

Chap. 2
Linear regression

I) Model specification

$x = (x_1, x_2, \dots, x_D) \rightarrow$ one input

output $y = \underbrace{\sum_{i=1}^D w_i x_i}_{= w^T x} + \underbrace{\varepsilon}_{\substack{\text{Gaussian noise, variance } \sigma^2}}$

If we want to model a non-linear relationship, we can replace x by $\phi(x)$ (this is called basis function expansion)

(x is 1-dim and $\phi(x)$ can be, ex., $\phi(x) = [1, x, x^2, \dots, x^{D-1}]$)

the model has σ^2 parameters $\theta = \begin{matrix} \sigma^2 \\ w \end{matrix}$ (in general, the parameters in a model is denoted by θ).

II) Maximum likelihood estimation (least squares)

A common way to estimate the parameters is the MLE, defined as:

$$\hat{\theta} := \underset{\theta}{\operatorname{argmax}} \log p(\mathcal{D} | \theta) \quad (\mathcal{D} = \text{data})$$

$$= \ell(\theta)$$

More precisely: if the $(x^{(i)}, y^{(i)})$ we have access to are iid, we write:

$$\ell(\theta) = \sum_{i=1}^N \log p(y^{(i)} | x^{(i)}, \theta)$$

where: $y \mapsto p(y | x, \theta)$ is the density of y knowing x and supposing we know the vector w and supposing ε is $\mathcal{N}(0, \sigma^2)$, that is

$$p(y | x, \theta) = \frac{\exp\left(-\frac{(w^T x - y)^2}{2\sigma^2}\right)}{\sqrt{2\pi\sigma^2}}$$

Instead of maximizing $l(\theta)$ (= log-likelihood), we could minimize the negative-log-likelihood (NLL):

$$NLL(\theta) := - \sum_{i=1}^N \log p(y^{(i)} | x^{(i)}, \theta)$$

(because some software packages are designed to find minima and not maxima)

$$\begin{aligned} \text{We have: } l(\theta) &= \sum_{i=1}^N \log \left[\left(\frac{1}{2\pi\sigma^2} \right)^{\frac{1}{2}} \exp \left(-\frac{1}{2\sigma^2} (y^{(i)} - w^T x^{(i)})^2 \right) \right] \\ &= -\frac{1}{2\sigma^2} \text{RSS}(w) - \frac{N}{2} \log(2\pi\sigma^2) \end{aligned}$$

where RSS is the residual sum of squares

$$\text{RSS}(w) := \sum_{i=1}^N (y_i - w^T x_i)^2$$

RSS is also called the sum of squared errors (or SSE)

and $\frac{\text{SSE}}{N}$ is called the mean square error (or MSE)

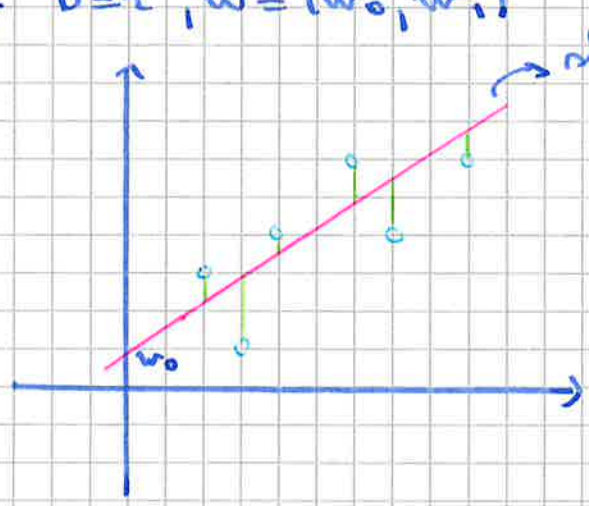
It can also be written as the square of the l^2 norm of the vector of residual errors:

$$\text{RSS}(w) = \sum_{i=1}^N (\epsilon^{(i)})^2 = \|\epsilon\|_2^2$$

↳ where $\epsilon^{(i)} := y^{(i)} - w^T x^{(i)}$

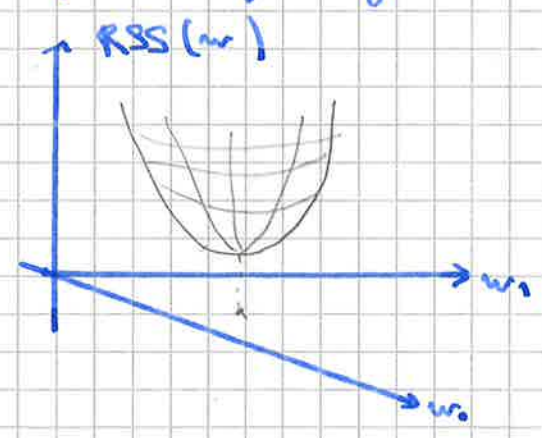
Suppose σ is known and we only look at the vector w minimizing the RSS → this method is known as least squares.

Ex: $D=2, w = (w_0, w_1)$



$x^{(i)} =$ blue points

We look for the pink line minimizing the sum of the squares of the green distances.



1) Derivation of the MLE

$$NLL(w) = \frac{1}{2} (y - Xw)^T (y - Xw) = \frac{1}{2} w^T (X^T X) w - w^T (X^T y) + \frac{1}{2} y^T y$$

$$\begin{cases} y = (y^{(1)}, \dots, y^{(N)}) \text{ (column vector)} \\ w = (w_0, w_1, \dots, w_D) \text{ (column vector)} \\ X = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & \dots & x_D^{(1)} \\ \vdots & \vdots & & \vdots \\ x_1^{(N)} & x_2^{(N)} & \dots & x_D^{(N)} \end{bmatrix} \end{cases}$$

where

$$\begin{cases} X^T X = \sum_{i=1}^N \begin{bmatrix} (x_1^{(i)})^2 & \dots & x_1^{(i)} x_D^{(i)} \\ \vdots & & \vdots \\ x_1^{(i)} x_D^{(i)} & \dots & (x_D^{(i)})^2 \end{bmatrix} \\ X^T y = \sum_{i=1}^N x^{(i)} \times y_i \end{cases}$$

We compute the gradient:

$$\nabla(NLL(w)) = (X^T X) w - X^T y$$

Indeed if A is a symmetric square matrix

$$A = \begin{bmatrix} a_{11} & \dots & a_{1D} \\ \vdots & & \vdots \\ a_{D1} & \dots & a_{DD} \end{bmatrix}$$

$$w^T A w = \sum_{i=1}^D \sum_{j=1}^D w_i a_{ij} w_j$$

whose gradient is a vector with line k:

$$\begin{aligned} \frac{\partial}{\partial w_k} \left(\sum_{i=1}^D \sum_{j=1}^D w_i a_{ij} w_j \right) &= \sum_{j=1}^D a_{kj} w_j + \sum_{i=1}^D w_i a_{ik} \\ &= 2 \sum_{j=1}^D a_{kj} w_j \end{aligned}$$

i.e. the vector 2Aw

and the gradient of $w^T \beta$ (for β a vector)

is a vector with line h :

$$\frac{\partial}{\partial w_i} \left(\sum_{i=1}^D w_i \beta_i \right) = \beta_i h$$

i.e. the vector β .

17

The critical point (the w for which $\nabla(\text{MSE}(w))=0$) is:

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

under the assumption that this matrix is invertible (explanation later)

i) Geometric interpretation

We assume $N > D$. We set $E = \text{Vect} \left(\begin{pmatrix} x_i^{(1)} \\ \vdots \\ x_i^{(N)} \end{pmatrix}, 1 \leq i \leq D \right)$.

E is a sub-space of dimension $\leq D$ of \mathbb{R}^N . Let us suppose $\dim(E) = D$. We set (for $1 \leq i \leq D$):

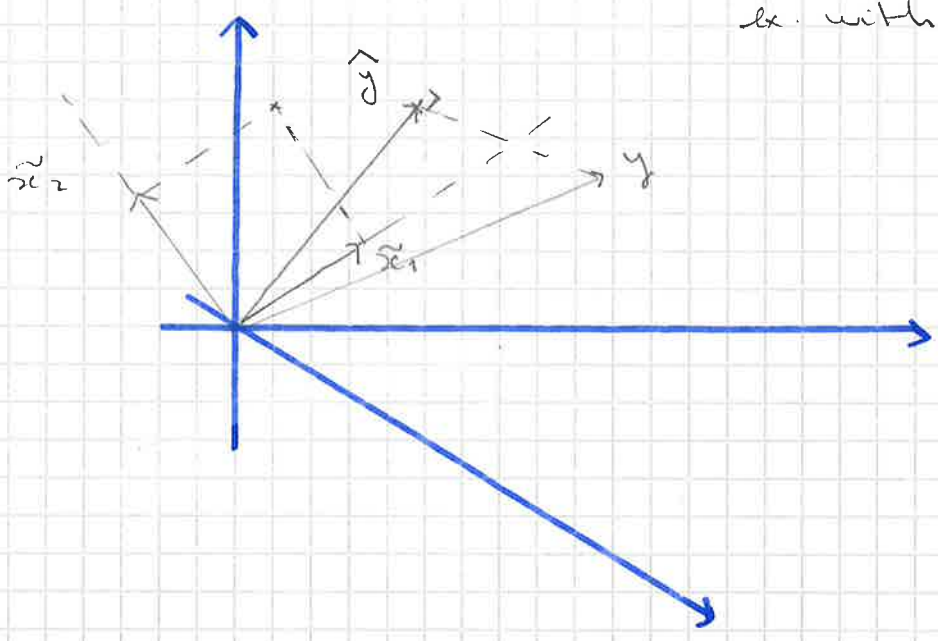
$$\tilde{x}_i = \begin{pmatrix} x_i^{(1)} \\ \vdots \\ x_i^{(N)} \end{pmatrix}$$

What we do in the linear regression is that we seek a vector $\hat{y} \in \mathbb{R}^N$ such that $\hat{y} \in E$ and is as close as possible to y : $\hat{y} = \text{argmin}_{y \in \text{span}(\tilde{x}_1, \dots, \tilde{x}_D)} \|y - \hat{y}\|_2$

Or to put it differently, we write $\hat{y} = w_1 \tilde{x}_1 + \dots + w_D \tilde{x}_D$
($\hat{y} \in \text{span}(\tilde{x}_1, \dots, \tilde{x}_D)$)
i.e. $\hat{y} = Xw$

To minimize the norm of $y - \hat{y}$, we want $y - \hat{y}$ to be orthogonal to $\text{span}(\tilde{x}_1, \dots, \tilde{x}_D) = E$

ex. with $\begin{cases} N=3 \\ D=2 \end{cases}$



So $\tilde{x}_j^T (y - \hat{y})$ for $j=1, 2, \dots, D$. Hence:

$$\tilde{x}_j^T (y - Xw) \text{ for } j=1, 2, \dots, D$$

$$X^T (y - Xw) = 0$$

And so $w = (X^T X)^{-1} X^T y$

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

this is the orthogonal projection of y onto $\text{span}(\tilde{x}_1, \dots, \tilde{x}_D)$.
the projection matrix $P := X(X^T X)^{-1} X^T$ is called the hat matrix since it "puts a hat on y ".

Why is $X^T X$ invertible in the case $N > D$?

For any $y \in \text{Span}(\tilde{x}_1, \dots, \tilde{x}_D)$, $\exists!$ w such that:
 $X^T (y - Xw) = 0$ (also true for $y \in \text{Span}(\tilde{x}_1, \dots, \tilde{x}_D)$)

i.e. $X^T X w = X^T y$

\downarrow
 a generic element of (x_1, \dots, x_N)
 remembers $x_i = \begin{pmatrix} x_1^{(i)} \\ \vdots \\ x_D^{(i)} \end{pmatrix}$
 \downarrow
 $w \in \mathbb{R}^D$

this means that $X^T X$ is invertible.

III) Ridge regression

The MLE can overfit because it picks the parameters that are good for the training data. But the training data can be noisy so there is no reason to pick the best parameters for the training data.

The alternative to linear regression is to use a Bayesian statistics approach. We suppose $y_i = w_0 + w^T x_i + \epsilon_i$

We assume an a priori distribution of the w_1, \dots, w_D (corresponding to a prior knowledge / our preconceived idea / ...) here, we have in mind that they should not be too big so we choose the following a priori:

$$p(w) = \prod_{j=1}^D \mathcal{N}(w_j | 0, \sigma^2)$$

↳ meaning each w_j is ind. of the others and of law $\mathcal{N}(0, \sigma^2)$

The likelihood is the density of y_1, \dots, y_D under the assumption that w_1, \dots, w_D is what it is:

$$p(y | x_i, \frac{\theta}{\sigma}) = \prod_{i=1}^N \mathcal{N}(y_i | w_0 + w^T x_i, \sigma^2)$$

$= w_1, \dots, w_D$

↳ knowing w_1, x_i , y_i has a Gaussian density

The observations y_1, \dots, y_D are fixed. The higher $p(y | x_i, \theta)$ is, the more likely it is that w_1, \dots, w_D has a certain value.

The a posteriori distribution is the product of the a priori distribution and the likelihood, renormalized to be a probability density (in w).

We want to find the MAP (Maximum A Posteriori),
 so we look for w maximizing

$$\prod_{i=1}^N \log W(y_i | w^T x_i, \sigma^2) + \prod_{j=1}^D \log \pi(w_j | 0, \tau^2)$$

(we took the log and threw away the normalization constant)
 or, equivalently, minimizing:

$$\frac{1}{N} \left[\sum_{i=1}^N \frac{(y_i - w_0 - w^T x_i)^2}{2\sigma^2} + \sum_{j=1}^D \frac{w_j^2}{2\tau^2} \right]$$

or (equ.) minimizing $\frac{1}{N} \left[\sum_{i=1}^N (y_i - w_0 - w^T x_i)^2 + \lambda w^T w \right]$
 $(\lambda = \frac{\sigma^2}{\tau^2})$

We can interpret the first term as the MSE of before and the second term as a complexity penalty (we do not want the w_j to be too big). The corresponding solution is given by:

$$\hat{w}_{ridge} = (\lambda I_D + X^T X)^{-1} X^T y$$

this technique is known as ridge regression or penalized least squares. In general, adding a Gaussian prior to the parameters of a model to encourage them to be small is called l_2 -regularization or weight-decay.

Projetés la figure 7.8 p.226 du livre.

IV) The LASSO

We define $\hat{w}_{lasso} = \operatorname{argmin} \left\{ \sum_{i=1}^N (y_i - w_0 - w^T x_i)^2, w \text{ such that } \sum_{j=1}^D |w_j| \leq t \right\}$
 $(t \text{ to be chosen})$

- Making t sufficiently small will cause some of the coefficients to be exactly 0.

- If t is chosen larger than $t_0 = \sum_{i=1}^D |\hat{w}_j^{LS}|$ (2)
↑
the ones we found in the least-square regression

then $\hat{w}_{\text{lasso}} = \hat{w}^{LS}$

- t should be adaptively chosen to minimize an estimate of expected prediction error

Def: LASSO and Ridge regressions are part of what is called shrinkage methods

Remarques pour le TP: p. 45 à 55 du O'Reilly

- * Le λ de la Ridge Regression est le λ du cours
- * Le t du LASSO est le t du cours.
- * L'algorithme de python trouve la solution du LASSO par une boucle (on doit converger vers la solution). On peut préciser le nombre d'itérations maximales (max_iter) autorisées (en particulier, on peut relever ce plafond pour s'assurer de la convergence).