

Mathematics of Data: From Theory to Computation

Prof. Volkan Cevher
volkan.cevher@epfl.ch

Lecture 1: Introduction to Continuous Optimization

Laboratory for Information and Inference Systems (LIONS)
École Polytechnique Fédérale de Lausanne (EPFL)

EE-556 (Fall 2019)



License Information for Mathematics of Data Slides

- ▶ This work is released under a [Creative Commons License](#) with the following terms:
- ▶ **Attribution**
 - ▶ The licensor permits others to copy, distribute, display, and perform the work. In return, licensees must give the original authors credit.
- ▶ **Non-Commercial**
 - ▶ The licensor permits others to copy, distribute, display, and perform the work. In return, licensees may not use the work for commercial purposes – unless they get the licensor's permission.
- ▶ **Share Alike**
 - ▶ The licensor permits others to distribute derivative works only under a license identical to the one that governs the licensor's work.
- ▶ [Full Text of the License](#)

Logistics

- ▶ **Credits:** 4
- ▶ **Prerequisites:** Previous coursework in calculus, linear algebra, and probability is required. Familiarity with optimization is useful.
- ▶ **Grading:** Continuous control via homework exercises & exam (cf., syllabus)
- ▶ **HW topics:** Support vector machines, compressive subsampling, neural networks, power flow...
- ▶ **Moodle:** My courses > Genie électrique et électronique (EL) > Master > EE-556
syllabus & course outline & HW exercises
- ▶ **TA's:** Thomas Sanchez (head TA), Paul Rolland, Maria Vladarean, Chaehwan Song, Ali Kavis, Mehmet Fatih Sahin, Fabian Latorre, and Ahmet Alacaoglu.

Outline

- ▶ This class:
 1. What is an optimization problem?
 2. Gradient descent: A basic introduction
 3. Common templates on convex/non-convex optimization
- ▶ Next class
 1. Review of probability, statistics and linear algebra

Recommended reading material

- ▶ Chapter 1 in S. Boyd, and L. Vandenberghe, *Convex Optimization*, Cambridge Univ. Press, 2009.
- ▶ Chapter 1 in Nocedal, Jorge, and Wright, Stephen J., *Numerical Optimization*, Springer, 2006.

What is optimization?

Problem (Mathematical formulation)

The optimization problem

$$\min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}) \quad (1)$$

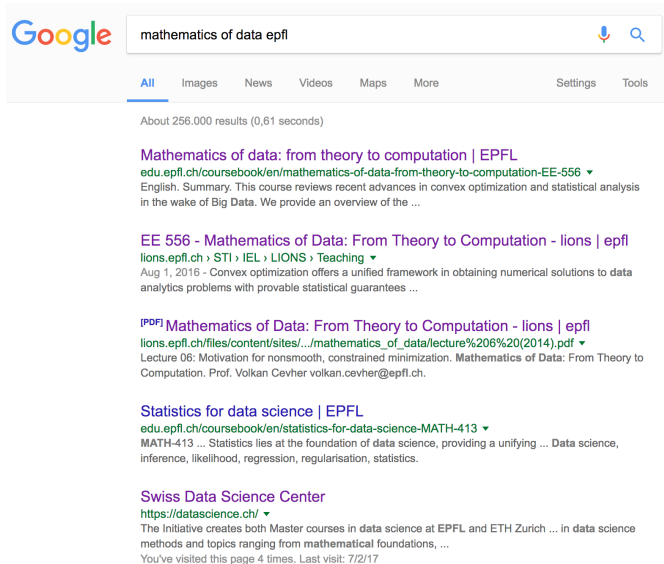
consists in finding a minimizer $\mathbf{x}^ \in S_f^*$. We say (1) has a solution if S_f^* is non-empty*

Definition (Set of minimizers)

Let $f : \mathcal{X} \rightarrow \mathbb{R}$ and $\delta > 0$

- ▶ $\mathbf{x}^* \in \mathcal{X}$ is a global minimizer of f if $f(\mathbf{x}^*) \leq f(\mathbf{x})$ for all $\mathbf{x} \in \mathcal{X}$
- ▶ $\mathbf{x}_{\text{loc}}^* \in \mathcal{X}$ is a local minimizer of f if $f(\mathbf{x}_{\text{loc}}^*) \leq f(\mathbf{x})$ for all $\mathbf{x} : \|\mathbf{x} - \mathbf{x}_{\text{loc}}^*\| \leq \delta$
- ▶ $f^* := f(\mathbf{x}^*)$ is the global minimum value of f
- ▶ S_f^* is the set of minimizers of f

Example: Google PageRank



The screenshot shows a Google search interface. The search bar contains the text "mathematics of data epfl". Below the search bar, there are tabs for "All", "Images", "News", "Videos", "Maps", and "More". The "All" tab is selected. To the right of the tabs are links for "Settings" and "Tools". Below the tabs, it says "About 256.000 results (0,61 seconds)". The search results are listed below, each with a title, a link, and a brief description.

Google

mathematics of data epfl

All Images News Videos Maps More Settings Tools

About 256.000 results (0,61 seconds)

Mathematics of data: from theory to computation | EPFL
edu.epfl.ch/coursebook/en/mathematics-of-data-from-theory-to-computation-EE-556 ▼
English. Summary. This course reviews recent advances in convex optimization and statistical analysis in the wake of Big Data. We provide an overview of the ...

EE 556 - Mathematics of Data: From Theory to Computation - lions | epfl
lions.epfl.ch › [STI](#) › [IEL](#) › [LIONS](#) › [Teaching](#) ▼
Aug 1, 2016 - Convex optimization offers a unified framework in obtaining numerical solutions to data analytics problems with provable statistical guarantees ...

[PDF] Mathematics of Data: From Theory to Computation - lions | epfl
[lions.epfl.ch/files/content/sites/.../mathematics_of_data/lecture%206%20\(2014\).pdf](http://lions.epfl.ch/files/content/sites/.../mathematics_of_data/lecture%206%20(2014).pdf) ▼
Lecture 06: Motivation for nonsmooth, constrained minimization. Mathematics of Data: From Theory to Computation. Prof. Volkan Cevher volkan.cevher@epfl.ch.

Statistics for data science | EPFL
edu.epfl.ch/coursebook/en/statistics-for-data-science-MATH-413 ▼
MATH-413 ... Statistics lies at the foundation of data science, providing a unifying ... Data science, inference, likelihood, regression, regularisation, statistics.

Swiss Data Science Center
<https://datascience.ch/> ▼
The Initiative creates both Master courses in data science at EPFL and ETH Zurich ... In data science methods and topics ranging from mathematical foundations, ...
You've visited this page 4 times. Last visit: 7/2/17

Modeling Google PageRank

- A basic model

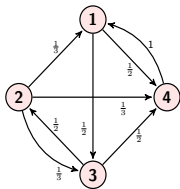


- Compute the conditional probabilities:

$$P(\text{The Washington Post} | \text{Google News}) = 2/8$$

$$P(\text{The Atlantic} | \text{Google News}) = 1/8$$

- A toy graph and transition matrix:

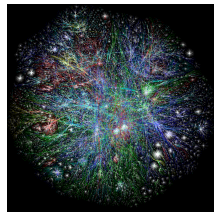


$$E = \begin{bmatrix} 0 & \frac{1}{3} & 0 & 1 \\ 0 & 0 & \frac{1}{2} & 0 \\ \frac{1}{2} & \frac{1}{3} & 0 & 0 \\ \frac{1}{2} & \frac{1}{3} & \frac{1}{2} & 0 \end{bmatrix}$$

Modeling Google PageRank

- Transition matrix for world wide web:

$$\mathbf{E} = \begin{bmatrix} c_{11} & c_{12} & \dots & c_{1n} \\ c_{21} & c_{22} & \dots & c_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ c_{n1} & c_{n2} & \dots & c_{nn} \end{bmatrix}$$

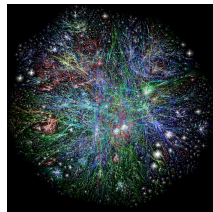


- $\sum_{i=1}^n c_{ij} = 1, \forall j \in \{1, 2, \dots, n\}$ ($n \approx 5.91\text{billion}$)
- Estimated memory to store \mathbf{E} : 10^{12} GB!

Modeling Google PageRank

- Transition matrix for world wide web:

$$\mathbf{E} = \begin{bmatrix} c_{11} & c_{12} & \dots & c_{1n} \\ c_{21} & c_{22} & \dots & c_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ c_{n1} & c_{n2} & \dots & c_{nn} \end{bmatrix}$$



- $\sum_{i=1}^n c_{ij} = 1, \forall j \in \{1, 2, \dots, n\}$ ($n \approx 5.91$ billion)
- Estimated memory to store \mathbf{E} : 10^{12} GB!
- A bit of mathematical modeling:
 - ▶ r_i^k : Probability of being at node i at k^{th} state. Let us define a state vector

$$\mathbf{r}^k = [r_1^k, r_2^k, \dots, r_n^k]^\top$$

- ▶ Multiplying \mathbf{r}^k by \mathbf{E} takes one random step along the edges of the graph:

$$r_i^1 = \sum_{j=1}^n c_{ij} r_j^0 = (\mathbf{E} \mathbf{r}^0)_i,$$

since $c_{ij} = P(i|j)$ (by the law of total probability).

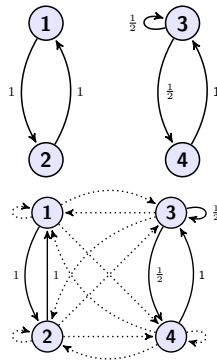
Towards a Formal Formulation for Google PageRank

Goal

Find the ranking vector \mathbf{r}^* after an infinite number of random steps.

- Disconnected web: Initial state vector affects the ranking vector.

A solution: Model the event that the surfer will quit the current webpage and open another.



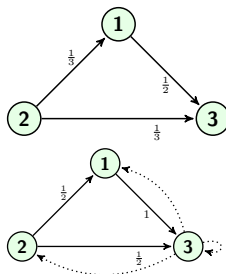
Towards a Formal Formulation for Google PageRank

Goal

Find the ranking vector \mathbf{r}^* after an infinite number of random steps.

- Sink nodes: Column of zeros in \mathbf{E} , moves \mathbf{r} to 0!

A solution: Create artificial links from sink nodes to all the nodes.



Towards a Formal Formulation for Google PageRank

Goal

Find the ranking vector \mathbf{r}^* after an infinite number of random steps.

- Disconnected web: Initial state vector affects the ranking vector.

A solution: Model the event that the surfer quits the current webpage to open another.

$$\mathbf{B} = \begin{bmatrix} 1 & 1 & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & 1 \end{bmatrix} = \frac{1}{n} \mathbf{1}\mathbf{1}^\top$$

- Sink nodes: Column of zeros in \mathbf{E} , moves \mathbf{r} to $\mathbf{0}$!

A solution: Create artificial links from sink nodes to all the nodes.

$$\lambda_i = \begin{cases} 1 & \text{if } i^{th} \text{ node is a sink node,} \\ 0 & \text{otherwise.} \end{cases}$$

Optimization formulation of Google PageRank

- Define the pagerank matrix \mathbf{M} as

$$\mathbf{M} = (1 - p)(\mathbf{E} + \frac{1}{n}\mathbf{1}\mathbf{1}^T) + p\mathbf{B}.$$

\mathbf{M} is a column stochastic matrix.

Problem Formulation

- We characterize the solution as
 - $\mathbf{M}\mathbf{r}^* = \mathbf{r}^*$.
 - \mathbf{r}^* is a probability state vector:

$$r_i \geq 0, \quad \sum_{i=1}^n r_i = 1.$$

- Find $\mathbf{r} \geq 0$ such that $\mathbf{M}\mathbf{r} = \mathbf{r}$ and $\mathbf{1}^\top \mathbf{r} = 1$.

Optimization formulation of Google PageRank

- Define the pagerank matrix \mathbf{M} as

$$\mathbf{M} = (1 - p)(\mathbf{E} + \frac{1}{n}\mathbf{1}\mathbf{1}^T) + p\mathbf{B}.$$

\mathbf{M} is a column stochastic matrix.

Problem Formulation

- We characterize the solution as
 - $\mathbf{M}\mathbf{r}^* = \mathbf{r}^*$.
 - \mathbf{r}^* is a probability state vector:

$$r_i \geq 0, \quad \sum_{i=1}^n r_i = 1.$$

- Find $\mathbf{r} \geq 0$ such that $\mathbf{M}\mathbf{r} = \mathbf{r}$ and $\mathbf{1}^\top \mathbf{r} = 1$.

Optimization formulation

$$\min_{\mathbf{x} \in \mathbb{R}^n} \left\{ f(\mathbf{x}) = \frac{1}{2} \|\mathbf{M}\mathbf{x} - \mathbf{x}\|^2 + \frac{\gamma}{2} (\mathbf{1}^\top \mathbf{x} - 1)^2 \right\}.$$

The general formulation: Least-squares

Optimization formulation (Least-squares estimator)

$$\min_{\mathbf{x} \in \mathbb{R}^d} \underbrace{\frac{1}{2} \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2^2}_{f(\mathbf{x})},$$

where $\mathbf{x} = \mathbf{r}$, $\mathbf{b} = \begin{bmatrix} \mathbf{r} \\ \frac{\gamma}{n} \mathbf{1} \end{bmatrix}$, $\mathbf{A} = \begin{bmatrix} \mathbf{M} \\ \frac{\gamma}{2n} \mathbf{1}\mathbf{1}^\top \end{bmatrix}$, $d = n$ in Google PageRank problem.

Linear regression problem

Let $\mathbf{x}^\natural \in \mathbb{R}^d$ and $\mathbf{A} \in \mathbb{R}^{n \times d}$ (full column rank). **Goal:** estimate \mathbf{x}^\natural , given \mathbf{A} and

$$\mathbf{b} = \mathbf{A}\mathbf{x}^\natural + \mathbf{w},$$



where \mathbf{w} denotes unknown noise.

- Many other examples:

Image reconstruction (MRI), stock market prediction, house pricing, etc.

Regression/Classification/Generation Examples

- Example: Taking a mortgage.
- Houses data (source: <https://www.homegate.ch>)

	Type	Apartment	Ecublens
	Rooms	5.5	1024 Ecublens VO
	Living space	200 m²	
	Year built	1991	
1,325,000.-			
	Type	Villa	1024 Ecublens VO
	Rooms	7.5	
	Living space	250 m²	
	Lot size	584 m²	
	Year built	1965	
1,390,000.-			

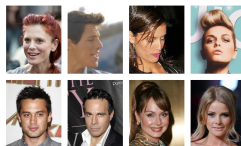
- Example: Image classification



- Imagenet: 1000 object classes.
1.2M/100K train/test images

source: <https://www.imagenet.org>

- Example: Image generation (GANs).
- Target distribution:



source: <http://mmlab.ie.cuhk.edu.hk/projects/CelebA.html>

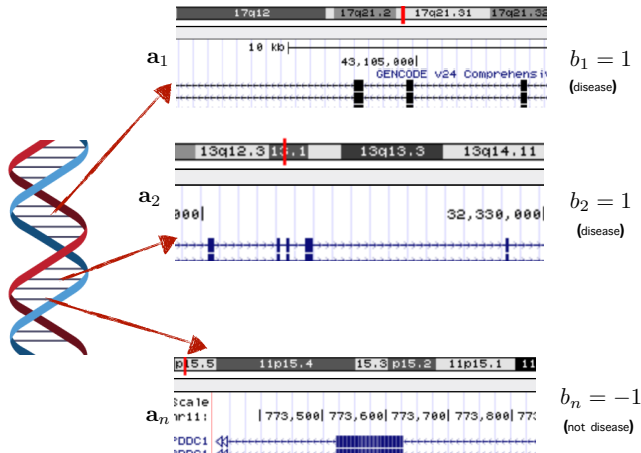
- How to generate (fake) images from the same distribution?



source: Progressive Growing of GANs for Improved Quality, Stability, and Variation Karras et al. ICLR 2018

Breast Cancer Detection

- Genome data for breast cancer (source: <http://genome.ucsc.edu>):



- A patient with genome data a_t : has he got breast cancer or not (i.e., $b_t = 1$ or -1)?

Score-based Classifiers (I)

Goal

Predict either $b = 1$ or $b = -1$ given \mathbf{a} .

- For a genome sequence \mathbf{a} compute a *score* $s_{\mathbf{x}}(\mathbf{a}) \in (-\infty, \infty)$:

Example: $\mathbf{a} \rightarrow s_{\mathbf{x}}(\mathbf{a}) = \underbrace{\mathbf{x}^{\top}}_{\text{weights = importance of genes}} \mathbf{a}$

- Use the **logistic function**

$$t \mapsto h(t) := \frac{1}{1 + \exp(-t)}.$$

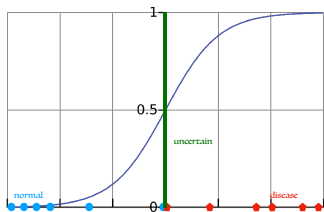
to transform $s_{\mathbf{x}}(\mathbf{a})$ into a probability of disease:

$$P(b = 1 | \mathbf{a}, \mathbf{x}) = h(s_{\mathbf{x}}(\mathbf{a})) \in (0, 1)$$

Score-based Classifiers (II)

- We have a model for the conditional probability of disease given \mathbf{a}

$$P(b = 1|\mathbf{a}, \mathbf{x}) = \frac{1}{1 + \exp(-s_{\mathbf{x}}(\mathbf{a}))}$$

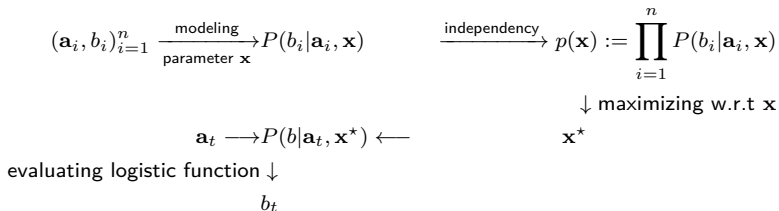


$$P(b = 1|\mathbf{a}, \mathbf{x}) \begin{cases} > 0.5, & \text{if } s_{\mathbf{x}}(\mathbf{a}) \text{ is positive,} \\ \leq 0.5, & \text{otherwise.} \end{cases}$$

$$\text{Prediction} = \begin{cases} \text{disease,} & \text{if } P(b = 1|\mathbf{a}, \mathbf{x}) > 0.5, \\ \text{normal,} & \text{if } P(b = 1|\mathbf{a}, \mathbf{x}) < 0.5. \\ \text{uncertain,} & \text{if } P(b = 1|\mathbf{a}, \mathbf{x}) = 0.5. \end{cases}$$

Score-based Classifiers: How do we choose \mathbf{x} ?

- Classification diagram:



- Maximizing $\log p(\mathbf{x})$ gives the **log-likelihood estimator**.

Logistic regression vs Neural networks

Optimization formulation (max Log-likelihood)

$$\min_{\mathbf{x} \in \mathbb{R}^p} \log p(\mathbf{x}) = \underbrace{\sum_{i=1}^n \log(1 + \exp(-b_i(s_{\mathbf{x}}(\mathbf{a}_i))))}_{f(\mathbf{x})} \quad (2)$$

Problem (Choice of score function)

- *Logistic Regression:*

$$s_{\mathbf{x}}(\mathbf{a}) = \mathbf{x}^{\top} \mathbf{a} \quad \Rightarrow (2) \text{ is convex} \quad (3)$$

- *Neural Networks* ($\mathbf{x} = [\mathbf{v}, \mathbf{W}]$):

$$s_{\mathbf{x}}(\mathbf{a}) = \mathbf{v}^{\top} \sigma(\mathbf{W}\mathbf{a}) \quad \Rightarrow (2) \text{ is non-convex} \quad (4)$$

Optimization landscapes

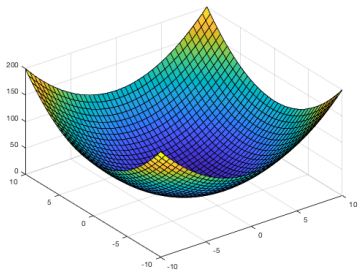


Figure: convex function

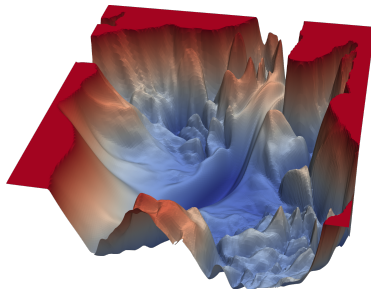


Figure: nonconvex function

A basic *iterative* strategy

General idea of an optimization algorithm

Guess a solution, and then *refine* it based on *oracle information*.

Repeat the procedure until the result is *good enough*.

Approximate vs. exact optimality and local optimality (I)

Is it possible to solve an optimization problem?

*"In general, optimization problems are **unsolvable**" - Y. Nesterov [1]*

- ▶ Even when a closed-form solution exists, numerical accuracy may still be an issue.
- ▶ We must be content with **approximately** optimal solutions.
- ▶ Sometimes we can only find **locally** optimal solutions.

Approximate vs. exact optimality and local optimality (II)

Definition

We say that \mathbf{x}_ϵ^* is ϵ -optimal in **objective value** if

$$f(\mathbf{x}_\epsilon^*) - f^* \leq \epsilon .$$

Definition

We say that \mathbf{x}_ϵ^* is ϵ -optimal in **sequence** if, for some norm $\|\cdot\|$,

$$\|\mathbf{x}_\epsilon^* - \mathbf{x}^*\| \leq \epsilon ,$$

- The latter approximation guarantee is considered stronger.

Definition

We say that $\mathbf{x}_{\text{loc}}^*$ is a local minimizer if for some $\delta > 0$,

$$f(\mathbf{x}_{\text{loc}}^*) \leq f(\mathbf{x}) \quad \text{for all } \mathbf{x} \text{ such that } \|\mathbf{x} - \mathbf{x}_{\text{loc}}^*\| \leq \delta$$

Necessary/Sufficient optimality conditions

Lemma (First-order necessary local optimality condition)

Let \mathbf{x}_{loc}^* be a local minimizer of a differentiable convex function f . It holds that

$$\nabla f(\mathbf{x}_{loc}^*) = \mathbf{0}.$$

Lemma (sufficient optimality condition for convex functions)

Let f be a **convex** function. It holds that

$$\mathbf{x}_{loc}^* \text{ is a local minimizer} \Rightarrow \mathbf{x}_{loc}^* \text{ is a global minimizer}$$

A gradient method

Gradient method

Choose a starting point \mathbf{x}^0 and iterate

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \alpha_k \nabla f(\mathbf{x}^k)$$

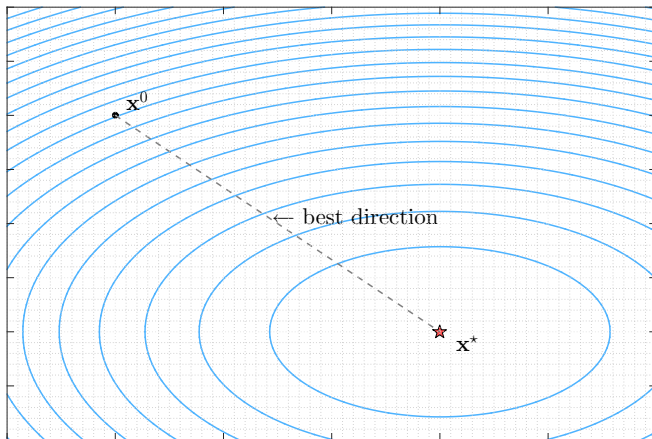
where α_k is a step-size to be chosen so that \mathbf{x}^k converges to some $\mathbf{x}_{\text{loc}}^*$.

Fixed-point characterization

Given that $\nabla f(\mathbf{x}_{\text{loc}}^*) = 0$, the following fixed point condition holds

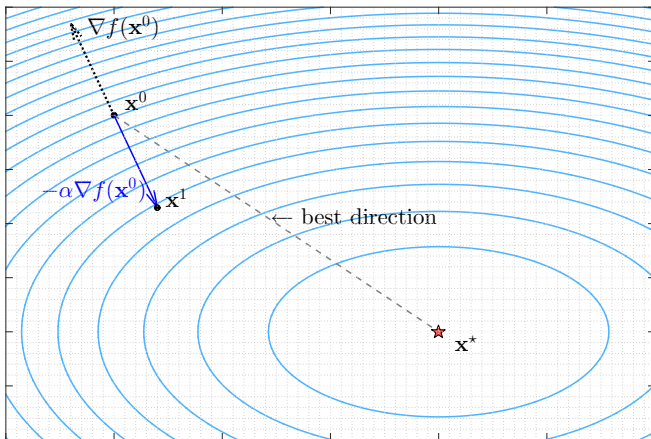
$$\mathbf{x}_{\text{loc}}^* = \mathbf{x}_{\text{loc}}^* - \alpha \nabla f(\mathbf{x}_{\text{loc}}^*) \quad \text{for all } \alpha \in \mathbb{R}$$

A simple example



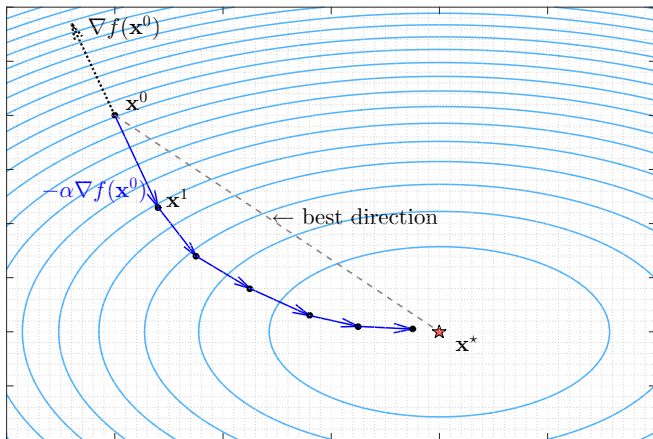
- Choose initial point: x^0 , and a step size $\alpha > 0$.

A simple example



- ▶ Choose initial point: x^0 , and a step size $\alpha > 0$.
- ▶ Take a step in the negative gradient direction: $x^{k+1} = x^k - \alpha \nabla f(x^k)$

A simple example



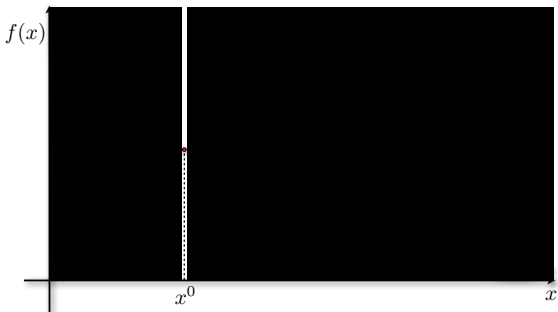
- ▶ Choose initial point: x^0 , and a step size $\alpha > 0$.
- ▶ Take a step in the negative gradient direction: $x^{k+1} = x^k - \alpha \nabla f(x^k)$
- ▶ Repeat this procedure until x^k is accurate enough.

Challenges for an iterative optimization algorithm

Problem

Find the minimizer x^* of $f(x)$, given starting point x^0 based on only local information.

► Fog of war

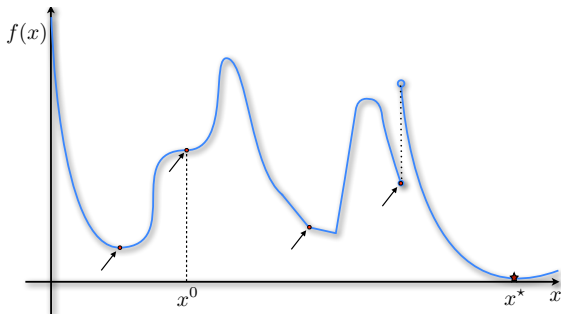


Challenges for an iterative optimization algorithm

Problem

Find the minimizer x^* of $f(x)$, given starting point x^0 based on only local information.

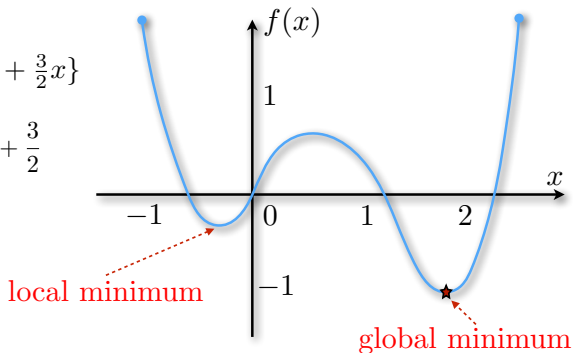
- Fog of war, non-differentiability, discontinuities, local minimizers, stationary points...



Local minimizers

$$\min_{x \in \mathbb{R}} \{x^4 - 3x^3 + x^2 + \frac{3}{2}x\}$$

$$\frac{df}{dx} = 4x^3 - 9x^2 + 2x + \frac{3}{2}$$



Choose $x^0 = 0$ and $\alpha = \frac{1}{6}$

$$x^1 = x^0 - \alpha \frac{df}{dx} \Big|_{x=x^0} = 0 - \frac{1}{6} \frac{3}{2} = -\frac{1}{4}$$

$$x^2 = -\frac{5}{16}$$

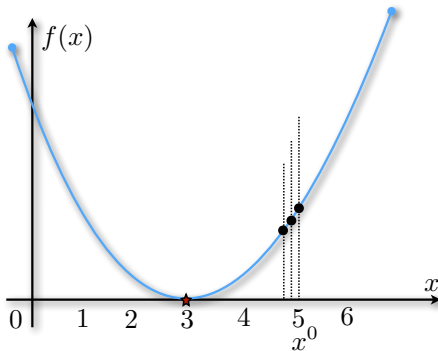
...

x^k is converging to **local minimizer**!

Effect of very small step-size α ...

$$\min_{x \in \mathbb{R}} \frac{1}{2}(x-3)^2$$

$$\frac{df}{dx} = x - 3$$



Choose $x^0 = 5$ and $\alpha = \frac{1}{10}$

$$x^1 = x^0 - \alpha \left. \frac{df}{dx} \right|_{x=x^0} = 5 - \frac{1}{10} 2 = 4.8$$

$$x^2 = x^1 - \alpha \left. \frac{df}{dx} \right|_{x=x^1} = 4.8 - \frac{1}{10} 1.8 = 4.62$$

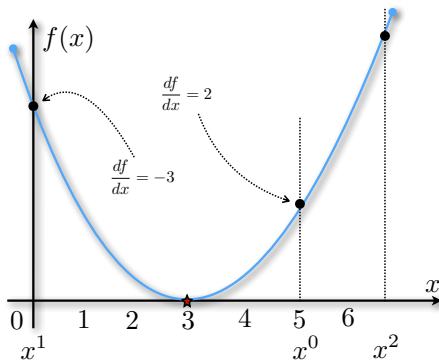
...

x^k converges **very slowly**.

Effect of very large step-size α ...

$$\min_{x \in \mathbb{R}} \frac{1}{2}(x-3)^2$$

$$\frac{df}{dx} = x - 3$$



Choose $x^0 = 5$ and $\alpha = \frac{5}{2}$

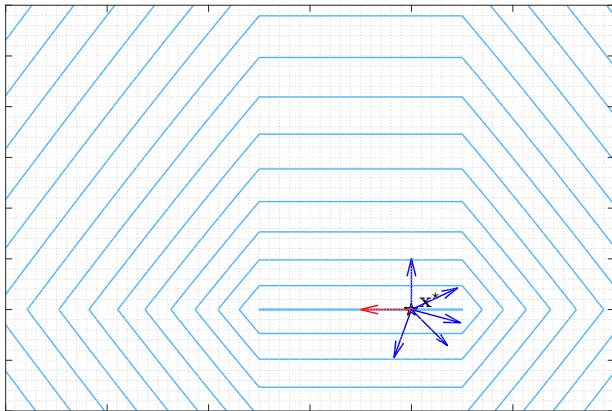
$$x^1 = x^0 - \alpha \frac{df}{dx} \Big|_{x=x^0} = 5 - \frac{5}{2} \cdot 2 = 0$$

$$x^2 = x^1 - \alpha \frac{df}{dx} \Big|_{x=x^1} = 0 - \frac{5}{2}(-3) = \frac{15}{2}$$

...

x^k diverges.

Nonsmooth optimization

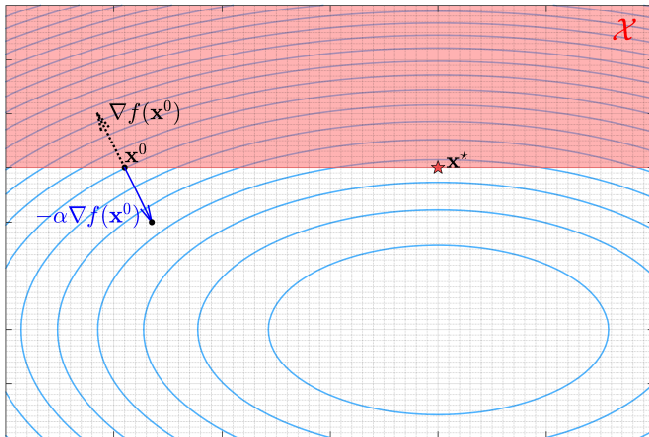


For nonsmooth optimization, the first order optimality condition

$$\nabla f(\mathbf{x}^*) = \mathbf{0}$$

does not hold for every descent direction.

Constrained optimization

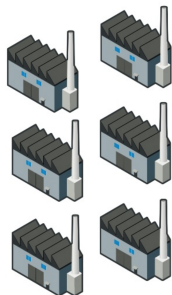


In many practical problems,
we need to **minimize** the cost **under some constraints**.

$$f^* := \min_{\mathbf{x} \in \mathbb{R}^p} \left\{ f(\mathbf{x}) : \mathbf{x} \in \mathcal{X} \right\}$$

Example: Facility Location Problem

Assign facilities to locations to minimize the total assignment cost.



Facilities



Locations

Example: Facility Location Problem

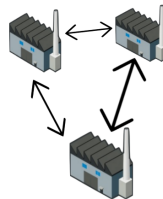
► **Goal:** To minimize the costs

► **Inputs:**

Distance between locations: $A = \begin{bmatrix} 0 & a_{12} & \dots & a_{1n} \\ a_{21} & \ddots & & \\ \vdots & & \ddots & \\ a_{n1} & & & 0 \end{bmatrix}$



Flow between facilities: $B = \begin{bmatrix} 0 & b_{12} & \dots & b_{1n} \\ b_{21} & \ddots & & \\ \vdots & & \ddots & \\ b_{n1} & & & 0 \end{bmatrix}$



Example: Facility Location Problem

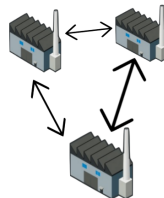
► **Goal:** To minimize the costs

► **Inputs:**

Distance between locations: $A = \begin{bmatrix} 0 & a_{12} & \dots & a_{1n} \\ a_{21} & \ddots & & \\ \vdots & & \ddots & \\ a_{n1} & & & 0 \end{bmatrix}$



Flow between facilities: $B = \begin{bmatrix} 0 & b_{12} & \dots & b_{1n} \\ b_{21} & \ddots & & \\ \vdots & & \ddots & \\ b_{n1} & & & 0 \end{bmatrix}$



► **Output:**

An assignment matrix $X \in \Pi_n$

Example: Quadratic Assignment Problem

Quadratic assignment problem, QAP, in the trace formulation

$$\mu^* := \min_{X \in \Pi_n} \text{Tr} \left(A X B X^\top \right)$$

Π_n : set of $n \times n$ permutation matrices

A and B : $n \times n$ real symmetric matrices

- ▶ **Non-convex, quadratic** objective over the (discrete) set of permutation matrices
- ▶ Convex relaxations exist

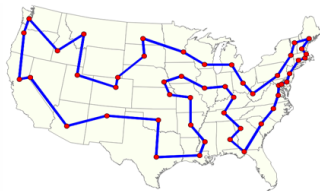
QAP example: Traveling Salesman Problem

Find a path passing from all vertices (e.g., cities) once to minimize the total trip time

$A = \frac{1}{2}D$, D : Matrix of edge weights such that $D_{ij} = D_{ji} \geq 0$ ($i \neq j$)

$B = C$ C : The adjacency matrix of the cities

$$TSP_{opt} := \min_{X \in \Pi_n} \text{Tr} \left(\frac{1}{2} D X C X^\top \right)$$



Implications of convexity

If f is convex,

- ▶ any local minimizer is also a **global minimizer**,
- ▶ we have a **principled step-size** selection,
- ▶ we can handle **non-smooth** problems like **constraints**.

Unfortunately, **convexity does not imply tractability...**

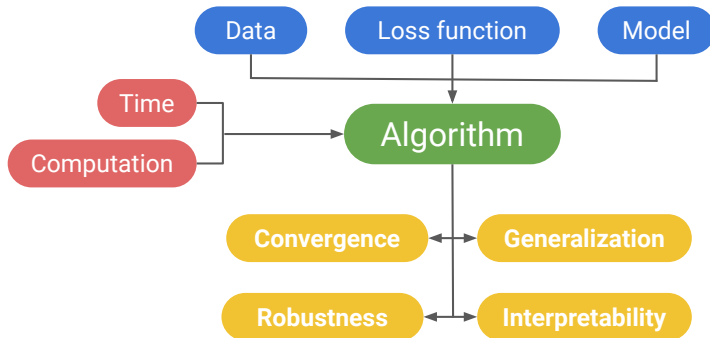
Implications of non-convexity

If f is not convex,

- ▶ Gradient descent converges only to **first order stationary points** ,
- ▶ Sometimes, they are **global minimizers**
- ▶ In certain applications, local minimizers can be **good enough**

Fortunately, **non-convexity does not imply intractability**
or **uselessness...**

Overview of mathematics of data



Do not forget!

- Lecture on Monday and recitation on Friday
 - ▶ Lecture: Basic probability theory and statistics.
 - ▶ Recitation: Terminology of optimization theory, gradient descent for logistic regression.

References

- [1] Yu. Nesterov.
Introductory Lectures on Convex Optimization: A Basic Course.
Kluwer, Boston, MA, 2004.