# Session 10: Linear regression

Optimization and Computational Linear Algebra for Data Science
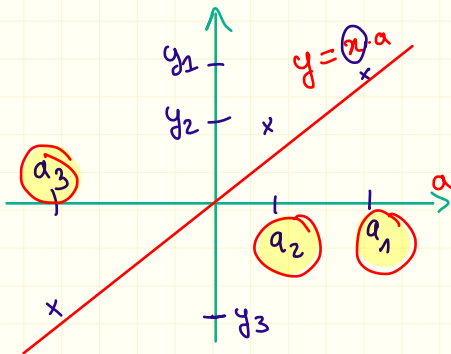
Léo Miolane

# Contents

- We have $n$ « feature vectors » $a_1, \ldots, a_n \in \mathbb{R}^d$. ← $d$ features
- Each point $a_i$ comes with a « target variable » $y_i \in \mathbb{R}$.

<u>Goal</u>: find a linear relation between the $a_i$'s and the $y_i$'s

→ find $x \in \mathbb{R}^d$ such that $\boxed{y_i \simeq \langle x, a_i \rangle}$ for all $i$.



Can we have some intercept, that is $y_i \simeq \langle x, a_i \rangle + \underline{\underline{c}}$ ?

<u>Yes</u> we can add a '1' coordinate to the $a_i$ → $\tilde{a}_i = \begin{pmatrix} a_i \\ 1 \end{pmatrix}$

$d$ first coords

$\langle x, \tilde{a}_i \rangle = \langle x, a_i \rangle + \boxed{x_{d+1}} - c$

# Solving $Ax = y$ is a bad idea

The system $Ax = y$ may have:

$$A = \begin{pmatrix} - a_1^T - \\ - a_n^T - \end{pmatrix} \in \mathbb{R}^{n \times d}$$

- No solution.

Example: if $A$ is a "tall matrix" ($n > d$)

→ $\dim \operatorname{Im}(A) \leq d < n$   $\operatorname{Im}(A) \subset \mathbb{R}^n$



→ $y$ is not very likely to belong to $\operatorname{Im}(A)$

in practice.

- Infinitely many solutions.   → no solution

Example: if $A$ is a "fat matrix" ← $d > n$



then   $\dim \ker(A) \geq d - n > 0$

$y \in \mathbb{R}^n$

→ infinitely many solutions.

# Ordinary least squares

# Least squares problem

**(LS)**   Minimize   $f(x) = \|Ax - y\|^2$   with respect to   $x \in \mathbb{R}^d$.

$f$ is convex ( HW 9 ) therefore

$x$ minimizes $f$ $\iff$ $\nabla f(x) = 0$.

$\iff 2 A^\top A x - 2 A^\top y = 0$

$\iff A^\top A x = A^\top y$.

<u>Conclusion</u> : the minimizers of $f$ are exactly the solutions of the linear system $A^\top A x = A^\top y$

If $A^\top A$ is invertible $\to$ $x = (A^\top A)^{-1} A^\top y$.

What if $A^\top A$ is not invertible?

# The Moore-Penrose pseudo-inverse

### Definition

Let $A = U\Sigma V^\mathsf{T}$ be the SVD of $A$. The matrix $A^\dagger \stackrel{\text{def}}{=} V\Sigma' U^\mathsf{T}$ is called the (Moore-Penrose) pseudo-inverse of $A$, where $\Sigma'$ is the $d \times n$ matrix given by

$$\Sigma'_{i,i} = \begin{cases} 1/\Sigma_{i,i} & \text{if } \Sigma_{i,i} \neq 0, \\ 0 & \text{otherwise,} \end{cases}$$

and $\Sigma'_{i,j} = 0$ for $i \neq j$.

$$A = U \begin{pmatrix} \sigma_1 & & & \\ & \ddots & & \\ & & \sigma_r & 0 \\ & & & \ddots \\ & & & & 0 \end{pmatrix} U^\mathsf{T} \qquad \in \mathbb{R}^{n \times d}$$

Exercise: check that if
$A$ invertible, then $A^{-1} = A^+$

$$A^+ = V \begin{pmatrix} 1/\sigma_1 & & & \\ & \ddots & & \\ & & 1/\sigma_r & 0 \\ & & & \ddots & 0 \end{pmatrix} U^\mathsf{T} \qquad \in \mathbb{R}^{d \times n}$$

# Solving $A^\mathsf{T} A x = A^\mathsf{T} y$

**Claim:** The vector $x^{\mathrm{LS}} \overset{\mathrm{def}}{=} A^\dagger y$ is a solution of $A^\mathsf{T} A x = A^\mathsf{T} y$

$$A^\mathsf{T} A\, x^{\mathrm{LS}} = V \Sigma^\mathsf{T} \cancel{U^\mathsf{T}} \cancel{U} \Sigma \cancel{V^\mathsf{T}} \cancel{V} \Sigma^\dagger U^\mathsf{T} y$$

$$= V \underbrace{\Sigma^\mathsf{T} \Sigma \Sigma^\dagger}_{=\ \Sigma^\mathsf{T}} U^\mathsf{T} y = V \Sigma^\mathsf{T} U^\mathsf{T} y = A^\mathsf{T} y.$$

## Theorem

The set of the minimizers of $f(x) = \|Ax - y\|^2$ is

$$\underbrace{A^\dagger y}_{\mathrm{Ker}(A^\mathsf{T} A)} + \mathrm{Ker}(A) = \left\{ x^{\mathrm{LS}} + v \,\middle|\, v \in \mathrm{Ker}(A) \right\}.$$

# Penalized least squares

# Ridge regression

Ridge regression consists in adding a « $\ell_2$ penalty » :

**(Ridge)**   Minimize   $f(x) = \frac{1}{2}\|Ax - y\|^2 + \frac{\lambda}{2}\|x\|^2$   w.r.t.   $x \in \mathbb{R}^d$

for some fixed $\lambda > 0$.

$$x_1^{\boxed{2}} + \cdots + x_d^2$$

- $f$ is <mark>strongly convex</mark>, it admits a unique minimizer 

$$\boxed{x^{Ridge} = \left(A^T A + \lambda \, \mathrm{Id}\right)^{-1} A^T y.}$$

- Why adding the $\ell_2$-penalty ?

Trade-off : 
- this promotes vectors of small norm $\|x^{Ridge}\| \leq \|x^{LS}\|$ (exercise!)

- issue : $\|Ax^{Ridge} - y\|^2 \geqslant \|Ax^{LS} - y\|^2$

$$\langle x, a_{new} \rangle = x_1 \cdot a_{new\,1} + \cdots + x_d \cdot a_{new\,d}$$

# Lasso

The Lasso adds a « $\ell_1$ penalty » :

**(Lasso)**   Minimize   $f(x) = \|Ax - y\|^2 + \lambda \|x\|_1$,   w.r.t.   $x \in \mathbb{R}^d$

for some fixed $\lambda > 0$.

$g(t) = \dfrac{t^2}{2}$   $g'(t) = t$

- $f$ is **not** strictly convex in general : there is not a unique minimizer a priori.

- In practice, the minimizer $x^{Lasso}$ is unique.

Why do we add this $\ell_1$-penalty ?

$\longrightarrow$ it promotes sparse vectors $\underline{x^{Lasso}}$ (lots of coefficients of $x^{Lasso}$ are likely to be zero).

$\longrightarrow$ Feature selection !

# Intuition behind feature selection

### Lemma

Let $x^{\mathrm{Lasso}}$ be a minimizer of the Lasso cost function and let $r = \|x^{\mathrm{Lasso}}\|_1$. Then $x^{\mathrm{Lasso}}$ is a solution to the constrained optimization problem:

$$\text{minimize} \quad \|Ax - y\|^2 \quad \text{subject to} \quad \|x\|_1 \leq r.$$

Proof: By contradiction, assume that there exists $x$ such that

$$\|Ax - y\|^2 < \|Ax^{\mathrm{Lasso}} - y\|^2$$

$$\|x\|_1 \leq r = \|x^{\mathrm{Lasso}}\|_1$$

$$\rightarrow \quad \|Ax - y\|^2 + \lambda \|x\|_1 \; < \; \|Ax^{\mathrm{Lasso}} - y\|^2 + \lambda \|x^{\mathrm{Lasso}}\|_1$$

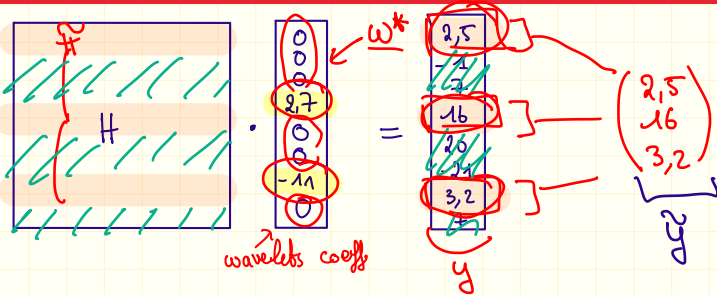$$\rightarrow \text{Contradiction}$$

# Application: compressed sensing

- In homework 4 we have seen that we can compress images very well.
- Most of the data can be thrown away !

$$H \cdot \begin{pmatrix} 0 \\ 0 \\ 0 \\ 2,7 \\ 0 \\ 0 \\ -11 \\ 0 \end{pmatrix} = \begin{pmatrix} 2,5 \\ -1 \\ 7 \\ 16 \\ 20 \\ -21 \\ 3,2 \\ 7 \end{pmatrix} \leftarrow \text{pixels of the image.}$$

wavelets coeff

Can we directly measure only the useful wavelet coefficients ?

① Measure only a small fraction of the pixels.

② We have $\tilde{H} \, \omega^* = \tilde{y}$

③ Minimize $f(\omega) = \| \tilde{H} \omega - \tilde{y} \|^2 + \lambda \| \omega \|_1$

to get $\omega^{Lasso}$ which should be a good estimate

of $\omega^*$

# Matrix norms

# Frobenius norm

## Definition

The Frobenius norm of a matrix $A \in \mathbb{R}^{n \times m}$ is defined as

$$\|A\|_F = \sqrt{\sum_{i=1}^{n} \sum_{j=1}^{m} A_{i,j}^2}$$

## Proposition

$$A = U \begin{pmatrix} \sigma_1 & \\ & \ddots \end{pmatrix} V^T$$

$$\|A\|_F = \sqrt{\sum_{i=1}^{\min(n,m)} \sigma_i(A)^2}$$

$\text{Tr}(AB)$
$= \text{Tr}(BA)$

$$\|A\|_F^2 = \text{Tr}(A A^T) = \text{Tr}(U \Sigma V^T V \Sigma^T U^T)$$
$$= \text{Tr}(U^T U \Sigma \Sigma^T) = \sigma_1^2 + \cdots + \sigma_r^2$$

# The spectral norm

## Definition

The spectral norm of a matrix $A \in \mathbb{R}^{n \times m}$ is defined as

$$\|A\|_{\mathrm{Sp}} = \max_{\|x\|=1} \|Ax\|.$$

## Proposition

$$\|A\|_{\mathrm{Sp}} = \sigma_1(A).$$

largest singular value of A.

Proof:
$$\|A\|_{\mathrm{Sp}}^2 = \max_{\|x\|=1} \|Ax\|^2$$

$$A = U \Sigma V^T$$
$$A^T A = V \begin{pmatrix} \sigma_1^2 & \\ & \ddots \end{pmatrix} V^T$$

$$= \max_{\|x\|=1} x^T A^T A x$$

$$= \lambda_1(A^T A) = \sigma_1(A)^2$$

## Definition

The nuclear norm of a matrix $A \in \mathbb{R}^{n \times m}$ is defined as

$$\|A\|_\star = \sum_{i=1}^{\min(n,m)} \sigma_i(A).$$

$\longrightarrow$ "$\ell_1$ – norm of the singular values"

# Application to matrix completion

We have a data matrix $M \in \mathbb{R}^{n \times m}$ that we only observe **partially**. That is we only have access to

$$M_{i,j} \quad \text{for } (i,j) \in \Omega,$$

where $\Omega \subset \{1, \ldots, n\} \times \{1, \ldots m\}$ is a **subset of the complete set of the entries.**

NP-HARD

$\hookrightarrow$ minimize $\quad \text{rank}(X) \quad$ with respect to $X \in \mathbb{R}^{n \times m}$

verifying $\quad X_{i,j} = M_{i,j}$

for all $(i,j) \in \Omega$

it has been proposed to solve instead:

minimize $\quad \|X\|_* \quad$ with respect to $X \in \mathbb{R}^{n \times m}$

verifying $\quad X_{i,j} = M_{i,j}$

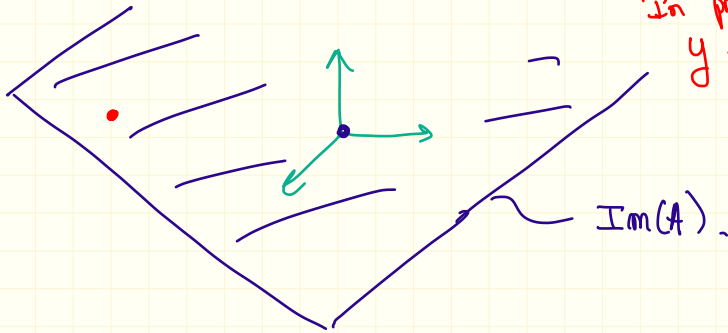for all $(i,j) \in \Omega$

# Application to matrix completion

$$A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{bmatrix}$$

2

3

$\mathbb{R}^3$

Im(A) is a subspace of $\mathbb{R}^3$ of dimension at most 2

Assume $\dim \text{Im}(A) = 2$

"In practice"

$y \notin \text{Im}(A)$



Im(A).

# Questions?