# KERNEL METHODS
# AND THE
# CURSE OF DIMENSIONALITY

**Stefano Spigler**

Jonas Paccolat, Mario Geiger, Matthieu Wyart

EPFL

# SUPERVISED DEEP LEARNING

- Why and how does deep **supervised** learning work?


- Learn from examples: **how many** are needed?


- Typical tasks:

  - Regression (fitting functions)

  - Classification

# LEARNING CURVES

- Performance is evaluated through the **generalization error** $\epsilon$

- Learning curves decay with number of examples $n$, often as

$$\epsilon \sim n^{-\beta}$$

- $\beta$ depends on the **dataset** and on the **algorithm**

Deep networks: $\beta \sim 0.07\text{-}0.35$ [Hestness et al. 2017]

*We lack a theory for $\beta$ for deep networks!*
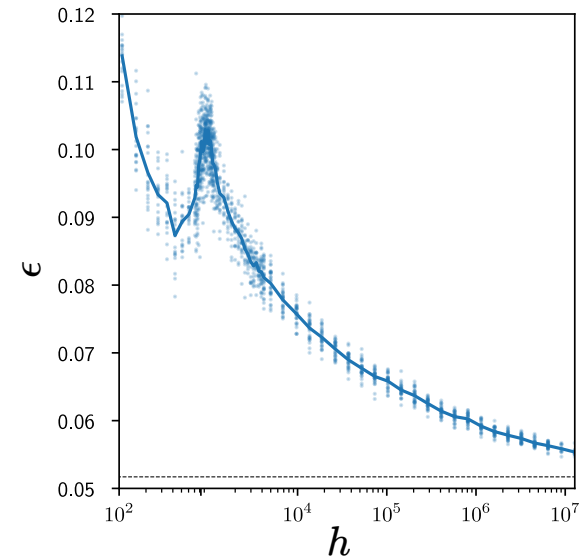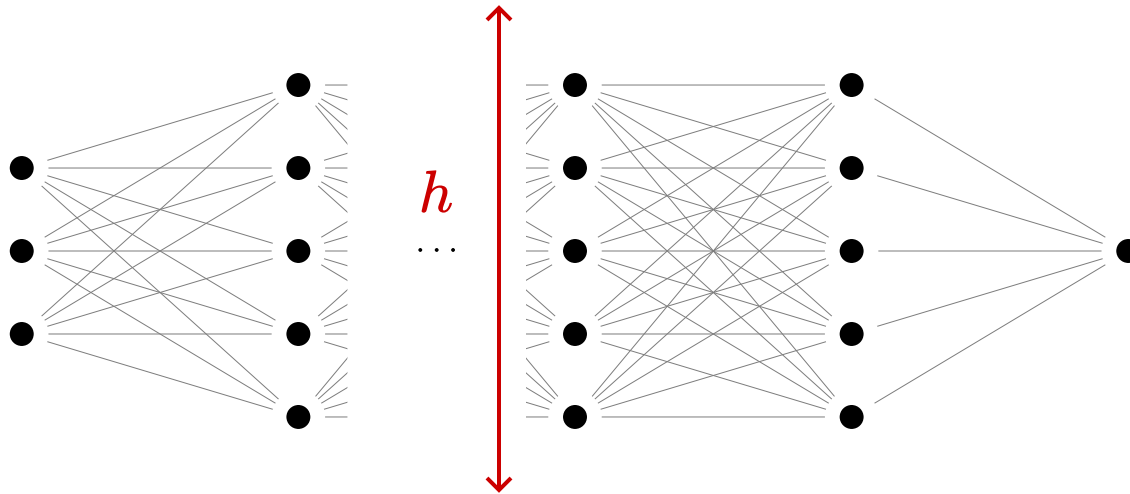
# LINK WITH KERNEL LEARNING

- Performance increases with **overparametrization**

[Neyshabur et al. 2017, 2018, Advani and Saxe 2017]
[Belkin et al. 2018, Spigler et al. 2018, Geiger et al. 2019]

$\longrightarrow$ study the infinite-width limit!

[Mei et al. 2017, Rotskoff and Vanden-Eijnden 2018, Jacot et al. 2018, Chizat and Bach 2018, ...]
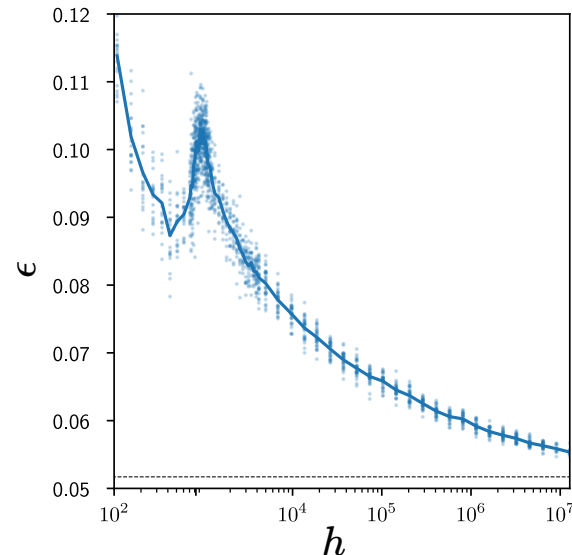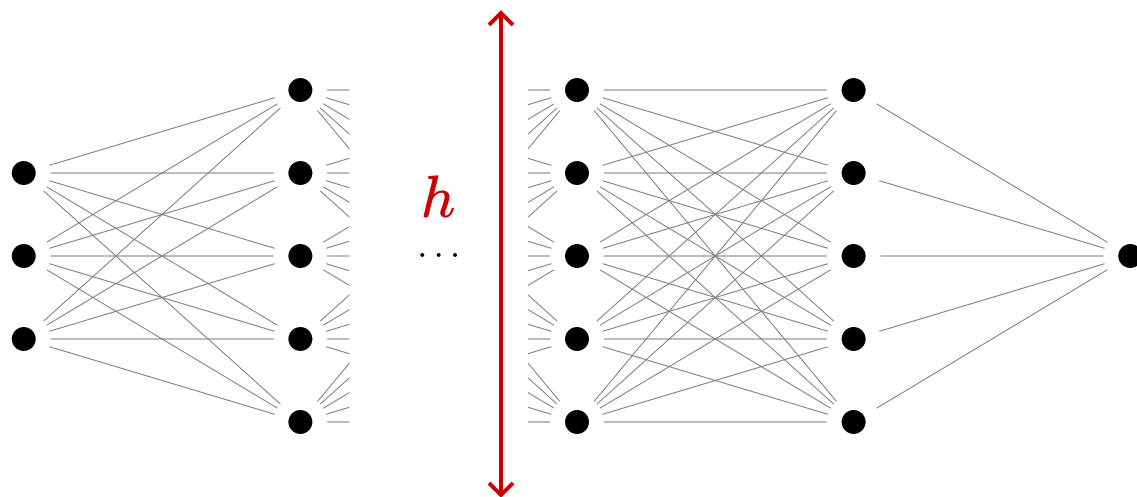
# LINK WITH KERNEL LEARNING

- Performance increases with **overparametrization**

[Neyshabur et al. 2017, 2018, Advani and Saxe 2017]
[Belkin et al. 2018, Spigler et al. 2018, Geiger et al. 2019]

$\longrightarrow$ study the infinite-width limit!

[Mei et al. 2017, Rotskoff and Vanden-Eijnden 2018, Jacot et al. 2018, Chizat and Bach 2018, ...]



- With a specific scaling, infinite-width limit $\rightarrow$ **kernel learning**

[Jacot et al. 2018]
Neural Tangent Kernel

(next slides)

**What are the learning curves of kernels like?**
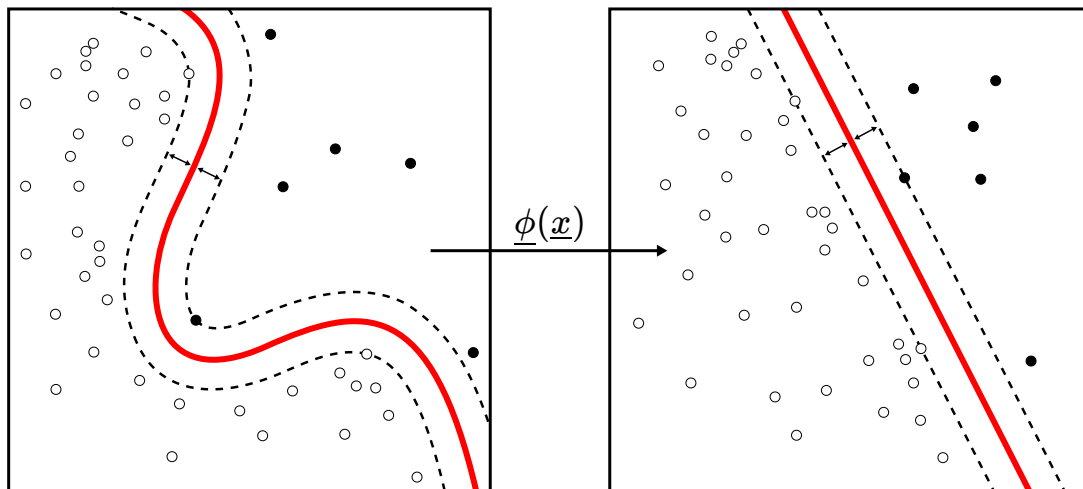
4

# OUTLINE

- Very brief introduction to kernel methods

- Performance of kernels on real data

- Gaussian data: Teacher-Student regression

- Gaussian approximation: smoothness and effective dimension

- Dimensional reduction via invariants in the task

# KERNEL METHODS

- Kernel methods learn non-linear functions or boundaries

- Map data to a **feature space**, where the problem is linear

data $\underline{x} \longrightarrow \underline{\phi}(\underline{x}) \longrightarrow$ use linear combination of features



$\underline{\phi}(\underline{x})$

# KERNEL METHODS

- Kernel methods learn non-linear functions or boundaries

- Map data to a **feature space**, where the problem is linear

data $\underline{x} \longrightarrow \underline{\phi}(\underline{x}) \longrightarrow$ use linear combination of features

only scalar products are needed: $\boxed{\underline{\phi}(\underline{x}) \cdot \underline{\phi}(\underline{x}')}$

kernel $K(\underline{x}, \underline{x}')$
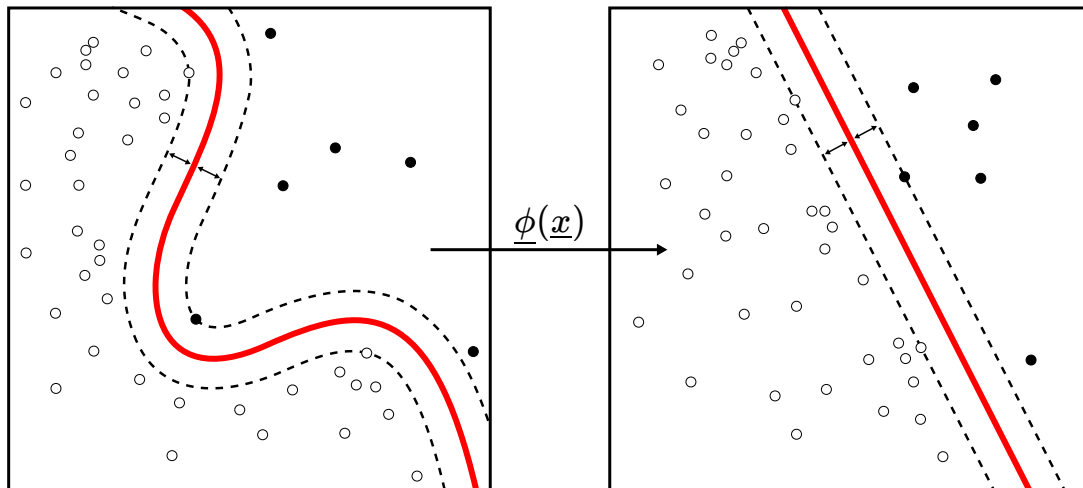


$\underline{\phi}(\underline{x})$

# KERNEL METHODS

- Kernel methods learn non-linear functions or boundaries

- Map data to a **feature space**, where the problem is linear

  data $\underline{x} \longrightarrow \underline{\phi}(\underline{x}) \longrightarrow$ use linear combination of features

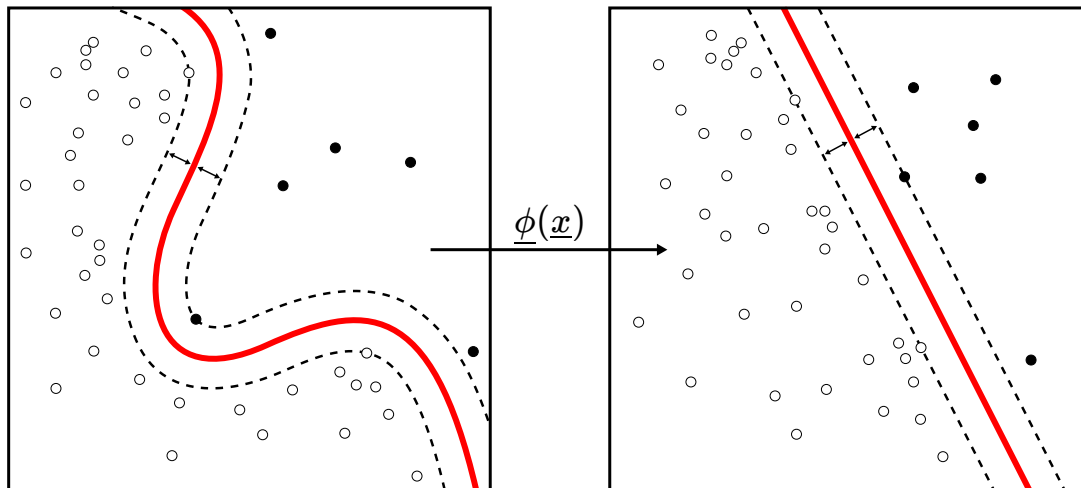  only scalar products are needed: $\boxed{\underline{\phi}(\underline{x}) \cdot \underline{\phi}(\underline{x}')}$

kernel $K(\underline{x}, \underline{x}')$

*Gaussian:*
$$K(\underline{x}, \underline{x}') = \exp\left(-\frac{\|\underline{x} - \underline{x}'\|^2}{\sigma^2}\right)$$

*Laplace:*
$$K(\underline{x}, \underline{x}') = \exp\left(-\frac{\|\underline{x} - \underline{x}'\|}{\sigma}\right)$$

$\underline{\phi}(\underline{x})$

6

# KERNEL REGRESSION

E.g. kernel regression:

- **Target function** $\underline{x}_\mu \to Z(\underline{x}_\mu), \ \mu = 1, \ldots, n$

# KERNEL REGRESSION

E.g. kernel regression:

- **Target function** $\underline{x}_\mu \to Z(\underline{x}_\mu), \ \mu = 1, \dots, n$

- Build an estimator $\hat{Z}_K(\underline{x}) = \sum_{\mu=1}^{n} c_\mu K(\underline{x}_\mu, \underline{x})$

# KERNEL REGRESSION

E.g. kernel regression:

- **Target function** $\underline{x}_\mu \to Z(\underline{x}_\mu), \ \ \mu = 1, \ldots, n$

- Build an estimator $\hat{Z}_K(\underline{x}) = \sum_{\mu=1}^{n} c_\mu K(\underline{x}_\mu, \underline{x})$

- Minimize training MSE $= \frac{1}{n} \sum_{\mu=1}^{n} \left[ \hat{Z}_K(\underline{x}_\mu) - Z(\underline{x}_\mu) \right]^2$

# KERNEL REGRESSION

E.g. kernel regression:

- **Target function** $\underline{x}_\mu \to Z(\underline{x}_\mu), \ \mu = 1, \ldots, n$

- Build an estimator $\hat{Z}_K(\underline{x}) = \sum_{\mu=1}^{n} c_\mu K(\underline{x}_\mu, \underline{x})$

- Minimize training MSE $= \frac{1}{n} \sum_{\mu=1}^{n} \left[ \hat{Z}_K(\underline{x}_\mu) - Z(\underline{x}_\mu) \right]^2$

- Estimate the **generalization error** $\epsilon = \mathbb{E}_{\underline{x}} \left[ \hat{Z}_K(\underline{x}) - Z(\underline{x}) \right]^2$

# REPRODUCING KERNEL HILBERT SPACE (RKHS)

A kernel $K$ induces a corresponding Hilbert space $\mathcal{H}_K$ with norm

$$\|Z\|_K = \int \mathrm{d}\underline{x}\mathrm{d}\underline{y}\, Z(\underline{x})K^{-1}(\underline{x},\underline{y})Z(\underline{y})$$

where $K^{-1}(\underline{x},\underline{y})$ is such that

$$\int \mathrm{d}\underline{y}\, K^{-1}(\underline{x},\underline{y})K(\underline{y},\underline{z}) = \delta(\underline{x},\underline{z})$$

$\mathcal{H}_K$ is called the **Reproducing Kernel Hilbert Space** (RKHS)

# PREVIOUS WORKS

*Regression:* **performance depends on the target function!**

# PREVIOUS WORKS

*Regression:* **performance depends on the target function!**

$d$ = dimension of the input space

- If only assumed to be **Lipschitz**, then $\beta = \frac{1}{d}$

  **Curse of dimensionality!**          [Luxburg and Bousquet 2004]

# PREVIOUS WORKS

*Regression:* **performance depends on the target function!**

$d$ = dimension of the input space

- If only assumed to be **Lipschitz**, then $\beta = \frac{1}{d}$

    **Curse of dimensionality!**                    [Luxburg and Bousquet 2004]

- If assumed to be in the **RKHS**, then $\beta \geq \frac{1}{2}$ does not depend on $d$

    [Smola et al. 1998, Rudi and Rosasco 2017]

*Regression:* **performance depends on the target function!**

- If only assumed to be **Lipschitz**, then $\beta = \frac{1}{d}$

  **Curse of dimensionality!**                    [Luxburg and Bousquet 2004]

- If assumed to be in the **RKHS**, then $\beta \geq \frac{1}{2}$ does not depend on $d$

  [Smola et al. 1998, Rudi and Rosasco 2017]

- Yet, RKHS is a very strong assumption on the smoothness of the target function (see later on)

  [Bach 2017]

# REAL DATA AND ALGORITHMS
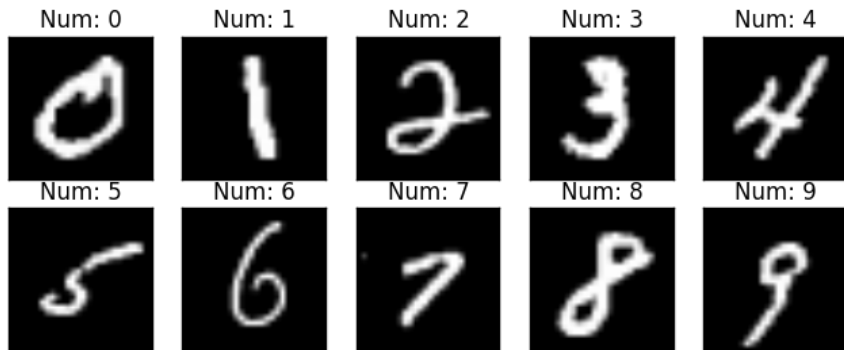
## We apply kernel methods on

| MNIST | CIFAR10 |
|---|---|
| 2 classes: even/odd | 2 classes: first 5/last 5 |
| 70000 28x28 b/w pictures | 60000 32x32 RGB pictures |
| ↓ | ↓ |
| dimension $d = 784$ | dimension $d = 3072$ |



We perform
$$\begin{cases} \textbf{regression} & \longrightarrow \quad \text{kernel regression} \\ \\ \textbf{classification} & \longrightarrow \quad \text{margin SVM} \end{cases}$$

- Same exponent for regression and classification

- Same exponent for Gaussian and Laplace kernel

- MNIST and CIFAR10 display exponents $\beta \gg \frac{1}{d}$ but $< \frac{1}{2}$

11

Regression on MNIST — Gaussian, Laplace, $n^{-0.37}$ — $\beta \approx 0.4$

Regression on CIFAR10 — Gaussian, Laplace, $n^{-0.08}$ — $\beta \approx 0.1$

Classification on MNIST

Classification on CIFAR10 — Gaussian, Laplace, $n^{-0.10}$

We need a **new framework**!

- Same exponent for regression and classification

- Same exponent for Gaussian and Laplace kernel

- MNIST and CIFAR10 display exponents $\beta \gg \frac{1}{d}$ but $< \frac{1}{2}$

11

# KERNEL TEACHER-STUDENT FRAMEWORK

- Controlled setting: **Teacher-Student regression**

# KERNEL TEACHER-STUDENT FRAMEWORK

- Controlled setting: **Teacher-Student regression**

- Training data are sampled from a Gaussian Process:

$$Z_T(\underline{x}_1), \ldots, Z_T(\underline{x}_n) \ \sim \ \mathcal{N}(0, K_T)$$

$\underline{x}_\mu$ are random on a $d$-**dim hypersphere**

# KERNEL TEACHER-STUDENT FRAMEWORK

- Controlled setting: **Teacher-Student regression**

- Training data are sampled from a Gaussian Process:

$$Z_T(\underline{x}_1), \ldots, Z_T(\underline{x}_n) \sim \mathcal{N}(0, K_T)$$

$\underline{x}_\mu$ are random on a $d$-**dim hypersphere**

$$\mathbb{E}Z_T(\underline{x}_\mu) = 0$$

$$\mathbb{E}Z_T(\underline{x}_\mu)Z_T(\underline{x}_\nu) = K_T(\|\underline{x}_\mu - \underline{x}_\nu\|)$$

# KERNEL TEACHER-STUDENT FRAMEWORK

- Controlled setting: **Teacher-Student regression**

- Training data are sampled from a Gaussian Process:

$$Z_T(\underline{x}_1), \ldots, Z_T(\underline{x}_n) \sim \mathcal{N}(0, K_T)$$

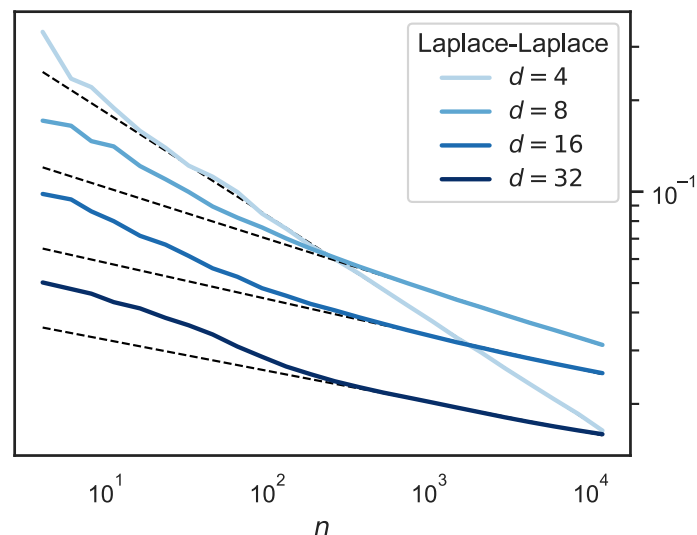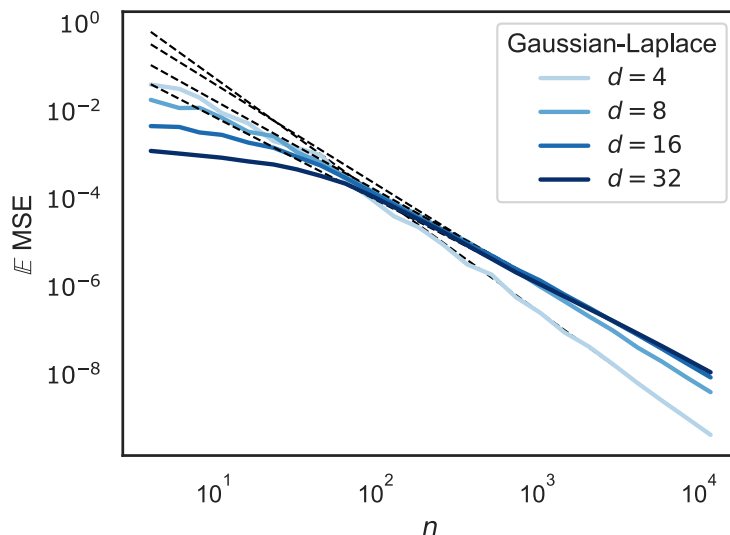  $\underline{x}_\mu$ are random on a **$d$-dim hypersphere**

$$\mathbb{E}Z_T(\underline{x}_\mu) = 0$$

$$\mathbb{E}Z_T(\underline{x}_\mu)Z_T(\underline{x}_\nu) = K_T(\|\underline{x}_\mu - \underline{x}_\nu\|)$$
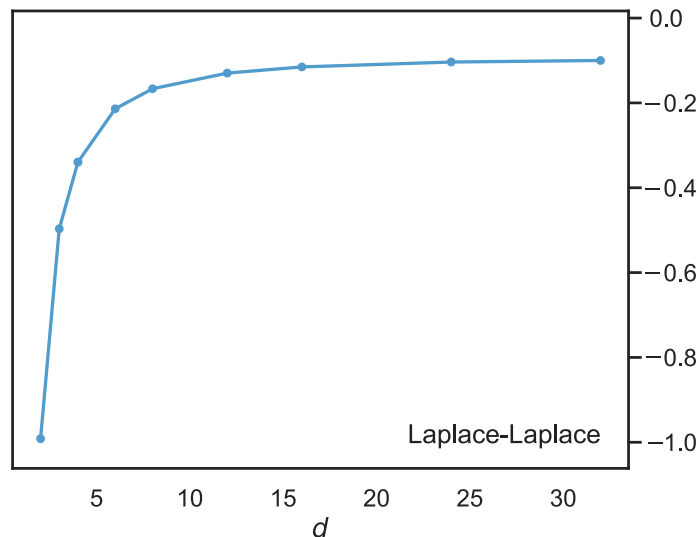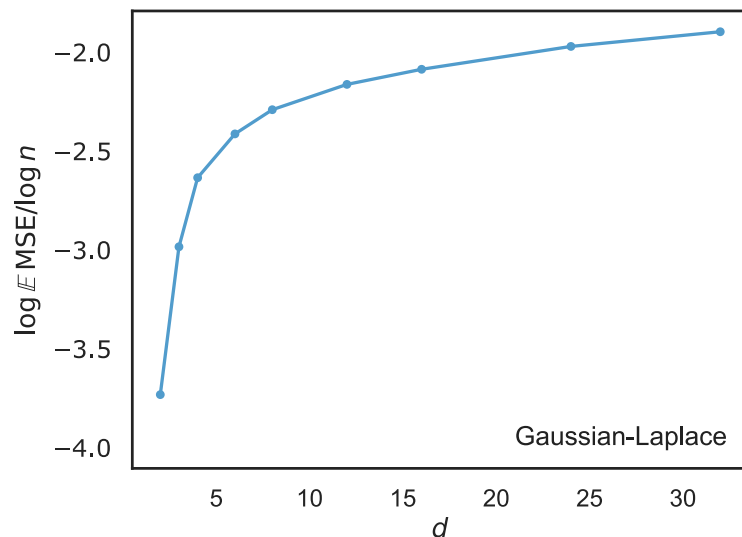
- Regression is done with another kernel $K_S$

Generalization error



Exponent $-\beta$

*Can we understand these curves?*

# TEACHER-STUDENT: REGRESSION

Regression:

$$\hat{Z}_S(\underline{x}) = \sum_{\mu=1}^{n} c_\mu K_S(\underline{x}_\mu, \underline{x})$$

$$\text{Minimize} = \frac{1}{n} \sum_{\mu=1}^{n} \left[ \hat{Z}_S(\underline{x}_\mu) - Z_T(\underline{x}_\mu) \right]^2$$

Explicit solution:

$$\hat{Z}_S(\underline{x}) = \underline{k}_S(\underline{x}) \cdot \mathbb{K}_S^{-1} \underline{Z} \qquad \text{where}$$

# TEACHER-STUDENT: REGRESSION

Regression:
$$\hat{Z}_S(\underline{x}) = \sum_{\mu=1}^{n} c_\mu K_S(\underline{x}_\mu, \underline{x})$$
$$\text{Minimize} = \frac{1}{n} \sum_{\mu=1}^{n} \left[ \hat{Z}_S(\underline{x}_\mu) - Z_T(\underline{x}_\mu) \right]^2$$

Explicit solution:

$$\hat{Z}_S(\underline{x}) = \underline{k}_S(\underline{x}) \cdot \mathbb{K}_S^{-1} \underline{Z}_T \qquad \text{where} \qquad \begin{cases} (\underline{k}_S(\underline{x}))_\mu = K_S(\underline{x}_\mu, \underline{x}) \\ \qquad\qquad\qquad \text{kernel overlap} \end{cases}$$

# TEACHER–STUDENT: REGRESSION

Regression:
$$\hat{Z}_S(\underline{x}) = \sum_{\mu=1}^{n} c_\mu K_S(\underline{x}_\mu, \underline{x})$$

$$\text{Minimize} = \frac{1}{n} \sum_{\mu=1}^{n} \left[ \hat{Z}_S(\underline{x}_\mu) - Z_T(\underline{x}_\mu) \right]^2$$

Explicit solution:

$$\hat{Z}_S(\underline{x}) = \underline{k}_S(\underline{x}) \cdot \mathbb{K}_S^{-1} \underline{Z}_T$$

where

$$(\underline{k}_S(\underline{x}))_\mu = K_S(\underline{x}_\mu, \underline{x})$$

kernel overlap

$$(\mathbb{K}_S)_{\mu\nu} = K_S(\underline{x}_\mu, \underline{x}_\nu)$$

Gram matrix

# TEACHER-STUDENT: REGRESSION

Regression:

$$\hat{Z}_S(\underline{x}) = \sum_{\mu=1}^{n} c_\mu K_S(\underline{x}_\mu, \underline{x})$$

$$\text{Minimize} = \frac{1}{n} \sum_{\mu=1}^{n} \left[ \hat{Z}_S(\underline{x}_\mu) - Z_T(\underline{x}_\mu) \right]^2$$

Explicit solution:

$$\hat{Z}_S(\underline{x}) = \underline{k}_S(\underline{x}) \cdot \mathbb{K}_S^{-1} \underline{Z}_T$$

where

$$\begin{cases} (\underline{k}_S(\underline{x}))_\mu = K_S(\underline{x}_\mu, \underline{x}) \\ \text{kernel overlap} \\ \\ (\mathbb{K}_S)_{\mu\nu} = K_S(\underline{x}_\mu, \underline{x}_\nu) \\ \text{Gram matrix} \\ \\ (\underline{Z}_T)_\mu = Z_T(\underline{x}_\mu) \\ \text{training data} \end{cases}$$

# TEACHER-STUDENT: REGRESSION

Regression:

$$\hat{Z}_S(\underline{x}) = \sum_{\mu=1}^n c_\mu K_S(\underline{x}_\mu, \underline{x})$$

$$\text{Minimize} = \frac{1}{n} \sum_{\mu=1}^n \left[ \hat{Z}_S(\underline{x}_\mu) - Z_T(\underline{x}_\mu) \right]^2$$

Explicit solution:

$$\hat{Z}_S(\underline{x}) = \underline{k}_S(\underline{x}) \cdot \mathbb{K}_S^{-1} \underline{Z}_T$$

where

$$(\underline{k}_S(\underline{x}))_\mu = K_S(\underline{x}_\mu, \underline{x})$$
kernel overlap

$$(\mathbb{K}_S)_{\mu\nu} = K_S(\underline{x}_\mu, \underline{x}_\nu)$$
Gram matrix

$$(\underline{Z}_T)_\mu = Z_T(\underline{x}_\mu)$$
training data

Compute the generalization error $\epsilon$ and how it scales with $n$

$$\epsilon = \mathbb{E}_T \int \mathrm{d}^d \underline{x} \left[ \hat{Z}_S(\underline{x}) - Z_T(\underline{x}) \right]^2 \sim n^{-\beta}$$

14

# TEACHER-STUDENT: THEOREM (1/2)

To compute the generalization error:

- We look at the problem in the **frequency domain**

- We assume that $\tilde{K}_S(\underline{w}) \sim \|\underline{w}\|^{-\alpha_S}$ and $\tilde{K}_T(\underline{w}) \sim \|\underline{w}\|^{-\alpha_T}$ as $\|\underline{w}\| \to \infty$

# TEACHER-STUDENT: THEOREM (1/2)

To compute the generalization error:

- We look at the problem in the **frequency domain**

- We assume that $\tilde{K}_S(\underline{w}) \sim \|\underline{w}\|^{-\alpha_S}$ and $\tilde{K}_T(\underline{w}) \sim \|\underline{w}\|^{-\alpha_T}$ as $\|\underline{w}\| \to \infty$

  E.g. Laplace has $\alpha = d + 1$ and Gaussian has $\alpha = \infty$

# TEACHER-STUDENT: THEOREM (1/2)

To compute the generalization error:

- We look at the problem in the **frequency domain**

- We assume that $\tilde{K}_S(\underline{w}) \sim \|\underline{w}\|^{-\alpha_S}$ and $\tilde{K}_T(\underline{w}) \sim \|\underline{w}\|^{-\alpha_T}$ as $\|\underline{w}\| \to \infty$

  E.g. Laplace has $\alpha = d + 1$ and Gaussian has $\alpha = \infty$

- **SIMPLIFYING ASSUMPTION:** We take the $n$ points $\underline{x}_\mu$ on a **regular** $d$-dim **lattice!**

# TEACHER-STUDENT: THEOREM (1/2)

To compute the generalization error:

- We look at the problem in the **frequency domain**

- We assume that $\tilde{K}_S(\underline{w}) \sim \|\underline{w}\|^{-\alpha_S}$ and $\tilde{K}_T(\underline{w}) \sim \|\underline{w}\|^{-\alpha_T}$ as $\|\underline{w}\| \to \infty$

  E.g. Laplace has $\alpha = d + 1$ and Gaussian has $\alpha = \infty$

- **SIMPLIFYING ASSUMPTION:** We take the $n$ points $\underline{x}_\mu$ on a **regular** $d$-dim **lattice!**

  (details: arXiv:1905.10843)

Then we can show that

for $n \gg 1$ $\qquad \epsilon \sim n^{-\beta}$ with $\boxed{\beta = \frac{1}{d} \min(\alpha_T - d, 2\alpha_S)}$

# TEACHER-STUDENT: THEOREM (2/2)

$$\beta = \frac{1}{d}\min(\alpha_T - d, 2\alpha_S)$$

- Large $\alpha \rightarrow$ fast decay at high freq $\rightarrow$ **indifference** to **local details**

- $\alpha_T$ is intrinsic to the **data** (T), $\alpha_S$ depends on the **algorithm** (S)

- If $\alpha_S$ is large enough, $\beta$ takes the largest possible value $\frac{\alpha_T - d}{d}$

  (optimal learning)

- As soon as $\alpha_S$ is small enough, $\beta = \frac{2\alpha_S}{d}$

# TEACHER-STUDENT: COMPARISON (1/2)

What is the prediction for our simulations?

$$\beta = \tfrac{1}{d} \min(\alpha_T - d, 2\alpha_S)$$

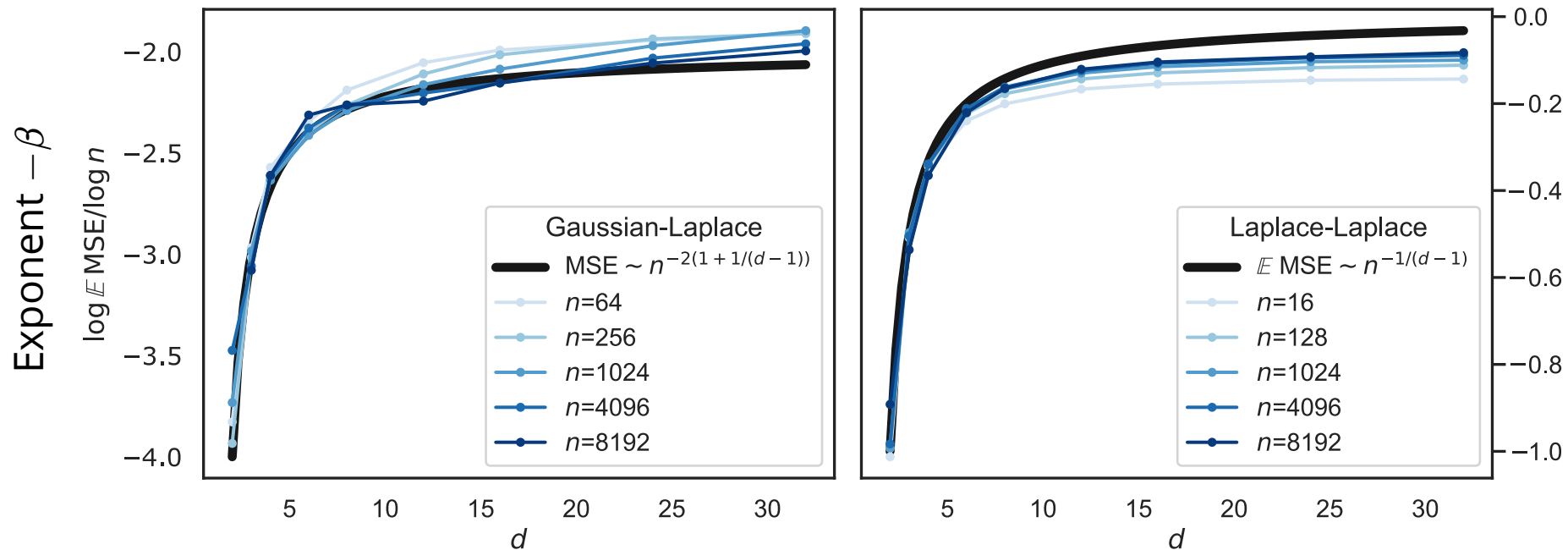- If Teacher=Student=Laplace $\qquad (\alpha_T = \alpha_S = d + 1)$

$$\beta = \frac{\alpha_T - d}{d} = \frac{1}{d} \qquad \text{(curse of dimensionality!)}$$

- If Teacher=Gaussian, Student=Laplace $\qquad (\alpha_T = \infty, \alpha_S = d + 1)$

$$\beta = \frac{2\alpha_S}{d} = 2 + \frac{2}{d}$$

- Our result matches the numerical simulations
<span style="color:red">(on hypersphere)</span>
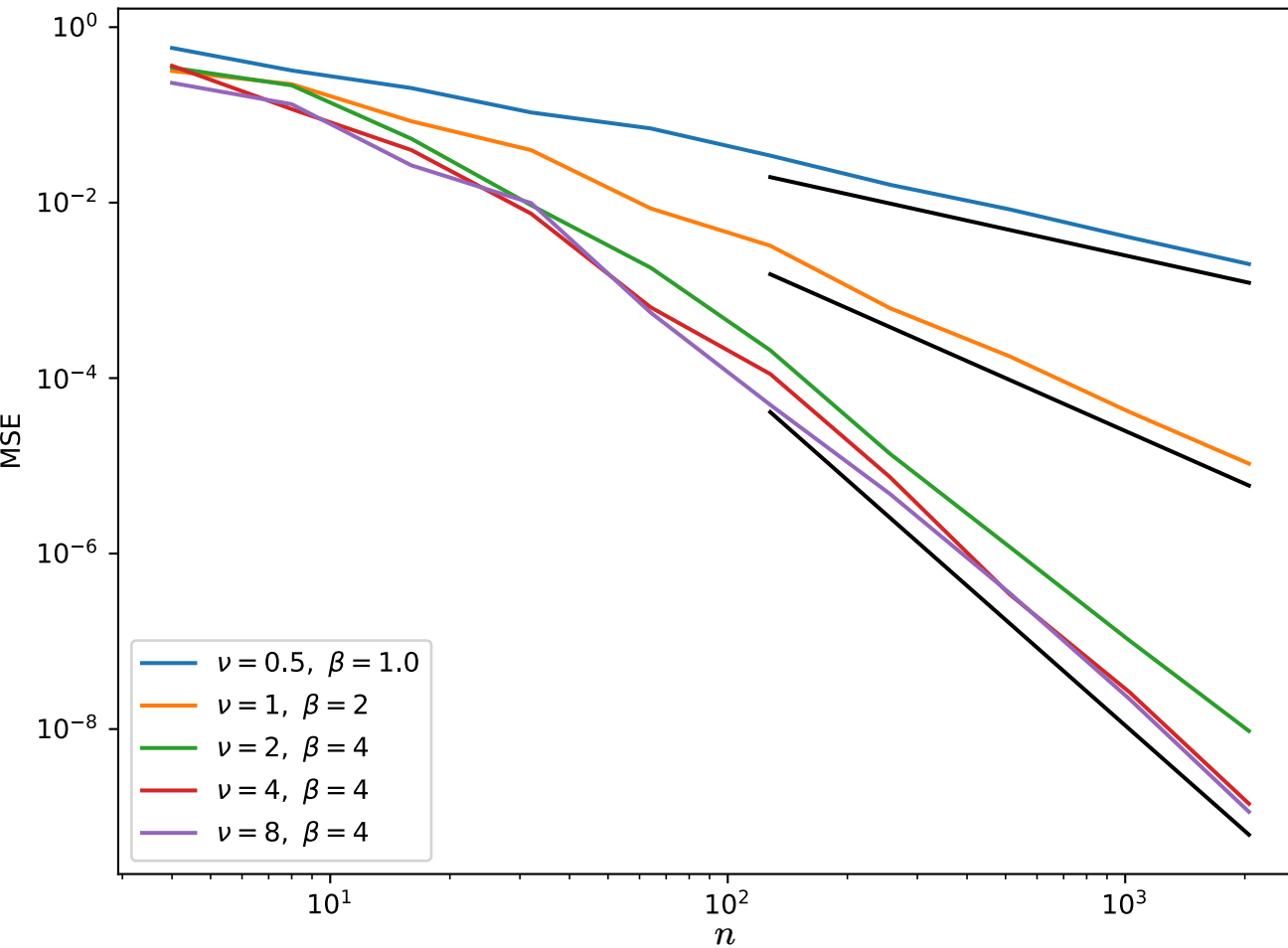
- There are finite size effects (small $n$)

# TEACHER–STUDENT: MATÉRN TEACHER

Matérn kernels: $K_T(\underline{x}) = \frac{2^{1-\nu}}{\Gamma(\nu)} z^\nu \mathcal{K}_\nu(z), \quad z = \sqrt{2\nu} \frac{\|x\|}{\sigma}, \quad \alpha = d + 2\nu$

Laplace student, $K_S(\underline{x}) = \exp\left(-\frac{\|x\|}{\sigma}\right)$

$d = 1$

$\beta = \min(2\nu, 4)$



Legend:
- $\nu = 0.5,\ \beta = 1.0$
- $\nu = 1,\ \beta = 2$
- $\nu = 2,\ \beta = 4$
- $\nu = 4,\ \beta = 4$
- $\nu = 8,\ \beta = 4$
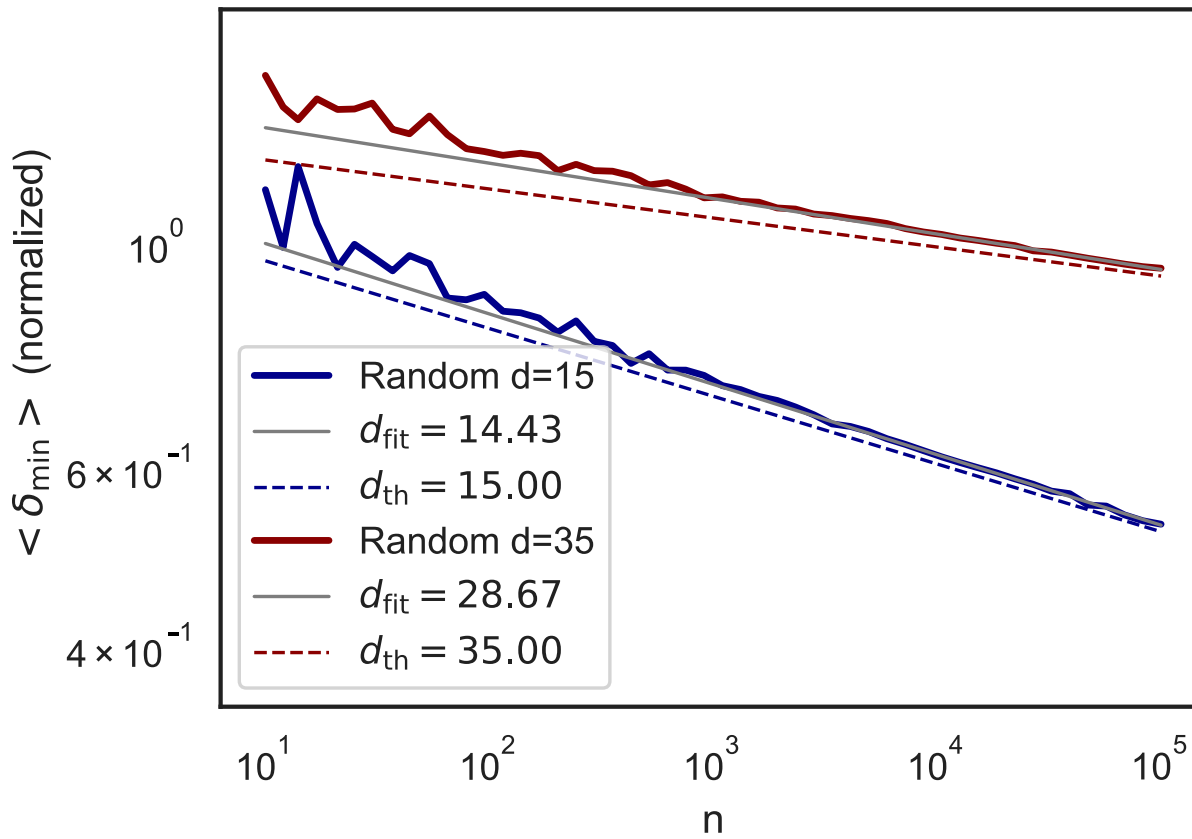
# NEAREST-NEIGHBOR DISTANCE

Same result with points on *regular lattice* or *random hypersphere*?

What matters is how **nearest-neighbor distance $\delta$ scales with $n$**

<span style="color:red">(conjecture)</span>



In both cases $\boxed{\delta \sim n^{\frac{1}{d}}}$

*Finite size effects: asymptotic scaling only when $n$ is large enough*

What about real data?

$\longrightarrow$ *second order approximation* with a Gaussian process $K_T$:

does it capture some aspects?

What about real data?

$\longrightarrow$ *second order approximation* with a Gaussian process $K_T$:

does it capture some aspects?

- Gaussian processes are $s$-times (mean-square) differentiable,
$$s = \frac{\alpha_T - d}{2}$$

What about real data?

$\longrightarrow$ *second order approximation* with a Gaussian process $K_T$:

does it capture some aspects?

- Gaussian processes are $s$-times (mean-square) differentiable,
$$s = \frac{\alpha_T - d}{2}$$

- Fitted exponents are $\beta \approx 0.4$ (MNIST) and $\beta \approx 0.1$ (CIFAR10), regardless of the Student $\longrightarrow \beta = \frac{\alpha_T - d}{d}$

  (since $\beta = \frac{1}{d}\min(\alpha_T - d, 2\alpha_S)$ indep. of $\alpha_S \longrightarrow \beta = \frac{\alpha_T - d}{d}$)

# BACK TOREAL DATA

<p align="center">What about real data?</p>

<p align="center">$\longrightarrow$ <em>second order approximation</em> with a Gaussian process $K_T$:</p>

<p align="center">does it capture some aspects?</p>

- Gaussian processes are $s$-times (mean-square) differentiable,
$$s = \frac{\alpha_T - d}{2}$$

- Fitted exponents are $\beta \approx 0.4$ (MNIST) and $\beta \approx 0.1$ (CIFAR10), regardless of the Student $\longrightarrow \beta = \frac{\alpha_T - d}{d}$

  (since $\beta = \frac{1}{d}\min(\alpha_T - d, 2\alpha_S)$ indep. of $\alpha_S \longrightarrow \beta = \frac{\alpha_T - d}{d}$)

$\longrightarrow s = \frac{1}{2}\beta d$, $s \approx 0.2d \approx 156$ (MNIST) and $s \approx 0.05d \approx 153$ (CIFAR10)

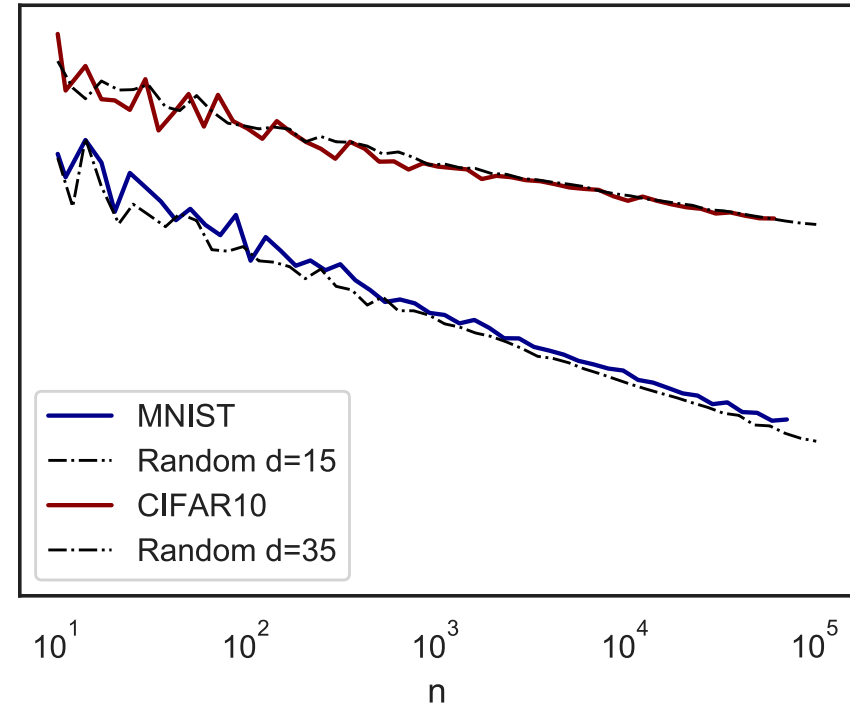<p align="center"><strong>This number is unreasonably large!</strong></p>

# EFFECTIVE DIMENSION

- Measure NN-distance $\delta$
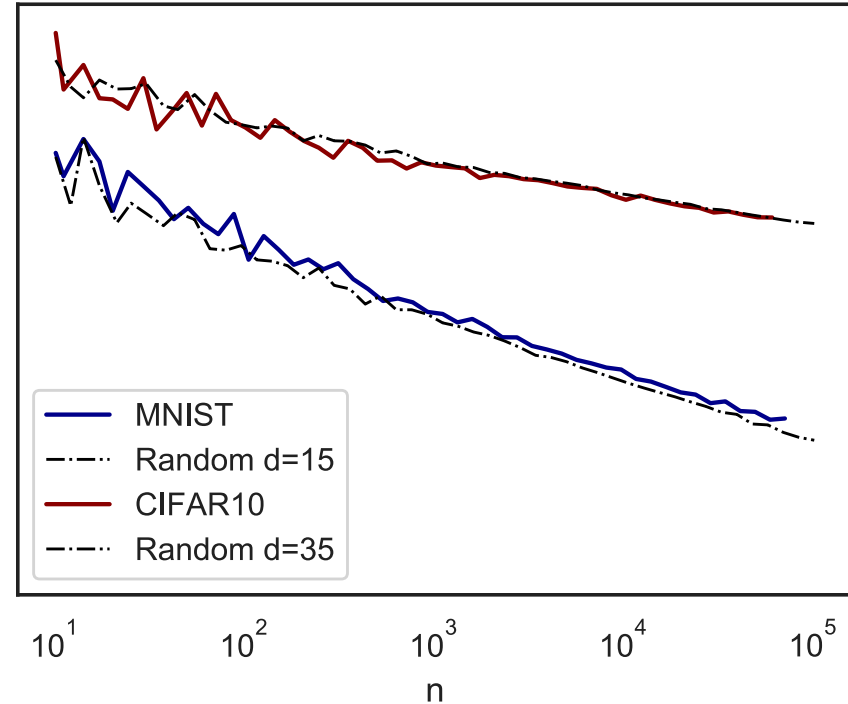
- $\delta \sim n^{-\text{some exponent}}$

# EFFECTIVE DIMENSION

- Measure NN-distance $\delta$

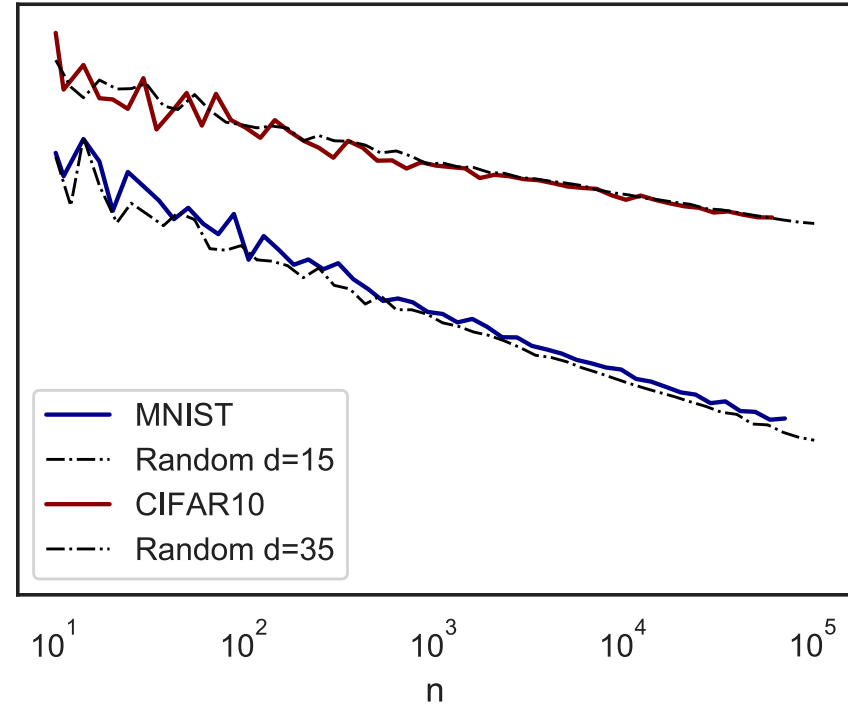- $\delta \sim n^{-\text{some exponent}}$

$\downarrow$



Define effective dimension as $\delta \sim n^{-\frac{1}{d_{\text{eff}}}}$

# EFFECTIVE DIMENSION

- Measure NN-distance $\delta$

- $\delta \sim n^{-\text{some exponent}}$

$\downarrow$



Define effective dimension as $\delta \sim n^{-\frac{1}{d_{\text{eff}}}}$

$d_{\text{eff}}$ is much smaller

| | $\beta$ | $d$ | $d_{\text{eff}}$ | $s = \left\lfloor \frac{1}{2}\beta d_{\text{eff}} \right\rfloor$ |
|---|---|---|---|---|
| MNIST | 0.4 | 784 | 15 | 3 |
| CIFAR10 | 0.1 | 3072 | 35 | 1 |

$s$ is more reasonable!

# CURSE OF DIMENSIONALITY (1/2)

- Loosely speaking, the (optimal) exponent is

$$\beta \approx \frac{\text{smoothness } \alpha_T - d = 2s}{\text{manifold dimension } d}$$

- To avoid the curse of dimensionality ($\beta \sim \frac{1}{d}$):

  - either the dimension of the manifold is small

  - or the data are extremely smooth

# RKHS & SMOOTHNESS

- Indeed, what happens if we consider a field $Z_T(\underline{x})$ that

    - is an instance of a Teacher $K_T$          $(\alpha_T)$
    - lies in the RKHS of a Student $K_S$     $(\alpha_S)$

# RKHS & SMOOTHNESS

- Indeed, what happens if we consider a field $Z_T(\underline{x})$ that

  - is an instance of a Teacher $K_T$ $\qquad\qquad$ $(\alpha_T)$
  - lies in the RKHS of a Student $K_S$ $\qquad\qquad$ $(\alpha_S)$

$$\mathbb{E}_T \|Z_T\|_{K_S} =$$
$$\mathbb{E}_T \int \mathrm{d}^d\underline{x}\,\mathrm{d}^d\underline{y}\ Z_T(\underline{x})K_S^{-1}(\underline{x},\underline{y})Z_T(\underline{y}) =$$
$$\int \mathrm{d}^d\underline{x}\,\mathrm{d}^d\underline{y}\ K_T(\underline{x},\underline{y})K_S^{-1}(\underline{x},\underline{y}) < \infty$$

$$\implies \qquad \alpha_T > \alpha_S + d$$

# RKHS & SMOOTHNESS

- Indeed, what happens if we consider a field $Z_T(\underline{x})$ that

  - is an instance of a Teacher $K_T$ $\quad\quad\quad$ ($\alpha_T$)
  - lies in the RKHS of a Student $K_S$ $\quad\quad$ ($\alpha_S$)

$$\mathbb{E}_T \|Z_T\|_{K_S} =$$
$$\mathbb{E}_T \int \mathrm{d}^d\underline{x}\,\mathrm{d}^d\underline{y}\, Z_T(\underline{x}) K_S^{-1}(\underline{x},\underline{y}) Z_T(\underline{y}) = \quad\quad \implies \quad\quad \alpha_T > \alpha_S + d$$
$$\int \mathrm{d}^d\underline{x}\,\mathrm{d}^d\underline{y}\, K_T(\underline{x},\underline{y}) K_S^{-1}(\underline{x},\underline{y}) < \infty$$

$$K_S(\underline{0}) \propto \int \mathrm{d}\underline{w}\, \tilde{K}_S(\underline{w}) < \infty \quad\quad\quad \implies \quad\quad\quad \alpha_S > d$$

# RKHS & SMOOTHNESS

- Indeed, what happens if we consider a field $Z_T(\underline{x})$ that

  - is an instance of a Teacher $K_T$       $(\alpha_T)$
  - lies in the RKHS of a Student $K_S$       $(\alpha_S)$

$$\mathbb{E}_T \| Z_T \|_{K_S} =$$
$$\mathbb{E}_T \int d^d\underline{x} d^d\underline{y} \, Z_T(\underline{x}) K_S^{-1}(\underline{x},\underline{y}) Z_T(\underline{y}) =$$
$$\int d^d\underline{x} d^d\underline{y} \, K_T(\underline{x},\underline{y}) K_S^{-1}(\underline{x},\underline{y}) < \infty$$
$$\implies \quad \alpha_T > \alpha_S + d$$

$$K_S(\underline{0}) \propto \int d\underline{w} \, \tilde{K}_S(\underline{w}) < \infty \qquad \implies \qquad \alpha_S > d$$

(it scales with $d$!)

Therefore the smoothness must be $s = \frac{\alpha_T - d}{2} > \frac{d}{2}$

$$\longrightarrow \beta > \tfrac{1}{2}$$

# CURSE OF DIMENSIONALITY (2/2)

- Assume that the data are not smooth enough and live in $d$ large

- **Dimensionality reduction** in the task rather than in the data?

- E.g. the $n$ points $\underline{x}_\mu$ live in $\mathbb{R}^d$, but the target function is such that

$$Z_T(\underline{x}) = Z_T(\underline{x}_\parallel) \equiv Z_T(x_1, \ldots, x_{d_\parallel}), \quad d_\parallel < d$$

Similar setting studied in Bach 2017

- Can kernels understand the lower dimensional structure?

# TASK INVARIANCE: KERNEL REGRESSION (1/2)

Theorem (informal formulation):

in the described setting with $d_\parallel \leq d$,

for $n \gg 1$
$\epsilon \sim n^{-\beta}$   with   $\boxed{\beta = \frac{1}{d} \min(\alpha_T - d, 2\alpha_S)}$

Regardless of $d_\parallel$!

Similar result in Bach 2017

Two reasons contribute to this result:

- the nearest-neighbor distance always scales as $\delta \sim n^{-\frac{1}{d}}$

- $\alpha_T(d) - d$ only depends on the function $K_T(z)$ and not on $d$

Teacher = Matérn (with parameter $\nu$),    Student = Laplace,    $d$=4

**Classification** with the **margin SVM** algorithm:

$$\hat{y}(\underline{x}) = \text{sign}\left[\sum_{\mu=1}^{n} c_\mu K\left(\frac{\|x - x^\mu\|}{\sigma}\right) + b\right]$$

find $\{c_\mu\}, b$ by minimizing some function

We consider a very simple setting:

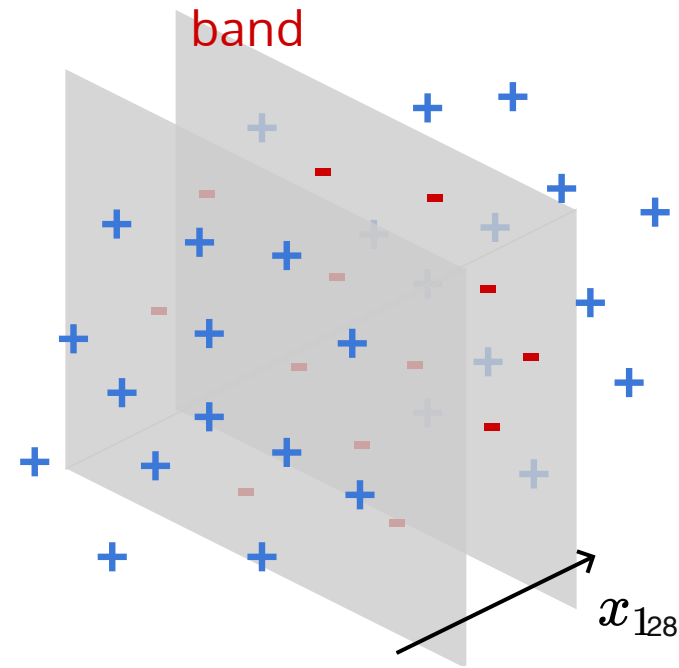- the label is $y(\underline{x}) = y(x_1) \longrightarrow d_\parallel = 1$

**Non-Gaussian data!**

$y(x_1):$

hyperplane

band



$x_1$

$x_{128}$

# TASK INVARIANCE: CLASSIFICATION (2/2)

Vary **kernel scale** $\sigma \longrightarrow$ **two regimes!**

- $\sigma \ll \delta$: then the estimator is tantamount to a **nearest-neighbor algorithm** $\longrightarrow$ curse of dimensionality $\beta = \frac{1}{d}$

- $\sigma \gg \delta$: important **correlations** in $c_\mu$ due to the **long-range kernel**. For the hyperplane with $d_\parallel = 1$ we find $\beta = \mathcal{O}(d^0)$!
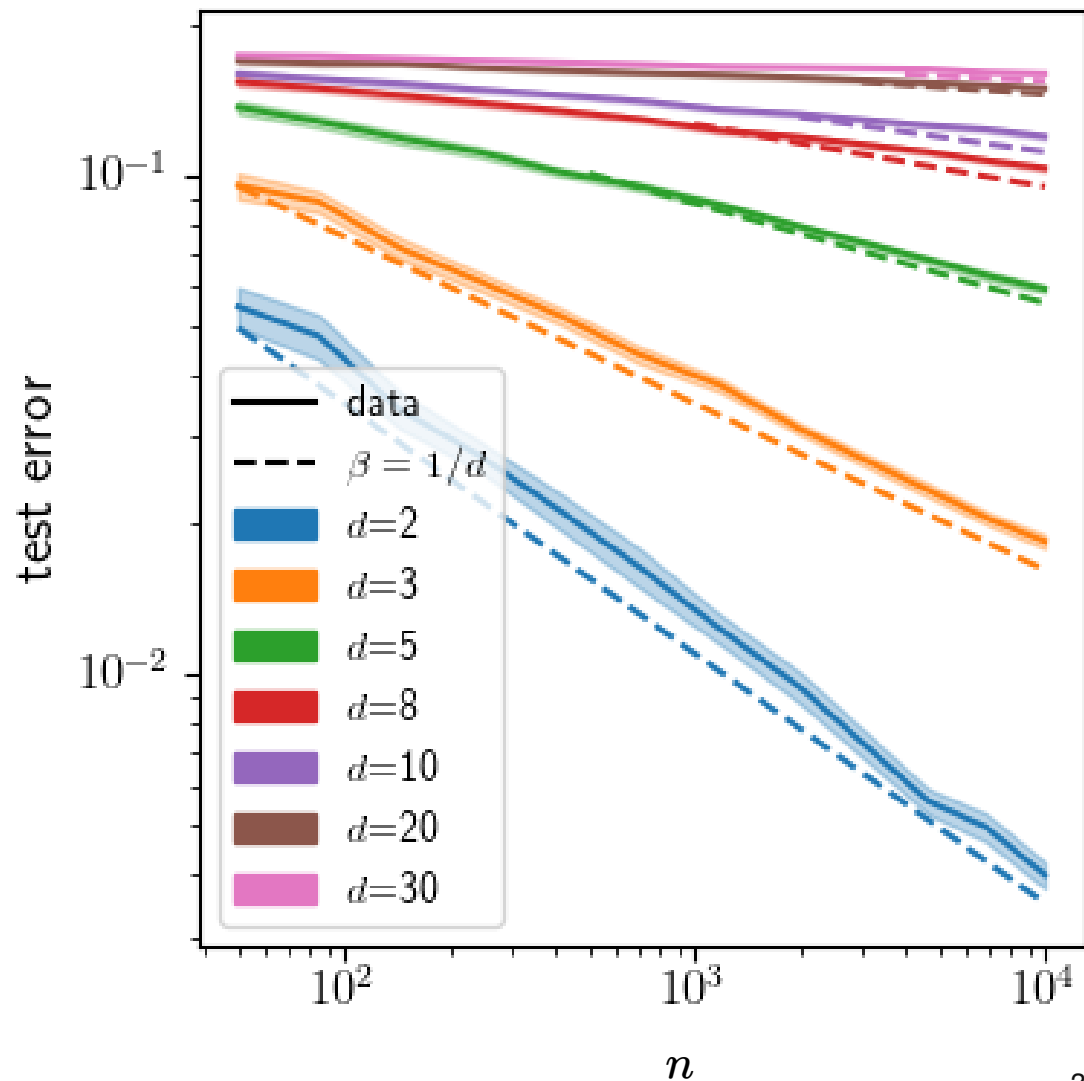
**No curse of dimensionality!**

# THE NEAREST-NEIGHBOR LIMIT

hyperplane
interface

using a Laplace kernel
and
varying the dimension $d$:

$$\beta = \frac{1}{d}$$

# KERNEL CORRELATIONS (1/2)

When $\sigma \gg \delta$ we can expand the kernel overlaps:

$$K\left(\tfrac{\|x-x^{\mu}\|}{\sigma}\right) \approx K(0) - \mathrm{const} \times \left(\tfrac{\|x-x^{\mu}\|}{\sigma}\right)^{\xi}$$

(the exponent $\xi$ is linked to the smoothness of the kernel)

We can derive some scaling arguments that lead to an exponent
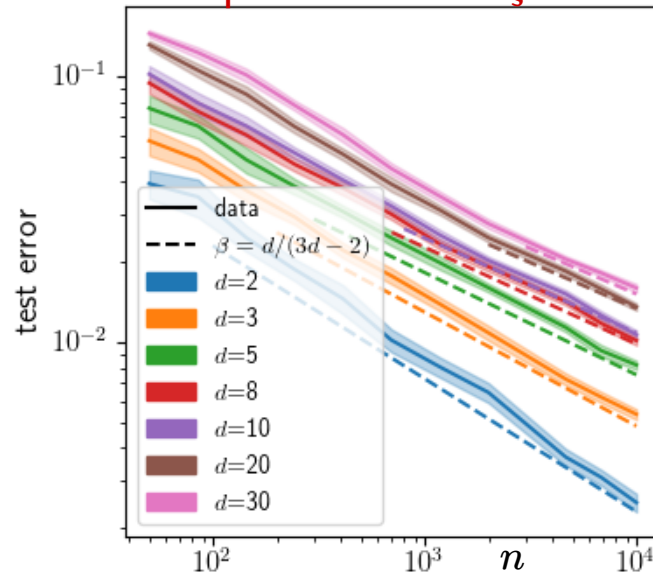
$$\beta = \tfrac{d+\xi-1}{3d+\xi-3}$$

# KERNEL CORRELATIONS (1/2)

When $\sigma \gg \delta$ we can expand the kernel overlaps:

$$K\left(\frac{\|x-x^\mu\|}{\sigma}\right) \approx K(0) - \text{const} \times \left(\frac{\|x-x^\mu\|}{\sigma}\right)^\xi$$

(the exponent $\xi$ is linked to the smoothness of the kernel)

We can derive some scaling arguments that lead to an exponent

$$\beta = \frac{d+\xi-1}{3d+\xi-3}$$

Idea:

- support vectors ($c_\mu \neq 0$) are close to the interface
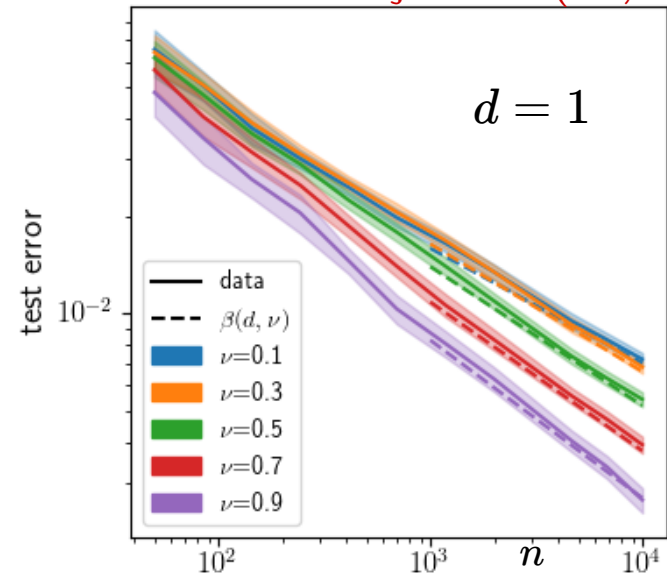- we impose that the decision boundary has $\mathcal{O}(1)$ spatial fluctuations on a scale proportional to $\delta$
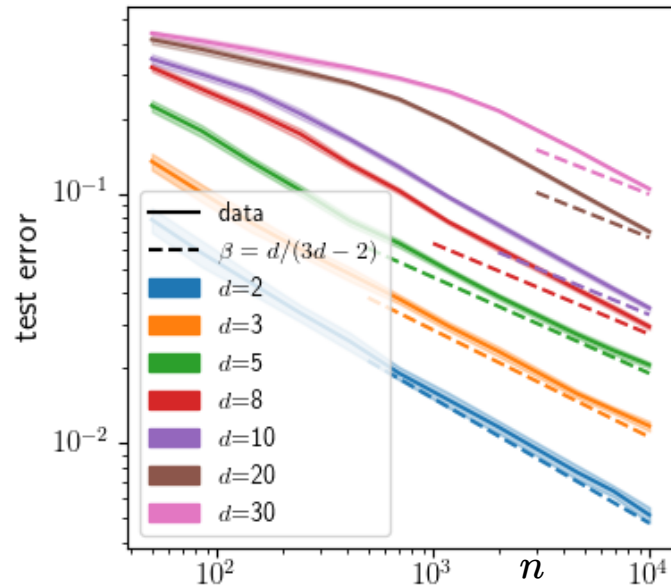
# KERNEL CORRELATIONS (2/2)



Laplace kernel $\xi = 1$

Matérn kernels $\xi = \min(2\nu, 2)$

hyperplane

$d = 1$

band

$$\beta = \frac{d+\xi-1}{3d+\xi-3}$$

in all these cases!

# KERNEL CORRELATIONS: HYPERSPHERE
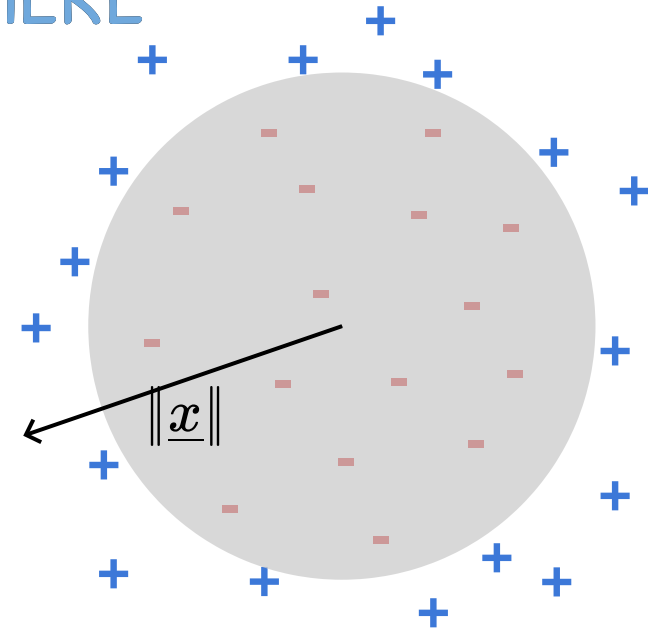
What about other interfaces?

boundary = hypersphere:

$$y(\underline{x}) = \text{sign}(\|\underline{x}\| - R)$$
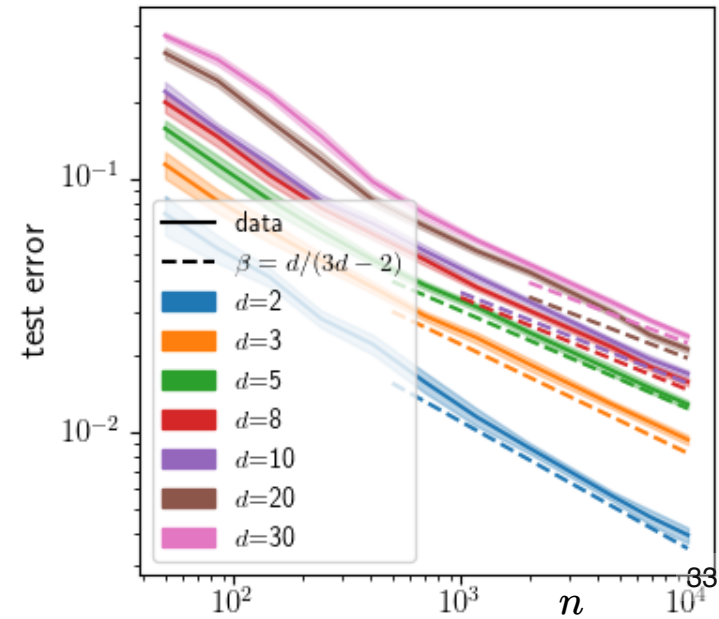
$$(d_\parallel = 1)$$

$$\beta = \frac{d+\xi-1}{3d+\xi-3}$$

(same exponent!)

(similar scaling arguments apply, provided $R \gg \delta$)



Laplace kernels ($\xi = 1$)

# CONCLUSION

- Learning curves of real data decay as **power laws** with exponents

$$\frac{1}{d} \ll \beta < \frac{1}{2}$$

- We introduce a **new framework** that links the exponent $\beta$ to the degree of smoothness of Gaussian random data

- We justify how different kernels can lead to the same exponent $\beta$

- We show that the **effective dimension** of real data is $\ll d$. It can be linked to a (small) **effective smoothness** $s$

- We show that kernel regression is not able to capture invariants in the task, while kernel classification can

(in some regime and for **smooth interfaces**)