

# Linear regression

R. Flamary, A. Rakotomamonjy

January 10, 2019

1/35

## Linear Prediction

### Linear function

Function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ , can be expressed as

$$f(\mathbf{x}) = \sum_{i=1}^d w_i x_i + b = \mathbf{x}^\top \mathbf{w} + b = [\mathbf{x}^\top \ 1] \boldsymbol{\alpha} \quad (1)$$

with  $\mathbf{w} \in \mathbb{R}^d$  a vector defining an hyperplane in  $\mathbb{R}^d$  et  $b \in \mathbb{R}$  a bias term displacing the function along the normal  $\mathbf{w}$  of the hyperplane. All parameters can be stored in a unique vector  $\boldsymbol{\alpha} = \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix}$  of dimensionality  $\mathbb{R}^{d+1}$  concatenating  $\mathbf{w}$  and  $b$ .

### Objective of linear prediction

- ▶ Regression:  $f(\cdot) \in \mathbb{R}$ .
- ▶ Classification:  $\text{sign}(f(\cdot)) \in \{-1, 1\}$ .

3/35

## Sommaire

### Introduction

- Linear function
- Training dataset
- Performance measures

### Least Square regression (LS)

- Optimization problem
- Geometrical interpretation
- Probabilistic interpretation

### Ridge regression

- Tikhonov regularization
- Optimization problem
- Probabilistic interpretation

### Variable selection with the Lasso

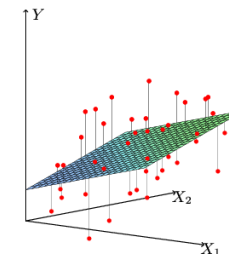
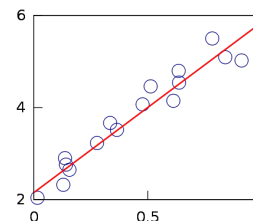
- Optimization problem
- Non smooth minimization

### Regularized least Square (RLS)

- General problem formulation

2/35

## Linear regression



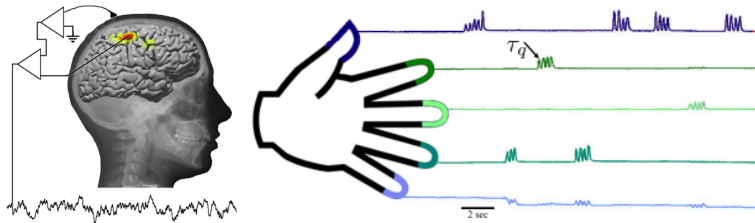
### Objective

Train a linear function  $f(\cdot)$  that can predict a continuous value  $y \in \mathbb{R}$  from an observation  $\mathbf{x} \in \mathbb{R}^d$ .

In practice we want to find the coefficients  $(\mathbf{w}, b)$  of  $f(\cdot)$  using a training dataset  $\{\mathbf{x}_i, y_i\}_{i=1, \dots, n}$ .

4/35

## Application for a Brain Computer Interface (BCI)



### BCI Competition IV, Dataset 4

- ▶ Data: Recordings of ECoG brain signals and of simultaneous finger flexion of a subject (using a glove).
- ▶ Objective of the competition: predict movement of the 5 fingers of the subject from its recorded ECoG.
- ▶ Best performances were obtained using a linear model.

## How do we store training data ?

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^\top & 1 \\ \mathbf{x}_2^\top & 1 \\ \vdots & \vdots \\ \mathbf{x}_i^\top & 1 \\ \vdots & \vdots \\ \mathbf{x}_n^\top & 1 \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1j} & \dots & x_{1d} & 1 \\ x_{21} & x_{22} & \dots & x_{2j} & \dots & x_{2d} & 1 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ x_{i1} & x_{i2} & \dots & x_{ij} & \dots & x_{id} & 1 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nj} & \dots & x_{nd} & 1 \end{bmatrix}, \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_i \\ \vdots \\ y_n \end{bmatrix}$$

### Training data

- ▶  $\mathbf{x}_i \in \mathbb{R}^d$  observations for  $i = 1, \dots, n$ .
- ▶  $y_i \in \mathbb{R}$  values to predict for  $i = 1, \dots, n$ .

Matrix form:

- ▶  $\mathbf{X} \in \mathbb{R}^{n \times (d+1)}$  such that  $\mathbf{x} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n, \mathbf{e}]^\top$  with  $\mathbf{e} \in \mathbb{R}^d$  and  $e_i = 1, \forall i$
- ▶  $\mathbf{y} \in \mathbb{R}^n$  such that  $\mathbf{y} = [y_1, y_2, \dots, y_n]^\top$ .
- ▶  $\alpha \in \mathbb{R}^{d+1}$  is a vector such that  $\alpha = \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix}$

5/35

6/35

## Performance measure

How to measure performance of a prediction?

Let  $\mathbf{y}$  be the values to predict and  $\hat{\mathbf{y}}$  the predictions.

### Mean square error

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

- ▶ 0 for a perfect prediction.
- ▶ Not normalized (depends on the variance of  $\mathbf{y}$ )

### Correlation coefficient

$$r = \frac{\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})(\hat{y}_i - \bar{\hat{y}})}{\sigma_y \sigma_{\hat{y}}}$$

- ▶  $\bar{y} = \frac{1}{n} \sum_i y_i$  mean of  $\mathbf{y}$  and  $\sigma_y$  its.
- ▶ 1 for perfect prediction.
- ▶ Normalized :  $r \in (-1, 1)$ .

### Warning

Always measure performance of a model on data that has NOT been used for training or else there is a risk of **over-fitting**

7/35

8/35

## Error and squared error

### Principle

We have the following model:

$$y = \mathbf{x}^\top \mathbf{w} + b = \tilde{\mathbf{x}}^\top \alpha \quad (2)$$

where  $\tilde{\mathbf{x}} = [x_1, \dots, x_d, 1]^\top$  is  $\mathbf{x}$  concatenated with 1.

We seek for parameters  $(\mathbf{w}, b) \equiv \alpha$  of function  $f(\cdot)$  that works well on training data.

### Residuals

The residual of sample  $i$  is the prediction error :

$$\epsilon_i = y_i - \mathbf{x}_i^\top \mathbf{w} - b = y_i - \tilde{\mathbf{x}}_i^\top \alpha \quad (3)$$

We want the residuals to be the smallest possible in average. To this end we can measure the error as the square error:

$$\epsilon_i^2 = (y_i - \mathbf{x}_i^\top \mathbf{w} - b)^2 = (y_i - \tilde{\mathbf{x}}_i^\top \alpha)^2 \quad (4)$$

## Least square optimization problem

We see the function  $f(\cdot)$  minimizing the squared error on the training samples :

$$\min_f \frac{1}{2} \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2 = \frac{1}{2} \sum_{i=1}^n \epsilon_i^2 \quad (5)$$

using the linear form of  $f(\cdot)$ , we obtain the following optimization problem :

$$\min_{\mathbf{w}, b} \frac{1}{2} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \mathbf{w} - b)^2 \quad (6)$$

that is equivalent to

$$\min_{\alpha} \frac{1}{2} \sum_{i=1}^n (y_i - \tilde{\mathbf{x}}_i^\top \alpha)^2 \quad (7)$$

## Interpretation of least squares

$$J(\alpha) = \frac{1}{2} \sum_{i=1}^n \underbrace{(y_i - \tilde{\mathbf{x}}_i^\top \alpha)}_{\epsilon_i}^2$$

The problem can be seen as finding an hyperplane  $\mathbf{x}^\top \mathbf{w} + b = y$  in a  $\mathbb{R}^{d+1}$  space that best fits a point cloud  $(\mathbf{x}_i, y_i) i = 1, n$  with respect to the  $y$  dimension.

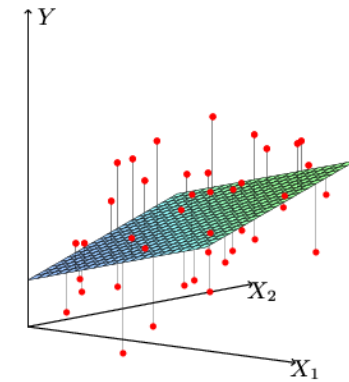


Figure: Residuals for the regression

9/35

10/35

## Matrix form of least squares (1)

### Residual as a vector

The residuals (error) on the samples can be expressed as:

$$\epsilon_i = y_i - \mathbf{x}_i^\top \mathbf{w} - b = y_i - \tilde{\mathbf{x}}_i^\top \alpha \quad (8)$$

Similarly to the training data, they can be stored in a vector  $\epsilon \in \mathbb{R}^n$  such that:

$$\epsilon = \mathbf{y} - \mathbf{X}\alpha \quad (9)$$

### Matrix form for least square

The least square optimization problem can be expressed as:

$$\min_{\alpha} \|\epsilon\|^2 = \|\mathbf{y} - \mathbf{X}\alpha\|^2 \quad (10)$$

where  $\|\cdot\|$  is the euclidean norm of a vector such that  $\|\epsilon\|^2 = \sum_{i=1}^n \epsilon_i^2$

## Matrix form of least squares (2)

The optimization problem can be expressed as:

$$\min_{\alpha} J(\alpha) \quad \text{avec} \quad J(\alpha) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\alpha\|^2$$

Using the properties of scalar product we get :

$$\begin{aligned} \min_{\alpha} J(\alpha) &= \frac{1}{2} \|\mathbf{y} - \mathbf{X}\alpha\|^2 \\ &= \frac{1}{2} (\mathbf{y} - \mathbf{X}\alpha)^\top (\mathbf{y} - \mathbf{X}\alpha) \\ &= \frac{1}{2} \mathbf{y}^\top \mathbf{y} - \frac{1}{2} \alpha^\top \mathbf{X}^\top \mathbf{y} - \frac{1}{2} \mathbf{y}^\top \mathbf{X}\alpha + \frac{1}{2} \alpha^\top \mathbf{X}^\top \mathbf{X}\alpha \\ &= \frac{1}{2} \mathbf{y}^\top \mathbf{y} - \alpha^\top \mathbf{X}^\top \mathbf{y} + \frac{1}{2} \alpha^\top \mathbf{X}^\top \mathbf{X}\alpha \end{aligned}$$

$\mathbf{y}^\top \mathbf{y}$  is a scalar,  $\mathbf{X}^\top \mathbf{y}$  is a vector  $\mathbb{R}^d$  and  $\mathbf{X}^\top \mathbf{X}$  is a  $(d+1) \times (d+1)$  matrix.

11/35

12/35

## Convex optimization basics

### Optimization problem

We want to solve

$$\min_{\alpha} J(\alpha) \quad \text{avec} \quad J(\alpha) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\alpha\|^2$$

where  $J(\alpha)$  is a convex function.

### Minimum of a convex function

Let  $J(\alpha)$  be a convex function  $\mathbb{R}^d \rightarrow \mathbb{R}$ .  $\alpha^*$  is a minimum  $J(\alpha)$  if and only if

$$\nabla J(\alpha^*) = \mathbf{0} \quad (11)$$

where  $\nabla J(\alpha) \in \mathbb{R}^d$  is the gradient of the function in  $\alpha$  such that

$$\nabla J(\alpha)_i = \frac{\partial J(\alpha)}{\partial \alpha_i} \quad \forall i$$

In order to find the minimum we need to find  $\alpha^*$  such that the gradient is  $\mathbf{0}$ .

13/35

## Least Squares solution

Minimizing the cost  $J(\alpha)$  corresponds to finding the parameter  $\alpha$  that leads to a null gradient:

$$\nabla J(\hat{\alpha}) = \mathbf{0} \quad \Leftrightarrow \quad -\mathbf{X}^T \mathbf{y} + \mathbf{X}^T \mathbf{X} \hat{\alpha} = \mathbf{0}$$

The solution of the minimization problem for Least Square is the vector  $\hat{\alpha}$  defined as

$$\hat{\alpha} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

### Hypothesis

$\mathbf{X}$  is a matrix of rank  $d + 1$  which means that  $\mathbf{X}^T \mathbf{X}$  is invertible.

In practice it means that  $n > d + 1$ , this method requires that we have more training samples than parameter to estimate.

15/35

## Gradient computation

$$\begin{aligned} J(\alpha) &= \frac{1}{2} \mathbf{y}^T \mathbf{y} - \alpha^T \mathbf{X}^T \mathbf{y} + \frac{1}{2} \alpha^T \mathbf{X}^T \mathbf{X} \alpha \\ \frac{\partial J(\alpha)}{\partial \alpha_i} &= 0 - p_i + \frac{1}{2} \sum_{j=1}^{d+1} (M_{ij} + M_{ji}) \alpha_j \end{aligned}$$

with  $\mathbf{p} = \mathbf{X}^T \mathbf{y}$  and  $\mathbf{M} = \mathbf{X}^T \mathbf{X}$

$$\blacktriangleright \frac{\partial \alpha^T \mathbf{p}}{\partial \alpha_i} = \frac{\partial \sum_{j=1}^{d+1} p_j \alpha_j}{\partial \alpha_i} = p_i$$

$$\blacktriangleright \frac{\partial \alpha^T \mathbf{M} \alpha}{\partial \alpha_i} = \frac{\partial \sum_{j=1}^{d+1} \sum_{k=1}^{d+1} \alpha_j \alpha_k M_{jk}}{\partial \alpha_i} = \sum_{j=1}^{d+1} \alpha_j M_{ji} + \sum_{k=1}^{d+1} \alpha_k M_{ik}$$

because  $(uv)' = uv' + u'v$  with  $u = \alpha_j$  et  $v = \sum_{k=1}^{d+1} \alpha_k M_{jk}$

$$\nabla J(\alpha) = -\mathbf{p} + \mathbf{M}\alpha = -\mathbf{X}^T \mathbf{y} + \mathbf{X}^T \mathbf{X} \alpha$$

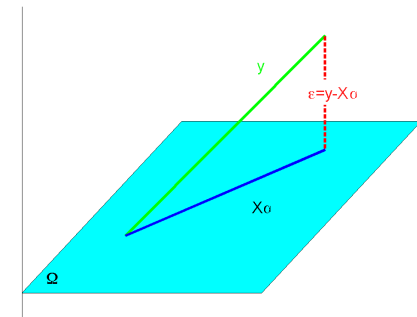
14/35

## Geometrical interpretation

$$\Omega = \text{span}\{\mathbf{X}\}$$

$\Omega$  is the linear subspace of  $\mathbb{R}^n$  generated by the columns of matrix  $\mathbf{X}$ .

$$\mathbf{z} \in \Omega \quad \Leftrightarrow \quad \exists \alpha \in \mathbb{R}^{d+1} \quad \mathbf{z} = \mathbf{X}\alpha$$



The least square is the projection of  $\mathbf{y}$  onto  $\Omega$ . We have :

$$\mathbf{X}\hat{\alpha} = \mathbf{H}\mathbf{y}$$

with the orthogonal projection operator  $\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$

16/35

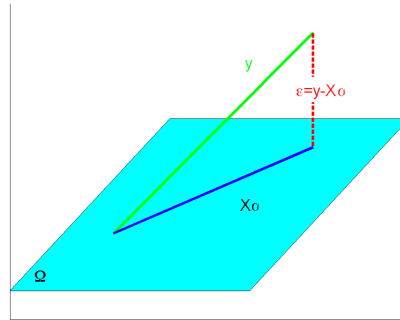
## Estimation and orthogonality

Objective value  $\|\epsilon\|^2$  is minimal for  $\alpha$  such that  $\mathbf{z} = \mathbf{X}\alpha$  is the orthogonal projection of  $\mathbf{y}$  on  $\Omega$ . This means that

$$\forall \mathbf{z} \in \Omega \quad \mathbf{z}^\top \epsilon = 0$$

Which means that the residual is orthogonal to all columns of  $\mathbf{X}$

$$\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\hat{\alpha}) = 0$$



$$\begin{aligned} \mathbf{X}^\top (\mathbf{y} - \mathbf{X}\hat{\alpha}) = 0 &\Leftrightarrow \mathbf{X}^\top \mathbf{y} - \mathbf{X}^\top \mathbf{X}\hat{\alpha} = 0 \\ &\Leftrightarrow \mathbf{X}^\top \mathbf{X}\hat{\alpha} = \mathbf{X}^\top \mathbf{y} \\ &\Leftrightarrow \hat{\alpha} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \end{aligned}$$

17/35

## Probabilistic interpretation of least squares (2)

### Maximum likelihood estimator

The MLE estimator is the solution of

$$\max_{\alpha} \mathcal{L}(\alpha)$$

In practice people often maximize the log-likelihood in order to have a simpler problem with the same solution.

### Maximizing the Log-likelihood

$$\begin{aligned} \log(\mathcal{L}(\alpha)) &= n \log \left( \frac{1}{\sqrt{2\pi\sigma_n^2}} \right) + \sum_{i=1}^n \log \left( \exp \left( -\frac{(y_i - \hat{\mathbf{x}}_i^\top \alpha)^2}{2\sigma_n^2} \right) \right) \\ &= n \log \left( \frac{1}{\sqrt{2\pi\sigma_n^2}} \right) - \frac{1}{2\sigma_n^2} \sum_{i=1}^n (y_i - \hat{\mathbf{x}}_i^\top \alpha)^2 = cst - J(\alpha) \end{aligned}$$

Maximizing the log-likelihood wrt  $\alpha$  is equivalent to minimizing  $J(\alpha)$ .

19/35

## Probabilistic interpretation of least squares (1)

### Observation model

We suppose that the observation model is the following:

$$y = \mathbf{x}^\top \mathbf{w} + b + \epsilon = \hat{\mathbf{x}}^\top \alpha + \epsilon$$

- ▶  $\epsilon$  is a centered random variable such that  $\epsilon \sim \mathcal{N}(0, \sigma_n^2)$ .
- ▶ The probability of a given observation  $(\mathbf{x}, y)$  when the parameters  $\alpha$  are known is then

$$p(\mathbf{x}, y | \alpha) = \frac{1}{\sqrt{2\pi\sigma_n^2}} \exp \left( -\frac{(y - \hat{\mathbf{x}}^\top \alpha)^2}{2\sigma_n^2} \right)$$

### Likelihood on the dataset

The likelihood for the whole dataset can be expressed as

$$\mathcal{L}(\alpha) = \prod_{i=1}^n p(\mathbf{x}_i, y_i | \alpha) = \left( \frac{1}{\sqrt{2\pi\sigma_n^2}} \right)^n \prod_{i=1}^n \exp \left( -\frac{(y_i - \hat{\mathbf{x}}_i^\top \alpha)^2}{2\sigma_n^2} \right)$$

18/35

## Why regularize ?

### Least Squares

We minimize the prediction error on the training data:

$$\min_{\alpha} J(\alpha) \quad \text{avec} \quad J(\alpha) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\alpha\|^2$$

Problem solution is

$$\hat{\alpha} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

### Numerical problems

- ▶ When  $n < d + 1$ , matrix  $\mathbf{X}^\top \mathbf{X}$  is non-invertible.
  - ▶ There exists an infinity of solutions, problem is ill-posed.
- ⇒ regularization (among all possible solutions, pick the simplest).

20/35

## Ridge Regression

### Optimization problem

$$\min_{\mathbf{w}, b} \frac{1}{2} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \mathbf{w} - b)^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2 \quad (12)$$

- ▶ We add a regularization term  $\|\mathbf{w}\|^2$  weighted by the regularization coefficient  $\lambda \geq 0$ .
- ▶ Parameter  $\lambda$  can be chosen to limit over-fitting on the data.
- ▶ This regularization promotes parameters  $\mathbf{w}$  of minimal norm.
- ▶ It makes the optimization problem strictly convex (a unique solution).
- ▶ When  $\lambda = 0$  the problem boils down to the least squares (special case).

### Gradient computation

$$J'(\alpha) = \frac{1}{2} \mathbf{y}^\top \mathbf{y} - \alpha^\top \mathbf{X}^\top \mathbf{y} + \frac{1}{2} \overbrace{\alpha^\top (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{S}) \alpha}^{\frac{1}{2} \alpha^\top (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{S}) \alpha} + \frac{\lambda}{2} \alpha^\top \mathbf{S} \alpha$$

$$\frac{\partial J'(\alpha)}{\partial \alpha_i} = 0 - p_i + \frac{1}{2} \sum_{j=1}^{d+1} (M_{ij} + M_{ji}) \alpha_j$$

with  $\mathbf{p} = \mathbf{X}^\top \mathbf{y}$  and  $\mathbf{M} = \mathbf{X}^\top \mathbf{X} + \lambda \mathbf{S}$

- ▶  $\frac{\partial \alpha^\top \mathbf{p}}{\partial \alpha_i} = \frac{\partial \sum_{j=1}^{d+1} p_j \alpha_j}{\partial \alpha_i} = p_i$
- ▶  $\frac{\partial \alpha^\top \mathbf{M} \alpha}{\partial \alpha_i} = \frac{\partial \sum_{j=1}^{d+1} \sum_{k=1}^{d+1} \alpha_j \alpha_k M_{jk}}{\partial \alpha_i} = \sum_{j=1}^{d+1} \alpha_j M_{ji} + \sum_{k=1}^{d+1} \alpha_k M_{ik}$

because  $(uv)' = uv' + u'v$  with  $u = \alpha_j$  and  $v = \sum_{k=1}^{d+1} \alpha_k M_{jk}$

$$\nabla J'(\hat{\alpha}) = -\mathbf{p} + \mathbf{M} \hat{\alpha} = -\mathbf{X}^\top \mathbf{y} + (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{S}) \hat{\alpha}$$

## Matrix form of ridge regression

$$\min_{\alpha} \left\{ J'(\alpha) = \frac{1}{2} \|\mathbf{y} - \mathbf{X} \alpha\|^2 + \frac{\lambda}{2} \alpha^\top \mathbf{S} \alpha \right\} \quad (13)$$

with  $\mathbf{S} \in \mathbb{R}^{(d+1) \times (d+1)}$  a matrix defined as

$$\mathbf{S} = \begin{bmatrix} 1 & 0 & \dots & 0 & 0 \\ 0 & 1 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & 1 & 0 \\ 0 & 0 & \dots & 0 & 0 \end{bmatrix} \quad S_{i,j} = \begin{cases} 1 & \text{si } i = j \text{ et } i \leq d \\ 0 & \text{sinon} \end{cases} \quad (14)$$

$\mathbf{S}$  is a diagonal matrix containing 1 on the diagonal except for the last term equal to 0. Problem (12) and (13) are equivalent because

$$\alpha^\top \mathbf{S} \alpha = \sum_{i,j=1}^{d+1} \alpha_i \alpha_j S_{i,j} = \sum_{i=1}^d \alpha_i^2 = \sum_{i=1}^d \mathbf{w}_i^2 = \|\mathbf{w}\|^2 \quad (15)$$

21/35

22/35

### Ridge regression solution

Minimizing the cost  $J'(\alpha)$  corresponds to finding the parameter  $\alpha$  that leads to a null gradient:

$$\nabla J'(\hat{\alpha}) = 0 \Leftrightarrow -\mathbf{X}^\top \mathbf{y} + (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{S}) \hat{\alpha} = 0$$

The solution of the minimization problem for ridge regression is the vector  $\hat{\alpha}$  defined as :

$$\hat{\alpha} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{S})^{-1} \mathbf{X}^\top \mathbf{y}$$

### Regularization

Matrix  $\mathbf{S}$  adds  $\lambda$  on the diagonal of  $\mathbf{X}^\top \mathbf{X}$ , making the matrix  $\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{S}$  invertible. The problem is now well posed and has a unique solution.

23/35

24/35

## Probabilistic interpretation of Ridge Regression

### Prior distribution for the parameters

- ▶ In LS we had no prior information about  $\alpha$ .
- ▶ We suppose that the  $\mathbf{w}$  parameter has been drawn from  $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \sigma_p^2 \mathbf{I})$ .
- ▶ Probability of a given  $\mathbf{w}$  is :  $p(\mathbf{w}) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma_p^2}} \exp\left(-\frac{(w_k)^2}{2\sigma_p^2}\right)$

### Maximum likelihood estimator

$$\begin{aligned} \log(\mathcal{L}(\alpha)) &= cst - \frac{1}{2\sigma_n^2} \sum_{i=1}^n (y_i - \hat{\mathbf{x}}_i^\top \alpha)^2 - \frac{1}{2\sigma_p^2} \sum_{k=1}^d (w_k)^2 \\ &= cst - \frac{1}{2\sigma_n^2} \sum_{i=1}^n (y_i - \hat{\mathbf{x}}_i^\top \alpha)^2 - \frac{1}{2\sigma_p^2} \|\mathbf{w}\|^2 \end{aligned}$$

- ▶ Maximizing the log-likelihood wrt  $\alpha$  is equivalent to minimizing  $J(\alpha)$ .
- ▶ Problems are equivalent when  $\lambda = \frac{\sigma_n^2}{\sigma_p^2}$ .

25/35

## Lasso estimator

$$\min_{\mathbf{w}, b} \frac{1}{2} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \mathbf{w} - b)^2 + \lambda \sum_{k=1}^d |w_k| \quad (16)$$

### Optimization problem

- ▶  $\|\mathbf{w}\|_1 = \sum_{k=1}^d |w_k|$  is the L1 norm of vector  $\mathbf{w}$ .
- ▶ Objective function is non differentiable in  $w_k = 0, \forall k$ .
- ▶ For a large enough  $\lambda$  the solution of the problem is sparse (some components of  $\mathbf{w}$  are exactly 0).
- ▶ The problem is equivalent to

$$\min_{\mathbf{w}, b, \|\mathbf{w}\|_1 \leq \mu} \frac{1}{2} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \mathbf{w} - b)^2 \quad (17)$$

I.e. there exists a  $\mu$  that leads to the same solution of the problem for a given  $\lambda$ .

26/35

## Lasso with no bias

### Data and model

- ▶ If the model has no bias  $b$  it means that the prediction can be expressed as

$$f(\mathbf{x}) = \sum_i x_i w_i = \mathbf{x}^\top \mathbf{w}$$

- ▶  $\mathbf{X}' \in \mathbb{R}^{n \times d}$  is the  $\mathbf{X}$  matrix without the last columns containing ones.
- ▶ Prediction can be done with  $\mathbf{X}'\mathbf{w}$  and  $\mathbf{w}$  is the only parameters

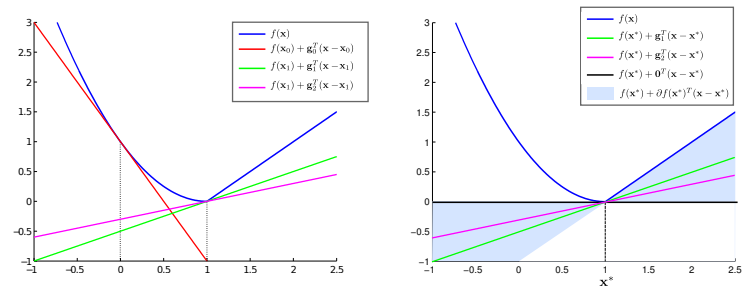
### Matrix form of the optimization problem

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}'\mathbf{w}\|^2 + \lambda \|\mathbf{w}\|_1 \quad (18)$$

- ▶ In the remaining we will focus on the Lasso with no bias for readability.
- ▶ Extension of the result when adding the bias is tedious but straightforward.

27/35

## Subgradients and subdifferential



### Non differentiable function

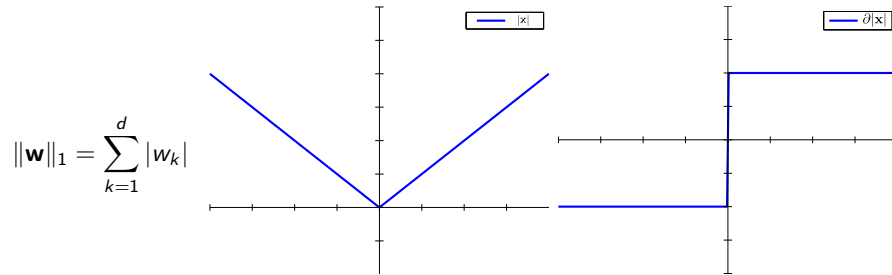
- ▶ A non differentiable function might not have a gradient (ex  $|\cdot|$  in 0).
- ▶ The tool used in place of gradient is the subdifferential and subgradients.
- ▶ For a convex function  $f(\mathbf{x})$ ,  $\mathbf{g}$  is a subgradient of  $f$  in  $\mathbf{x}_0$  if

$$f(\mathbf{x}) \geq f(\mathbf{x}_0) + \mathbf{g}^\top (\mathbf{x} - \mathbf{x}_0) \quad (19)$$

- ▶ The set of all subgradients at  $\mathbf{x}_0$  is the subdifferential  $\partial f(\mathbf{x}_0)$ .
- ▶  $\mathbf{x}_0$  is a minimum of the convex function  $f$  if  $\mathbf{0} \in \partial f(\mathbf{x}_0)$ .

28/35

## Subdifferential for the L1 norm



### L1 norm

- ▶ The subdifferential is of the form

$$\partial\|\mathbf{w}\|_1 = \begin{cases} \mathbf{g} : \|\mathbf{g}\|_\infty \leq 1 & \text{si } \mathbf{w} = 0 \\ \mathbf{g} : \|\mathbf{g}\|_\infty \leq 1 \text{ et } \mathbf{g}^T \mathbf{w} = \|\mathbf{w}\|_1 & \text{si } \mathbf{w} \neq 0 \end{cases}$$

- ▶ Which give

$$\partial\|\mathbf{w}\|_1 = \begin{cases} \mathbf{g} : g_k \in [-1, 1] & \text{si } \mathbf{w} = 0 \\ \mathbf{g} : g_k = \text{sign}(w_k) & \text{si } \mathbf{w} \neq 0 \end{cases}$$

## Interpreting optimality conditions

### Correlation with the residue

$$c_k = \mathbf{X}'_k{}^T (\mathbf{y} - \mathbf{X}'\mathbf{w}^*) = \|\mathbf{X}'_{\cdot,k}\| \|\mathbf{y} - \mathbf{X}'\mathbf{w}^*\| \cos(\theta)$$

- ▶  $c_k$  is the scalar product between the feature  $k$  and the residuals  $\epsilon = \mathbf{y} - \mathbf{X}'\mathbf{w}^*$ .
- ▶  $\theta$  is the angle between the two vectors.
- ▶  $c_k = 0, \forall k$  for Least Squares regression (optimality condition).

### Effect of the regularization parameter $\lambda$

- ▶  $\lambda = 0$  boils down to Least Squares (no sparsity).
- ▶ If  $\lambda$  is small we have  $w_k = 0$  only for variable  $k$  where

$$|\mathbf{X}'_k{}^T (\mathbf{y} - \mathbf{X}'\mathbf{w}^*)| \leq \lambda$$

- ▶ If  $\lambda$  is very large at some point we have for all  $k$ ,

$$|\mathbf{X}'_k{}^T \mathbf{y}| \leq \lambda \quad \text{which means} \quad w_k = 0, \forall k.$$

## Optimality conditions of the Lasso

### Optimality conditions

$\mathbf{w}^*$  is a solution of the optimization problem if

$$\mathbf{0} \in \partial J_{\text{lasso}}(\mathbf{w}^*) \quad \text{with} \quad J_{\text{lasso}}(\mathbf{w}) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}'\mathbf{w}\|^2 + \lambda \|\mathbf{w}\|_1$$

This can be reformulated as the following condition

$$-\mathbf{X}'^T (\mathbf{y} - \mathbf{X}'\mathbf{w}^*) + \lambda \mathbf{g} = \mathbf{0} \quad \text{with} \quad \mathbf{g} \in \partial\|\mathbf{w}^*\|_1$$

### Conditions on the components of $\mathbf{w}^*$

$$\begin{aligned} w_k^* \neq 0 &\Rightarrow -\mathbf{X}'_k{}^T (\mathbf{y} - \mathbf{X}'\mathbf{w}^*) + \lambda \text{sign}(w_k^*) = 0 \\ w_k^* = 0 &\Rightarrow |\mathbf{X}'_k{}^T (\mathbf{y} - \mathbf{X}'\mathbf{w}^*)| \leq \lambda \end{aligned}$$

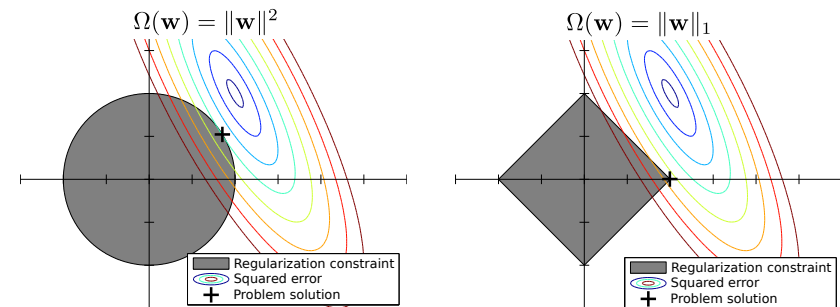
- ▶  $\mathbf{X}'_k$  is the  $k$ th column of  $\mathbf{X}'$  (feature  $k$ ).

- ▶ There is no closed-form solution except for special cases of  $\mathbf{X}'$  (orthogonality).

29/35

30/35

## L2 VS L1 regularization



### Tikhonov regularization

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}'\mathbf{w}\|^2 + \lambda \Omega(\mathbf{w})$$

$\Leftrightarrow$

### Ivanov regularization

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{y} - \mathbf{X}'\mathbf{w}\|^2 \\ \text{s.t.} \quad & \Omega(\mathbf{w}) \leq \mu \end{aligned}$$

The two optimization problems are equivalent for a strictly convex function  $\Omega$ .

31/35

32/35



## Solving the optimization problem

### Least-angle regression (LARS)

- ▶ Algorithm developed by B. Efron, T. Hastie, I. Johnstone and R. Tibshirani.
- ▶ Allows to find efficiently the whole regularization path (all solution for all  $\lambda$ )
- ▶ Potential problems with highly correlated variables.

### Proximal gradient descent (PGD)

- ▶ Subgradient descent is known to converge slowly.
- ▶ Proximal gradient descent allows for acceleration of the resolution.
- ▶ Can be seen as a Majoration-Minimization method.
- ▶ Each iteration is a simple soft thresholding of the parameter.
- ▶ Can be coupled for active sets to speedup sparse solutions.

### Coordinate descent algorithm

- ▶ Optimize each components of  $\mathbf{w}$  independently until convergence.
- ▶ Very fast for sparse solutions.

33/35

## Coordinate descent for the Lasso

### Algorithm

- ▶ Select an initial vector  $\mathbf{w}$  (usually  $\mathbf{w} = \mathbf{0}$ ).
- ▶ For all  $k$  :  $w_k \leftarrow \min_{w_k} J_{\text{lasso}}(\mathbf{w})$  with all  $w_j, j \neq k$  fixed
- ▶ Repeat until optimality conditions are satisfied.

### Iteration for $w_k$

$$\begin{aligned} \min_{w_k} \quad & \frac{1}{2} \|\mathbf{y} - \mathbf{X}'\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_1 \\ \min_{w_k} \quad & \frac{1}{2} \|\mathbf{y} - \sum_{j \neq k} \mathbf{X}'_j \mathbf{w}_j - \mathbf{X}'_k w_k\|_2^2 + \lambda |w_k| \\ \min_{w_k} \quad & \frac{1}{2} \|\mathbf{s} - \mathbf{X}'_k w_k\|_2^2 + \lambda |w_k| \end{aligned}$$

were  $\mathbf{s} = \mathbf{y} - \sum_{j \neq k} \mathbf{X}'_j \mathbf{w}_j$  is the residual wrt  $w_k$ . The last problem is a Lasso with only one variable, its solution is

$$w_k^* = \text{sign}(\mathbf{X}'_k^T \mathbf{s}) (|\mathbf{X}'_k^T \mathbf{s}| - \lambda)_+$$

This operator is called the soft thresholding.

34/35

## Regularized linear regression

General problem formulation:

$$\min_{\mathbf{w}, b} \sum_{i=1}^n L(y_i, \mathbf{w}^T \mathbf{x}_i + b) + \lambda \Omega(\mathbf{w}) \quad (20)$$

With

- ▶  $L(\dots)$  a loss function.
- ▶  $\Omega(\cdot)$  a regularization term.

Examples:

### Loss function $L(y, \hat{y})$

- ▶  $(y - \hat{y})^2$ , quadratic (this course).
- ▶  $|y - \hat{y}|$ , absolute value.
- ▶  $\min(0, |y - \hat{y}| - \epsilon)$  epsilon insensitive

### Regularizations $\Omega(\mathbf{w})$

- ▶  $\|\mathbf{w}\|_2^2$ , quadratic.
- ▶  $\|\mathbf{w}\|_1$ ,  $\ell_1$  norm.
- ▶  $\mathbf{w}^T \Sigma \mathbf{w}$ , Mahalanobis.

35/35