

# Optimization and Computational Linear Algebra for Data Science

## Lecture 11: Linear regression, matrix completion

Léo MIOLANE · leo.miolane@gmail.com

October 26, 2019

**Warning:** *This material is not meant to be lecture notes. It only gathers the main concepts and results from the lecture, without any additional explanation, motivation, examples, figures...*

## 1 Least squares

Assume that we are given point  $a_i = (a_{i,1}, \dots, a_{i,d}) \in \mathbb{R}^d$  with labels  $y_i \in \mathbb{R}$  for  $i = 1 \dots n$ . We aim at finding a vector  $x \in \mathbb{R}^d$  such that

$$y_i \simeq \langle a_i, x \rangle = \sum_{j=1}^d a_{i,j} x_j, \quad \text{for } i = 1 \dots n.$$

If we denote by  $A$  the  $n \times d$  matrix whose rows are  $a_1, \dots, a_n$ , i.e.  $A_{i,j} = a_{i,j}$ , we are looking for some  $x$  such that  $Ax \simeq y$ .

### 1.1 Solving the system $Ax = y$

As we have seen in Lecture 2, we can distinguish two cases:

- If  $y \notin \text{Im}(A)$  then the equation  $Ax = y$  does not admit any solution (by definition of  $\text{Im}(A)$ ).
- If  $y \in \text{Im}(A)$  then the equation  $Ax = y$  admits at least a solution  $x_0$  (by definition of  $\text{Im}(A)$ ). Moreover, the set of (all) solutions is

$$x_0 + \text{Ker}(A) = \{x_0 + v \mid v \in \text{Ker}(A)\}.$$

In particular, if  $\text{Ker}(A) = \{0\}$  then the equation admits a unique solution.

In the second case, one can obtain an expression for a particular solution  $x_0$  using the SVD of  $A$ . Let  $r = \text{rank}(A)$ ,  $\sigma_1, \sigma_2, \dots, \sigma_r > 0$  be the non-zero singular values of  $A$  and  $\Sigma = \text{Diag}(\sigma_1, \dots, \sigma_r)$ . Finally, let  $A = U\Sigma V^T$  be the SVD of  $A$ , where  $V \in \mathbb{R}^{n \times r}$  and  $U \in \mathbb{R}^{d \times r}$  are matrices that have orthonormal columns.

Notice that  $V^T V = \text{Id}$  and that  $UU^T$  is the orthogonal projection on  $\text{Im}(A)$ . Hence, if we let  $x_0 = V\Sigma^{-1}U^T y$ , we have

$$Ax_0 = U\Sigma V^T V\Sigma^{-1}U^T y = UU^T y = y$$

because we assumed that  $y \in \text{Im}(A)$ . This motivates the following definition:

#### **Definition 1.1 (Moore-Penrose pseudo-inverse)**

The matrix  $A^\dagger \stackrel{\text{def}}{=} V\Sigma^{-1}U^T$  is called the (Moore-Penrose) pseudo-inverse of  $A$ .

Notice that in the case where  $A$  is invertible,  $A^\dagger = A^{-1}$ . From the analysis above, we deduce:

### Proposition 1.1

The set of solution of the linear system  $Ax = y$  is

- $\emptyset$  if  $y \notin \text{Im}(A)$ .
- $A^\dagger y + \text{Ker}(A)$  otherwise.

## 1.2 Least squares

In general, there is no reason for  $y$  to belong to  $\text{Im}(A)$ , especially when  $n > d$ . (Exercise: why?) Therefore one is rather interested by solving

$$\min_{x \in \mathbb{R}^d} \|Ax - y\|^2. \quad (1)$$

The function  $f : x \mapsto \|Ax - y\|^2$  is convex (Exercise: why?) and differentiable. Hence  $x$  is solution of (1) if and only if  $\nabla f(x) = 0$ . Compute

$$f(x) = (Ax - y)^\top (Ax - y) = x^\top A^\top Ax - 2y^\top Ax + \|y\|^2.$$

Hence  $\nabla f(x) = 2A^\top Ax - 2A^\top y$ . We conclude

$$x \text{ is solution of (1)} \iff A^\top Ax = A^\top y.$$

If  $A^\top A$  is invertible there is a unique minimizer  $x^* = (A^\top A)^{-1} A^\top y$ . In the general case, we see that the solutions of (1) are the solutions of the linear system  $A^\top Ax = A^\top y$ . From Proposition 1.1 we get that the solutions of (1) are

$$(A^\top A)^\dagger A^\top y + \text{Ker}(A^\top A).$$

This expression simplifies a lot. First (exercise!) we have  $\text{Ker}(A^\top A) = \text{Ker}(A)$ . Then if we let  $A = U\Sigma V^\top$  be the SVD of  $A$ , we have

$$A^\top A = V\Sigma^2 V^\top.$$

$V\Sigma^2 V^\top$  is therefore the SVD of  $A^\top A$ . Hence  $(A^\top A)^\dagger = V\Sigma^{-2} V^\top$ . This gives  $(A^\top A)^\dagger A^\top = V\Sigma^{-2} V^\top V\Sigma U^\top = A^\dagger$ . We conclude:

### Proposition 1.2 (Least squares)

The set of solution of the minimization problem  $\min_{x \in \mathbb{R}^n} \|Ax - y\|^2$  is

$$A^\dagger y + \text{Ker}(A).$$

## 2 Penalized least squares: Ridge regression and Lasso

When  $\text{Ker}(A) \neq \emptyset$  the least squares problem (1) has an infinite number of solutions: which one should we pick?

## 2.1 Ridge regression

The Ridge regression adds a  $\ell_2$  penalty to the least square problem in order to “select” a solution of smaller norm, and minimizes

$$\min_{x \in \mathbb{R}^d} \left\{ \|Ax - y\|^2 + \lambda \|x\|^2 \right\}, \quad (2)$$

for some penalization parameter  $\lambda > 0$ .

**Exercise 2.1.** Show that (2) admits a unique solution given by

$$x^{\text{Ridge}} = (A^\top A + \lambda \text{Id})^{-1} A^\top y.$$

## 2.2 Lasso

The Lasso adds a  $\ell_1$  penalty to the least square problem, and minimizes

$$\min_{x \in \mathbb{R}^d} \left\{ \frac{1}{2} \|Ax - y\|^2 + \lambda \|x\|_1 \right\}, \quad (3)$$

for some penalization parameter  $\lambda > 0$ . The Lasso has the wonderful property of *feature selection*: the solution  $x^*$  of (3) is likely to be *sparse* (many coordinates  $x_j^*$  will be set to 0). In words, the Lasso estimator discards the “useless features” by setting its corresponding coefficient to 0. This is particularly nice for the interpretability of the results: in many application, each data point  $a_i$  has an enormous number of features ( $d$  very large), but only a small number of them are useful to predict the label  $y_i$ .

We can gain some intuition on this phenomena on Figure.

**Lasso for orthonormal design** In general there is no closed form formula for the Lasso estimator. It is however possible to derive one in the case where the matrix  $A$  has orthonormal columns:  $A^\top A = \text{Id}$ . The least square estimator is then given by  $x^{\text{LS}} = A^\top y$  and the Lasso cost function becomes

$$\frac{1}{2} \|Ax - y\|^2 + \lambda \|x\|_1 = \frac{1}{2} \|x\|^2 - \langle x^{\text{LS}}, x \rangle + \frac{1}{2} \|y\|^2 + \lambda \|x\|_1. \quad (4)$$

The minimizer of (4) is therefore the minimizer of

$$\frac{1}{2} \|x\|^2 - \langle x^{\text{LS}}, x \rangle + \lambda \|x\|_1 = \sum_{j=1}^d \frac{x_j^2}{2} - x_j^{\text{LS}} x_j + \lambda |x_j| = \sum_{j=1}^d f_{x_j^{\text{LS}}}(x_j),$$

where  $f_{x_0}(x) \stackrel{\text{def}}{=} \frac{1}{2} x^2 - x_0 x + \lambda |x|$ .

**Exercise 2.2.** Show that the function  $f_{x_0}$  admits a unique minimizer given by

$$x^* = \eta(x_0; \lambda),$$

where  $\eta$  denotes the “soft-thresholding” function:

$$\eta(x_0; \lambda) = \begin{cases} x_0 - \lambda & \text{if } x_0 \geq \lambda \\ 0 & \text{if } -\lambda \leq x_0 \leq \lambda \\ x_0 + \lambda & \text{if } x_0 \leq -\lambda. \end{cases}$$

We conclude that the Lasso estimator is given by

$$x_j^{\text{Lasso}} = \eta(x_j^{\text{LS}}; \lambda) \quad \text{for } j = 1, \dots, d.$$

This formula confirms the intuition of Figure: the Lasso estimator translates the coefficients of the least-square solution by  $\lambda$ , truncating at 0.

### 3 Norms for matrices

Before looking at low-rank matrix estimation and matrix completion, we need to look at different norms we can have for matrices. Recall that  $\mathbb{R}^{n \times m}$ , the set of  $n \times m$  matrices, is a vector space (of dimension  $nm$ ).

**The Frobenius norm.** The most obvious norm to consider is the equivalent of the  $\ell_2$  norm, called the Frobenius norm:

**Definition 3.1 (Frobenius norm)**

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

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

**Remark 3.1.** The Frobenius norm is the norm induced by the following inner-product on matrices:

$$\langle A, B \rangle = \text{Tr}(A^T B) = \sum_{i=1}^n \sum_{j=1}^m A_{i,j} B_{i,j}, \quad \text{for } A, B \in \mathbb{R}^{n \times m}.$$

**Proposition 3.1**

Let  $A \in \mathbb{R}^{n \times n}$  and let  $\sigma_1, \dots, \sigma_{\min(n,m)}$  be the singular values of  $A$ . Then

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

**The spectral norm.** Another extremely common norm for matrices is the spectral norm, also called the operator norm:

**Definition 3.2 (Spectral norm)**

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

$$\|A\|_{\text{Sp}} = \max_{\substack{x \in \mathbb{R}^m \\ \|x\|=1}} \|Ax\|.$$

**The nuclear norm.** We have seen above that the Frobenius norm of a matrix correspond to the  $\ell_2$  norm of its singular values and that the spectral norm correspond to the infinity norm of the singular values. The nuclear norm is simply the  $\ell_1$  norm of its singular values:

**Definition 3.3 (Nuclear norm)**

Let  $A \in \mathbb{R}^{n \times n}$  and let  $\sigma_1, \dots, \sigma_{\min(n,m)}$  be the singular values of  $A$ . The nuclear norm of  $A$  is defined as:

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

## 4 Low-rank matrix estimation and matrix completion

Further reading



References