

Session 10: Linear regression

Optimization and Computational Linear Algebra for Data Science

Textbook: Boyd & Vandenberghe: Introduction to applied linear algebra - Chaps 12 & 13

Marylou Gabrié (based on material by Léo Miolane)

Contents

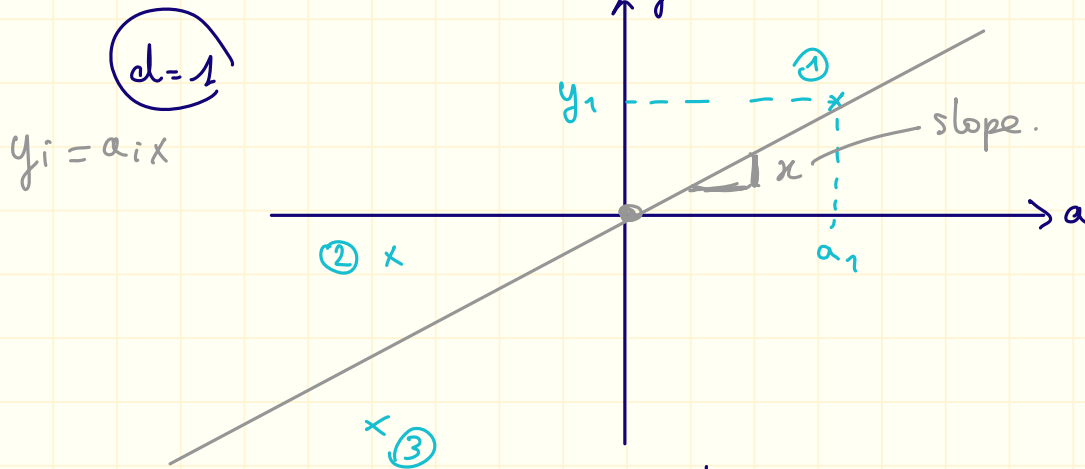
1. Ordinary least squares
2. Penalized linear regression
3. Matrix norms

Introduction

example $\left\{ \begin{array}{l} d \text{ physiological measures } a_i \\ 1 \text{ output } y_i - \text{diabetes stage} \end{array} \right. \rightarrow d \text{ features}$

- ❑ We have n « feature vectors » $a_1, \dots, a_n \in \mathbb{R}^d$.
- ❑ Each point a_i comes with a « target variable » $y_i \in \mathbb{R}$.
- ❑ **Goal.** Find a linear relation between the a_i s and the y_i s:

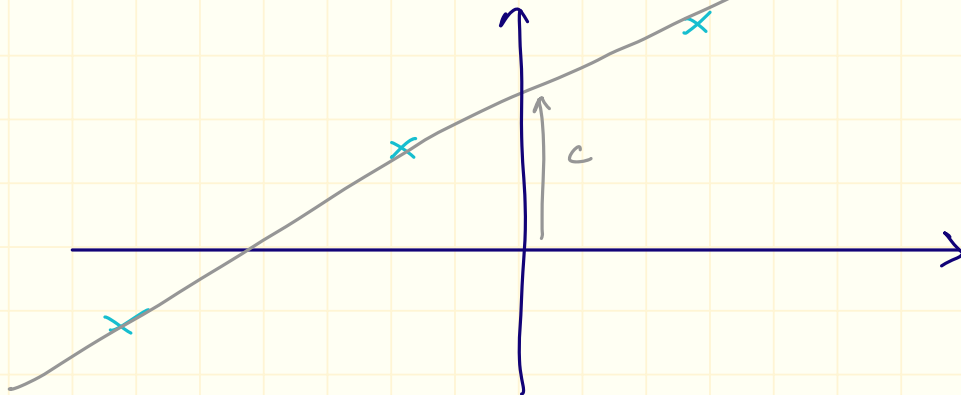
↳ Find $x \in \mathbb{R}^d$ such that $y_i \approx \langle x, a_i \rangle$ for all $i = 1 \dots n$



- ❑ **Prediction:** New $a \in \mathbb{R}^d$ $\hat{y} = \langle x, a \rangle$

Can we have an intercept?

Find $x \in \mathbb{R}^d$ such that $y_i = \langle x, a_i \rangle + c. = \langle \tilde{x}, \tilde{a}_i \rangle$



$$\tilde{a}_i \in \mathbb{R}^{d+1} \quad - \quad \tilde{a}_i = \begin{pmatrix} a_{i1} \\ \vdots \\ a_{id} \\ 1 \end{pmatrix}$$

$$\tilde{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_d \\ c \end{pmatrix}$$

\Rightarrow Write: $A = \begin{pmatrix} -a_1^T \\ \vdots \\ -a_n^T \end{pmatrix} \in \mathbb{R}^{n \times (d+1)}$ $y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \in \mathbb{R}^n$ Find x in \mathbb{R}^{d+1}
 " $Ax = y$ "

Solving $Ax = y$ is a bad idea

From now on "no intercept"

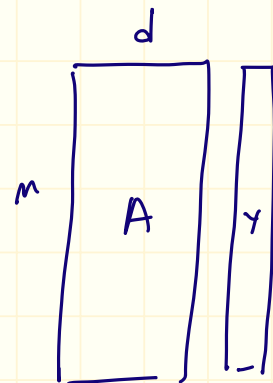
The system $Ax = y$ may have:

❑ No solution.

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

$$\Rightarrow \dim \text{Im}(A) \leq d < n$$

$\Rightarrow y \in \mathbb{R}^n$ is unlikely to belong to $\text{Im}(A) \subseteq \mathbb{R}^n$.

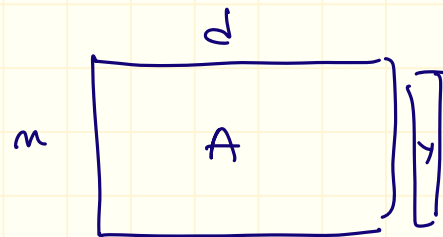


❑ Infinitely many solutions.

example: A is a "fat" matrix ($d > n$)

$$\dim(\text{Ker}(A)) \geq d - n > 0$$

\Rightarrow typically infinitely many solutions



1. Ordinary least squares

Least squares problem

Euclidean, ℓ_2 -norm

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

f is convex (HW9) therefore:

$$x \text{ minimizes } f \iff \nabla f(x) = 0$$

$$\iff 2(A^T A x - A^T y) = 0$$

$$\iff A^T A x = A^T y$$

Conclusion: the minimizers of f are the solutions $x \in \mathbb{R}^d$ of the linear system

$$A^T A x = A^T y.$$

CASE 1: $A^T A$ is invertible:

$$x = (A^T A)^{-1} A^T y.$$

The Moore-Penrose pseudo-inverse

What if $A^T A$ is not invertible?

Definition

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

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

$$A = U \begin{pmatrix} \sigma_1 & & & \\ & \ddots & & \\ & & \sigma_n & \\ & & & 0 & \dots & 0 \end{pmatrix} V^T$$
$$A^\dagger = V \begin{pmatrix} 1/\sigma_1 & & & \\ & \ddots & & \\ & & 1/\sigma_n & \\ & & & 0 & \dots & 0 \end{pmatrix} U^T$$

Exercise: Check that if A is invertible then $A^{-1} = A^\dagger$.

"dagger"



Solving $A^T A x = A^T y$

A

Claim: The vector $x^{\text{LS}} \stackrel{\text{def}}{=} A^\dagger y$ is a minimize $\|y - Ax\|^2$
solution of $A^T A x = A^T y$

$$\begin{aligned} A^T A x^{\text{LS}} &= V \cancel{\Sigma^T U^T} U \cancel{\Sigma V^T} V \Sigma' U^T y \\ &= V \underbrace{\Sigma^T \Sigma \Sigma'}_{\Sigma^T} U^T y = V \Sigma^T U^T y = A^T y. \end{aligned}$$

Theorem

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

$$\left\{ x^{\text{LS}} + v \mid v \in \underbrace{\text{Ker}(A)}_{\text{Ker}(A^T A)} \right\}.$$

2. Penalized least squares

Ridge regression

Ridge regression consists in adding a « ℓ_2 penalty » :

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

for some fixed $\lambda > 0$.

- CONVEXITY: f is strongly convex \Rightarrow strictly convex
 f admits a unique minimizer.

$$x^{\text{Ridge}} = (A^T A + \lambda I_d)^{-1} A^T y. \quad \text{+w!}$$

- WHY ℓ_2 penalty?

Tradeoff: • promotes solution x with small norm

$$\mathbb{E}_x \quad \|x^{\text{Ridge}}\| \leq \|x^{\text{LS}}\|.$$

\Rightarrow small norm solution can be robust.

- issue $\|Ax^{\text{LS}} - y\| \leq \|Ax^{\text{Ridge}} - y\|$

Lasso

$$\|x\|_1 = \sum_{i=1}^d |x_i|$$

The Lasso adds a « ℓ_1 penalty » :

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

for some fixed $\lambda > 0$.

- CONVEXITY: f is not strictly convex, no unique minimizer in general.

↘ In practice, LASSO minimizer is still unique.

- WHY ℓ_1 PENALTY:

- promotes sparse vectors: x^{Lasso} will have "a lot" of coordinates equal to 0

$$x^{\text{Lasso}} = \begin{pmatrix} x_1 \\ 0 \\ x_2 \\ 0 \\ 0 \\ \vdots \end{pmatrix}$$

- Tradeoff.

→ Features selection.

Intuition behind feature selection

Lemma

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

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

Proof: By contradiction: assume that there exist x

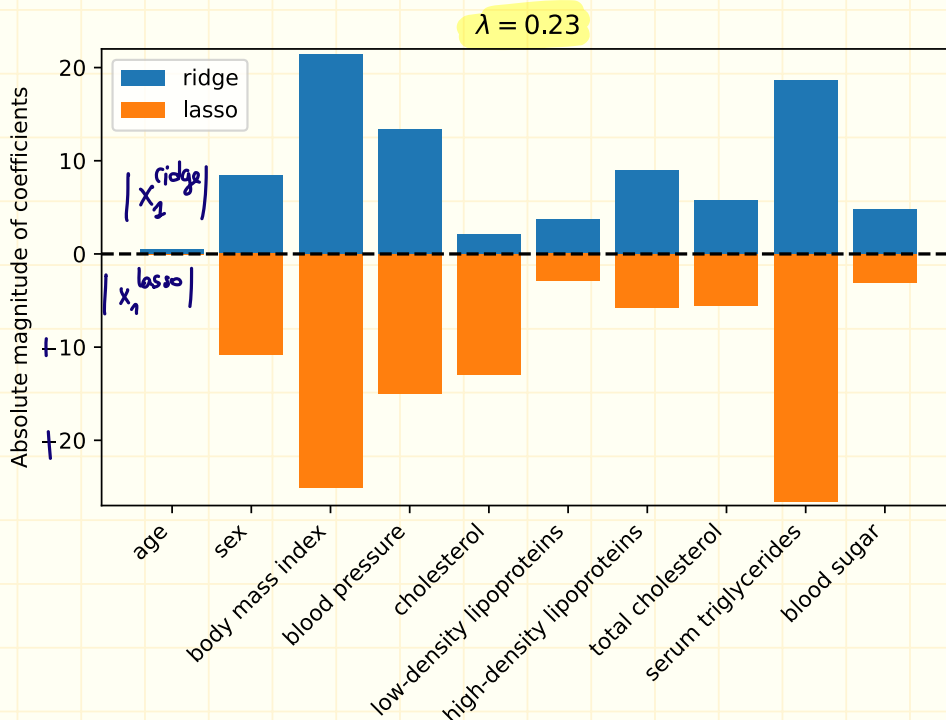
$$\text{such that } \begin{cases} \|Ax - y\|^2 < \|Ax^{\text{Lasso}} - y\|^2 \\ \|x\|_1 \leq \|x^{\text{Lasso}}\|_1 = r. \end{cases}$$

$$\Rightarrow \|Ax - y\|^2 + \lambda \|x\|_1 < \underbrace{\|Ax^{\text{Lasso}} - y\|^2 + \lambda \|x^{\text{Lasso}}\|_1}_{f(x^{\text{Lasso}})}.$$

→ Contradiction since x^{Lasso} minimizer of $f(\cdot)$

Effect of Regularization

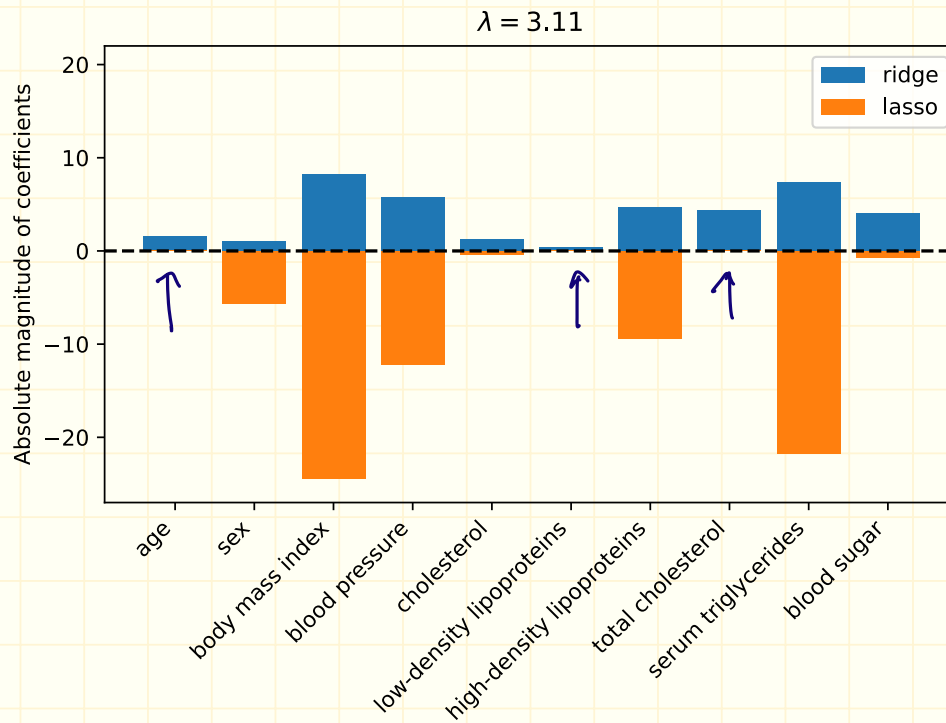
Consider a data set with $n = 442$ patients, with $d = 10$ dimensions feature vectors and the prediction of diabetes disease progression



https://scikit-learn.org/stable/datasets/toy_dataset.html

Effect of Regularization

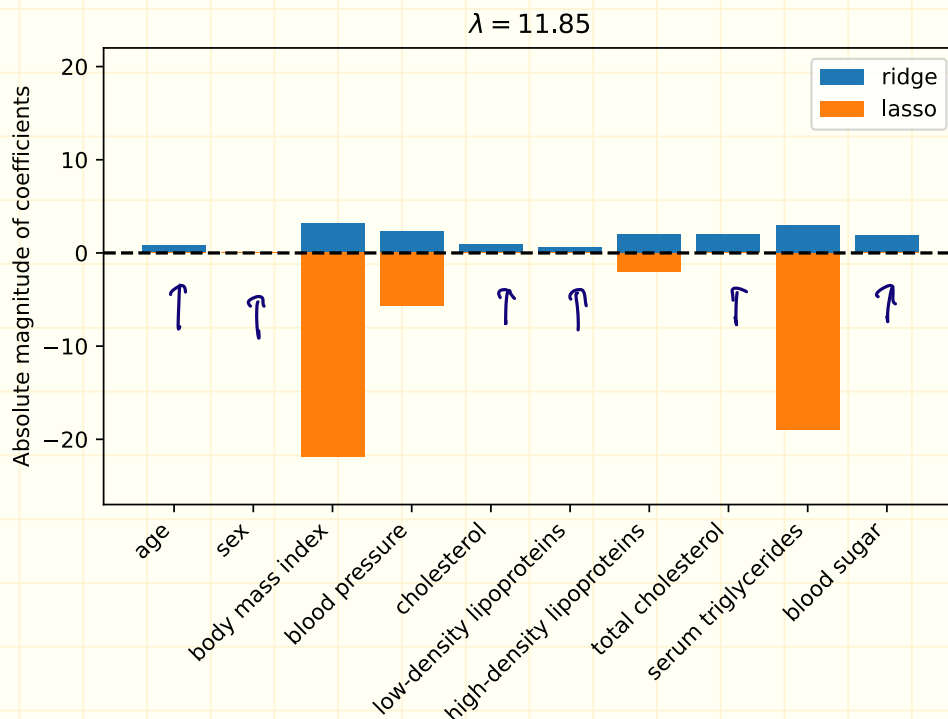
Consider a data set with $n = 442$ patients, with $d = 10$ dimensions feature vectors and the prediction of diabetes disease progression



https://scikit-learn.org/stable/datasets/toy_dataset.html

Effect of Regularization

Consider a data set with $n = 442$ patients, with $d = 10$ dimensions feature vectors and the prediction of diabetes disease progression

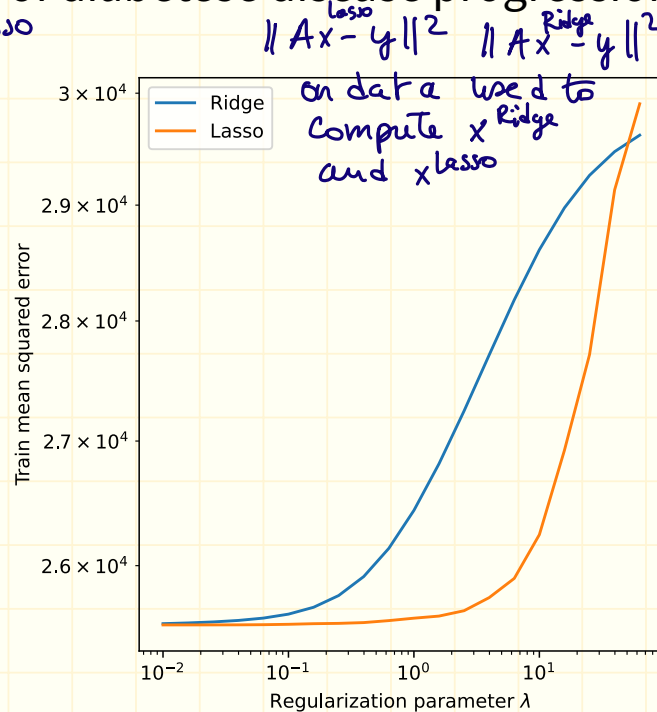
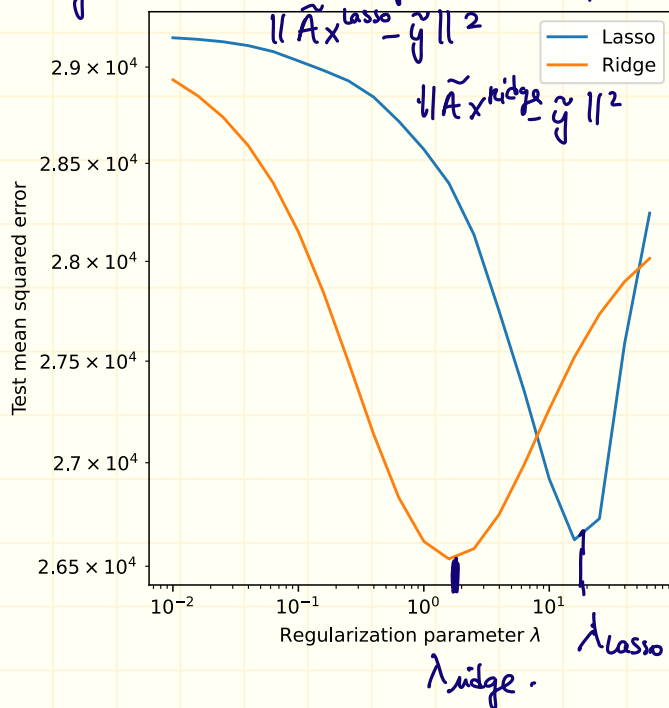


https://scikit-learn.org/stable/datasets/toy_dataset.html

Effect of Regularization

Consider a data set with $n = 442$ patients, with $d = 10$ dimensions feature vectors and the prediction of diabetes disease progression

\tilde{A} \tilde{y} not used to compute x^{Ridge} , x^{Lasso}



https://scikit-learn.org/stable/datasets/toy_dataset.html

3. Matrix norms

Frobenius norm

Definition

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

$$\|A\|_F = \sqrt{\sum_{i=1}^n \sum_{j=1}^m A_{i,j}^2} = \sqrt{\text{Tr}(A^T A)}$$

Proposition

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

$$\|A\|_F^2 = \text{Tr}(AA^T) = \text{Tr}\left(\cancel{U} \Sigma \cancel{V^T} \cancel{V} \Sigma^T \cancel{U^T}\right) = \text{Tr}(\Sigma \Sigma^T) = \sigma_1^2 + \dots + \sigma_n^2$$

The spectral norm

Definition

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

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

Proposition

largest singular values

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

Proof:

$$\begin{aligned} \|A\|_{\text{Sp}}^2 &= \max_{\|x\|=1} \|Ax\|^2 \\ &= \max_{\|x\|=1} x^T A^T A x \end{aligned}$$

max eigenvalue of $A^T A$

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

The nuclear norm

Definition

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

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

" ℓ_1 norm of matrices"

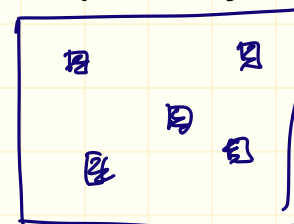
Application to matrix completion

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

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

observed entries
indices.

Ω



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

→ minimize

$$\text{rank}(X) \text{ w.r.t } X \in \mathbb{R}^{n \times m} \text{ (such that)} \\ \text{s.t. } \pi_{i,j} = X_{i,j} \\ \text{for all } (i,j) \in \Omega$$

→ NP-HARD

→ Instead
Solve

$$\text{minimize } \|X\|_* \text{ w.r.t } X \in \mathbb{R}^{n \times m} \text{ s.t. } \pi_{i,j} = X_{i,j} \\ \text{for all } (i,j) \in \Omega.$$

Application to matrix completion

Questions?

Questions?