

## Foundations of Data Science, Fall 2020

### 7b. Optimisation I

Prof. Dan Olteanu

**DaST**  
Data • (Systems+Theory)

Oct 16, 2020



University of  
Zurich

<https://lms.uzh.ch/url/RepositoryEntry/16830890400>

<https://uzh.zoom.us/j/96690150974?pwd=cnZmMTduWUtCeWoxYW85Z3RMfnpTZz09>

## Solving Machine Learning Problems

Most machine learning methods can be cast as optimisation problems.

- So far: Closed-form solutions  
e.g., minimisation of least squares and ridge regression objectives
- Most interesting learning problems do not admit closed-form solutions :(

Two approaches to solving the problems beyond closed-form solutions:

1. Frame the objective of the ML problem as a mathematical problem

Use **existing blackbox solver** for such problems

When objectives can be formulated as **convex optimisation problems**

2. **Gradient-based optimisation methods**

They are **not blackbox**: optimisation hyper-parameters affect performance

1

## A Crash Course in Optimisation

Today:

- Convex optimisation

Next time:

- Recap: Gradients, Hessians
- Gradient Descent
- Stochastic Gradient Descent
- Constrained optimisation

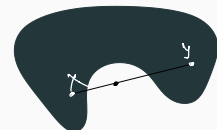
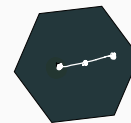
Most machine learning packages, e.g., scikit-learn, tensorflow, octave, torch, have optimisation methods readily implemented.

You need to understand the basics of optimisation to use them effectively.

2

## Convex Sets

A set  $C \subseteq \mathbb{R}^D$  is **convex** if for any  $\mathbf{x}, \mathbf{y} \in C$  and  $\lambda \in [0, 1]$ , it holds  $\lambda \mathbf{x} + (1 - \lambda) \mathbf{y} \in C$



3

## Examples of Convex Sets

- **Set  $\mathbb{R}^D$**

$\lambda \mathbf{x} + (1 - \lambda) \mathbf{y} \in \mathbb{R}^D$  for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^D$

- **Intersections of convex sets**

Given convex sets  $C_1, \dots, C_n$ , the set  $\bigcap_{i=1}^n C_i$  is convex

- **Norm balls**

For any  $L$ -norm  $\|\cdot\|$ , the set  $B = \{\mathbf{x} \in \mathbb{R}^D : \|\mathbf{x}\| \leq 1\}$  is convex

- **Polyhedra**

Given  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $\mathbf{b} \in \mathbb{R}^m$ , the polyhedron  $\{\mathbf{x} \in \mathbb{R}^n : \mathbf{A} \mathbf{x} \leq \mathbf{b}\}$  is convex

$$\mathbf{A} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \leq \begin{bmatrix} b_1 \\ \vdots \\ b_m \end{bmatrix}$$

- **Positive semidefinite cone**: The set of positive semi-definite matrices

4

## Showing the Set of PSD Matrices is Convex

$\mathbf{A} \in \mathbb{R}^{D \times D}$  is PSD :  $\forall \mathbf{x} \in \mathbb{R}^D : \mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0$ .

To show: Any combination  $\lambda \mathbf{A} + (1 - \lambda) \mathbf{B} \in \mathbb{S}_+^D$  if  $\mathbf{A}, \mathbf{B} \in \mathbb{S}_+^D$   
PSD cone.

$$\begin{aligned} \mathbf{x}^T (\lambda \mathbf{A} + (1 - \lambda) \mathbf{B}) \mathbf{x} &= \underbrace{\mathbf{x}^T \lambda \mathbf{A} \mathbf{x}}_{\geq 0} + \underbrace{\mathbf{x}^T (1 - \lambda) \mathbf{B} \mathbf{x}}_{\geq 0} \\ &= \underbrace{\lambda \cdot \mathbf{x}^T \mathbf{A} \mathbf{x}}_{\geq 0} + \underbrace{(1 - \lambda) \cdot \mathbf{x}^T \mathbf{B} \mathbf{x}}_{\geq 0} \\ &\geq 0. \end{aligned}$$

5

### Showing the Norm Balls Form Convex Sets

$$B = \{x \in \mathbb{R}^d \mid \|x\| \leq 1\}$$

Take  $x, y \in B$  and  $\lambda \in [0, 1]$ .

To show:  $\lambda x + (1-\lambda)y \in B$ . — triangle inequality.

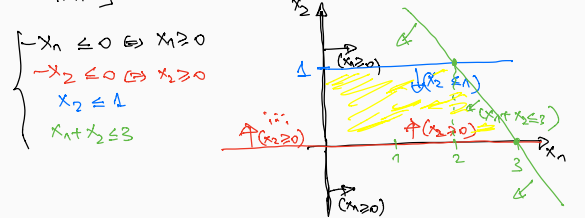
$$\begin{aligned} \|\lambda x + (1-\lambda)y\| &\leq \lambda \|x\| + (1-\lambda)\|y\| \\ &\leq \lambda \cdot 1 + (1-\lambda) \cdot 1 \\ &\leq 1. \end{aligned}$$

6

### Showing the Polyhedron is Convex + Example

Given  $A \in \mathbb{R}^{m \times n}$  and  $b \in \mathbb{R}^m$ , the polyhedron  $P = \{x \in \mathbb{R}^n : Ax \leq b\}$  is convex

$$\begin{aligned} \text{Take } x, y \in P : A(\lambda x + (1-\lambda)y) &= \lambda Ax + (1-\lambda)Ay \\ &\leq \lambda b + (1-\lambda)b \\ &\leq b. \end{aligned}$$



7

### Convex Functions

A function  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  defined on a convex domain is **convex** if:

for all  $x, y \in \mathbb{R}^n$  where  $f$  is defined and  $0 \leq \lambda \leq 1$ ,

$$f(\lambda x + (1-\lambda)y) \leq \lambda f(x) + (1-\lambda)f(y)$$

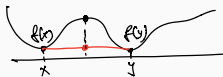


$$\lambda = 0 \Rightarrow f(y) \leq f(y)$$

$$\lambda = 1 \Rightarrow f(x) \leq f(x)$$

$$\lambda = 0.5 \Rightarrow f\left(\frac{x+y}{2}\right) \leq \frac{f(x)+f(y)}{2}$$

Not convex:



8

### Examples of Convex Functions

• **Affine functions:**  $f(x) = b^T x + c$

• **Quadratic functions:**  $f(x) = 1/2 x^T A x + b^T x + c$ , where  $A$  is symmetric positive semidefinite

• **Nonnegative weighted sums of convex functions:** Given convex functions  $f_1, \dots, f_n$  and  $w_1, \dots, w_n \in \mathbb{R}_{\geq 0}$ , the following is a convex function

$$f(x) = \sum_{i=1}^n w_i f_i(x)$$

• **Norms:**  $\|\cdot\|_p$  except  $p=0$  (0 is the number of non-zero entries in the vector).

Counterexample for convexity:

$x = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, y = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$ , Take  $\theta \in (0, 1)$ :

$$\begin{aligned} \theta \|x\|_0 + (1-\theta)\|y\|_0 &= \theta \cdot 1 + (1-\theta) \cdot 1 = 1 \\ \|\theta x + (1-\theta)y\|_0 &= \left\| \begin{bmatrix} \theta \\ 1-\theta \end{bmatrix} \right\|_0 = 2 \end{aligned}$$

9

### Convex Optimisation

Given convex functions  $f(x), g_1(x), \dots, g_m(x)$  and affine functions  $h_1(x), \dots, h_n(x)$ ,

a **convex optimisation problem** is of the form:

$$\begin{aligned} &\text{minimize } f(x) \\ &\text{subject to } g_i(x) \leq 0 \quad i \in \{1, \dots, m\} \\ &\quad \quad \quad h_j(x) = 0 \quad j \in \{1, \dots, n\} \end{aligned}$$

Goal is to find an **optimal value** of a convex optimisation problem:

$$v^* = \min\{f(x) : g_i(x) \leq 0, i \in \{1, \dots, m\}, h_j(x) = 0, j \in \{1, \dots, n\}\}$$

Whenever  $f(x^*) = v^*$  then  $x^*$  is a (not necessarily unique) **optimal point**

- $v^* \stackrel{\text{def}}{=} +\infty$  for **infeasible instances** (feasible = fulfils all constraints  $g_i$  and  $h_j$ )
- $v^* \stackrel{\text{def}}{=} -\infty$  for **unbounded instances** (unbounded = the set of feasible instances has no infimum)

10

### Local Optima are Global Optima for Convex Optimisation Problems

$x$  is **locally optimal** if:

- $x$  is feasible and
- There is  $B > 0$  s.t.  $f(x) \leq f(y)$  for all feasible  $y$  with  $\|x - y\|_2 \leq B$ .

$x$  is **globally optimal** if:

- $x$  is feasible and
- $f(x) \leq f(y)$  for all feasible  $y$ .

**Theorem:** For a convex optimisation problem, all locally optimal points are globally optimal.

*Proof by contradiction.* If locally optimal point  $x$  that is not globally optimal:  $\exists y \neq x$  s.t.  $f(y) < f(x)$

$x$  locally optimal  $\Rightarrow \exists B$  s.t.  $\forall$  feasible  $z$   $f(x) \leq f(z)$  s.t.  $\|x - z\|_2 \leq B$ .

11

## Local Optima are Global Optima for Convex Optimisation Problems: Proof

Let  $z = \lambda y + (1-\lambda)x$  where  $\lambda = \frac{\beta}{2\|x-y\|_2}$ .  
 To get  $\lambda \in [0, 1]$ , we may choose  $\beta$  as we like.  
 With this  $\lambda$ , it holds that  $\|x-z\|_2 \leq \beta$  since:  
 $\|x-z\|_2 = \|x - (\lambda y + (1-\lambda)x)\|_2 = \|(1-\lambda)x - \lambda y\|_2$   
 $= \lambda \|x-y\|_2 = \frac{\beta}{2\|x-y\|_2} \cdot \|x-y\|_2 = \frac{\beta}{2} \leq \beta$ .  
 Let us expand  $f(z)$ .  $f$  is convex  
 $f(z) = f(\lambda y + (1-\lambda)x) \leq \lambda f(y) + (1-\lambda)f(x)$   
 hypothesis:  
 $f(y) \leq \lambda f(x) + (1-\lambda)f(y)$   
 $= f(x)$ .  
 We obtained:  $f(y) \leq f(x)$ . Contradiction (it means that  $x$  is not locally optimal).

12

## Classes of Convex Optimisation Problems

### Linear Programming:

$$\begin{aligned} &\text{minimize } \mathbf{c}^T \mathbf{x} + d \\ &\text{subject to } \mathbf{A} \mathbf{x} \leq \mathbf{e} \\ &\quad \mathbf{B} \mathbf{x} = \mathbf{f} \end{aligned}$$

### Quadratically Constrained Quadratic Programming:

$$\begin{aligned} &\text{minimize } \frac{1}{2} \mathbf{x}^T \mathbf{B} \mathbf{x} + \mathbf{c}^T \mathbf{x} + d \\ &\text{subject to } \frac{1}{2} \mathbf{x}^T \mathbf{Q}_i \mathbf{x} + \mathbf{r}_i^T \mathbf{x} + s_i \leq 0 \quad i \in \{1, \dots, m\} \\ &\quad \mathbf{A} \mathbf{x} = \mathbf{b} \end{aligned}$$

### Semidefinite Programming:

$$\begin{aligned} &\text{minimize } \text{tr}(\mathbf{C} \mathbf{X}) \\ &\text{subject to } \text{tr}(\mathbf{A}_i \mathbf{X}) = b_i \quad i \in \{1, \dots, m\} \\ &\quad \mathbf{X} \text{ positive semidefinite} \end{aligned}$$

$\text{tr}(\mathbf{A})$  is the **trace** of the matrix  $\mathbf{A}$

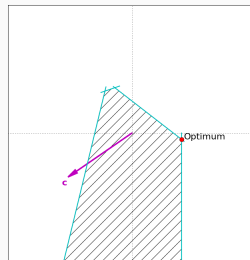
13

## Linear Programming

Looking for solutions  $\mathbf{x} \in \mathbb{R}^n$  to the following optimisation problem

$$\begin{aligned} &\text{minimize } \mathbf{c}^T \mathbf{x} + d \\ &\text{subject to } \mathbf{A} \mathbf{x} \leq \mathbf{e} \\ &\quad \mathbf{B} \mathbf{x} = \mathbf{f} \end{aligned}$$

- No closed-form solution
- Efficient algorithms exist, both in theory and practice (for tens of thousands of variables)



14

## Linear Model with Absolute Loss

Suppose we have data  $(\mathbf{X}, \mathbf{y})$  and that we want to minimise the objective:

$$\mathcal{L}(\mathbf{w}) = \sum_{i=1}^N |\mathbf{x}_i^T \mathbf{w} - y_i|$$

We would like to transform this optimisation problem into a linear program.

We introduce one  $\zeta_i$  for each datapoint.

The linear program in the  $D + N$  variables  $w_1, \dots, w_D, \zeta_1, \dots, \zeta_N$

$$\begin{aligned} &\text{minimize } \sum_{i=1}^N \zeta_i \\ &\text{subject to:} \\ &\quad \mathbf{w}^T \mathbf{x}_i - y_i \leq \zeta_i, \quad i = 1, \dots, N \\ &\quad y_i - \mathbf{w}^T \mathbf{x}_i \leq \zeta_i, \quad i = 1, \dots, N \end{aligned}$$

The solution to this linear program gives  $\mathbf{w}$  that minimises the objective  $\mathcal{L}$ .

15

## Recall: Likelihood of Linear Regression (Gaussian Noise Model)

### Likelihood

$$p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \sigma) = \left( \frac{1}{2\pi\sigma^2} \right)^{N/2} \exp \left( -\frac{1}{2\sigma^2} (\mathbf{Xw} - \mathbf{y})^T (\mathbf{Xw} - \mathbf{y}) \right)$$

Maximise Likelihood = Maximise Log-Likelihood ( $\log: \mathbb{R}^+ \rightarrow \mathbb{R}$  is increasing)

$$\text{LL}(\mathbf{y} | \mathbf{X}, \mathbf{w}, \sigma) = -\frac{N}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (\mathbf{Xw} - \mathbf{y})^T (\mathbf{Xw} - \mathbf{y})$$

Maximise Log-Likelihood = Minimise Negative Log-Likelihood

$$\begin{aligned} \text{NLL}(\mathbf{y} | \mathbf{X}, \mathbf{w}, \sigma) &= \frac{N}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2} (\mathbf{Xw} - \mathbf{y})^T (\mathbf{Xw} - \mathbf{y}) \\ &= \underbrace{\frac{N}{2} \log(2\pi\sigma^2)}_{\text{constant}} + \frac{1}{2\sigma^2} \left( \underbrace{\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w}}_{\mathbf{w}^T \mathbf{B} \mathbf{w}} - \underbrace{2\mathbf{y}^T \mathbf{X} \mathbf{w}}_{\mathbf{c}^T \mathbf{w}} + \underbrace{\mathbf{y}^T \mathbf{y}}_{\text{constant}} \right) \end{aligned}$$

This is a **convex quadratic optimisation problem with no constraints!**

18

## Minimising the Lasso Objective

For the Lasso objective, i.e., linear model with  $\ell_1$ -regularisation, we have

$$\mathcal{L}_{\text{lasso}}(\mathbf{w}) = \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}_i - y_i)^2 + \lambda \sum_{i=1}^D |w_i| = \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{y}^T \mathbf{X} \mathbf{w} + \mathbf{y}^T \mathbf{y} + \lambda \sum_{i=1}^D |w_i|$$

- Quadratic part of the loss function cannot be framed as linear programming
- Lasso regularisation does not allow for closed-form solutions
- Can be rephrased as quadratic programming problem
- Alternatively resort to general optimisation methods

19