

Optimization for Machine Learning

Optimisation pour l'apprentissage automatique

Clément Royer

Université Paris-Dauphine

Master 2 IASD/ID Apprentissage



A warning

- This course will be given in **English**.
- The slides will be in English.
- The instructor is...French.

A warning

- This course will be given in **English**.
- The slides will be in English.
- The instructor is...French.

Why?

- It's in the syllabus!
- Latest advances in Machine Learning/Optimization are international;
- Both academic and industrial research are produced in English.

French touch

A warning

- This course will be given in **English**.
- The slides will be in English.
- The instructor is...French.

Why?

- It's in the syllabus!
- Latest advances in Machine Learning/Optimization are international;
- Both academic and industrial research are produced in English.

Aims of the course

- Present the main optimization tools used in ML;
- Motivate the use of these methods;
- Illustrate on typical ML problems.

Regarding this course

Lecturer : Clément Royer

- *Maître de conférences* at Dauphine since September 2019;
- From 2016 to 2019 : University of Wisconsin-Madison (USA);
- Research : Continuous optimization.

Useful information

- `clement.royer@dauphine.psl.eu`.
- [Link to these slides \(updated as we go\)](#).

URL:

<https://www.lamsade.dauphine.fr/~croyer/docs/courseOptiML.pdf>

Three-hour slots

- Week 1: 09/25 (8.30am-11.45am), 09/27 (8.30am-11.45am);
- Week 2: 10/02 (1.45pm-5pm), 10/04 (8.30am-11.45am);
- Week 3: 10/07 (1.45pm-5pm), 10/10 (1.45pm-5pm);
- Week 4: 11/06 (1.45pm-5pm), 11/08 (1.45pm-5pm);
- Week 5: **Exam on 11/15.**

Lab sessions

- On 10/10 and 11/08 for ID (instructor: Clément Royer);
- On 10/10 and 11/07 for IASD (instructor: Laurent Meunier).

Schedule

Three-hour slots

- Week 1: 09/25 (8.30am-11.45am), 09/27 (8.30am-11.45am);
- Week 2: 10/02 (1.45pm-5pm), 10/04 (8.30am-11.45am);
- Week 3: 10/07 (1.45pm-5pm), 10/10 (1.45pm-5pm);
- Week 4: 11/06 (1.45pm-5pm), 11/08 (1.45pm-5pm);
- Week 5: **Exam on 11/15.**

Lab sessions

- On 10/10 and 11/08 for ID (instructor: Clément Royer);
- On 10/10 and 11/07 for IASD (instructor: Laurent Meunier).

I expect to...

- Start/finish on time;
- Be able to hear everyone;
- **Get feedback from you.**

- 1 Introduction
- 2 Basics of optimization
- 3 Unconstrained optimization
- 4 Constrained optimization
- 5 Stochastic optimization
- 6 Nonsmooth optimization
- 7 Advanced topics

- 1 Introduction (Some examples)
- 2 Basics of optimization (Gradients and convexity)
- 3 Unconstrained optimization (Gradient descent)
- 4 Constrained optimization (ADMM)
- 5 Stochastic optimization (Stochastic gradient)
- 6 Nonsmooth optimization (Proximal methods)
- 7 Advanced topics (Second-order methods?)

- 1 Introduction
 - Optimization and ML
 - An example: text classification via Support Vector Machine
- 2 Basics of optimization
- 3 Unconstrained optimization

What you may have heard of/read about

- Data Analysis;
- Data Mining;
- Machine Learning (ML);
- Artificial Intelligence;
- Big Data;
- ...

What you may have heard of/read about

- Data Analysis;
- Data Mining;
- Machine Learning (ML);
- Artificial Intelligence;
- Big Data;
- ...

What this course is about

- Optimization for ML...
- ...and for all of **data science**.
- We will focus on **generic principles**.

Main goals

- Extract meaning/information from data:
Statistics, main features and structures;
- Use this information to predict behavior of yet unseen data.

Main goals

- Extract meaning/information from data:
Statistics, main features and structures;
- Use this information to predict behavior of yet unseen data.

Components of ML

- Statistics;
- Computer Science (data management, parallel computing, etc);
- **Optimization** for modeling and algorithms.

Optimization $\not\subset$ Machine Learning

- Optimization is a mathematical tool;
- Used in many areas: Economics, Chemistry, Physics, Social sciences,...
- Appears in other branches of (applied) mathematics: Linear Algebra, PDEs, Statistics, etc.

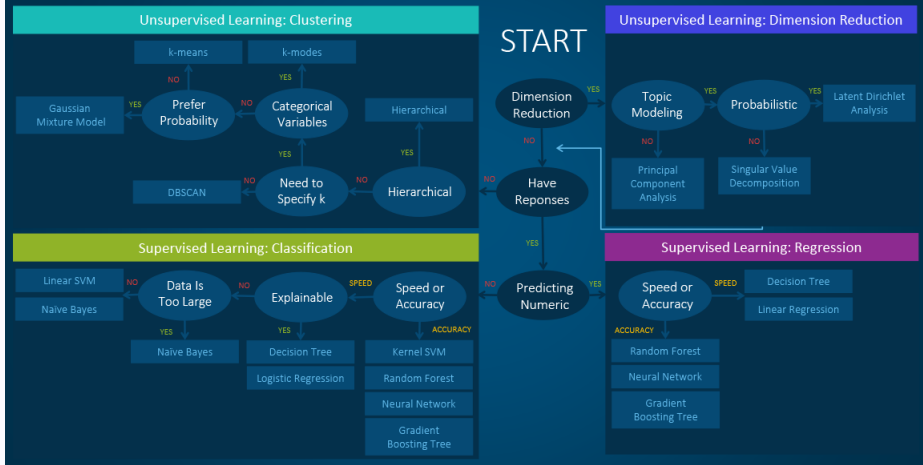
Optimization $\not\subset$ Machine Learning

- Optimization is a mathematical tool;
- Used in many areas: Economics, Chemistry, Physics, Social sciences,...
- Appears in other branches of (applied) mathematics: Linear Algebra, PDEs, Statistics, etc.

Machine Learning $\not\subset$ Optimization

- Optimization targets a certain problem;
- ML is not just about this problem;
- Other features of ML (data cleaning, hardware,...) will not appear in the optimization.

Machine Learning Algorithms Cheat Sheet



Source: <https://blogs.sas.com/content/subconsciousmusings/2017/04/12/machine-learning-algorithm-use/>

At first there was optimization...

- The rise of optimization: 1970-1980;
- Many algorithms have proven effective in various fields;
- *Standard practice for a physics-motivated problem*: run an interior-point Newton-type method (developed in the 2000s).

...then came ML!

From an optimization point of view:

- ML problems have challenging characteristics;
- The usual solvers are not so efficient in ML problems;
- But other (old) methods have regained interest.

Ubiquitous practice in ML: Run Stochastic Gradient with Momentum (1950s + a 1983 theoretical paper).

What changed?

Big data setting

- Very expensive to compute full derivatives/look at the entire data set;
- First-order methods have proven very effective to reach low accuracies.

What changed?

Big data setting

- Very expensive to compute full derivatives/look at the entire data set;
- First-order methods have proven very effective to reach low accuracies.

Community has changed

- The optimization problem is not everything;
- Interest in statistical properties of the solutions;
- Different analyzes and theoretical results.

- 1 Introduction
 - Optimization and ML
 - An example: text classification via Support Vector Machine
- 2 Basics of optimization
- 3 Unconstrained optimization

Given: A dataset $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$.

- \mathbf{x}_i is a **feature** vector in \mathbb{R}^d ;
- y_i is a **label**.

Statistical machine learning approach

Given: A dataset $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$.

- \mathbf{x}_i is a **feature** vector in \mathbb{R}^d ;
- y_i is a **label**.

Example: text classification

Using d words for classification:

- \mathbf{x}_i represents the words contained in a text document:

$$[\mathbf{x}_i]_j = \begin{cases} 1 & \text{if word } j \text{ is in document } i, \\ 0 & \text{otherwise.} \end{cases}$$

- y_i is equal to $+1$ if the document addresses a certain topic of interest, to -1 otherwise.

Learning process

- Given $\{(\mathbf{x}_i, y_i)\}_i$, discover a function $h : \mathbb{R}^d \rightarrow \mathbb{R}$ such that $h(\mathbf{x}_i) \approx y_i \forall i = 1, \dots, n$.
- Choose the predictor function h among a set \mathcal{H} parameterized by a vector $\mathbf{w} \in \mathbb{R}^d$: $\mathcal{H} = \{h \mid h = h(\cdot; \mathbf{w}), \mathbf{w} \in \mathbb{R}^{\hat{d}}\}$;

Learning process

- Given $\{(\mathbf{x}_i, y_i)\}_i$, discover a function $h : \mathbb{R}^d \rightarrow \mathbb{R}$ such that $h(\mathbf{x}_i) \approx y_i \forall i = 1, \dots, n$.
- Choose the predictor function h among a set \mathcal{H} parameterized by a vector $\mathbf{w} \in \mathbb{R}^d$: $\mathcal{H} = \{h \mid h = h(\cdot; \mathbf{w}), \mathbf{w} \in \mathbb{R}^{\hat{d}}\}$;

Linear model for text classification

- We seek a hyperplane in \mathbb{R}^d separating the feature vectors associated with $y_i = +1$ and those associated with $y_i = -1$;
- This corresponds to a linear model $h(\mathbf{x}) = \mathbf{x}^T \mathbf{u} - v$, and we want to choose $\mathbf{w}_1, \mathbf{w}_0$ such that:

$$\forall i = 1, \dots, n, \quad \begin{cases} \mathbf{x}_i^T \mathbf{u} - v \geq 1 & \text{if } y_i = +1 \\ \mathbf{x}_i^T \mathbf{u} - v \leq -1 & \text{if } y_i = -1. \end{cases}$$

Objective of the problem

An objective to optimize over

- Our goal is to penalize values of $\mathbf{w} = (\mathbf{u}, v)$ for which $h(\mathbf{x}_i) \neq y_i$.
- One possibility: the **hinge loss function**

$$\forall (h, y) \in \mathbb{R}^2, \quad \ell(h, y) = \max \{1 - yh, 0\}.$$

About the hinge loss

- $hy > 1 \Rightarrow \ell(h, y) = 0$: no penalty (h and y are of the same sign, $|h| > 1$ so this is a good prediction);
- $hy < -1 \Rightarrow \ell(h, y) > 2$: large penalty (h and y are of opposite sign and $|h| > 1$, this is a bad prediction);
- $|hy| \leq 1 \Rightarrow \ell(h, y) \in [0, 2]$: small penalty (h and y can be of the same sign, but the value of $|h|$ makes the prediction less certain).

An optimization problem

$$\min_{\mathbf{u}, v} \frac{1}{n} \sum_{i=1}^n \max \{1 - y_i(\mathbf{x}_i^T \mathbf{u} - v), 0\} \quad .$$

An optimization problem

$$\min_{\mathbf{u}, v} \frac{1}{n} \sum_{i=1}^n \max \{1 - y_i(\mathbf{x}_i^T \mathbf{u} - v), 0\} \quad .$$

- Minimize the sum of the losses for all examples;

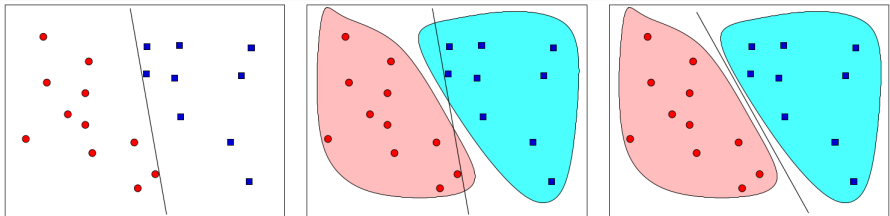
An optimization problem

$$\min_{\mathbf{u}, v} \frac{1}{n} \sum_{i=1}^n \max \{1 - y_i(\mathbf{x}_i^T \mathbf{u} - v), 0\} + \frac{\lambda}{2} \|\mathbf{u}\|_2^2.$$

for $\lambda \geq 0$.

- Minimize the sum of the losses for all examples;
- A regularizing term is usually added (more on that later).

Different solutions



Source: B. Recht and S. J. Wright, Nonlinear Optimization for Machine Learning (forthcoming).

- Red/Blue dots: data points labeled $+1/-1$;
- Red/Blue clouds: distribution of the text documents;
- Two linear classifiers;
- Rightmost plot: maximal-margin solution.

$$\min_{\mathbf{u}, v} \frac{1}{n} \sum_{i=1}^n \max \{1 - y_i(\mathbf{x}_i^T \mathbf{u} - v), 0\} + \frac{\lambda}{2} \|\mathbf{u}\|_2^2$$

Reformulation

- Add variables to replace the max \Rightarrow Convex quadratic program;
- Use duality \Rightarrow Convex quadratic program:

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^n} \frac{1}{2} \boldsymbol{\alpha}^T \mathbf{Q} \boldsymbol{\alpha} - \mathbf{1}^T \boldsymbol{\alpha} \quad \text{subject to} \quad 0 \leq \boldsymbol{\alpha} \leq \frac{1}{\lambda} \mathbf{1}, \mathbf{y}^T \boldsymbol{\alpha} = 0$$

with $\mathbf{Q} = [y_i y_j h(x_i) h(x_j)]_{ij}$, $\mathbf{y} = [y_1 \dots y_n]^T$, $\mathbf{1} = [1 \dots 1]^T \in \mathbb{R}^n$.

$$\min_{\mathbf{u}, v} \frac{1}{n} \sum_{i=1}^n \max \{1 - y_i(\mathbf{x}_i^T \mathbf{u} - v), 0\} + \frac{\lambda}{2} \|\mathbf{u}\|_2^2$$

Reformulation

- Add variables to replace the max \Rightarrow Convex quadratic program;
- Use duality \Rightarrow Convex quadratic program:

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^n} \frac{1}{2} \boldsymbol{\alpha}^T \mathbf{Q} \boldsymbol{\alpha} - \mathbf{1}^T \boldsymbol{\alpha} \quad \text{subject to} \quad 0 \leq \boldsymbol{\alpha} \leq \frac{1}{\lambda} \mathbf{1}, \mathbf{y}^T \boldsymbol{\alpha} = 0$$

with $\mathbf{Q} = [y_i y_j h(x_i) h(x_j)]_{ij}$, $\mathbf{y} = [y_1 \dots y_n]^T$, $\mathbf{1} = [1 \dots 1]^T \in \mathbb{R}^n$.

$$\min_{\mathbf{u}, v} \frac{1}{n} \sum_{i=1}^n \max \{1 - y_i(\mathbf{x}_i^T \mathbf{u} - v), 0\} + \frac{\lambda}{2} \|\mathbf{u}\|_2^2$$

Reformulation

- Add variables to replace the max \Rightarrow Convex quadratic program;
- Use duality \Rightarrow Convex quadratic program:

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^n} \frac{1}{2} \boldsymbol{\alpha}^T \mathbf{Q} \boldsymbol{\alpha} - \mathbf{1}^T \boldsymbol{\alpha} \quad \text{subject to} \quad 0 \leq \boldsymbol{\alpha} \leq \frac{1}{\lambda} \mathbf{1}, \quad \mathbf{y}^T \boldsymbol{\alpha} = 0$$

with $\mathbf{Q} = [y_i y_j h(x_i) h(x_j)]_{ij}$, $\mathbf{y} = [y_1 \dots y_n]^T$, $\mathbf{1} = [1 \dots 1]^T \in \mathbb{R}^n$.

Optimizers know how to solve this efficiently.

$$\min_{\mathbf{u}, v} \underbrace{\frac{1}{n} \sum_{i=1}^n \max \{1 - y_i(\mathbf{x}_i^T \mathbf{u} - v), 0\}}_{\text{loss}} + \underbrace{\frac{\lambda}{2} \|\mathbf{u}\|_2^2}_{\text{regularizer}} .$$

The key questions

- Are all solutions with “zero loss” equally good?
- We want to do good not only on our **training set** $\{(\mathbf{x}_i, y_i)\}$...
- ...but also on yet unseen data (from a similar distribution)!
- **In our example**, we want our classifier to apply to new text documents.

$$\min_{\mathbf{u}, v} \underbrace{\frac{1}{n} \sum_{i=1}^n \max \{1 - y_i(\mathbf{x}_i^T \mathbf{u} - v), 0\}}_{\text{loss}} + \underbrace{\frac{\lambda}{2} \|\mathbf{u}\|_2^2}_{\text{regularizer}} .$$

The key questions

- Are all solutions with “zero loss” equally good?
- We want to do good not only on our **training set** $\{(\mathbf{x}_i, y_i)\}$...
- ...but also on yet unseen data (from a similar distribution)!
- **In our example**, we want our classifier to apply to new text documents.

Optimizers may not be able to do that efficiently.

Takeaways from the example

- We formulate the optimization problem based on **observed data**;
- We want the solution to have properties with respect to **unseen data**;
- Optimization may help but is not the ultimate answer.

Takeaways from the example

- We formulate the optimization problem based on **observed data**;
- We want the solution to have properties with respect to **unseen data**;
- Optimization may help but is not the ultimate answer.

Other issues with ML problems

- What if the feature space is large (*all French/English words*)?
- What if the parameter space \mathbb{R}^n is huge (*all Wikipedia articles*)?
- What if linear models do not give good results?

Takeaways from the example

- We formulate the optimization problem based on **observed data**;
- We want the solution to have properties with respect to **unseen data**;
- Optimization may help but is not the ultimate answer.

Other issues with ML problems

- What if the feature space is large (*all French/English words*)?
Reduce dimensionality, look for sparse solutions.
- What if the parameter space \mathbb{R}^n is huge (*all Wikipedia articles*)?
Sampling/Batch/Stochastic methods.
- What if linear models do not give good results?
Nonlinear optimization (kernel SVM**).**

- Methodologies to solve given optimization problems;
- Focus on common structures in ML: finite sum, regularization;
- Discussion on properties of various formulations.

Focus on the optimization side

- Main algorithms and characteristics;
- Some applications, but always from an optimization perspective;
- Plenty of other data science courses in these Master programs!

1 Introduction

2 Basics of optimization

- Notation and background
- Optimization problem and optimality
- Convexity
- Optimization algorithms

3 Unconstrained optimization

- 1 Introduction
- 2 Basics of optimization
 - Notation and background
 - Optimization problem and optimality
 - Convexity
 - Optimization algorithms
- 3 Unconstrained optimization

For simplicity

- Optimization on real variables;
- Finite dimension;
- Canonical vector space structure.

I will use the following

- Scalars: a, b, c, \dots
- Vectors: $\mathbf{a}, \mathbf{b}, \mathbf{c}, \dots$
- Matrices: $\mathbf{A}, \mathbf{B}, \mathbf{C}, \dots$
- Sets: $\mathcal{A}, \mathcal{B}, \mathcal{C}, \dots$

- \mathbb{R}^d : set of vectors with $d \geq 1$ real components;
- For any $\mathbf{w} \in \mathbb{R}^d$ and $i \in \{1, \dots, d\}$, $w_i \in \mathbb{R}$ is the i -component of \mathbf{w} :
 $\mathbf{w} = [w_i]_{1 \leq i \leq d}$;
- Any $\mathbf{w} \in \mathbb{R}^d$ will be represented columnwise: $\mathbf{w} = \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}$;
- We will use row vectors as “transposed” (from column to row) of their column vectors counterpart: $\mathbf{w}^T := [w_1 \cdots w_d]$;

- \mathbb{R}^d : set of vectors with $d \geq 1$ real components;
- For any $\mathbf{w} \in \mathbb{R}^d$ and $i \in \{1, \dots, d\}$, $w_i \in \mathbb{R}$ is the i -component of \mathbf{w} :
 $\mathbf{w} = [w_i]_{1 \leq i \leq d}$;
- Any $\mathbf{w} \in \mathbb{R}^d$ will be represented columnwise: $\mathbf{w} = \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}$;
- We will use row vectors as “transposed” (from column to row) of their column vectors counterpart: $\mathbf{w}^T := [w_1 \cdots w_d]$;

Vector operations

- *Addition in \mathbb{R}^d : $\mathbf{w} + \mathbf{z} := [w_i + z_i]_{1 \leq i \leq d}$*
- *Multiply a vector in \mathbb{R}^d by a real number: $\lambda \mathbf{w} := [\lambda w_i]_{1 \leq i \leq d}$.*

Linear algebra (2)

Euclidean norm on \mathbb{R}^d

The Euclidean norm (or ℓ_2 norm) of a vector $\mathbf{w} \in \mathbb{R}^d$ is given by:

$$\|\mathbf{w}\| := \sqrt{\sum_{i=1}^d w_i^2}.$$

Scalar product on \mathbb{R}^d

The scalar product is defined for every $\mathbf{w}, \mathbf{z} \in \mathbb{R}^d$ by:

$$\mathbf{w}^T \mathbf{z} := \sum_{i=1}^d w_i z_i.$$

One thus has $\mathbf{w}^T \mathbf{z} = \mathbf{z}^T \mathbf{w} = \|\mathbf{w}\|^2$.

Matrices

- $\mathbb{R}^{n \times d}$: set of n -by- d matrices;
- $\mathbb{R}^{d \times 1} \simeq \mathbb{R}^d$.

Transposed matrix

Let $\mathbf{A} = [\mathbf{A}_{ij}] \in \mathbb{R}^{n \times d}$ be a matrix with n rows and d columns.

The *transposed matrix* of \mathbf{A} , denoted by \mathbf{A}^T , is the matrix with n rows and m columns such that

$$\forall i = 1, \dots, n, \forall j = 1, \dots, d, \quad [\mathbf{A}^T]_{ij} = \mathbf{A}_{ji}.$$

Matrices

- $\mathbb{R}^{n \times d}$: set of n -by- d matrices;
- $\mathbb{R}^{d \times 1} \simeq \mathbb{R}^d$.

Transposed matrix

Let $\mathbf{A} = [\mathbf{A}_{ij}] \in \mathbb{R}^{n \times d}$ be a matrix with n rows and d columns.

The *transposed matrix* of \mathbf{A} , denoted by \mathbf{A}^T , is the matrix with n rows and m columns such that

$$\forall i = 1, \dots, n, \forall j = 1, \dots, d, \quad [\mathbf{A}^T]_{ij} = \mathbf{A}_{ji}.$$

Squared matrix case

- $\mathbf{A}^T \in \mathbb{R}^{d \times d}$;
- \mathbf{A} is called a *symmetric matrix* if $\mathbf{A} = \mathbf{A}^T$.

Matrix inversion

A matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$ is *invertible* if it exists $\mathbf{B} \in \mathbb{R}^{d \times d}$ such that $\mathbf{BA} = \mathbf{AB} = \mathbf{I}_d$, where \mathbf{I}_d is the identity matrix of $\mathbb{R}^{d \times d}$.

In this case, \mathbf{B} is the unique matrix with this property: \mathbf{B} is called the *inverse matrix of \mathbf{A}* , and is denoted by \mathbf{A}^{-1} .

Positive (semi-)definiteness

A matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$ is *positive semidefinite* if

$$\forall \mathbf{x} \in \mathbb{R}^n, \quad \mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0.$$

It is called *positive definite* when $\mathbf{x}^T \mathbf{A} \mathbf{x} > 0$ for every nonzero vector \mathbf{x} .

Eigenvalues and eigenvectors

Let $\mathbf{A} \in \mathbb{R}^{d \times d}$. A real λ is called an *eigenvalue of \mathbf{A}* if

$$\exists \mathbf{v} \in \mathbb{R}^d, \|\mathbf{v}\| \neq 0, \quad \mathbf{A}\mathbf{v} = \lambda\mathbf{v}.$$

The vector \mathbf{v} is then called an *eigenvector of \mathbf{A}* (associated to the eigenvalue λ).

Eigenvalues and eigenvectors

Let $\mathbf{A} \in \mathbb{R}^{d \times d}$. A real λ is called an *eigenvalue of \mathbf{A}* if

$$\exists \mathbf{v} \in \mathbb{R}^d, \|\mathbf{v}\| \neq 0, \quad \mathbf{A}\mathbf{v} = \lambda\mathbf{v}.$$

The vector \mathbf{v} is then called an *eigenvector of \mathbf{A}* (associated to the eigenvalue λ).

Any symmetric matrix in $\mathbb{R}^{d \times d}$ possesses d real eigenvalues. Given two symmetric matrices $(\mathbf{A}, \mathbf{B}) \in \mathbb{R}^{d \times d}$, we introduce the following notations:

- $\lambda_{\min}(\mathbf{A})/\lambda_{\max}(\mathbf{A})$: smallest/largest eigenvalue of \mathbf{A} ;
- $\mathbf{A} \stackrel{n}{\succeq} \mathbf{B} \Leftrightarrow \lambda_{\min}(\mathbf{A}) \geq \lambda_{\max}(\mathbf{B})$;
- $\mathbf{A} \stackrel{n}{\succ} \mathbf{B} \Leftrightarrow \lambda_{\min}(\mathbf{A}) > \lambda_{\max}(\mathbf{B})$.

With these notations, \mathbf{A} is positive semi-definite (resp. positive definite) if and only if $\mathbf{A} \succeq 0$ (resp. $\mathbf{A} \succ 0$).

We consider a **smooth** function $f : \mathbb{R}^d \rightarrow \mathbb{R}$.

We consider a **smooth** function $f : \mathbb{R}^d \rightarrow \mathbb{R}$.

First-order derivative

If f is continuously differentiable on \mathbb{R}^d , one defines for any $\mathbf{w} \in \mathbb{R}^d$ the **gradient of f at \mathbf{w}** by

$$\nabla f(\mathbf{w}) := \left[\frac{\partial f}{\partial w_i} \right]_{1 \leq i \leq d} \in \mathbb{R}^d.$$

The set of continuously differentiable functions will be denoted by $\mathcal{C}^1 \stackrel{n}{=} \mathcal{C}^1(\mathbb{R}^d, \mathbb{R})$.

We consider a **smooth** function $f : \mathbb{R}^d \rightarrow \mathbb{R}$.

We consider a **smooth** function $f : \mathbb{R}^d \rightarrow \mathbb{R}$.

Second-order derivative

If f is *twice* continuously differentiable on \mathbb{R}^d , one defines for any $\mathbf{w} \in \mathbb{R}^d$ the **Hessian of f at \mathbf{w}** by

$$\nabla^2 f(\mathbf{w}) := \left[\frac{\partial^2 f}{\partial w_i \partial w_j} \right]_{1 \leq i, j \leq d} \in \mathbb{R}^{d \times d}.$$

This matrix is symmetric.

The set of twice continuously differentiable functions will be denoted by \mathcal{C}^2 .

First-order Taylor expansions

If $f \in \mathcal{C}^1$, for any $\mathbf{w}, \mathbf{h} \in \mathbb{R}^d$,

$$\begin{cases} f(\mathbf{w} + \mathbf{h}) = f(\mathbf{w}) + \nabla f(\mathbf{w} + t\mathbf{h})^\top \mathbf{h} & \text{for some } t \in (0, 1) \\ f(\mathbf{w} + \mathbf{h}) = f(\mathbf{w}) + \int_0^1 \nabla f(\mathbf{w} + t\mathbf{h})^\top \mathbf{h} dt. \end{cases}$$

First-order Taylor expansions

If $f \in \mathcal{C}^1$, for any $\mathbf{w}, \mathbf{h} \in \mathbb{R}^d$,

$$\begin{cases} f(\mathbf{w} + \mathbf{h}) = f(\mathbf{w}) + \nabla f(\mathbf{w} + t\mathbf{h})^\top \mathbf{h} & \text{for some } t \in (0, 1) \\ f(\mathbf{w} + \mathbf{h}) = f(\mathbf{w}) + \int_0^1 \nabla f(\mathbf{w} + t\mathbf{h})^\top \mathbf{h} dt. \end{cases}$$

Second-order Taylor expansions

If $f \in \mathcal{C}^2$, for any $\mathbf{w}, \mathbf{h} \in \mathbb{R}^d$,

$$\begin{cases} f(\mathbf{w} + \mathbf{h}) = f(\mathbf{w}) + \nabla f(\mathbf{w})^\top \mathbf{h} + \frac{1}{2} \mathbf{h}^\top \nabla^2 f(\mathbf{w} + t\mathbf{h}) \mathbf{h} \\ \quad \text{for some } t \in (0, 1) \\ f(\mathbf{w} + \mathbf{h}) = f(\mathbf{w}) + \nabla f(\mathbf{w})^\top \mathbf{h} + \frac{1}{2} \int_0^1 \mathbf{h}^\top \nabla^2 f(\mathbf{w} + t\mathbf{h}) \mathbf{h} dt. \end{cases}$$

Definition

A function $g : \mathbb{R}^d \rightarrow \mathbb{R}^m$ is L -Lipschitz continuous if it exists $L > 0$ such that

$$\forall (\mathbf{w}, \mathbf{z}) \in (\mathbb{R}^d)^2, \quad \|g(\mathbf{w}) - g(\mathbf{z})\| \leq L \|\mathbf{w} - \mathbf{z}\|.$$

The value L is called a Lipschitz constant for g .

- Ex) Any linear function is Lipschitz continuous;
- $\mathcal{C}_L^{1,1}$: set continuously differentiable functions with L -Lipschitz continuous first-order derivative;
- $\mathcal{C}_L^{2,2}$: set of twice continuously differentiable functions with L -Lipschitz continuous second-order derivative.

First-order Taylor bound

Let $f \in \mathcal{C}_L^{1,1}$. For any $\mathbf{w}, \mathbf{h} \in \mathbb{R}^d$,

$$f(\mathbf{w} + \mathbf{h}) \leq f(\mathbf{w}) + \nabla f(\mathbf{w})^\top \mathbf{h} + \frac{L}{2} \|\mathbf{h}\|^2.$$

First-order Taylor bound

Let $f \in \mathcal{C}_L^{1,1}$. For any $\mathbf{w}, \mathbf{h} \in \mathbb{R}^d$,

$$f(\mathbf{w} + \mathbf{h}) \leq f(\mathbf{w}) + \nabla f(\mathbf{w})^\top \mathbf{h} + \frac{L}{2} \|\mathbf{h}\|^2.$$

⇒ One of the two key inequalities in optimization.

First-order Taylor bound

Let $f \in \mathcal{C}_L^{1,1}$. For any $\mathbf{w}, \mathbf{h} \in \mathbb{R}^d$,

$$f(\mathbf{w} + \mathbf{h}) \leq f(\mathbf{w}) + \nabla f(\mathbf{w})^\top \mathbf{h} + \frac{L}{2} \|\mathbf{h}\|^2.$$

\Rightarrow **One of the two key inequalities in optimization.**

Second-order Taylor bound

Let $f \in \mathcal{C}_L^{2,2}$. For any $\mathbf{w}, \mathbf{h} \in \mathbb{R}^d$,

$$f(\mathbf{w} + \mathbf{h}) \leq f(\mathbf{w}) + \nabla f(\mathbf{w})^\top \mathbf{h} + \frac{1}{2} \mathbf{h}^\top \nabla^2 f(\mathbf{w}) \mathbf{h} + \frac{L}{6} \|\mathbf{h}\|^3,$$

Some references

- Plenty of lecture notes, courses freely available;
- **Appendix material** of many optimization (and some ML) textbooks!

Examples (subject to updates)

- In French:
<https://www.lpsm.paris/pageperso/bolley/poly-cdiff.pdf>
<https://www.lpsm.paris/pageperso/bolley/poly-algebre3.pdf>
- In English:
<http://vmls-book.stanford.edu/vmls.pdf> (Chapters 1-3)
https://sebastianraschka.com/pdf/books/dlb/appendix_d_calculus.pdf.

- 1 Introduction
- 2 Basics of optimization
 - Notation and background
 - Optimization problem and optimality
 - Convexity
 - Optimization algorithms
- 3 Unconstrained optimization

What's optimization?

- Operations research;
- Decision-making;
- Decision sciences;
- Mathematical programming;
- Mathematical optimization.

⇒ All of these can be considered as optimization.

What's optimization?

- Operations research;
- Decision-making;
- Decision sciences;
- Mathematical programming;
- Mathematical optimization.

⇒ All of these can be considered as optimization.

My definition

The purpose of optimization is to make the best decision out of a set of alternatives.

Formulation of an optimization problem

A **minimization problem** of d real parameters is written as follows :

$$\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{x}) \quad \text{subject to } \mathbf{w} \in \mathcal{F}$$

Formulation of an optimization problem

A **minimization problem** of d real parameters is written as follows :

$$\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{x}) \quad \text{subject to } \mathbf{w} \in \mathcal{F}$$

- \mathbf{w} represents the optimization variable(s);
- d is the dimension of the problem (we will assume $d \geq 1$);
- $f(\cdot)$ is the **objective/cost/loss** function;
- \mathcal{F} is the constraint/feasible set.

Formulation of an optimization problem

A **minimization problem** of d real parameters is written as follows :

$$\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{x}) \quad \text{subject to } \mathbf{w} \in \mathcal{F}$$

- \mathbf{w} represents the optimization variable(s);
- d is the dimension of the problem (we will assume $d \geq 1$);
- $f(\cdot)$ is the **objective/cost/loss** function;
- \mathcal{F} is the constraint/feasible set.

Maximizing f is equivalent to minimizing $-f$.

Local and global solutions

$$\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) \quad \text{subject to } \mathbf{w} \in \mathcal{F}$$

Local minimum (also called minimizer)

- A point \mathbf{w}^* is a **local minimum** of the problem if there exists a neighborhood \mathcal{N} of \mathbf{w}^* such that $f(\mathbf{w}^*) \leq f(\mathbf{w}) \forall \mathbf{w} \in \mathcal{N} \cap \mathcal{F}$;
- A local minimum such that $f(\mathbf{w}^*) < f(\mathbf{w}) \forall \mathbf{w} \in \mathcal{N} \cap \mathcal{F}, \mathbf{w} \neq \mathbf{w}^*$ is called a strict local minimum.

Global minimum

A point \mathbf{w}^* is a **global minimum** of the problem if $f(\mathbf{w}^*) \leq f(\mathbf{w}) \forall \mathbf{w} \in \mathcal{F}$.

Local and global solutions (2)

- In general, finding global solutions is hard;
- Local solutions can also be hard to find.

Local and global solutions (2)

- In general, finding global solutions is hard;
- Local solutions can also be hard to find.

Tractable cases

- When the objective function behaves nicely;
- Suitable properties of the constraint set (more on that in the constrained optimization lecture).

Unconstrained problem: $\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$,
 f continuously differentiable.

Unconstrained problem: $\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$,
 f continuously differentiable.

First-order necessary condition

If \mathbf{w}^* is a local minimum of the problem, then

$$\|\nabla f(\mathbf{w}^*)\| = 0.$$

Unconstrained problem: $\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$,
 f continuously differentiable.

First-order necessary condition

If \mathbf{w}^* is a local minimum of the problem, then

$$\|\nabla f(\mathbf{w}^*)\| = 0.$$

- This condition is only necessary;
- A point such that $\|\nabla f(\mathbf{w}^*)\| = 0$ can also be a local maximum or a saddle point.

Unconstrained problem: $\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$,
 f twice continuously differentiable.

Optimality conditions for unconstrained optimization (2)

Unconstrained problem: $\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$,
 f twice continuously differentiable.

Second-order necessary condition

If \mathbf{w}^* is a local minimum of the problem, **then**

$$\|\nabla f(\mathbf{w}^*)\| = 0 \quad \text{and} \quad \nabla^2 f(\mathbf{w}^*) \succeq 0.$$

Optimality conditions for unconstrained optimization (2)

Unconstrained problem: $\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$,
 f twice continuously differentiable.

Second-order necessary condition

If \mathbf{w}^* is a local minimum of the problem, **then**

$$\|\nabla f(\mathbf{w}^*)\| = 0 \quad \text{and} \quad \nabla^2 f(\mathbf{w}^*) \succeq 0.$$

Second-order sufficient condition

If \mathbf{w}^* is such that

$$\|\nabla f(\mathbf{w}^*)\| = 0 \quad \text{and} \quad \nabla^2 f(\mathbf{w}^*) \succ 0,$$

then it is a local minimum of the problem.

- 1 Introduction
- 2 Basics of optimization
 - Notation and background
 - Optimization problem and optimality
 - Convexity
 - Optimization algorithms
- 3 Unconstrained optimization

Convex set

A set $\mathcal{C} \in \mathbb{R}^d$ is called **convex** if

$$\forall (\mathbf{u}, \mathbf{v}) \in \mathcal{C}^2, \forall t \in [0, 1], \quad t\mathbf{u} + (1 - t)\mathbf{v} \in \mathcal{C}.$$

Convex set

A set $\mathcal{C} \in \mathbb{R}^d$ is called **convex** if

$$\forall (\mathbf{u}, \mathbf{v}) \in \mathcal{C}^2, \forall t \in [0, 1], \quad t\mathbf{u} + (1 - t)\mathbf{v} \in \mathcal{C}.$$

Examples:

- \mathbb{R}^d ;
- Line segment: $\{t\mathbf{w} | t \in \mathbb{R}\}$ for some $\mathbf{w} \in \mathbb{R}^d$;
- Sphere: $\{\mathbf{w} \in \mathbb{R}^d | \sum_i [\mathbf{w}_i]^2 \leq 1\}$.

Generic definition

A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is **convex** if

$$\forall (\mathbf{u}, \mathbf{v}) \in (\mathbb{R}^d)^2, \forall t \in [0, 1], \quad f(t\mathbf{u} + (1 - t)\mathbf{v}) \leq t f(\mathbf{u}) + (1 - t) f(\mathbf{v}).$$

Generic definition

A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is **convex** if

$$\forall (\mathbf{u}, \mathbf{v}) \in (\mathbb{R}^d)^2, \forall t \in [0, 1], \quad f(t\mathbf{u} + (1-t)\mathbf{v}) \leq tf(\mathbf{u}) + (1-t)f(\mathbf{v}).$$

Examples:

- Linear function: $f(\mathbf{w}) = \mathbf{a}^T \mathbf{w} + b$;
- Squared Euclidean norm: $f(\mathbf{w}) = \|\mathbf{w}\|^2 = \mathbf{w}^T \mathbf{w}$.

Convexity and gradient

A continuously differentiable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex if and only if

$$\forall \mathbf{u}, \mathbf{v} \in \mathbb{R}^d, \quad f(\mathbf{v}) \geq f(\mathbf{u}) + \nabla f(\mathbf{u})^T (\mathbf{v} - \mathbf{u}).$$

Convexity and gradient

A continuously differentiable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex if and only if

$$\forall \mathbf{u}, \mathbf{v} \in \mathbb{R}^d, \quad f(\mathbf{v}) \geq f(\mathbf{u}) + \nabla f(\mathbf{u})^T (\mathbf{v} - \mathbf{u}).$$

The other key inequality in optimization.

Convexity and gradient

A continuously differentiable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex if and only if

$$\forall \mathbf{u}, \mathbf{v} \in \mathbb{R}^d, \quad f(\mathbf{v}) \geq f(\mathbf{u}) + \nabla f(\mathbf{u})^T (\mathbf{v} - \mathbf{u}).$$

The other key inequality in optimization.

Convexity and Hessian

A twice continuously differentiable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex if and only if for every $\mathbf{w} \in \mathbb{R}^d$, $\nabla^2 f(\mathbf{w}) \succeq 0$.

Convex optimization problem

$$\min_{\mathbf{w} \in \mathcal{X}} f(\mathbf{w}), f \text{ convex}, \mathcal{X} \subset \mathbb{R}^d \text{ closed+convex.}$$

$$\min_{\mathbf{w} \in \mathcal{X}} f(\mathbf{w}), f \text{ convex}, \mathcal{X} \subset \mathbb{R}^d \text{ closed+convex.}$$

Theorem

Every local minimum of a f is a global minimum.

$$\min_{\mathbf{w} \in \mathcal{X}} f(\mathbf{w}), f \text{ convex}, \mathcal{X} \subset \mathbb{R}^d \text{ closed+convex.}$$

Theorem

Every local minimum of a f is a global minimum.

Corollary

If f is continuously differentiable, every point \mathbf{w}^* such that $\|\nabla f(\mathbf{w}^*)\| = 0$ is a global minimum.

Definition

A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ in \mathcal{C}^1 is μ -strongly convex (or *strongly convex of modulus $\mu > 0$*) if for all $(\mathbf{u}, \mathbf{v}) \in (\mathbb{R}^d)^2$ and $t \in [0, 1]$,

$$f(t\mathbf{u} + (1-t)\mathbf{v}) \leq t f(\mathbf{u}) + (1-t)f(\mathbf{v}) - \frac{\mu}{2}t(1-t)\|\mathbf{v} - \mathbf{u}\|^2.$$

Definition

A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ in \mathcal{C}^1 is μ -strongly convex (or *strongly convex of modulus $\mu > 0$*) if for all $(\mathbf{u}, \mathbf{v}) \in (\mathbb{R}^d)^2$ and $t \in [0, 1]$,

$$f(t\mathbf{u} + (1-t)\mathbf{v}) \leq t f(\mathbf{u}) + (1-t)f(\mathbf{v}) - \frac{\mu}{2}t(1-t)\|\mathbf{v} - \mathbf{u}\|^2.$$

Theorem

Any strongly convex function in \mathcal{C}^1 has a unique global minimizer.

Strong convexity (2)

Gradient and strong convexity

Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$, $f \in \mathcal{C}^1$. Then,

$$\forall \mathbf{u}, \mathbf{v} \in \mathbb{R}^d, \quad f(\mathbf{v}) \geq f(\mathbf{u}) + \nabla f(\mathbf{u})^T (\mathbf{v} - \mathbf{u}) + \frac{\mu}{2} \|\mathbf{v} - \mathbf{u}\|^2.$$

Hessian and strong convexity

Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$, $f \in \mathcal{C}^2$. Then,

$$f \text{ is } \mu\text{-strongly convex} \iff \nabla^2 f(\mathbf{w}) \succeq \mu \mathbf{I} \quad \forall \mathbf{w} \in \mathbb{R}^d.$$

Examples of (strongly) convex problems

Minimize a convex quadratic

$$\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) := \frac{1}{2} \mathbf{w}^T \mathbf{A} \mathbf{w} + \mathbf{b}^T \mathbf{w}, \quad \mathbf{A} \succeq 0.$$

- $\nabla^2 f(\mathbf{w}) = \mathbf{A}$;
- Strongly convex if $\mathbf{A} \succ 0$, with $\mu = \lambda_{\min}(\mathbf{A})$.

Examples of (strongly) convex problems

Minimize a convex quadratic

$$\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) := \frac{1}{2} \mathbf{w}^T \mathbf{A} \mathbf{w} + \mathbf{b}^T \mathbf{w}, \quad \mathbf{A} \succeq 0.$$

- $\nabla^2 f(\mathbf{w}) = \mathbf{A}$;
- Strongly convex if $\mathbf{A} \succ 0$, with $\mu = \lambda_{\min}(\mathbf{A})$.

Projection onto a closed, convex set

$$\min_{\mathbf{w} \in \mathcal{X}} \frac{1}{2} \|\mathbf{w} - \mathbf{a}\|^2, \quad \mathcal{X} \text{ closed, convex.}$$

- The objective is 1-strongly convex \Rightarrow the problem has a unique solution;
- Generalization of the case $\mathcal{X} = \mathbb{R}^d$.

- 1 Introduction
- 2 Basics of optimization
 - Notation and background
 - Optimization problem and optimality
 - Convexity
 - Optimization algorithms
- 3 Unconstrained optimization

Three ways to study optimization problems

- **Mathematical** : Prove existence of solutions, well-posedness of a problem. Study complex optimization formulations.
- **Computational** : Write a piece of software to solve specific or generic optimization problems in practice.
- **Algorithmic** : Design algorithms, establish theoretical guarantees and validate their practical implementation.

Three ways to study optimization problems

- **Mathematical** : Prove existence of solutions, well-posedness of a problem. Study complex optimization formulations.
- **Computational** : Write a piece of software to solve specific or generic optimization problems in practice.
- **Algorithmic** : Design algorithms, establish theoretical guarantees and validate their practical implementation.

This course is about the third category.

How to solve an optimization problem?

The ideal approach

- Find the solutions of $\|\nabla f(\mathbf{w})\| = 0$;
- Choose the one with the lowest function value.

How to solve an optimization problem?

The ideal approach

- Find the solutions of $\|\nabla f(\mathbf{w})\| = 0$;
- Choose the one with the lowest function value.

What's wrong with that?

- Solving a nonlinear equation directly is hard;
- There can be infinitely many solutions;
- The procedure has to be implemented eventually.

Iterative procedures

- Driving principle : given the current solution, move towards a (potentially) better point;
- Requires a certain amount of calculation at every iteration.

Our goal in the rest of the course

- Propose several algorithms;
- Analyze their theoretical behavior and guarantees;
- Check their practical appeal (lab sessions).

What do we expect?

In order to solve $\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$, we hope to achieve one of the following:

- 1 The iterates should get close to a solution;
- 2 The function values should get close to the optimum;
- 3 The optimality conditions should get close to be satisfied.

What do we expect?

In order to solve $\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$, we hope to achieve one of the following:

- 1 The iterates should get close to a solution;
- 2 The function values should get close to the optimum;
- 3 The optimality conditions should get close to be satisfied.

Convergence of iterates

The method generates a sequence of points (iterates) $\{\mathbf{w}_k\}_k$ such that

$$\|\mathbf{w}_k - \mathbf{w}^*\| \rightarrow 0 \quad \text{when } k \rightarrow \infty,$$

where \mathbf{w}^* is an optimal value of the problem.

(Typical of (strongly) convex functions.)

What do we expect? ('ed)

Convergence in function value

$$f(\mathbf{w}_k) \rightarrow f^* \quad \text{when } k \rightarrow \infty,$$

where f^* is the optimal value of the problem.

(Typical of (strongly) convex functions.)

What do we expect? ('ed)

Convergence in function value

$$f(\mathbf{w}_k) \rightarrow f^* \quad \text{when } k \rightarrow \infty,$$

where f^* is the optimal value of the problem.

(Typical of (strongly) convex functions.)

Convergence to a stationary point for differentiable f

$$\|\nabla f(\mathbf{w}_k)\| \rightarrow 0 \quad \text{when } k \rightarrow \infty.$$

More generic condition.

Why these conditions?

Unlike in theory, in practice:

- We do not know the optimal solution(s);
- We do not know the optimal value.

Why these conditions?

Unlike in theory, in practice:

- We do not know the optimal solution(s);
- We do not know the optimal value.

From an algorithmic standpoint,

- We can measure the behavior of the iterates;
- We can evaluate the objective and try to decrease it iteratively;
- We can evaluate/estimate the gradient norm and measure its decrease to zero.

- In optimization, classical results are asymptotic:

$$\|\nabla f(\mathbf{w}_k)\| \rightarrow 0 \quad \text{when } k \rightarrow \infty.$$

Remark: Convergence and convergence rates

- In optimization, classical results are asymptotic:

$$\|\nabla f(\mathbf{w}_k)\| \rightarrow 0 \quad \text{when } k \rightarrow \infty.$$

- **Global convergence rates** are now very popular:

$$\|\nabla f(\mathbf{w}_k)\| = \mathcal{O}\left(\frac{1}{k}\right) \quad \Leftrightarrow \quad \exists C > 0, \|\nabla f(\mathbf{w}_k)\| \leq \frac{C}{k} \quad \forall k.$$

- Common in convex optimization;
- Standard in theoretical computer science/statistics.

Optimizers code in...

- C/C++/Fortran (high-performance computing)
- Matlab, Python (prototyping);
- Julia.

Optimizers code in...

- C/C++/Fortran (high-performance computing)
- Matlab, Python (prototyping);
- Julia.

Specific optimization modeling languages

- GAMS, AMPL, CVX are broad-spectrum languages;
- MATPOWER, PyTorch are domain-oriented;
- Can be interfaced with the languages above.

Modeling framework

- Objective, constraints;
- Characterization of the solutions.

Conclusions: basics of optimization

Modeling framework

- Objective, constraints;
- Characterization of the solutions.

Important tools

- Derivatives and Taylor expansion;
- Convexity.

Conclusions: basics of optimization

Modeling framework

- Objective, constraints;
- Characterization of the solutions.

Important tools

- Derivatives and Taylor expansion;
- Convexity.

Algorithmic principle

- Iterative process: find a sequence of points that leads to a solution;
- Quantify how fast.

- 1 Introduction
- 2 Basics of optimization
- 3 Unconstrained optimization**
 - Linear least squares
 - Gradient descent method

Our problem today

$$\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}).$$

Assumptions

- f bounded below by f^* ;
- f smooth \Rightarrow derivatives can be used to solve this problem.

Two categories

Least squares

- Heavily relies on **linear algebra**;
- Can precisely characterize the solution(s);
- Application: Linear regression problems.

Generic smooth unconstrained problems

- One tool from analysis: the gradient;
- Goal: Converge iteratively towards a solution;
- Application(s): Logistic regression (among others).

Aims of this lecture

- Survey classical techniques;
 - Illustrate how they may be used in ML problems.
-
- Highlight the role of the gradient and that of convexity;
 - **Show** how convergence rates are obtained.

- 1 Introduction
- 2 Basics of optimization
- 3 **Unconstrained optimization**
 - Linear least squares
 - Gradient descent method

Data

- Dataset with n elements (individuals, trials, samples, etc);
- Every element i is characterized by a vector $\mathbf{x}_i \in \mathbb{R}^d$ of **features** and a label $y_i \in \mathbb{R}$.

$$\Rightarrow \text{Matrix } \mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix} \in \mathbb{R}^{n \times d} \text{ and vector } \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}.$$

Data

- Dataset with n elements (individuals, trials, samples, etc);
- Every element i is characterized by a vector $\mathbf{x}_i \in \mathbb{R}^d$ of **features** and a label $y_i \in \mathbb{R}$.

$$\Rightarrow \text{Matrix } \mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix} \in \mathbb{R}^{n \times d} \text{ and vector } \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}.$$

Goal

We seek a **linear** predictor function $h : \mathbf{x} \mapsto \mathbf{x}^T \mathbf{w}$ that correctly predicts y_i from \mathbf{x}_i .

- Linear models often provide a good first approximation;
- Relies on linear algebra, a rich area both theoretically and computationally.

Ideal predictor

- Would achieve $h(\mathbf{x}_i) = \mathbf{x}_i^T \mathbf{w} = y_i$ for every i ;
- These n equations can be written under the form of a linear system:
 $\mathbf{X}\mathbf{w} = \mathbf{y}$.

Ideal predictor

- Would achieve $h(\mathbf{x}_i) = \mathbf{x}_i^T \mathbf{w} = y_i$ for every i ;
- These n equations can be written under the form of a linear system:
 $\mathbf{X}\mathbf{w} = \mathbf{y}$.

Solving linear systems of equations

- A purely linear algebra problem;
- The solution is completely characterized by the properties of \mathbf{X} and \mathbf{y} .

Here's the catch

A dataset

- $\mathbf{x}_1 = \mathbf{x}_2 = \dots = \mathbf{x}_n = \mathbf{1}$ ($d = 1$);
- y_1, \dots, y_n are distinct (typical of **noisy** measurements).

Here's the catch

A dataset

- $\mathbf{x}_1 = \mathbf{x}_2 = \dots = \mathbf{x}_n = 1$ ($d = 1$);
- y_1, \dots, y_n are distinct (typical of **noisy** measurements).

Fitting a linear model

- We seek $\mathbf{w} = w \in \mathbb{R}$ such that $\mathbf{x}_i^T \mathbf{w} = x_i w = y_i \forall i$;
- The corresponding linear system is:

$$\begin{cases} w = y_1 \\ w = y_2 \\ \vdots \\ w = y_n \end{cases}$$

Here's the catch

A dataset

- $\mathbf{x}_1 = \mathbf{x}_2 = \dots = \mathbf{x}_n = 1$ ($d = 1$);
- y_1, \dots, y_n are distinct (typical of **noisy** measurements).

Fitting a linear model

- We seek $\mathbf{w} = w \in \mathbb{R}$ such that $\mathbf{x}_i^T \mathbf{w} = x_i w = y_i \forall i$;
- The corresponding linear system is:

$$\begin{cases} w = y_1 \\ w = y_2 \\ \vdots \\ w = y_n \end{cases}$$

- This system does not have a solution!
- Yet it is possible to compute a solution to the “data fitting” problem.

Problem formulation

Given a data set $\{(\mathbf{x}_i, y_i)\}_{1 \leq i \leq n}$ where $\mathbf{x}_i \in \mathbb{R}^d$, compute $\mathbf{w}^* \in \mathbb{R}^d$ as a solution of

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 = \frac{1}{2} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}),$$

$$\text{where } \mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix} \in \mathbb{R}^{n \times d} \text{ and } \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \in \mathbb{R}^n.$$

Problem formulation

Given a data set $\{(\mathbf{x}_i, y_i)\}_{1 \leq i \leq n}$ where $\mathbf{x}_i \in \mathbb{R}^d$, compute $\mathbf{w}^* \in \mathbb{R}^d$ as a solution of

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 = \frac{1}{2}(\mathbf{X}\mathbf{w} - \mathbf{y})^T(\mathbf{X}\mathbf{w} - \mathbf{y}),$$

where $\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix} \in \mathbb{R}^{n \times d}$ and $\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \in \mathbb{R}^n$.

Characteristics

- Unconstrained optimization problem;
- Nonnegative objective function (values bounded below by 0);
- Smooth: polynomial in the coefficients of \mathbf{w} .

How to solve linear least squares

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{X} \mathbf{w} - \mathbf{y}\|^2.$$

How to solve linear least squares

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2.$$

- If \mathbf{w}^* is a solution of the linear system $\mathbf{X}\mathbf{w} = \mathbf{y}$, then it is a solution of the least-squares problem!

How to solve linear least squares

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2.$$

- If \mathbf{w}^* is a solution of the linear system $\mathbf{X}\mathbf{w} = \mathbf{y}$, then it is a solution of the least-squares problem!
- What happens when the system has no solution?

Solving a linear system: the nice case

Squared linear system

$\mathbf{X} \mathbf{w} = \mathbf{y}$, avec $\mathbf{X} \in \mathbb{R}^{n \times d}$ and $n = d$.

Case 1: \mathbf{X} possesses an inverse

$$\mathbf{X} \mathbf{w} = \mathbf{y} \Leftrightarrow \mathbf{w} = \mathbf{X}^{-1} \mathbf{y}.$$

The system possesses a unique solution $\mathbf{w}^* = \mathbf{X}^{-1} \mathbf{y}$, which is also the global minimum of the least-squares problem $\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{X} \mathbf{w} - \mathbf{y}\|^2$.

Example with $d = n = 2$

$$\begin{cases} w_1 + w_2 &= 0, \\ 3w_1 + 2w_2 &= 1. \end{cases}$$

The unique solution is $\mathbf{w} = [1 \ -1]^T$.

Solving a linear system: the other cases

Squared linear system

$\mathbf{X} \mathbf{w} = \mathbf{y}$, avec $\mathbf{X} \in \mathbb{R}^{n \times d}$ and $n = d$.

Case 2: \mathbf{X} is not invertible

- There could be no solution;
- There could be infinitely many.

In both cases, we can compute a solution in the least-squares sense!

Solving a linear system: the other cases

Squared linear system

$\mathbf{X} \mathbf{w} = \mathbf{y}$, avec $\mathbf{X} \in \mathbb{R}^{n \times d}$ and $n = d$.

Case 2: \mathbf{X} is not invertible

- There could be no solution;
- There could be infinitely many.

In both cases, we can compute a solution in the least-squares sense!

Other cases

- $\mathbf{X} \mathbf{w} = \mathbf{y}$, avec $\mathbf{X} \in \mathbb{R}^{n \times d}$, $n \neq d$;
- Can have no solution, one or infinitely many!

What we want

- An analogous of the inverse;
- Provides a solution when there exists one (or infinitely many).

What we want

- An analogous of the inverse;
- Provides a solution when there exists one (or infinitely many).

Pseudo-inverse

Given a matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$, there exists a matrix $\mathbf{A} \in \mathbb{R}^{d \times n}$ that satisfies the Moore-Penrose equations:

$$\begin{cases} \mathbf{AXA} = \mathbf{A} \\ \mathbf{XAX} = \mathbf{X} \end{cases} \quad \text{and} \quad \begin{cases} (\mathbf{AX})^T = \mathbf{AX} \\ (\mathbf{XA})^T = \mathbf{XA} \end{cases}$$

This matrix is called the pseudo-inverse of \mathbf{X} , and we note $\mathbf{A} = \mathbf{X}^\dagger$.
If \mathbf{X} is invertible, $\mathbf{X}^\dagger = \mathbf{X}^{-1}$.

$$\mathbf{X} \in \mathbb{R}^{n \times d}, \quad \mathbf{y} \in \mathbb{R}^n.$$

Theorem

For any $\mathbf{y} \in \mathbb{R}^n$, $\mathbf{X}^\dagger \mathbf{y}$ is the solution of the least-squares problem

$$\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) = \frac{1}{2} \|\mathbf{X} \mathbf{w} - \mathbf{y}\|^2.$$

with minimal norm. That is, for any $\hat{\mathbf{w}}$ solution of the problem:

$$\mathbf{X} \in \mathbb{R}^{n \times d}, \quad \mathbf{y} \in \mathbb{R}^n.$$

Theorem

For any $\mathbf{y} \in \mathbb{R}^n$, $\mathbf{X}^\dagger \mathbf{y}$ is the solution of the least-squares problem

$$\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) = \frac{1}{2} \|\mathbf{X} \mathbf{w} - \mathbf{y}\|^2.$$

with minimal norm. That is, for any $\hat{\mathbf{w}}$ solution of the problem:

- $f(\mathbf{X}^\dagger \mathbf{y}) = f(\hat{\mathbf{w}})$;
- $\|\mathbf{X}^\dagger \mathbf{y}\| \leq \|\hat{\mathbf{w}}\|$;
- $\mathbf{X}^\dagger \mathbf{y}$ can be represented using less information than $\hat{\mathbf{w}}$.

$$\mathbf{X}\mathbf{w} = \mathbf{y}, \quad \mathbf{X} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}.$$

Finding a solution

- The least-squares problem $\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$ has infinitely many solutions;
- Among them, $\mathbf{w}^* = \mathbf{X}^\dagger \mathbf{y}$ is the one with minimal norm;
- This solution turns out to be the mean $\mathbf{w}^* = \frac{1}{n} \sum_{i=1}^n y_i!$

Key points

- Computing the pseudo-inverse;
- Or an approximation thereof!

What linear algebra solvers can do

- Put the data matrix \mathbf{X} in a nicer form (QR, LU, SVD, etc) easier to (pseudo-)invert;
- Use iterative linear algebra routines (LSQR, LSLQ, etc) to compute an approximate solution, paying attention to round-off errors;
- **Run in parallel/distributed environments.**

Two applications among many more

Least-squares formulations

- Naturally arise in a plurality of fields that try to minimize the error between a **model** and some **data**
Ex) weather forecasting, statistics, economy.
- Some problems can also be formulate or **reformulated** as linear least squares.

Two applications among many more

Least-squares formulations

- Naturally arise in a plurality of fields that try to minimize the error between a **model** and some **data**
Ex) weather forecasting, statistics, economy.
- Some problems can also be formulate or **reformulated** as linear least squares.

Two illustrations

- 1 Rewrite an optimization problem as a linear least-squares problem;
- 2 Form a linear least-squares formulation of a problem.

Illustration: Minimization of quadratic functions

Problem

Given a symmetric matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$ and a vector $\mathbf{b} \in \mathbb{R}^d$, solve

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \mathbf{w}^T \mathbf{A} \mathbf{w} - \mathbf{b}^T \mathbf{w}.$$

Solving this problem

- No solution if \mathbf{A} has negative eigenvalues! (In that case, the problem is unbounded);
- We will always assume that $\mathbf{A} \succeq 0$.

Illustration: Minimization of quadratic functions

Problem

Given a symmetric matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$ and a vector $\mathbf{b} \in \mathbb{R}^d$, solve

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \mathbf{w}^T \mathbf{A} \mathbf{w} - \mathbf{b}^T \mathbf{w}.$$

- Unconstrained problem;
- Smooth: objective is a degree-2 polynomial in the components of \mathbf{w} .

Solving this problem

- No solution if \mathbf{A} has negative eigenvalues! (In that case, the problem is unbounded);
- We will always assume that $\mathbf{A} \succeq 0$.

Illustration: Minimization of quadratic functions (2)

Square root of a matrix matrix

For any **symmetric positive semidefinite matrix** \mathbf{A} , there exists a matrix \mathbf{B} such that $\mathbf{B}^2 = \mathbf{B} \times \mathbf{B} = \mathbf{A}$.

The matrix \mathbf{B} is called the square root of \mathbf{A} ; we write $\mathbf{B} = \mathbf{A}^{1/2}$.

Illustration: Minimization of quadratic functions (2)

Square root of a matrix matrix

For any **symmetric positive semidefinite matrix** \mathbf{A} , there exists a matrix \mathbf{B} such that $\mathbf{B}^2 = \mathbf{B} \times \mathbf{B} = \mathbf{A}$.

The matrix \mathbf{B} is called the square root of \mathbf{A} ; we write $\mathbf{B} = \mathbf{A}^{1/2}$.

Reformulations

The problem $\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \mathbf{w}^T \mathbf{A} \mathbf{w} - \mathbf{b}^T \mathbf{w}$ is equivalent to

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \left\| \mathbf{A}^{1/2} \mathbf{w} - \mathbf{c} \right\|^2,$$

where $\mathbf{b} = \mathbf{A}^{1/2} \mathbf{c}$ (\mathbf{c} may not be unique).

Illustration: Minimization of quadratic functions (3)

A useful example

$$\min_{\mathbf{z} \in \mathbb{R}^d} \varphi(\mathbf{z}) = \mathbf{g}^T \mathbf{z} + \frac{m}{2} \|\mathbf{z} - \mathbf{w}\|^2 \quad \text{where } \mathbf{g} \in \mathbb{R}^d, m \geq 0.$$

Illustration: Minimization of quadratic functions (3)

A useful example

$$\min_{\mathbf{z} \in \mathbb{R}^d} \varphi(\mathbf{z}) = \mathbf{g}^T \mathbf{z} + \frac{m}{2} \|\mathbf{z} - \mathbf{w}\|^2 \quad \text{where } \mathbf{g} \in \mathbb{R}^d, m \geq 0.$$

Reformulation

- 1 Expand the least-squares formula:

$$\varphi(\mathbf{z}) = \frac{1}{2} \mathbf{z}^T \underbrace{(\mathbf{m} \mathbf{I}_d)}_{\text{Identity matrix}} \mathbf{z} + (\mathbf{g} - m\mathbf{w})^T \mathbf{z} + \underbrace{\frac{m}{2} \|\mathbf{w}\|^2}_{\text{Independent of } \mathbf{z}}.$$

Illustration: Minimization of quadratic functions (3)

A useful example

$$\min_{\mathbf{z} \in \mathbb{R}^d} \varphi(\mathbf{z}) = \mathbf{g}^T \mathbf{z} + \frac{m}{2} \|\mathbf{z} - \mathbf{w}\|^2 \quad \text{where } \mathbf{g} \in \mathbb{R}^d, m \geq 0.$$

Reformulation

- 1 Expand the least-squares formula:

$$\varphi(\mathbf{z}) = \frac{1}{2} \mathbf{z}^T \underbrace{(\mathbf{m} \mathbf{I}_d)}_{\text{Identity matrix}} \mathbf{z} + (\mathbf{g} - m\mathbf{w})^T \mathbf{z} + \underbrace{\frac{m}{2} \|\mathbf{w}\|^2}_{\text{Independent of } \mathbf{z}}.$$

- 2 Equivalent linear least-squares problem:

$$\min_{\mathbf{z} \in \mathbb{R}^d} \frac{1}{2} \left\| \mathbf{z} - \left(\mathbf{w} - \frac{1}{m} \mathbf{g} \right) \right\|^2.$$

Illustration: Minimization of quadratic functions (3)

A useful example

$$\min_{\mathbf{z} \in \mathbb{R}^d} \varphi(\mathbf{z}) = \mathbf{g}^T \mathbf{z} + \frac{m}{2} \|\mathbf{z} - \mathbf{w}\|^2 \quad \text{where } \mathbf{g} \in \mathbb{R}^d, m \geq 0.$$

Reformulation

- 1 Expand the least-squares formula:

$$\varphi(\mathbf{z}) = \frac{1}{2} \mathbf{z}^T \underbrace{(\mathbf{m} \mathbf{I}_d)}_{\text{Identity matrix}} \mathbf{z} + (\mathbf{g} - m\mathbf{w})^T \mathbf{z} + \underbrace{\frac{m}{2} \|\mathbf{w}\|^2}_{\text{Independent of } \mathbf{z}}.$$

- 2 Equivalent linear least-squares problem:

$$\min_{\mathbf{z} \in \mathbb{R}^d} \frac{1}{2} \left\| \mathbf{z} - \left(\mathbf{w} - \frac{1}{m} \mathbf{g} \right) \right\|^2.$$

- 3 The global minimum of the problem is $\mathbf{z}^* = \mathbf{w} - \frac{1}{m} \mathbf{g}$.

Illustration: Linear regression

- Data $\{(\mathbf{x}_i, y_i)\}_i$, $\mathbf{x}_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}$.
- Goal: compute a linear model $h(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ such that $h(\mathbf{x}_i) \approx y_i$ for $i = 1, \dots, n$.
- Objective: Minimize the squares of the errors $|h(\mathbf{x}_i) - y_i|$.

Illustration: Linear regression

- Data $\{(\mathbf{x}_i, y_i)\}_i$, $\mathbf{x}_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}$.
- Goal: compute a linear model $h(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ such that $h(\mathbf{x}_i) \approx y_i$ for $i = 1, \dots, n$.
- Objective: Minimize the squares of the errors $|h(\mathbf{x}_i) - y_i|$.

Linear regression

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2n} \sum_{i=1}^n (h(\mathbf{x}_i) - y_i)^2 = \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2.$$

- This is a linear least-squares problem!

Illustration: Linear regression (2)

Quality of the least-squares solution

- Best possible approximation in terms of errors;
- Fixed solution when $\{(\mathbf{x}_i, y_i)\}_i$ are deterministic.

Illustration: Linear regression (2)

Quality of the least-squares solution

- Best possible approximation in terms of errors;
- Fixed solution when $\{(\mathbf{x}_i, y_i)\}_i$ are deterministic.

In presence of random data

- True linear regression: consider the distribution of the data;
- Statistical interpretation of the least-squares solution: **maximum likelihood estimator**.

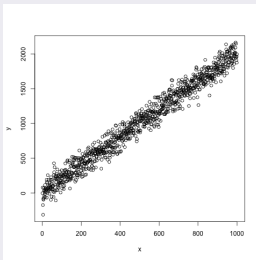
Illustration: Linear regression (2)

Quality of the least-squares solution

- Best possible approximation in terms of errors;
- Fixed solution when $\{(\mathbf{x}_i, y_i)\}_i$ are deterministic.

In presence of random data

- True linear regression: consider the distribution of the data;
- Statistical interpretation of the least-squares solution: **maximum likelihood estimator**.



In short: linear least squares

Aims

- Find a linear relationship between features and labels in your data;
- Work even when an exact linear model does not exist!

Techniques

- Look for solutions of the associated linear system;
- In practice, can use direct linear algebra solvers (even better when you choose one matching your problem characteristics);
- If too costly, can think of iterative methods.

- 1 Introduction
- 2 Basics of optimization
- 3 Unconstrained optimization
 - Linear least squares
 - Gradient descent method

$$\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}).$$

Assumptions: f smooth, bounded below.

Back to the general problem

$$\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}).$$

Assumptions: f smooth, bounded below.

Key properties

- Smoothness: We will exploit the gradient of f ;
- Convexity: Will allow for **fast convergence** (with the right method).

Example: Logistic regression

Context

- Data set $\{(\mathbf{x}_i, y_i)\}_i$, $\mathbf{x}_i \in \mathbb{R}^d$, $y_i \in \{-1, +1\}$;
- Goal: Classification through a linear classifier $\mathbf{x} \mapsto \mathbf{w}^T \mathbf{x}$;
- **Difference with SVM:** we want probabilities of belonging to a class!

Example: Logistic regression

Context

- Data set $\{(\mathbf{x}_i, y_i)\}_i$, $\mathbf{x}_i \in \mathbb{R}^d$, $y_i \in \{-1, +1\}$;
- Goal: Classification through a linear classifier $\mathbf{x} \mapsto \mathbf{w}^T \mathbf{x}$;
- **Difference with SVM:** we want probabilities of belonging to a class!

A probabilistic measure

- We define an odds-like function

$$p(\mathbf{x}; \mathbf{w}) = (1 + e^{\mathbf{x}^T \mathbf{w}})^{-1} \in (0, 1).$$

- The parameters \mathbf{w} should be chosen such that

$$\begin{cases} p(\mathbf{x}_i; \mathbf{w}) \approx 1 & \text{if } y_i = +1; \\ p(\mathbf{x}_i; \mathbf{w}) \approx 0 & \text{if } y_i = -1. \end{cases}$$

Example: Logistic regression (2)

Towards an objective function

$$p(\mathbf{x}; \mathbf{w}) = (1 + e^{\mathbf{x}^T \mathbf{w}})^{-1},$$

- Penalize cases where
 - $y_i = +1$ and $p(\mathbf{x}_i; \mathbf{w})$ is small;
 - $y_i = +1$ and $p(\mathbf{x}_i; \mathbf{w})$ is close to 1;
- Use logarithm of the $p(\mathbf{x}_i; \mathbf{w})$ in the cost function:
 - Motivation: Statistical interpretation (joint distribution);
 - Mathematical interest for gradient calculations.

Example: Logistic regression (2)

Towards an objective function

$$p(\mathbf{x}; \mathbf{w}) = (1 + e^{\mathbf{x}^T \mathbf{w}})^{-1},$$

- Penalize cases where
 - $y_i = +1$ and $p(\mathbf{x}_i; \mathbf{w})$ is small;
 - $y_i = +1$ and $p(\mathbf{x}_i; \mathbf{w})$ is close to 1;
- Use logarithm of the $p(\mathbf{x}_i; \mathbf{w})$ in the cost function:
 - Motivation: Statistical interpretation (joint distribution);
 - Mathematical interest for gradient calculations.

Resulting function: logistic loss

$$f(\mathbf{w}) = \frac{1}{n} \left\{ \sum_{y_i=-1} \ln(1 + e^{-\mathbf{x}_i^T \mathbf{w}}) + \sum_{y_i=+1} \ln(1 + e^{\mathbf{x}_i^T \mathbf{w}}) \right\}.$$

Logistic loss problem

$$\min_{\mathbf{w} \in \mathbb{R}^n} \frac{1}{n} \left\{ \sum_{y_i = -1} \ln \left(1 + e^{-\mathbf{x}_i^T \mathbf{w}} \right) + \sum_{y_i = +1} \ln \left(1 + e^{\mathbf{x}_i^T \mathbf{w}} \right) \right\}$$

Logistic loss problem

$$\min_{\mathbf{w} \in \mathbb{R}^n} \frac{1}{n} \left\{ \sum_{y_i=-1} \ln \left(1 + e^{-\mathbf{x}_i^T \mathbf{w}} \right) + \sum_{y_i=+1} \ln \left(1 + e^{\mathbf{x}_i^T \mathbf{w}} \right) \right\}$$

- The logistic loss is convex (but not strongly);
- To make it convex, possible to add a regularizing term $\frac{\mu}{2} \|\mathbf{w}\|^2$
 \Rightarrow The problem becomes μ -strongly convex!

Example: Nonlinear regression

Context

- Data set $\{(\mathbf{x}_i, y_i)\}_i$, $\mathbf{x}_i \in \mathbb{R}^d$, $y_i \in \{0, 1\}$;
- Goal: Classification through a linear classifier $\mathbf{x} \mapsto \mathbf{w}^T \mathbf{x}$.

Example: Nonlinear regression

Context

- Data set $\{(\mathbf{x}_i, y_i)\}_i$, $\mathbf{x}_i \in \mathbb{R}^d$, $y_i \in \{0, 1\}$;
- Goal: Classification through a linear classifier $\mathbf{x} \mapsto \mathbf{w}^T \mathbf{x}$.

A loss function

- We use a sigmoid function: $\phi(\mathbf{x}_i; \mathbf{w}) = \left(1 + e^{-\mathbf{x}_i^T \mathbf{w}}\right)^{-1}$;
- Our goal is now to penalize the squared error $(y_i - \phi(\mathbf{x}_i; \mathbf{w}))^2$.

Example: Nonconvex loss function (2)

The optimization problem

$$\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \left(y_i - \frac{1}{1 + e^{-x_i^T \mathbf{w}}} \right)^2.$$

Example: Nonconvex loss function (2)

The optimization problem

$$\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \left(y_i - \frac{1}{1 + e^{-x_i^T \mathbf{w}}} \right)^2.$$

- Nonconvex problem;
- **Nonlinear** least-squares structure;
- Smooth: can apply **gradient descent**.