

---

# Information Theory and Signal Processing (for Data Science)

*Lecture Notes — Fall 2019*

---

Gastpar, Telatar, Urbanke

EPFL



# Contents

<b>1</b>	<b>Introduction</b>	<b>7</b>
1.1	Foreword . . . . .	7
1.2	Acknowledgments . . . . .	8
1.3	Practical Information, Fall 2019, EPFL . . . . .	8
1.4	Lecture Schedule, Fall 2019, EPFL . . . . .	10
<b>2</b>	<b>Information Measures</b>	<b>11</b>
2.1	Hypothesis testing . . . . .	11
2.1.1	Hypothesis testing with repeated independent observations . . . . .	12
2.2	Large deviations via types . . . . .	13
2.2.1	Example . . . . .	18
2.3	Problems . . . . .	19
<b>3</b>	<b>Compression and Quantization</b>	<b>27</b>
3.1	Data compression . . . . .	27
3.2	Universal data compression with the Lempel-Ziv algorithm . . . . .	31
3.2.1	Finite state information lossless encoders . . . . .	34
3.3	Quantization . . . . .	36
3.4	Problems . . . . .	39
<b>4</b>	<b>Exponential Families and Maximum Entropy Distributions</b>	<b>45</b>
4.1	Definition . . . . .	46
4.2	Examples . . . . .	46
4.3	Convexity of $A(\theta)$ . . . . .	48
4.4	Derivatives of $A(\theta)$ . . . . .	49
4.5	Application to Parameter Estimation and Machine Learning . . . . .	49
4.6	Conjugate Priors . . . . .	50
4.7	Maximum Entropy Distributions . . . . .	51
4.8	Application To Physics . . . . .	53
4.9	I-Projections . . . . .	55
4.10	Relationship between $\theta$ and $\mathbb{E}[\phi(x)]$ . . . . .	56
4.10.1	The forward map $\nabla A(\theta)$ . . . . .	56
4.10.2	The backward map . . . . .	58

---

4.11 Problems . . . . .	58
<b>5 Multi-Arm Bandits . . . . .</b>	<b>61</b>
5.1 Introduction . . . . .	61
5.2 Some References . . . . .	62
5.3 Stochastic Bandits with a Finite Number of Arms . . . . .	62
5.3.1 Set-Up . . . . .	62
5.3.2 Explore then Exploit . . . . .	62
5.3.3 The Upper Confidence Bound Algorithm . . . . .	67
5.3.4 Information-theoretic Lower Bound . . . . .	71
5.4 Further Topics . . . . .	75
5.4.1 Asymptotic Optimality . . . . .	75
5.4.2 Adversarial Bandits . . . . .	75
5.4.3 Contextual Bandits . . . . .	77
5.5 Problems . . . . .	77
<b>6 Distribution Estimation, Property Testing and Property Estimation . . . . .</b>	<b>79</b>
6.1 Distribution Estimation . . . . .	79
6.1.1 Notation and Basic Task . . . . .	79
6.1.2 Empirical Estimator . . . . .	80
6.1.3 Loss Functions . . . . .	80
6.1.4 Min-Max Criterion . . . . .	80
6.1.5 Risk of Empirical Estimator in $\ell_2^2$ . . . . .	81
6.1.6 Risk of “Add Constant” Estimator in $\ell_2^2$ . . . . .	82
6.1.7 Matching lower bound for $\ell_2^2$ . . . . .	83
6.1.8 Risk in $\ell_1$ . . . . .	84
6.1.9 Risk in KL-Divergence . . . . .	85
6.1.10 The problem with the min-max formulation . . . . .	86
6.1.11 Competitive distribution estimation . . . . .	86
6.1.12 Multi-set genie estimator . . . . .	86
6.1.13 Natural Genie and Good-Turing Estimator . . . . .	86
6.2 Property Testing . . . . .	88
6.2.1 General Idea . . . . .	89
6.2.2 Testing Against a Uniform Distribution . . . . .	90
6.3 Property Estimation . . . . .	94
6.3.1 Entropy Estimation . . . . .	94
6.3.2 Symmetric Properties . . . . .	95
6.3.3 Profiles and Natural Estimators . . . . .	95
6.4 Problems . . . . .	95
<b>7 Information Measures and Generalization Error . . . . .</b>	<b>99</b>
7.1 Exploration Bias and Information Measures . . . . .	99
7.1.1 Definitions and Problem Statement . . . . .	99
7.1.2 $L_1$ -Distance Bound . . . . .	101

---

7.1.3	Mutual Information Bound . . . . .	102
7.2	Information Measures and Generalization Error . . . . .	105
7.2.1	Setup and Problem Statement . . . . .	105
7.2.2	Mutual Information Bound . . . . .	106
7.2.3	Differential Privacy Bound . . . . .	107
7.3	Problems . . . . .	107
<b>8</b>	<b>Elements of Statistical Signal Processing</b>	<b>109</b>
8.1	Optimum Estimation . . . . .	109
8.1.1	MMSE Estimation . . . . .	109
8.1.2	Linear MMSE Estimation . . . . .	110
8.2	Wiener Filtering, Smoothing, Prediction . . . . .	111
8.3	Adaptive Filters . . . . .	112
8.4	Problems . . . . .	113
<b>9</b>	<b>Signal Representation</b>	<b>117</b>
9.1	Review : Notions of Linear Algebra . . . . .	117
9.2	Fourier Representations . . . . .	121
9.2.1	DFT and FFT . . . . .	121
9.2.2	The Other Fourier Representations . . . . .	122
9.3	The Hilbert Space Framework for Signal Representation . . . . .	123
9.4	General Bases, Frames, and Time-Frequency Analysis . . . . .	125
9.4.1	The General Transform . . . . .	125
9.4.2	The Heisenberg Box Of A Signal . . . . .	126
9.4.3	The Uncertainty Relation . . . . .	127
9.4.4	The Short-time Fourier Transform . . . . .	127
9.5	Multi-Resolution Concepts and Wavelets . . . . .	130
9.5.1	The Haar Wavelet . . . . .	130
9.5.2	Multiresolution Concepts . . . . .	134
9.5.3	Wavelet Design — A Fourier Technique . . . . .	136
9.5.4	Wavelet Algorithms . . . . .	137
9.5.5	Wavelet Design — Further Considerations . . . . .	138
9.6	Data-adaptive Signal Representations . . . . .	140
9.6.1	<i>Example : word2vec</i> . . . . .	142
9.7	Problems . . . . .	143



# Chapter 1

## Introduction

### 1.1 Foreword

This is a set of lecture notes for a MS level class called “*Information Theory and Signal Processing (for Data Science)*” (COM-406) at EPFL. The class was first designed for the Fall Semester 2017.

Lausanne, Switzerland, September 2019

M. Gastpar, E. Telatar, R. Urbanke

## 1.2 Acknowledgments

The authors thank Dr. Ibrahim Issa for contributions to the class development as well as to the Lecture Notes.

## 1.3 Practical Information, Fall 2019, EPFL

### Instructors:

Michael Gastpar, [michael.gastpar@epfl.ch](mailto:michael.gastpar@epfl.ch), Office: INR 130

Emre Telatar, [emre.telatar@epfl.ch](mailto:emre.telatar@epfl.ch), Office: INR 1

Rüdiger Urbanke, [ruediger.urbanke@epfl.ch](mailto:ruediger.urbanke@epfl.ch), Office: INR 1

### Teaching Assistants:

Amedeo Esposito, [amedeo.esposito@epfl.ch](mailto:amedeo.esposito@epfl.ch), Office: INR 031

Pierre Quinton, [pierre.quinton@epfl.ch](mailto:pierre.quinton@epfl.ch), Office: INR 030

### Administrative Assistants:

Muriel Bardet, [muriel.bardet@epfl.ch](mailto:muriel.bardet@epfl.ch), Office: INR 137

France Faille, [france.faille@epfl.ch](mailto:france.faille@epfl.ch), Office: INR 311

### Class Meetings:

Mondays, 9:15-11:00, INM 200

Fridays, 8:15-10:00, GR A3 31

Fridays, 10:15-12:00, GR A3 31 (Exercises)

**Course Web Page:** We will use <https://ipg.epfl.ch/cms/lang/en/pid/147664>

### Official Prerequisites:

COM-300 “Modèles stochastiques pour les communications” (or equivalent)

COM-303 “Signal processing for communications” (or equivalent)

**Homework:** Some Homework will be graded....

**Final Exam:** The Final Exam for the course will take place at some point between January 13 and February 1, 2019. The precise date will be decided by EPFL some time in November 2019.

**Grading:** If you do not hand in your final exam your overall grade will be NA. Otherwise, your grade will be determined based on the following weighted average: 10% for the Homework, 90% for the Final Exam.





## 1.4 Lecture Schedule, Fall 2019, EPFL

Date	Topics	Reading
Sept 20	General Introduction ; Review Probability <i>Exercise: Review Session (Probability)</i>	Handout
Sept 23 Sept 27  Sept 30 Oct 4	Basic Information Measures  <i>Exercise: HW 1</i>  <i>Exercise: HW 1</i>	Chapter 2
Oct 7 Oct 11  Oct 14 Oct 18	Compression and Quantization Compression and Quantization <i>Exercise: HW 4</i> Compression and Quantization Compression and Quantization <i>Exercise: HW 4</i>	Chapter 3
Oct 21 Oct 25	Exponential families ; Max Entropy problems Boltzmann distribution ; Exponential families <i>Exercise: HW 3</i>	Chapter 4
Oct 28 Nov 1  Nov 4 Nov 8	Multi-armed Bandits : Explore & Exploit Multi-armed Bandits : UCB algorithm <i>Exercise: HW 5</i> Multi-armed Bandits : Converse bound Multi-armed Bandits : Variations <i>Exercise: HW 5</i>	Chapter 5
Nov 11 Nov 15  Nov 18 Nov 22	Distribution Estimation ; Property Testing and Estimation Distribution Estimation ; Property Testing and Estimation <i>Exercise: HW 7</i> Distribution Estimation ; Property Testing and Estimation Distribution Estimation ; Property Testing and Estimation <i>Exercise: HW 7</i>	Chapter 6
Nov 25 Nov 29  Dec 2 Dec 6	Information measures, Learning and Generalization Information measures, Learning and Generalization <i>Exercise: HW 2</i> Optimum Detection and Estimation ; MMSE Wiener Filter, LMS Adaptive Filter <i>Exercise: HW 2</i>	Chapter 7  Chapter 8
Dec 9 Dec 13  Dec 16 Dec 20	Review Linear Algebra (SVD, Eckart–Young) ; Fourier Sparse Fourier ; Hilbert space perspective <i>Exercise: HW 6</i> Time–Frequency ; Wavelets Wavelets ; Data-adaptive Signal Representations <i>Exercise: HW 6</i>	Chapter 9

## Chapter 2

# Information Measures

### 2.1 Hypothesis testing

Consider the problem of deciding which of two hypotheses, hypothesis 0 or hypothesis 1, is true, based on an observation  $U$ . The observation  $U$  is a random variable taking values in an alphabet  $\mathcal{U}$  — a finite set of  $K = |\mathcal{U}|$  letters — and under hypothesis  $j$  it has distribution  $P_j$ . To avoid trivial cases we will assume that for each  $u \in \mathcal{U}$  both  $P_0(u)$  and  $P_1(u)$  are strictly positive. Otherwise, if we observe a  $u$  with, say,  $P_0(u) = 0$ , we would know for sure that hypothesis 1 is true.

A deterministic decision rule associates to each  $u \in \mathcal{U}$  a binary value — i.e., the rule is a function  $\phi : \mathcal{U} \rightarrow \{0, 1\}$  — and we decide in favor of hypothesis  $\phi(u)$  if the observation  $U$  equals  $u$ . In general, we will allow for randomized decision rules: such a rule is characterized by a function  $\phi : \mathcal{U} \rightarrow [0, 1]$  that associates to each  $u \in \mathcal{U}$  a value in the *interval*  $[0, 1]$ , that gives the probability of deciding in favor of hypothesis 1. If our observation  $U$  equals  $u$ , we flip a coin that comes heads with probability  $\phi(u)$  and tails with probability  $1 - \phi(u)$ , and decide accordingly: 1 if heads, 0 if tails. We will identify a decision rule with the function  $\phi$ .

In this set up there are two kinds of error: deciding 1 when the true hypothesis is 0, and deciding 0 when the true hypothesis is 1. Letting  $\pi_\phi(i|j)$  for rule  $\phi$  denote the probability of deciding  $i$  when the truth is  $j$ , we see that

$$\pi_\phi(0|1) = \sum_u P_1(u)[1 - \phi(u)], \quad \pi_\phi(1|0) = \sum_u P_0(u)\phi(u).$$

Given  $P_0$  and  $P_1$  and a positive real number  $\eta > 0$ , let  $\Phi_\eta$  to be the set of decision rules  $\phi$  of the form

$$\phi(u) = \begin{cases} 1 & \text{if } P_1(u) > \eta P_0(u) \\ 0 & \text{if } P_1(u) < \eta P_0(u). \end{cases} \quad (2.1)$$

Note that if there is no  $u$  for which  $P_1(u) = \eta P_0(u)$ , the test  $\phi$  is uniquely specified and  $\Phi_\eta$  contains only this test.

**Lemma 2.1.** *The rules in  $\Phi_\eta$  are minimizers of  $\pi(0|1) + \eta\pi(1|0)$ .*

*Proof.* For any rule  $\phi \in \Phi_\eta$ , as a consequence of (2.1), for every  $u \in \mathcal{U}$

$$P_1(u)[1 - \phi(u)] + \eta P_0(u)\phi(u) = \min\{P_1(u), \eta P_0(u)\}.$$

Thus for any rule  $\phi \in \Phi_\eta$

$$\pi_\phi(0|1) + \eta\pi_\phi(1|0) = \sum_u P_1(u)[1 - \phi(u)] + \eta P_0(u)\phi(u) = \sum_u \min\{P_1(u), \eta P_0(u)\}.$$

Suppose now  $\psi$  is any decision rule. The lemma follows by noting that

$$\pi_\psi(0|1) + \eta\pi_\psi(1|0) = \sum_u P_1(u)[1 - \psi(u)] + \eta P_0(u)\psi(u) \geq \sum_u \min\{P_1(u), \eta P_0(u)\}. \quad \square$$

**Theorem 2.2.** For any  $\alpha \in [0, 1]$ , (i) there is a rule  $\phi$  of the form (2.1) such that  $\pi_\phi(0|1) = \alpha$ , and (ii) for any decision rule  $\psi$  either  $\pi_\psi(0|1) \geq \pi_\phi(0|1)$  or  $\pi_\psi(1|0) \geq \pi_\phi(1|0)$ .

*Proof.* Assertion (ii) follows from the lemma above: a  $\psi$  that violates both the inequalities would contradict the lemma. It thus suffices to prove (i), the existence of a rule  $\phi$  of the form (2.1) with  $\pi_\phi(0|1) = \alpha$ . To that end, define  $L(u) = P_1(u)/P_0(u)$ , and label the elements of  $\mathcal{U}$  as  $\mathcal{U} = \{u_1, \dots, u_K\}$  such that  $L(u_1) \geq L(u_2) \geq \dots \geq L(u_K)$ . Now define,  $a_i = \sum_{j=1}^i P_1(u_j)$  for  $i = 0, \dots, K$ . We then have  $0 = a_0 < a_1 < \dots < a_K = 1$ . Given  $0 \leq \alpha \leq 1$ , we can find  $1 \leq i \leq K$  for which  $a_{i-1} \leq 1 - \alpha \leq a_i$ , so that  $1 - \alpha = (1 - \rho)a_{i-1} + \rho a_i$  for some  $\rho \in [0, 1]$ . Then, the rule

$$\phi(u) = \begin{cases} 1 & u \in \{u_1, \dots, u_{i-1}\} \\ \rho & u = u_i \\ 0 & u \in \{u_{i+1}, \dots, u_K\} \end{cases}$$

is of the form (2.1) with  $\eta = L(u_i)$ , and  $\pi_\phi(0|1) = \alpha$ .  $\square$

Rules of the form (2.1) are based on a *likelihood ratio test*: they compare the likelihood ratio  $P_1(u)/P_0(u)$  to a threshold  $\eta$  to make a decision. If the likelihood ratio is larger than the threshold, decide 1; if less, decide 0. Equivalently one may compare the *log likelihood ratio*,  $\log(P_1(u)/P_0(u))$  to the threshold  $\log \eta$ .

The theorem stated just above shows the dominant nature of likelihood ratio tests in making decisions: given any decision rule  $\psi$ , we can find a (log) likelihood ratio test  $\phi$  which is ‘as good or better’ — in the sense that the two error probabilities satisfy  $\pi_\phi(0|1) \leq \pi_\psi(0|1)$  and  $\pi_\phi(1|0) \leq \pi_\psi(1|0)$ .

### 2.1.1 Hypothesis testing with repeated independent observations

Suppose now that we make repeated independent observations of  $U$ . That is, we observe a sequence  $U_1, \dots, U_n$  of independent and identically distributed (i.i.d.) random variables, with common distribution  $P_i$  under hypothesis  $i$ , for  $i = 0, 1$ .

The log likelihood ratio tests for this scenario are of the form

$$\phi(u_1, \dots, u_n) = \begin{cases} 1 & \Lambda_n(u_1, \dots, u_n) > t \\ 0 & \Lambda_n(u_1, \dots, u_n) < t \end{cases}$$

where

$$\Lambda_n(u_1, \dots, u_n) = \frac{1}{n} \sum_{i=1}^n \log \frac{P_1(u_i)}{P_0(u_i)}$$

is the normalized log likelihood ratio for the observation  $u_1, \dots, u_n$ .

If hypothesis 0 is true, then  $U_1, \dots, U_n$  are i.i.d. random variables with distribution  $P_0$ , and, by the law of large numbers

$$\Lambda_n(U_1, \dots, U_n) \rightarrow E_0 \left[ \log \frac{P_1(U_1)}{P_0(U_1)} \right] = \sum_u P_0(u) \log \frac{P_1(u)}{P_0(u)}$$

as  $n$  gets large. In the expression above, the subscript 0 to the expectation operator indicates that the expectation is taken with the distribution of the  $U_i$ 's given by  $P_0$ . Similarly, if hypothesis 1 is true,

$$\Lambda_n(U_1, \dots, U_n) \rightarrow E_1 \left[ \log \frac{P_1(U_1)}{P_0(U_1)} \right] = \sum_u P_1(u) \log \frac{P_1(u)}{P_0(u)}$$

as  $n$  gets large.

In the following section we will show that for any two probability distributions  $P$  and  $Q$  on an alphabet  $\mathcal{U}$ , the quantity  $D(P\|Q) = \sum_u P(u) \log[P(u)/Q(u)]$  is non-negative, and equals zero if and only if  $P = Q$ .

Thus, as  $n$  gets large  $\Lambda_n(U_1, \dots, U_n)$  concentrates around  $-D(P_0\|P_1) \leq 0$  under hypothesis 0 and, concentrates around  $D(P_1\|P_0) \geq 0$  under hypothesis 1. One expects that the threshold  $t$  will be chosen to lie between  $-D(P_0\|P_1)$  and  $D(P_1\|P_0)$  so that under either hypothesis, making a wrong decision becomes a large deviations event — an event that the empirical average of a collection of i.i.d. random variables deviates significantly from its expected value.

We will shortly see that  $D(\cdot\|\cdot)$  plays a central role in estimating the probabilities of large deviation events.

## 2.2 Large deviations via types

Let  $\Pi := \Pi(\mathcal{U})$  denote the set of all probability distributions on  $\mathcal{U}$ . With  $K = |\mathcal{U}|$ , we can identify  $\Pi$  with the *simplex* in  $\mathbb{R}^K$ : the set of all  $(p_1, \dots, p_K) \in \mathbb{R}^K$  with  $\sum_k p_k = 1$ , and  $p_k \geq 0$ .

**Definition 2.1.** For  $P$  and  $Q$  in  $\Pi$ , we call  $D(P\|Q) = \sum_u P(u) \log[P(u)/Q(u)]$  the *divergence of  $P$  from  $Q$* .

In the sum above, we skip the terms with  $P(u) = 0$ , and we set  $D(P\|Q) = +\infty$  if there is a  $u$  for which  $P(u) > 0$  and  $Q(u) = 0$ .

**Definition 2.2.** Given two alphabets  $\mathcal{U}$  and  $\mathcal{V}$ , a *probability kernel from  $\mathcal{U}$  to  $\mathcal{V}$*  is a matrix  $W = [W(v|u) : u \in \mathcal{U}, v \in \mathcal{V}]$  such that  $W(v|u) \geq 0$ , and for each  $u \in \mathcal{U}$ ,  $\sum_v W(v|u) = 1$ . We will write  $W : \mathcal{U} \rightarrow \mathcal{V}$  to indicate that  $W$  is such a kernel. The set of probability kernels describes all possible conditional probabilities on  $\mathcal{V}$ , conditional on elements of  $\mathcal{U}$ .

**Lemma 2.3.** Given  $P$  and  $Q$  in  $\Pi(\mathcal{U})$  and  $W : \mathcal{U} \rightarrow \mathcal{V}$ , let  $\tilde{P}(v) = \sum_u P(u)W(v|u)$  and  $\tilde{Q}(v) = \sum_u Q(u)W(v|u)$ . Then  $\tilde{P}$  and  $\tilde{Q}$  are in  $\Pi(\mathcal{V})$ , and

$$D(\tilde{P}||\tilde{Q}) \leq D(P||Q).$$

The inequality is strict, unless  $Q(u)/P(u) = \tilde{Q}(v)/\tilde{P}(v)$  for all  $u, v$  with  $P(u)W(v|u) > 0$ .

*Proof.* That  $\tilde{P}$  and  $\tilde{Q}$  are probability distributions is an easy consequence of  $W$  being a probability kernel. To prove the claimed inequality between the divergences let us first show that  $\log$  is a strictly concave function. I.e., for any non-negative  $\lambda_1, \dots, \lambda_K$  for which  $\sum_k \lambda_k = 1$ , and any positive  $x_1, \dots, x_K$ , we have, with  $\bar{x} = \sum_k \lambda_k x_k$ ,

$$\sum_k \lambda_k \log x_k \leq \log \bar{x},$$

and equality happens if and only if for all  $k$  with  $\lambda_k > 0$ , we have  $x_k = \bar{x}$ . It suffices to prove this statement with  $\ln$  instead of  $\log$ . To that end, first note that with  $f(x) = \ln x$  we have  $f'(x) = 1/x$  and  $f''(x) = -1/x^2 < 0$ . Thus, Taylor expansion of  $\ln x$  around 1 yields  $\ln x = (x - 1) - (x - 1)^2/(2\xi^2)$  for some  $\xi$  between 1 and  $x$ , and we see that  $\ln x \leq x - 1$ , with equality if and only if  $x = 1$ . Consequently

$$\sum_k \lambda_k \ln x_k - \ln \bar{x} = \sum_k \lambda_k \ln [x_k/\bar{x}] \leq \sum_k \lambda_k [x_k/\bar{x} - 1] = 1 - 1 = 0,$$

with the inequality being strict if there is a  $k$  for which  $\lambda_k > 0$  and  $x_k \neq \bar{x}$ .

Having thus proved the strict concavity of  $\log$ , now observe (with  $P(u)W(v|u)$ 's cast in the role of  $\lambda_k$ 's) that

$$\begin{aligned} D(\tilde{P}||\tilde{Q}) - D(P||Q) &= \sum_v \tilde{P}(v) \log \frac{\tilde{P}(v)}{\tilde{Q}(v)} - \sum_u P(u) \log \frac{P(u)}{Q(u)} \\ &= \sum_{u,v} W(v|u) P(u) \log \frac{\tilde{P}(v)}{\tilde{Q}(v)} - \sum_{u,v} W(v|u) P(u) \log \frac{P(u)}{Q(u)} \\ &= \sum_{u,v} W(v|u) P(u) \log \frac{\tilde{P}(v) Q(u)}{\tilde{Q}(v) P(u)} \\ &\leq \log \left[ \sum_{u,v} W(v|u) \frac{\tilde{P}(v) Q(u)}{\tilde{Q}(v)} \right] = \log \left[ \sum_v \tilde{P}(v) \right] = 0. \quad \square \end{aligned}$$

**Corollary 2.4.**  $D(P||Q) \geq 0$  with equality if and only if  $P = Q$ .

*Proof.* Take  $\mathcal{V} = \{0\}$  and set  $W(0|u) = 1$ . Then  $\tilde{P}(0) = \tilde{Q}(0) = 1$  and  $D(\tilde{P}||\tilde{Q}) = 0$ .  $\square$

**Corollary 2.5.**  $D(P\|Q)$  is a convex function of the pair  $(P, Q)$ .

*Proof.* Suppose  $P_0, Q_0, P_1, Q_1$  are in  $\Pi(\mathcal{U})$  and suppose  $0 \leq \lambda \leq 1$ . We need to show that

$$D((1 - \lambda)P_0 + \lambda P_1 \| (1 - \lambda)Q_0 + \lambda Q_1) \leq (1 - \lambda)D(P_0 \| Q_0) + \lambda D(P_1 \| Q_1).$$

To that end consider the distributions  $P$  and  $Q$  on the set  $\{0, 1\} \times \mathcal{U}$  with

$$P(z, u) = \begin{cases} (1 - \lambda)P_0(u) & \text{if } z = 0 \\ \lambda P_1(u) & \text{if } z = 1, \end{cases} \quad \text{and} \quad Q(z, u) = \begin{cases} (1 - \lambda)Q_0(u) & \text{if } z = 0 \\ \lambda Q_1(u) & \text{if } z = 1, \end{cases}$$

Consider also the channel  $W : \{0, 1\} \times \mathcal{U} \rightarrow \mathcal{U}$  with  $W(u'|(z, u)) = \mathbb{1}\{u' = u\}$ . It is easily checked that  $D(P\|Q) = (1 - \lambda)D(P_0\|Q_0) + \lambda D(P_1\|Q_1)$  and also that

$$\tilde{P} = (1 - \lambda)P_0 + \lambda P_1 \quad \text{and} \quad \tilde{Q} = (1 - \lambda)Q_0 + \lambda Q_1.$$

The conclusion now follows from Lemma 2.3.  $\square$

**Definition 2.3.** For  $P$  in  $\Pi$ , we call  $H(P) = \sum_u P(u) \log[1/P(u)]$  the *entropy* of  $P$ .

**Lemma 2.6.**  $0 \leq H(P) \leq \log |\mathcal{U}|$ , with equality on the left if and only if there is a  $u_0 \in \mathcal{U}$  with  $P(u_0) = 1$ , and equality on the right if and only if  $P$  is the uniform distribution on  $\mathcal{U}$ .

*Proof.* The non-negativity of  $H(P)$  follows from  $P(u) \geq 0$  and  $\log[1/P(u)] \geq 0$ , so that each term in the sum defining  $H(P)$  is non-negative. Moreover, the sum equals zero only if each term is zero, which yields the condition for  $H(P)$  to equal 0. The right hand side inequality and the condition for equality follows from noting that  $\log |\mathcal{U}| - H(P) = D(P\|\text{unif}_{\mathcal{U}})$  where  $\text{unif}_{\mathcal{U}}$  is the uniform distribution on  $\mathcal{U}$  with  $\text{unif}_{\mathcal{U}}(u) = 1/|\mathcal{U}|$ .  $\square$

**Notation.** We will use  $x^n$  as a short-hand to denote the sequence  $(x_1, \dots, x_n)$ .

**Notation.** For  $P \in \Pi$ , we will let  $P^n$  denote the distribution of the i.i.d. sequence  $U^n$ , each  $U_i$  with distribution  $P$ . I.e.,  $P^n(u^n) = \prod_{i=1}^n P(u_i)$ .

Any *empirical average* based on a sequence  $(u_1, \dots, u_n)$  — a quantity of the form  $\frac{1}{n} \sum_{i=1}^n f(u_i)$  for some function  $f : \mathcal{U} \rightarrow \mathbb{R}$  — depends on  $u^n = (u_1, \dots, u_n)$  only via its *empirical distribution*:

**Definition 2.4.** The *empirical distribution* (also called the *type*) of a sequence  $u^n \in \mathcal{U}^n$  is the probability distribution  $\hat{P}$  on  $\mathcal{U}$  defined by

$$\hat{P}(u) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{u_i = u\}, \quad u \in \mathcal{U}.$$

We will also write  $\hat{P} = \hat{P}_{u^n}$  to emphasize that  $\hat{P}$  is the type of the sequence  $u^n$ .

With  $\hat{P}$  denoting the type of  $u^n$ , observe that  $\frac{1}{n} \sum_{i=1}^n f(u_i)$  equals  $\sum_u \hat{P}(u) f(u)$ . As a particular case  $\Lambda_n(u^n)$  is an empirical average with  $f(u) = \log[P_1(u)/P_0(u)]$ .

Furthermore, if  $U^n = (U_1, \dots, U_n)$  is a collection of i.i.d. random variables with common distribution  $P$ , then  $\Pr(U^n = u^n) = P^n(u^n)$ , and

$$\frac{1}{n} \log P^n(u^n) = \frac{1}{n} \sum_{i=1}^n \log P(u_i) = \sum_u \hat{P}_{u^n}(u) \log P(u) = -H(\hat{P}_{u^n}) - D(\hat{P}_{u^n} \| P).$$

We state this formally as:

**Lemma 2.7.** *For  $P \in \Pi$  and  $Q$  denoting the type of  $u^n$ ,  $P^n(u^n) = \exp[-n(D(Q \| P) + H(Q))]$ .*

The set of types of sequences of length  $n$  form a subset  $\Pi_n$  of  $\Pi$ :

$$\Pi_n = \{P \in \Pi : nP(u) \text{ is an integer for all } u \in \mathcal{U}\}.$$

**Lemma 2.8.** *With  $K = |\mathcal{U}|$ , we have  $|\Pi_n| = \binom{n+K-1}{K-1} \leq (n+1)^K$ .*

So, even though the number of sequences of length  $n$  is exponential in  $n$ , the number of types is only polynomial in  $n$ .

*Proof.* Without loss of generality, let  $\mathcal{U} = \{1, \dots, K\}$ . An element  $P$  of  $\Pi_n$  can be identified with a  $K$ -tuple of non-negative integers  $(n_1, \dots, n_K)$  that sum to  $n$  via  $n_i = nP(i)$ . But the set of such integers are in one-to-one correspondence with set of binary sequences that contain exactly  $n$  ones and  $K-1$  zeros: with  $1^m$  denoting a repetition of  $m$  1's,  $(n_1, \dots, n_K)$  is identified with the sequence  $1^{n_1}01^{n_2}0 \dots 01^{n_K}$ . This yields the size of  $\Pi_n$  as the binomial coefficient. The upper bound on  $|\Pi_n|$  follows from noting that each  $n_i$  can take on  $n+1$  possible values (0 up to  $n$ ).  $\square$

**Remark.** It should be clear that  $\cup_{n>0} \Pi_n$  is dense in  $\Pi$ . Indeed, for any  $P \in \Pi$  we can find  $P_n \in \Pi_n$  with  $\|P_n - P\|_\infty := \max_u |P(u) - P_n(u)| < 1/n$ . To see this, suppose  $\mathcal{U} = \{1, \dots, K\}$ , and let  $nP(i) = n_i + f_i$  with  $n_i = \lfloor nP(i) \rfloor$  and  $0 \leq f_i < 1$ . Since  $\sum_i nP(i) = n$ , the sum  $r = \sum_i f_i$  is an integer between 0 and  $K-1$ . Assume  $f_1 \geq \dots \geq f_K$  and define  $P_n$  by  $nP_n(i) = n_i + \mathbb{1}\{i \leq r\}$ . One can check that  $\|P_n - P\|_\infty \leq \frac{1}{n} \frac{K-1}{K}$  and  $\|P_n - P\|_1 \leq \frac{K}{2n}$ .

**Definition 2.5.** For  $Q \in \Pi_n(\mathcal{U})$ , define  $T^n(Q) = \{u^n \in \mathcal{U}^n : \hat{P}_{u^n} = Q\}$ , i.e., the set of all sequences of length  $n$  with type  $Q$ .

For  $Q \in \Pi_n$ , the set  $T^n(Q)$  is the set of all sequences of length  $n$  with exactly  $n_u = nQ(u)$  occurrences of the letter  $u \in \mathcal{U}$ . Thus, with  $\mathcal{U} = \{1, \dots, K\}$ ,

$$|T^n(Q)| = \binom{n}{nQ(1) \dots nQ(K)} = \frac{n!}{\prod_u (nQ(u))!}.$$

**Lemma 2.9.** *For  $P \in \Pi_n$ ,  $Q \in \Pi_n$  we have  $P^n(T^n(Q)) \leq P^n(T^n(P))$ , with equality if and only if  $Q = P$ .*



*Proof.* We had already noted in Lemma 2.7 that the value of  $P^n(u^n)$  is determined by the type of  $u^n$  and is thus constant over the set  $T^n(Q)$ ; let  $c(P, Q)$  denote this constant. Thus,  $P^n(T^n(Q)) = |T^n(Q)|c(P, Q)$ .

The lemma states that  $P$  is a global maximizer of  $Q \in \Pi_n \mapsto P^n(T^n(Q))$ . We will prove the lemma by showing that any  $Q \neq P$  cannot even be a local maximizer. Let  $\mathcal{U} = \{1, \dots, K\}$ , let  $n_k = nP(k)$  and  $m_k = nQ(k)$ , and suppose  $Q \neq P$ . We may assume that for any  $k$  with  $n_k = 0$  we have  $m_k = 0$ : otherwise  $c(P, Q) = 0$  and clearly  $Q$  is not a maximizer.

Since  $\sum_k m_k$  and  $\sum_k n_k$  are both equal to  $n$ , and  $Q \neq P$ , there will be indices  $k_1$  and  $k_2$  such that  $m_{k_1} > n_{k_1}$  and  $m_{k_2} < n_{k_2}$ . (By the assumption above  $n_{k_1} > 0$ .) Assume, without loss of generality, that  $k_1 = 1$  and  $k_2 = 2$ . Consider now the type  $\tilde{Q}$  for which  $n\tilde{Q}(1) = m_1 - 1$ ,  $n\tilde{Q}(2) = m_2 + 1$ , and  $n\tilde{Q}(k) = m_k$  for  $k > 2$ . Observe that for any  $u^n \in T^n(Q)$  and  $\tilde{u}^n \in T^n(\tilde{Q})$ ,

$$\frac{c(P, \tilde{Q})}{c(P, Q)} = \frac{P^n(\tilde{u}^n)}{P^n(u^n)} = \frac{P(2)}{P(1)} = \frac{n_2}{n_1}.$$

Thus,

$$\frac{P^n(T^n(\tilde{Q}))}{P^n(T^n(Q))} = \frac{|T^n(\tilde{Q})|}{|T^n(Q)|} \frac{n_2}{n_1} = \frac{m_1}{n_1} \frac{n_2}{m_2 + 1}.$$

But  $m_1 > n_1$  and  $n_2 \geq m_2 + 1$ . Thus  $P^n(T^n(\tilde{Q})) > P^n(T^n(Q))$ , so  $Q$  is not a global maximum.  $\square$

**Corollary 2.10.** For  $P \in \Pi_n$ ,  $|\Pi_n|^{-1} \leq P^n(T^n(P)) \leq 1$ .

*Proof.* The right hand inequality is trivial. For the left, note that the collection  $\{T^n(Q) : Q \in \Pi_n\}$  is a partition of  $\mathcal{U}^n$  into disjoint sets. (A given sequence  $u^n$  has one, and exactly one, type  $\hat{P}_{u^n}$ .) Consequently

$$1 = \sum_{Q \in \Pi_n} P^n(T^n(Q)) \leq \sum_{Q \in \Pi_n} P^n(T^n(P)) = |\Pi_n| P^n(T^n(P)). \quad \square$$

**Corollary 2.11.** For  $P \in \Pi_n$ ,  $|\Pi_n|^{-1} \exp(nH(P)) \leq |T^n(P)| \leq \exp(nH(P))$ .

*Proof.* By Lemma 2.7, for  $u^n \in T^n(P)$ ,  $P^n(u^n) = \exp(-nH(P))$ . The conclusion follows from the previous corollary by noting that  $P^n(T^n(P)) = |T^n(P)| \exp(-nH(P))$ .  $\square$

**Lemma 2.12.** For  $P \in \Pi$ ,  $Q \in \Pi_n$ ,  $|\Pi_n|^{-1} \exp(-nD(Q\|P)) \leq P^n(T^n(Q)) \leq \exp(-nD(Q\|P))$ .

*Proof.* By Lemma 2.7, for  $u^n \in T^n(Q)$ ,  $P^n(u^n) = \exp[-n(H(Q) + D(Q\|P))]$ . Using Corollary 2.11 to bound the size of  $T^n(Q)$  yields the bounds on  $P^n(T^n(Q))$ . Note that just as in Lemma 2.7,  $P$  is assumed to be in  $\Pi$ , not necessarily in  $\Pi_n$ .  $\square$

**Theorem 2.13.** Suppose  $U_1, U_2, \dots$  is an i.i.d. sequence of random variables with common distribution  $P$ . Let  $\hat{P}_n$  denote the type of  $U^n$ . Suppose  $A \subset \Pi$  is a set of distributions with  $G \subset A \subset F$  with  $G$  open and  $F$  closed. Then

$$-\inf_{Q \in G} D(Q\|P) \leq \liminf_{n \rightarrow \infty} \frac{1}{n} \log \Pr(\hat{P}_n \in A) \leq \limsup_{n \rightarrow \infty} \frac{1}{n} \log \Pr(\hat{P}_n \in A) \leq -\min_{Q \in F} D(Q\|P).$$

If  $A$  is such that its closure is equal to the closure of its interior, then,

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log \Pr(\hat{P}_n \in A) = - \inf_{Q \in A} D(Q \| P).$$

*Proof.* The last claim follows from setting  $F$  to be the closure of  $A$  and  $G$  to be the interior of  $A$ . It thus suffices to prove the upper and lower bounds.

For the upper bound, let  $D_* = \min_{Q \in F} D(Q \| P)$ . Note that

$$\Pr(\hat{P}_n \in A) = \sum_{Q \in A \cap \Pi_n} \Pr(\hat{P}_n = Q) = \sum_{Q \in A \cap \Pi_n} P^n(T^n(Q)) \leq \sum_{Q \in A \cap \Pi_n} \exp(-nD(Q \| P)).$$

Since each term in the sum is upper bounded by  $\exp(-nD_*)$ , and since there are at most  $|\Pi_n|$  terms,

$$\frac{1}{n} \log \Pr(\hat{P}_n \in A) \leq -D_* + \frac{1}{n} \log |\Pi_n|,$$

and the upper bound follows by noting that  $\lim_{n \rightarrow \infty} \frac{1}{n} \log |\Pi_n| = 0$ .

For the lower bound, let  $D^* = \inf_{Q \in G} D(Q \| P)$ , fix  $\epsilon > 0$  and find  $Q_0 \in G$  with  $D(Q_0 \| P) < D^* + \epsilon$ . Since  $G$  is open and  $Q \mapsto D(Q \| P)$  is continuous, we can find  $\delta = \delta(\epsilon, Q_0) > 0$  such that whenever  $Q$  satisfies  $\|Q - Q_0\|_\infty < \delta$ , we will have (i)  $Q \in G$ , and (ii)  $|D(Q \| P) - D(Q_0 \| P)| < \epsilon$ . By Remark 2.2, for  $n > 1/\delta$ , we can find  $Q_n \in \Pi_n$  such that  $\|Q_n - Q_0\|_\infty < \delta$ . Consequently, such a  $Q_n$  belongs to  $G$  (and thus to  $A$ ), and  $D(Q_n \| P) < D^* + 2\epsilon$ . So, for  $n > 1/\delta$  and sufficiently large to ensure that  $\frac{1}{n} \log |\Pi_n| < \epsilon$ ,

$$\frac{1}{n} \log \Pr(\hat{P}_n \in A) \geq \frac{1}{n} \log \Pr(\hat{P}_n = Q_n) \geq -\frac{1}{n} \log |\Pi_n| - D(Q_n \| P) > -D^* - 3\epsilon.$$

As  $\epsilon > 0$  is arbitrary, the lower bound follows.  $\square$

### 2.2.1 Example

Consider the setting of hypothesis testing with repeated independent observations. Let  $f(u) = \log[P_1(u)/P_0(u)]$ , and define  $A = \{Q \in \Pi : \sum_u Q(u)f(u) \geq t\}$ , and  $B = \{Q \in \Pi : \sum_u Q(u)f(u) \leq t\}$ . With these, the events

$$\{\hat{P}_n \in A\} \quad \text{and} \quad \{\hat{P}_n \in B\}$$

are exactly the events  $\{\Lambda_n(U_1, \dots, U_n) \geq t\}$  and  $\{\Lambda_n(U_1, \dots, U_n) \leq t\}$ . Furthermore,  $A$  and  $B$  are both equal to the closure of their interiors. Thus,

$$D_0 = \min_{Q \in A} D(Q \| P_0) \quad \text{and} \quad D_1 = \min_{Q \in B} D(Q \| P_1)$$

are the exponents of the rates of decay of the two error probabilities. For the ‘interesting case’ when  $-D(P_0 \| P_1) < t < D(P_1 \| P_0)$ , we see that  $P_0 \notin A$  and  $P_1 \notin B$ . Consequently, both  $D_0$  and  $D_1$  are positive.

One can guess the form of the minimizers by forming the Lagrangians  $J_0$  and  $J_1$  for the minimization problems with Lagrange multiplier  $s$ . For the first minimization this gives

$$J_0(Q, s) = D(Q\|P_0) - s \sum_u Q(u) \log \frac{P_1(u)}{P_0(u)}.$$

Notice that  $J_0(Q, s) = D(Q\|P_s) - \log Z(s)$  where  $Z(s) = \sum_u P_0(u)^{1-s} P_1(u)^s$  and  $P_s(u) = P_0(u)^{1-s} P_1(u)^s / Z(s)$ . Thus,  $P_s$  is the minimizer for  $J_0(\cdot, s)$ , and this suggests that the minimizer for  $D_0$  will be among  $\{P_s : s \in \mathbb{R}\}$ . The same conclusion holds for the minimizer for  $D_1$ .

One might guess that one should choose  $s = s^*$  so that  $Q = P_{s^*}$  is on the boundary of  $A$  and  $B$ , i.e., to find the  $s$  for which  $\sum_u P_s(u) \log[P_1(u)/P_0(u)] = t$ . As  $s$  ranges from 0 to 1,  $\sum_u P_s(u) \log[P_1(u)/P_0(u)]$  ranges from  $-D(P_0\|P_1)$  to  $D(P_1\|P_0)$ , so we see that  $s^*$  will be in the open interval  $(0, 1)$ .

Having made our guess, let us now verify that  $Q^* = P_{s^*}$  is indeed the minimizer for finding  $D_0$ . First observe that

$$\log \frac{Q^*(u)}{P_0(u)} = s^* \log \frac{P_1(u)}{P_0(u)} - \log Z(s^*),$$

so that  $D(Q^*\|P_0) = s^*t - \log Z(s^*)$ . On the other hand, for any  $Q \in A$ ,

$$\sum_u Q(u) \log \frac{Q^*(u)}{P_0(u)} \geq s^*t - \log Z(s^*) = D(Q^*\|P_0).$$

Thus

$$D(Q\|P_0) = \sum_u Q(u) \log \frac{Q(u)}{P_0(u)} = D(Q\|Q^*) + \sum_u Q(u) \log \frac{Q^*(u)}{P_0(u)} \geq D(Q^*\|P_0),$$

verifying that  $Q^*$  minimizes  $D(Q\|P_0)$  among all  $Q$  in  $A$ , and  $D_0 = D(Q^*\|P_0)$ . An analogous computation shows that for the same  $Q^*$ , and any  $Q$  in  $B$  we have  $D(Q\|P_1) \geq D(Q^*\|P_1)$ .

If we summarize the conclusions above parametrically in  $0 \leq s \leq 1$ , we get

$$D_0 = D(P_s\|P_0), \quad D_1 = D(P_s\|P_1), \quad t = \sum_u P_s(u) \log \frac{P_1(u)}{P_0(u)}.$$

As  $s$  changes from 0 to 1, it is easy to check that  $D_0$  increases from 0 to  $D(P_1\|P_0)$ , and  $D_1$  decreases from  $D(P_0\|P_1)$  to 0. One natural choice for  $s$  (and thus the threshold  $t$ ) is the choice that makes  $D_0 = D_1$  so that  $\min\{D_0, D_1\}$  is as large as possible.

## 2.3 Problems

**Problem 2.1** (Random Variables). Let  $X$  and  $Y$  be discrete random variables defined on some probability space with a joint pmf  $p_{XY}(x, y)$ . Let  $a, b \in \mathbb{R}$  be fixed.

(a) Prove that  $\mathbb{E}[aX + bY] = a\mathbb{E}[X] + b\mathbb{E}[Y]$ . Do not assume independence.

- (b) Prove that if  $X$  and  $Y$  are independent random variables, then  $\mathbb{E}[X \cdot Y] = \mathbb{E}[X] \cdot \mathbb{E}[Y]$ .  
 (c) Assume that  $X$  and  $Y$  are not independent. Find an example where  $\mathbb{E}[X \cdot Y] \neq \mathbb{E}[X] \cdot \mathbb{E}[Y]$ , and another example where  $\mathbb{E}[X \cdot Y] = \mathbb{E}[X] \cdot \mathbb{E}[Y]$ .  
 (d) Prove that if  $X$  and  $Y$  are independent, then they are also uncorrelated, i.e.,

$$\text{Cov}(X, Y) := \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] = 0. \quad (2.2)$$

- (e) Find an example where  $X$  and  $Y$  are uncorrelated but dependent.  
 (f) Assume that  $X$  and  $Y$  are uncorrelated and let  $\sigma_X^2$  and  $\sigma_Y^2$  be the variances of  $X$  and  $Y$ , respectively. Find the variance of  $aX + bY$  and express it in terms of  $\sigma_X^2, \sigma_Y^2, a, b$ .  
**Hint:** First show that  $\text{Cov}(X, Y) = \mathbb{E}[X \cdot Y] - \mathbb{E}[X] \cdot \mathbb{E}[Y]$ .

**Problem 2.2** (Gaussian Random Variables). A random variable  $X$  with probability density function

$$p_X(x) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-m)^2}{2\sigma^2}} \quad (2.3)$$

is called a *Gaussian* random variable.

- (a) Explicitly calculate the mean  $\mathbb{E}[X]$ , the second moment  $\mathbb{E}[X^2]$ , and the variance  $\text{Var}[X]$  of the random variable  $X$ .  
 (b) Let us now consider events of the following kind:

$$\mathbb{P}(X < \alpha). \quad (2.4)$$

Unfortunately for Gaussian random variables this cannot be calculated in closed form. Instead, we will rewrite it in terms of the standard Q-function:

$$Q(x) = \int_x^\infty \frac{1}{\sqrt{2\pi}} e^{-\frac{u^2}{2}} du \quad (2.5)$$

Express  $\mathbb{P}(X < \alpha)$  in terms of the Q-function and the parameters  $m$  and  $\sigma^2$  of the Gaussian pdf.

Like we said, the Q-function cannot be calculated in closed form. Therefore, it is important to have *bounds* on the Q-function. In the next 3 subproblems, you derive the most important of these bounds, learning some very general and powerful tools along the way:

- (c) Derive the Markov inequality, which says that for any non-negative random variable  $X$  and positive  $a$ , we have

$$\mathbb{P}(X \geq a) \leq \frac{\mathbb{E}[X]}{a}. \quad (2.6)$$

- (d) Use the Markov inequality to derive the Chernoff bound: the probability that a real random variable  $Z$  exceeds  $b$  is given by

$$\mathbb{P}(Z \geq b) \leq \mathbb{E}[e^{s(Z-b)}], \quad s \geq 0. \quad (2.7)$$

- (e) Use the Chernoff bound to show that

$$Q(x) \leq e^{-\frac{x^2}{2}} \quad \text{for } x \geq 0. \quad (2.8)$$

**Problem 2.3** (Moment Generating Function). In the class we had considered the logarithmic moment generating function

$$\phi(s) := \ln E[\exp(sX)] = \ln \sum_x p(x) \exp(sx)$$

of a real-valued random variable  $X$  taking values on a finite set, and showed that  $\phi'(s) = E[X_s]$  where  $X_s$  is a random variable taking the same values as  $X$  but with probabilities  $p_s(x) := p(x) \exp(sx) \exp(-\phi(s))$ .

(a) Show that

$$\phi''(s) = \text{Var}(X_s) := E[X_s^2] - E[X_s]^2$$

and conclude that  $\phi''(s) \geq 0$  and the inequality is strict except when  $X$  is deterministic.

(b) Let  $x_{\min} := \min\{x : p(x) > 0\}$  and  $x_{\max} := \max\{x : p(x) > 0\}$  be the smallest and largest values  $X$  takes. Show that

$$\lim_{s \rightarrow -\infty} \phi'(s) = x_{\min}, \quad \text{and} \quad \lim_{s \rightarrow \infty} \phi'(s) = x_{\max}.$$

**Problem 2.4** (Divergence and  $L_1$ ). Suppose  $p$  and  $q$  are two probability mass functions on a finite set  $\mathcal{U}$ . (I.e., for all  $u \in \mathcal{U}$ ,  $p(u) \geq 0$  and  $\sum_{u \in \mathcal{U}} p(u) = 1$ ; similarly for  $q$ .)

(a) Show that the  $L_1$  distance  $\|p - q\|_1 := \sum_{u \in \mathcal{U}} |p(u) - q(u)|$  between  $p$  and  $q$  satisfies

$$\|p - q\|_1 = 2 \max_{\mathcal{S} \subset \mathcal{U}} p(\mathcal{S}) - q(\mathcal{S})$$

with  $p(\mathcal{S}) = \sum_{u \in \mathcal{S}} p(u)$  (and similarly for  $q$ ), and the maximum is taken over all subsets  $\mathcal{S}$  of  $\mathcal{U}$ .

For  $\alpha$  and  $\beta$  in  $[0, 1]$ , define the function  $d_2(\alpha\|\beta) := \alpha \log \frac{\alpha}{\beta} + (1 - \alpha) \log \frac{1 - \alpha}{1 - \beta}$ . Note that  $d_2(\alpha\|\beta)$  is the divergence of the distribution  $(\alpha, 1 - \alpha)$  from the distribution  $(\beta, 1 - \beta)$ .

(b) Show that the first and second derivatives of  $d_2$  with respect to its first argument  $\alpha$  satisfy  $d_2'(\beta\|\beta) = 0$  and  $d_2''(\alpha\|\beta) = \frac{\log e}{\alpha(1 - \alpha)} \geq 4 \log e$ .

(c) By Taylor's theorem conclude that

$$d_2(\alpha\|\beta) \geq 2(\log e)(\alpha - \beta)^2.$$

(d) Show that for any  $\mathcal{S} \subset \mathcal{U}$

$$D(p\|q) \geq d_2(p(\mathcal{S})\|q(\mathcal{S}))$$

[Hint: use the data processing theorem for divergence.]

(e) Combine (a), (c) and (d) to conclude that

$$D(p\|q) \geq \frac{\log e}{2} \|p - q\|_1^2.$$

- (f) Show, by example, that  $D(p\|q)$  can be  $+\infty$  even when  $\|p - q\|_1$  is arbitrarily small. [Hint: considering  $\mathcal{U} = \{0, 1\}$  is sufficient.] Consequently, there is no generally valid inequality that upper bounds  $D(p\|q)$  in terms of  $\|p - q\|_1$ .

**Problem 2.5** (Other Divergences). Suppose  $f$  is a convex function defined on  $(0, \infty)$  with  $f(1) = 0$ . Define the  $f$ -divergence of a distribution  $p$  from a distribution  $q$  as

$$D_f(p\|q) := \sum_u q(u) f(p(u)/q(u)).$$

In the sum above we take  $f(0) := \lim_{t \rightarrow 0} f(t)$ ,  $0f(0/0) := 0$ , and  $0f(a/0) := \lim_{t \rightarrow 0} tf(a/t) = a \lim_{t \rightarrow 0} tf(1/t)$ .

- (a) Show that for any non-negative  $a_1, a_2, b_1, b_2$  and with  $A = a_1 + a_2$ ,  $B = b_1 + b_2$ ,

$$b_1 f(a_1/b_1) + b_2 f(a_2/b_2) \geq B f(A/B);$$

and that in general, for any non-negative  $a_1, \dots, a_k, b_1, \dots, b_k$ , and  $A = \sum_i a_i$ ,  $B = \sum_i b_i$ , we have

$$\sum_i b_i f(a_i/b_i) \geq B f(A/B).$$

[Hint: since  $f$  is convex, for any  $\lambda \in [0, 1]$  and any  $x_1, x_2 > 0$   $\lambda f(x_1) + (1 - \lambda)f(x_2) \geq f(\lambda x_1 + (1 - \lambda)x_2)$ ; consider  $\lambda = b_1/B$ .]

- (b) Show that  $D_f(p\|q) \geq 0$ .
- (c) Show that  $D_f$  satisfies the data processing inequality: for any transition probability kernel  $W(v|u)$  from  $\mathcal{U}$  to  $\mathcal{V}$ , and any two distributions  $p$  and  $q$  on  $\mathcal{U}$

$$D_f(p\|q) \geq D_f(\tilde{p}\|\tilde{q})$$

where  $\tilde{p}$  and  $\tilde{q}$  are probability distributions on  $\mathcal{V}$  defined via  $\tilde{p}(v) := \sum_u W(v|u)p(u)$ , and  $\tilde{q}(v) := \sum_u W(v|u)q(u)$ ,

- (d) Show that each of the following are  $f$ -divergences.

- i.  $D(p\|q) := \sum_u p(u) \log(p(u)/q(u))$ . [Warning:  $\log$  is not the right choice for  $f$ .]
- ii.  $R(p\|q) := D(q\|p)$ .
- iii.  $1 - \sum_u \sqrt{p(u)q(u)}$
- iv.  $\|p - q\|_1$ .
- v.  $\sum_u (p(u) - q(u))^2 / q(u)$

**Problem 2.6** (Entropy and pairwise independence). Suppose  $X, Y, Z$  are pairwise independent fair flips, i.e.,  $I(X; Y) = I(Y; Z) = I(Z; X) = 0$ .

- (a) What is  $H(X, Y)$ ?

- (b) Give a lower bound to the value of  $H(X, Y, Z)$ .
- (c) Give an example that achieves this bound.

**Problem 2.7** (Generating fair coin flips from biased coins). Suppose  $X_1, X_2, \dots$  are the outcomes of independent flips of a biased coin. Let  $\Pr(X_i = 1) = p$ ,  $\Pr(X_i = 0) = 1 - p$ , with  $p$  unknown. By processing this sequence we would like to obtain a sequence  $Z_1, Z_2, \dots$  of *fair* coin flips.

Consider the following method: We process the  $X$  sequence in successive pairs,  $(X_1X_2)$ ,  $(X_3X_4)$ ,  $(X_5X_6)$ , mapping (01) to 0, (10) to 1, and the other outcomes (00) and (11) to the empty string. After processing  $X_1, X_2$ , we will obtain either nothing, or a bit  $Z_1$ .

- (a) Show that, if a bit is obtained, it is fair, i.e.,  $\Pr(Z_1 = 0) = \Pr(Z_1 = 1) = 1/2$ .

In general we can process the  $X$  sequence in successive  $n$ -tuples via a function  $f : \{0, 1\}^n \rightarrow \{0, 1\}^*$  where  $\{0, 1\}^*$  denote the set of all finite length binary sequences (including the empty string  $\lambda$ ). [The case in (a) is the function  $f(00) = f(11) = \lambda$ ,  $f(01) = 0$ ,  $f(10) = 1$ . The function  $f$  is chosen such that  $(Z_1, \dots, Z_K) = f(X_1, \dots, X_n)$  are i.i.d., and fair (here  $K$  may depend on  $(X_1, \dots, X_n)$ ).

- (b) With  $h_2(p) = -p \log p - (1 - p) \log(1 - p)$ , prove the following chain of (in)equalities.

$$\begin{aligned} nh_2(p) &= H(X_1, \dots, X_n) \\ &\geq H(Z_1, \dots, Z_K, K) \\ &= H(K) + H(Z_1, \dots, Z_K | K) \\ &= H(K) + E[K] \\ &\geq E[K]. \end{aligned}$$

Consequently, on the average no more than  $nh_2(p)$  fair bits can be obtained from  $(X_1, \dots, X_n)$ .

- (c) Find a good  $f$  for  $n = 4$ .

**Problem 2.8** (Extremal characterization for Rényi entropy). Given  $s \geq 0$ , and a random variable  $U$  taking values in  $\mathcal{U}$ , with probabilities  $p(u)$ , consider the distribution  $p_s(u) = p(u)^s / Z(s)$  with  $Z(s) = \sum_u p(u)^s$ .

- (a) Show that for any distribution  $q$  on  $\mathcal{U}$ ,

$$(1 - s)H(q) - sD(q||p) = -D(q||p_s) + \log Z(s).$$

- (b) Given  $s$  and  $p$ , conclude that the left hand side above is maximized by the choice by  $q = p_s$  with the value  $\log Z(s)$ ,

The quantity

$$H_s(p) := \frac{1}{1 - s} \log Z(s) = \frac{1}{1 - s} \log \sum_u p(u)^s$$

is known as the *Rényi entropy of order  $s$  of the random variable  $U$* . When convenient, we will also write  $H_s(U)$  instead of  $H_s(p)$ .

- (c) Show that if  $U$  and  $V$  are independent random variables

$$H_s(UV) := H_s(U) + H_s(V).$$

[Here  $UV$  denotes the pair formed by the two random variables — not their product. E.g., if  $\mathcal{U} = \{0, 1\}$  and  $\mathcal{V} = \{a, b\}$ ,  $UV$  takes values in  $\{0a, 0b, 1a, 1b\}$ .]

**Problem 2.9** (Guessing and Rényi entropy). Suppose  $X$  is a random variable taking values  $K$  values  $\{a_1, \dots, a_K\}$  with  $p_i = \Pr\{X = a_i\}$ . We wish to guess  $X$  by asking a sequence of binary questions of the type ‘Is  $X = a_i$ ?’ until we are answered ‘yes’. (Think of guessing a password).

A *guessing strategy* is an ordering of the  $K$  possible values of  $X$ ; we first ask if  $X$  is the first value; then if it is the second value, etc. Thus the strategy is described by a function  $G(x) \in \{1, \dots, K\}$  that gives the position (first, second, ...  $K$ th) of  $x$  in the ordering. I.e., when  $X = x$ , we ask  $G(x)$  questions to guess the value of  $X$ . Call  $G$  the guessing function of the strategy.

For the rest of the problem suppose  $p_1 \geq p_2 \geq \dots \geq p_K$ .

- (a) Show that for any guessing function  $G$ , the probability of asking fewer than  $i$  questions satisfies

$$\Pr(G(X) \leq i) \leq \sum_{j=1}^i p_j$$

and equality holds for the guessing function  $G^*$  with  $G^*(a_i) = i$ ,  $i = 1, \dots, K$ ; this is the strategy that first guesses the most probable value  $a_1$ , then the next most probable value  $a_2$ , etc.

- (b) Show that for any increasing function  $f : \{1, \dots, K\} \rightarrow \mathbb{R}$ ,  $E[f(G(X))]$  is minimized by choosing  $G = G^*$ . [Hint:  $E[f(G(X))] = \sum_{i=1}^K f(i) \Pr(G = i)$ . Write  $\Pr(G = i) = \Pr(G \leq i) - \Pr(G \leq i - 1)$ , to write the expectation in terms of  $\sum_i [f(i) - f(i + 1)] \Pr(G \leq i)$ , and use (a).]
- (c) For any  $i$  and  $s \geq 0$  prove the inequalities

$$i \leq \sum_{j=1}^i (p_j/p_i)^s \leq \sum_j (p_j/p_i)^s$$

- (d) For any  $\rho \geq 0$ , show that

$$E[G^*(X)^\rho] \leq \left( \sum_i p_i^{1-s\rho} \right) \left( \sum_j p_j^s \right)^\rho.$$

for any  $s \geq 0$ . [Hint: write  $E[G^*(X)^\rho] = \sum_i p_i i^\rho$ , and use (c) to upper bound  $i^\rho$ ]

- (e) By choosing  $s$  carefully, show that

$$E[G^*(X)^\rho] \leq \left( \sum_i p_i^{1/(1+\rho)} \right)^{1+\rho} = \exp[\rho H_{1/(1+\rho)}(X)].$$



- (f) Suppose  $U_1, \dots, U_n$  are i.i.d., each with distribution  $p$ , and  $X = (U_1, \dots, U_n)$ . (I.e., we are trying to guess a password that is made of  $n$  independently chosen letters.) Show that

$$\frac{1}{n\rho} \log E[G^*(U_1, \dots, U_n)^\rho] \leq H_{1/(1+\rho)}(U_1)$$

[Hint: first observe that  $H_\alpha(X) = nH_\alpha(U_1)$ . In other words, the  $\rho$ -th moment of the number of guesses grows exponentially in  $n$  with a rate upper bounded by in terms of the Rényi entropy of the letters.

It is possible a lower bound to  $E[G^*(U_1, \dots, U_n)^\rho]$  that establishes that the exponential upper bound we found here is asymptotically tight.

**Problem 2.10** (Gaussian variance estimation). Consider estimating the mean  $\mu$  and variance  $\sigma^2$  from  $n$  independent samples  $(X_1, \dots, X_n)$  of a Gaussian with this mean and variance.

- (a) Show that  $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$  is an unbiased estimator of  $\mu$ .  
 (b) Show that

$$S_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$$

is a biased estimator of  $\sigma^2$  whereas

$$S_{n-1}^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

is an unbiased estimator of  $\sigma^2$ .

- (c) Show that  $S_n^2$  has a lower mean squared error than  $S_{n-1}^2$ . Thus it is possible that a biased estimator may be better than an unbiased one.



## Chapter 3

# Compression and Quantization

### 3.1 Data compression

**Notation.** Given a set  $A$  we denote by  $A^*$  the set of all finite sequences  $\{(a_1, \dots, a_n) : n \geq 0, a_i \in A\}$  (including the null sequence  $\lambda$  of length 0). In particular  $\{0, 1\}^* = \{\lambda, 0, 1, 00, 01, 10, 11, 000, \dots\}$ .

Consider the problem of assigning binary sequences (also called binary strings) to elements of a finite set  $\mathcal{U}$ . Such an assignment  $c : \mathcal{U} \rightarrow \{0, 1\}^*$  is called a *binary code* for the set  $\mathcal{U}$ . The binary string  $c(u)$  is called the *codeword* for  $u$ . The collection  $\{c(u) : u \in \mathcal{U}\}$  is thus the set of codewords.

**Definition 3.1.** A code  $c$  is called *injective* if for all  $u \neq v$  we have  $c(u) \neq c(v)$ .

**Definition 3.2.** A code  $c$  is called *prefix-free* if  $c(u)$  is not a prefix of  $c(v)$  for all  $u \neq v$ . In particular, if  $c$  is prefix-free then  $c$  is injective. (To be clear: a string  $a_1 \dots a_m$  is a prefix of a string  $b_1 \dots b_n$  if  $m \leq n$  and  $a_i = b_i$  for  $i = 1, \dots, m$ . Thus, the null string is a prefix of any string, and each string is a prefix of itself.)

**Lemma 3.1.** Suppose  $c : \mathcal{U} \rightarrow \{0, 1\}^*$  is injective. Then,  $\sum_u 2^{-\text{length}(c(u))} \leq \log_2(1 + |\mathcal{U}|)$ .

*Proof.* Without loss of generality, we can assume that whenever  $k = \text{length}(c(u))$  for some  $u$ , then for every binary string  $b$  of length  $i < k$  there is a  $v$  with  $b = c(v)$ . (Otherwise, there is a  $b$  with  $\text{length}(b) < k$  which is not a codeword, and replacing  $c(u)$  with  $b$  will preserve the injectiveness of  $c$  and increase the left hand side of the inequality.)

For such a code  $c$ , with  $k$  denoting the length of the longest codeword, the set of codewords is the union of  $\bigcup_{i=0}^{k-1} \{0, 1\}^i$  with a non-empty subset of  $\{0, 1\}^k$ . With  $1 \leq r \leq 2^k$  denoting the cardinality of this last subset, we have  $|\mathcal{U}| = 2^k - 1 + r$  and  $\sum_u 2^{-\text{length}(c(u))} = k + r2^{-k}$ . As  $\log_2(1 + |\mathcal{U}|) = k + \log_2(1 + r2^{-k})$  and  $0 < r2^{-k} \leq 1$ , all we need to show is  $x \leq \log_2(1 + x)$  for  $0 < x \leq 1$ . As equality obtains for  $x = 0$  and  $x = 1$ , the inequality follows from the concavity of  $\log$ .  $\square$

**Lemma 3.2.** Suppose  $c : \mathcal{U} \rightarrow \{0, 1\}^*$  is prefix-free. Then,  $\sum_u 2^{-\text{length}(c(u))} \leq 1$ . Conversely, if  $\ell : \mathcal{U} \rightarrow \{0, 1, 2, \dots\}$  with  $\sum_u 2^{-\ell(u)} \leq 1$ , then there exists a prefix-free code  $c : \mathcal{U} \rightarrow \{0, 1\}^*$  with  $\text{length}(c(u)) = \ell(u)$ .

*Proof.* Given a binary sequence  $a = a_1 \dots a_m$ , let  $p(a) = \sum_{i=1}^m a_i 2^{-i}$  denote the rational number whose binary expansion is  $0.a_1 \dots a_m$ . With this notation, a binary sequence  $a = a_1 \dots a_m$  is a prefix of a binary sequence  $b = b_1 \dots b_n$  if and only if  $p(b)$  lies in the interval  $I(a) = [p(a), p(a) + 2^{-m})$ .

For the first claim, observe that  $c$  being prefix-free thus implies that the intervals  $I(c(u))$  are disjoint. As  $I(c(u))$  is of size  $2^{-\text{length}(c(u))}$  and all of the intervals are included in  $[0, 1)$ , the inequality follows.

For the second claim, order the elements of  $\mathcal{U}$  as  $u_1, \dots, u_K$  such that  $\ell_1 := \ell(u_1) \leq \dots \leq \ell_K := \ell(u_K)$ . Let  $p_k = \sum_{i < k} 2^{-\ell_i}$  and set  $I_k = [p_k, p_k + 2^{-\ell_k})$ . Observe that the intervals  $I_1, \dots, I_K$  are disjoint, and  $I_k \subset [0, 1)$ . Furthermore, for each  $k$ ,  $2^{\ell_k} p_k$  is an integer, thus  $p_k$  can be expressed in binary as  $0.b^{(k)}$  with  $b^{(k)}$  a binary string of length  $\ell_k$ . The code  $c(u_k) = b^{(k)}$  now has the required properties — it being prefix free a consequence of the disjointness of the collection intervals  $I_k$ .  $\square$

**Lemma 3.3.** *Suppose  $P \in \Pi(\mathcal{U})$  is a probability distribution on  $\mathcal{U}$  and  $U$  is random variable with distribution  $P$ . Then, with  $H(U) = -\sum_u P(u) \log_2 P(u)$  denoting the entropy of  $U$ ,*

- (i) *for any prefix-free  $c : \mathcal{U} \rightarrow \{0, 1\}^*$ ,  $E[\text{length}(c(U))] \geq H(U)$ ;*
- (ii) *there exists a prefix-free  $c : \mathcal{U} \rightarrow \{0, 1\}^*$  with  $E[\text{length}(c(U))] \leq H(U) + 1$ ;*
- (iii) *for any injective  $c : \mathcal{U} \rightarrow \{0, 1\}^*$ ,  $E[\text{length}(c(U))] \geq H(U) - \log_2 \log_2(1 + |\mathcal{U}|)$ ,*
- (iv) *there exists an injective  $c : \mathcal{U} \rightarrow \{0, 1\}^*$  with  $E[\text{length}(c(U))] \leq H(U)$ .*

*Proof.* For (i) and (iii) let  $Q(u) = 2^{-\text{length}(c(u))}$  and observe that

$$H(U) - E[\text{length}(c(U))] = \sum_u P(u) \log_2 \frac{Q(u)}{P(u)} \leq \log_2 \sum_u Q(u),$$

where the inequality is because  $\log$  is concave. When  $c$  is prefix-free  $\sum_u Q(u) \leq 1$  by Lemma 3.2, and when  $c$  is injective  $\sum_u Q(u) \leq \log_2(1 + |\mathcal{U}|)$  by Lemma 3.1. The inequalities (i) and (iii) thus follow.

For (ii) set  $\ell(u) = \lceil -\log_2 P(u) \rceil$ . As  $2^{-\ell(u)} \leq P(u)$ , we see that  $\sum_u 2^{-\ell(u)} \leq 1$  and by Lemma 3.2 there exists a prefix-free code  $c$  with  $\text{length}(c(u)) = \ell(u)$ . As  $\ell(u) < -\log_2 P(u) + 1$ , (ii) follows.

For (iv) order the elements of  $\mathcal{U}$  as  $u_1, \dots, u_K$  with  $P(u_1) \geq \dots \geq P(u_K)$ . Let  $c(u_k) = b_k$  where  $b_k$  is the  $k$ th element of the sequence  $\lambda, 0, 1, 00, 01, 10, 11, 000, 001, \dots$ , (e.g.,  $b_1 = \lambda$ ,  $b_2 = 0$ ,  $b_3 = 1$ ,  $b_4 = 00$ ,  $\dots$ ,  $b_9 = 001$ ,  $\dots$ ). Observe that  $\text{length}(b_k) = \lfloor \log_2 k \rfloor \leq \log_2 k$ . Also note that  $1 \geq \sum_{i=1}^k P(u_i) \geq kP(u_k)$ , and thus  $\log_2 k \leq -\log_2 P(u_k)$ . Consequently, for this  $c$ ,  $E[\text{length}(c(U))] \leq -\sum_k P(u_k) \log_2 P(u_k) = H(U)$ .  $\square$

**Corollary 3.4.** *Suppose  $U_1, U_2, \dots$  is a stochastic process. Then for any sequence  $c_n : \mathcal{U}^n \rightarrow \{0, 1\}^*$  of injective codes*

$$\liminf_n \frac{1}{n} E[\text{length}(c_n(U^n))] \geq \liminf_n \frac{1}{n} H(U^n),$$

and there exists a sequence  $c_n$  of prefix-free codes for which

$$\limsup_n \frac{1}{n} E[\text{length}(c_n(U^n))] \leq \limsup_n \frac{1}{n} H(U^n).$$

In particular, if  $r = \lim_n \frac{1}{n} H(U^n)$  exists, all faithful representations of the process  $U_1, U_2, \dots$  with bits will asymptotically require at least  $r$  bits per letter, and there is a representation that asymptotically requires exactly as much.

*Proof.* The first inequality follows from noting that

$$E[\text{length}(c_n(U^n))] \geq H(U^n) - \log_2 \log_2(1 + |\mathcal{U}|^n)$$

and observing that  $\lim_n \frac{1}{n} \log_2 \log_2(1 + |\mathcal{U}|^n) = 0$ . The second inequality follows from noting that there exist prefix-free  $c_n$  with

$$E[\text{length}(c_n(U^n))] \leq H(U^n) + 1$$

and that  $\lim_n 1/n = 0$ . □

**Remark.** Lemma 3.2 gives evidence of a strong connection between prefix-free codes and probability distributions. On the one hand, given a prefix-free code  $c$ , one can construct a probability distribution  $Q$  that assigns the letter  $u$  the probability  $Q(u) = 2^{-\text{length}(c(u))}$ . By the lemma,  $\sum_u Q(u) \leq 1$ ; if equality holds  $Q$  is indeed a probability distribution, otherwise, we can assign  $1 - \sum_u Q(u)$  as the probability  $Q(u_0)$  of a fictitious symbol  $u_0 \notin \mathcal{U}$ . If  $U$  is a random variable with distribution  $P$ , we then have (by assigning  $P(u_0) = 0$  if necessary),

$$E[\text{length}(c(U))] - H(U) = \sum_u P(u)[\text{length}(c(U)) + \log P(u)] = D(P\|Q).$$

On the other hand, given a distribution  $Q \in \Pi(\mathcal{U})$ , by Lemma 3.2 we can construct a prefix-free code  $c : \mathcal{U} \rightarrow \{0, 1\}^*$  with  $\text{length}(c(u)) = \lceil -\log_2 Q(u) \rceil$ . As  $-\log_2 Q(u) \leq \text{length}(c(u)) < -\log_2 Q(u) + 1$ , we see that

$$E[\text{length}(c(U))] - H(U) = \sum_u P(u)[\text{length}(c(u)) + \log_2 P(u)]$$

is bounded from below by  $D(P\|Q)$ , and from above  $D(P\|Q) + 1$ .

These observations give the divergence  $D(P\|Q)$  an interpretation as the expected number of “excess” bits (beyond the minimum possible  $H(U)$ ) a code based on  $Q$  requires when describing a random variable with distribution  $P$ .

Consequently, if we are given  $S \subset \Pi$  and told that the distribution  $P$  of a random variable  $U$  belongs to  $S$ , a reasonable strategy to design a code  $c$  is to look for a distribution  $Q \in \Pi$  such that

$$\sup_{P \in S} D(P\|Q)$$

is small (e.g., by finding the  $Q$  that minimizes this quantity) and construct a code  $c$  based on  $Q$  as above.

**Example 3.1.** To illustrate the remark above, suppose we are told that  $U_1, U_2, \dots$  are binary and i.i.d. random variables. The distribution of  $U^n$  can be parametrized by  $\theta = \Pr(U_1 = 1)$ , and is given by

$$\Pr(U^n = u^n) = P_\theta^n(u^n) = (1 - \theta)^{n_0(u^n)} \theta^{n_1(u^n)}$$

where  $n_0(u^n)$  and  $n_1(u^n)$  are the number of zeros and ones in the sequence  $u_1 \dots u_n$ . With this notation,  $S_n = \{P_\theta^n : 0 \leq \theta \leq 1\}$  is the class of distributions that we are told the distribution of  $U^n$  belongs to.

Consider now a sequence of conditional distributions

$$Q_{U_{k+1}|U^k}(u|u^k) = \frac{n_u(u^k) + 1}{k + 2}$$

where  $n_u(u^k)$  is as above, denoting the number of  $u$ 's in  $u_1 \dots u_k$ . Note that  $Q_{U_1}(0) = Q_{U_1}(1) = 1/2$ . Define

$$Q_n(u^n) = \prod_{i=1}^n Q_{U_i|U^{i-1}}(u_i|u^{i-1}).$$

One can prove by induction on  $n$ , that for any  $n \geq 1$  and any  $u^n \in \{0, 1\}^n$ ,

$$Q_n(u^n) \geq \frac{1}{n+1} \left( \frac{n_0(u^n)}{n} \right)^{n_0(u^n)} \left( \frac{n_1(u^n)}{n} \right)^{n_1(u^n)}.$$

If  $U_1, \dots, U_n$  are i.i.d. with common distribution  $P_\theta$ ,

$$\begin{aligned} D(P_\theta^n \| Q_n) &= E \left[ \log \frac{P_\theta^n(U^n)}{Q_n(U^n)} \right] \\ &\leq \log(n+1) + E \left[ \log \frac{P_\theta^n(U^n)}{(n_0(U^n)/n)^{n_0(U^n)} (n_1(U^n)/n)^{n_1(U^n)}} \right] \\ &= \log(n+1) + E \left[ n_0(U^n) \log \frac{n(1-\theta)}{n_0(U^n)} + n_1(U^n) \log \frac{n\theta}{n_1(U^n)} \right] \\ &\leq \log(n+1) + n(1-\theta) \log \frac{n(1-\theta)}{n(1-\theta)} + n\theta \log \frac{n\theta}{n\theta} = \log(n+1), \end{aligned}$$

where the inequality in the last line is because  $x \mapsto x \log[1/x]$  is concave and  $E[n_0(U^n)] = n(1-\theta)$ , and  $E[n_1(U^n)] = n\theta$ .

Consequently, we see that  $\sup_{P^n \in S_n} D(P^n \| Q_n) \leq \log(n+1)$ . If  $Q_n$  were used to construct a prefix-free code  $c_n : \{0, 1\}^n \rightarrow \{0, 1\}^*$ , by the remark above,  $c_n$  will satisfy

$$\frac{1}{n} E[\text{length}(c_n(U^n))] - H(P) \leq \frac{1}{n} [\log(n+1) + 1]$$

whenever  $U^n$  is i.i.d. with distribution  $P$ . As the right hand side vanishes as  $n$  gets large, it would be appropriate to call the sequence of codes  $c_n$  “asymptotically universal for the class of binary i.i.d. data”. In the exercises we will see another choice of  $Q_n$  which improves the upper bound on  $D(P^n \| Q_n)$  to  $\frac{1}{2} \log n$ .

Note that, had we chosen  $Q_n$  to be a member of  $S_n$ , say  $Q_n = P_{\theta_0}^n$  for some  $\theta_0$ , then  $D(P_{\theta}^n \| Q_n)$  would have grown linearly in  $n$  for any  $\theta \neq \theta_0$ . Thus, even if we know that the true distribution  $P$  is in  $S$ , choosing  $Q$  outside of  $S$  (as we have done above) may lead to a better code construction.

**Remark.** The example above also illustrates a connection between compression and prediction. (One can also use the term ‘learning’ instead of prediction.) Suppose we have a family  $S_n$  of distributions on  $\mathcal{U}^n$ , and we are given a prefix-free code  $c_n : \mathcal{U}^n \rightarrow \{0,1\}^*$  performs well, in the sense that

$$\sup_{P \in S_n} \frac{1}{n} E_P[\text{length}(c_n(U^n))] - \frac{1}{n} H(U^n)$$

is small. Construct the distribution  $Q$  associated with the code  $c$ , i.e.,  $Q(u^n) = 2^{-\text{length}(c_n(u^n))}$  and factorize it as  $Q(u^n) = \prod_{i=1}^n Q(u_i | u^{i-1})$ . As the code  $c$  performs well,  $\frac{1}{n} D(P \| Q)$  is small for all  $P \in S_n$ . But

$$\begin{aligned} \frac{1}{n} D(P \| Q) &= \frac{1}{n} \sum_{u^n} P(u^n) \log \frac{P(u^n)}{Q(u^n)} \\ &= \frac{1}{n} \sum_{i=1}^n \sum_{u^n} P(u^n) \log \frac{P(u_i | u^{i-1})}{Q(u_i | u^{i-1})} \\ &= \frac{1}{n} \sum_{i=1}^n \sum_{u^i} P(u^i) \log \frac{P(u_i | u^{i-1})}{Q(u_i | u^{i-1})} \\ &= \frac{1}{n} \sum_{i=1}^n \sum_{u^{i-1}} P(u^{i-1}) \sum_{u_i} P(u_i | u^{i-1}) \log \frac{P(u_i | u^{i-1})}{Q(u_i | u^{i-1})} \\ &= \frac{1}{n} \sum_{i=1}^n \sum_{u^{i-1}} P(u^{i-1}) D(P(\cdot | u^{i-1}) \| Q(\cdot | u^{i-1})), \end{aligned}$$

so we conclude that for a large fraction of  $i$ 's in  $1, \dots, n$ , and for a set of  $u^{i-1}$ 's with large  $P$  probability, the quantity  $D(P(\cdot | u^{i-1}) \| Q(\cdot | u^{i-1}))$  is small.<sup>1</sup> Which is to say, no matter what  $P$  from  $S_n$  is the true distribution of the data, if after observing  $u^{i-1}$  we predicted the distribution of the next symbol  $u_i$  to be  $Q(\cdot | u^{i-1})$ , our prediction will be close to the true distribution  $P(\cdot | u^{i-1})$  for most  $i$ 's and for a high probability set of  $u^{i-1}$ 's.

## 3.2 Universal data compression with the Lempel-Ziv algorithm

In the example in the previous section we saw a compression method that was universal over the class of binary i.i.d. processes. We will now see a much more powerful method that

<sup>1</sup>To be concrete, if  $\frac{1}{n} D(P \| Q)$  is less than  $\epsilon$ , then, except for a  $\epsilon^{1/3}$  fraction of the  $i$ 's, we have  $\sum_{u^{i-1}} P(u^{i-1}) D(P(\cdot | u^{i-1}) \| Q(\cdot | u^{i-1})) < \epsilon^{2/3}$ , and except for a set of  $P$  probability  $\epsilon^{1/3}$  of  $u^{i-1}$ 's, we have  $D(P(\cdot | u^{i-1}) \| Q(\cdot | u^{i-1})) < \epsilon^{1/3}$ .

is universal over all stationary processes. The method was invented by Ziv and Lempel in 1977, the version we present here is a variant due to Welch from 1984.

Given an alphabet  $\mathcal{U} = \{a_1, \dots, a_K\}$ , the method encodes an infinite sequence  $u_1 u_2 \dots$  from this alphabet to binary as follows:

1. Set a dictionary  $\mathcal{D} = \mathcal{U}$ . Denote the dictionary entries as  $d(0) = a_1, \dots, d(s-1) = a_K$ , with  $s = K$  being the size of the dictionary. Set  $i = 0$  (the number of input letters read so far).
2. Find the largest  $l$  such that  $w = u_{i+1} \dots u_{i+l}$  is in  $\mathcal{D}$ .
3. With  $0 \leq j < s$  denoting the index of  $w$  in  $\mathcal{D}$ , output the  $\lceil \log_2 s \rceil$  bit binary representation of  $j$ .
4. Add the word  $wu_{i+l+1}$  to  $\mathcal{D}$ , i.e., set  $d(s) = wu_{i+l+1}$ , and increment  $s$  by 1. Increment  $i$  by  $l$ . Goto step 2.

For example, with  $\mathcal{U} = \{a, b\}$ , the input string `abbbbbaaab...` will lead to the execution steps

$\mathcal{D}$ at 2	$w$	output at 3	added-word at 4
a b	a	0	ab
a b ab	b	01	bb
a b ab bb	bb	11	bbb
a b ab bb bbb	b	001	ba
a b ab bb bbb ba	a	000	aa
a b ab bb bbb ba aa	aa	110	aab

The first question we need to answer is if we can recover the input sequence  $u_1 u_2 \dots$  from the output of the algorithm. The question is answered in the affirmative in Lemma 3.5 below.

Note that the algorithm parses the sequence  $u_1 u_2 \dots$  into a sequence of words  $w_1, w_2, \dots$  found at step 2 of the algorithm. So, recovery of  $u_1 u_2 \dots$  is equivalent to the recovery of these words. Let  $j_1, j_2, \dots$  the dictionary indices that appear at step 3, and  $d_1, d_2, \dots$  the words added to the dictionary in step 4. As the dictionary size  $s$  increases by 1 each time a dictionary word is parsed, the bitstream that is output by the algorithm can be parsed into the indices  $j_1, j_2, \dots$ .

**Lemma 3.5.** *From  $j_1, \dots, j_i$  we can determine  $w_1, \dots, w_i$ . In other words, we can recover the input  $u_1 u_2 \dots$  from the output of the algorithm.*

*Proof.* From  $j_1$ , we can determine  $w_1$ , and so the claim is true for  $i = 1$ . We proceed by induction. Suppose now we observe  $j_1, \dots, j_{i+1}$ . By the induction hypothesis,  $j_1, \dots, j_i$  determines  $w_1, \dots, w_i$ . Since  $d_k$  is the concatenation of  $w_k$  with the first letter of  $w_{k+1}$ , we know  $d_1, \dots, d_{i-1}$ , and, except for its last letter,  $d_i$ . We need to show that we can reconstruct  $w_{i+1}$  from the additional information obtained by  $j_{i+1}$ . Note that  $j_{i+1}$  refers to a word in the dictionary formed by augmenting  $\mathcal{U}$  with the words  $d_1, \dots, d_i$ . If  $j_{i+1}$  refers



to any word other than  $d_i$ , we already determine  $w_{i+1}$ . Otherwise  $j_{i+1}$  refers to  $d_i$ . But in this case the last letter of  $d_i$  equals the first letter of  $d_i$  which is already known, and we can again determine  $w_{i+1}$ .  $\square$

Next, we will obtain an upper bound on the number of bits per letter the algorithm uses to describe a sequence  $u_1 u_2 \dots$ . The answer will be given in Theorem 3.8 below.

**Lemma 3.6.** *A word  $w$  can appear in the sequence  $w_1, w_2, \dots$  at most  $|\mathcal{U}|$  times.*

*Proof.* As the algorithm always looks for the longest word  $w$  in the dictionary that matches the start of the as-yet-unprocessed segment of the input,  $wu_{i+l+1}$  is not in the dictionary before its addition to the dictionary in step 4. Thus the words added to the dictionary are distinct. For each occurrence of a word  $w$  in the parsing a word of the form  $wu$  with  $u \in \mathcal{U}$  is added to the dictionary. Since these are distinct,  $w$  cannot appear more than  $|\mathcal{U}|$  times in the parsing.  $\square$

**Lemma 3.7.** *Suppose  $u^n = u_1 \dots u_n$  is parsed into  $m(u^n)$  words  $w_1 \dots w_m$  by the algorithm. Then  $\lim_n m(u^n)/n = 0$ .*

*Proof.* There are  $|\mathcal{U}|^i$  words of length  $i$ , and by the previous lemma, each can appear at most  $|\mathcal{U}|$  times in the list  $w_1, \dots, w_m$ . As the algorithm does not parse the null string, at most  $F(k) = |\mathcal{U}| \sum_{i=1}^{k-1} |\mathcal{U}|^i$  words in the list are of length  $k-1$  or less, and each of the remaining words in the list has length  $k$  or more. Thus  $n \geq k[m - F(k)]$ . Consequently,

$$\limsup_n \frac{m(u^n)}{n} \leq \limsup_n \frac{n/k + F(k)}{n} = 1/k.$$

As  $k$  is arbitrary, the lemma follows.  $\square$

**Theorem 3.8.** *Let  $\ell(u^n)$  denote the number of bits produced by the algorithm after reading  $u^n$ . Then,  $\limsup_n \ell(u^n)/n \leq \limsup_n \frac{m(u^n)}{n} \log m(u^n)$ .*

*Proof.* As the dictionary size increases by 1 at each iteration of the algorithm, with  $m = m(u^n)$ ,  $\ell(u^n) = \sum_{i=0}^{m-1} \lceil \log_2(|\mathcal{U}| + i) \rceil$ . Thus

$$\ell(u^n) \leq m \log_2(|\mathcal{U}| + m - 1) + m = m \log_2 m + m \log_2(1 + (|\mathcal{U}| - 1)/m) + m,$$

and the lemma follows from  $\lim_n \frac{m(u^n)}{n} = 0$ .  $\square$

We now have an upper bound to the number of bits per letter the LZW algorithm requires to describe a sequence  $u_1 u_2 \dots$ . Next, we will derive a lower bound to the number of bits per letter produced by *any* information loss finite state machine that maps  $u_1 u_2 \dots$  to a sequence of bits (the terms ‘finite state machine’ and ‘information lossless’ will be defined formally in the paragraphs that follow). The lower bound will even apply to machines that may have been designed with prior knowledge of the sequence  $u_1 u_2 \dots$ , and most remarkably, will match the upper bound just derived above. That is to say, for any  $u_1 u_2 \dots$ , LZW competes well against any information lossless finite state machine. In particular, each prefix-free encoder  $c_n$  that appears in Corollary 3.4 can be implemented by a finite state information lossless machine. Consequently, once we obtain the lower bound, we will have proved:

**Theorem 3.9.** *If  $U_1, U_2$  is a stationary and ergodic stochastic process, then the number of bits per letter emitted by LZW when its input is  $U_1 U_2 \dots$  approaches  $\lim_n H(U^n)/n$  with probability one.*

### 3.2.1 Finite state information lossless encoders

For our purposes, a finite state machine is a device that reads the input sequence one symbol at a time. Each symbol of the input sequence belongs to a finite alphabet  $\mathcal{U}$  with  $|\mathcal{U}|$  symbols. The machine is in one of a finite number  $s$  of states before it reads a symbol, and goes to a new state determined by the old state and the symbol read. We will assume that the machine is in a fixed, known state  $z_1$  before it reads the first input symbol. The machine also produces a finite string of binary digits (possibly the null string) after each input. This output string is again a function of the old state and the input symbol. That is, when the infinite sequence  $u = u_1 u_2 \dots$  is given as the input, the encoder produces  $y = y_1 y_2 \dots$ , while visiting an infinite sequence of states  $z = z_1 z_2 \dots$ , given by

$$\begin{aligned} y_k &= f(z_k, u_k), \quad k \geq 1 \\ z_{k+1} &= g(z_k, u_k), \quad k \geq 1 \end{aligned}$$

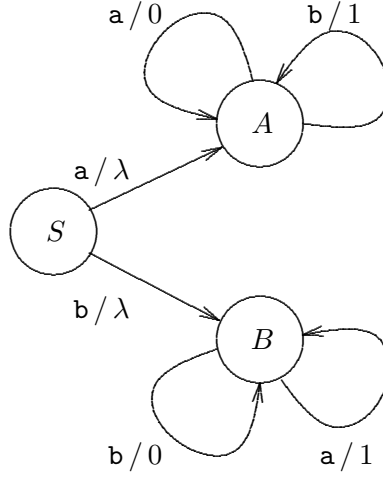
where the function  $f$  takes values on the set  $\{0, 1\}^*$  of finite binary strings, so that each  $y_k$  is a (perhaps null) binary string. A finite segment  $x_k x_{k+1} \dots x_j$  of a sequence  $x = x_1 x_2 \dots$  will be denoted by  $x_k^j$ , and by an abuse of the notation, the functions  $f$  and  $g$  will be extended to indicate the output sequence and the final state. Thus,  $f(z_k, u_k^j)$  will denote  $y_k^j$  and  $g(z_k, u_k^j)$  will denote  $z_{j+1}$ . Without loss of generality we will assume that any state  $z$  is reachable from the initial state  $z_1$  — i.e., that some input sequence will take the machine from state  $z_1$  to  $z$ .

To make the question of compressibility meaningful one has to require some sort of an ‘invertibility’ condition on the finite state encoders. Given the description of the finite state machine that encoded the string, and the starting state  $z_1$ , but (of course) without the knowledge of the input string, it should be possible to reconstruct the input string  $u$  from the output of the encoder  $y$ . A weaker requirement than this is the following: for any state  $z$  and two distinct input sequences  $v^m$  and  $\tilde{v}^n$ , either  $f(z, v^m) \neq f(z, \tilde{v}^n)$  or  $g(z, v^m) \neq g(z, \tilde{v}^n)$ . An encoder satisfying this condition will be called *information lossless* (IL). It is clear that if an encoder is not IL, then there is no hope to recover the input from the output, and thus every ‘invertible’ encoder is IL. <sup>2</sup>

We will need the following fact:

**Lemma 3.10.** *Suppose  $v_1, \dots, v_m$  are binary strings, with no string occurring more than  $k$  times. Then, writing  $m = \sum_{i=0}^{j-1} k2^i + r$  with  $0 \leq r < k2^j$ , we have  $\sum_{i=1}^m \text{length}(v_i) \geq k \sum_{i=0}^{j-1} i2^i + rj$ .*

<sup>2</sup>However, as illustrated in Figure 3.1, an IL encoder is not necessarily uniquely decodable. Starting from state  $S$ , two distinct input sequences will leave the encoder in distinct states if they have different first symbols, otherwise they will lead to different output sequences. Thus, the above encoder is IL. Nevertheless, no decoder can distinguish between the input sequences  $\mathbf{aaaa} \dots$  and  $\mathbf{bbbb} \dots$  by observing the output  $000 \dots$ .



A finite state machine with three states  $S$ ,  $A$  and  $B$ . The notation  $i / \text{output}$  means that the machine produces *output* in response to the input  $i$ .  $\lambda$  denotes the null output.

Figure 3.1: An IL encoder which is not uniquely decodable.

*Proof.* The set of binary strings ordered in increasing length consists of: 1 string of length 0, 2 strings of length 1,  $\dots$ ,  $2^i$  strings of length  $i$ ,  $\dots$ . The shortest total length for the  $v_i$ 's will be attained if the  $v_i$ 's are chosen by traversing the set of all binary strings in increasing length, each string repeated  $k$  times, until all strings of length  $j - 1$  or less are repeated  $k$  times, and we are left to find  $0 \leq r < k2^j$  strings, which are chosen from the set of strings of length  $j$ . The lower bound in the lemma is precisely the total length of this optimal collection.  $\square$

**Lemma 3.11.** Suppose  $v_1, \dots, v_m$  are binary strings, with no string occurring more than  $k$  times. Then,

$$\sum_{i=1}^m \text{length}(v_i) \geq m \log_2 \frac{m}{8k}. \quad (3.1)$$

*Proof.* Noting that  $\sum_{i=0}^{j-1} 2^i = 2^j - 1$  and  $\sum_{i=0}^{j-1} i2^i = (j-2)2^j + 2$ , the previous lemma states: writing  $m = k2^j - k + r$  with  $0 \leq r < k2^j$ , the total length of the  $v_i$ 's is lower bounded by

$$k((j-2)2^j + 2) + rj = (j-2)m + kj + 2r \geq (j-2)m.$$

As  $r < k2^j$ , we have  $m < k(2^{j+1} - 1)$ . Rearranging, we get  $2^{j+1} > 1 + m/k > m/k$ , and thus  $j - 2 > \log \frac{m}{8k}$ .  $\square$

Now we can state the following

**Lemma 3.12.** For any IL-encoder with  $s$  states,

$$\text{length}(y^n) \geq m(u^n) \log_2 \frac{m(u^n)}{8|U|s^2}. \quad (3.2)$$

where, as before,  $m(u^n)$  is the number of words in the parsing of  $u^n$  by LZW.

*Proof.* Let  $u^n = w_1 \dots w_m$  be the parsing of the input by LZW. Let  $z_i$  be the state the IL machine is in at just before it reads  $w_i$ , and  $z_{i+1}$  be the state just after it has read  $w_i$ . Let  $t_i$  be the binary string output by the IL machine while it reads  $w_i$ , so that  $y^n = t^m$ . No binary string  $t$  can occur in  $t_1, \dots, t_m$  more than  $k = s^2|\mathcal{U}|$  times. If it did, there will be a state-pair  $(z, z')$  that occurs among  $(z_i, z_{i+1})$  more than  $|\mathcal{U}|$  times with  $t_i = t$ . By Lemma 3.6 then there will be  $w_i \neq w_j$  with  $t_i = t_j = t$  and  $(z_i, z_{i+1}) = (z_j, z_{j+1}) = (z, z')$ . But this contradicts the IL property of the machine. Thus the output of the IL machine  $t_1 \dots t_m$  is a concatenation of  $m$  binary strings with no string occurring more than  $k = s^2|\mathcal{U}|$  times, so their total length is at least  $m \log \frac{m}{8k}$ .  $\square$

Using the lemma just proved, and by Lemma 3.7 we have

**Theorem 3.13.** *For any finite state information lossless encoder, the number of output bits  $\ell(u^n)$  produced by the encoder after reading  $u^n$  satisfies*

$$\limsup_n \ell(u^n)/n \geq \limsup_n \frac{m(u^n)}{n} \log m(u^n).$$

### 3.3 Quantization

Often it is not necessary to represent data exactly; an approximate representation suffices, e.g., we may be content if a sequence  $u^n$  of bits are reproduced as a sequence  $v^n$  of bits as long as  $v^n$  and  $u^n$  differ only at a few indices. We formulate this state of affairs as follows: our data is a sequence of random variables  $U^n$ , each  $U_i$  taking values in an alphabet  $\mathcal{U}$ . We are given a representation alphabet  $\mathcal{V}$  and also a distortion measure  $d : \mathcal{U} \times \mathcal{V} \rightarrow \mathbb{R}$  that gives the distortion caused by representing a data letter  $u$  by  $v$ . For  $n \geq 1$  define  $d_n(u^n, v^n) = \frac{1}{n} \sum_{i=1}^n d(u_i, v_i)$ . A quantizer consists of two maps:  $f_n : \mathcal{U}^n \rightarrow \{1, \dots, M\}$  and  $g_n : \{1, \dots, M\} \rightarrow \mathcal{V}^n$ . The data sequence  $u^n$  is mapped by the quantizer to  $f_n(u^n)$ , which is a  $\log M$  bit representation of the data, and reconstructed from this representation as  $v^n = g_n(f_n(u^n))$ . We measure the quality of the the quantizer by two criteria: its *rate*  $R = (\log M)/n$  (the number of bits per data letter), and its expected *distortion*  $\Delta = E[d_n(U^n, V^n)]$  where  $V^n = g_n(f_n(U^n))$ .

The following lemma says that one cannot have too small rate for a given level of distortion.

**Lemma 3.14.** *Suppose  $d$  is a distortion measure,  $U_1, U_2, \dots$  are i.i.d. with distribution  $P_U$ , and we have a quantizer with rate  $R$  and expected distortion  $\Delta$ . Then, there exists a distribution  $P_{UV}$  such that  $R \geq I(U; V)$  and  $\Delta = E[d(U, V)]$ .*

*Proof.* Let  $P_{U_i V_i}$  denote the joint distribution of  $(U_i, V_i)$  and set  $P_{UV} = \frac{1}{n} \sum_{i=1}^n P_{U_i V_i}$ . We will show that  $P_{UV}$  satisfies the conditions claimed in the lemma. First note that the  $U$ -marginal of each  $P_{U_i V_i}$  equals  $P_U$ , so  $P_{UV}$  indeed has  $U$  marginal  $P_U$ . Also,  $P_V = \frac{1}{n} \sum_i P_{V_i}$ . Now, by the data processing theorem for mutual information,  $nR \geq I(U^n; V^n)$ . As the  $U_i$ 's

are independent,  $H(U^n) = \sum_i H(U_i)$ . Also by the chain rule and removing terms from the conditioning  $H(U^n|V^n) \leq \sum_i H(U_i|V_i)$ . Thus

$$I(U^n; V^n) = H(U^n) - H(U^n|V^n) \geq \sum_{i=1}^n H(U_i) - H(U_i|V_i) = \sum_{i=1}^n I(U_i; V_i),$$

and consequently

$$R \geq \frac{1}{n} \sum_{i=1}^n I(U_i; V_i) = \frac{1}{n} \sum_{i=1}^n D(P_{U_i V_i} \| P_U P_{V_i}) \geq D(P_{UV} \| P_U P_V) = I(U; V)$$

where the last inequality is due to the convexity of  $D(\cdot \| \cdot)$ . Furthermore,

$$E[d(U, V)] = \sum_{u,v} P_{UV}(u, v) d(u, v) = \frac{1}{n} \sum_{i=1}^n P_{U_i V_i}(u, v) d(u, v) = E[d_n(U^n, V^n)] = \Delta. \quad \square$$

The lemma says that the possible rate distortion pairs when quantizing an i.i.d. source with distribution  $P$  lie in the following region of  $\mathbb{R}^2$ :

$$\bigcup_{P_{UV}: P_U=P} \{(R, \Delta) : R \geq I(U; V), \Delta = E[d(U, V)]\}.$$

We will now show that there are quantizers whose performance closely approximates the bound given in the lemma.

**Theorem 3.15.** *Given  $P_U$ , a distortion measure  $d$  and a joint distribution  $P_{UV}$ , for any  $\epsilon > 0$ , there is a quantizer  $(f_n, g_n)$  with rate  $R < I(U; V) + \epsilon$  and expected distortion  $\Delta$  satisfying  $|\Delta - E[d(U, V)]| < \epsilon$  when it is used to quantize i.i.d. data  $U^n$  with distribution  $P_U$ .*

To prove the theorem we will make use of the following lemma.

**Lemma 3.16.** *For all  $\delta > 0$ , for all large enough  $n$ , for every joint type  $P_{UV} \in \Pi_n(\mathcal{U} \times \mathcal{V})$ , there is a subset  $S := S(P_{UV}) \subset T^n(P_V)$  with  $|S| \leq \lceil 2^{n(I(U; V) + \delta)} \rceil$  so that for every  $u^n \in T^n(P_U)$  there is an  $v^n \in S$  for which  $(u^n, v^n) \in T^n(P_{UV})$ .*

*Proof.* Let  $M = \lceil 2^{n(I(U; V) + \delta)} \rceil$ , and choose  $V^n(1), \dots, V^n(M)$  independently and uniformly from  $T^n(P_V)$ . Let  $S$  be the (random) set  $\{V^n(1), \dots, V^n(M)\}$ . For each  $u^n \in T^n(P_U)$ , let

$$Z(u^n) = \prod_{m=1}^M \mathbb{1}\{(u^n, V^n(m)) \notin T^n(P_{UV})\}.$$

Note that the  $S$  has the claimed property if and only if  $Z = \sum_{u^n \in T^n(P_U)} Z(u^n) = 0$ . We will prove the existence of the set  $S$  with the property, by showing that  $E[Z] < 1$ . To that end, note that for a given  $u^n \in T^n(P_U)$ , the  $M$  random variables  $\mathbb{1}\{(u^n, V^n(m)) \in T^n(P_{UV})\}$  are independent, with expectation

$$\frac{|T^n(P_{UV})|}{|T^n(P_U)||T^n(P_V)|} \geq \left( \frac{n + |\mathcal{U} \times \mathcal{V}| - 1}{|\mathcal{U} \times \mathcal{V}| - 1} \right)^{-1} 2^{n[H(UV) - H(U) - H(V)]} \geq 2^{-n[I(U; V) + \delta/2]}.$$

Consequently,

$$E[Z(u^n)] \leq (1 - 2^{-n(I(U;V)+\delta/2)})^M \leq \exp(-M2^{-n(I(U;V)+\delta/2)}) \leq \exp(-2^{n\delta/2}).$$

As this expectation is doubly exponentially small in  $n$ , and since  $|T^n(P_U)|$  is only exponentially large, we conclude that for large enough  $n$  we will have  $E[Z] < 1$ .  $\square$

*Proof of the theorem.* Given  $P_{UV}$ , set  $R_0 = I(U;V)$  and  $\Delta_0 = E[d(U,V)]$ . Fix  $\delta > 0$ , and let  $\Omega$  the set of all joint distributions  $Q$  on  $\mathcal{U} \times \mathcal{V}$  such that  $I(U;V) < R_0 + \delta$  and  $|E[d(U,V)] - \Delta_0| < \delta$  when  $(U,V)$  has joint distribution  $Q$ . The distribution  $P_{UV}$  belongs to  $\Omega$ , thus  $\Omega$  is a non-empty open set of joint distributions. Let  $\Omega_0$  be the set of  $U$  marginals of the distributions in  $\Omega$ . Note that  $\Omega_0$  is also open and non-empty as it contains  $P_U$ . As a consequence  $D^* := \inf_{P \notin \Omega_0} D(P||P_U) = D(P^*||P_U)$  for some  $P^* \notin \Omega_0$ , and thus  $D^* > 0$  since  $P^*$  can't equal  $P_U \in \Omega_0$ .

Let  $S(Q)$  be as in the lemma above and set  $S = \bigcup S(Q)$  where the union is over all  $Q \in \Omega \cap \Pi_n(\mathcal{U} \times \mathcal{V})$ . For  $n$  large enough  $M = |S| \leq 2^{n(R_0+3\delta)}$ , since each  $S(Q)$  has  $|S(Q)| \leq \lceil 2^{n(R_0+2\delta)} \rceil$  and the union contains only  $|\Pi_n|$  (which is a polynomial in  $n$ ) such sets. Assign to each  $v^n$  in  $S$  a unique index in  $\{1, \dots, M\}$ , and for  $i = 1, \dots, M$ , let  $g_n(i)$  be the element of  $S$  with index  $i$ .

By the construction of the set  $S$ , for every  $P \in \Omega_0 \cap \Pi_n(\mathcal{U})$  and every  $u^n \in T^n(P)$  we can find a  $v^n$  in  $S$  such that the type of  $(u^n, v^n)$  belongs to  $\Omega$ , and thus  $|d_n(u^n, v^n) - \Delta_0| < \delta$ . For  $u^n \in \bigcup_{P \in \Omega_0} T^n(P)$  define  $f_n(u^n)$  to be the index of such a  $v^n$  in  $S$ . For  $u^n$  not in this union define  $f_n(u^n)$  in any arbitrary fashion, e.g., by setting  $f_n(u^n) = 1$ . For those  $u^n$ , with  $v^n = g_n(f_n(u^n))$ , we do not necessarily have a small value for  $|d_n(u^n, v^n) - \Delta_0|$ , but nevertheless we can say  $|d_n(u^n, v^n) - \Delta_0| < 2 \max_{u,v} |d(u,v)|$ .

When  $U_1, \dots, U_n$  are i.i.d. with distribution  $P_U$ , the probability that the type of  $U^n$  does not belong to  $\Omega_0$  decays exponentially to zero as  $2^{-nD^*}$ . Consequently, with  $V^n = g_n(f_n(U^n))$

$$\eta_n = \Pr(|d_n(U^n, V^n) - \Delta_0| \geq \delta)$$

decays to zero like  $2^{-nD^*}$ . Thus the quantizer  $(f_n, g_n)$  has rate  $\frac{1}{n} \log M \leq R_0 + 3\delta$ , and distortion  $\Delta = E[d_n(U^n, V^n)]$  satisfying

$$|\Delta - \Delta_0| \leq E[|d(U^n, V^n) - \Delta_0|] < \delta + \eta_n 2 \max_{u,v} |d(u,v)| < 2\delta$$

for  $n$  large enough. As  $\delta > 0$  is arbitrary, the theorem follows.  $\square$

In the light of the Theorem above and Lemma 3.14 that precedes it, it makes sense to define

**Definition 3.3.** Given source and reproduction alphabets  $\mathcal{U}$  and  $\mathcal{V}$ , a distortion function  $d : \mathcal{U} \times \mathcal{V} \rightarrow \mathbb{R}$ , and a distribution  $P_U$  on  $\mathcal{U}$ , define the *rate distortion function*  $R(\Delta)$  and *distortion rate function*  $\Delta(R)$  as

$$R(\Delta) = \inf\{I(U;V) : E[d(U,V)] \leq \Delta\}, \quad \Delta(R) = \inf\{E[d(U,V)] : I(U;V) \leq R\}$$

where the infima are over all joint distributions  $P_{UV}$  with  $U$  marginal  $P_U$ .

Lemma 3.14 says that any quantizer with rate at most  $R$  must have expected distortion at least  $\Delta(R)$ , and any quantizer with expected distortion at most  $\Delta$  must have rate at least  $R(\Delta)$ . The theorem says that these bounds are achievable arbitrarily closely.

As an example consider  $\mathcal{U} = \mathcal{V} = \{0, 1\}$ ,  $d(u, v) = \mathbb{1}\{u \neq v\}$ , and  $P_U(1) = p \leq 1/2$ . It is easy to show that

$$R(\Delta) = \begin{cases} 0 & \Delta \geq p \\ h_2(p) - h_2(\Delta) & 0 \leq \Delta < p \end{cases}$$

where  $h_2(q) = -q \log_2 q - (1 - q) \log_2 (1 - q)$  is the binary entropy function. In particular, for  $p = 1/2$  and  $\Delta = 0.1$ ,  $R(\Delta) \approx 0.531$ , so it is possible to compress binary data by almost a factor half and still be 90 percent accurate. The naive approach for achieving the same accuracy would have retained 80 percent of the data bits (the remaining 20 percent can then be guessed in any fashion; the guesses would be right with probability  $1/2$ , resulting in 90 percent overall accuracy).

### 3.4 Problems

**Problem 3.1** (Elias coding). Let  $0^n$  denote a sequence of  $n$  zeros. Consider the code (the subscript  $U$  a mnemonic for ‘Unary’),  $\mathcal{C}_U : \{1, 2, \dots\} \rightarrow \{0, 1\}^*$  for the positive integers defined as  $\mathcal{C}_U(n) = 0^{n-1}$ .

- (a) Is  $\mathcal{C}_U$  injective? Is it prefix-free?

Consider the code (the subscript  $B$  a mnemonic for ‘Binary’),  $\mathcal{C}_B : \{1, 2, \dots\} \rightarrow \{0, 1\}^*$  where  $\mathcal{C}_B(n)$  is the binary expansion of  $n$ . I.e.,  $\mathcal{C}_B(1) = 1$ ,  $\mathcal{C}_B(2) = 10$ ,  $\mathcal{C}_B(3) = 11$ ,  $\mathcal{C}_B(4) = 100$ , .... Note that

$$\text{length} \mathcal{C}_B(n) = \lceil \log_2(n+1) \rceil = 1 + \lfloor \log_2 n \rfloor.$$

- (b) Is  $\mathcal{C}_B$  injective? Is it prefix-free?

With  $k(n) = \text{length} \mathcal{C}_B(n)$ , define  $\mathcal{C}_0(n) = \mathcal{C}_U(k(n))\mathcal{C}_B(n)$ .

- (c) Show that  $\mathcal{C}_0$  is a prefix-free code for the positive integers. To do so, you may find it easier to describe how you would recover  $n_1, n_2, \dots$  from the concatenation of their codewords  $\mathcal{C}_0(n_1)\mathcal{C}_0(n_2)\dots$ .
- (d) What is  $\text{length}(\mathcal{C}_0(n))$ ?

Now consider  $\mathcal{C}_1(n) = \mathcal{C}_0(k(n))\mathcal{C}_B(n)$ .

- (e) Show that  $\mathcal{C}_1$  is a prefix-free code for the positive integers, and show that  $\text{length}(\mathcal{C}_1(n)) = 2 + 2\lfloor \log(1 + \lfloor \log n \rfloor) \rfloor + \lfloor \log n \rfloor \leq 2 + 2\log(1 + \log n) + \log n$ .

Suppose  $U$  is a random variable taking values in the positive integers with  $\Pr(U = 1) \geq \Pr(U = 2) \geq \dots$ .

- (f) Show that  $E[\log U] \leq H(U)$ , [Hint: first show  $i \Pr(U = i) \leq 1$ ], and conclude that

$$E[\text{length}\mathcal{C}_1(U)] \leq H(U) + 2\log(1 + H(U)) + 2.$$

**Problem 3.2** (Universal codes). Suppose we have an alphabet  $\mathcal{U}$ , and let  $\Pi$  denote the set of distributions on  $\mathcal{U}$ . Suppose we are given a family of  $S$  of distributions on  $\mathcal{U}$ , i.e.,  $S \subset \Pi$ . For now, assume that  $S$  is finite.

Define the distribution  $Q_S \in \Pi$

$$Q_S(u) = Z^{-1} \max_{P \in S} P(u)$$

where the normalizing constant  $Z = Z(S) = \sum_u \max_{P \in S} P(u)$  ensures that  $Q_S$  is a distribution.

- (a) Show that  $D(P\|Q) \leq \log Z \leq \log |S|$  for every  $P \in S$ .
- (b) For any  $S$ , show that there is a prefix-free code  $\mathcal{C} : \mathcal{U} \rightarrow \{0,1\}^*$  such that for any random variable  $U$  with distribution  $P \in S$ ,

$$E[\text{length}\mathcal{C}(U)] \leq H(U) + \log Z + 1.$$

(Note that  $\mathcal{C}$  is designed on the knowledge of  $S$  alone, it cannot change on the basis of the choice of  $P$ .) [Hint: consider  $L(u) = -\log_2 Q_S(u)$  as an ‘almost’ length function.]

- (c) Now suppose that  $S$  is not necessarily finite, but there is a finite  $S_0 \subset \Pi$  such that for each  $u \in \mathcal{U}$ ,  $\sup_{P \in S} P(u) \leq \max_{P \in S_0} P(u)$ . Show that  $Z(S) \leq |S_0|$ .

Now suppose  $\mathcal{U} = \{0,1\}^m$ . For  $\theta \in [0,1]$  and  $(x_1, \dots, x_m) \in \mathcal{U}$ , let

$$P_\theta(x_1, \dots, x_m) = \prod_i \theta^{x_i} (1 - \theta)^{1-x_i}.$$

(This is a fancy way to say that the random variable  $U = (X_1, \dots, X_m)$  has i.i.d. Bernoulli  $\theta$  components). Let  $S = \{P_\theta : \theta \in [0,1]\}$ .

- (d) Show that for  $u = (x_1, \dots, x_m) \in \{0,1\}^m$

$$\max_{\theta} P_\theta(x_1, \dots, x_m) = P_{k/m}(x_1, \dots, x_m)$$

where  $k = \sum_i x_i$ .

- (e) Show that there is a prefix-free code  $\mathcal{C} : \{0,1\}^m \rightarrow \{0,1\}^*$  such that whenever  $X_1, \dots, X_m$  are i.i.d. Bernoulli,

$$\frac{1}{m} E[\text{length}\mathcal{C}(X_1, \dots, X_m)] \leq H(X_1) + \frac{1 + \log_2(1 + m)}{m}.$$



**Problem 3.3** (Prediction and coding). After observing a binary sequence  $u_1, \dots, u_i$ , that contains  $n_0(u^i)$  zeros and  $n_1(u^i)$  ones, we are asked to estimate the probability that the next observation,  $u_{i+1}$  will be 0. One class of estimators are of the form

$$\hat{P}_{U_{i+1}|U^i}(0|u^i) = \frac{n_0(u^i) + \alpha}{n_0(u^i) + n_1(u^i) + 2\alpha} \quad \hat{P}_{U_{i+1}|U^i}(1|u^i) = \frac{n_1(u^i) + \alpha}{n_0(u^i) + n_1(u^i) + 2\alpha}.$$

We will consider the case  $\alpha = 1/2$ , this is known as the Krichevsky–Trofimov estimator. Note that for  $i = 0$  we get  $\hat{P}_{U_1}(0) = \hat{P}_{U_1}(1) = 1/2$ .

Consider now the joint distribution  $\hat{P}(u^n)$  on  $\{0, 1\}^n$  induced by this estimator,

$$\hat{P}(u^n) = \prod_{i=1}^n \hat{P}_{U_i|U^{i-1}}(u_i|u^{i-1}).$$

(a) Show, by induction on  $n$  that, for any  $n$  and any  $u^n \in \{0, 1\}^n$ ,

$$\hat{P}(u_1, \dots, u_n) \geq \frac{1}{2\sqrt{n}} \left(\frac{n_0}{n}\right)^{n_0} \left(\frac{n_1}{n}\right)^{n_1},$$

where  $n_0 = n_0(u^n)$  and  $n_1 = n_1(u^n)$ .

(b) Conclude that there is a prefix-free code  $\mathcal{C} : \mathcal{U} \rightarrow \{0, 1\}^*$  such that

$$\text{length}\mathcal{C}(u_1, \dots, u_n) \leq nh_2\left(\frac{n_0(u^n)}{n}\right) + \frac{1}{2} \log n + 2,$$

with  $h_2(x) = -x \log x - (1-x) \log(1-x)$ .

(c) Show that if  $U_1, \dots, U_n$  are i.i.d. Bernoulli, then

$$\frac{1}{n} E[\text{length}\mathcal{C}(U_1, \dots, U_n)] \leq H(U_1) + \frac{1}{2n} \log n + \frac{2}{n}$$

**Problem 3.4** (Lempel Ziv 78). Suppose  $\dots, U_{-1}, U_0, U_1, \dots$  is a stationary process, i.e., for any  $k = 1, 2, \dots$ , any  $u_0, \dots, u_{k-1}$ , and any  $n = \dots, -1, 0, 1, \dots$

$$\Pr(U_n \dots U_{n+k-1} = u_0 \dots u_{k-1}) = \Pr(U_0 \dots U_{k-1} = u_0 \dots u_{k-1}).$$

Suppose also that  $U$  is a recurrent process, i.e., any letter  $u_0$  with  $\Pr(U_0 = u_0) > 0$ , the event  $A = \{\text{there exists } i \geq 0 \text{ and } j > 0 \text{ such that } U_i = U_{-j} = u_0\}$  has  $\Pr(A) = 1$ . (That is, a positive probability letter  $u_0$  will occur infinitely often.)

Fix  $u_0$  with  $\Pr(U_0 = u_0) > 0$ . For  $i \geq 0$  and  $j < 0$ , let

$$A_{ij} = \{U_i = u_0\} \cap \{U_{-j} = u_0\} \cap \bigcap_{k=-j+1}^{i-1} \{U_k \neq u_0\}$$

denote the event that  $j$  is the last time before time 0 that  $u_0$  was seen and  $i$  was the first time after time 0 that  $u_0$  is seen.

(a) Show that  $\sum_{i \geq 0, j > 0} \Pr(A_{ij}) = \Pr(A) = 1$ .

(b) Show that  $\Pr(A_{ij}) = f(i + j)$ , where

$$f(k) = \Pr(U_{-k} = u_0, U_{-l} \neq u_0 \text{ for } l = 1, \dots, k-1, U_0 = u_0).$$

(c) Using (a) and (b), show that

$$1 = \sum_{k \geq 1} k f(k) = 1.$$

(d) Let  $K = \inf\{k > 0 : U_{-k} = u_0\}$  (i.e., the negative index of the most recent time before time 0  $u_0$  was seen). Observe that the event  $\{K = k, U_0 = u_0\}$  is the event whose probability is  $f(k)$ . Using (c) show that

$$E[K | U_0 = u_0] = 1 / \Pr(U_0 = u_0)$$

and that  $E[\log K] \leq H(U_0)$ .

Suppose we have a stationary and ergodic source  $\dots, X_{-1}, X_0, X_1, \dots$ . This means, in particular, that for any  $n > 0$ , the process  $\{U_i\}$  defined by  $U_i = (X_i, X_{i+1}, X_{i+n-1})$  is stationary and recurrent.

Fix a sequence  $x_0, \dots, x_{n-1}$  with  $\Pr((X_0 \dots X_{n-1}) = (x_0 \dots x_{n-1})) > 0$ . Let

$$K = \inf\{k > 0 : (X_{-k} \dots X_{-k+n-1}) = (x_0 \dots x_{n-1})\}.$$

(e) Show that  $E[\log K] \leq H(X_0 \dots X_{n-1})$ .

(f) Consider the following data compression method. Assuming that the encoder has already described the infinite past  $\dots, X_{-2}, X_{-1}$  to the decoder, he describes  $X_0, \dots, X_{n-1}$  by (i) finding the most recent occurrence  $X_0 \dots X_{n-1}$  in the past, (ii) describing the index  $K$  of this occurrence by the method of problem 1(f). Now that the decoder knows  $\dots, X_{n-1}$ , the encoder describes  $X_n \dots X_{2n-1}$  in the same way, etc. Show that this method uses fewer than

$$\frac{1}{n} H(X_0 \dots X_{n-1}) + \frac{2}{n} \log(1 + H(X_0 \dots X_{n-1})) + \frac{2}{n}$$

bits per letter on the average.

**Problem 3.5** (Quantization with two criteria). Suppose  $U^n$  has i.i.d. components with distribution  $P$ . We want to describe  $U^n$  at rate  $R$ , i.e., we want to design a function  $f : \mathcal{U}^n \rightarrow \{1, \dots, 2^{nR}\}$ .

We are given two distortion measures  $d_1 : \mathcal{U} \times \mathcal{V}_1 \rightarrow \mathbb{R}$  and  $d_2 : \mathcal{U} \times \mathcal{V}_2 \rightarrow \mathbb{R}$ , and we wish to ensure that from  $i = f(U^n)$  we can reconstruct  $V_1^n = g_1(i) \in \mathcal{V}_1^n$  and  $V_2^n = g_2(i) \in \mathcal{V}_2^n$  so that

$$E[d_1(U^n, V_1^n)] \leq D_1 \quad \text{and} \quad E[d_2(U^n, V_2^n)] \leq D_2$$

with given distortion criteria  $D_1$  and  $D_2$ . (As in class  $d(U^n, V^n) = \frac{1}{n} \sum_{i=1}^n d(U_i, V_i)$ .)

- 
- (a) What is the rate distortion function  $R(D_1, D_2)$ ?
- (b) Suppose  $R_1(D_1)$  is the rate distortion function with the first distortion criterion alone, and  $R_2(D_2)$  is the rate distortion function with the second criterion alone. What relationship exists between  $R(D_1, D_2)$  and  $R_1(D_1) + R_2(D_2)$ ?



## Chapter 4

# Exponential Families and Maximum Entropy Distributions

Exponential families are a class of parametrized distributions. They are important for several reasons. First, many “standard” distributions we are well acquainted with (like the Gaussian distribution) are members of this family. Therefore, they appear frequently in applications. Second, all members of this family have nice theoretical properties. Hence, rather than discussing these properties for each member of this family, it is convenient to discuss them for the whole family at once.

To give some “practical motivations,” in machine learning exponential families are used in the context of classification, giving rise to so-called *generalized linear models*. Further, these are the distributions that maximize the *entropy* given constraints on the moments. E.g., we will see that the Gaussian has the maximum entropy of any family with a given second moment constraint. It is therefore natural to consider such distributions as prior distributions since they make “the least assumptions” if all we know are constraints on moments.

In the following it will be convenient to treat continuous and discrete cases together. So we will assume that we have a space  $\mathcal{X}$  and a measure  $\nu$ . Let us list our most important examples:

1. Reals: Let  $\mathcal{X} = \mathbb{R}$  and let  $\nu$  be the Lebesgue measure on  $\mathbb{R}$ ; recall that  $\nu$  assigns to intervals  $[a, b]$ ,  $a \leq b$ , the measure  $\nu([a, b]) = b - a$ .
2. Bernoulli: Let  $\mathcal{X} = \{0, 1\}$  and let  $\nu$  be the counting measure on  $\{0, 1\}$ ; i.e.,  $\nu(\emptyset) = 0$ ,  $\nu(\{0\}) = \nu(\{1\}) = 1$ ,  $\nu(\{0, 1\}) = 2$
3. Poisson: Let  $\mathcal{X} = \mathbb{Z}$  and let  $\nu$  be the counting measure on  $\mathbb{Z}$ ; i.e., for  $S \subseteq \mathbb{Z}$ ,  $\nu(S) = |S|$ , the cardinality of the set  $S$ .

In the sequel it will hopefully be clear what the base measure is and so we will typically not include it in our notation.

We will be relatively terse. If you want to dig deeper, we recommend the lecture notes by John Duchi [1], Chapter 6 and 7, or the extensive monograph by Wainwright and Jordan [2].

## 4.1 Definition

**Definition 4.1.** Let  $\mathcal{X}$  be a given alphabet and let  $\phi : \mathcal{X} \rightarrow \mathbb{R}^d$ ,  $d \in \mathbb{N}$ . The *exponential family* associated with  $\phi$  is the set of distributions parametrized by  $\theta \in \mathbb{R}^d$  with densities given by

$$p_\theta(x) = h(x)e^{\langle \theta, \phi(x) \rangle - A(\theta)}.$$

Note that  $A(\theta)$  is a *normalizing* constant. As such it might not seem to play an important role. But, as we will discuss soon, it in fact encodes (in its derivatives) crucial information. The function  $A(\theta)$  is some-times called the *log-partition function* (the partition function is a term used in statistical physics for the normalization constant and  $A$  is the log of this). In statistics it is known as the *cumulant* function.

In our definition of an exponential family we included the term  $h(x)$ . In principle this term can be absorbed by the underlying measure  $\nu(x)$  and is in this sense redundant. But it might sometimes be more “natural” to represent a distribution in this way. For all our subsequent computations and proofs of properties it does not really matter which point of view we take, i.e., if we explicitly write out the term  $h(x)$  think of it as being included in the underlying measure.

## 4.2 Examples

**Example 4.1** (Gaussian). Let  $\mathcal{X} = \mathbb{R}$  and let  $\nu$  be the Lebesgue measure on  $\mathbb{R}$ . Then the density of the normal distribution with mean  $\mu$  and variance  $\sigma^2$  can be written as

$$\begin{aligned} p(x) &= \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \\ &= e^{x\frac{\mu}{\sigma^2} - x^2\frac{1}{2\sigma^2} - [\frac{\mu^2}{2\sigma^2} + \frac{1}{2}\ln(2\pi\sigma^2)]} \\ &= e^{\langle \theta, \phi(x) \rangle - A(\theta)}, \end{aligned}$$

where  $h(x) = 1$ ,  $\theta = (\frac{\mu}{\sigma^2}, -\frac{1}{2\sigma^2})^\top$ ,  $\phi(x) = (x, x^2)^\top$ , and

$$A = \frac{\mu^2}{2\sigma^2} + \frac{1}{2}\ln(2\pi\sigma^2) = -\frac{\theta_1^2}{4\theta_2} - \frac{1}{2}\ln(-\theta_2/\pi).$$

Note that  $\phi(x)$  is a vector of dimension 2, reflecting the fact that the Gaussian has two degrees of freedom. Further, we have the following bijective relationships

$$\begin{aligned} \theta &= (\theta_1, \theta_2)^\top = \left(\frac{\mu}{\sigma^2}, -\frac{1}{2\sigma^2}\right)^\top, \\ (\mu, \sigma^2)^\top &= \left(-\frac{\theta_1}{2\theta_2}, -\frac{1}{2\theta_2}\right)^\top. \end{aligned}$$

**Example 4.2** (Bernoulli). Let  $\mathcal{X} = \{0, 1\}$  and let  $\nu$  be the counting measure on  $\mathcal{X}$ . We can represent the Bernoulli distribution with  $P(X = 1) = p$  in the form

$$\begin{aligned} P(X = x) &= p^x (1 - p)^{1-x} \\ &= e^{x \ln p + (1-x) \ln(1-p)} \\ &= e^{x \ln \frac{p}{1-p} + \ln(1-p)} \\ &= h(x) e^{\langle \theta, \phi(x) \rangle - A(\theta)}, \end{aligned}$$

where  $h(x) = 1$ ,  $\theta = \ln \frac{p}{1-p}$ , so that  $p = \frac{e^\theta}{1+e^\theta}$ ,  $\phi(x) = x$ , and  $A(\theta) = -\ln(1-p) = -\ln(1 - \frac{e^\theta}{1+e^\theta}) = \ln(1 + e^\theta)$ .

**Example 4.3** (Poisson). Let  $\mathcal{X} = \mathbb{N}$  and let  $\nu$  be the counting measure on  $\mathcal{X}$ . We can represent the Poisson distribution with parameter  $\lambda$  in the form

$$\begin{aligned} P(X = x) &= \frac{\lambda^x e^{-\lambda}}{x!} \\ &= \frac{1}{x!} e^{x \ln(\lambda) - \lambda} \\ &= \frac{1}{x!} e^{\theta x - e^\theta} \\ &= h(x) e^{\langle \theta, \phi(x) \rangle - A(\theta)}, \end{aligned}$$

where  $h(x) = 1/x!$ ,  $\theta = \ln(\lambda)$ ,  $\phi(x) = x$ , and  $A(\theta) = e^\theta$ .

**Example 4.4** (Multinomial). Let  $\mathcal{X} = \{0, \dots, n\}^d$  and let  $\nu$  be the counting measure on  $\mathcal{X}$ . Then the multinomial distribution with parameter  $\alpha = (\alpha_1, \dots, \alpha_d)$  can be expressed as

$$\begin{aligned} P(X_1 = x_1, \dots, X_d = x_d) &= \binom{n}{x_1, \dots, x_d} \prod_{i=1}^d \alpha_i^{x_i} \\ &= \binom{n}{x_1, \dots, x_d} e^{\sum_{i=1}^d \ln(\alpha_i) x_i} \\ &= h(x) e^{\langle \theta, \phi(x) \rangle - A(\theta)}, \end{aligned}$$

where

$$h(x) = \binom{n}{x_1, \dots, x_d},$$

$\theta = (\ln(\alpha_1), \dots, \ln(\alpha_d))$ ,  $\phi(x) = x = (x_1, \dots, x_d)$ , and  $A(\theta) = 0$ .

**Example 4.5** (Dirichlet). The Dirichlet distribution of order  $d \geq 2$  with parameter  $\alpha = (\alpha_1, \dots, \alpha_d)$ ,  $\alpha_i > 0$ , has a density with respect to the Lebesgue measure on  $\mathbb{R}^{d-1}$  of the

form

$$\begin{aligned} p_\theta(x) &= \frac{1}{B(\alpha)} \prod_{i=1}^d x_i^{\alpha_i-1} \\ &= h(x) e^{\langle \theta, \phi(x) \rangle - A(\theta)} \end{aligned}$$

where  $x$  belongs to the  $(d-1)$ -dimensional simplex, i.e.,  $\sum_{i=1}^d x_i = 1$ ,  $x_i \geq 0$ , and where

$$B(\alpha) = \frac{\prod_{i=1}^d \Gamma(\alpha_i)}{\Gamma(\sum_{i=1}^d \alpha_i)}.$$

Further,  $h(x) = 1$ ,  $\theta = (\alpha_1 - 1, \dots, \alpha_d - 1)$ ,  $\phi(x) = (\ln(x_1), \dots, \ln(x_d))$ , and  $A(\theta) = \ln(B(\alpha))$ .

### 4.3 Convexity of $A(\theta)$

**Theorem 4.1.** *Let  $\Theta = \{\theta \in \mathbb{R}^d : A(\theta) < \infty\}$ . The log-partition function  $A(\theta)$  is convex in  $\theta$  on  $\Theta$ .*

*Proof.* Let  $\theta_\lambda = \lambda\theta_1 + (1-\lambda)\theta_2$ ,  $\theta_1, \theta_2 \in \Theta$ . Let  $p = \frac{1}{\lambda}$  and  $q = \frac{1}{1-\lambda}$  so that  $p\lambda = q(1-\lambda) = 1$ . Note that  $1/p + 1/q = \lambda + (1-\lambda) = 1$  and that  $p, q \in [1, \infty]$ . Hölder's inequality states that  $\|fg\|_1 \leq \|f\|_p \|g\|_q$ . In more detail,

$$\left( \int |f(x)g(x)| d\nu(x) \right) \leq \left( \int |f(x)|^p d\nu(x) \right)^{\frac{1}{p}} \left( \int |g(x)|^q d\nu(x) \right)^{\frac{1}{q}}.$$

Recall that  $p_\theta(x) = h(x) e^{\langle \theta, \phi(x) \rangle - A(\theta)}$  so that  $1 = \int p_\theta(x) d\nu(x) = \left[ \int h(x) e^{\langle \theta, \phi(x) \rangle} d\nu(x) \right] e^{-A(\theta)}$ , or,  $A(\theta) = \ln \left[ \int h(x) e^{\langle \theta, \phi(x) \rangle} d\nu(x) \right]$ . We have

$$\begin{aligned} A(\theta_\lambda) &= \ln \left[ \int h(x) e^{\langle \theta_\lambda, \phi(x) \rangle} d\nu(x) \right] \\ &= \ln \left[ \int \underbrace{\left( h(x) e^{\langle \theta_1, \phi(x) \rangle} \right)^\lambda}_{f(x)} \underbrace{\left( h(x) e^{\langle \theta_2, \phi(x) \rangle} \right)^{(1-\lambda)}}_{g(x)} d\nu(x) \right] \\ &\stackrel{\text{Hölder}}{\leq} \ln \left[ \left( \int \left( h(x) e^{\langle \theta_1, \phi(x) \rangle} \right)^{p\lambda} d\nu(x) \right)^{\frac{1}{p}} \left( \int \left( h(x) e^{\langle \theta_2, \phi(x) \rangle} \right)^{q(1-\lambda)} d\nu(x) \right)^{\frac{1}{q}} \right] \\ &= \ln \left[ \left( \int \left( h(x) e^{\langle \theta_1, \phi(x) \rangle} \right)^{p\lambda} d\nu(x) \right)^{\frac{1}{p}} \right] + \ln \left[ \left( \int \left( h(x) e^{\langle \theta_2, \phi(x) \rangle} \right)^{q(1-\lambda)} d\nu(x) \right)^{\frac{1}{q}} \right] \\ &\stackrel{p\lambda=q(1-\lambda)=1}{=} \frac{1}{p} \ln \left[ \left( \int h(x) e^{\langle \theta_1, \phi(x) \rangle} d\nu(x) \right) \right] + \frac{1}{q} \ln \left[ \left( \int h(x) e^{\langle \theta_2, \phi(x) \rangle} d\nu(x) \right) \right] \\ &= \lambda A(\theta_1) + (1-\lambda) A(\theta_2). \end{aligned}$$

□



## 4.4 Derivatives of $A(\theta)$

Without proof we state that  $A(\theta)$  is infinitely often differentiable on  $\Theta$ . In particular the first two derivatives are of interest to us.

Let us compute the first derivative (gradient). We have

$$\begin{aligned}\nabla A(\theta) &= \nabla \ln \int h(x) e^{\langle \theta, \phi(x) \rangle} d\nu(x) \\ &= \frac{\int \nabla h(x) e^{\langle \theta, \phi(x) \rangle} d\nu(x)}{\int h(x) e^{\langle \theta, \phi(x) \rangle} d\nu(x)} \\ &= \frac{\int h(x) e^{\langle \theta, \phi(x) \rangle} \phi(x) d\nu(x)}{e^{A(\theta)}} \\ &= \int h(x) e^{\langle \theta, \phi(x) \rangle - A(\theta)} \phi(x) d\nu(x) \\ &= \mathbb{E}[\phi(x)].\end{aligned}$$

For future reference, let us record that

$$\nabla A(\theta) = \mathbb{E}[\phi(x)]. \quad (4.1)$$

In a similar manner we have

$$\nabla^2 A(\theta) = \mathbb{E}[\phi(x)\phi(x)^\top] - \mathbb{E}[\phi(x)]\mathbb{E}[\phi(x)^\top].$$

Note that this gives us a second proof that  $A(\theta)$  is convex since we see that the Hessian of  $A(\theta)$  is a covariance matrix and hence positive semidefinite.

**Example 4.6** (Bernoulli). For the Bernoulli distribution we have seen that  $A(\theta) = \ln(1+e^\theta)$  and  $\theta = \ln \frac{p}{1-p}$ . Therefore,

$$\begin{aligned}\frac{dA(\theta)}{d\theta} &= \frac{d \ln(1+e^\theta)}{d\theta} = \frac{e^\theta}{1+e^\theta} = \sigma(\theta) = p, \\ \frac{d^2 A(\theta)}{d\theta^2} &= \frac{d\sigma(\theta)}{d\theta} = \sigma(\theta)(1-\sigma(\theta)) = p(1-p).\end{aligned}$$

## 4.5 Application to Parameter Estimation and Machine Learning

The convexity of  $A(\theta)$  is one of the main reason why this family of distributions is so convenient to work with.

Assume that we have a set of samples  $\mathbf{x}_1, \dots, \mathbf{x}_n$  and that we assume that they are iid according to an exponential family with an unknown parameter  $\theta$ . We want to estimate this parameter.

We can then write down the likelihood as

$$p_\theta(\mathbf{x}_1, \dots, \mathbf{x}_N) = \prod_{n=1}^N h(\mathbf{x}_n) e^{\langle \theta, \phi(\mathbf{x}_n) \rangle - A(\theta)}. \quad (4.2)$$

Instead of maximizing this likelihood we can equivalently take the log of this expression, multiply by minus one, and minimize instead. This gives us

$$-\ln p_\theta(\mathbf{x}_1, \dots, \mathbf{x}_N) = \sum_{n=1}^N [-\ln(h(\mathbf{x}_n)) - \langle \theta, \phi(\mathbf{x}_n) \rangle + A(\theta)] \quad (4.3)$$

Now note that the function on the right is convex – it is the sum of the constant (with respect to  $\theta$ )  $-\sum_{n=1}^N \ln(h(\mathbf{x}_n))$ , the linear function  $-\sum_{n=1}^N \langle \theta, \phi(\mathbf{x}_n) \rangle$  and the convex function  $NA(\theta)$ . Greedy, local algorithms are therefore expected to work well in locating the optimal parameter  $\theta$ .

A word of caution is in order here. Just because a function is convex, it does not mean that it is easy to minimize. We need in addition that the function itself (and perhaps its derivative) is easy to compute. If you look ahead at the *Ising model* described in Example 4.13, then you will see that even though the function  $A(\theta)$  to be minimized is convex, there is no low-complexity algorithm known to accomplish this minimization since the computation of  $A(\theta)$  requires in general exponential effort.

## 4.6 Conjugate Priors

In the previous section we considered an application in ML where we estimated the underlying parameter  $\theta$ , given some iid samples from the distribution by maximizing the likelihood. We have seen that for exponential distributions the underlying maximization problem is “simple” since the underlying function is convex. The justification for maximizing the likelihood is that under some technical conditions this leads to a consistent estimator (see Section 4.10.2).

Alternatively, in the Bayesian setting, we assume that there is a prior on the set of parameters and we will then maximize the posterior instead. In this case the question is what prior we should pick. One part of the question is what priors are “meaningful” or “appropriate.” Leaving out this question for the moment, there is still the question what priors lead to “manageable” computational tasks, e.g., convex functions to be minimized. Here is where conjugate priors enter. If we start with a likelihood that is a member of an exponential family and use as a conjugate prior then we end up again with an element of the exponential family. Rather than discussing this in the abstract, let us look at some important examples.

**Example 4.7** (Bernoulli). Consider a Bernoulli distribution with parameter  $p$ ,

$$P_p(X = x) = p^x(1 - p)^{1-x},$$

where we recall  $x \in \{0, 1\}$ . Assume that the parameter  $p \in [0, 1]$  is unknown and has a prior  $q(p)$  of the form

$$q(p) = K(\alpha_1, \alpha_2)p^{\alpha_1-1}(1 - p)^{\alpha_2-1},$$

where  $\alpha, \alpha_2 > 0$  so that the density can be normalized and where  $K(\alpha_1, \alpha_2)$  is the normalization constant,  $K(\alpha_1, \alpha_2) = \frac{\Gamma(\alpha_1 + \alpha_2)}{\Gamma(\alpha_1)\Gamma(\alpha_2)}$ . That we write the parameters in the form  $\alpha_i - 1$

(instead of  $\alpha_i$ ) is purely for convenience. This is called the *beta* distribution. Note that the beta distribution is an member of the exponential family itself.

Let us now quickly discuss, why this prior is convenient. Assume that we have a set of iid samples  $\mathbf{x}_1, \dots, \mathbf{x}_N$ . As in Section 4.5 we assume that they are iid according to a Bernoulli distribution with an unknown parameter  $p$ . In addition we assume that the parameter itself is distributed according to  $q(p)$  with the parameters  $\alpha_1$  and  $\alpha_2$  fixed. Let us then write down the posterior distribution for the parameter  $p$  given the samples. We have

$$p(p \mid \mathbf{x}_1, \dots, \mathbf{x}_N) \propto p^{\alpha_1-1} (1-p)^{\alpha_2-1} \prod_{n=1}^N p^{\mathbf{x}_n} (1-p)^{1-\mathbf{x}_n} \quad (4.4)$$

$$\propto p^{\alpha_1-1+\sum_{n=1}^N \mathbf{x}_n} (1-p)^{\alpha_2-1+N-\sum_{n=1}^N \mathbf{x}_n}. \quad (4.5)$$

The key is to notice that this is again a beta distribution but now with parameters  $\sum_{n=1}^N \mathbf{x}_n + \alpha_1$  and  $N - \sum_{n=1}^N \mathbf{x}_n + \alpha_2$ . In particular, this is again an exponential distribution and so the optimization of this expression is again “simple.”

**Example 4.8** (Multinomial). The conjugate prior of a Multinomial is a Dirichlet distribution.

**Example 4.9** (Gaussian). The conjugate prior of a Gaussian is a Gaussian.

## 4.7 Maximum Entropy Distributions

Assume that we are given a function  $\phi : \mathcal{X} \rightarrow \mathbb{R}^d$  and a vector  $\alpha \in \mathbb{R}^d$ . What distribution on  $\mathcal{X}$  maximizes the entropy subject to condition that the expected value of  $\phi(X)$  is equal to  $\alpha$ ? Mathematically, we are looking for

$$P^* = \operatorname{argmax}_{P: \mathbb{E}_P[\phi(X)] = \alpha} H(P)$$

Why are we interested in this problem? If all that we know is a constraint on the mean it makes sense to look at the “most random” distribution that fulfills this constraint. This is the distribution that makes the least “assumptions” if we use it as a prior.

When looking at maximum entropy distributions we will drop the factor  $h(x)$  from exponential families. As we have mentioned earlier, any specific factor  $h(x)$  can be absorbed into the underlying measure  $\nu(x)$  and this is indeed the natural view point for our current purpose.

**Theorem 4.2.** For  $\theta \in \mathbb{R}^d$ , let  $P_\theta$  have density

$$p_\theta(x) = \exp\{\langle \theta, \phi(x) \rangle - A(\theta)\}$$

with respect to the measure  $\nu$ . If  $\mathbb{E}_{P_\theta}[\phi(X)] = \alpha$ , then  $P_\theta$  maximizes  $H(P)$  over  $\{P : \mathbb{E}_P[\phi(X)] = \alpha\}$  and it is the unique distribution with this property.

*Proof.* Let  $\theta$  be a parameter so that  $\mathbb{E}_{p_\theta}[\phi(X)] = \alpha$  and let  $P$  be any other distribution so that  $\mathbb{E}_P[\phi(X)] = \alpha$ . Then

$$\begin{aligned}
 H(P) &= - \int p(x) \log p(x) d\nu(x) \\
 &= - \int p(x) \log p_\theta(x) d\nu(x) + \int p(x) \log p_\theta(x) d\nu(x) - \int p(x) \log p(x) d\nu(x) \\
 &= - \int p(x) \log p_\theta(x) d\nu(x) - \int p(x) \log \frac{p(x)}{p_\theta(x)} d\nu(x) \\
 &= - \int p(x) [\langle \theta, \phi(x) \rangle - A(\theta)] d\nu(x) - D(p(x) \| p_\theta(x)) \\
 &= - \int p_\theta(x) [\langle \theta, \phi(x) \rangle - A(\theta)] d\nu(x) - D(p(x) \| p_\theta(x)) \\
 &= H(P_\theta) - \underbrace{D(p(x) \| p_\theta(x))}_{\geq 0} \\
 &\leq H(P_\theta).
 \end{aligned}$$

□

In all three of the following examples we pick  $\phi(x) = x^2$  and  $\alpha = 1$ , i.e., we are constraining the distribution  $P$  by asking that the second moment is equal to 1,  $\mathbb{E}_P[X^2] = 1$ . The general form of the density that maximizes the entropy is then

$$p_\theta(x) = \frac{\exp\{x^2\theta\}}{Z}, \quad (4.6)$$

where  $Z$  is the normalizing constant.

**Example 4.10.** If the measure  $\nu$  is the counting measure on  $\{-1, 1\}$  then the distribution  $P$  is of the form

$$(P(x = -1) = \frac{e^\theta}{Z}, P(x = 1) = \frac{e^\theta}{Z}),$$

where we have the condition  $\mathbb{E}_P[X^2] = 2\frac{e^\theta}{Z} = 1$ . Hence the maximum entropy distribution  $P(x)$  is  $(P(x = -1) = \frac{1}{2}, P(x = 1) = \frac{1}{2})$ , the uniform distribution.

**Example 4.11.** If the measure  $\nu$  is the counting measure on  $\mathbb{Z}$  then the distribution  $P$  is of the form

$$p_\theta(x) = \frac{e^{-\theta x^2}}{\sum_i e^{-\theta i^2}}, x \in \mathbb{Z},$$

where  $\theta$  is chosen so that

$$\sum_{x \in \mathbb{Z}} x^2 \frac{e^{-\theta x^2}}{\sum_i e^{-\theta i^2}} = 1.$$

**Example 4.12.** If the measure  $\nu$  is the Lebesgue measure on  $\mathbb{R}$  then we recognize from the basic form of the density given in (4.6) that the density that maximizes the entropy is the Gaussian distribution with mean 0 and variance 1,

$$p(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}.$$

In the proof of Theorem 4.2 we have seen that *if* an exponential distribution exists that yields the right moment then it is the maximum entropy distribution with this moment. Assume for a moment that we did not already *know* the form of this distribution. It is then perhaps insightful to “derive” the form of the distribution from first principles. Consider the following Lagrangian:

$$L = \int p(x) \log p(x) d\nu(x) + \theta^\top (\mu - \int p(x) \phi(x) d\nu(x)) + \kappa (1 - \int p(x) d\nu(x)).$$

Our aim is to minimize this Lagrangian. The first term is equal to  $-H(P)$ . Indeed, we want to maximize  $H(P)$ , i.e., equivalently we want to minimize  $-H(P)$ . The second term corresponds to all the constraints on the moments. Here,  $\theta$  is a vector of length  $d$ . And the third term corresponds to the normalization constraint on the density. Note that we have not included any term to ensure that the “density” is non-negative. We will see in a second that even without adding this constraint the solution will automatically fulfill this constraint, hence there is no need to add this constraint explicitly.

If we now take the “derivative” with respect to  $p(x)$  and set it to 0 we get

$$0 = 1 + \log(p(x)) - \langle \theta, \phi(x) \rangle - \kappa.$$

Solving for  $p(x)$ ,

$$p(x) = e^{\langle \theta, \phi(x) \rangle + \kappa - 1}.$$

This is of course an exponential distribution as expected. Note that due to the special structure of the solution  $p(x) \geq 0$  is automatically fulfilled.

## 4.8 Application To Physics

Let us re-derive one of the basic laws of physics – the Maxwell-Boltzmann distribution.

Assume that we have particles in  $\mathbb{R}^3$ . They each have a position and a velocity vector associated to them. We will not be interested in the position but we are asking how the velocity vectors are distributed. Let  $\mathbf{v} = (v_1, v_2, v_3)$  be the velocity vector associated to a particular particle.

We associate an average “kinetic energy”  $E$  (per particle) to the distribution

$$\int p(\mathbf{v}) \frac{1}{2} m (\mathbf{v}_1^2 + \mathbf{v}_2^2 + \mathbf{v}_3^2) d\mathbf{v} = E, \quad (4.7)$$

where  $m$  is the mass of a particle (all are assumed to have equal mass).

Let  $s = \sqrt{\mathbf{v}_1^2 + \mathbf{v}_2^2 + \mathbf{v}_3^2}$ , the speed. What is the maximum entropy distribution  $p(s)$ ? Note that in this case  $\phi(\mathbf{v}) = \mathbf{v}_1^2 + \mathbf{v}_2^2 + \mathbf{v}_3^2$ . Therefore, the form of the maximizing distribution is

$$p(\mathbf{v}) = e^{\theta(\mathbf{v}_1^2 + \mathbf{v}_2^2 + \mathbf{v}_3^2) - A(\theta)}.$$

We recognize this to be a three-dimensional zero-mean Gaussian distribution with independent and identically distributed components. We conclude that each component is distributed according to

$$p(v) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{v^2}{2\sigma^2}},$$

for some value of  $\sigma^2$ . Going back to our original constraint (4.7) we see that  $\sigma^2 = \frac{2E}{3m}$ . Summarizing, the velocity distribution of each component has the form

$$p(v) = \sqrt{\frac{3m}{4\pi E}} e^{-\frac{3mv^2}{4E}}.$$

What is the induced distribution of the overall speed  $s$ ? Recall that the surface of a sphere (in 3 D) of radius  $s$  has area  $4\pi s^2$ . Hence,

$$\mathbb{P}\{s \leq S \leq s + ds\} = \left(\frac{3m}{4\pi E}\right)^{\frac{3}{2}} e^{-\frac{3ms^2}{4E}} 4\pi s^2 ds,$$

so that

$$p(s) = \sqrt{\frac{27m^3}{4\pi E^3}} s^2 e^{-\frac{3ms^2}{4E}}. \quad (4.8)$$

Appealing to thermodynamics, we write  $E$  as  $E = 3/2kT$ , where  $k$  is the Boltzmann constant and  $T$  is the temperature. The factor 3 accounts for the three degrees of freedom and  $\frac{1}{2}kT$  is the kinetic energy per degree of freedom. Then we get the usual form

$$p(s) = \sqrt{\frac{2m^3}{\pi k^3 T^3}} s^2 e^{-\frac{ms^2}{2kT}}.$$

As an alternative derivation, we could have found  $p(s)$  by directly finding the maximum entropy distribution on  $\mathcal{X} = \mathbb{R}^+$  with measure  $d\nu(s) = 4\pi s^2 ds$ , for  $s \geq 0$ . In this case we know that

$$p(s) = e^{-\theta s^2 - A(\theta)} \mathbb{1}_{\{s \geq 0\}}.$$

The normalizing condition reads

$$\int_{s \geq 0} p(s) d\nu(s) = \int_{s \geq 0} e^{-\theta s^2 - A(\theta)} 4\pi s^2 ds = \frac{\pi^{3/2}}{\theta^{3/2}} e^{-A(\theta)} = 1.$$

This tells us that  $A(\theta) = \frac{3}{2} \ln \frac{\pi}{\theta}$ . The second moment requires that

$$\mathbb{E}_{p(s)}[s^2] = \int_{s \geq 0} \left(\frac{\theta}{\pi}\right)^{3/2} e^{-\theta s^2} s^2 d\nu(s) = \frac{3}{2\theta} = \frac{2E}{m},$$

so that  $\theta = \frac{3m}{4E}$ . It follows that the distribution that maximizes this entropy when written with respect to the Lebesgue measure on  $\mathbb{R}^+$  is equal to

$$p(s) = \left(\frac{\theta}{\pi}\right)^{3/2} 4\pi s^2 e^{-\theta s^2} \Big|_{\theta = \frac{3m}{4E}},$$

which is equal to what we got in (4.8).

## 4.9 I-Projections

In previous lectures we have discussed at length Sanov's theorem. Recall that if we have given a family of distributions, call the family  $\Pi$ , and a fixed distribution  $P$ , then the chance that we will mistake samples from  $P$  for samples from one of the elements of  $\Pi$ , is exponentially small and the exponent is asymptotically equal to  $\argmin_{Q \in \Pi} D(Q \| P)$ . The operation of finding the "closest" element of  $\Pi$  is called an I-projection. In general it is difficult to compute this projection. But if the family is linear then the projection is again easy to compute as we will see now.

**Theorem 4.3.** *Let  $P$  be a fixed distribution with density  $p(x)$  and let  $\Pi$  be the set of all distributions so that  $\mathbb{E}_q[\phi(x)] = \mu$  for  $q \in \Pi$ . If  $P_\theta$  has density*

$$p_\theta = p(x) e^{\langle \theta, \phi(x) \rangle - A(\theta)}$$

*and  $\mathbb{E}_{P_\theta}[\phi(x)] = \mu$  then*

$$P_\theta = \argmin_{Q \in \Pi} D(Q \| P).$$

*In words,  $P_\theta$  is the I-projection of  $P$  onto  $\Pi$ .*

*Proof.* The proof uses the same idea as we used to show that exponential distributions solve the maximum entropy problem. We have

$$\begin{aligned} D(Q \| P) &= \int q(x) \log \frac{q(x)}{p(x)} d\nu(x) \\ &= \int q(x) \log \frac{p_\theta(x)}{p(x)} d\nu(x) + \int q(x) \log \frac{q(x)}{p_\theta(x)} d\nu(x) \\ &= \int q(x) [\langle \theta, \phi(x) \rangle - A(\theta)] d\nu(x) + D(Q \| P_\theta) \\ &= \int p_\theta(x) [\langle \theta, \phi(x) \rangle - A(\theta)] d\nu(x) + D(Q \| P_\theta) \\ &= \int p_\theta(x) \log \frac{p_\theta(x)}{p(x)} d\nu(x) + D(Q \| P_\theta) \\ &= D(P_\theta \| P) + D(Q \| P_\theta) \\ &\geq D(P_\theta \| P). \end{aligned}$$

□

## 4.10 Relationship between $\theta$ and $\mathbb{E}[\phi(x)]$

### 4.10.1 The forward map $\nabla A(\theta)$

Assume that we fix  $h(x)$  and  $\phi(x)$ . Then for every  $\theta \in \Theta$  there is a distribution  $p_\theta(x)$  and an associated “mean”  $\mathbb{E}_{P_\theta}[\phi(x)]$ . This mapping  $\theta \mapsto \mu = \mathbb{E}_{P_\theta}[\phi(x)]$  is called the “forward” map. We have seen in (4.1) that it is given by  $\nabla_\theta A(\theta)$ .

Clearly this mapping is important. For simple distributions as for Bernoulli, Poisson, or Gaussian this map is simple to state and simple to compute. But there are important classes of distributions in the exponential family where this map is computationally difficult. Let us give one such example.

**Example 4.13** (Ising Model). The *Ising* model is a classical example from statistical physics which was initially introduced in order to study magnetism. The associated exponential distribution has the form

$$p_\theta(x) = e^{\sum_{s \in V} \theta_s x_s + \sum_{(s,t) \in E} \theta_{st} x_s x_t - A(\theta)}.$$

Here, the  $x_s$  take values in  $\{\pm 1\}$  and they are called *spins*. The set of spins is  $V$  and there is an underlying undirected graph with vertex set  $V$  and edge set  $E$ . The strength of the “interactions” between two spins that are connected by an edge  $(s, t)$  is  $\theta_{st}$ . There are also the “local fields”  $\theta_s$  for every spin  $s \in V$ .

Note that we are here in a much “higher dimensional” setting –  $\phi(x)$  has dimension  $|V| + |E|$ . Given the local fields and the strength of the interactions we are typically interested in the marginals and the pairwise correlations, i.e., we are exactly interested in  $\mu = \mathbb{E}_{P_\theta}[\phi(x)]$ . In particular we are interested if, e.g., some marginals become strongly “biased” or pairs become strongly correlated. If we stay with the physical interpretation of this model then such an “emergent” bias would represent the emergence of a global magnetic field given the local interaction rules. “Emergent” here means that we envision that we change the parameters of the model (the  $\theta_s$  and  $\theta_{st}$ ) and that for some such parameters even a small extra change might suddenly lead to biased marginals.

In summary, we are interested in the forward map  $\theta \mapsto \mu$ . But this map is in general exponentially complex to compute. E.g., if you look at the expression of  $A(\theta)$ , it has the form

$$A(\theta) = \ln \sum_{x \in \{0,1\}^{|V|}} e^{\sum_{s \in V} \theta_s x_s + \sum_{(s,t) \in E} \theta_{st} x_s x_t},$$

requiring a sum over an exponential number of terms. And the same is true if you look at the gradient of this expression since  $A(\theta)$  is part of this gradient computation. So even though the map is well-defined and mathematically simple to describe, it might be difficult to compute.



**Definition 4.2** (Set of Feasible Means). For a fixed  $\phi(x)$  let

$$\mathcal{M} = \{\mu \in \mathbb{R}^d : \exists P \text{ so that } \mathbb{E}_P[\phi(x)] = \mu\}.$$

In words, for a fixed sufficient statistics,  $\mathcal{M}$  is the set of all means that can be achieved by *some* distribution. It is important here that  $P$  is *not* assumed to be an element of the exponential family.

**Definition 4.3** (Regular Families). Let  $\phi(x)$  be given. We say that the associated exponential family is *regular* if  $\Theta$  is *open*.

**Definition 4.4** (Minimal Families). Let  $\phi(x)$  be given. We say that the associated exponential family is *minimal* if there *does not exist* a vector  $\eta$  so that

$$\eta^\top \phi(x) = \text{const},$$

$\nu(x)$ -almost everywhere.

**Theorem 4.4.** For a regular family the gradient  $\nabla_\theta A(\theta) : \Theta \rightarrow \mathcal{M}$  is one-to-one if and only if the exponential representation is minimal.

*Proof.* Assume at first that the family is not minimal, i.e., there does exist an  $\eta$  so that  $\eta^\top \phi(x) = c$ , a constant,  $\nu(x)$ -almost everywhere. Pick a  $\theta_1 \in \Theta$ . Then for a sufficiently small  $\epsilon$ ,  $\theta_2 = \theta_1 + \epsilon\eta \in \Theta$  since we assumed that  $\Theta$  was open (the family is regular). Note that  $A(\theta_2) = A(\theta_1) + \epsilon c$ . Therefore

$$\nabla_\theta A(\theta_1) = \nabla_\theta A(\theta_2).$$

Conversely, assume that the family is minimal. We claim that in this case  $A(\theta)$  is strictly convex. This implies that

$$\begin{aligned} A(\theta_2) &> A(\theta_1) + \langle \nabla_\theta A(\theta_1), \theta_2 - \theta_1 \rangle, \\ A(\theta_1) &> A(\theta_2) + \langle \nabla_\theta A(\theta_2), \theta_1 - \theta_2 \rangle. \end{aligned}$$

We therefore have

$$\langle \nabla_\theta A(\theta_1), \theta_1 - \theta_2 \rangle > A(\theta_1) - A(\theta_2) > \langle \nabla_\theta A(\theta_2), \theta_1 - \theta_2 \rangle.$$

This implies that

$$\langle \nabla_\theta A(\theta_1) - \nabla_\theta A(\theta_2), \theta_1 - \theta_2 \rangle > 0.$$

It remains to explain why  $A(\theta)$  is strictly convex for a minimal family. Recall that  $\nabla_\theta^2 A(\theta)$  is the covariance matrix of  $\phi(x)$ . So  $A(\theta)$  is always convex. If the family is minimal then for no  $\eta$  is  $\langle \eta, \phi(x) \rangle$  a constant. We conclude that the covariance of  $\langle \eta, \phi(x) \rangle$  is strictly positive. But this covariance is equal to  $\eta^\top \nabla_\theta^2 A(\theta) \eta$ , and so this quantity is strictly positive for any  $\eta \in \Theta$ .  $\square$

### 4.10.2 The backward map

When we discussed the maximum entropy problem we had to assume that for a given mean vector  $\mu$  there exists a parameter  $\theta$  so that  $\mathbb{E}_{P_\theta}[\phi(x)] = \mu$ . Only then could we conclude that the maximum entropy solution is an element of the exponential distribution. As we will see now, this is not really much of a restriction as long as there is *some* distribution that has this mean.

**Theorem 4.5.** *In a minimal exponential family, the gradient map  $\nabla_\theta A(\theta) : \Theta \rightarrow \mathcal{M}$  is onto the interior of  $\mathcal{M}$ .*

We will not provide a proof here but refer the reader to [2].

We can therefore define a backward map from the interior of  $\mathcal{M}$  onto  $\Theta$ . This has the pleasing consequence that if we are looking for a maximum entropy distribution then as long as we pick a mean vector from the interior of  $\mathcal{M}$  then the solution will be an element of the exponential family.

This has another important consequence. Let us go back to the parameter estimation problem discussed in Section 4.5. Assume that the samples do come from a minimal exponential family with sufficient statistic  $\phi(x)$  and that the parameter  $\theta_0$  is such that  $\nabla A(\theta_0) = \mu$  is in the interior of  $\mathcal{M}$ . Assume that we compute the empirical mean

$$\hat{\mu} = \frac{1}{N} \sum_{n=1}^N \phi(x).$$

We know that  $\hat{\mu} \xrightarrow{N \rightarrow \infty} \mu$  almost surely. Therefore, if we estimate the parameter by applying the inverse map to  $\hat{\mu}$  then this estimate will converge almost surely to the true parameter  $\theta_0$ . In other words, this estimator is *consistent*. This gives us an well-founded justification for using the ML estimator in the first place.

## 4.11 Problems

**Problem 4.1.** Find the maximum entropy density  $f$ , defined for  $x \geq 0$ , satisfying  $E[X] = \alpha_1$ ,  $E[\ln X] = \alpha_2$ . That is, maximize  $-\int f \ln f$  subject to  $\int x f(x) dx = \alpha_1$ ,  $\int (\ln x) f(x) dx = \alpha_2$ , where the integral is over  $0 \leq x < \infty$ . What family of densities is this?

**Problem 4.2.** What is the maximum entropy distribution  $p(x, y)$  that has the following marginals?

$\begin{array}{c} y \\ \diagdown \\ x \end{array}$	1	2	3	
1	$p_{11}$	$p_{12}$	$p_{13}$	$\frac{1}{2}$
2	$p_{21}$	$p_{22}$	$p_{23}$	$\frac{1}{4}$
3	$p_{31}$	$p_{32}$	$p_{33}$	$\frac{1}{4}$
	$\frac{2}{3}$	$\frac{1}{6}$	$\frac{1}{6}$	

**Problem 4.3.** (a) What is the parametric-form maximum entropy density  $f(x)$  satisfying the two conditions

$$E[X^8] = a \quad E[X^{16}] = b$$

(b) What is the maximum entropy density satisfying the condition

$$E[X^8 + X^{16}] = a + b$$

(c) Which entropy is higher?

**Problem 4.4.** Find the parametric form of the maximum entropy density  $f$  satisfying the Laplace transform condition

$$\int f(x)e^{-x}dx = \alpha,$$

and give the constraints on the parameter.

**Problem 4.5.** Let  $Y = X_1 + X_2$ . Find the maximum entropy of  $Y$  under the constraint  $E[X_1^2] = P_1$ ,  $E[X_2^2] = P_2$  :

- (a) If  $X_1$  and  $X_2$  are independent.
- (b) If  $X_1$  and  $X_2$  are allowed to be dependent.

**Problem 4.6.** We learned in the course that as long as the set of feasible means is open then every such mean can be realized by an element of the exponential family. In the following verify this explicitly (by not referring to the above statement for the following scenario).

- (i) Let  $\phi(x) = x^2$ .
- (ii) Let  $\phi(x)$  consist of all elements  $x_i x_j$ , where  $i$  and  $j$  go from 1 to  $K$ .

**Problem 4.7.** (*Exponential Families*) What is the maximum entropy distribution, call it  $p(x, i)$ , on  $[0, \infty] \times \mathbb{N}$ , both of whose marginals have mean  $\mu > 0$ . (I.e., in one axis the distribution is over the positive reals, whereas in the other one it is over the natural numbers.)



## Chapter 5

# Multi-Arm Bandits

This is a HUGE topic and we will only be able to cover basic material. If you want to know more (also about the historical development) we highly recommend <http://banditalgs.com>

### 5.1 Introduction

The basic model is the following. For each *round*  $t = 1, 2, \dots, n$ , the *learner* chooses an *action*  $A_t$  from a set of available actions  $\mathcal{A}$ . To each action  $A \in \mathcal{A}$  corresponds a probability distribution  $P_A$ . The *environment* receives the chosen action  $A_t$  from the learner and in response generates the random variable  $X_t$  that is distributed according to  $P_{A_t}$ .

The *reward* up to and including time  $n$  is  $\sum_{t=1}^n X_t$ . The decision by the learner at time  $t$  is in general a function of the history  $H_{t-1} = \{A_1, X_1, \dots, A_{t-1}, X_{t-1}\}$ . Our aim is to maximize the reward by employing an appropriate learning algorithm.

More precisely, we typically try to *minimize* the *regret* rather than maximize the reward. The regret with respect to a particular action  $a \in \mathcal{A}$  is the difference of what we could have gotten if we had used action  $a$  in all  $n$  rounds versus the actual reward we got. The advantage of this *competitive* view (comparing to some other action) is that this measure is invariant to e.g. shifting all rewards by a constant amount. We typically compute the *worst case* regret, i.e., the regret with respect to the best action we could have taken, and since the reward is a random variable it is common to first average over the reward for each action. This means we compute

$$R_n = \max_{A \in \mathcal{A}} n\mu_A - \mathbb{E}\left[\sum_{t=1}^n X_t\right].$$

It is probably not surprising that a good learner will be able to achieve a sublinear worst-case regret, i.e.,  $R_n = o(n)$ . Let us quickly go over the argument. We will do a much more thorough analysis later on. Assume that  $|\mathcal{A}| = K$ . If we take  $m$  samples from any of these  $K$  distributions we can compute each mean with an additive error bounded by  $c/\sqrt{m}$  with high probability.

Assume that we spend a fraction  $\epsilon$  of the total time  $n$  on learning the  $K$  actions and afterwards always play the “best” one according to the derived estimates. In this way we

will achieve a regret that behaves like  $n\mu^*(\epsilon + cK/\sqrt{\epsilon n})$ . The term  $n\mu^*\epsilon$  is an upper bound on the regret that we get since for a fraction  $\epsilon$  of the time (when we are learning) we might have a regret as large as  $\mu^*$ . The second term, namely  $n\mu^*cK/\sqrt{\epsilon n}$  accounts for the fact that during the remaining fraction  $1 - \epsilon \leq 1$  of the time, we always play the “best” arm according to our estimates but for each arm the estimate can be off by  $c/\sqrt{\epsilon n}$  with a fixed probability and so each arm in expectation will contribute a term of this order to the expected regret and we have  $K$  arms. We are still free to optimize over the choice of  $\epsilon$ . If we choose  $\epsilon = (\frac{cK}{2\sqrt{n}})^{\frac{2}{3}}$  then we get  $3(n\frac{cK}{2})^{\frac{2}{3}}$ , which vanishes as a function of  $n$ . So the interesting question is how *fast* we can make the normalized regret converge to 0.

In the above paragraph we have assume that we know the *time horizon*  $n$ . This is often the case. But also the setting where the time horizon is not known a priori is of interest.

The setting we described, where the rewards come from a distribution that only depends on the chosen action and this distribution is fixed over time is called the *stochastic stationary bandit* problem. We will limit ourselves to this setting.

## 5.2 Some References

Besides the web page pointed out at the beginning there are many very good references for this topic.

The paper that started the whole area is William R. Thompson. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, 1933.

You can find a good survey at <https://arxiv.org/pdf/1204.5721v2.pdf>. If you are looking for a book, we recommend Multi-armed Bandit Allocation Indices, 2nd Edition, by John Gittins, Kevin Glazebrook, Richard Weber, or Bandit problems – Sequential Allocation of Experiments, by Berry, Donald A. and Fristedt, Bert.

## 5.3 Stochastic Bandits with a Finite Number of Arms

### 5.3.1 Set-Up

Let us analyze the simplest strategy that we already mentioned in a little bit more detail. It is somewhat easier to think of problems with infinite horizons, i.e., there is no fixed  $n$ , but we assume that the game goes on forever.

### 5.3.2 Explore then Exploit

A sub-optimal but very natural strategy is the following. First, get sufficiently many samples from every bandit in order to determine its mean sufficiently accurately. This is the exploring stage. Then exploit the so gained knowledge and play according to these empirical means. If we have a bandit with a finite number of arms then it is not very surprising that this strategy achieves a sub-linear regret.

Let us do the calculations. Let  $X_1, \dots, X_n$  be a sequence of iid random variables with mean  $\mu = \mathbb{E}[X_i]$ . Given the sequence  $X_1, \dots, X_m$  the empirical estimator for  $\mu$ , call it

$\hat{\mu}(X_1, \dots, X_m)$  is

$$\hat{\mu}(X_1, \dots, X_m) = \frac{1}{m} \sum_{t=1}^m X_t.$$

The above estimate is itself a random variable. It's mean is unbiased, i.e.,

$$\mathbb{E}[\hat{\mu}(X_1, \dots, X_m)] = \frac{1}{m} \mathbb{E}\left[\sum_{t=1}^m X_t\right] = \mu.$$

But of course we have a variance. Let us therefore give a bound on the probability that this estimator deviates significantly from the true value. The natural avenue is to use Chernov bounds. Indeed, this is what we do. But rather than optimizing over the degree of freedom that we have in this bound, we will make an assumption on the underlying random variable that will make this optimization unnecessary.

**Definition 5.1.** A random variable  $X$  is  $\sigma^2$ -subgaussian if for all  $\lambda \in \mathbb{R}$  it holds that  $\mathbb{E}[e^{\lambda X}] \leq e^{\lambda^2 \sigma^2 / 2}$ .

The quantity  $\mathbb{E}[e^{\lambda X}]$  is called the *moment-generating function*.

**Lemma 5.1** (Basic Properties of subgaussians). *Let  $X_i$ ,  $i = 1, 2$ , be two  $\sigma_i^2$ -subgaussian independent random variables. Then*

- (i)  $\mathbb{E}[X_i] = 0$ ;  $\mathbb{E}[X_i^2] \leq \sigma_i^2$ .
- (ii) For all  $\alpha \in \mathbb{R}$ ,  $\alpha X_i$  is  $(\alpha^2 \sigma_i^2)$ -subgaussian.
- (iii)  $X_1 + X_2$  is  $(\sigma_1^2 + \sigma_2^2)$ -subgaussian.

*Proof.* Pick  $\lambda$  so that  $\lambda \mathbb{E}[X_i] \geq 0$ . Then using our assumption in the second step,

$$1 + \frac{1}{2} \lambda^2 \sigma_i^2 + O(\lambda^4) \geq \mathbb{E}[e^{\lambda^2 \sigma_i^2 / 2}] \geq \mathbb{E}[e^{\lambda X_i}] \geq 1 + \lambda \mathbb{E}[X_i] + \frac{1}{2} \lambda^2 \mathbb{E}[X_i^2] + O(\lambda^3).$$

Claim (i) follows by letting  $\lambda$  tend to 0. Claim (ii) is true since  $\mathbb{E}[e^{\lambda X_i}] \leq e^{\lambda^2 \sigma_i^2 / 2}$  implies  $\mathbb{E}[e^{\lambda(\alpha X_i)}] = \mathbb{E}[e^{(\lambda \alpha) X_i}] \leq \mathbb{E}[e^{\lambda^2 \alpha^2 \sigma_i^2 / 2}] = e^{\lambda^2 (\alpha \sigma_i)^2 / 2}$ . And to prove claim (iii), note that  $\mathbb{E}[e^{\lambda(X_1 + X_2)}] = \mathbb{E}[e^{\lambda X_1} e^{\lambda X_2}] = \mathbb{E}[e^{\lambda X_1}] \mathbb{E}[e^{\lambda X_2}] \leq e^{\lambda^2 \sigma_1^2 / 2} e^{\lambda^2 \sigma_2^2 / 2} = e^{\lambda^2 (\sigma_1^2 + \sigma_2^2) / 2}$ .  $\square$

**Lemma 5.2** (Zero-Mean Gaussian is subgaussian). *Let  $X$  be a Gaussian random variable with mean zero and variance  $\sigma^2$ . Then  $X$  is  $\sigma^2$ -subgaussian.*

*Proof.* By Lemma 5.1 we can assume that  $\sigma^2 = 1$ . We then have

$$\begin{aligned} \mathbb{E}[e^{\lambda X}] &= \frac{1}{\sqrt{2\pi}} \int e^{\lambda x} e^{-x^2/2} dx \\ &= \frac{1}{\sqrt{2\pi}} \int e^{\lambda x - x^2/2} dx \\ &= e^{\lambda^2/2} \frac{1}{\sqrt{2\pi}} \int e^{-\frac{1}{2}(\lambda - x)^2} dx \\ &= e^{\lambda^2/2}. \end{aligned}$$

$\square$

**Lemma 5.3** (Zero-Mean Bernoulli is subgaussian). *Let  $X$  be a zero-mean Bernoulli random variable,  $\mathbb{P}\{X = 1-p\} = p$  and  $\mathbb{P}\{X = -p\} = 1-p$ ,  $0 \leq p \leq 1$ . Then  $X$  is  $1/4$ -subgaussian.*

*Proof.* We have

$$\begin{aligned}\mathbb{E}[e^{\lambda X}] &= (1-p)e^{-\lambda p} + pe^{\lambda(1-p)} \\ &\leq e^{\frac{1}{4}\lambda^2/2}.\end{aligned}$$

The last step is a special form of the so-called *Hoeffding* inequality: Note that we only have to proof the inequality for  $\lambda \geq 0$  and  $p \in [0, 1]$ . This is true since the pair  $(\lambda, p)$  results in the same value as the pair  $(-\lambda, 1-p)$  and the bound only depends on  $\lambda^2$ . Consider

$$\begin{aligned}\phi(\lambda) &= \ln[(1-p)e^{-\lambda p} + pe^{\lambda(1-p)}] \\ &= -\lambda p + \ln[(1-p) + pe^{\lambda}].\end{aligned}$$

Write it as  $\phi(\lambda) = \phi(0) + \phi'(0)\lambda + \frac{1}{2}\phi''(\xi)\lambda^2$ , where  $0 \leq \xi \leq \lambda$ . Note that  $\phi(0) = \ln(1) = 0$ . Next,

$$\phi'(\lambda) = -p + \frac{pe^{\lambda}}{(1-p) + pe^{\lambda}},$$

so that  $\phi'(0) = -p + p = 0$ . Finally,

$$\phi''(\lambda) = \frac{e^{\lambda}p(1-p)}{((1-p) + pe^{\lambda})^2}.$$

Now note that by the convexity of the function  $x \mapsto x^2$ ,  $((1-p) + pe^{\lambda})^2 \geq 4e^{\lambda}p(1-p)$ . From this it follows that  $\phi''(\xi) \leq \frac{1}{4}$  for all  $0 \leq \xi \leq \lambda$ . We therefore have  $\phi(\lambda) \leq \frac{1}{8}\lambda^2$ , which is equivalent to the claim.  $\square$

**Lemma 5.4** (Zero-Mean RV with finite range is subgaussian). *Let  $X$  be a zero-mean random variable with  $X \in [a, b]$ . Then  $X$  is  $(b-a)^2/4$ -subgaussian.*

*Proof.* The proof follows along the same lines of the previous proof and we skip the details.  $\square$

**Lemma 5.5** (Hoeffding's Bound). *Assume that  $X_1 - \mu, \dots, X_m - \mu$  are zero-mean independent  $\sigma^2$ -subgaussian random variables. Then the empirical estimator  $\hat{\mu}$  satisfies*

$$\begin{aligned}\mathbb{P}\{\hat{\mu} \geq \mu + \epsilon\} &\leq \exp\left\{-\frac{m\epsilon^2}{2\sigma^2}\right\}, \\ \mathbb{P}\{\hat{\mu} \leq \mu - \epsilon\} &\leq \exp\left\{-\frac{m\epsilon^2}{2\sigma^2}\right\}.\end{aligned}$$



*Proof.* Assume at first that  $X$  is  $\sigma^2$ -subgaussian. Then

$$\begin{aligned} \mathbb{P}\{X \geq \epsilon\} &\stackrel{\lambda \geq 0}{\leq} \mathbb{P}\{e^{\lambda X} \geq e^{\lambda \epsilon}\} \\ &\stackrel{\text{Markov inequality}}{\leq} \frac{\mathbb{E}[e^{\lambda X}]}{e^{\lambda \epsilon}} \\ &\stackrel{\sigma^2\text{-subgaussian}}{\leq} \mathbb{E}[e^{\frac{1}{2}\lambda^2 \sigma^2 - \lambda \epsilon}]. \end{aligned}$$

Therefore,

$$\begin{aligned} \mathbb{P}\left\{\frac{1}{m} \sum_{t=1}^m X_t \geq \mu + \epsilon\right\} &\stackrel{\lambda \geq 0}{\leq} \mathbb{P}\left\{e^{\frac{\lambda}{m} \sum_{t=1}^m (X_t - \mu)} \geq e^{\lambda \epsilon}\right\} \\ &\stackrel{\text{Markov inequality}}{\leq} \frac{\mathbb{E}[e^{\frac{\lambda}{m} \sum_{t=1}^m (X_t - \mu)}]}{e^{\lambda \epsilon}} \\ &\stackrel{\sigma^2\text{-subgaussian}}{\leq} \mathbb{E}[e^{\frac{1}{2m}\lambda^2 \sigma^2 - \lambda \epsilon}] \\ &\stackrel{\lambda = m\epsilon/\sigma^2}{=} \mathbb{E}[e^{-m\epsilon^2/(2\sigma^2)}]. \end{aligned}$$

In the one before-last step we have used property (ii) and (iii) of Lemma 5.1 to conclude that  $\frac{1}{m} \sum_{t=1}^m (X_t - \mu)$  is  $\frac{\sigma^2}{m}$ -subgaussian.  $\square$

Let us now get back to analyzing the explore-then-exploit (or commit) strategy. Assume that at the start we get  $m$  samples from each of the  $K$  bandit arms. Let the expected gain from arm  $k$  be  $\mu_k$  and let  $\mu^* = \max_{1 \leq k \leq K} \mu_k$ . To simplify notation, let us assume that  $\mu^* = \mu_1$ . Define  $\Delta_k = \mu^* - \mu_k$ . Finally, assume that each of the  $K$  arms corresponds to a random variable that is 1-subgaussian.

After the initial exploration stage we choose that bandit with the largest empirical pay-off for the remaining  $n - Km$  steps. This gives us an expected regret of

$$R_n = m \sum_{k=1}^K \Delta_k + (n - mK) \sum_{k=1}^K \Delta_k \mathbb{P}\{k = \operatorname{argmax}_j \hat{\mu}_j\}.$$

This expression is easy to explain. The first sum on the right is the expected regret due to the exploration stage – we get  $m$  samples from each arm, and in so doing, accumulate for each arm an expected regret of  $m\Delta_k$ .

The second sum on the right accounts for the regret that we accumulate over the remaining  $n - mK$  steps in case we choose a sub-optimal arm in the exploitation stage. This second term we can now bound using our tail-bound inequalities. Recall that  $\Delta_k$  is the regret if we use arm  $k$ , instead of the optimum arm 1. We get  $m$  samples. What is the probability that the average of the  $m$  samples of arm  $k$  look better than the average of the  $m$  samples of arm 1? This is equivalent to asking for the probability that

$$\mathbb{P}\left\{\frac{1}{m} \sum_{t=1}^m (X_t^{(1)} - X_t^{(k)}) \leq 0\right\},$$

where  $X_t^{(1)}$  denotes the  $m$  independent samples from arm 1 and  $X_t^{(k)}$  denotes the  $m$  independent samples from arm  $k$ . Now note that by assumption  $X_t^{(1)} - \mu_1$  is 1-subgaussian and so is  $X_t^{(k)} - \mu_k$ . Therefore  $X_t^{(1)} - \mu_1 + X_t^{(k)} - \mu_k$  is 2-subgaussian by Lemma 5.1. We therefore have

$$\begin{aligned} \mathbb{P}\left\{\frac{1}{m} \sum_{t=1}^m (X_t^{(1)} - X_t^{(k)}) \leq 0\right\} &= \mathbb{P}\left\{\frac{1}{m} \sum_{t=1}^m (X_t^{(1)} - \mu_1 - X_t^{(k)} + \mu_k) \leq \mu_k - \mu_1\right\} \\ &= \mathbb{P}\left\{\frac{1}{m} \sum_{t=1}^m (X_t^{(1)} - \mu_1 - X_t^{(k)} + \mu_k) \leq -\Delta_k\right\} \\ &\leq e^{-m\Delta_k^2/4}. \end{aligned}$$

Therefore, our regret can be upper bounded as

$$\begin{aligned} R_n &= m \sum_{k=1}^K \Delta_k + (n - mK) \sum_{k=1}^K \Delta_k \mathbb{P}\{k = \operatorname{argmax}_j \hat{\mu}_j\} \\ &\leq m \sum_{k=1}^K \Delta_k + (n - mK) \sum_{k=1}^K \Delta_k \exp\left\{-\frac{m\Delta_k^2}{4}\right\}. \end{aligned}$$

It is instructive to consider the special case of  $K = 2$ . Let  $\Delta$  be the regret of the second best arm (compared to the best one). Our expression for the regret is then

$$R_n = m\Delta + (n - 2m)\Delta \exp\left\{-\frac{m\Delta^2}{4}\right\} \leq m\Delta + n\Delta \exp\left\{-\frac{m\Delta^2}{4}\right\}.$$

If we assume that we know  $n$  and  $\Delta$  a priori we can find the optimal value of  $m$ . This leads to the equation

$$\begin{aligned} \frac{dR_n}{dm} &= \Delta(1 - ne^{-\frac{m\Delta^2}{4}} \Delta^2/4) = 0, \\ e^{-\frac{m\Delta^2}{4}} \frac{\Delta^2}{4} &= \frac{1}{n}. \end{aligned}$$

We see from this equation that the optimum choice (ignoring integer constraints) is

$$m \sim \frac{4}{\Delta^2} \ln\left(\frac{n\Delta^2}{4}\right).$$

This gives us

$$R_n \sim \frac{4}{\Delta} \left(1 + \ln\left(\frac{n\Delta^2}{4}\right)\right).$$

At first this looks pretty promising. The bound on the right is only logarithmic in  $n$ . But there is a slight problem with our bound. So far we have implicitly assumed that  $\Delta$  is relatively large. But what if  $\Delta$  is small? E.g., assume that  $\Delta = \frac{1}{\sqrt{n}}$ . Then the term  $\frac{1}{\Delta}$  is

equal to  $\sqrt{n}$  and so our regret is now much larger. Indeed, what if  $\Delta$  is even smaller? It seems that the regret has no bound – the smaller the gap the larger the regret. This seems counter-intuitive. Should a small gap not be good for us?

We can easily fix this bound by noting that the regret can never be larger than  $n\Delta$ . Hence, we have a bound of the form

$$R_n \leq \min\left\{n\Delta, \frac{4}{\Delta} \left(1 + \ln\left(\frac{n\Delta^2}{4}\right)\right)\right\}.$$

Now it is easy to see that the worst case is to have a gap of order  $1/\sqrt{n}$ . In this case the regret is of order  $\sqrt{n}$ .

To summarize. If the gap is large (a constant) then we only need  $\ln(n)$  samples to figure out which of the arms is best with high probability. After that we will always use the best arm. This gives us a regret of order  $\ln(n)$ . But if the gap is small then even though pulling the wrong trigger is less costly we need a considerably larger exploration phase. And once the gap becomes of size  $1/\sqrt{n}$ , we will spend all the time in exploring and never reach the exploitation phase.

There is another issue. All our previous discussion is based on the assumption that we *know* the horizon  $n$  and the gap  $\Delta$  (just look at the expression for the optimum  $m$  – it depends both on  $n$  and  $\Delta$ ). Perhaps it is realistic to assume that  $n$  is known. But it is not realistic to assume that we know  $\Delta$ . So is there a way to choose  $m$  that is *universal*? We will explore this in the exercises. We will see that there is. But in this case the worst-case regret is of order  $n^{\frac{2}{3}}$ . Note that this is the same order that we got in our very first back of the envelope calculation.

### 5.3.3 The Upper Confidence Bound Algorithm

The upper confidence bound (UCB) algorithm is a celebrated algorithm that overcomes the shortcomings of the explore-then-exploit algorithm. Rather than separating the exploring phase from the exploiting phase these two phases are mixed and the algorithm learns continuously. The idea is simple: At any point the algorithm gets a sample from that arm that, according to optimistic estimates, looks best.

Recall our upper bound of

$$\mathbb{P}\{\hat{\mu}(X_1, \dots, X_m) - \mu \geq \epsilon\} \leq \exp(-m\epsilon^2/2).$$

If we set the right-hand side to  $\delta > 0$  and then solve for  $\delta$  we get

$$\mathbb{P}\{\hat{\mu}(X_1, \dots, X_m) - \mu \geq \sqrt{\frac{2}{m} \ln\left(\frac{1}{\delta}\right)}\} \leq \delta.$$

If we think of  $\delta$  as small then this suggest that, at time  $t-1$ , it is unlikely that our empirical estimator  $\hat{\mu}_{k,t-1}$  of the  $k$ -th bandit arm overestimates its mean by more than  $\sqrt{\frac{2}{T_k(t-1)} \ln\left(\frac{1}{\delta}\right)}$ . Here  $T_k(t-1)$  denotes the number of times we have chosen arm  $k$  in the first  $t-1$  steps.

The idea of the UCB algorithm is to take these upper bounds on the individual confidence intervals as our estimates and to choose as an action  $A_t$  at time  $t$  that arm  $i$  that maximizes this upper bound.

To specify the algorithm it remains to specify the *confidence* level  $\delta_t$  that is used at time  $t$ . We will choose

$$\delta_t = \frac{1}{f(t)} = \frac{1}{1 + t \ln^2(t)}. \quad (5.1)$$

Note that the above algorithm has the following property. Once all arms have been explored at depth all the upper bounds on the confidence intervals will be very close to the true means and so we will likely explore further only arms whose mean is very close to the maximum mean.

Let us now formally specify the algorithm. We have

$$A_t = \begin{cases} t, & t \leq K, \\ \operatorname{argmax}_k \hat{\mu}_k(t-1) + \sqrt{\frac{2 \ln f(t)}{T_k(t-1)}}, & t > K. \end{cases}$$

This algorithm is pretty intuitive. Even if a genie had given us the correct mean of the best arm for free, in order to verify that indeed this is the best arm, what we would do is to compute its confidence interval. And how confident should we be, how should we choose  $\delta$ ? If we make a mistake we will pay linear regret for the remainder of the running time. Therefore  $\delta$  should be smaller than  $\frac{1}{n}$ . If we think now of  $n$  as  $t$  then (5.1) makes sense.

**Lemma 5.6.** *The regret of the UCB algorithm is bounded by*

$$R_n \leq \sum_{k: \Delta_k > 0} \inf_{\epsilon \in (0, \Delta_k)} \Delta_k \left( 1 + \frac{7}{\epsilon^2} + \frac{2}{(\Delta_k - \epsilon)^2} (\ln f(n) + \sqrt{\pi \ln f(n)} + 1) \right).$$

Let us compare this result to what we have seen for the explore-then-exploit algorithm. If we pick  $\epsilon$  small but not too small then the dominant terms in this expression are of the form  $\frac{2 \ln f(n)}{\Delta_k} \sim \frac{2 \ln(n)}{\Delta_k}$ . This is essentially the same as what we derived for the explore-then-exploit algorithm. But this time we neither required the knowledge of  $n$  nor of  $\Delta_k$ . Of course, we have the same issue when one of the  $\Delta_k$  becomes small. The worst case is again when one of these gaps is of order  $1/\sqrt{n}$ . This will, as before, result in a regret of order  $\sqrt{n} \ln(n)$ . In summary, we have

$$R_n \leq \sum_{k: \Delta_k > 0} \min \left\{ n \Delta_k, \inf_{\epsilon \in (0, \Delta_k)} \Delta_k \left( 1 + \frac{7}{\epsilon^2} + \frac{2}{(\Delta_k - \epsilon)^2} (\ln f(n) + \sqrt{\pi \ln f(n)} + 1) \right) \right\}.$$

*Proof.* Let  $\hat{\mu}_t$  be the empirical (natural) estimator of the mean of a 1-subgaussian random variable based on  $t$  independent observations. Let  $a \in \mathbb{R}^+$  and  $\epsilon > 0$ . Consider the quantity

$$\mathbb{P} \left\{ \hat{\mu}_t + \sqrt{\frac{2a}{t}} \geq \epsilon \right\}.$$

For  $t$  no more than  $\frac{2a}{\epsilon^2}$  this probability is very close to 1. But for larger  $t$  we can use our tail bound to conclude that

$$\mathbb{P}\{\hat{\mu}_t + \sqrt{\frac{2a}{t}} \geq \epsilon\} \leq e^{-\frac{1}{2}t(\epsilon - \sqrt{\frac{2a}{t}})^2}.$$

Therefore,

$$\begin{aligned} \mathbb{E}\left[\sum_{t=1}^n \mathbb{I}_{\{\hat{\mu}_t + \sqrt{\frac{2a}{t}} \geq \epsilon\}}\right] &= \sum_{t=1}^n \mathbb{P}\{\hat{\mu}_t + \sqrt{\frac{2a}{t}} \geq \epsilon\} \\ &\leq \frac{2a}{\epsilon^2} + \sum_{t \geq \frac{2a}{\epsilon^2}}^n \mathbb{P}\{\hat{\mu}_t + \sqrt{\frac{2a}{t}} \geq \epsilon\} \\ &\leq \frac{2a}{\epsilon^2} + \sum_{t \geq \frac{2a}{\epsilon^2}}^n e^{-\frac{1}{2}t(\epsilon - \sqrt{\frac{2a}{t}})^2} \\ &\stackrel{(a)}{\leq} \frac{2a}{\epsilon^2} + 1 + \int_{\frac{2a}{\epsilon^2}}^{\infty} e^{-\frac{1}{2}t(\epsilon - \sqrt{\frac{2a}{t}})^2} dt \\ &\stackrel{(b)}{=} \frac{2a}{\epsilon^2} + 1 + \frac{2}{\epsilon^2} \int_0^{\infty} e^{-\frac{1}{2}x^2} (x + \sqrt{2a}) dx \\ &= 1 + \frac{2}{\epsilon^2} (a + \sqrt{\pi a} + 1). \end{aligned} \tag{5.2}$$

In step (a) we note that the terms in the sum are decreasing and that each term is less than 1. We can hence bound the sum by the corresponding integral plus the constant 1 (the maximum value that the function can take on at the left boundary). In step (b) we made two substitutions. First we set  $z = \epsilon\sqrt{t}$  so that  $dt = 2z/\epsilon^2 dz$ . This will change the lower bound to  $\sqrt{2a}$  and the argument in the exponent to  $-\frac{1}{2}(z - \sqrt{2a})^2$ . Then we shift the integration boundaries by defining  $x = z - \sqrt{2a}$ . This gives us the indicated integral.

Let us now bound the regret, which has the form  $R_n = \sum_{k: \Delta_k > 0} \Delta_k \mathbb{E}[T_k(n)]$ . The key is to find a good bound on  $\mathbb{E}[T_k(n)]$ . Note that

$$T_k(n) = \sum_{t=1}^n \mathbb{I}_{\{A_t=k\}} \leq \sum_{t=1}^n \mathbb{I}_{\{\hat{\mu}_1(t-1) + \sqrt{\frac{2 \ln f(t)}{T_1(t-1)}} \leq \mu_1 - \epsilon\}} + \sum_{t=1}^n \mathbb{I}_{\{\hat{\mu}_k(t-1) + \sqrt{\frac{2 \ln f(t)}{T_k(t-1)}} \geq \mu_1 - \epsilon \wedge A_t=k\}} \tag{5.3}$$

The idea of this bound is the following. Rather than counting how often the upper confidence bound of arm  $k$  is larger than the upper confidence bounds of all other arms, we count how often it is larger than the upper confidence bound of arm 1.

Clearly, this count is an upper bound. Further, rather than comparing the upper confidence bound of arm  $k$  and arm 1 directly, we compare each individually to a third quantity. This quantity is chosen to be slightly below the true mean of arm 1. We increase our count if either the upper confidence bound of arm  $k$  is above this threshold, or if the upper confidence bound of arm 1 is below this threshold. Again, this leads to an upper bound. This explains the two terms on the right in (5.3).

We start with the first one,

$$\begin{aligned}
\mathbb{E}\left[\sum_{t=1}^n \mathbb{I}_{\{\hat{\mu}_1(t-1) + \sqrt{\frac{2 \ln f(t)}{T_1(t-1)}} \leq \mu_1 - \epsilon\}}\right] &= \sum_{t=1}^n \mathbb{P}\left(\hat{\mu}_1(t-1) + \sqrt{\frac{2 \ln f(t)}{T_1(t-1)}} \leq \mu_1 - \epsilon\right) \\
&\stackrel{(a)}{\leq} \sum_{t=1}^n \sum_{s=1}^t \mathbb{P}\left(\hat{\mu}_{1,s} + \sqrt{\frac{2 \ln f(t)}{s}} \leq \mu_1 - \epsilon\right) \\
&\leq \sum_{t=1}^n \sum_{s=1}^t e^{-\frac{s}{2}(\sqrt{\frac{2 \ln f(t)}{s}} + \epsilon)^2} \\
&= \sum_{t=1}^n \sum_{s=1}^t e^{-\ln(f(t)) - \sqrt{2s \ln f(t)} - \frac{s}{2}\epsilon^2} \\
&\leq \sum_{t=1}^n \frac{1}{f(t)} \sum_{s=1}^t e^{-\frac{s}{2}\epsilon^2} \\
&= \sum_{t=1}^n \frac{1}{f(t)} \frac{e^{-\frac{\epsilon^2}{2}}}{1 - e^{-\frac{\epsilon^2}{2}}} \\
&= \sum_{t=1}^n \frac{1}{f(t)} \underbrace{\frac{1}{e^{\frac{\epsilon^2}{2}} - 1}}_{\substack{\text{take Taylor series} \\ \text{all terms are positive; keep only first two}}} \\
&\leq \sum_{t=1}^n \frac{1}{f(t)} \frac{2}{\epsilon^2} \\
&\leq \frac{2}{\epsilon^2} \sum_{t=1}^n \frac{1}{1 + t \ln(t)^2} \\
&\stackrel{(b)}{\leq} \frac{2}{\epsilon^2} (2 + \int_2^\infty \frac{1}{x \ln(x)^2} dx) \\
&= \frac{2}{\epsilon^2} (2 + \frac{1}{\ln(2)}) \\
&\leq \frac{7}{\epsilon^2}
\end{aligned}$$

□

In step (a) we argue as follows. We do not know how many samples of arm 1 we have taken at time  $t$ . Hence we bound this probability via a union bound, where we sum over all possibilities. In step (b) we used the fact that  $\frac{1}{1+t \ln(t)^2}$  is a decreasing function so that the sum can be bounded by an appropriately chosen integral. In particular, we note that each term is upper bounded by 1. We hence bound the first two terms by 2 and then we bound the remainder of the sum by the corresponding integral starting at 2 (the sum starts at 3). Finally, we drop the 1 from the denominator and we extend the integral to infinity. This further upper bounds the sum and leads to a simple expression.

It remains to bound the second term in (5.3),

$$\begin{aligned}
\mathbb{E}\left[\sum_{t=1}^n \mathbb{I}_{\{\hat{\mu}_k(t-1) + \sqrt{\frac{2 \ln f(t)}{T_k(t-1)}} \geq \mu_1 - \epsilon \wedge A_t = k\}}\right] &\stackrel{(a)}{\leq} \mathbb{E}\left[\sum_{t=1}^n \mathbb{I}_{\{\hat{\mu}_k(t-1) + \sqrt{\frac{2 \ln f(n)}{T_k(t-1)}} \geq \mu_1 - \epsilon \wedge A_t = k\}}\right] \\
&\stackrel{(b)}{\leq} \mathbb{E}\left[\sum_{s=1}^n \mathbb{I}_{\{\hat{\mu}_{k,s} + \sqrt{\frac{2 \ln f(n)}{s}} \geq \mu_1 - \epsilon\}}\right] \\
&\leq \mathbb{E}\left[\sum_{s=1}^n \mathbb{I}_{\{\hat{\mu}_{k,s} - \mu_k + \sqrt{\frac{2 \ln f(n)}{s}} \geq \Delta_k - \epsilon\}}\right] \\
&\stackrel{(c)}{\leq} 1 + \frac{2}{(\Delta_k - \epsilon)^2} (\ln f(n)) + \sqrt{\pi \ln f(n)} + 1.
\end{aligned}$$

In step (a) we replaced  $f(t)$  by the larger quantity  $f(n)$ . This gives us an upper bound.

Step (b) also warrants some explanation. As in the previous case, we do not know what  $T_k(t-1)$  is, other than that it must be in the range from 1 to  $t-1$ . When we derived a bound on the first term in (5.3) we got around this problem by taking a union bound over all possible such values.

We could do the same thing here, but this bound would be loose. The trick is to realize that whatever value of  $s$  we have for a particular step  $t$ , the same value cannot appear again in a later step due to the condition that  $A_t = k$ . Therefore, it suffices to take the sum over all possible values of  $s$  *once*, i.e., *instead* of the sum over  $t$ .

Finally, in step (c) we have used (5.2) with  $a = \ln(f(n))$  and  $\epsilon$  replaced by  $\Delta_k = \epsilon$ .

### 5.3.4 Information-theoretic Lower Bound

We have seen that the UCB algorithm has a worst-case (worst-case over the choice of gaps) regret of order  $\sqrt{n} \ln(n)$ . Could there be an algorithm that is much better than that. We will now see that we cannot hope to do better than  $\sqrt{n}$ .

So far we have discussed two concrete algorithms. In general, an algorithm is specified by a *policy*  $\pi$ . A policy is a sequence of conditional probabilities that specify the probability of the action at time  $t$  given the history  $H_{t-1} = \{A_1, X_1, \dots, A_{t-1}, X_{t-1}\}$ . The policies of the explore-then-exploit as well as the UCB algorithm were deterministic (other than perhaps when breaking ties). But in general a policy might be randomized. Recall also our notion of *environment*  $\nu$ . The environment is the set of  $K$  probability distributions  $\nu = (\mathbb{P}_1, \dots, \mathbb{P}_K)$ .

**Lemma 5.7** (Lower Bound on Worst-Case Regret). *Let  $K > 1$  and  $n \geq K-1$ . Then for any policy  $\pi$  there exists an environment  $\nu$  so that the regret  $R_n(\pi, \nu) \geq \frac{1}{27} \sqrt{(K-1)n}$ . Further, this environment can be chosen to be a Gaussian environment, where all distributions are unit-variance Gaussians.*

*Proof.* The idea of the proof is the following. We are given a policy  $\pi$ . Based on this policy we construct two Gaussian environments that are quite similar and differ only in a single distribution. We then show that the given policy  $\pi$  cannot do well on *both* of these

environments. Note that the “bad” environment that we prove to exist depends in general not only on the policy but also on  $n$ .

Let  $K$  be the number of arms,  $K > 1$ , and let  $\pi$  be given policy. Our first environment is Gaussian with unit-variance distributions and a mean vector of the form  $(\Delta, 0, \dots, 0)$ , where  $\Delta > 0$  is a parameter. We will chose it to be  $\sqrt{(K-1)/(4n)}$ .

Let  $\mathbb{E}_\nu[T_k(n)]$  denote the expected number of times we choose arm  $k$  for this environment  $\nu$  under the policy  $\pi$  (since the policy  $\pi$  is fixed we do not explicitly denote it). Let  $i$ ,  $1 \leq i \leq K$ , be an arm that we choose the least often under this policy. More formally,  $i = \operatorname{argmin}_k \mathbb{E}_\nu[T_k(n)]$ . If  $i = 1$  then at least  $(1 - 1/K)n$  times we do not choose arm 1 and so our regret is at least  $(1 - 1/K)n\Delta$ , which, for our choice of  $\Delta$  gives a regret of  $(1 - 1/K)n\sqrt{(K-1)/(4n)} \geq \frac{1}{4}\sqrt{(K-1)n} \geq \frac{1}{27}\sqrt{(K-1)n}$ .

So let us assume that  $i \neq 1$ . In this case the second environment  $\nu'$  is again Gaussian with unit-variance distributions and a mean vector of

$$(\Delta, 0, \dots, 0, \underbrace{2\Delta}_{i\text{-th component}}, 0, \dots, 0).$$

Let  $p_\nu(A_1, X_1, \dots, A_n, X_n)$  denote the joint distribution under policy  $\pi$  in environment  $\nu$  and let  $p_{\nu'}(A_1, X_1, \dots, A_n, X_n)$  denote the joint distribution under policy  $\pi$  in environment  $\nu'$ . How different are these distributions? Let us compute their KL divergence,

$$D(p_\nu, p_{\nu'}) = \int p_\nu(A_1, X_1, \dots, A_n, X_n) \ln \frac{p_\nu(A_1, X_1, \dots, A_n, X_n)}{p_{\nu'}(A_1, X_1, \dots, A_n, X_n)}.$$

Note that these distributions can be factorized in the following form

$$p_\nu(A_1, X_1, \dots, A_n, X_n) = \pi(A_1)\pi_\nu(X_1 | A_1) \cdots \pi_{H_{n-1}}(A_n)\pi_\nu(X_n | A_n).$$



Therefore

$$\begin{aligned}
D(p_\nu, p_{\nu'}) &= \int p_\nu(A_1, X_1, \dots, A_n, X_n) \ln \frac{p_\nu(A_1, X_1, \dots, A_n, X_n)}{p_{\nu'}(A_1, X_1, \dots, A_n, X_n)} \\
&= \int p_\nu(A_1, X_1, \dots, A_n, X_n) \ln \frac{\pi(A_1)\pi_\nu(X_1 | A_1) \cdots \pi_{H_{n-1}}(A_n)\pi_\nu(X_n | A_n)}{\pi(A_1)\pi_{\nu'}(X_1 | A_1) \cdots \pi_{H_{n-1}}(A_n)\pi_{\nu'}(X_n | A_n)} \\
&= \int p_\nu(A_1, X_1, \dots, A_n, X_n) \ln \frac{\pi_\nu(X_1 | A_1) \cdots \pi_\nu(X_n | A_n)}{\pi_{\nu'}(X_1 | A_1) \cdots \pi_{\nu'}(X_n | A_n)} \\
&= \sum_{t=1}^n \int p_\nu(A_1, X_1, \dots, A_t, X_t) \ln \frac{\pi_\nu(X_t | A_t)}{\pi_{\nu'}(X_t | A_t)} \\
&= \sum_{t=1}^n \int \sum_{k=1}^K p_\nu(A_t = k) p_\nu(X_t | A_t = k) \ln \frac{\pi_\nu(X_t | A_t = k)}{\pi_{\nu'}(X_t | A_t = k)} \\
&= \sum_{t=1}^n \sum_{k=1}^K p_\nu(A_t = k) D(P_k, P'_k) \\
&= \sum_{k=1}^K \mathbb{E}_\nu[T_k(n)] D(P_k, P'_k) \\
&\stackrel{(a)}{=} \mathbb{E}_\nu[T_i(n)] \frac{4\Delta^2}{2} \\
&\stackrel{(b)}{\leq} \frac{n}{K-1} \frac{4\Delta^2}{2} = \frac{2n\Delta^2}{K-1}.
\end{aligned}$$

In step (a) we have used the fact that the two environments only differ in position  $i$  and that in this position we have two unit-variance Gaussians, one with mean 0 and one with mean  $2\Delta$ . As you will show in your homework, if  $P_i$ ,  $i = 1, 2$ , are two Gaussians with means  $\mu_i$  and variances  $\sigma_i^2$ , then

$$D_{KL}(P_1 || P_2) = \ln(\sigma_2/\sigma_1) + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2}.$$

To see step (b) note that amongst the  $K-1$  arms  $2, \dots, K$ , the one that is chosen the least cannot be chosen more than  $n/(K-1)$  times. We now have

$$\begin{aligned}
R_n(\pi, \nu) + R_n(\pi, \nu') &\stackrel{(a)}{\geq} \mathbb{P}_\nu\{T_1(n) \leq n/2\} \frac{n\Delta}{2} + \mathbb{P}_{\nu'}\{T_1(n) > n/2\} \frac{n\Delta}{2} \\
&\stackrel{(b)}{\geq} \frac{n\Delta}{4} e^{-D(\mathbb{P}_\nu, \mathbb{P}_{\nu'})} \\
&\geq \frac{n\Delta}{4} e^{-\frac{2n\Delta^2}{K-1}} \\
&\stackrel{(c)}{=} \frac{\sqrt{n(K-1)}}{8} e^{-\frac{1}{2}} \\
&\stackrel{(d)}{=} \frac{2\sqrt{n(K-1)}}{27}.
\end{aligned}$$

Before we justify any of these steps note that this inequality completes our proof: We have shown that for a given policy there are two environments so that the sum of their regrets at time  $n$  is at least  $2\frac{\sqrt{n(K-1)}}{27}$ . So at least one of these environments must have a regret of at least  $\frac{\sqrt{n(K-1)}}{27}$ .

Step (a) is easy to explain. If we choose the arm 1 in environment  $\mu$  at most half of the time then for  $n/2$  time steps we have a regret of  $\Delta$  for each step. And if we choose the arm 1 in environment  $\nu'$  more than half of the time then for  $n/2$  time steps we have a regret of at least  $\Delta$ .

In step (c) we made the choice  $\Delta = \sqrt{(K-1)/(4n)}$  and in step (d) we lower bounded  $e^{-\frac{1}{2}}/8 \sim 0.0758$  by  $2/27 \sim 0.07407$ .

To see step (b) let  $P$  and  $Q$  be two distributions,  $p$  and  $q$  be their densities, and let  $A$  be an event. We have

$$\begin{aligned}
P(A) + Q(A^c) &= \int_{x \in A} p(x) + \int_{x \in A^c} q(x) \\
&\geq \int_{x \in A} \min\{p(x), q(x)\} + \int_{x \in A^c} \min\{p(x), q(x)\} \\
&= \int \min\{p(x), q(x)\} \\
&\stackrel{(a)}{\geq} \frac{1}{2} \left( \int \min\{p(x), q(x)\} \right) \left( \int \max\{p(x), q(x)\} \right) \\
&\stackrel{(a)}{\geq} \frac{1}{2} \left( \int \sqrt{\min\{p(x), q(x)\} \max\{p(x), q(x)\}} \right)^2 \\
&= \frac{1}{2} \left( \int \sqrt{p(x)q(x)} \right)^2 \\
&= \frac{1}{2} e^{2 \ln \int \sqrt{p(x)q(x)}} \\
&= \frac{1}{2} e^{2 \ln \int p \sqrt{\frac{q(x)}{p(x)}}} \\
&\stackrel{\text{Jensen}}{\geq} \frac{1}{2} e^{2 \int \frac{1}{2} p \ln \frac{q(x)}{p(x)}} \\
&= \frac{1}{2} e^{\int p \ln \frac{q(x)}{p(x)}} \\
&= \frac{1}{2} e^{-D(p,q)}.
\end{aligned}$$

□

For step (a) note that  $\frac{1}{2} \left( \int \max\{p(x), q(x)\} \right) \leq 1$  with equality if the two distributions have disjoint support. Step (b) follows by Cauchy-Schwartz.<sup>1</sup>

---

<sup>1</sup>  $|fg|_1^2 \leq |f|_2^2 |g|_2^2$  with  $f = \sqrt{\min\{p(x), q(x)\}}$  and  $g = \sqrt{\max\{p(x), q(x)\}}$ .

## 5.4 Further Topics

There are many extensions and variations of this topic. Let us quickly mention a few without proofs or details.

### 5.4.1 Asymptotic Optimality

Assume we are given a fixed policy  $\pi$ .

In Section 5.3.4 proved that if we first fix the time  $n$  and then choose an environment we can make the regret as large as  $\sqrt{n}$ .

But what if we first fix the environment and then let  $n$  tend to infinity. We have seen that for the UCB algorithm the asymptotic regret scales logarithmically. Can we do better? It turns out that we cannot as long as we stick to policies that have an asymptotic regret that is upper bounded by  $n^\alpha$ , for every  $\alpha > 0$ , for all environments (E.g., the UCB fulfills this condition). So the UCB algorithm is optimal also in this sense.

### 5.4.2 Adversarial Bandits

So far we have assumed that the environment consists of  $K$  distributions that are unknown but fixed. This is a relatively strong assumption. In the *adversarial bandit* setting we do not assume that the rewards are iid samples from a distribution. Rather, we allow them to be arbitrary numbers  $x_{t,k}$  in  $[0, 1]$ .

Assume at first that the policy is deterministic. Then it is clear that we can make the regret equal to  $n$ . For any time  $t$ , once  $A_t$  has been chosen by the policy, let's say it has value  $k$ , pick  $j \neq k$ , and set  $x_{t,j} = 1$ , and  $x_{t,i} = 0$ ,  $i \neq j$ .

We can remedy this problem by making the following two changes. First, clearly we need a randomized strategy. Second, we will compare ourselves to a genie who knows all rewards  $x_{t,k}$  and who picks that arm whose average reward is largest up to time  $n$  (so we do NOT compare to a genie who is allowed to pick at every time  $t$  that arm that contains the highest reward at this time).

Perhaps surprisingly, for a randomized strategy and this proper choice of genie we can make the regret almost as small as for the stochastic case. Here, the regret is the expected regret, where the expectation is over the choice of randomness of the algorithm.

### Exponential-Weight Algorithm for Exploration and Exploitation

The most common algorithm in this setting is the *Exponential-weight algorithm for Exploration and Exploitation* (Exp3 for short). It is defined as follows.

We start with a uniform distribution on the set of actions,  $\mathbb{P}_{t=1,k} = 1/K$ ,  $k = 1, \dots, K$ . At time  $t$ , we have computed the distribution  $\mathbb{P}_{t,k}$ . Sample an action  $A_t$  from this distribution, assume it is  $k$ . Reveal the sample. It is the number  $x_{t,k}$  and we call it  $X_t$ . Estimate the rewards for all arms based on  $X_t$  and then compute  $P_{t+1,k}$  by updating  $P_{t,k}$ .

### Reward Estimation

Recall that in round  $t$  we chose some action  $A_t$  according to the distribution  $\mathbb{P}_{t,k}$  and then we observed the reward  $X_t$  which is the  $t$ -reward of arm  $A_t$ . Based on this number we would like to estimate the reward of all arms. We use the estimator

$$\hat{x}_{t,k} = \frac{\mathbb{I}_{\{A_t=k\}}}{\mathbb{P}_{t,k}} X_t.$$

This makes sense. We scale each number by one over the probability that we sample it. We have<sup>2</sup>

$$\begin{aligned} \mathbb{E}[\hat{x}_{t,k} \mid A_1, X_1, \dots, A_{t-1}, X_{t-1}] &= \mathbb{E}\left[\frac{\mathbb{I}_{\{A_t=k\}}}{\mathbb{P}_{t,k}} X_t \mid A_1, X_1, \dots, A_{t-1}, X_{t-1}\right] \\ &= \mathbb{E}\left[\frac{\mathbb{I}_{\{A_t=k\}}}{\mathbb{P}_{t,k}} x_{t,k} \mid A_1, X_1, \dots, A_{t-1}, X_{t-1}\right] \\ &= \frac{x_{t,k}}{\mathbb{P}_{t,k}} \mathbb{E}[\mathbb{I}_{\{A_t=k\}} \mid A_1, X_1, \dots, A_{t-1}, X_{t-1}] \\ &= x_{t,k}. \end{aligned}$$

Note that the conditional expectation  $\mathbb{E}[\hat{x}_{t,k} \mid A_1, X_1, \dots, A_{t-1}, X_{t-1}]$  is a random variable (it depends on the history) that is constant on the “partitions” given by the conditioning. What the above says is that independent of the history this random variable is in fact equal to  $x_{t,k}$ . So irrespective of the history, if we repeated the same history many times, the expected regret we would see is always  $x_{t,k}$ . This makes of course sense since the history only changes  $\mathbb{P}_{t,k}$  but by definition we sampled exactly according to this distribution. (In order for things to be well-defined we need to ensure that the probability of sampling a particular  $k$  is never 0.)

In the same way we can compute the variance of this estimator. We have

$$\begin{aligned} &= \mathbb{E}[\hat{x}_{t,k}^2 \mid A_1, X_1, \dots, A_{t-1}, X_{t-1}] - \mathbb{E}[\hat{x}_{t,k} \mid A_1, X_1, \dots, A_{t-1}, X_{t-1}]^2 \\ &= \mathbb{E}\left[\left(\frac{\mathbb{I}_{\{A_t=k\}}}{\mathbb{P}_{t,k}} X_t\right)^2 \mid A_1, X_1, \dots, A_{t-1}, X_{t-1}\right] - x_{t,k}^2 \\ &= \mathbb{E}\left[\frac{\mathbb{I}_{\{A_t=k\}}}{\mathbb{P}_{t,k}^2} x_{t,k}^2 \mid A_1, X_1, \dots, A_{t-1}, X_{t-1}\right] - x_{t,k}^2 \\ &= \frac{x_{t,k}^2}{\mathbb{P}_{t,k}^2} \mathbb{E}[\mathbb{I}_{\{A_t=k\}} \mid A_1, X_1, \dots, A_{t-1}, X_{t-1}] - x_{t,k}^2 \\ &= x_{t,k}^2 \frac{1 - \mathbb{P}_{t,k}}{\mathbb{P}_{t,k}}. \end{aligned}$$

<sup>2</sup>Recall the definition of conditional expectation. Let  $X$  be a random variable defined on a  $\sigma$ -algebra  $\mathcal{F}$ . Let  $\mathcal{G}$  be a *sub  $\sigma$ -algebra*. Then the conditional  $\mathbb{E}[X \mid \mathcal{G}]$  is the unique random variable that is  $\mathcal{G}$ -measurable and so that for all  $G \in \mathcal{G}$

$$\int_G \mathbb{E}[X \mid \mathcal{G}] d\mathbb{P} = \int_G X d\mathbb{P}$$

From this we see that the variance can be substantial if  $\mathbb{P}_{t,k}$  is very small. This can of course cause trouble.

### Updating the Probability Distribution

Now where we have discussed how we can estimate the total reward of each arm up to time  $t - 1$ , call this quantity  $\hat{S}_{t-1,k}$ , we still need to discuss how we convert this estimate into the probability distribution  $\mathbb{P}_{t,k}$ . One standard way of doing this is to set

$$\mathbb{P}_{t,k} = \frac{e^{\eta \hat{S}_{t-1,k}}}{\sum_j e^{\eta \hat{S}_{t-1,j}}},$$

where  $\eta$  is a parameter that we can choose freely.

**Lemma 5.8** (Regret of Exp3). *For any assignment of the rewards  $x_{t,k} \in [0, 1]$  the expected rewards of the Exp3 algorithm is bounded as*

$$R_n \leq 2\sqrt{nK \ln(K)}.$$

### 5.4.3 Contextual Bandits

In many scenarios we have some side information available. How can we use this information to improve our choice. One idea is to define a *context*. E.g., perhaps we built a movie recommendation site. In this case we might have a prior classification of various user “types.” This could be the context.

Assume that there are a finite number of contexts. Then we could define one bandit algorithm for each of the finite number of contexts and run them independently. But we pay a price. Now each algorithm only sees a fraction of the examples!

## 5.5 Problems

**Problem 5.1.** Compute the KL Divergence of two scalar Gaussians  $p(x) = \mathcal{N}(\mu_1, \sigma_1^2)$  and  $q(x) = \mathcal{N}(\mu_2, \sigma_2^2)$ .

**Problem 5.2.** In the course we analyzed the Upper Confidence Bound algorithm. As was suggested in the course, we should get something similar if instead we use the Lower Confidence Bound algorithm. It is formally defined as follows.

$$A_t = \begin{cases} t, & t \leq K, \\ \arg \max_k \hat{\mu}_k(t-1) - \sqrt{\frac{2 \ln f(t)}{T_k(t-1)}}, & t > K. \end{cases}$$

Analyze the performance of this algorithm in the same way as we did this in the course for the UCB algorithm.

**Problem 5.3.** Recall our original *explore-then-exploit* strategy. We had a fixed time horizon  $n$ . For some  $m$ , a function of  $n$  and the gaps  $\{\Delta_k\}$ , we explore each of the  $K$  arms  $m$  times initially. Then we pick the best arm according to their empirical gains and play this arm until we reach round  $n$ . We have seen that this strategy achieves an asymptotic regret of order  $\ln(n)$  if the environment is fixed and we think of  $n$  tending to infinity but a worst-case regret of order  $\sqrt{n}$  if we use the gaps when determining  $m$  and of order  $n^{\frac{2}{3}}$  if we do not use the gaps in order to determine  $m$ .

Here is a slightly different algorithm. Let  $\epsilon_t = t^{-\frac{1}{3}}$ . For each round  $t = 1, \dots$ , toss a coin with success probability  $\epsilon_t$ . If success, then explore arms uniformly at random. If not success, then pick in this round the arm that currently has the highest empirical average.

Show that for this algorithm the expected regret at *any* time  $t$  is upper bounded by  $t^{\frac{2}{3}}$  times terms in  $t$  and  $K$  of lower order. This is similar to the worst-case of the explore-then-exploit strategy but here we do not need to know the horizon a priori. Assume that the rewards are in  $[0, 1]$ .

**Problem 5.4.** Prove Lemma 7.4 in the lecture notes. In other words, prove that if  $X$  is a zero-mean random variable taking values in  $[a, b]$  then

$$\mathbb{E}[e^{\lambda X}] \leq e^{\frac{\lambda^2}{2}[(a-b)^2/4]}.$$

Expressed differently,  $X$  is  $[(a-b)^2/4]$ -subgaussian.

## Chapter 6

# Distribution Estimation, Property Testing and Property Estimation

Assume that we are given iid samples from an unknown distribution. How many samples do we need before we can estimate the distribution with an “acceptable” accuracy? And what if we are interested in only particular properties of the distribution, such as its support size, or perhaps it’s entropy. These are the questions that we will discuss in this chapter.

This chapter closely follows the tutorial by Acharya, Orlitsky, and Suresh, see

<https://people.ece.cornell.edu/acharya/papers/isit-tutorial-acharya-orlitsky-suresh.pdf>

## 6.1 Distribution Estimation

### 6.1.1 Notation and Basic Task

Consider a random variable  $X$  taking values in the discrete set  $\mathcal{X}$  and let  $p(x), x \in \mathcal{X}$ , describe the distribution of  $X$ . Of course we have  $p : \mathcal{X} \rightarrow \mathbb{R}_{\geq 0}$  and  $\sum_{x \in \mathcal{X}} p(x) = 1$ . We assume that  $\mathcal{X}$  has a finite support,  $|\mathcal{X}| = k$ . Without loss of generality we will assume that  $\mathcal{X} = \{1, \dots, k\}$ . In this case we can think of  $p(x)$  also as a vector of length  $k$  written as  $p = (p_1, \dots, p_k)$ . Note that in this way  $p$  is an element of  $\Delta_k$ , the simplex in  $\mathbb{R}^k$ .

In the sequel we will assume that we are given a sample of  $n$  elements from  $\mathcal{X}$ , call it  $x^n = x_1, \dots, x_n$ , chosen iid according to  $p(x)$ , so that

$$p(x^n) = \prod_{i=1}^n p(x_i).$$

Given this sample  $x^n$  our task is to find a distribution  $q = q(x^n)$ ,  $q \in \Delta_k$ , that is “close” to the distribution  $p$ .

### 6.1.2 Empirical Estimator

The perhaps most “natural” estimator is the *empirical* estimator  $q^{\text{emp}}$ . Given a sequence  $x^n$  let

$$t_i(x^n) = |\{j \in \{1, \dots, n\} : x_j = i\}|.$$

In words,  $t_i(x^n)$  counts how many times the symbol  $i$  appears in  $x^n$ . The empirical estimator is then defined as

$$q_i^{\text{emp}}(x^n) = t_i(x^n)/n.$$

Clearly,

$$\begin{aligned} q_i^{\text{emp}}(x^n) &\geq 0, \\ \sum_{i=1}^k q_i^{\text{emp}}(x^n) &= 1. \end{aligned}$$

In other words,  $q^{\text{emp}}(x^n) \in \Delta_k$ , so this estimator is well-defined.

**Example 6.1.** Let’s assume that  $n = 4$ ,  $k = 3$ , and  $x^4 = 3112$ . Then

$$q^{\text{emp}}(3112) = (q_1^{\text{emp}}(3112), q_2^{\text{emp}}(3112), q_3^{\text{emp}}(3112)) = \left(\frac{2}{4}, \frac{1}{4}, \frac{1}{4}\right).$$

The empirical estimator will play a prominent role in the following.

### 6.1.3 Loss Functions

Before we can analyse how a given estimator does we need to specify how we will measure the quality of the estimator. More precisely, given  $p$  and  $q$ , how do we measure their distance? The three most common choices are the  $\ell_1$  distance, the  $\ell_2$  distance, and the Kullback-Leibler divergence. Generically we call the loss  $L(p, q)$ . We will start by looking at  $\ell_2^2$  since it is mathematically the most convenient.

### 6.1.4 Min-Max Criterion

So assume that we have fixed the loss function  $L$ . We then still have various degrees of freedom. Let us go over these options.

- We are given a fixed distribution  $p$  and a fixed estimator  $q$ . It is then natural that we compute the expected loss, where the expectation is over the sample:

$$\mathbb{E}_{X^n \sim p}[L(p, q(X^n))].$$

- We are given a fixed distribution  $p$ . For each estimator we can compute an expected loss as we discuss in the previous scenario. It is then natural to find that estimator that minimizes this expected loss:

$$q^* = \operatorname{argmin}_q \mathbb{E}_{X^n \sim p}[L(p, q(X^n))].$$



- We are given a fixed estimator  $q$ . For each fixed distribution we can compute an expected loss as we discussed in the first scenario. It is then natural that we find that distribution  $p^*$  that maximizes this expected loss:

$$p^* = \operatorname{argmax}_{p \in \Delta_k} \mathbb{E}_{X^n \sim p}[L(p, q(X^n))].$$

What if we are neither given  $p$  nor  $q$ ? We get a *robust* definition if we choose the estimator  $q$  in such a way that we *minimize* the expected risk for the *worst* distribution  $p$ . This is called the *min-max* criterion and in formulae it reads

$$r_{k,n}^L = \min_q \max_{p \in \Delta_k} \mathbb{E}_{X^n \sim p}[L(p, q(X^n))].$$

Let us emphasize: For each  $q$  (estimator) we look at that  $p$  (distribution) that gives the worst result. We then pick that  $q$  that minimizes the worst case.

We will be interested in finding what this min-max optimal estimator is and how it performs. Our strategy will be the following. We first compute the risk of the empirical estimator. By computing a lower bound on the risk, we will then see that a small variant of this natural estimator is min-max optimal.

### 6.1.5 Risk of Empirical Estimator in $\ell_2^2$

We start by looking at the case where the difference between the true distribution and the estimated one is measured in  $\ell_2^2$  distance, i.e.,

$$\|p - q\|_2^2 = \sum_i (p_i - q_i)^2.$$

We want to compute

$$r_{k,n}^{q^{\text{emp}}} = \max_{p \in \Delta_k} \mathbb{E}_{X^n \sim p} \left[ \sum_i (p_i - q_i^{\text{emp}}(X^n))^2 \right],$$

where  $q_i^{\text{emp}}(x^n) = t_i(x^n)/n$ .

Recall that the components of  $X^n$  are iid, chosen from  $\mathcal{X}$  according to the distribution  $p$ . This (the fact that they are iid) implies that for all  $i \in \mathcal{X}$ ,  $t_i(X^n)$  is a Binomial with parameters  $\text{Binom}(p_i, n)$ . Therefore,

$$\begin{aligned} \mathbb{E}_{X^n \sim p}[t_i(X^n)] &= np_i, \\ \mathbb{E}_{X^n \sim p}[(t_i(X^n) - np_i)^2] &= np_i(1 - p_i). \end{aligned}$$

Hence,

$$\mathbb{E}_{X^n \sim p} \left[ \sum_{i=1}^k (t_i(X^n)/n - p_i)^2 \right] = \sum_{i=1}^k \frac{p_i(1 - p_i)}{n} = \frac{1 - \sum_{i=1}^k p_i^2}{n} \stackrel{(a)}{\leq} \frac{1 - \frac{1}{k}}{n}.$$

In step we have used the Cauchy-Schwartz inequality,

$$|\langle x, y \rangle|^2 \leq \langle x, x \rangle \cdot \langle y, y \rangle.$$

To get our inequality, pick  $x = (1, \dots, 1)$  and  $y = p$ . Note that this inequality is achieved by the uniform distribution  $p_i = 1/k$ .

Note that the above bound is *universal* with respect to the underlying distribution  $p$ . This is good news. So we have shown that for  $\ell_2^2$  loss

$$\max_{p \in \Delta_k} \mathbb{E}_{X^n \sim p} [L(p, q^{\text{emp}}(X^n))] = \frac{1 - \frac{1}{k}}{n}.$$

### 6.1.6 Risk of “Add Constant” Estimator in $\ell_2^2$

Is the empirical estimator min-max optimal? Not quite. Here is a slightly better estimator:

$$q_i^{+\sqrt{n}/k}(x^n) = \frac{t_i(x^n) + \frac{\sqrt{n}}{k}}{n + \sqrt{n}}.$$

This is an instance of an “add constant” estimator. We will soon see that this estimator is min-max optimal. Intuitively, adding a constant to our observations makes sense. If the number of samples is small, we cannot possibly have seen all elements, not because their probability is zero, but because it is small and there is randomness in sampling. What is the risk of this estimator? We claim that

$$\max_{p \in \Delta_k} \mathbb{E}_{X^n \sim p} [L(p, q^{+\sqrt{n}/k}(X^n))] = \frac{1 - \frac{1}{k}}{(\sqrt{n} + 1)^2}$$

To see this claim note that

$$\mathbb{E}_{X^n \sim p} [q_i^{+\sqrt{n}/k}(x^n)] = \mathbb{E}_{X^n \sim p} \left[ \frac{t_i(X^n) + \frac{\sqrt{n}}{k}}{n + \sqrt{n}} \right] = \frac{np_i + \frac{\sqrt{n}}{k}}{n + \sqrt{n}}.$$

Note that this is a *biased* estimator. We get

$$\begin{aligned} \mathbb{E}_{X^n \sim p} [(q_i^{+\sqrt{n}/k}(X^n) - p_i)^2] &= \frac{\mathbb{E}_{X^n \sim p} [(t_i(X^n) - np_i - \frac{\sqrt{n}}{k}(kp_i - 1))^2]}{(n + \sqrt{n})^2} \\ &= \frac{\text{Var}[t_i(X^n)] + \frac{n}{k^2}(kp_i - 1)^2}{(n + \sqrt{n})^2} \\ &= \frac{np_i(1 - p_i) + \frac{n}{k^2}(kp_i - 1)^2}{(n + \sqrt{n})^2} \\ &= \frac{np_i(1 - \frac{2}{k}) + \frac{n}{k^2}}{(n + \sqrt{n})^2}. \end{aligned}$$

To compute the worst-case loss of the  $q^{+\sqrt{n}/k}(x^n)$  estimator we need to sum this expression over all  $k$  components and then maximize with respect to the distribution  $p$ . This gives us

$$\max_{p \in \Delta_k} \sum_{i=1}^k \mathbb{E}_{X^n \sim p} [(q_i^{+\sqrt{n}/k}(X^n) - p_i)^2] = \frac{n(1 - \frac{1}{k})}{(n + \sqrt{n})^2} = \frac{(1 - \frac{1}{k})}{(\sqrt{n} + 1)^2},$$

as claimed.

Note: This calculation shows that for this estimator the loss *does not* depend on the underlying distribution  $p$ . This will become important soon.

Note: If we had done this calculation with a general additive term  $\beta$  instead of the specific term  $\beta^* = \frac{\sqrt{n}}{k}$  then it is easy to see that the choice  $\beta^*$  is optimal.

### 6.1.7 Matching lower bound for $\ell_2^2$

We will now derive a matching lower bound. We proceed as follows. Let  $\pi$  be a prior distribution on  $\Delta_k$ . We then have

$$\begin{aligned} r_{k,n}^{\ell_2^2} &= \min_q \max_{p \in \Delta_k} \mathbb{E}_{X^n \sim p} \left[ \sum_i (p_i - q_{X^n,i})^2 \right] \geq \min_q \mathbb{E}_{P \sim \pi; X^n \sim P} \left[ \sum_i (P_i - q_i(X^n))^2 \right] \\ &= \mathbb{E}_{P \sim \pi; X^n \sim P} \left[ \sum_i (P_i - \mathbb{E}_{P \sim \pi; X^n \sim P} [P_i | X^n])^2 \right]. \end{aligned}$$

where in the last step we have used the fact that the minimum-mean squared error estimator is given by the conditional expectation.

Therefore, if we can guess a “good” prior then we will get a good bound. It turns out that a suitable Dirichlet prior gives us the matching lower bound. The Dirichlet distribution on  $\Delta_k$  is characterized by a vector  $\alpha \in (\mathbb{R}_+)^k$ . Let  $(x_1, \dots, x_k) \in \Delta_k$ . The associated density is given by

$$f(x_1, \dots, x_k; \alpha) = \frac{\prod_{i=1}^k \Gamma(\alpha_i)}{\Gamma(\sum_{i=1}^k \alpha_i)} \prod_{i=1}^k x_i^{\alpha_i - 1}.$$

Note that this is an exponential distribution, i.e., it can be written in a form

$$p_{\Theta}(x) = e^{\langle \Theta, \phi(x) \rangle - A(\Theta)},$$

where  $x = (x_1, \dots, x_k)$ ,  $\phi(x) = (\ln(x_1), \dots, \ln(x_k))$ , and  $\Theta = (\alpha_1 - 1, \dots, \alpha_k - 1)$ .

The important property of a Dirichlet distribution for our application is that it is the conjugate distribution of a multi-nomial distribution. So assume that the parameters  $p_1, \dots, p_k$  of a multi-nomial distribution are themselves random with prior  $\text{Dir}(\alpha)$ . Assume that we sample from this Dirichlet distribution and then sample from the multi-nomial according to the chosen parameter. We get  $n$  samples and their counts are  $T_1, \dots, T_k$ . Given this observation what is the posterior of the parameters? I.e., we want to determine

$$f(p_1, \dots, p_k \mid T_1 = t_1, \dots, T_k = t_k).$$

We claim that this is again a Dirichlet distribution, namely with parameters  $\alpha + T$  (both vectors). Using Bayes' rule we get

$$\begin{aligned}
 f(p_1, \dots, p_k \mid t_1, \dots, t_k) &= \frac{f(t_1, \dots, t_k \mid p_1, \dots, p_k) f(p_1, \dots, p_k)}{Z(t_1, \dots, t_k)} \\
 &= \frac{\binom{n}{t_1, \dots, t_k} [\prod_{i=1}^k p_i^{t_i}] \frac{\prod_{i=1}^k \Gamma(\alpha_i)}{\Gamma(\sum_{i=1}^k \alpha_i)} [\prod_{i=1}^k p_i^{\alpha_i - 1}]}{Z(t_1, \dots, t_k)} \\
 &= \frac{\prod_{i=1}^k p_i^{\alpha_i + t_i - 1}}{Z} \\
 &= \text{Dir}(p_1, \dots, p_k; \alpha + t).
 \end{aligned}$$

The Dirichlet distribution has been well studied. In particular, its mean is known. If we assume that  $\alpha_i = \alpha$  for all  $i = 1, \dots, k$ , then we get

$$\mathbb{E}_{P \sim \pi; X^n \sim P}[P_i \mid X^n] = \frac{t_i(X^n) + \alpha}{n + k\alpha}.$$

Now note that if we pick  $\alpha = \sqrt{n}/k$  then our previous calculations have shown that the loss does not in fact depend on the distribution  $p$  but is always equal to  $\frac{1-1/k}{(\sqrt{n}+1)^2}$ . This finishes our claim.

Although  $\ell_2$  is nice and easy from an analysis perspective it has also downsides. Perhaps the biggest one is that it is not a good measure if  $k$  is very large. Assume that  $p$  has  $2/k$  in its first  $k/2$  components and 0 in its remaining ones and assume that for  $q$  the roles of these two parts are exactly reversed. Their  $\ell_2^2$  distance is then equal to  $4/k$  which quickly converges to 0 as  $k$  gets large. But it is hard to think of distributions that are more different! If instead we looked at their  $\ell_1$  distance we would get 2 and their KL “distance” is in fact infinity.

### 6.1.8 Risk in $\ell_1$

This motivates us to look at  $\ell_1$ , i.e., now we look at

$$\|p - q\|_1 = \sum_{i=1}^k |p_i - q_i|.$$

Note that this has a probabilistic interpretation. E.g., we can couple the two random variables so that up to a fraction of time  $\|p - q\|_1$  they take the same value.

**Lemma 6.1.** *For  $k$  fixed and as  $n$  tends to infinity, the worst case min-max loss behaves like*

$$r_{k,n} \leq \sqrt{\frac{2(k-1)}{\pi n}} + O(n^{-\frac{3}{4}}).$$

*Further, this is achieved by the empirical estimator.*

The idea of the proof is similar to the technique we used before. We first compute the loss of the empirical estimator. We show that this loss is highest for the uniform distribution and this gives us an upper bound. Then we derive a lower bound that matches the dominant terms by computing again the Bayes loss with a proper Dirichlet prior. We skip the details.

If we are content with a slightly looser upper bound we can proceed as follows. Consider the empirical estimator  $q^{\text{emp}}(X^n)$ . We have

$$\begin{aligned} \mathbb{E}[\|p - q^{\text{emp}}(X^n)\|_1] &= \sum_{i=1}^k \mathbb{E}[|p_i - \frac{T_i(X^n)}{n}|] \\ &\stackrel{(a)}{\leq} \sum_{i=1}^k \sqrt{\mathbb{E}[|p_i - \frac{T_i(X^n)}{n}|^2]} \\ &\stackrel{(b)}{\leq} \sqrt{\frac{k-1}{n}}. \end{aligned}$$

Step (a) follows by Jensen's inequality and in step (b) we have used our results for  $\ell_2^2$  and the Cauchy-Schwartz inequality. We see that we loose a factor  $2/\pi$  compared to the previous result.

### 6.1.9 Risk in KL-Divergence

If we are using the KL divergence as our loss metric then we need to make sure that none of our estimated probabilities are 0 since otherwise our metric will be  $\infty$ . It is therefore natural to use an “add constant” estimator.

When the number of samples becomes large compared to the alphabet size one can show that the best “add constant” estimator is of the form  $q_i^{+0.509}$  and this gives us an expected worst case loss of

$$\max_{p \in \Delta_k} \mathbb{E}[D(p \| q^{+0.509})] \sim 0.509 \frac{k-1}{n}. \quad (6.1)$$

One can do slightly better ( $\frac{1}{2}$  instead of 0.509) by using different constants depending on the observed frequency.

Note that by the Pinsker inequality we have

$$\sqrt{1.2 \frac{k-1}{n}} \sim \max_{p \in \Delta_k} \sqrt{2\mathbb{E}[D(p \| q^{+0.509})]} \geq \max_{p \in \Delta_k} \mathbb{E}[\sqrt{2D(p \| q^{+0.509})}] \geq \max_{p \in \Delta_k} \mathbb{E}[\|p - q^{+0.509}\|_1].$$

The first step is (6.1). Note that this step is approximate (valid for large ratios  $n/k$ ). In the second we have used Jensen's inequality. The final step is Pinsker's inequality  $\|p - q\|_1 \leq \sqrt{2D(p \| q)}$ . We see from this sequence of inequalities that these two results are related roughly as we would expect. (Compared to the result in Lemma 6.1 we loose only a factor 0.64 inside the square root.)

### 6.1.10 The problem with the min-max formulation

We have now surveyed how the min-max estimator behaves for various risks. One problem we encounter is that min-max is really quite pessimistic. Yes, the worst case is as good as we could hope for but the estimator could be quite bad for pretty much any case. In fact, we have seen that the whole trick of proving what the min-max estimator was for the  $\ell_2$  case was to come up with an estimator that was uniformly bad (so to speak). This brings us to a slightly different point of view.

### 6.1.11 Competitive distribution estimation

The new point of view is that we want an estimator  $q$  that is close to optimal for every distribution  $p$ . The key question is here: What is our comparison group? Let us make this more formal.

Look at the probability simplex  $\Delta_k$ . Partition this space into groups  $P_1, P_2, \dots$  so that  $\Delta_k = \cup_j P_j$ . Call this partition  $\mathcal{P}$ . Let  $L(\cdot, \cdot)$  be the loss as usual. Then we are interested in

$$r_{n,k}^{\mathcal{P}} = \min_q \max_j \left[ \max_{p \in P_j} L(p, q) - \min_{q'} \max_{p' \in P_j} L(p', q') \right].$$

This is easy to interpret. Assume at first that we pick the partition so that every element of  $\Delta_k$  forms a group on its own. Within each group we compare to  $\min_{q'} \max_{p' \in P_j} L(p', q')$ . This is the min-max estimator for that group. But since the group only exists of a single distribution we can use an estimator that *knows* that particular distribution. This measure then collapses to our original min-max formulation and, as we have discussed, this is often simply to pessimistic.

On the other hand, if our partition consists only of a single group, i.e., the “genie” we compare ourselves to has no more knowledge than we have ourselves, then our loss is 0. This is not very useful either.

The key hence is to find for every case a suitable partition. And for each partition we compare ourselves in each group to a genie who knows that group a priori. E.g., the genie might know the entropy of the distribution a priori. Or perhaps the genie knows the *set* of probabilities but not which component has which of these probabilities.

### 6.1.12 Multi-set genie estimator

Assume that the genie knows the probability multi-set. I.e., the genie knows the set  $\{p_i\}$  but not the order. E.g., if  $p_1 = 0.4$ ,  $p_2 = 0.2$ , and  $p_3 = 0.4$  then the genie would be given the set (ordered in decreasing order of values)  $\{0.4, 0.4, 0.2\}$ .

### 6.1.13 Natural Genie and Good-Turing Estimator

A slightly more powerful genie is one that in fact knows the distribution but is forced to give the same estimate to any group of symbols that appear the same number of times. E.g.,

$X^5 = 12213$ . Then this genie will use the estimates

$$q_1 = q_2 = \frac{p_1 + p_2}{2},$$

$$q_3 = p_3.$$

Let  $M(t)$  denote the total probability of all symbols that appear exactly  $t$  times. And let us assume that we are using the KL divergence as loss function. Further, let  $\phi(t)$  denote the number of symbols that appeared  $t$  times.

The so-called Good-Turing estimator is then

$$q_i^{\text{GT}} = \frac{T_i + 1}{n} \frac{\phi(T_i + 1)}{\phi(T_i)}.$$

This estimator has a fabled history and was supposedly one of the tools used in breaking the Enigma code. Good published it in 1953 based on an unpublished note by Turing.

Let us do an example. Assume that  $X^9 = 121234555$ . We then have  $\phi(1) = 2$  since 3 and 4 appeared once,  $\phi(2) = 2$  since 1 and 2 appeared twice, and  $\phi(3) = 1$  since 5 appeared three times. We then have

$$q_3 = \frac{T_3 + 1}{9} \frac{\phi(2)}{\phi(1)} = \frac{2}{9} \frac{2}{2} = \frac{2}{9}.$$

What is the intuition for this estimator? The intuition comes by looking at what the natural genie will do. Recall, it will give the probability  $M(t)/\phi(t)$  to each of the  $\phi(t)$  symbols that appear  $t$  times. We are trying to compete against this genie. So it makes sense that we use an expression motivated by this estimate. Of course, we do not know what  $M(t)$  is since we do not know the probabilities. We claim that

$$\mathbb{E}[M(t)] = \frac{t+1}{n} \mathbb{E}[\phi(t+1)]. \quad (6.2)$$

To be slightly more precise. We claim that we have this identity if we use *Poisson* sampling. You will explore this more in the homework. It is a standard trick to get rid of the dependency that you get between the various coefficients when you sample a fixed number of samples.

Assume that we have given a distribution  $p$  on  $\mathcal{X} = \{1, \dots, k\}$ . Let  $X^n$  denote a sequence of  $n$  iid samples. Let  $T_i = T_i(X^n)$  be the number of times symbol  $i$  appears in  $X^n$ . Then

$$\mathbb{P}\{T_i = t_i\} = \binom{n}{t_i} p_i^{t_i} (1 - p_i)^{n-t_i}.$$

Note that the random variables  $T_i$  are *dependent*, since  $\sum_i T_i = n$ . This dependence can cause difficulties if we are using this distribution in a scheme and want to analyse its performance.

There is a convenient way of getting around this problem. This is called *Poisson* sampling. Let  $N$  be a random variable distributed according to a Poisson distribution with mean  $n$ . Let  $X^N$  be an iid sequence of  $N$  variables distributed according to  $p$ .

Then the following statements are true.

- $T_i(X^N)$  is distributed according to a Poisson random variable with mean  $p_i n$ .
- The  $T_i(X^N)$  are independent.
- Conditioned on  $N = n$ , the induced distribution of the Poisson sampling scheme is equal to the distribution of the *original scheme*.

We will verify (6.2) in a moment. Equation (6.2) does not completely solve our problem since we do not know  $\mathbb{E}[\phi(t+1)]$ . But we do have its “instantaneous” value  $\phi(t+1)$ . Hence define  $\hat{M}(t) = \frac{t+1}{n} \phi(t+1)$ . If we now use  $\hat{M}(t)/\phi(t)$  instead of  $M(t)/\phi(t)$  then we get our Good-Turing estimator. Let us now verify (6.2). Note that

$$M(t) = \sum_{i=1}^k p_i \mathbb{1}_{\{T_i(X^N)=t\}},$$

where in the notation  $M(t)$  we omit the fact that this quantity depends on  $X^n$ . We now have

$$\begin{aligned} \mathbb{E}[\hat{M}(t)] &= \sum_i \frac{t+1}{n} \mathbb{E}[\mathbb{1}_{\{T_i(X^N)=t+1\}}] \\ &= \sum_i \frac{t+1}{n} e^{-(np_i)} \frac{(np_i)^{t+1}}{(t+1)!} \\ &= \sum_i \frac{np_i}{n} e^{-(np_i)} \frac{(np_i)^t}{(t)!} \\ &= \sum_i p_i e^{-(np_i)} \frac{(np_i)^t}{(t)!} \\ &= \sum_i p_i \mathbb{E}[\mathbb{1}_{\{T_i(X^N)=t\}}] \\ &= \mathbb{E}[M(t)]. \end{aligned}$$

We can now write down the competitive loss. We have.

$$r_{n,k}^{\text{nat}} = \min_{\hat{M}} \max_p \mathbb{E} \left[ \sum_{t=0}^n M(t) \ln \frac{M(t)}{\hat{M}(t)} \right].$$

One can show that an estimator based on the Good-Turing estimator and the empirical estimator achieves

$$r_{n,k}^{\text{nat}} \sim \min \left\{ n^{-\frac{1}{3}}, \frac{k}{n} \right\}.$$

## 6.2 Property Testing

We are given iid samples from an unknown distribution and we want to know if this distribution has a particular property or if it is at least an  $\epsilon$  away from having this property. Here are a few examples.



1. We want to test if the distribution is uniform.
2. We want to test if the distribution is equal to a given distribution. This is called identity testing.
3. We want to test if a distribution over  $\mathcal{X} \times \mathcal{X}$  is the product of two marginal distributions.
4. We want to test if the pdf is monotone.
5. We want to test if the pdf is log-concave.<sup>1</sup>

We can frame all these questions in the following manner. Let  $\mathcal{P}$  and  $\mathcal{Q}$  be two *families* of distributions with  $\mathcal{P} \cap \mathcal{Q} = \emptyset$ . Let  $P \in \mathcal{P} \cup \mathcal{Q}$  and let  $X^n$  be  $n$  iid samples according to  $P$ . We are given  $X^n$  but do not know  $P$ . We are asked to decide whether the samples were drawn according to a distribution in  $\mathcal{P}$  or a distribution in  $\mathcal{Q}$ . More formally we are asked to design an estimator,  $C(X^n) \rightarrow \{\mathcal{P}, \mathcal{Q}\}$  in such a way as to minimize the maximum probability of error,

$$p^n = \max_{P \in \mathcal{P} \cup \mathcal{Q}; X^n \sim P} \max\{\mathbb{P}\{C(X^n) = \mathcal{Q} \mid P \in \mathcal{P}\}, \mathbb{P}\{C(X^n) = \mathcal{P} \mid P \in \mathcal{Q}\}\}.$$

We could ask now how this probability of error behaves as  $n$  tends to infinity. This would bring us back to questions of large deviations. But in the current context we are more interested in how the quantity behaves for a small sample size; a sample size that is just big enough so that our error probability is bounded away from  $\frac{1}{2}$ .

### 6.2.1 General Idea

Assume that we can design a so-called *test statistics*  $T(X^n) \rightarrow \mathbb{R}$  with the following properties: there exists a *threshold*  $\tau$  so that

1. if  $P \in \mathcal{P}$  then  $\mathbb{P}\{T(X^n) > \tau\} < 0.1$ ,
2. if  $P \in \mathcal{Q}$  then  $\mathbb{P}\{T(X^n) < \tau\} < 0.1$ .

In this case, given the sample  $X^n$  we simply evaluate this test statistics and make our decision accordingly. I.e., we define

$$C(X^n) = \begin{cases} \mathcal{P}, & T(X^n) < \tau, \\ \mathcal{Q}, & T(X^n) > \tau. \end{cases}$$

In the above description we have assumed that the number of samples is fixed. As we have seen this for distribution estimation it is sometimes useful to allow this number to be itself a random variable distributed according to a Poisson distribution. We will then write  $X^N$ ,  $N \sim \text{Poi}(n)$ .

---

<sup>1</sup>For discrete contiguous distributions we say that it is log-concave if  $P_i^2 \geq P_{i-1}P_{i+1}$ .

### 6.2.2 Testing Against a Uniform Distribution

Even though there are many questions that fall under the category of property testing we will consider only one – namely the question of testing whether samples come from a uniform distribution.

We will assume that the alphabet size  $k$  is known and we will ask whether the samples come from a uniform distribution with support on the *whole* alphabet size. It is important to note that, even though this is a meaningful question, and it is mathematically simpler, perhaps an even more meaningful question would be to allow distributions whose support is not all of  $\mathcal{X}$ .

#### Learning Approach

The first approach is obvious. Let us learn the distribution reasonably accurately and then compute the distance of this learned distribution to the uniform one. In the following let us assume that we measure the distance according to  $\ell_1$ . Let  $U$  denote the uniform distribution on  $\mathcal{X} = \{1, \dots, k\}$ . Then  $\mathcal{P} = \{U\}$  and  $\mathcal{Q}$  is the set of distributions that have  $\ell_1$  distance at least  $\epsilon$  from  $U$ .

We then have the following algorithm.

1. Given  $X^n$  learn  $\hat{P}$  so that  $\|\hat{P} - P\|_1 \leq \epsilon/2$  with probability at least 0.9.
2. Output decision according to

$$C(X^n) = \begin{cases} \mathcal{P}, & \|\hat{P} - U\|_1 < \epsilon/2, \\ \mathcal{Q}, & \text{otherwise.} \end{cases}$$

Let us quickly check that this scheme works as intended. If  $P \in \mathcal{P}$ , i.e.,  $P = U$  then by assumption  $\|\hat{P} - U\|_1 \leq \epsilon/2$  with probability at least 0.9. So we make a mistake with probability at most 0.1. And if  $P \in \mathcal{Q}$ , then by assumption  $\|P - U\|_1$  is at least  $\epsilon$ . Since further by assumption  $\|\hat{P} - P\|_1 < \epsilon/2$ , it follows by the triangle inequality that  $\|\hat{P} - U\|_1 > \epsilon/2$ . So we see that indeed we have constructed an appropriate decision statistics and threshold for this case.

We have seen in Section 6.1.8 that in expectation the  $\ell_1$ -distance is upper bounded by  $\sqrt{\frac{k-1}{n}}$ . This means that if we want the  $\ell_1$  risk to be bounded by  $\epsilon/2$  with some fixed probability then  $O(k/\epsilon^2)$  samples will suffice.

#### A Better Approach

We can do better than that. There is no reason we first have to learn the whole distribution if at the end we are only interested in this one bit of information. We will now see that  $\sqrt{k}/\epsilon^2$  samples suffice. This might not seem a big deal if the alphabet size is small but for large alphabet sizes this is significant.

### A Lower Bound

We claim that we need at least  $\Omega(\sqrt{k})$  samples for any fixed  $\epsilon$ . This bound does not give the correct scaling with respect to  $\epsilon$  but it does tell us that we cannot hope to do better than  $\sqrt{k}$  with respect to the alphabet size.

Recall that  $\mathcal{P} = \{U\}$ , where  $U$  is the uniform distribution on  $\{1, \dots, k\}$ . Let  $X^N$  be iid samples according to  $U$ , where  $N$  is chosen according to a  $\text{Poi}(n)$  distribution, and assume that  $n \leq k$ . Recall that in this setting  $T_i(X^N)$  has distribution  $\text{Poi}(n/k)$ . The probability that symbol  $i$  is chosen 2 or more times is equal to

$$\sum_{j \geq 2} e^{-n/k} \frac{(n/k)^j}{j!}.$$

But since for  $\lambda = k/n \leq 1$ ,  $\sum_{j \geq 2} \frac{\lambda^j}{j!} \leq \lambda^2 \sum_{j \geq 2} \frac{1}{j!} \leq \lambda^2$ ,  $\sum_{j \geq 2} e^{-n/k} \frac{(n/k)^j}{j!} \leq e^{-n/k} (n/k)^2$ . Therefore, the expected number of symbols that appear more than once is upper bounded by  $ke^{-n/k} (n/k)^2 = e^{-n/k} n^2/k$ .

Assume that we pick  $n < \sqrt{k}/10$ . Then this expected value is upper bounded by  $1/100$ . So let us recap. We have a random variable, call it  $Z$ , that is integer-valued and non-negative and whose expected value is upper bounded by  $1/100$ . So

$$1/100 \geq \mathbb{E}[Z] = \sum_{i \geq 0} \mathbb{P}\{Z = i\}i \geq \sum_{i \geq 1} \mathbb{P}\{Z = i\} = \mathbb{P}\{Z \geq 1\}.$$

Therefore the probability that none of the symbols appear at least twice is upper bounded by  $1/100$ . This is called the first moment method.

Now consider a distribution, call it  $\tilde{U}$  that is also uniform, but uniform on a subset of  $\{1, \dots, k\}$  of size  $k/2$ . By exactly the same argument, replacing  $k$  with  $k/2$  everywhere, we have that with probability at most  $1/50$  we see any of the symbols repeated more than once. We conclude that, under these conditions, we cannot hope to be able to distinguish between those two distributions. And clearly these two distributions are quite different. In fact, their  $\ell_1$  distance is 1!

We recognize that  $n \sim \sqrt{k}$  is not an arbitrary threshold. This is the threshold that we know from the *birthday* paradox. This is not a co-incidence. As we will discuss in more depth when we discuss property estimation, essentially the only information that is contained in the samples *are* the overlaps. So  $n \sim \sqrt{k}$  is when we start getting useful information.

### An Upper Bound

Now where we have the “right” (we don’t know this yet, but soon ...) scaling in  $k$  let us look at an actual algorithm that gives us the desired result of  $\sqrt{k}/\epsilon^2$  samples. Let us first relate  $\ell_1$  to  $\ell_2$ .

**Lemma 6.2.** *Let  $P, Q \in \Delta_k$ . If  $\|P - Q\|_1 \geq \epsilon$  then  $\|P - Q\|_2^2 \geq \epsilon^2/k$ .*

*Proof.* By Cauchy-Schwarz  $|\langle u, v \rangle|^2 \leq \langle u, u \rangle \langle v, v \rangle$ , with  $u^\top = (|P_1 - Q_1|, \dots, |P_k - Q_k|)$ , and  $v^\top = (1, \dots, 1)$  we have

$$\underbrace{\sum_i (P_i - Q_i)^2}_{\langle u, u \rangle} \underbrace{k}_{\langle v, v \rangle} \geq \left( \sum_i |P_i - Q_i| \right)^2 \geq \epsilon^2.$$

□

Lemma 6.2 might give you pause. Why do we go via an  $\ell_2$  route? Did we not claim a few pages ago that  $\ell_2$  is not a good metric when it comes to large alphabet sizes? Indeed this is the case. In the current context we assumed that “uniform” means that the support of  $U$  is *all* of the alphabet  $\mathcal{X} = \{1, \dots, k\}$ . If we considered a slightly more general scenario where we allowed  $U$  to have a support that was strictly smaller than  $k$ , as long as all symbols with non-negative weight have equal weight, then we would have to proceed in a different manner.

Before we proceed let us quickly recall some facts about Poisson distributions.

**Lemma 6.3.** *Let  $X$  be a random variable with Poisson distribution  $\text{Poi}(\lambda)$ . Then for  $l \geq 1$*

$$\mathbb{E}[X(X-1)\cdots(X-l+1)] = \lambda^l.$$

*Further, if  $\mu$  is any real number then*

$$\text{Var}((X - \mu)^2 - X) = 2\lambda^2 + 4\lambda(\lambda - \mu)^2.$$

*Proof.* Consider the first statement. The generating function of the Poisson distribution is  $e^{\lambda(x-1)}$ . Taking the derivative with respect to  $x$  and then setting  $x = 1$  gives us the mean since this corresponds to the weighted sum with weight  $i$ . We get  $\lambda e^{\lambda(x-1)}|_{x=1} = \lambda$ . More generally, taking the  $l$ -th derivative of  $e^{\lambda(x-1)}$  and then setting  $x = 1$  gives us  $\lambda^l$  and this corresponds to the stated expression.

Now consider the second statement. Expanding  $\mathbb{E}[(X - \mu)^2 - X]$  and using the previous result we see that  $\mathbb{E}[(X - \mu)^2 - X] = (\mu - \lambda)^2$ . To compute the variance write down the corresponding expected value and expand in terms of  $X(X-1)\cdots(X-l+1)$  for  $l = 0$  up to  $l = 3$ . Use the previous trick to get the result. □

Recall that we need a test statistics. We claim that

$$T(X^n) = \sum_i (T_i(X^n) - \frac{n}{k})^2 - T_i(X^n)$$

is a good candidate.

NOTE: This is perhaps not the best of notation.  $T_i(X^n)$  refers to the count of the symbol  $i$  in the sample  $X^n$ , whereas  $T(X^n)$  refers to the test statistics.

**Lemma 6.4.** *Let  $P, Q \in \Delta_k$ . Let  $N$  be chosen according to  $\text{Poi}(n)$  and let  $X^N$  be  $N$  iid samples according to  $P$ . Then*

$$\mathbb{E}\left[\sum_i (T_i(X^N) - nQ_i)^2 - T_i(X^N)\right] = n^2 \sum_i (P_i - Q_i)^2.$$

*Proof.* Recall that  $T_i(X^N)$  has distribution  $\text{Poi}(nP_i)$ . Therefore,

$$\begin{aligned} \mathbb{E}\left[\sum_i (T_i(X^N) - nQ_i)^2 - T_i(X^N)\right] &= \mathbb{E}\left[\sum_i T_i(X^N)(T_i(X^N) - 1) - 2nT_i(X^N)Q_i + n^2Q_i^2\right] \\ &\stackrel{\text{Lemma 6.3}}{=} \sum_i [n^2P_i^2 - 2nP_iQ_i + n^2Q_i^2] \\ &= \sum_i n^2(P_i - Q_i)^2. \end{aligned}$$

□

Assume now that  $P \in \mathcal{P} = \{U\}$ , i.e.,  $P = U$ . Then Lemma 6.4 tells us that

$$\mathbb{E}[T(X^N)] = 0.$$

But if  $P$  is such that  $\|P - U\|_1 \geq \epsilon$  then

$$\mathbb{E}[T(X^N)] \stackrel{\text{Lemma 6.4}}{=} n^2 \sum_i (P_i - \frac{n}{k})^2 \stackrel{\text{Lemma 6.2}}{\geq} \frac{n^2\epsilon^2}{k}.$$

We are now ready to state the algorithm that has the claimed performance.

1. Obtain  $N \sim \text{Poi}(n)$  iid samples  $X^N$  from  $P$ , where  $P$  is unknown.
2. Output decision according to

$$\begin{cases} \mathcal{P}, & T(X^N) < \tau = \frac{n^2\epsilon^2}{2k}, \\ \mathcal{Q}, & T(X^N) > \tau. \end{cases}$$

**Lemma 6.5.** *Consider the previous algorithm and assume that  $n > \sqrt{80k}/\epsilon^2$ . Then*

1. *if  $P \in \mathcal{P}$  then  $\mathbb{P}\{T(X^N) > \tau\} < 0.1$ ,*
2. *if  $P \in \mathcal{Q}$  then  $\mathbb{P}\{T(X^N) < \tau\} < 0.1$ .*

*Proof.* Let us look at the two cases separately. If  $P = U$  then we know that  $\mathbb{E}[T(X^N)] = 0$ . From Lemma 6.3 with  $\mu = \lambda = n/k$ , and taking into account that the  $T_i(X^N)$  are independent we get

$$\text{Var}(T(X^N)) = k2(n/k)^2 = 2\frac{n^2}{k}.$$

By the Chebyshev inequality

$$\mathbb{P}\{T(X^N) > \underbrace{\sqrt{10\frac{2n^2}{k}}}_{\text{will see in a second where this comes from}}\} \leq \frac{\text{Var}(T(X^N))}{10\frac{2n^2}{k}} = \frac{2n^2}{k} \frac{k}{10 \cdot 2n^2} = \frac{1}{10}.$$

will see in a second where this comes from

Recall that we have  $\tau = \epsilon^2 n^2 / (2k)$ . Therefore, if  $\tau > \sqrt{10 \frac{2n^2}{k}}$  then  $T(X^N)$  will not exceed  $\tau$  with probability 0.9 as desired. This is true if  $n > \sqrt{80k}/\epsilon^2$ .

The second case follows in a similar manner. Recall that if  $P$  is such that  $\|P - U\|_1 \geq \epsilon$  then

$$\mathbb{E}[T(X^N)] \geq \frac{n^2 \epsilon^2}{k}.$$

We skip the remaining details. □

### 6.3 Property Estimation

We now get the last topic. We have seen how to estimate distributions and how to test properties of distributions. Let us now look how we can estimate properties of distributions.

There are plenty of properties that one might be interested in: entropy, support size, or perhaps the mutual information between two densities.

The simplest approach is to use so-called *plug-in* estimators. This means, estimate the distribution and then plug in these estimates into the functional that computes the desired quantity. But the question is if it is really necessary (and optimal) first to learn the whole distribution if all that we are interested in is one number.

#### 6.3.1 Entropy Estimation

The set-up is very similar than what we used for in the distribution estimation scenario. We have an alphabet  $\mathcal{X} = \{1, \dots, k\}$ . We get iid samples  $X^n = X_1, \dots, X_n$  that are drawn according to a fixed but unknown distribution  $p$ .

We are given a functional  $f(p)$ ,

$$f(p) = \sum_{i=1}^k f(p_i).$$

E.g., if we are interested in the entropy then we want to compute  $f(p) = \sum_i p_i \log_2 \frac{1}{p_i}$ .

Our aim therefore is to design a functional  $\hat{f} : \mathcal{X}^n \rightarrow \mathbb{R}$  that is best in our usual min-max sense,

$$\min_{\hat{f}} \max_{p \in \Delta_k} \mathbb{E}[(f(p) - \hat{f}(X^n))^2]$$

As always one of the most natural such estimators is to use the *empirical one*

$$\hat{f}^{\text{emp}}(X^n) = \sum_i f\left(\frac{T_i(X^n)}{n}\right).$$

### 6.3.2 Symmetric Properties

Many of the properties that one is interested are *symmetric*. This means that they only depend on the set of probabilities but are invariant to a permutation to these probabilities. E.g., the entropy has this property. In the sequel we will limit ourselves to estimating such symmetric properties.

### 6.3.3 Profiles and Natural Estimators

When we talked about the competitive framework for distribution estimation we encountered the notion of a natural genie who knew the distribution but was forced to assign the same probability to symbols that were seen the same number of times.

We will use a similar notion here. We will insist that the estimator for a symmetric property only depends on the *profile* of the sequence we receive. What is this profile? It is a multi-set of multiplicities. More formally we have

$$\phi(X^n) = \{T_i(X^n) : 1 \leq i \leq k\}.$$

**Example 6.2.**  $\Phi(1, 1, 2) = \Phi(1, 2, 1) = \Phi(2, 2, 1) = \{1, 2\}$ .

Assume that the underlying distribution  $P$  is known. For a sequence  $x^n$  let  $P(x^n)$  denote the probability of this sequence. Then, given a profile  $\Phi$  we can compute the *profile probability*

$$p(\Phi) = \sum_{x^n: \Phi(x^n) = \Phi} P(x^n)$$

**Example 6.3.** Let  $\mathcal{X} = \{1, 2, 3\}$ . Let  $P_1 = \frac{1}{2}, P_2 = P_3 = \frac{1}{4}$ . Let  $n = 3$  and consider the profile  $\{1, 2\}$ . Then

$$\begin{aligned} p(\{1, 2\}) &= P(1, 2, 2) + P(1, 3, 3) + P(2, 3, 3) + P(1, 1, 2) + P(1, 1, 3) + P(2, 2, 3) \\ &= 1 - P(1, 2, 3). \end{aligned}$$

## 6.4 Problems

**Problem 6.1** ( $\ell_1$  versus Total Variation). In class we defined the  $\ell_1$  distance as

$$\|p - q\|_1 = \sum_{i=1}^k |p_i - q_i|.$$

Another important distance is the total variation distance  $d_{\text{TV}}(p, q)$ . It is defined as

$$d_{\text{TV}}(p, q) = \max_{S \subseteq \{1, \dots, k\}} \left| \sum_{i \in S} (p_i - q_i) \right|.$$

Show that  $d_{\text{TV}}(p, q) = \frac{1}{2} \|p - q\|_1$ .

**Problem 6.2** (Poisson Sampling). Assume that we have given a distribution  $p$  on  $\mathcal{X} = \{1, \dots, k\}$ . Let  $X^n$  denote a sequence of  $n$  iid samples. Let  $T_i = T_i(X^n)$  be the number of times symbol  $i$  appears in  $X^n$ . Then

$$\mathbb{P}\{T_i = t_i\} = \binom{n}{t_i} p_i^{t_i} (1 - p_i)^{n-t_i}.$$

Note that the random variables  $T_i$  are *dependent*, since  $\sum_i T_i = n$ . This dependence can sometimes be inconvenient.

There is a convenient way of getting around this problem. This is called *Poisson sampling*. Let  $N$  be a random variable distributed according to a Poisson distribution with mean  $n$ . Let  $X^N$  be then an iid sequence of  $N$  variables distributed according to  $p$ .

Show that

- $T_i(X^N)$  is distributed according to a Poisson random variable with mean  $p_i n$ .
- The  $T_i(X^N)$  are independent.
- Conditioned on  $N = n$ , the induced distribution of the Poisson sampling scheme is equal to the distribution of the *original scheme*.

**Problem 6.3** (Add- $\beta$  Estimator). The add- $\beta$  estimator  $q_{+\beta}$  over  $[k]$ , assigns to symbol  $i$  a probability proportional to its number of occurrences plus  $\beta$ , namely,

$$q_i \stackrel{\text{def}}{=} q_i(X^n) \stackrel{\text{def}}{=} q_{+\beta,i}(X^n) \stackrel{\text{def}}{=} \frac{T_i + \beta}{n + k\beta}$$

where  $T_i \stackrel{\text{def}}{=} T_i(X^n) \stackrel{\text{def}}{=} \sum_{j=1}^n \mathbf{1}(X_j = i)$ . Prove that for all  $k \geq 2$  and  $n \geq 1$ ,

$$\min_{\beta \geq 0} r_{k,n}^{\ell_2^2}(q_{+\beta}) = r_{k,n}^{\ell_2^2}(q_{+\sqrt{n}/k}) = \frac{1 - \frac{1}{k}}{(\sqrt{n} + 1)^2}$$

Furthermore,  $q_{+\sqrt{n}/k}$  has the same expected loss for every distribution  $p \in \Delta_k$ .

**Problem 6.4** (Uniformity Testing). Let us reconsider the problem of testing against uniformity. In the lecture we saw a particular *test statistics* that required only  $O(\sqrt{k}/\epsilon^2)$  samples where  $\epsilon$  was the  $\ell_1$  distance.

Let us now derive a test from scratch. To make things simple let us consider the  $\ell_2^2$  distance. Recall that the alphabet is  $\mathcal{X} = \{1, \dots, k\}$ , where  $k$  is known. Let  $U$  be the uniform distribution on  $\mathcal{X}$ , i.e.,  $u_i = 1/k$ . Let  $P$  be a given distribution with components  $p_i$ . Let  $X^n$  be a set of  $n$  iid samples. A pair of samples  $(X_i, X_j)$ ,  $i \neq j$ , is said to *collide* if  $X_i = X_j$ , if they take on the same value.

1. Show that the expected number of collisions is equal to  $\binom{n}{2} \|p\|_2^2$ .
2. Show that the uniform distribution minimizes this quantity and compute this minimum.



3. Show that  $\|p - u\|_2^2 = \|p\|_2^2 - \frac{1}{k}$ .

*NOTE:* In words, if we want to distinguish between the uniform distribution and distributions  $P$  that have an  $\ell_2^2$  distance from  $U$  of at least  $\epsilon$ , then this implies that for those distributions  $\|p\|_2^2 \geq 1/k + \epsilon$ . Together with the first point this suggests the following test: compute the number of collisions in a sample and compare it to  $\binom{n}{2}(1/k + \epsilon/2)$ . If it is below this threshold decide on the uniform one. What remains is to compute the variance of the collision number as a function of the sample size. This will tell us how many samples we need in order for the test to be reliable.

4. Let  $a = \sum_i p_i^2$  and  $b = \sum_i p_i^3$ . Show that the variance of the collision number is equal to

$$\begin{aligned} & \binom{n}{2}a + \binom{n}{2} \left[ \binom{n}{2} - \left( 1 + \binom{n-2}{2} \right) \right] b + \binom{n}{2} \binom{n-2}{2} a^2 - \binom{n}{2}^2 a^2 \\ &= \binom{n}{2} [b2(n-2) + a(1 + a(3-2n))] \end{aligned}$$

by giving an interpretation of each of the terms in the above sum.

*NOTE:* If you don't have sufficient time, skip this step and go to the last point.

For the uniform distribution this is equal to

$$\binom{n}{2} \frac{(k-1)(2n-3)}{k^2} \leq \frac{n^2}{2k}.$$

*NOTE:* You don't have to derive this from the previous result. Just assume it.

5. Recall that we are considering the  $\ell_2^2$  distance which becomes generically small when  $k$  is large. Therefore, the proper scale to consider is  $\epsilon = \kappa/k$ . Use the Chebyshev inequality and conclude that if we have  $\Theta(\sqrt{k}/\kappa^2)$  samples then with high probability the empirical number of collisions will be less than  $\binom{n}{2}(1/k + \kappa/(2k))$  assuming that we get samples from a uniform distribution.

*NOTE:* The second part, namely verifying that the number of collisions is with high probability no smaller than  $\binom{n}{2}(1/k + \kappa/(2k))$  when we get  $\Theta(\sqrt{k}/\kappa^2)$  samples from a distribution with  $\ell_2^2$  distance at least  $\kappa/k$  away from a uniform distribution follows in a similar way.

*HINT:* Note that if  $p$  represents a vector with components  $p_i$  then  $\|p\|_1 = \sum_i |p_i|$  and  $\|p\|_2^2 = \sum_i p_i^2$ .



## Chapter 7

# Information Measures and Generalization Error

### 7.1 Exploration Bias and Information Measures

**Example 7.1.** *It's bad to request too many unnecessary medical tests.* You may have heard this advice before, and we can show why this is indeed sound advice. Suppose we run 15 (unnecessary) medical tests on a healthy patient, each corresponding to the detection of a certain disease. Suppose the tests are accurate so that the rate of false positives is only 2%, i.e.,  $\Pr(\text{test } i \text{ positive} \mid \text{no disease}) = 0.02$ . What is the chance that all 15 tests are negative given that the patient doesn't have any disease?  $\Pr(\text{all negative} \mid \text{no diseases}) = 1 - (1 - 0.02)^{15} \approx 0.26$ .

In the following section, we investigate the problem of “exploration bias”, which arises in data analysis. That is, suppose we collect a large amount of data and perform many measurements/tests on it. Based on these measurements, we decide to report (seemingly) interesting/significant features of the data. In this scenario, not only does the value of the measurement depend on the data (evidently), but also *the choice of which measurement to report depends on the data*. This can bias our results, i.e., if we repeat the measurement on a fresh set of data, the results could differ significantly.

#### 7.1.1 Definitions and Problem Statement

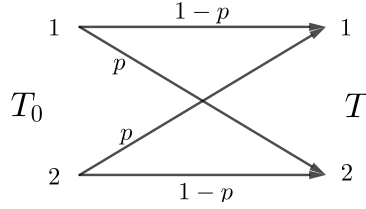
Let  $\mathcal{X}$  denote the sample space, and  $D \in \mathcal{X}^n$  denote the data set. Let  $\phi_i(D)$ ,  $i \in \{1, 2, \dots, m\}$ , denote hypothesis tests' statistics, indexed by  $i$ . (Since  $D$  is random, so is  $\phi_i(D)$ .) The true mean is  $\mu_i = \mathbb{E}[\phi_i(D)]$ , where the expectation is over the randomness of the dataset. On a particular dataset  $D$ , if  $T(D) = i$  is the selected test, the output of data exploration is the value  $\phi_{T(D)}(D)$ . The reported value is thus  $\mathbb{E}[\phi_{T(D)}(D)]$ , resulting in a bias of  $\mathbb{E}[\phi_{T(D)}(D)] - \mu_{T(D)}$ . Hence, we would like to bound

$$|\mathbb{E}[\phi_{T(D)}(D)] - \mu_{T(D)}|. \quad (7.1)$$

Note that  $T$  is not necessarily a deterministic function of  $D$ , rather there exists a conditional distribution  $P_{T|D}$ . In the remainder, we will assume that  $T$  is chosen based on the measurements  $\phi = (\phi_1, \phi_2, \dots, \phi_m)$  and suppress  $D$  in the notation. That is, we can rewrite (7.1) as

$$|\mathbb{E}[\phi_T - \mu_T]|. \quad (7.2)$$

**Example 7.2.** Let  $\phi_1$  and  $\phi_2 \sim \mathcal{N}(\mu, \sigma^2)$  i.i.d. Let  $T_0 = \operatorname{argmax}_{i \in \{1,2\}} \phi_i$  and generate  $T$  as follows:  $T = \begin{cases} T_0, & \text{with probability } 1-p \\ 3-T_0, & \text{with probability } p \end{cases}$  for some  $p \in [0, 1]$ . Now to compute the



exploration bias, note that  $\mathbb{E}[\mu_T] = \mu$ . On the other hand,

$$\begin{aligned} \mathbb{E}[\phi_T] &= \mathbf{Pr}(T = T_0) \mathbb{E}[\phi_{T_0}] + \mathbf{Pr}(T = 3 - T_0) \mathbb{E}[\phi_{(3-T_0)}] \\ &= (1-p) \mathbb{E}[\max\{\phi_1, \phi_2\}] + p \mathbb{E}[\min\{\phi_1, \phi_2\}]. \end{aligned}$$

Now let  $S = \phi_1 + \phi_2$  and  $\Delta = \phi_1 - \phi_2$ . It is straightforward to check that  $S \sim \mathcal{N}(2\mu, 2\sigma^2)$ ,  $\Delta \sim \mathcal{N}(0, 2\sigma^2)$ ,  $\max\{\phi_1, \phi_2\} = \frac{S+|\Delta|}{2}$ , and  $\min\{\phi_1, \phi_2\} = \frac{S-|\Delta|}{2}$ . Then,

$$\begin{aligned} \mathbb{E}[\phi_T] &= \frac{1}{2} \left( (1-p) \mathbb{E}[S + |\Delta|] + p \mathbb{E}[S - |\Delta|] \right) \\ &= \frac{1}{2} \left( \mathbb{E}[S] + (1-2p) \mathbb{E}[|\Delta|] \right) \\ &= \frac{1}{2} \left( 2\mu + (1-2p) \sqrt{\frac{4\sigma^2}{\pi}} \right) \\ &= \mu + (1-2p) \sigma \sqrt{\frac{1}{\pi}}. \end{aligned}$$

Hence, the exploration bias is given by

$$|\mathbb{E}[\phi_T] - \mathbb{E}[\mu_T]| = |1-2p| \sigma \sqrt{\frac{1}{\pi}}.$$

Note that for  $p = \frac{1}{2}$ , the bias is zero. Indeed, for  $p = 1/2$ ,  $T$  is independent of  $(\phi_1, \phi_2)$ ; hence the index does not depend on the data, so we are not introducing any bias. As we decrease  $p$  from  $\frac{1}{2}$  to 0, we “increase the dependence” between  $T_0$  and  $(\phi_1, \phi_2)$ , and the exploration bias increases accordingly.

As we saw in the above example, the exploration bias depends on the degree to which  $T$  depends on  $\phi$ . Hence, we will use dependence measures to find good bounds on the bias.

We can rewrite (7.2) in a more abstract way. In particular, suppose we have an alphabet  $\mathcal{Z}$ , two distributions on  $\mathcal{Z}$  denoted by  $P$  and  $Q$ , and a function  $f : \mathcal{Z} \rightarrow \mathbb{R}$ . We want to bound

$$|\mathbb{E}_P[f(Z)] - \mathbb{E}_Q[f(Z)]| \quad (7.3)$$

In the above setup  $\mathcal{Z} = \mathbb{R}^m \times \{1, 2, \dots, m\}$  (where  $\mathbb{R}^m$  and  $\{1, 2, \dots, m\}$  represent the sets in which  $\phi$  and  $T$  live, respectively),  $f(\phi, t) = \phi_t$ ,  $P = P_{\phi T}$  (i.e., the joint distribution of  $T$  and  $\phi$ ), and  $Q = P_{\phi} P_T$  (i.e., the product of the marginals of  $T$  and  $\phi$ ).

### 7.1.2 $L_1$ -Distance Bound

We have already seen a result somewhat similar to a bound on (7.3). In particular,

**Lemma 7.1.** *Let  $P$  and  $Q$  be two probability mass functions on a finite set  $\mathcal{Z}$ . Then,*

$$\|P - Q\|_1 = 2 \max_{S \subseteq \mathcal{Z}} P(S) - Q(S).$$

Note that for any subset  $S$ ,  $P(S)$  can also be seen as  $\mathbb{E}_P[f_S(Z)]$  where  $f_S(z) = \begin{cases} 1, & z \in S, \\ 0, & z \notin S. \end{cases}$

Moreover, the proof of Lemma 7.1 can be simply modified to show the following:

**Lemma 7.2.**

$$\|P - Q\|_1 = 2 \max_{f: \mathcal{Z} \rightarrow [0,1]} \mathbb{E}_P[f(Z)] - \mathbb{E}_Q[f(Z)].$$

*Proof:* Let  $A = \{z \in \mathcal{Z} : P(z) \geq Q(z)\}$ . Then,

$$\begin{aligned} \mathbb{E}_P[f(Z)] - \mathbb{E}_Q[f(Z)] &= \sum_{z \in A} f(z) (P(z) - Q(z)) + \sum_{z \notin A} f(z) (P(z) - Q(z)) \\ &\leq \sum_{z \in A} (P(z) - Q(z)) \\ &= \frac{\|P - Q\|_1}{2}. \end{aligned}$$

Equality can be achieved if we choose  $f(z) = \begin{cases} 1, & z \in A, \\ 0, & z \notin A. \end{cases}$  ■

**Remark.** The form of the equality in Lemma 7.2 is called the *variational representation* of  $L_1$ -distance. More generally, we can represent any convex function (such as the  $L_1$ -distance) as the supremum of affine functions.

We now have a bound on (7.3) for bounded  $f$ :

**Corollary 7.3.** *For any distributions  $P$  and  $Q$  of a finite set  $\mathcal{Z}$ , and any function  $f : \mathcal{Z} \rightarrow [0, 1]$ , we have*

$$|\mathbb{E}_P[f(Z)] - \mathbb{E}_Q[f(Z)]| \leq \frac{\|P - Q\|_1}{2}.$$

**Exercise 7.1.** *The statement of Lemma 7.2 does not include the absolute value. Verify that the corollary follows by applying Lemma 7.2 twice: once using  $f$ , and another using  $g = 1 - f$ .*

As noted earlier, our initial setup corresponds to choosing  $P$  to be some joint distribution  $P_{XY}$ , and  $Q$  to be the product of the marginals  $P_X P_Y$ . Then, the closer  $P_{XY}$  is to  $P_X P_Y$ , the closer they are to independence (i.e., the less  $Y$  depends on  $X$ ), which makes the exploration bias smaller, as captured in the corollary.

One disadvantage of the above bound is that it is restricted to bounded functions. And as noted in the remark, the main property that allowed us to derive such bound is the convexity of the  $L_1$ -distance. Hence, we can derive similar bounds using other convex dependence measures. In particular, we will turn to the KL divergence.

### 7.1.3 Mutual Information Bound

The following lemma is called the Donsker-Varadhan variational representation of KL divergence.

**Lemma 7.4.** *Let  $p$  and  $q$  be two probability density functions. Then,*

$$D(p||q) = \sup_{\substack{f: \mathbb{R} \rightarrow \mathbb{R} \\ \mathbb{E}_q[e^f] < +\infty}} \left\{ \mathbb{E}_p[f(Z)] - \log \mathbb{E}_q \left[ e^{f(Z)} \right] \right\}$$

**Remark.**  $\log$  is taken to the base  $e$  (both in the computation of  $D(p||q)$  and in the right-hand side).

*Proof:* 1) Suppose  $D(p||q) < +\infty$  and consider any function  $f$ . Then,

$$\begin{aligned} \mathbb{E}_p[f(Z)] - \log \mathbb{E}_q \left[ e^{f(Z)} \right] &= \mathbb{E}_p \left[ \log e^{f(Z)} \right] - \log \mathbb{E}_q \left[ e^{f(Z)} \right] \\ &= \mathbb{E}_p \left[ \log \frac{e^{f(Z)}}{\mathbb{E}_q \left[ e^{f(Z)} \right]} \right] \\ &= \mathbb{E}_p \left[ \log \left( \frac{e^{f(Z)}}{\mathbb{E}_q \left[ e^{f(Z)} \right]} \frac{p}{q} \frac{q}{p} \right) \right] \\ &= D(p||q) + \mathbb{E}_p \left[ \log \left( \frac{q e^{f(Z)}}{p \mathbb{E}_q \left[ e^{f(Z)} \right]} \right) \right] \\ &= D(p||q) - \mathbb{E}_p \left[ \log \frac{p}{q \frac{e^{f(Z)}}{\mathbb{E}_q \left[ e^{f(Z)} \right]}} \right] \end{aligned}$$

Now note that  $\int q \frac{e^{f(Z)}}{\mathbb{E}_q[e^{f(Z)}]} dz = \frac{1}{\mathbb{E}_q[e^{f(Z)}]} \int q e^{f(Z)} dz = 1$ . Hence, the second term is also a KL divergence. Then,

$$\mathbb{E}_p[f(Z)] - \log \mathbb{E}_q[e^{f(Z)}] = D(p||q) - D(p||\tilde{q}) \leq D(p||q).$$

Equality is achieved if we set  $f = \log \frac{p}{q}$ .

2) Suppose  $D(p||q) = +\infty$ . Then, we need to show that the supremum is also  $+\infty$ .  $D(p||q) = +\infty$  implies that there exists a set  $A$  such that  $p(A) > 0$  and  $q(A) = 0$ . Choose  $f = \lambda \mathbb{1}\{z \in A\}$ . Then,  $\mathbb{E}_p[f(Z)] - \log \mathbb{E}_q[e^{f(Z)}] = \lambda p(A)$ . Taking  $\lambda \rightarrow +\infty$  yields the result. ■

Before introducing the main bound, we need to introduce the concept of subgaussian distributions.

**Definition 7.1.**  $Z$  is called  $\sigma^2$ -subgaussian if  $\log \mathbb{E}[e^{\lambda Z}] \leq \frac{\lambda^2 \sigma^2}{2}$  for all  $\lambda \in \mathbb{R}$ .

**Lemma 7.5.** If  $Z$  is  $\sigma^2$ -subgaussian, then  $\mathbb{E}[Z] = 0$  and  $\mathbb{E}[Z^2] \leq \sigma^2$ .

*Proof:* Pick  $\lambda$  such that  $\lambda \mathbb{E}[Z] > 0$ . Using the Taylor expansion of the exponential, we get

$$1 + \frac{1}{2} \lambda^2 \sigma^2 + O(\lambda^4) \geq e^{\lambda^2 \sigma^2 / 2} \geq \mathbb{E}[e^{\lambda Z}] \geq 1 + \lambda \mathbb{E}[Z] + \frac{1}{2} \lambda^2 \mathbb{E}[Z^2] + O(\lambda^3).$$

The result follows by taking  $\lambda \rightarrow 0$ . ■

**Exercise 7.2.** For  $Z \sim \mathcal{N}(0, \sigma^2)$ , show that  $\mathbb{E}[e^{\lambda Z}] = e^{\lambda^2 \sigma^2 / 2}$ . Hence,  $Z$  is  $\sigma^2$ -subgaussian.

We are now ready to prove the main bound for this section on the exploration bias  $|\mathbb{E}[\phi_T - \mu_T]|$ , where  $\phi = (\phi_1, \phi_2, \dots, \phi_m)$  and  $T \in \{1, 2, \dots, m\}$ .

**Theorem 7.6.** Suppose for each  $i \in \{1, 2, \dots, m\}$ ,  $\phi_i - \mu_i$  is  $\sigma^2$ -subgaussian. Then,

$$|\mathbb{E}[\phi_T - \mu_T]| \leq \sigma \sqrt{2I(T; \phi)}.$$

**Remark.** As expected, if  $T$  is independent of  $\phi$ , then the exploration bias is zero. If  $T$  does not depend “too much” on  $\phi$ , as captured by mutual information, then we can guarantee a small bias.

*Proof:* Fix  $T = i$ . We will use Lemma 7.4 with the distributions  $P_{\phi_i|T=i}$  (the distribution of  $\phi_i$  conditioned on the choice of  $T$  being  $i$ ) and  $P_{\phi_i}$  (the prior distribution of  $\phi_i$ ). For some  $\lambda \in \mathbb{R}$ , let  $f = \lambda(\phi_i - \mu_i)$ . Then, it follows from Lemma 7.4 that

$$\begin{aligned} D(P_{\phi_i|T=i}||P_{\phi_i}) &\geq \lambda \left( \mathbb{E}_{P_{\phi_i|T=i}}[\phi_i] - \mu_i \right) - \log \mathbb{E}_{P_{\phi_i}} \left[ e^{\lambda(\phi_i - \mu_i)} \right] \\ &\geq \lambda \left( \mathbb{E}_{P_{\phi_i|T=i}}[\phi_i] - \mu_i \right) - \lambda^2 \sigma^2 / 2, \end{aligned}$$

where the second inequality follows from the  $\sigma^2$ -subgaussianity assumption. Since  $\lambda$  was arbitrary, we get

$$\begin{aligned} D(P_{\phi_i|T=i}||P_{\phi_i}) &\geq \sup_{\lambda \in \mathbb{R}} \left\{ \lambda \left( \mathbb{E}_{P_{\phi_i|T=i}}[\phi_i] - \mu_i \right) - \lambda^2 \sigma^2 / 2 \right\} \\ &= \frac{\left( \mathbb{E}_{P_{\phi_i|T=i}}[\phi_i] - \mu_i \right)^2}{2\sigma^2}. \end{aligned}$$

Hence,

$$\left| \mathbb{E}_{P_{\phi_i|T=i}}[\phi_i] - \mu_i \right| \leq \sigma \sqrt{2D(P_{\phi_i|T=i}||P_{\phi_i})}.$$

Finally,

$$\begin{aligned} |\mathbb{E}[\phi_T - \mu_T]| &= \left| \sum_{i=1}^m \Pr(T=i) \mathbb{E}_{P_{\phi_i|T=i}}[\phi_i] - \mu_i \right| \\ &\stackrel{(a)}{\leq} \sum_{i=1}^m \Pr(T=i) \left| \mathbb{E}_{P_{\phi_i|T=i}}[\phi_i] - \mu_i \right| \\ &\leq \sum_{i=1}^m \Pr(T=i) \sigma \sqrt{2D(P_{\phi_i|T=i}||P_{\phi_i})} \\ &\stackrel{(b)}{\leq} \sum_{i=1}^m \Pr(T=i) \sigma \sqrt{2D(P_{\phi|T=i}||P_{\phi})} \\ &\stackrel{(c)}{\leq} \sigma \sqrt{2 \sum_{i=1}^m \Pr(T=i) D(P_{\phi|T=i}||P_{\phi})} \\ &= \sigma \sqrt{2D(P_{\phi T}||P_{\phi} P_T)} = \sigma \sqrt{2I(T; \phi)}, \end{aligned}$$

where (a) and (c) follow from Jensen's inequality, and (b) follows from the data processing inequality.  $\blacksquare$

**Exercise 7.3.** Show that, if  $\phi_i - \mu_i$  is  $\sigma_i^2$ -subgaussian for each  $i \in \{1, 2, \dots, m\}$ , then

$$|\mathbb{E}[\phi_T - \mu_T]| \leq \sqrt{\mathbb{E}[\sigma_T^2]} \sqrt{2I(T; \phi)}.$$

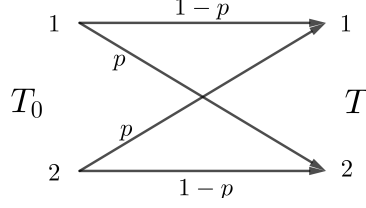
Let's revisit the initial example:

**Example 7.3.** Let  $\phi_1$  and  $\phi_2 \sim \mathcal{N}(\mu, \sigma^2)$  i.i.d. Let  $T_0 = \operatorname{argmax}_{i \in \{1, 2\}} \phi_i$  and generate  $T$  as follows:  $T = \begin{cases} T_0, & \text{with probability } 1-p \\ 3 - T_0, & \text{with probability } p \end{cases}$  for some  $p \in [0, 1]$ .

Since  $\phi_i - \mu_i \sim \mathcal{N}(0, \sigma^2)$ , it is  $\sigma^2$ -subgaussian, thus satisfying the assumption of Theorem 7.6. To compute  $I(T; \phi)$ :

$$I(T; \phi) = H(T) - H(T|\phi).$$





Since  $\phi_1$  and  $\phi_2$  are i.i.d, then  $\Pr(T_0 = 1) = \Pr(T_0 = 2) = \frac{1}{2}$ . Hence,  $H(T) = \log 2$ . Since both  $\phi - T_0 - T$  and  $T_0 - \phi - T$  are Markov chains, we get  $H(T|\phi, T_0) = H(T|\phi)$  and  $H(T|\phi, T_0) = H(T|T_0)$ . Hence,

$$I(T; \phi) = H(T) - H(T|\phi) = \log 2 - H(T|T_0) = \log 2 - H(p).$$

Hence, by the above theorem,

$$|\mathbb{E}[\phi_T - \mu_T]| \leq \sigma \sqrt{2(\log 2 - H(p))}.$$

**Example 7.4.** Suppose  $\phi_i \sim \mathcal{N}(0, \sigma^2)$  i.i.d. for  $i \in \{1, 2, \dots, m\}$ , and  $T = \operatorname{argmax}_i \phi_i$ . Then,

$$I(T; \phi) = H(T) = \log m,$$

and

$$\mathbb{E}[\max\{\phi_1, \phi_2, \dots, \phi_m\}] \leq \sigma \sqrt{2 \log m}.$$

## 7.2 Information Measures and Generalization Error

In the previous section, we studied the exploration bias. Herein, we will consider a related concept that arises in machine learning, called *generalization error*. Roughly speaking, it tries to capture the idea of model “overfitting”.

### 7.2.1 Setup and Problem Statement

The standard setup of statistical learning theory is as follows: we have an instance space  $\mathcal{X}$ , a hypothesis space  $\mathcal{W}$ , and a loss function  $\ell : \mathcal{W} \times \mathcal{X} \rightarrow \mathbb{R}_+$ .

We observe  $D = (X_1, X_2, \dots, X_n)$  i.i.d samples from some *unknown* distribution  $P$ . Using these samples, we want to pick a hypothesis in  $\mathcal{W}$ . We can think of the hypotheses as models we use to explain our data, and we use the loss function to evaluate the performance of our chosen model.

**Example 7.5.**  $\mathcal{X} = \mathbb{R} \times \mathbb{R}$ ,  $\mathcal{W} = \{\text{affine functions from } \mathbb{R} \text{ to } \mathbb{R}\}$ , and  $\ell(w, (x_1, x_2)) = (x_2 - w(x_1))^2$ . That is, we observe pairs of values  $\{(x_1, x_2)\}_{i=1}^n$ , and we want to find the best linear approximation of  $x_2$  in terms of  $x_1$ . This setup with the choice of the squared error loss is referred to as linear regression.

The algorithm to choose a hypothesis  $w \in \mathcal{W}$  based on  $D$  does not need to be deterministic. Therefore, we model the *learning algorithm* as a conditional distribution  $P_{W|D}$ .

**Definition 7.2.** Given  $w \in \mathcal{W}$ , the *population risk* of  $w$  is defined as:

$$L_P(w) := \mathbb{E}_P[\ell(w, X)].$$

The population risk indicates how well the hypothesis  $w$  models the data. We would like this risk to be small, but we cannot evaluate directly since  $P$  is unknown. On the other hand, given a data set  $D$ , we can evaluate the empirical risk.

**Definition 7.3.** Given a data set  $D = (X_1, X_2, \dots, X_n)$  and a hypothesis  $w$ , the *empirical risk* is defined as:

$$L_D(w) = \frac{1}{n} \sum_{i=1}^n \ell(w, X_i).$$

**Definition 7.4.** Given a learning algorithm  $P_{W|D}$ , the *generalization error* is defined as

$$\text{gen}(P, P_{W|D}) = \mathbb{E}_{P_{DW}}[L_P(W) - L_D(W)].$$

Thus, the generalization error is the difference between the performance of the chosen hypothesis on the given data versus its performance on a fresh data sample (given by the population risk), averaged over the choice of the hypothesis. Note that

$$\mathbb{E}_{P_D P_W}[L_D(W)] = \mathbb{E}_{P_W}[L_P(W)],$$

however the expectation is taken with respect to the joint  $P_{DW}$  in the definition of the generalization error.

Similarly to the exploration bias setup, if the chosen hypothesis “depends too much” on the given data, then the generalization error can be large, i.e., we are “overfitting”. We can bound the error by controlling the degree of dependence, and we will use again information measures to do so. Note that, in this setting, if  $W$  is chosen independently from the data, then the generalization error will be zero but the population risk will be large, which we ultimately want to be small. So we have a tension where on the one hand, if  $W$  is independent from  $D$  so the learning algorithm is not “learning” anything; and on the other hand, if  $W$  depends too much on the data, it will be overfitting.

### 7.2.2 Mutual Information Bound

**Theorem 7.7.** *If for all  $w \in \mathcal{W}$ ,  $\ell(w, X)$  is  $\sigma^2$ -subgaussian, then*

$$|\text{gen}(P, P_{W|D})| = \sqrt{\frac{2\sigma^2}{n} I(D; W)}.$$

The proof is the same as that of Theorem 7.6. The only “new” component is to show that if  $\ell(w, X)$  is  $\sigma^2$ -subgaussian, then  $\frac{1}{n} \sum_{i=1}^n \ell(w, X_i)$  is  $\frac{\sigma^2}{n}$ -subgaussian, which is left as an exercise.

### 7.2.3 Differential Privacy Bound

On a high level, to get a bound on the errors, we have been trying to control the degree of dependence of the output (the hypothesis) of our algorithm on its input (the data). Unsurprisingly, this notion of controlling the dependence arises in the privacy and security literature.

One such notion is *differential privacy* which tries to capture the following intuition: small changes in the input should lead to small changes in the output. That is, no small change should have a large effect on the output. So we will define “neighboring” data sets to represent small changes in the input.

**Definition 7.5.** We say two datasets  $D = (X_1, X_2, \dots, X_n)$  and  $D' = (X'_1, X'_2, \dots, X'_n)$  are neighbors if they differ only in one entry, i.e.,  $|\{i : X_i \neq X'_i\}| = 1$ .

**Definition 7.6.** Suppose  $\mathcal{W}$  is discrete. Given  $\epsilon > 0$ , we say  $P_{W|D}$  is  $\epsilon$ -differentially private if for any neighboring datasets  $D_1$  and  $D_2$ ,

$$\max_{w \in \mathcal{W}} \frac{P(w|D_1)}{P(w|D_2)} \leq e^\epsilon.$$

If  $\mathcal{W}$  is not discrete, we simply replace pmfs by pdfs, i.e.,

$$\sup_w \frac{p(w|D_1)}{p(w|D_2)} \leq e^\epsilon.$$

This definition is equivalent to the following: for any neighboring datasets  $D_1$  and  $D_2$ , and any subset  $S \subseteq \mathcal{W}$ ,

$$P_{W|D_1}(W \in S) \leq e^\epsilon P_{W|D_2}(W \in S).$$

**Theorem 7.8.** If  $\ell \in [0, 1]$  and  $P_{W|D}$  is  $\epsilon$ -differentially private, then

$$|\text{gen}(P, P_{W|D})| \leq e^\epsilon - 1.$$

## 7.3 Problems

**Problem 7.1** (Exploration Bias). Given a real random variable  $X$  taking values on a finite set  $\mathcal{X} \subset \mathbb{R}$ , define  $\psi(\lambda) = \log \mathbb{E}[e^{\lambda X}]$ . Show that

(a)  $\psi'(\lambda) = \mathbb{E}[X_\lambda]$  where  $\mathbb{E}[X_\lambda]$  is a random variable taking values on  $\mathcal{X}$ , with distribution  $p_\lambda(x) = p(x) \exp(\lambda x) \exp(-\psi(\lambda))$ . Hence  $\psi'(0) = \mathbb{E}[X]$ .

(b)  $\psi''(\lambda) = \text{var}(X_\lambda)$ . Conclude that  $\psi$  is convex.

**Problem 7.2** (Exploration Bias). (a) Let  $X_1, X_2, \dots, X_n \sim \text{i.i.d. } \mathcal{N}(0, 1)$ . Let  $Y = \arg\max_i X_i$  and  $T \in \{1, 2, \dots, n\}$  is such that

$$P_{T|Y}(t|y) = \begin{cases} p, & t = y \\ \frac{1-p}{n-1}, & t \neq y \end{cases} \quad \text{for some } p \in [0, 1]. \quad (7.4)$$

1. Compute  $I(X; T)$  where  $X = (X_1, X_2, \dots, X_n)$ . (Hint: write  $I(X; T) = H(T) - H(T|X)$ . What is the marginal distribution of  $T$ ?)
- (b) Let  $X_1, \dots, X_4 \sim \text{i.i.d. } \mathcal{N}(0, 1)$  and  $X_5 \sim \mathcal{N}(0, 4)$ . Let  $Y$  and  $T$  be as in part (a) with  $p = 0.3$ .

1. Show that  $\Pr(Y = 5) = \int_{-\infty}^{\infty} \frac{1}{\sqrt{8\pi}} (1 - Q(x))^4 e^{-x^2/8} dx$ , where we are using  $Q(x) = \int_x^{\infty} \frac{1}{\sqrt{2\pi}} e^{-u^2/2} du$ , and find a corresponding numerical approximation (using Mathematica, for example).
2. Using the previous numerical approximation, find the marginal distributions  $P_Y$  and  $P_T$ .

**Problem 7.3** (Exploration Bias). Let  $\mathcal{X}$  be the sample space,  $\mathcal{W}$  the hypothesis space, and let  $\ell : \mathcal{W} \times \mathcal{X} \rightarrow \mathbb{R}_+$  be a corresponding loss function. On a dataset  $D = (X_1, X_2, \dots, X_n)$ , the empirical risk for a hypothesis  $w$  is given by  $L_D(w) = \frac{1}{n} \sum_{i=1}^n \ell(w, X_i)$ . We saw in class that  $I(D; W)$  can be used to bound the generalization error. Hence, we can use it as a *regularizer* in empirical risk minimization.

- (a) First, show that given any joint distribution  $P_{XY}$  on  $\mathcal{X} \times \mathcal{Y}$  and marginal distribution  $Q$  on  $\mathcal{Y}$ ,  $D(P_{XY} \| P_X P_Y) \leq D(P_{XY} \| P_X Q)$ .

Since we cannot directly compute  $D(P_{DW} \| P_D P_W)$ , we will use  $D(P_{DW} \| P_D Q)$  as a proxy, where  $Q$  is a distribution on  $\mathcal{W}$ .

- (b) Let

$$P_{W|D}^* = \operatorname{argmin}_{P_{W|D}} \left( \mathbb{E}[L_D(W)] + \frac{1}{\beta} D(P_{DW} \| P_D Q) \right).$$

1. Show that

$$\min_{P_{W|D}} \left( \mathbb{E}[L_D(W)] + \frac{1}{\beta} D(P_{DW} \| P_D Q) \right) = \mathbb{E}_D \left[ \min_{P_{W|D=d}} \left( \mathbb{E}[L_d(W)] + \frac{1}{\beta} D(P_{W|D=d} \| Q) \right) \right].$$

2. Show that the minimizer on the right-hand side  $P_{W|D=d}^*$  is given by

$$P_{W|D=d}^* = \frac{e^{-\beta L_d(w)} Q(w)}{\mathbb{E}_Q [e^{-\beta L_d(W)}]}.$$

This is known in the literature as the Gibbs algorithm. (Hint: Write  $\mathbb{E}[\beta L_d(W)] = \mathbb{E}[\log e^{\beta L_d(W)}]$ , combine with the KL divergence term and use non-negativity of KL divergence.)

3. Show that  $P_{W|D=d}^*$  is  $2\beta/n$ -differential private if  $\ell \in [0, 1]$ .

## Chapter 8

# Elements of Statistical Signal Processing

### 8.1 Optimum Estimation

#### 8.1.1 MMSE Estimation

Consider two (real- or complex-valued) random vectors  $\mathbf{D}$  and  $\mathbf{X}$  with known joint probability density function  $f_{\mathbf{D},\mathbf{X}}$ . Suppose that using only  $\mathbf{X}$ , we are tasked to construct an estimate of  $\mathbf{D}$ . This estimate is thus a function  $g(\mathbf{x})$ , to be selected optimally. A natural criterion is the standard mean-squared error

$$\mathbb{E} \left[ \|\mathbf{D} - g(\mathbf{X})\|^2 \middle| \mathbf{X} = \mathbf{x} \right], \quad (8.1)$$

and the goal is to find a function  $g(\cdot)$  that minimizes this.

This problem has a nice and intuitively pleasing solution:

$$g(\mathbf{x}) = \mathbb{E} [\mathbf{D} | \mathbf{X} = \mathbf{x}], \quad (8.2)$$

and we will use the shorthand notation  $\hat{\mathbf{D}}_{MMSE}(\mathbf{X} = \mathbf{x})$  for this optimal estimator (optimal in the mean-squared error sense).

**Example 8.1** (Gaussian signal and noise). Let  $D$  be a real-valued zero-mean unit-variance Gaussian random variable. Let  $X = D + Z$ , where  $Z$  is a zero-mean Gaussian random variable of variance  $\sigma^2$ . Then, it is straightforward to evaluate

$$\hat{D}_{MMSE}(X = x) = \mathbb{E}[D | X = x] = \int_{-\infty}^{\infty} d f_{D|X}(d|x) dd = \int_{-\infty}^{\infty} d \frac{f_{X|D}(x|d) p_D(d)}{f_X(x)} dd \quad (8.3)$$

$$= \int_{-\infty}^{\infty} d \frac{\frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-d)^2}{2\sigma^2}) \frac{1}{\sqrt{2\pi}} \exp(-\frac{d^2}{2})}{\frac{1}{\sqrt{2\pi(1+\sigma^2)}} \exp(-\frac{x^2}{2(1+\sigma^2)})} dd = \frac{1}{1+\sigma^2} x. \quad (8.4)$$

Moreover, the mean-squared error incurred by this optimum estimator is

$$\mathbb{E} \left[ \left| D - \hat{D}_{MMSE}(X) \right|^2 \right] = \frac{\sigma^2}{1+\sigma^2}. \quad (8.5)$$

### 8.1.2 Linear MMSE Estimation

In a slight variation of the consideration, let us now assume that the estimator must be linear (with fixed coefficients, independent of the data). Let us first consider the case where the desired data  $D$  is scalar. That is, we seek to find

$$\hat{D}_{LMMSE}(\mathbf{X}) = \mathbf{w}^T \mathbf{X}, \quad (8.6)$$

where  $\mathbf{w}$  is a fixed vector of coefficients. In the MMSE perspective, we strive to select this vector such as to minimize

$$\mathbb{E} \left[ \left| D - \hat{D}_{LMMSE}(\mathbf{X}) \right|^2 \right]. \quad (8.7)$$

To express the solution, it is convenient to introduce the notation

$$R_{\mathbf{X}} = \mathbb{E}[\mathbf{X}\mathbf{X}^H] \quad (8.8)$$

for the *covariance matrix* of the data (recall that we are assuming zero-mean signals throughout this Module), and

$$\mathbf{r}_{D\mathbf{X}} = \mathbb{E}[D\mathbf{X}^*] \quad (8.9)$$

for the covariance between the desired and the observed data. With this, the optimal Wiener coefficients are (assuming that the matrix  $R_{\mathbf{X}}$  is invertible — for the more general case, see the homework)

$$\mathbf{w} = R_{\mathbf{X}}^{-1} \mathbf{r}_{D\mathbf{X}}, \quad (8.10)$$

and the corresponding mean-squared error can be expressed as

$$\mathbb{E} \left[ \left| D - \hat{D}_{LMMSE}(\mathbf{X}) \right|^2 \right] = \sigma_D^2 - \mathbf{r}_{D\mathbf{X}}^H R_{\mathbf{X}}^{-1} \mathbf{r}_{D\mathbf{X}}, \quad (8.11)$$

where  $\sigma_D^2$  denotes the variance of the desired data  $D$ .

To prove this, it is instructive to observe that with the optimal coefficient vector  $\mathbf{w}$ , the error must be orthogonal to the observed data, which (here) means that

$$\mathbb{E}[(D - \mathbf{w}^T \mathbf{X})\mathbf{X}^H] = \mathbf{0}^H. \quad (8.12)$$

This can be established either by observing that the objective function is convex and finding the gradient with respect to the coefficient vector, or by observing the Hilbert space structure and invoking the projection theorem. The orthogonality condition can be rewritten as

$$\mathbb{E}[D\mathbf{X}^H] - \mathbf{w}^T \mathbb{E}[\mathbf{X}\mathbf{X}^H] = \mathbf{0}^H. \quad (8.13)$$

Assuming that the matrix  $\mathbb{E}[\mathbf{X}\mathbf{X}^H]$  is invertible, this implies the claimed formula.

The corresponding incurred mean-squared error can be calculated as follows:

$$\begin{aligned} \mathbb{E} \left[ \left\| D - \hat{D}_{LMMSE}(\mathbf{X}) \right\|^2 \right] &= \mathbb{E} \left[ (D - \mathbf{w}^T \mathbf{X})^* (D - \mathbf{w}^T \mathbf{X}) \right] \\ &= \mathbb{E} \left[ (D - \mathbf{w}^T \mathbf{X})^* D \right] - \mathbf{w}^T \underbrace{\mathbb{E} \left[ (D - \mathbf{w}^T \mathbf{X})^* \mathbf{X} \right]}_{=\mathbf{0}, \text{ due to orthogonality}} \\ &= \mathbb{E} [|D|^2] - \mathbf{w}^H \mathbb{E} [\mathbf{X}^* D], \end{aligned} \quad (8.14)$$

and if we plug in the formula for the optimal Wiener solution for  $\mathbf{w}$ , we obtain the claimed formula.

## 8.2 Wiener Filtering, Smoothing, Prediction

The tools of signal processing are often most powerful if we consider (long) sequences of data. That is, we now suppose that we have a time-domain signal  $D[n]$  (where  $n$  ranges over integers), and the observed data is  $X[n]$ . In the world view of signal processing, we would then form an estimate of the form

$$\hat{D}[n] = \sum_{k=-p_0}^{p_1} w[k]X[n-k], \quad (8.15)$$

where  $p_0$  and  $p_1$  are non-negative integers. Defining the vector (of length  $p_0 + p_1 + 1$ )

$$\mathbf{X}[n] = (X[n+p_1], X[n+p_1-1], \dots, X[n], \dots, X[n-p_0+1], X[n-p_0])^T \quad (8.16)$$

and the vector  $\mathbf{w}$  containing the corresponding  $p_0 + p_1 + 1$  filter coefficients, namely,

$$\mathbf{w} = (w[p_1], w[p_1-1], \dots, w[0], \dots, w[-p_0+1], w[-p_0])^T, \quad (8.17)$$

we can express the optimal coefficients as

$$\mathbf{w}^T = \mathbb{E}[D[n]\mathbf{X}[n]^H] (\mathbb{E}[\mathbf{X}[n]\mathbf{X}[n]^H])^{-1}, \quad (8.18)$$

which, in general, depends on  $n$ . This motivates the definition of *wide-sense stationary* random processes that you have encountered in earlier classes. For such processes, we have that

$$\mathbb{E}[X[n]X^*[n-k]] = R_X[k], \quad (8.19)$$

$$\mathbb{E}[D[n]X^*[n-k]] = R_{DX}[k], \quad (8.20)$$

that is, these expectations do not depend on  $n$ , but only on the “lag”  $k$  between the two arguments. With this, it can easily be verified that the above formula for the optimal coefficient vector  $\mathbf{w}$  does not depend on  $n$ .

An equally enlightening but alternative view is to allow  $p_0$  and  $p_1$  to be infinite. In this case, the orthogonality principle stipulates that the optimum filter coefficients must satisfy

$$\mathbb{E} \left[ \left( D[n] - \sum_{k=-\infty}^{\infty} w[k]X[n-k] \right) X^*[n-\ell] \right] = 0, \quad (8.21)$$

for all integers  $\ell$ . Rewriting,

$$\mathbb{E}[D[n]X^*[n-\ell]] - \sum_{k=-\infty}^{\infty} w[k]\mathbb{E}[X[n-k]X^*[n-\ell]] = 0, \quad (8.22)$$

or

$$R_{DX}[\ell] - \sum_{k=-\infty}^{\infty} w[k]R_X[\ell-k] = 0, \quad (8.23)$$

where we observe that the sum is a convolution. This suggests that it may be instructive to take Fourier transforms:

$$S_{DX}(e^{j\omega}) - W(e^{j\omega})S_{XX}(e^{j\omega}) = 0. \quad (8.24)$$

### 8.3 Adaptive Filters

Let us consider the scenario where

$$\hat{D}[n] = \sum_{k=0}^p w[k]X[n-k]. \quad (8.25)$$

Suppose that we pick an arbitrary initial choice of filter coefficients  $\mathbf{w}_0$ . Let us take the perspective that we gradually update these filter coefficients so as to make them better. A classical choice is called *gradient descent*. Here, we consider the gradient of the error (with respect to the filter coefficients), which is easily found to be

$$\nabla_{\mathbf{w}_n} \mathbb{E} \left[ |D[n] - \mathbf{w}_n^T \mathbf{X}[n]|^2 \right] = -2\mathbb{E} \left[ (D[n] - \mathbf{w}_n^T \mathbf{X}[n]) \mathbf{X}^*[n] \right] \quad (8.26)$$

The idea is to take a (“small”) step *against* the gradient, i.e.,

$$\mathbf{w}_{n+1} = \mathbf{w}_n + \mu \mathbb{E} \left[ (D[n] - \mathbf{w}_n^T \mathbf{X}[n]) \mathbf{X}^*[n] \right], \quad (8.27)$$

where the step-size parameter  $\mu$  is to be chosen wisely.

To gain some understanding of what this algorithm does, let us consider the special case where the signals  $D[n]$  and  $X[n]$  are jointly wide-sense stationary, and hence, all expected values above do not depend on  $n$ . For this special case, the update equation becomes

$$\mathbf{w}_{n+1} = \mathbf{w}_n + \mu \mathbb{E} \left[ (D - \mathbf{w}_n^T \mathbf{X}) \mathbf{X}^* \right] \quad (8.28)$$

$$= \mathbf{w}_n + \mu (\mathbf{r}_{D\mathbf{X}} - R_{\mathbf{X}} \mathbf{w}_n) \quad (8.29)$$

If we plug in  $\mathbf{w}_n = R_{\mathbf{X}}^{-1} \mathbf{r}_{D\mathbf{X}}$  (which is the optimal Wiener solution), then the algorithm will not move any further, and thus, will stay at the globally optimal solution, which is a first important sanity check. In more detail, we can also suppose that we start the algorithm with an arbitrary  $\mathbf{w}_0$ . Let us denote the optimal Wiener solution by  $\bar{\mathbf{w}}$ . Then, we can express

$$\begin{aligned} \mathbf{w}_{n+1} - \bar{\mathbf{w}} &= \mathbf{w}_n - \bar{\mathbf{w}} + \mu \mathbf{r}_{D\mathbf{X}} - R_{\mathbf{X}} \mathbf{w}_n \\ &= (I - \mu R_{\mathbf{X}}) (\mathbf{w}_n - \bar{\mathbf{w}}), \end{aligned} \quad (8.30)$$

where, for the last step, we have used the fact that the Wiener solution satisfies  $\mathbf{r}_{D\mathbf{X}} = R_{\mathbf{X}} \bar{\mathbf{w}}$ . It is then instructive to express the matrix in terms of its spectral decomposition  $R_{\mathbf{X}} = U \Lambda U^H$ , leading to  $p+1$  independent recursions. Specifically, defining the notation  $\mathbf{u}_n = U^H (\mathbf{w}_n - \bar{\mathbf{w}})$ , we obtain

$$\begin{aligned} \mathbf{u}_{n+1} &= U^H (\mathbf{w}_{n+1} - \bar{\mathbf{w}}) \\ &= U^H (I - \mu R_{\mathbf{X}}) (\mathbf{w}_n - \bar{\mathbf{w}}) \\ &= U^H (U (I - \mu \Lambda) U^H) (\mathbf{w}_n - \bar{\mathbf{w}}) \\ &= (I - \mu \Lambda) \mathbf{u}_n, \end{aligned} \quad (8.31)$$

and since  $\Lambda$  is a diagonal matrix, each of the  $p+1$  components of the vector  $\mathbf{u}_n$  follows a separate recursion, independently of the others. Clearly, the overall sequence converges if



and only if all  $p+1$  so-called *modes* converge individually. But each one of them is simply an exponential series governed by  $(1 - \mu\lambda_i(R_{\mathbf{X}}))^n$  (times the initial value). Such an exponential series converges (to zero) if and only if  $|1 - \mu\lambda_i(R_{\mathbf{X}})| < 1$ , or, equivalently,  $\lambda_{\max}(R_{\mathbf{X}}) < 2/\mu$ . Here, we are also using the fact that for covariance matrices  $R_{\mathbf{X}}$ , eigenvalues must be non-negative.

To make the algorithm useful, we cannot use the shape given in Equation (8.27) since in any realistic scenario, we would not know the involved expected values. Instead, we may estimate them from the data. In the extreme case, we could estimate the expectation from just a single sample, hence use  $\mathbb{E}[(D[n] - \mathbf{w}_n^T \mathbf{X}[n]) \mathbf{X}^*[n]] \approx (D[n] - \mathbf{w}_n^T \mathbf{X}[n]) \mathbf{X}^*[n]$ . Of course, this should not be expected to be a particularly good estimate of said expectation, but in exchange, we can actually calculate this simply based on the data at hand (at least for all times  $n$  for which we have “training data,” i.e., where we know the true desired outcome  $D[n]$ ). This is called the LMS adaptive algorithm and was discovered in 1960 [3]. It is thus characterized by the update equation

$$\mathbf{w}_{n+1} = \mathbf{w}_n + \mu (D[n] - \mathbf{w}_n^T \mathbf{X}[n]) \mathbf{X}^*[n]. \quad (8.32)$$

The full analysis of the convergence of this algorithm, even for the special case of wide-sense stationary data, is not feasible in closed form. It is important to observe that here, the filter coefficients  $\mathbf{w}_n$  are random vectors, induced by the data. Hence, a first order of business would be to prove convergence of the mean  $\mathbb{E}[\mathbf{w}_n]$ , perhaps starting with the case of wide-sense stationary data. Unfortunately, it is quickly seen that the resulting vector sequence  $\{\mathbb{E}[\mathbf{w}_n]\}_{n \geq 0}$  depends on higher-order statistics of the data and is thus out of reach. A common alternative (approximate) consideration is to assume that the filter taps  $\mathbf{w}_n$  are (*statistically*) *independent* of the data vector corresponding to the same time slot,  $\mathbf{X}[n]$ . While this is not exactly true, it may hold approximately if  $\mu$  is sufficiently small. Under this so-called *independence assumption* for the LMS, one easily finds that convergence is again determined by the vector sequence from Equation (8.30), and thus, by the eigenvalues of the covariance matrix of the observed data. A well written account on adaptive filters can be found, e.g., in [4], and an exhaustive compendium in [5].

## 8.4 Problems

**Problem 8.1** (MMSE Estimation). Consider the scenario where  $f_{X|D}(x|d) = de^{-dx}$ , for  $x \geq 0$  (and zero otherwise), that is, the observed data  $x$  is distributed according to an exponential with mean  $1/d$ . Moreover, the desired variable  $d$  itself is also exponentially distributed, with mean  $1/\mu$ .

(a) Find the MMSE estimator of  $d$  given  $x$ , and calculate the corresponding mean-squared error incurred by this estimator.

(b) Find the MAP estimator of  $d$  given  $x$ .

**Problem 8.2** (Tweedie’s Formula). For the special case where  $X = D + N$ , where  $N$  is Gaussian noise of mean zero and variance  $\sigma^2$ , *Tweedie’s formula* says that the conditional mean (that is, the MMSE estimator) can be expressed as

$$\mathbb{E}[D|X = x] = x + \sigma^2 \ell'(x), \quad (8.33)$$

where

$$\ell'(x) = \frac{d}{dx} \log f_X(x), \quad (8.34)$$

where  $f_X(x)$  denotes the marginal PDF of  $X$ . In this exercise, we derive this formula.

(a) Assume that  $f_{X|D}(x|d) = e^{\alpha dx - \psi(d)} f_0(x)$  for some functions  $\psi(d)$  and  $f_0(x)$  and some constant  $\alpha$  (such that  $f_{X|D}(x|d)$  is a valid PDF for every value of  $d$ ). Define

$$\lambda(x) = \log \frac{f_X(x)}{f_0(x)}, \quad (8.35)$$

where  $f_X(x)$  is the marginal PDF of  $X$ , i.e.,  $f_X(x) = \int f_{X|D}(x|\delta) f_D(\delta) d\delta$ . With this, establish that

$$\mathbb{E}[D|X=x] = \frac{1}{\alpha} \frac{d}{dx} \lambda(x). \quad (8.36)$$

(b) Show that the case where  $X = D + N$ , where  $N$  is Gaussian noise of mean zero and variance  $\sigma^2$ , is indeed of the form required in Part (a) by finding the corresponding  $\psi(d)$ ,  $f_0(x)$ , and  $\alpha$ . Show that in this case, we have

$$\frac{f_0'(x)}{f_0(x)} = -\frac{x}{\sigma^2}, \quad (8.37)$$

and use this fact in combination with Part (a) to establish Tweedie's formula.

**Problem 8.3** (Wiener Filter). Consider a (discrete-time) signal that satisfies the difference equation  $d[n] = 0.5d[n-1] + v[n]$ , where  $v[n]$  is a sequence of uncorrelated zero-mean unit-variance random variables. We observe  $x[n] = d[n] + w[n]$ , where  $w[n]$  is a sequence of uncorrelated zero-mean random variables with variance 0.5.

(a) (you may skip this at first and do it later — it is conceptually straightforward) Show that for this signal model, the autocorrelation function of the signal  $d[n]$  is

$$\mathbb{E}[d[n]d[n+k]] = \frac{4}{3} \left(\frac{1}{2}\right)^{|k|}, \quad (8.38)$$

and thus the autocorrelation function of the signal  $x[n]$  is

$$\mathbb{E}[x[n]x[n+k]] = \begin{cases} \frac{11}{6}, & \text{for } k = 0, \\ \frac{4}{3} \left(\frac{1}{2}\right)^{|k|}, & \text{otherwise.} \end{cases} \quad (8.39)$$

(b) We would like to find an (approximate) linear predictor  $\hat{d}[n+3]$  using only the observations  $x[n], x[n-1], x[n-2], \dots, x[n-p]$ . Using the Wiener Filter framework, determine the optimal coefficients for the linear predictor. Find the corresponding mean-squared error for your predictor.

(c) We would like to find a linear denoiser  $\hat{d}[n]$  using *all* of the samples  $\{x[k]\}_{k=-\infty}^{\infty}$ . Find the filter coefficients and give a formula for the incurred mean-squared error.

**Problem 8.4** (Wiener Filter and Irrelevant Data). As we have seen in class, the (FIR) Wiener filter is given by

$$\mathbf{w} = R_x^{-1} \mathbf{r}_{dx}, \quad (8.40)$$

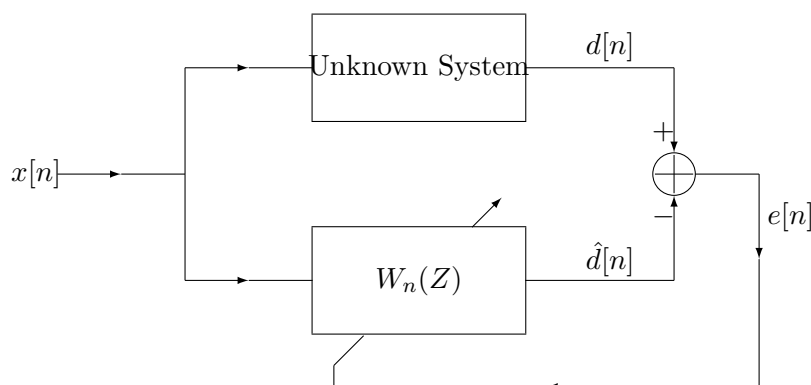
where  $R_x$  is the autocorrelation matrix of the data that's being used, and  $\mathbf{r}_{dx}$  is the cross-correlation between the data and the desired output. For this to be well defined,  $R_x$  should be full rank. In this problem, we study this question in more detail.

(a) In many applications, the signal acquisition process is noisy. That is, the data  $x[n] = s[n] + w[n]$ , where  $s[n]$  is an *arbitrary* signal, and  $w[n]$  is white noise. Prove that in this case, the  $p$ -dimensional autocorrelation matrix  $R_x$  is full rank (i.e., invertible) for any  $p$ . (Note: Be careful not to make *any* assumptions about the signal  $s[n]$ .)

(b) In some other cases,  $R_x$  could be rank-deficient. To study this, prove first that if the (FIR) Wiener filter based on the data  $\mathbf{x} = \{x[n]\}_{n=0}^{p-1}$  is  $\mathbf{w}$ , then the (FIR) Wiener filter based on the modified data  $A\mathbf{x}$  (where  $A$  is an invertible matrix) is  $A^{-H}\mathbf{w}$ , (where we use the relatively common notation  $A^{-H} = (A^{-1})^H$ ).

(c) Explain how to find the (FIR) Wiener filter when  $R_x$  is rank-deficient. Discuss existence and uniqueness. *Hint:* Use Part (b) to transform your data to a more convenient basis.

**Problem 8.5** (Adaptive Filters). One of the many uses of adaptive filters is for system identification as shown in the figure below. In this configuration, the same input is applied to an adaptive filter and to an unknown system, and the coefficients of the adaptive filter are adjusted until the difference between the outputs of the two systems is as small as possible.



Let the unknown system that is to be characterized by

$$d[n] = x[n] + 1.8x[n-1] + 0.81x[n-2] \quad (8.41)$$

With an input  $x[n]$  consisting of 1000 samples of unit variance white Gaussian noise, create the reference signal  $d[n]$ .

(a) Determine the range of values for the step size  $\mu$  in the LMS algorithm for the convergence in the mean.

(b) Implement an adaptive filter of order  $p = 4$  using the LMS algorithm. Set the initial weight vector equal to zero, and use a step size of  $\mu = 0.1\mu_{max}$ , where  $\mu_{max}$  is the largest step size allowed for convergence in the mean. Let the adaptive filter adapt and record the final set of coefficients.

(c) Repeat part (b) using the normalized LMS algorithm with  $\beta = 0.1$ , and compare your results.

(d) Make a plot of the learning curve by repeating the experiment described in part (b) for 100 different realizations of  $d[n]$ , and plotting the average of the plots of  $e^2[n]$  versus  $n$ . How many iterations are necessary for the mean-square error to fall to 10% of its peak value? Calculate the theoretical value for the excess mean-square error and compare it to what you observe in your plot of the learning curve.

**Problem 8.6** (Canonical Correlation Analysis). Let  $\mathbf{X}$  and  $\mathbf{Y}$  be zero-mean real-valued random vectors with covariance matrices  $R_{\mathbf{X}}$  and  $R_{\mathbf{Y}}$ , respectively. Moreover, let  $R_{\mathbf{XY}} = \mathbb{E}[\mathbf{XY}^T]$ . Our goal is to find vectors  $\mathbf{u}$  and  $\mathbf{v}$  such as to maximize the correlation between  $\mathbf{u}^T\mathbf{X}$  and  $\mathbf{v}^T\mathbf{Y}$ , that is,

$$\max_{\mathbf{u}, \mathbf{v}} \frac{\mathbb{E}[\mathbf{u}^T \mathbf{X} \mathbf{Y}^T \mathbf{v}]}{\sqrt{\mathbb{E}[\mathbf{u}^T \mathbf{X} \mathbf{X}^T \mathbf{u}]} \sqrt{\mathbb{E}[\mathbf{v}^T \mathbf{Y} \mathbf{Y}^T \mathbf{v}]}}. \quad (8.42)$$

Show how we can find the optimizing choices of the vectors  $\mathbf{u}$  and  $\mathbf{v}$  from the problem parameters  $R_{\mathbf{X}}$ ,  $R_{\mathbf{Y}}$ , and  $R_{\mathbf{XY}}$ .

*Hint:* Recall for the singular value decomposition that

$$\max_{\mathbf{v}} \frac{\|A\mathbf{v}\|}{\|\mathbf{v}\|} = \max_{\|\mathbf{v}\|=1} \|A\mathbf{v}\| = \sigma_1(A), \quad (8.43)$$

where  $\sigma_1(A)$  denotes the maximum singular value of the matrix  $A$ . The corresponding maximizer is the right singular vector  $\mathbf{v}_1$  (i.e., eigenvector of  $A^T A$ ) corresponding to  $\sigma_1(A)$ .

## Chapter 9

# Signal Representation

*It is generally assumed that the student has basic familiarity with the early topics in this chapter, such as bases, projections, and so on. Canonical references include [6, 7, 8].*

### Introduction

“Signal representation” refers to the act of representing an object (a “signal”)  $x$  as a linear combination of elements in a basic “dictionary” composed of elements  $\phi_k$ , where  $k$  runs over integers :

$$x = \sum_{\ell \in \mathbb{Z}} X_{\ell} \phi_{\ell}, \quad (9.1)$$

where the equality may be exact or approximate. The coefficients  $X_{\ell}$  are (for the purpose of this class) real or complex numbers. Then, instead of working with the original object (signal)  $x$ , we may work with its representation, given by the coefficients  $X_{\ell}$ .

The primary object of study is to find good and suitable dictionaries  $\{\phi_{\ell}\}_{\ell \in \mathbb{Z}}$ . In this module, we discuss the main methods and arguments relating to this quest. In guise of a motivational speech, such representations will satisfy one or more of these:

- Sparse representation
- Efficient processing possible
- etc.

### 9.1 Review : Notions of Linear Algebra

A key set of tools here (and throughout engineering and computer science) is linear algebra.

We will denote (column) vectors by bold symbols  $\mathbf{x}$ . The transpose of a vector  $\mathbf{x}$  is the row vector  $\mathbf{x}^T$ . The Hermitian transpose (transpose and complex conjugate) of a complex-valued vector  $\mathbf{x}$  is the row vector  $\mathbf{x}^H$ . The inner product (or dot product) of two vectors (of equal length) will be denoted as  $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{y}^H \mathbf{x}$ . (That is, following standard notational

conventions, in the  $\langle \cdot, \cdot \rangle$  notation, it is the *second* argument that is complex-conjugated.) The 2-norm of a vector is  $\|\mathbf{x}\| = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle}$ . More generally, the  $p$ -norm of a vector is  $\|\mathbf{x}\|_p = (\sum_{i=1}^n |x_i|^p)^{1/p}$ . These are genuine norms for all real numbers  $p \geq 1$ . For  $0 \leq p < 1$ , they are not norms since they violate the triangle inequality, but they are nonetheless of interest in applications.

Matrices are denoted by upper-case symbols  $A$ . They are of dimensions  $m \times n$ , meaning that they have  $m$  rows and  $n$  columns. The entry in row  $i$ , column  $j$  is denoted by  $\{A\}_{ij}$ , or simply  $A_{ij}$  when there is no confusion possible. The identity matrix is denoted by  $I$ . The matrix-vector product  $\mathbf{y} = A\mathbf{x}$ , where  $\mathbf{x}$  is of length  $n$ , is the vector  $\mathbf{y}$  of length  $m$  with entries  $y_i = \sum_{j=1}^n A_{ij}x_j$ . We use  $A^T$  to denote the transpose and  $A^H$  to denote the Hermitian transpose of the matrix  $A$ . Consider two matrices  $A$  and  $B$  with columns denoted by  $\mathbf{a}_i$  and  $\mathbf{b}_i$ , respectively. Then, the matrix product  $B^H A$  (if dimension-compatible) has as its entry in row  $i$ , column  $j$  the inner product  $\langle \mathbf{a}_j, \mathbf{b}_i \rangle$ . An alternative and equally useful expression for matrix multiplication is that  $AB^H = \sum_i \mathbf{a}_i \mathbf{b}_i^H$  (if the matrices are dimension-compatible). A *unitary* matrix is a matrix  $U$  satisfying  $UU^H = U^H U = I$ . The *trace* of a square matrix,  $\text{trace}(A)$ , is the sum of its diagonal entries. A property of many uses states that for any two (dimension-compatible) matrices  $A$  and  $B$ , we have  $\text{trace}(AB) = \text{trace}(BA)$ . Another useful property is that  $(AB)^H = B^H A^H$ .

### Symmetric Matrices : Spectral Decomposition

Perhaps the most important class of matrices are the symmetric (hence square) matrices. That is, matrices  $A$  for which we have  $A = A^H$  (and thus *a fortiori*,  $m = n$ ). Such matrices always admit a *spectral decomposition*,<sup>1</sup> i.e., they can be written as

$$A = U\Lambda U^H, \quad (9.2)$$

where  $\Lambda$  is a real-valued diagonal matrix and  $U$  is a unitary matrix. The  $n$  columns of  $U$  are called the *eigenvectors* of  $A$ , denoted by  $\mathbf{u}_1, \dots, \mathbf{u}_n$ . The  $n$  diagonal entries of  $\Lambda$ , usually denoted by  $\lambda_i$ , are called the corresponding *eigenvalues*. By inspection, Formula (9.2) can also be expressed as

$$A = \sum_{i=1}^n \lambda_i \mathbf{u}_i \mathbf{u}_i^H, \quad (9.3)$$

a shape that will be of interest to us in this class. A property with many uses is the fact that  $\text{trace}(A) = \sum_{i=1}^n \lambda_i$ , which follows simply from  $\text{trace}(A) = \text{trace}(U\Lambda U^H) = \text{trace}(\Lambda U^H U)$ .

### General Matrices : Singular Value Decomposition

For general matrices  $A$ , one can construct two instructive symmetric matrices, namely,  $AA^H$  and  $A^H A$ . Both of these admit spectral decompositions:

$$AA^H = U\Lambda'U^H \quad \text{and} \quad A^H A = V\Lambda''V^H, \quad (9.4)$$

<sup>1</sup>In general, the eigendecomposition is expressed as  $A = Q\Lambda Q^{-1}$ . When  $Q$  turns out to be a unitary matrix (thus,  $Q^{-1} = Q^H$ ), then one often refers to the eigendecomposition as a *spectral decomposition*.

and it is straightforward to show that, (i), the non-zero entries in  $\Lambda'$  and  $\Lambda''$  are the same, i.e.,  $AA^H$  and  $A^HA$  have the same eigenvalues, and (ii), all eigenvalues are non-negative. As you have seen, from these, one can construct the *singular value decomposition*

$$A = U\Sigma V^H = \sum_{i=1}^{\min(m,n)} \sigma_i \mathbf{u}_i \mathbf{v}_i^H, \quad (9.5)$$

where  $\Sigma$  is an  $m \times n$  diagonal matrix whose entries  $\sigma_i$  are simply the square roots of the eigenvalues of  $AA^H$  (or, equivalently,  $A^HA$ ). The values  $\sigma_i$  are thus non-negative and are referred to as the *singular values* of the matrix  $A$ .

### Rank of a matrix ; Norm of a matrix

The *rank* of a matrix  $A$  is the number of non-zero singular values it has in its singular value decomposition and will be denoted by  $\text{rank}(A)$ . An important relationship is that for any two (dimension-compatible) matrices, we have  $\text{rank}(AB) \leq \min\{\text{rank}(A), \text{rank}(B)\}$ . Note that no similarly useful and non-trivial relationship can be given for the rank of the *sum* of two matrices.

We will also find it useful to define *norms* for matrices. First, let us introduce the so-called *operator norms* that are derived from standard vector norms as follows:

$$\|A\|_p = \sup_{\mathbf{x} \neq \mathbf{0}} \frac{\|A\mathbf{x}\|_p}{\|\mathbf{x}\|_p} \quad (9.6)$$

An interesting special case is when  $p = 2$ , which is often called the *spectral norm* of the matrix  $A$ , and is easily seen to be equal to the largest singular value of the matrix  $A$ .

Of equal importance is the Frobenius norm

$$\|A\|_F = \sqrt{\sum_{i,j} |A_{ij}|^2}. \quad (9.7)$$

A first interesting observation (which can be proved by elementary manipulations) is that  $\|A\|_F^2 = \text{trace}(A^HA)$ . This also implies that  $\|A\|_F^2 = \sum_{i=1}^r \sigma_i^2$ . Another interesting property is that for any two (dimension-compatible) matrices  $A$  and  $B$ , we have that  $\|AB\|_F \leq \|A\|_F \|B\|_F$ , a consequence of the Cauchy-Schwarz inequality.

### Low-rank Matrix Approximation

Let us consider the following intuitively pleasing problem: Given a matrix  $A \in \mathbb{R}^{k \times n}$ , we seek to find a matrix  $B$  of rank no larger than  $p$  such that  $B$  is as close as possible to  $A$ , i.e., such that the norm  $\|A - B\|$  is as small as possible. If we use as the norm the spectral or the Frobenius norm, then this problem has an intuitively pleasing solution, given in the following theorem.

**Theorem 9.1** (Eckart-Young). *Let the SVD of the rank- $r$  matrix  $A$  be*

$$A = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^H, \quad \text{with } \sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_r. \quad (9.8)$$

*For integers  $p$  between 1 and  $r-1$ , let  $\hat{A}_p$  denote the truncated sum*

$$\hat{A}_p = \sum_{i=1}^p \sigma_i \mathbf{u}_i \mathbf{v}_i^H. \quad (9.9)$$

*Then, we have*

$$\min_{B: \text{rank}(B) \leq p} \|A - B\|_2 = \sigma_{p+1} \quad (9.10)$$

$$\min_{B: \text{rank}(B) \leq p} \|A - B\|_F = \sqrt{\sum_{k=p+1}^r \sigma_k^2}, \quad (9.11)$$

*and a minimizer of each of the two is  $B = \hat{A}_p$ . For the Frobenius norm,  $\hat{A}_p$  is the unique minimizer if and only if  $\sigma_p > \sigma_{p+1}$  (strict inequality).*

*Proof.* First, observe that  $\text{rank}(\hat{A}_p) \leq p$  and that we can write  $A - \hat{A}_p = \sum_{k=p+1}^r \sigma_k \mathbf{u}_k \mathbf{v}_k^H$ . Therefore :

- $\|A - \hat{A}_p\|_2 = \sigma_{p+1}$ , thus  $\min_{B: \text{rank}(B) \leq p} \|A - B\|_2 \leq \sigma_{p+1}$ , and
- $\|A - \hat{A}_p\|_F = \sqrt{\sum_{k=p+1}^r \sigma_k^2}$ , thus  $\min_{B: \text{rank}(B) \leq p} \|A - B\|_F \leq \sqrt{\sum_{k=p+1}^r \sigma_k^2}$ .

The more interesting part is the converse. We provide the converse proof only for the spectral norm in these notes. Consider any matrix  $B$  with  $\text{rank}(B) \leq p$ . Its null space has dimension no smaller than  $n - p$ , and thus, the dimension of the intersection  $\text{null}(B) \cap \text{span}\{\mathbf{v}_1, \dots, \mathbf{v}_{p+1}\}$  is at least one. For all vectors  $\mathbf{x}$  of norm one in this intersection, we have

$$(A - B)\mathbf{x} = A\mathbf{x} = \sum_{k=1}^{p+1} \sigma_k (\mathbf{v}_k^H \mathbf{x}) \mathbf{u}_k, \quad (9.12)$$

thus,

$$\|(A - B)\mathbf{x}\|^2 = \sum_{k=1}^{p+1} \sigma_k^2 |\mathbf{v}_k^H \mathbf{x}|^2. \quad (9.13)$$

But since the unit-norm vector  $\mathbf{x}$  lies in the span of  $\{\mathbf{v}_1, \dots, \mathbf{v}_{p+1}\}$ , we must have  $\sum_{k=1}^{p+1} |\mathbf{v}_k^H \mathbf{x}|^2 = 1$ . Then, since  $\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_{p+1}$ , we have  $\|(A - B)\mathbf{x}\|^2 \geq \sigma_{p+1}^2$ . Therefore, for every matrix  $B$  with  $\text{rank}(B) \leq p$ , a unit-norm vector  $\mathbf{x}$  can be found such that  $\|(A - B)\mathbf{x}\|^2 \geq \sigma_{p+1}^2$ . Thus,  $\|A - B\|_2 \geq \sigma_{p+1}$ .  $\square$



## 9.2 Fourier Representations

It is assumed that you have come across Fourier representations at least three times in your previous education (specifically, in your classes *Analysis III*, *Circuits & Systems II*, and *Signal Processing for Communications*). We here present a very brief overview, emphasizing some of the more advanced aspects.

Among all signal representations, Fourier representations are arguably the most important ones. This is due to several important reasons. First of all, they represent eigenvectors of LTI systems. Further reasons include important connections to wide-sense stationary signals and observations that many naturally occurring signal classes (audio, images, etc) have specific characteristics in the frequency domain. Moreover, Fourier representations can be calculated efficiently and have many desirable properties.

### 9.2.1 DFT and FFT

For the *discrete Fourier transform* (DFT), we follow the notation used in your prerequisite class, see [9, Section 4.2]. In line with this, let

$$W_N = e^{-j\frac{2\pi}{N}}. \quad (9.14)$$

The Fourier matrix  $W$  is the matrix whose entry in row  $k$ , column  $n$ , is given by

$$\{W\}_{kn} = W_N^{(k-1)(n-1)}, \text{ for } n, k \in \{1, 2, \dots, N\}. \quad (9.15)$$

and the DFT of the vector  $\mathbf{x}$  is the vector  $\mathbf{X}$  defined as

$$\mathbf{X} = W\mathbf{x}. \quad (9.16)$$

With this, the inverse transform is

$$\mathbf{x} = \frac{1}{N}W^H\mathbf{X}. \quad (9.17)$$

Explicitly, we can express the  $(k+1)$  entry of the vector  $\mathbf{X}$  (for  $k = 0, 1, \dots, N-1$ ) as

$$X[k] = \langle \mathbf{x}, \mathbf{w}_k \rangle = \mathbf{w}_k^H \mathbf{x} = \sum_{n=0}^{N-1} x[n]e^{-j2\pi\frac{kn}{N}}, \quad (9.18)$$

where the very last expression illustrates a slight notational anomaly, namely, we denote the elements of the signal vector  $\mathbf{x}$  by  $x[0], x[1], \dots, x[N-1]$ , that is, we number them from 0 to  $N-1$ . With this, we choose to follow the terminology used in [9, Section 4.2]. Also, we have used the notation  $\mathbf{w}_k = (1, W_N^k, W_N^{2k}, \dots, W_N^{(N-1)k})^H$ .

### Properties of the DFT

One of the most important reasons for the importance of the DFT is the wealth of useful properties it has. You have encountered these in detail in [9, Chapter 4].

*Cyclic shifts.* Consider the signal vector  $\mathbf{x}$  of length  $N$  and with entries denoted  $x[0], \dots, x[N-1]$ . Let  $\mathbf{y}$  be the signal vector  $\mathbf{x}$ , cyclically shifted to the right by  $n_0$  positions. We have the following DFT pair:

$$y[n] = x[(n - n_0) \bmod N] \quad \circ\!\!-\!\!\bullet \quad Y[k] = W_N^{kn_0} X[k], \quad (9.19)$$

where  $X[0], \dots, X[N-1]$  denote the entries of the DFT vector  $\mathbf{X} = W\mathbf{x}$ .

To establish this property, it is more convenient to use the sum representation than the matrix-vector representation. Namely,

$$Y[k] = \sum_{n=0}^{N-1} y[n] e^{-j\frac{2\pi}{N}kn} = \sum_{n=0}^{N-1} x[(n - n_0) \bmod N] e^{-j\frac{2\pi}{N}kn} \quad (9.20)$$

and change summation variables by defining  $m = n - n_0$ , which yields

$$Y[k] = \sum_{m=n_0}^{N-1+n_0} x[m \bmod N] e^{-j\frac{2\pi}{N}k(m+n_0)}. \quad (9.21)$$

We can rewrite the exponent as  $e^{-j\frac{2\pi}{N}k(m+n_0)} = e^{-j\frac{2\pi}{N}k(m \bmod N)} e^{-j\frac{2\pi}{N}kn_0}$  and introduce  $\ell = m \bmod N$  to obtain  $Y[k] = e^{-j\frac{2\pi}{N}kn_0} \sum_{\ell=0}^{N-1} x[\ell] e^{-j\frac{2\pi}{N}k\ell}$ , which completes the proof.

*Modulation property.* Consider the signal vector  $\mathbf{x}$  of length  $N$  and with entries denoted  $x[0], \dots, x[N-1]$ . Let  $\mathbf{y}$  be the signal vector with entries

$$y[n] = W_N^{-k_0 n} x[n] \quad \circ\!\!-\!\!\bullet \quad Y[k] = X[(k - k_0) \bmod N], \quad (9.22)$$

and the proof can be done following exactly the same steps as for the cyclic shift property.

*Duality.* A key observation is that these two properties are essentially one and the same. This is a reflection of the fact that DFT and inverse DFT are essentially the same (up to a complex-conjugate), and thus, the time and frequency variables can be exchanged.

### 9.2.2 The Other Fourier Representations

In your prerequisite classes, you have encountered several Fourier representations. For the theoretical understanding of the underpinnings and underlying ideas, the most important is the Fourier transform,

$$X(j\omega) = \int_{-\infty}^{\infty} x(t) e^{-j\omega t} dt, \quad (9.23)$$

whose inverse is given by

$$x(t) = \frac{1}{2\pi} \int_{-\infty}^{\infty} X(j\omega) e^{j\omega t} d\omega. \quad (9.24)$$

As you recall, there are several subtle aspects as to whether this inversion formula will indeed give back the original signal. Those will not be of particular interest to our class.

## 9.3 The Hilbert Space Framework for Signal Representation

Perhaps the most powerful framework to understand signal representation and approximation is that of *Hilbert space* which you have briefly encountered in the class *Signal Processing for Communications* [9, Chapter 3].<sup>2</sup>

A (real or complex) *vector space* is a set of vectors  $\mathbf{x} \in E$  with an addition for vectors, denoted by  $+$ , and a scalar multiplication (i.e., multiplication of a vector  $\mathbf{x}$  by a real- or complex-valued scalar  $\alpha$ ) such that for all  $\mathbf{x}, \mathbf{y}, \mathbf{z} \in E$  and all scalars  $\alpha, \beta$ :

- Commutativity :  $\mathbf{x} + \mathbf{y} = \mathbf{y} + \mathbf{x}$ .
- Associativity :  $(\mathbf{x} + \mathbf{y}) + \mathbf{z} = \mathbf{x} + (\mathbf{y} + \mathbf{z})$ , and  $\alpha(\beta\mathbf{x}) = (\alpha\beta)\mathbf{x}$ .
- Distributive laws :  $\alpha(\mathbf{x} + \mathbf{y}) = \alpha\mathbf{x} + \alpha\mathbf{y}$ , and  $(\alpha + \beta)\mathbf{x} = \alpha\mathbf{x} + \beta\mathbf{x}$ .
- There exists a vector  $\mathbf{0} \in E$  such that  $\mathbf{x} + \mathbf{0} = \mathbf{x}$  for all  $\mathbf{x} \in E$ .
- For all  $\mathbf{x} \in E$ , there exists an element  $-\mathbf{x} \in E$  such that  $\mathbf{x} + (-\mathbf{x}) = \mathbf{0}$ .
- For all  $\mathbf{x} \in E$ ,  $1 \cdot \mathbf{x} = \mathbf{x}$ .

A (real or complex) *inner product space* is a (real or complex) vector space together with an inner product  $\langle \mathbf{x}, \mathbf{y} \rangle \in \mathbb{R}$  or  $\mathbb{C}$  satisfying, for all  $\mathbf{x}, \mathbf{y}, \mathbf{z} \in E$  and scalars  $\alpha$ ,

- $\langle \mathbf{x} + \mathbf{z}, \mathbf{y} \rangle = \langle \mathbf{x}, \mathbf{y} \rangle + \langle \mathbf{z}, \mathbf{y} \rangle$ .
- $\langle \alpha\mathbf{x}, \mathbf{y} \rangle = \alpha\langle \mathbf{x}, \mathbf{y} \rangle$ .
- $\langle \mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{y}, \mathbf{x} \rangle^*$ .
- $\langle \mathbf{x}, \mathbf{x} \rangle \geq 0$ , with equality if and only if  $\mathbf{x} = \mathbf{0}$ .

The induced norm of the inner product space is defined as  $\|\mathbf{x}\| \stackrel{\text{def}}{=} \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle}$ . This definition directly implies the following important and useful facts:

- Cauchy-Schwarz inequality:  $|\langle \mathbf{x}, \mathbf{y} \rangle| \leq \|\mathbf{x}\| \|\mathbf{y}\|$  for all  $\mathbf{x}, \mathbf{y} \in E$ , with equality if and only if  $\mathbf{x} = \alpha\mathbf{y}$  for some scalar  $\alpha$ .
- Triangle inequality:  $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$ , with equality if and only if  $\mathbf{x} = \alpha\mathbf{y}$  for some real-valued non-negative scalar  $\alpha$ .
- Parallelogram identity:  $\|\mathbf{x} + \mathbf{y}\|^2 + \|\mathbf{x} - \mathbf{y}\|^2 = 2(\|\mathbf{x}\|^2 + \|\mathbf{y}\|^2)$ .

For example, to establish the Cauchy-Schwarz inequality, we may start by observing that  $\|\mathbf{x} - \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\|\mathbf{y}\|^2} \mathbf{y}\|^2 \geq 0$ , which holds by the definition of the norm. Writing this norm in terms of inner products and repeatedly applying the properties of the inner product leads to the Cauchy-Schwarz inequality.

The final key ingredient pertains to the *convergence* of sequences of vectors  $\mathbf{x}_n \in E$ . Quite naturally, we say that such a sequence converges to  $\mathbf{x} \in E$  if  $\lim_{n \rightarrow \infty} \|\mathbf{x}_n - \mathbf{x}\| = 0$ . A sequence  $\mathbf{x}_n \in E$  is called a *Cauchy sequence* if  $\lim_{m, n \rightarrow \infty} \|\mathbf{x}_m - \mathbf{x}_n\| = 0$ . Then, a *Hilbert space* is an inner product space with the additional property that *every* Cauchy sequence converges to a vector  $\mathbf{x} \in E$ . (This can be thought of as a technical condition which for the purpose of our class will not matter too much since it is satisfied for all examples of interest to us.)

---

<sup>2</sup>Our treatment here closely follows the development in the excellent textbook by Pierre Brémaud [7]. Alternatively, you may consult [8, Chapter 2] and/or follow the class COM-514 *Mathematical Foundations of Signal Processing*.

**Example 9.1** ( $n$ -dimensional complex vector space). This is the usual vector space with inner product  $\langle \mathbf{x}, \mathbf{y} \rangle = \sum_{i=1}^n x_i y_i^*$ . The induced norm is  $\|\mathbf{x}\| = \sqrt{\sum_{i=1}^n |x_i|^2}$ .

**Example 9.2** (Square-integrable functions (often denoted as  $L^2(\mathbb{R})$  or  $L_2(\mathbb{R})$ )). The set of all functions  $f(t)$  satisfying  $\int_{-\infty}^{\infty} |f(t)|^2 dt < \infty$ , with inner product  $\langle f, g \rangle = \int_{-\infty}^{\infty} f(t) g^*(t) dt$ , is a Hilbert space. The induced norm is  $\|f\| = \sqrt{\int_{-\infty}^{\infty} |f(t)|^2 dt}$ .

### Projection Theorem

For signal representation problems, the main reason why the Hilbert space framework is powerful is the projection theorem. This theorem tackles the following question : Given a Hilbert space  $H$  and a subspace  $G$  of  $H$  such that  $G$  is also a Hilbert space (meaning that  $G$  is also closed). Then, a very common task is that of representing any element  $\mathbf{x} \in H$  only using elements from the subspace  $G$  “in the best possible way.” Of course, since  $G$  is smaller than  $H$ , this leads to a more compact (approximate) representation of  $x$ , and is thus of obvious interest for many applications. More precisely, for any  $\mathbf{x} \in H$ , we are looking for an approximation  $\hat{\mathbf{x}} \in G$  such that  $\|\mathbf{x} - \hat{\mathbf{x}}\|$  is as small as possible. The projection theorem guarantees existence and uniqueness of this minimizer, and it establishes that the minimizer has the very useful property that it is *orthogonal* to the approximation error, i.e.,  $\langle \hat{\mathbf{x}}, \mathbf{x} - \hat{\mathbf{x}} \rangle = 0$ . This last relationship is often called the *orthogonality principle* and considerably simplifies the problem of finding the best approximation  $\hat{\mathbf{x}}$ . For any Hilbert subspace  $G$ , let us define  $G^\perp = \{\mathbf{z} \in H : \langle \mathbf{z}, \mathbf{x} \rangle = 0, \forall \mathbf{x} \in G\}$ . Then, we have the following statement:

**Theorem 9.2.** *Let  $\mathbf{x} \in H$ . There exists a unique element  $\mathbf{y} \in G$  such that  $\mathbf{x} - \mathbf{y} \in G^\perp$ . Moreover,  $\|\mathbf{y} - \mathbf{x}\| = \inf_{\mathbf{u} \in G} \|\mathbf{u} - \mathbf{x}\|$ .*

A proof can be found e.g. in [7, Sec. C1] or in [8, Ch.2]. We will also explore it to some extent in the Homework.

### Orthonormal Basis

A collection of vectors  $\{\mathbf{e}_n\}_{n \geq 0}$  in a Hilbert space  $H$  is called an *orthonormal system* if  $\langle \mathbf{e}_n, \mathbf{e}_k \rangle = 0$  for all  $n \neq k$ , and  $\|\mathbf{e}_n\| = 1$ , for all  $n \geq 0$ .

**Theorem 9.3** (Hilbert Basis Theorem).  *$\{\mathbf{e}_n\}_{n \geq 0}$  is an orthonormal system in  $H$ . Then, the following statements are equivalent:*

- $\{\mathbf{e}_n\}_{n \geq 0}$  generates the Hilbert space  $H$ .
- For all  $\mathbf{x} \in H$ , we have  $\|\mathbf{x}\|^2 = \sum_n |\langle \mathbf{x}, \mathbf{e}_n \rangle|^2$ .
- For all  $\mathbf{x} \in H$ , we have  $\mathbf{x} = \sum_n \langle \mathbf{x}, \mathbf{e}_n \rangle \mathbf{e}_n$ .

**Theorem 9.4** (Projection theorem, revisited). *Suppose  $G$  is spanned by the orthonormal basis  $\{\mathbf{g}_n\}_{n \geq 0}$ . Then, the element  $\mathbf{y} \in G$  that attains  $\min_{\mathbf{u} \in G} \|\mathbf{u} - \mathbf{x}\|$  is given by  $\mathbf{y} = \sum_n \langle \mathbf{x}, \mathbf{g}_n \rangle \mathbf{g}_n$ .*

## 9.4 General Bases, Frames, and Time-Frequency Analysis

### 9.4.1 The General Transform

#### Definition

A useful general way of thinking of transforms is in the shape of inner products with a set of “basis” functions:

$$T_x(\gamma) = \langle x(t), \phi_\gamma(t) \rangle \quad (9.25)$$

$$= \int_{-\infty}^{\infty} x(t) \phi_\gamma^*(t) dt, \quad (9.26)$$

where  $*$  denotes the complex conjugate.

The idea here is that ‘T’ denotes what kind of “basis” functions are being used and  $\gamma$  is the index of a basis function. The basis functions are  $\phi_\gamma(t)$  for all values of  $\gamma$ .

A good way of thinking about this is that for a fixed  $\gamma$ , the transform coefficient  $T_x(\gamma)$  is the result of projecting the original signal  $x(t)$  onto the “basis” element  $\phi_\gamma(t)$ .

An example is the Fourier transform, where instead of the letter  $\gamma$ , we more often use the letter  $\Omega$ , and where  $\phi_\Omega(t) = e^{j\Omega t}$ . Hence, in line with the above general notation, we could write

$$FT_x(\Omega) = \langle x(t), \phi_\Omega(t) \rangle \quad (9.27)$$

$$= \int_{-\infty}^{\infty} x(t) e^{-j\Omega t} dt, \quad (9.28)$$

Of course, we more often simply write  $X(\Omega)$  (or  $X(j\Omega)$ ) in place of  $FT_x(\Omega)$ .

#### Alternative Formulation

For our next step, we need the (general) Parseval/Plancherel formula, which asserts that

$$\int_{-\infty}^{\infty} f(t) g^*(t) dt = \frac{1}{2\pi} \int_{-\infty}^{\infty} F(j\Omega) G^*(j\Omega) d\Omega. \quad (9.29)$$

Using this, we can rewrite the general transform as

$$T_x(\gamma) = \langle x(t), \phi_\gamma(t) \rangle \quad (9.30)$$

$$= \int_{-\infty}^{\infty} x(t) \phi_\gamma^*(t) dt \quad (9.31)$$

$$= \frac{1}{2\pi} \int_{-\infty}^{\infty} X(j\Omega) \Phi_\gamma^*(j\Omega) d\Omega \quad (9.32)$$

$$= \langle X(j\Omega), \frac{1}{2\pi} \Phi_\gamma(j\Omega) \rangle \quad (9.33)$$

Hence, we now have *two* good ways of thinking about transforms: For a fixed  $\gamma$ , the transform coefficient  $T_x(\gamma)$  is the result of projecting the original signal  $x(t)$  onto the “basis”

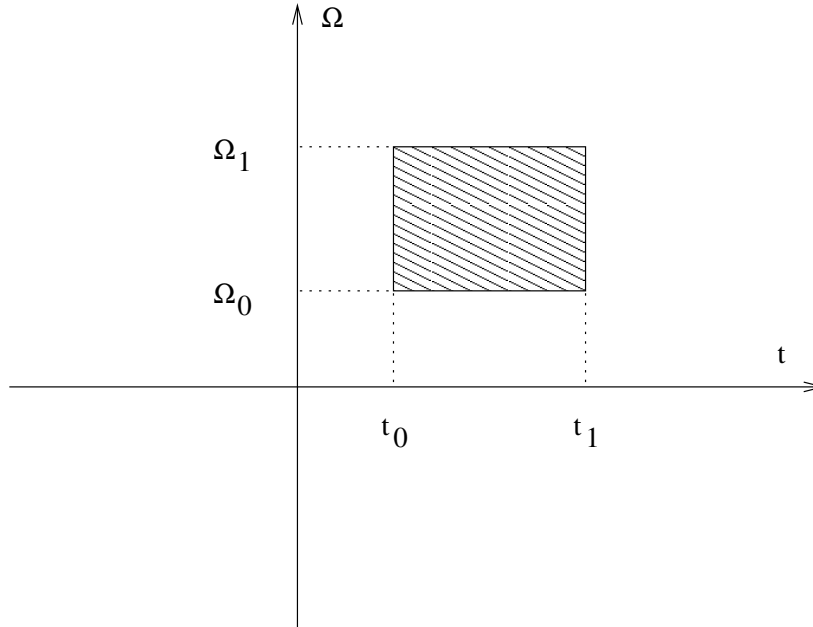


Figure 9.1: A conceptual picture: We imagine that the basis element  $\phi_\gamma(t)$  only lives in the shaded box, i.e., that the signal is very small outside the interval  $t_0 \leq t \leq t_1$ , and that its spectrum  $\Phi_\gamma(j\Omega)$  is very small outside of the interval  $\Omega_0 \leq \Omega \leq \Omega_1$ .

element  $\phi_\gamma(t)$ , and *equivalently*, of projecting the original spectrum  $X(j\Omega)$  onto the spectrum of the “basis” element  $\phi_\gamma(t)$ , which is  $\frac{1}{2\pi}\Phi_\gamma(j\Omega)$ .

Consider Figure 9.1: Merely as a thought experiment, let us think of a “basis” element  $\phi_\gamma(t)$  that lives<sup>3</sup> only inside the box illustrated in Figure 9.1. Then, a great way of thinking about the transform coefficient  $T_x(\gamma)$  is that it tells us “how much” of the original signal  $x(t)$  sits inside that box.

In line with this intuition, for the Fourier transform, the transform coefficient  $T_x(\Omega)$  tells us “how much” of the original signal  $x(t)$  sits at frequency  $\Omega$ , and the “box” shown in Figure 9.1 is infinitesimally thin in frequency and infinitely long in time.

### 9.4.2 The Heisenberg Box Of A Signal

Reconsider the conceptual picture given in Figure 9.1. Now, we want to make this precise. In order to do so, consider any signal  $\phi(t)$ . For simplicity (and without loss of generality), we assume that the signal is “normalized” such that

$$\int_{-\infty}^{\infty} |\phi(t)|^2 dt = 1. \quad (9.34)$$

Note that by Parseval, this also means that  $\frac{1}{2\pi} \int_{-\infty}^{\infty} |\Phi(j\Omega)|^2 d\Omega = 1$ .

<sup>3</sup>In the next section, we will make precise what “lives” means.

We define the following quantities. The “middle” of the signal  $\phi(t)$  is given by

$$m_t = \int_{-\infty}^{\infty} t |\phi(t)|^2 dt. \quad (9.35)$$

If you have taken a class in probability, you will recognize this to be the mean value of the distribution  $|\phi(t)|^2$ .

Similarly, we define the “middle” of the spectrum  $\Phi(j\Omega)$  to be

$$m_\Omega = \int_{-\infty}^{\infty} \Omega \frac{1}{2\pi} |\Phi(j\Omega)|^2 d\Omega, \quad (9.36)$$

with a similar probability interpretation.

Moreover, we define:

$$\sigma_t^2 = \int_{-\infty}^{\infty} (t - m_t)^2 |\phi(t)|^2 dt, \quad (9.37)$$

$$\sigma_\Omega^2 = \int_{-\infty}^{\infty} (\Omega - m_\Omega)^2 \frac{1}{2\pi} |\Phi(j\Omega)|^2 d\Omega. \quad (9.38)$$

Again, these can be understood as the respective *variances* of the two “probability distributions.”

With these definitions, we can now draw a more precise picture of the time-frequency box of the signal  $\phi(t)$ , as given in Figure 9.2.

We should also point out that for the Fourier transform, the basis functions are of the form  $\phi(t) = e^{j\Omega_0 t}$ , and for those, the above integrals do not all converge, so special care is required mathematically. However, the right intuition is to say that the Heisenberg box (the term appears in [6], and perhaps earlier) of the function  $\phi(t) = e^{j\Omega_0 t}$  is a horizontal line at frequency  $\Omega_0$ .

### 9.4.3 The Uncertainty Relation

So, what are the possible Heisenberg boxes?

**Theorem 9.5** (uncertainty relation). *For any function  $\phi(t)$ , the Heisenberg box must satisfy*

$$\sigma_t \sigma_\Omega \geq \frac{1}{2}. \quad (9.39)$$

That is, Heisenberg boxes cannot be too small. Or: transforms cannot have a very high time resolution *and* a very high frequency resolution at the same time. (Proof: see class.)

### 9.4.4 The Short-time Fourier Transform

It has long been recognized that one of the most significant drawbacks of the Fourier transform is its lack of *time localization*: An event that is localized in time (such as a signal discontinuity) affects all of the frequencies (remember the Gibbs phenomenon). This feature is clearly undesirable for many engineering tasks, including compression and classification.

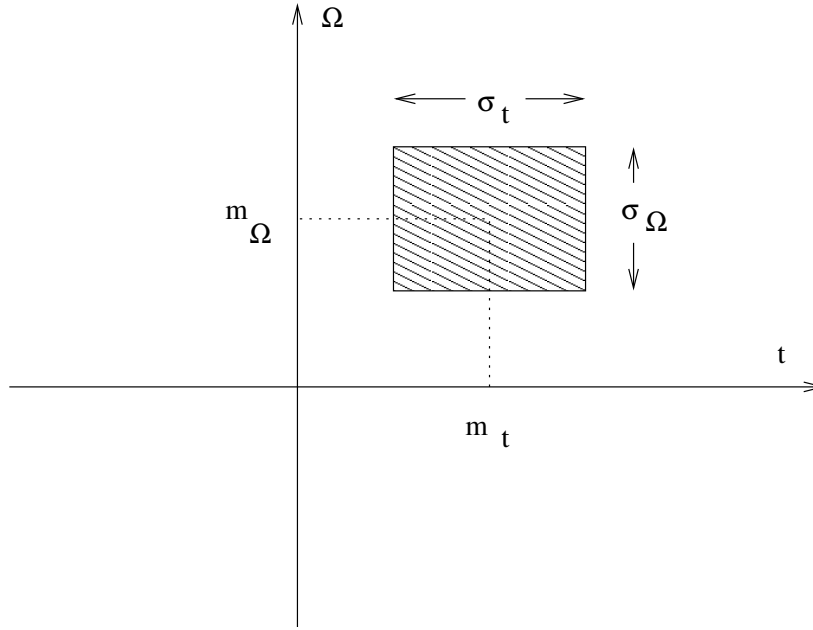


Figure 9.2: The Heisenberg box of the function  $\phi(t)$  (i.e., the place in time and frequency where the function  $\phi(t)$  is really alive).

To regain some of the time localization, one could do a “short-time” Fourier transform, essentially chopping up the signal into “short” pieces and taking Fourier transforms separately for each piece. Kind of trivially, this gives back some time localization.

More generally, the following form can be given:

$$STFT_x(\tau, \Omega) = \int_{-\infty}^{\infty} x(t) g^*(t - \tau) e^{-j\Omega t} dt, \quad (9.40)$$

where the function  $g(t)$  is an appropriate “window” function that cuts out a piece of the signal  $x(t)$ . With the parameter  $\tau$ , we can place the window wherever we want.

With regard to the general transform, here, instead of the letter  $\gamma$ , we use the pair  $(\tau, \Omega)$ , and

$$\phi_{\tau, \Omega}(t) = g(t - \tau) e^{j\Omega t}. \quad (9.41)$$

Many different window functions  $g(t)$  are being used, but one of the easiest to understand is the Gaussian window:

$$g(t) = \frac{1}{\sqrt[4]{\pi\sigma^2}} e^{-\frac{t^2}{2\sigma^2}}. \quad (9.42)$$

Note that strictly speaking, this window is never zero, so it does not really “cut” the signal. However, if  $|t|$  is large,  $g(t)$  is tiny, so this is “almost the same as zero,” but much easier to



analyze. With this window, we find the “basis” elements to be

$$\phi_{\tau_0, \Omega_0}(t) = \frac{1}{\sqrt[4]{\pi\sigma^2}} e^{-\frac{(t-\tau_0)^2}{2\sigma^2}} e^{j\Omega_0 t}. \quad (9.43)$$

Now, we want to find explicitly the Heisenberg box of this “basis” function. To this end, we need the Fourier transform of the Gaussian window, which is known to be

$$G(j\Omega) = \sqrt[4]{4\pi\sigma^2} e^{-\frac{\Omega^2\sigma^2}{2}}, \quad (9.44)$$

and thus, using the standard time- and frequency-shift properties of the Fourier transform,

$$\Phi_{\tau_0, \Omega_0}(j\Omega) = \sqrt[4]{4\pi\sigma^2} e^{-\frac{(\Omega-\Omega_0)^2\sigma^2}{2}} e^{-j\Omega\tau_0}. \quad (9.45)$$

Now, we can find the corresponding parameters of the Heisenberg box as:

$$m_t = \tau_0, \quad (9.46)$$

$$m_\Omega = \Omega_0 \quad (9.47)$$

$$\sigma_t^2 = \frac{\sigma^2}{2} \quad (9.48)$$

$$\sigma_\Omega^2 = \frac{1}{2\sigma^2}, \quad (9.49)$$

and so, we can draw the corresponding Figure 9.2. It is also interesting to note that for the Gaussian window, the Heisenberg uncertainty relation (Theorem 9.5) is satisfied with equality. It can be shown that the Gaussian window is (essentially) the only function that satisfies the uncertainty relation with equality, see e.g. [6, p.31].

## 9.5 Multi-Resolution Concepts and Wavelets

This is an interesting case of the general transform where we start from a *single* function

$$\psi(t), \quad (9.50)$$

sometimes called the *mother wavelet*.

Then, we build up our “dictionary” by shifting and scaling the mother wavelet, specifically,

$$\psi_{m,n}(t) = 2^{-m/2} \psi(2^{-m}t - n), \quad (9.51)$$

where  $n$  and  $m$  are arbitrary *integers* (positive, negative, or zero). That is, in place of the parameter  $\gamma$ , we will use the pair of integers  $(m, n)$ , where  $m$  is the scale (the bigger the coarser) and  $n$  is the shift. We will often denote the transform coefficient as

$$a_{m,n} = WT_x(m, n) = \langle x(t), \psi_{m,n}(t) \rangle. \quad (9.52)$$

The following key questions are of obvious interest:

- What are the conditions such that we can recover the original signal  $x(t)$  from the wavelet coefficients  $a_{m,n}$ ?
- How do we design good mother wavelets  $\psi(t)$ ?
- How do we efficiently compute the wavelet coefficients  $a_{m,n}$  for a given signal  $x(t)$ ?
- and of course many more...

### 9.5.1 The Haar Wavelet

We will start with the Haar wavelet:

$$\psi(t) = \begin{cases} 1, & \text{for } 0 \leq t < \frac{1}{2}, \\ -1, & \text{for } \frac{1}{2} \leq t < 1, \\ 0, & \text{otherwise.} \end{cases} \quad (9.53)$$

*Exercise:* Sketch  $\psi_{0,0}(t)$ ,  $\psi_{1,0}(t)$ ,  $\psi_{-1,0}(t)$ , and  $\psi_{-1,3}(t)$ .

Facts:

1. The functions  $\psi_{m,n}(t)$ , taken over all integers  $m$  and  $n$ , are an orthonormal set. (Easy to verify.)
2. The functions  $\psi_{m,n}(t)$ , taken over all integers  $m$  and  $n$ , are in fact an orthonormal basis for  $L^2(\mathbb{R})$ , the space of all functions  $x(t)$  for which  $\int_{-\infty}^{\infty} |x(t)|^2 dt$  is finite. (This is more difficult to prove, see e.g. [6, ch.7].)

Due to these facts, we can express any square-integrable function  $x(t)$  in the following form:

$$x(t) = \sum_{m=-\infty}^{\infty} \sum_{n=-\infty}^{\infty} a_{m,n} \psi_{m,n}(t), \quad (9.54)$$

which will be called the *Haar expansion* (or more generally, the wavelet expansion) of the signal  $x(t)$ . Moreover, due to the orthogonality, we also know that

$$a_{m,n} = \langle x(t), \psi_{m,n}(t) \rangle = \int_{-\infty}^{\infty} x(t) \psi_{m,n}(t) dt. \quad (9.55)$$

To understand how this wavelet works, it is instructive to consider a piecewise constant function, as in Figure 9.3. We here follow the development in [10, p.211]. Specifically, consider the function  $f^{(0)}(t)$  which is piecewise constant over intervals of length  $2^0 = 1$ , and assumes the values  $\dots, b_{-1}, b_0, b_1, b_2, \dots$ . As shown in Figure 9.3, we can write  $f^{(0)}(t)$  as the sum of two components: A sequence of shifted versions of the Haar wavelet at scale  $m = 1$  (i.e., of the functions  $\psi_{1,n}(t)$ ) and a “residual” function  $f^{(1)}(t)$ , which is piecewise constant over intervals of length  $2^1 = 2$ .

Specifically, we can write

$$d^{(1)}(t) = \sum_{n=-\infty}^{\infty} \underbrace{\frac{b_{2n} - b_{2n+1}}{\sqrt{2}}}_{a_{1,n}} \psi_{1,n}(t), \quad (9.56)$$

and

$$f^{(1)}(t) = \sum_{n=-\infty}^{\infty} \underbrace{\frac{b_{2n} + b_{2n+1}}{\sqrt{2}}}_{c_n} \varphi_{1,n}(t), \quad (9.57)$$

where we use

$$\psi_{1,n}(t) = \begin{cases} \frac{1}{\sqrt{2}}, & \text{for } 2n \leq t < 2n+1, \\ -\frac{1}{\sqrt{2}}, & \text{for } 2n+1 \leq t < 2n+2, \\ 0, & \text{otherwise.} \end{cases} \quad (9.58)$$

$$\varphi_{1,n}(t) = \begin{cases} \frac{1}{\sqrt{2}}, & \text{for } 2n \leq t < 2n+2, \\ 0, & \text{otherwise.} \end{cases} \quad (9.59)$$

The real boost now comes from observing that we can just continue along the same lines and decompose  $f^{(1)}(t)$  into two parts.

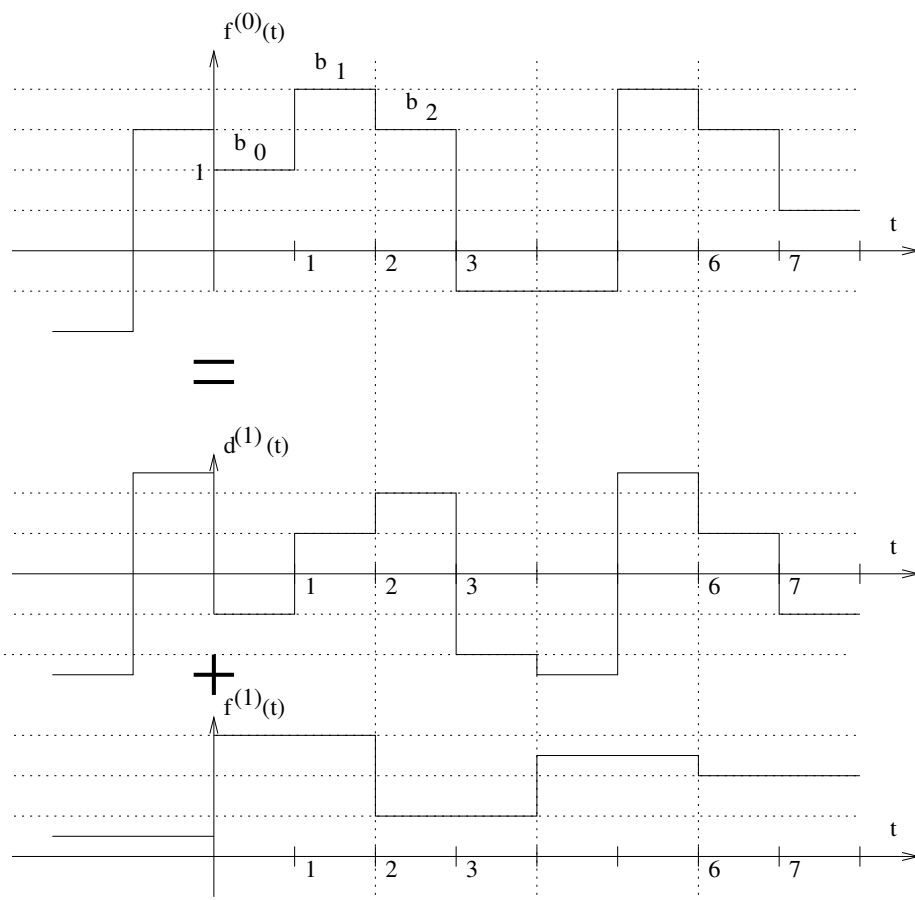


Figure 9.3: The Haar wavelet at work.

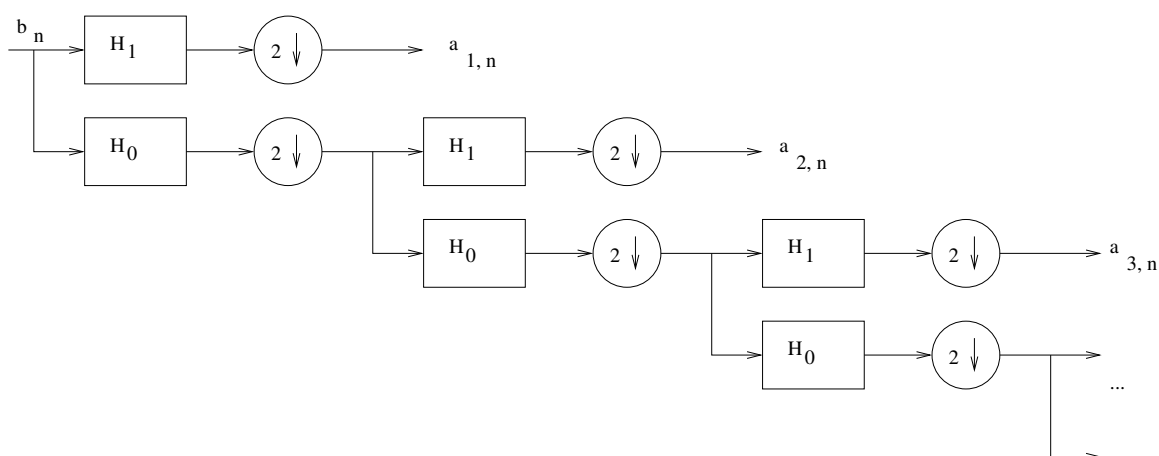


Figure 9.4: The Haar filter bank.

### A Filter Bank Companion To The Haar Wavelet

Let us now reconsider Equations (9.56) and (9.57). We can rewrite them purely in terms of the coefficients as

$$a_{1,n} = \frac{b_{2n} - b_{2n+1}}{\sqrt{2}} \quad (9.60)$$

$$c_n = \frac{b_{2n} + b_{2n+1}}{\sqrt{2}}, \quad (9.61)$$

which we easily recognize as filtering the sequence  $b_n$  with two different filters and then downsampling by a factor of two! In other words, starting from the coefficients  $b_n$ , there is a simple filter bank structure that computes all the wavelet coefficients, illustrated in Figure 9.4. For the Haar example, the filters are

$$H_1(z) = \frac{1}{\sqrt{2}}(1 - z) \quad (9.62)$$

$$H_0(z) = \frac{1}{\sqrt{2}}(1 + z). \quad (9.63)$$

Note that these filters are not quite causal, but since they are FIR, this is not a problem at all.

### Tilings Of The Time-frequency Plane

An interesting observation follows by observing that the Haar filter  $H_1(z)$  is (a crude version of) a highpass filter, and  $H_0(z)$  a (no less crude version of a) lowpass filter. Nevertheless, if we merely go with this highpass/lowpass picture, and merge it with the downsampling in time, we can draw a figure like the one given in Figure 9.5: consider for example the wavelet coefficient  $a_{1,0}$ . It results from highpass-filtering, so it pertains to the upper half of

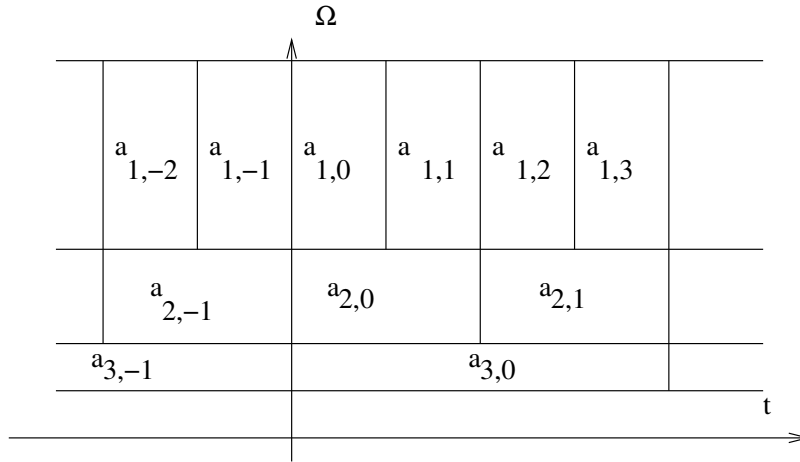


Figure 9.5: A tiling of the time frequency plane that somewhat mirrors what the Haar wavelet is doing.

the frequencies. Moreover, it pertains only to the time interval from 0 to 2. This is how we found the rectangle labelled  $a_{1,0}$  in Figure 9.5, and this is how you can find all the other rectangles in the figure.

### 9.5.2 Multiresolution Concepts

We discuss the “axiomatic” way of thinking about wavelets, as introduced by Mallat and Meyer, see e.g. [10, Section 4.2] or [6, Section 7.1].

The basic object is an *embedded* sequence of closed subspaces, which we will denote as

$$\cdots \subset V_2 \subset V_1 \subset V_0 \subset V_{-1} \subset V_{-2} \subset \cdots \quad (9.64)$$

In the Haar example,  $V_0$  is the space of functions that are piece-wise constant over intervals of length 1, and  $V_1$  the space of functions that are piece-wise constant over intervals of length 2. So, with respect to Figure 9.3,  $f^{(0)} \in V_0$  and  $f^{(1)}(t) \in V_1$ . However, due to the multiresolution embedding, we also have  $f^{(1)}(t) \in V_0$ , which is easily verified in Figure 9.3.

The key to wavelets is that the *difference* between subsequent spaces  $V_m$  are precisely the wavelet spaces  $W_m$ . That is,

$$V_m = V_{m+1} \oplus W_{m+1}, \quad (9.65)$$

and the wavelet space  $W_{m+1}$  is *orthogonal* to  $V_{m+1}$ . For the Haar example, with respect to Figure 9.3,  $d^{(1)}(t) \in W_1$  and it is easy to verify that  $W_1$  is orthogonal to  $V_1$ .

More generally, here are the basic requirements on the spaces  $V_m$  :

A1.  $\bigcup_{m \in \mathbb{Z}} V_m = L^2(\mathbb{R})$

A2.  $\bigcap_{m \in \mathbb{Z}} V_m = \emptyset$

- A3.  $f(t) \in V_m \iff f(2^m t) \in V_0$
- A4.  $f(t) \in V_0 \implies f(t - n) \in V_0, \forall n \in \mathbb{Z}$
- A5. There exists a function  $\varphi(t)$  (called the *scaling function*) such that  $\{\varphi(t - n)\}_{n \in \mathbb{Z}}$  is an orthonormal<sup>4</sup> basis for  $V_0$ .

The big question is how to find interesting scaling functions  $\varphi(t)$ , and the corresponding wavelet  $\psi(t)$ . The following consequences of the axioms of multiresolution will help us in doing so.

- C1. Because we require  $\{\varphi(t - n)\}_{n \in \mathbb{Z}}$  be an orthonormal basis, we must have

$$\sum_{k \in \mathbb{Z}} |\Phi(\omega + 2k\pi)|^2 = 1. \quad (9.66)$$

- C2. Because  $\varphi(t)$  also lives in  $V_{-1}$  and  $\{\varphi(2t - n)\}_{n \in \mathbb{Z}}$  is an orthonormal basis of  $V_{-1}$ , it must be possible to express

$$\varphi(t) = \sqrt{2} \sum_{n \in \mathbb{Z}} g_0[n] \varphi(2t - n), \quad (9.67)$$

or, equivalently,

$$\Phi(\omega) = \frac{1}{\sqrt{2}} G_0(e^{j\omega/2}) \Phi(\omega/2). \quad (9.68)$$

This is sometimes called the *two-scale equation*, a name that I find rather cute.

- C3. Combining points C1 and C2, we can also infer that the “filter”  $g_0[n]$  must satisfy

$$|G_0(e^{j\omega})|^2 + |G_0(e^{j(\omega+\pi)})|^2 = 2. \quad (9.69)$$

Namely, use (9.66) for  $2\omega$  and plug in (9.68).

- C4. Now, we want to find a function  $\psi(t)$  (the wavelet) such that  $\{\psi(t - n)\}_{n \in \mathbb{Z}}$  is an orthonormal basis for  $W_0$ . Note that it is *not trivial* that such a function even exists! However, one can prove this (see e.g. [10, Theorem 4.3]). If we merely assume that such a  $\psi(t)$  exists, then we immediately know that since it lies in  $V_{-1}$ , it must be possible to express

$$\psi(t) = \sqrt{2} \sum_{n \in \mathbb{Z}} g_1[n] \varphi(2t - n), \quad (9.70)$$

or, equivalently,

$$\Psi(\omega) = \frac{1}{\sqrt{2}} G_1(e^{j\omega/2}) \Phi(\omega/2). \quad (9.71)$$

---

<sup>4</sup>The non-orthogonal case has also been studied extensively in the literature.

Moreover,  $\psi(t)$  must also be orthogonal to any function in  $V_0$ , which is satisfied if we select

$$g_1[n] = (-1)^n g_0[-n+1], \quad (9.72)$$

or equivalently,

$$G_1(e^{j\omega}) = -e^{-j\omega} G_0^*(e^{j(\omega+\pi)}). \quad (9.73)$$

*Note:* One can also show that there are not too many more choices.

The main upshot is that we have reduced the wavelet design problem to one of finding a scaling function  $\varphi(t)$  and a “filter”  $g_0[n]$ . This is precisely what we will exploit next.

### 9.5.3 Wavelet Design — A Fourier Technique

We discuss Meyer’s Fourier-based wavelet design technique, see e.g. [10, Section 4.3] or [6, Section 7.2.2].

Specifically, we now exploit the insights from the previous section. In particular, the minute we have a valid scaling function  $\varphi(t)$  along with the “filter”  $g_0[n]$ , we are done. The problem is that the two are not nicely “ordered”. That is, we cannot just select any  $\varphi(t)$  that satisfies orthogonality and then select any filter. Here is a design procedure due to Meyer:

Y1. Select a function  $\theta(x)$  such that  $\theta(x) = 0$  for  $x \leq 0$  and  $\theta(x) = 1$  for  $x \geq 1$  and

$$\theta(x) + \theta(1-x) = 1, \text{ for } 0 \leq x \leq 1. \quad (9.74)$$

Y2. Set

$$\Phi(\omega) = \begin{cases} \sqrt{\theta(2 + \frac{3\omega}{2\pi})}, & \omega \leq 0, \\ \sqrt{\theta(2 - \frac{3\omega}{2\pi})}, & \omega > 0. \end{cases} \quad (9.75)$$

It is easy to verify that this satisfies point C1 above. This is the scaling function. Hence,  $V_0 = \text{span}\{\varphi(t-n), n \in \mathbb{Z}\}$ .

Y3. Set

$$G_0(e^{j\omega}) = \sqrt{2} \sum_{k \in \mathbb{Z}} \Phi(2\omega + 4k\pi). \quad (9.76)$$

You may need a quick sketch to convince yourself that this satisfies the two-scale equation C2.

This is it - the wavelet now follows directly from C4.



### 9.5.4 Wavelet Algorithms

We discuss Mallat's algorithm, see e.g. [10, Section 4.5.3] or [6, Section 7.3].

Suppose that the signal of interest  $f^{(0)}(t)$  lives at scale  $V_0$ . Hence, it can be written as

$$f^{(0)}(t) = \sum_{n=-\infty}^{\infty} b_n \varphi(t - n). \quad (9.77)$$

Our goal is to express that function in terms of the wavelet, i.e.,

$$f^{(0)}(t) = \sum_{m=1}^{\infty} \sum_{n=-\infty}^{\infty} a_{m,n} \psi_{m,n}(t). \quad (9.78)$$

Note that the scale parameter  $m$  only starts at 1. This is because our signal  $f^{(0)}(t)$  lives at scale  $V_0$ , and hence, only the wavelets that are coarser than scale  $V_0$  are needed.

We now describe an algorithm that takes the coefficients  $b_n$  as inputs and produces the coefficients  $a_{m,n}$  as outputs. The key step is, at each scale  $m = 1, 2, \dots$ , to split the remaining signal into two parts, namely  $V_m$  and  $W_m$ . More specifically:

- M1. *Project  $f^{(0)}(t)$  into  $W_1$ .* This is easy because we have an orthonormal basis for  $W_1$ , namely, the wavelets  $\{\frac{1}{\sqrt{2}}\psi(t/2 - n)\}_{n \in \mathbb{Z}}$ . Therefore, the projection is given by

$$d^{(1)}(t) \stackrel{\text{def}}{=} \sum_{n=-\infty}^{\infty} \underbrace{\langle f^{(0)}(t), \frac{1}{\sqrt{2}}\psi(t/2 - n) \rangle}_{a_{1,n}} \psi(t/2 - n). \quad (9.79)$$

But the coefficients  $a_{1,n}$  can be computed easily from the coefficients  $b_n$  as

$$a_{1,n} = \sum_{\ell=-\infty}^{\infty} \tilde{g}_1[2n - \ell] b_{\ell}, \quad (9.80)$$

where  $\tilde{g}_1[n] = g_1[-n]$ . Note that this is merely a filtering operation, followed by downsampling by a factor of two.

- M2. *Project  $f^{(0)}(t)$  into  $V_1$ .* This is easy because we have an orthonormal basis for  $V_1$ , namely,  $\{\frac{1}{\sqrt{2}}\varphi(2t - n)\}_{n \in \mathbb{Z}}$ . Therefore, the projection is given by

$$f^{(1)}(t) \stackrel{\text{def}}{=} \sum_{n=-\infty}^{\infty} \underbrace{\langle f^{(0)}(t), \frac{1}{\sqrt{2}}\varphi(t/2 - n) \rangle}_{b_{1,n}} \varphi(t/2 - n). \quad (9.81)$$

But the coefficients  $b_{1,n}$  can be computed easily from the coefficients  $b_n$  as

$$b_{1,n} = \sum_{\ell=-\infty}^{\infty} \tilde{g}_0[2n - \ell] b_{\ell}, \quad (9.82)$$

where  $\tilde{g}_0[n] = g_0[-n]$ . Note that this is merely a filtering operation, followed by downsampling by a factor of two.

The main trick is to realize that we can now repeat these steps for  $f^{(1)}(t)$ , projecting this signal into the spaces  $W_2$  and  $V_2$ , respectively. This leads exactly to the same filtering and downsampling. In the end, this leads exactly to the filter bank structure of Figure 9.4, with the filters being

$$h_1[n] = \tilde{g}_1[n] \quad (= g_1[-n]) \quad (9.83)$$

$$h_0[n] = \tilde{g}_0[n] \quad (= g_0[-n]) \quad (9.84)$$

### 9.5.5 Wavelet Design — Further Considerations

We discuss the vanishing moments versus support size question, leading to Daubechies' wavelet, see e.g. [10, Section 4.4.4] or [6, Section 7.2].

One of the key goals of wavelet analysis is to ensure a *compact* signal description. That is, we would like many wavelet coefficients to be small or zero, and only as few as possible to be large. There are two wavelet “dimensions,” namely, shift and scale, and we discuss them in turn:

- D1. Along the *time shifts*, such a property will be enabled by *short time support* of the wavelet: That way, local behavior only affects a few wavelets. It can be verified that  $\varphi(t)$  has compact support if and only if the filter  $g_0[n]$  also has compact support (i.e., is an FIR filter), and those supports are equal. Let those supports be  $[N_1, N_2]$ . Then, the wavelet  $\psi(t)$  also has compact support, namely,  $[(N_1 - N_2 + 1)/2, (N_2 - N_1 + 1)/2]$ .
- D2. Along the *scale dimension*, we will consider the so-called vanishing moments. A wavelet  $\psi(t)$  is said to have  $p$  vanishing moments if

$$\int_{-\infty}^{\infty} t^k \psi(t) dt = 0, \text{ for } 0 \leq k < p. \quad (9.85)$$

*Intuition:* Wherever the signal  $x(t)$  is smooth, this ensures that wavelet coefficients corresponding to small scales will be small. (This can be understood by expressing  $x(t)$  locally as a Taylor series.) With a small loss of generality, the wavelet having  $p$  vanishing moments is that same as the filter  $G_0(e^{j\omega})$  having  $p$  zeros at  $\omega = \pi$ . That is, the filter  $G_0(e^{j\omega})$  can be written as

$$G_0(e^{j\omega}) = \sqrt{2} \left( \frac{1 + e^{-j\omega}}{2} \right)^p R(e^{-j\omega}), \quad (9.86)$$

where  $R(e^{j\omega})$  is arbitrary but has no poles at  $\omega = \pi$ .

Now, what we would like is to have a wavelet of compact support and having as many vanishing moments as possible. However, unfortunately, there is a kind of uncertainty principle, given by Daubechies' theorem:

$$\text{time support of wavelet} \geq 2(\# \text{ of vanishing moments}) - 1. \quad (9.87)$$

The most astonishing thing about this insight is that the functions that attain this bound with equality look rather crazy...!

A bit more formally, the theorem is stated as follows (see e.g. [6, Theorem 7.5]):

**Theorem 9.6.** *A real filter  $g_0[n]$  satisfying*

$$|G_0(e^{j\omega})|^2 + |G_0(e^{j(\omega+\pi)})|^2 = 2. \quad (9.88)$$

*and having  $p$  zeros at  $\omega = \pi$  must have at least  $2p$  non-zero coefficients. Daubechies filters have exactly  $2p$  non-zero coefficients.*

Using point D1., this means that also the corresponding scaling function as well as the wavelet must have support at least  $2p - 1$ . The theorem follows from a deep result about polynomials called *Bezout theorem*.

*Proof.* Recall that we can write

$$G_0(e^{j\omega}) = \sqrt{2} \left( \frac{1 + e^{-j\omega}}{2} \right)^p R(e^{-j\omega}), \quad (9.89)$$

or, equivalently,

$$|G_0(e^{j\omega})|^2 = 2 \left( \cos \frac{\omega}{2} \right)^{2p} |R(e^{j\omega})|^2. \quad (9.90)$$

But now, since  $g_0[n]$  is real valued,  $|G_0(e^{j\omega})|^2$  is an even function of  $\omega$  and can therefore be expressed as a polynomial in  $\cos \omega$ . Using the trigonometric identity  $\sin^2(\frac{\omega}{2}) = (1 - \cos \omega)/2$ , we can equivalently write  $|R(e^{j\omega})|^2$  as a polynomial in  $\sin^2(\frac{\omega}{2})$ , and hence, we obtain

$$|G_0(e^{j\omega})|^2 = 2 \left( \cos \frac{\omega}{2} \right)^{2p} P(\sin^2 \frac{\omega}{2}), \quad (9.91)$$

for some yet to be determined polynomial  $P(\cdot)$ . The next trick is to recall that  $\sin^2(\frac{\omega}{2}) + \cos^2(\frac{\omega}{2}) = 1$ , and thus,

$$|G_0(e^{j\omega})|^2 = 2 \left( 1 - \sin^2 \frac{\omega}{2} \right)^p P(\sin^2 \frac{\omega}{2}). \quad (9.92)$$

To make notation simpler, we define  $y = \sin^2(\omega/2)$ . Plugging this back into Equation (9.88), this leads to the condition

$$(1 - y)^p P(y) + y^p P(1 - y) = 1. \quad (9.93)$$

In order to minimize the coefficients of the filter  $g_0[n]$ , we want to find the polynomial  $P(y)$  with the smallest degree that satisfies this equation. The Bezout theorem tells us that the smallest-degree polynomial  $P(y)$  has degree  $p - 1$ . This implies that  $G_0(e^{j\omega})$  has degree  $2p - 1$ , which means it has  $2p$  coefficients.  $\square$

## 9.6 Data-adaptive Signal Representations

Fourier, short-time Fourier and wavelet bases are picked based on fundamental considerations (physics, mathematics, general structure). By contrast, let us now study a scenario where we have ample data (and computational power). Then, it is tempting to conjecture that appropriate signal representations can be found simply by looking at the data.

Specifically, let us now consider the setting where we are given a (large) set of  $n$  data points in potentially high dimension  $k$ . Denote the data points by  $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(n)}$ , where each data point  $\mathbf{x}^{(j)}$  is a  $k$ -dimensional vector. The canonical approach to finding an (orthonormal) basis is the Gram-Schmidt procedure : We arbitrarily select a first sample, normalize it, and take this to be our first basis element. Then, we arbitrarily select a second sample and subtract from it its projection onto the first basis element, and so on. It should be intuitively clear that practically speaking, this is not a very desirable procedure. The outcome will in general strongly depend on the order in which the samples are processed, and it will lack any interesting structural properties. The procedure will work fine if the data samples all lie exactly in a  $p$ -dimensional subspace (with  $p < k$ ), but if this is only approximately true, the outcome will generally suffer.

Instead, we will now consider an alternative procedure. Here, we start by fixing a certain  $p < k$ . The goal is to find a good  $p$ -dimensional basis such that most of the data points can be represented quite accurately in terms of this basis. It is not initially clear what “most” and “quite accurately” should mean. One intuitively pleasing metric is to select the basis (and corresponding coefficients for each data sample) so as to minimize the overall mean-squared error, that is:

$$\sum_{j=1}^n \|\mathbf{x}^{(j)} - \widehat{\mathbf{x}^{(j)}}\|^2, \quad (9.94)$$

where  $\widehat{\mathbf{x}^{(j)}}$  represents the best approximation to  $\mathbf{x}^{(j)}$  within the chosen basis. In spite of appearance, this problem actually has a clean solution : This is precisely the Eckart-Young theorem.

To see this, let us denote the (yet unknown) basis vectors by  $\phi_1, \phi_2, \dots, \phi_p$ , and collect them (as column vectors) into the  $k \times p$  matrix

$$\Phi = \begin{pmatrix} \phi_1 & \phi_2 & \cdots & \phi_p \end{pmatrix}. \quad (9.95)$$

Then, we can express

$$\widehat{\mathbf{x}^{(j)}} = \Phi \mathbf{f}^{(j)}, \quad (9.96)$$

where  $\mathbf{f}^{(j)}$  is the feature vector corresponding to data sample  $\mathbf{x}^{(j)}$ . Hence, we are looking for

$$\min_{\text{feature vectors } \{\mathbf{f}^{(j)}\}_{j=1}^n \in \mathbb{C}^p} \left\{ \min_{\text{basis vectors } \{\phi_i\}_{i=1}^p \in \mathbb{C}^k} \sum_{j=1}^n \|\mathbf{x}^{(j)} - \Phi \mathbf{f}^{(j)}\|^2 \right\}. \quad (9.97)$$

To see how to proceed, we can rewrite

$$\begin{aligned}
\sum_{j=1}^n \|\mathbf{x}^{(j)} - \Phi \mathbf{f}^{(j)}\|^2 &= \sum_{j=1}^n \text{trace} \left( (\mathbf{x}^{(j)} - \Phi \mathbf{f}^{(j)})(\mathbf{x}^{(j)} - \Phi \mathbf{f}^{(j)})^H \right) \\
&= \text{trace} \left( \sum_{j=1}^n (\mathbf{x}^{(j)} - \Phi \mathbf{f}^{(j)})(\mathbf{x}^{(j)} - \Phi \mathbf{f}^{(j)})^H \right) \\
&= \text{trace} \left( (X - \Phi F)(X - \Phi F)^H \right) \\
&= \|X - \Phi F\|_F^2,
\end{aligned} \tag{9.98}$$

where we have collected all the data samples into the  $k \times n$  matrix  $X$  and all the feature vectors into the  $p \times n$  matrix  $F$ .

This problem is precisely addressed by the Eckart–Young theorem that we have discussed earlier. The answer is simply to determine the SVD of the matrix  $X$ , and retain only the  $p$  largest singular values along with their corresponding singular vectors. Explicitly:

$$X = U \Sigma V^H = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^H, \tag{9.99}$$

where  $r$  is the rank of the matrix  $X$  and where, as always, we assume that the singular values are ordered in decreasing order. Then, from the Eckart–Young theorem, we know that our error criterion is minimized if

$$\Phi F = \sum_{i=1}^p \sigma_i \mathbf{u}_i \mathbf{v}_i^H. \tag{9.100}$$

In other words, we may select our basis vectors of length  $k$  to be the left singular vectors of  $X$  (that is, the eigenvectors of  $XX^H$ ),

$$\phi_1 = \mathbf{u}_1, \quad \phi_2 = \mathbf{u}_2, \quad \dots, \quad \phi_p = \mathbf{u}_p, \tag{9.101}$$

in which case the matrix of feature vectors (of dimension  $p \times n$ ) is given by

$$F = \begin{pmatrix} \sigma_1 \mathbf{v}_1^H \\ \sigma_2 \mathbf{v}_2^H \\ \vdots \\ \sigma_p \mathbf{v}_p^H \end{pmatrix}. \tag{9.102}$$

For example, the first column of this matrix is the feature vector corresponding to the first data sample,  $\mathbf{x}^{(1)}$ . Of course, this feature vector (of length  $p$ ) can also be found by projecting the data sample successively into the  $p$  basis elements  $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_p$ . It is left as an exercise to the reader to show that this indeed leads to the same answer.

Let us also remark that this is, quite obviously, not the unique basis — we can always rotate the basis to find an alternative basis (spanning exactly the same space). This will change the feature vectors, but it will not change the approximation quality.

We should note that this is precisely the PCA (Principal Components Analysis) or Karhunen–Loève transform (KLT) that you have quite possibly encountered elsewhere. In the PCA, we first find the covariance matrix of the data samples (using the above notation, this is the matrix  $XX^H$ ), and then find its eigendecomposition. Clearly, the eigenvectors of the matrix  $XX^H$  are precisely the vectors  $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k$  above, and the PCA stipulates to use the eigenvectors corresponding to the  $p$  largest eigenvalues of  $XX^H$  as the approximate basis — exactly the same solution as the one we found above. However, we believe that the above is a much more insightful (and convincing) derivation than the “derivation” given in many texts on Machine Learning. Of course, that’s just our own biased opinion.

### 9.6.1 Example : word2vec

As a concrete (rather advanced) example, we now study the **word2vec** machinery developed by Mikolov *et al.* in [11]. This is a way of representing words by real-valued vectors.

#### The word2vec construction

We suppose that we have a large corpus of sentences in the English language (or your language of preference). We isolate individual words. For English, this leads to about 1’000’000 words. Every word  $j$  will be represented by a (column) vector  $\mathbf{x}^{(j)}$  of length 1’000’000: The entry  $x_i^{(j)}$  captures the relationship from word  $j$  to word  $i$ . Specifically, the value  $x_i^{(j)}$  is the *count* of how many times word  $i$  is observed in “close vicinity” of word  $j$  in the corpus of sentences. The definition of “close vicinity” can be tweaked in a number of ways. For example, we may say that word  $i$  is observed in close vicinity of word  $j$  if it is found within a window of 5 words around word  $j$ . (Asymmetric windows can also be defined, as can be weighted.)

Clearly, most vectors  $\mathbf{x}^{(j)}$  are rather sparsely populated, and may have a few large entries and several much smaller ones (representing the rarer occasions of use of word  $j$ ).

Next, we form the matrix (of dimensions 1’000’000 by 1’000’000) :

$$X = \begin{pmatrix} \mathbf{x}^{(1)} & \mathbf{x}^{(2)} & \dots & \mathbf{x}^{(1'000'000)} \end{pmatrix} \quad (9.103)$$

This matrix is sparsely populated.

The next step is to form the so-called *PMI matrix*, for *pointwise mutual information*. The usefulness of this step is rather mysterious and not fully understood. Specifically, we form the matrix  $P$  which is of the same dimensions as  $X$  and whose entry in the  $i$ th row and  $j$ th column is

$$P_{ij} = \log \frac{X_{ij}}{(\sum_{\ell} X_{i\ell}) (\sum_k X_{kj})}. \quad (9.104)$$

Empirically, it is observed that this matrix is *low-rank*.

## 9.7 Problems

**Problem 9.1** (Some review problems on linear algebra). (a) (*Frobenius norm*) Prove that  $\|A\|_F^2 = \text{trace}(A^H A)$ .

(b) (*Singular Value Decomposition*) Let  $\sigma_i(A)$  denote the  $i^{\text{th}}$  singular value of an  $m \times n$  matrix  $A$ . Prove that  $\|A\|_F^2 = \sum_{i=1}^{\min\{m,n\}} \sigma_i^2(A)$ .

(c) (*Projection Matrices*) Consider a set of  $k$  orthonormal vectors in  $\mathbb{C}^n$ , denoted by  $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k$ . The projection matrix (that projects an arbitrary vector into the subspace spanned by these orthonormal vectors) is given by

$$P = \sum_{i=1}^k \mathbf{u}_i \mathbf{u}_i^H. \quad (9.105)$$

- Prove that this matrix is *Hermitian*, i.e.,  $P^H = P$ .
- Prove that this matrix is *idempotent*, i.e.,  $P^2 = P$ . (In words, projecting twice into the same subspace is the same as projecting only once.)
- Prove that  $\text{trace}(P) = k$ , i.e., equal to the dimension of the subspace.
- Prove that the diagonal entries of  $P$  must be real-valued and non-negative. Then, prove that the diagonal entries of  $P$  cannot be larger than 1 (this is a little more tricky).

**Problem 9.2** (Eckart–Young Theorem). In these lecture notes, we show the proof of the converse part of the Eckart–Young theorem for the spectral norm. In this problem, you do the same for the case of the Frobenius norm.

(a) For any matrix  $A$  of dimension  $m \times n$  and an arbitrary orthonormal basis  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  of  $\mathbb{C}^n$ , prove that

$$\|A\|_F^2 = \sum_{k=1}^n \|A\mathbf{x}_k\|^2. \quad (9.106)$$

(b) Consider any  $m \times n$  matrix  $B$  with  $\text{rank}(B) \leq p$ . Clearly, its null space has dimension no smaller than  $n - p$ . Therefore, we can find an orthonormal set  $\{\mathbf{x}_1, \dots, \mathbf{x}_{n-p}\}$  in the null space of  $B$ . Prove that for such vectors, we have

$$\|A - B\|_F^2 \geq \sum_{k=1}^{n-p} \|A\mathbf{x}_k\|^2. \quad (9.107)$$

(c) (*This requires slightly more subtle manipulations.*) For any matrix  $A$  of dimension  $m \times n$  and any orthonormal set of  $n - p$  vectors in  $\mathbb{C}^n$ , denoted by  $\{\mathbf{x}_1, \dots, \mathbf{x}_{n-p}\}$ , prove that

$$\sum_{k=1}^{n-p} \|A\mathbf{x}_k\|^2 \geq \sum_{j=p+1}^r \sigma_j^2. \quad (9.108)$$

*Hint:* Consider the case  $m \geq n$  and the set of vectors  $\{\mathbf{z}_1, \dots, \mathbf{z}_{n-p}\}$ , where  $\mathbf{z}_k = V^H \mathbf{x}_k$ . Express your formulas in terms of these and the SVD representation  $A = U\Sigma V^H$ .

(d) Briefly explain how (a)-(c) imply the desired statement.

**Problem 9.3** (The Fourier matrix diagonalizes all circulant matrices). The discrete Fourier transform (DFT)  $\mathbf{X}$  of the vector  $\mathbf{x}$  is given by

$$\mathbf{X} = W\mathbf{x} \quad \text{and} \quad \mathbf{x} = \frac{1}{N}W^H\mathbf{X}. \quad (9.109)$$

In this homework problem, you will prove that the Fourier matrix diagonalizes all *circulant* matrices.

(a) To cut the derivation into two simpler steps, we introduce an auxiliary matrix  $M$ , defined as

$$M = WA = W \underbrace{\begin{pmatrix} b_0 & b_{N-1} & b_{N-2} & b_{N-3} & \dots & b_1 \\ b_1 & b_0 & b_{N-1} & b_{N-2} & \dots & b_2 \\ b_2 & b_1 & b_0 & b_{N-1} & \dots & b_3 \\ b_3 & b_2 & b_1 & b_0 & \dots & b_4 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ b_{N-1} & b_{N-2} & b_{N-3} & b_{N-4} & \dots & b_0 \end{pmatrix}}_{\text{This is a circulant matrix}}. \quad (9.110)$$

Let us denote the unitary DFT of the sequence  $\{b_0, b_1, \dots, b_{N-1}\}$  by  $\{B_0, B_1, \dots, B_{N-1}\}$ . Write out the matrix  $M$  in terms of  $\{B_0, B_1, \dots, B_{N-1}\}$ . *Hint:* The first column of the matrix  $M$  is simply given by

$$W \begin{pmatrix} b_0 \\ b_1 \\ b_2 \\ b_3 \\ \vdots \\ b_{N-1} \end{pmatrix} = \begin{pmatrix} B_0 \\ B_1 \\ B_2 \\ B_3 \\ \vdots \\ B_{N-1} \end{pmatrix} \quad (9.111)$$

To find the second column, you will need to use some Fourier properties.

(b) Using the matrix  $M$  from above, compute the full matrix product

$$WAW^H = MW^H. \quad (9.112)$$

*Hint:* Handle every *row* of the matrix  $M$  separately. Define the vector  $\mathbf{m}$  such that  $\mathbf{m}^H$  is simply the first row of the matrix  $M$ . But the product  $\mathbf{m}^H W^H$  is easily computed, recalling that  $\mathbf{m}^H W^H = (W\mathbf{m})^H$ .

**Problem 9.4** (Inner Products). Consider the standard  $n$ -dimensional vector space  $\mathbb{R}^n$ .

1. Characterize the set of matrices  $W$  for which  $\mathbf{y}^T W \mathbf{x}$  is a valid inner product for any  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ .
2. Prove that *every* inner product  $\langle \mathbf{x}, \mathbf{y} \rangle$  on  $\mathbb{R}^n$  can be expressed as  $\mathbf{y}^T W \mathbf{x}$  for an appropriately chosen matrix  $W$ .
3. For a subspace of dimension  $k < n$ , spanned by the basis  $\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_k \in \mathbb{R}^n$ , express the orthogonal projection operator (matrix) with respect to the general inner product  $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{y}^T W \mathbf{x}$ . *Hint:* For any vector  $\mathbf{x} \in \mathbb{R}^n$ , express its projection as  $\hat{\mathbf{x}} = \sum_{j=1}^k \alpha_j \mathbf{b}_j$ .



# Bibliography

- [1] J. Duchi, *Lecture Notes for Statistics 311/Electrical Engineering 377*. Stanford, 2016.
- [2] M. J. Wainwright and M. I. Jordan, *Graphical Models, Exponential Families and Variational Inference*. NOW Publishers, 2008.
- [3] B. Widrow and M. E. Hoff, “Adaptive switching circuits,” in *IRE WESON Convention Record*, vol. 4, (Los Angeles, CA), pp. 96–104, 1960.
- [4] M. H. Hayes, *Statistical Digital Signal Processing and Modelling*. Wiley, 1996.
- [5] S. O. Haykin, *Adaptive Filter Theory*. Prentice Hall, 5th ed., 2013.
- [6] S. Mallat, *A Wavelet Tour of Signal Processing*. Academic Press, 1998.
- [7] P. Brémaud, *Mathematical Principles of Signal Processing*. New York: Springer-Verlag, 2002.
- [8] M. Vetterli, J. Kovacevic, and V. Goyal, *Foundations of Signal Processing*. Cambridge University Press, 2014.
- [9] P. Prandoni and M. Vetterli, *Signal Processing for Communications*. EPFL Press, 2008. Available online at <http://www.sp4comm.org>.
- [10] M. Vetterli and J. Kovacevic, *Wavelets and Subband Coding*. Prentice Hall, 1995.
- [11] T. Mikolov, K. Chen, G. Corrado, and J. Dean, “Efficient estimation of word representations in vector space.” Preprint available at <https://arxiv.org/abs/1301.3781>, 2013.