



**Mondragon  
Unibertsitatea**

Escuela Politécnica  
Superior

# Regression

# Outline

- Simple regression
  - Parametric. Simple Linear Regression
  - Non-parametric. K-Nearest Neighbour
- Multiple regression
  - Non-additive
  - Non-linear
    - Polynomial regression
    - Splines regression
- Fitting procedures
  - Dimensionality Reduction
    - Principal Components Regression (PCR)
    - Partial Least Squares (PLS)
  - Shrinkage
    - Ridge regression
    - Lasso
    - ElasticNet
- Validation and model selection
  - Explained variance
  - Mean absolute, quadratic and quadratic logarithmic errors
  - Coefficient of determination

# Simple linear regression

$$Y = f(X) + \epsilon = a + bX + \epsilon$$

- A **linear relation** between the predictor  $X$  and the response  $Y$  is assumed.
- The **coefficients** to be optimized are  $a$  and  $b$ . The **error** would be  $\epsilon$ .
- It is not efficient by itself, but it is the basis for many sophisticated models.
- Given pairs of values  $(x_i, y_i), i = 1, \dots, N$  the coefficients  $a$  and  $b$  must be estimated in such a way that the predictions of the  $y_i$  from the  $x_i$  are as good as possible.
- For instance:

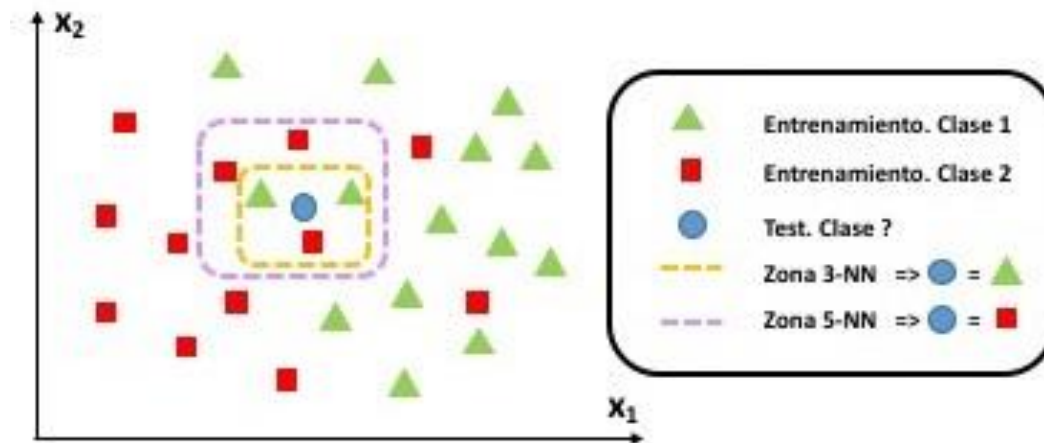
$$\hat{a}, \hat{b} = \arg \min_{a,b} \sum_{i=1}^N (y_i - a - bx_i)^2$$

Therefore, the coefficients that minimize the sum of the quadratic error.

- Note: the minimum above is located at the same point  $(\hat{a}, \hat{b})$  as the minimum mean square error, then it is equivalent to using that score.

# K-Nearest Neighbors regression

- Given a training set and a sample, for which we want to obtain a prediction:
  - Calculate all pairwise “distances” between the test samples and the train samples
  - Select the K smaller distances
  - Predict the class for the test sample by means of averaging among its K nearest neighbors



- Variants:
  - Use a weighted mean, which weights are inversely proportional to the distance
  - Without prefixing K, choose a threshold for defining the neighborhood (radius) and consider as neighbors all training samples that are not further than that radius
- **Prefixing K** is a tricky task to solve
- We must design an strategy for **breaking ties**

# Multiple linear regression

$$Y = f(\mathbf{X}) + \epsilon = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \epsilon$$

- A **linear relation** between the predictors  $X_1, \dots, X_p$ , and the response  $Y$  is assumed.
- The **coefficients** to be optimized are  $p+1$  coefficients  $\beta_j$ . The **error** would be  $\epsilon$ .
- It suffers in situations de co-linearity of the predictions.
- Given pairs of values  $(x_i, y_i), i = 1, \dots, N$  the coefficients  $\beta_j$  estimated in such a way that the predictions of the  $y_i$  from the  $x_i$  are as good as possible.
- For instance:

$$\hat{\beta} = \arg \min_{\beta} \sum_{i=1}^N (y_i - \beta_0 + \beta_1 X_{1,i} + \dots + \beta_p X_{p,i})^2$$

Therefore, the coefficients that minimize the sum of the quadratic error.

- Note: the minimum above is located at the same point  $\hat{\beta}$  as the minimum mean square error, then it is equivalent to using that score.

# Linearity and additivity

$$Y = f(\mathbf{X}) + \epsilon = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \epsilon$$

- Linearity. The predictors contribute linearly in the prediction.
- Additivity. The effects over the prediction of the perturbations in a predictor are independent of the value of the rest of predictors.
- Avoiding linearity implies, for instance, include quadratic effects.
- Avoiding additivity implies introducing interaction terms between predictors. Moreover, it diminishes **interpretability** of the model. For example:

$$Y = f(\mathbf{X}) + \epsilon = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2 + \epsilon$$

- The effect of  $X_1$  is no longer constant, as changes in  $X_1$  modify the impact of  $X_1$  over  $Y$

$$Y = f(\mathbf{X}) + \epsilon = \beta_0 + (\beta_1 + \beta_3 X_2) X_1 + \beta_2 X_2 + \epsilon$$

# Non-linear regression

- We skip the need for linear relations
- Certain functions  $b_k(\cdot)$  acting over  $X$  in a non-linear way are chosen

$$Y \sim \beta_0 + \beta_1 b_1(X) + \dots + \beta_k b_k(X)$$

- Notice that  $X$  could be a single predictor (simple) or several ones (multiple)
- In case of multiple, the functions would affect matrixally by columns, being even possible the consideration of the  $b_k(\cdot)$  having  $p$  components, one per predictor
- For instance, simple quadratic regression, which is a particular case of polynomial regression, is given by

$$Y \sim \beta_0 + \beta_1 X + \dots + \beta_2 X^2$$

- The base functions would be in this case:  $b_1(x) = x$  y  $b_2(x) = x^2$ .

# Polynomial regression

- The regression function has the shape

$$Y \sim \beta_0 + \beta_1 X + \beta_2 X^2 + \dots + \beta_d X^d$$

- Base functions would be in this case:  $b_k(x) = x^k, k = 0, \dots, d$ .
- Polynomial curves (simple) or surfaces (multiple) of degree  $d$  are generated.



- Instead of a polynomial of high degree for the whole prediction space, it is possible that we could be interested on using locally different polynomials of lower degrees in certain concrete regions
- Those regions are defined by mean of cut-points called *knots* that indicate the points where the model changes
- For instance, a cubic polynomial in two pieces with a knot in a point  $c$  would be:

$$Y = \begin{cases} \beta_{01} + \beta_{11}X + \beta_{21}X^2 + \beta_{31}X^3, si\ X < c \\ \beta_{02} + \beta_{12}X + \beta_{22}X^2 + \beta_{32}X^3, si\ X \geq c \end{cases}$$

# Cubic spline regression

$$Y \sim \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3 + \beta_{3+1} h(X, c_1) + \dots + \beta_{3+k} h(X, c_k)$$

- where

$$h(X, c) = \begin{cases} (X - c)^3, & \text{si } X > c \\ 0 & , \text{si } X \leq c \end{cases}$$

- Defines piecewise polynomials of degree 3 that is continuous and with smooth curvature
- We could enforce certain conditions at the knots: continuity, or first order differentiability and/or second order continuous derivatives, etc.
- The total number of coefficients is: 1 ( $\beta_0$ ) + 3 (degree) + k (knots)

# Cubic spline regression

- Cubic splines are the ones with lower order so that the discontinuity at the knots is not perceptible by human eye
- Splines vs polynomial. In polynomial regression we need really high degrees in many problems. This causes issues in low represented zones of the space as well as in extreme extrapolation areas
- Splines vs polynomial. The flexibility in splines is achieved by incrementing the number of knots, keeping the degree fixed

# Fitting procedures

- Fitting a multiple linear regression model

$$Y = f(\mathbf{X}) + \epsilon = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \epsilon$$

is an optimization problem

$$(\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p) = \arg \min_{\beta} \sum_{i=1}^N e_i^2 \rightarrow \hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \dots + \hat{\beta}_p X_p$$

- Instead of optimizing only the error, we could modify the objective function in order to look for fitting alternatives.
- Why?
  - Getting better fitting, i.e. lower error
  - Improving interpretability, e.g. reducing the use of certain variables internally (embedded feature selection)

- Available alternatives:
  - **Feature subset selection.** Identify those variables we consider more related with the response variable
  - **Dimensionality reduction.** Projecting the  $p$  predictors over a  $M$ -dimensional space, with  $M < p$ , using the projections as the new predictors (feature extraction). An example would be the application of PCA, previously to the multiple linear regression. This algorithm is called *principal components regression* (PCR).
  - **Regularization (or *shrinkage*).** Fitting the model with all predictors, forcing, during the optimization, that the absolute values of the coefficients become small, even null. For instance, Lasso.

# Principal components regression

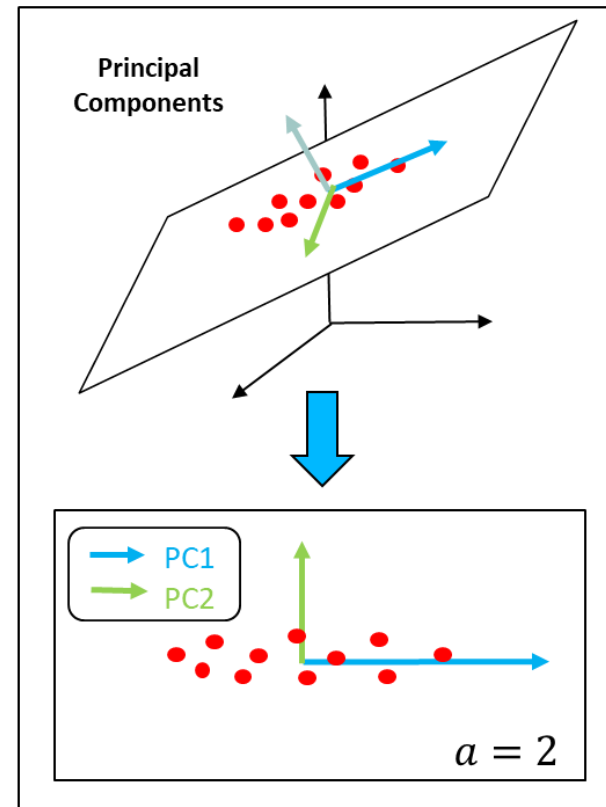
- We apply PCA, fixing the right number,  $a$ , of PCs
- We project on the  $a$ -dimensional space, using the scores as new predictors
- We apply MLR using them

## Principal component analysis

- Based on the decomposition given by:

$$X = T_a P_a^T + E_a \quad a \equiv \#PCs$$

- The PCs (LVs) are selected by maximizing scores variance
- $P \equiv$  Loadings  $\Rightarrow$  Coefficients of the linear combination that defines the PCs
- $T \equiv$  Scores  $\Rightarrow$  Coordinates of the data in the new axes
- $E \equiv$  Residuals (null if  $a = N$ )



# Partial least squares

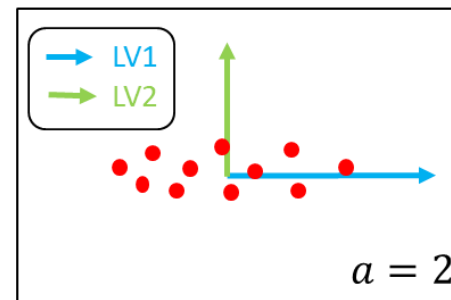
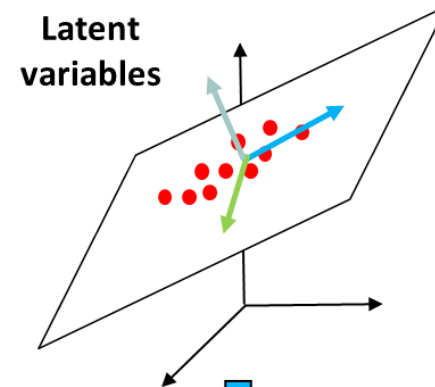
- We apply PLS, choosing the right number,  $a$ , of latent variables (LVs)
- The prediction now is obtained directly, without the need of MLR, using

$$Y = T_a Q_a$$

## Partial least squares

$$\begin{aligned}
 X &= T_a P_a^T + E_a \\
 Y &= T_a Q_a + F_a \\
 a &\equiv \#LVs
 \end{aligned}$$

- The LVs are selected by maximizing scores variance and correlation with the target
- $P, Q \equiv$  Loadings
- $T \equiv$  Scores
- $E, F \equiv$  Residuals
- Supervised and linear



# Lasso & Elastic Net

- We modify the objective function of the multiple linear regression model

$$Y = f(\mathbf{X}) + \epsilon = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \epsilon ,$$

appearing in the optimization problem

$$(\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p) = \arg \min_{\beta} \sum_{i=1}^N e_i^2 \rightarrow \hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \dots + \hat{\beta}_p X_p.$$

- The new objective function is not only considering the quadratic errors  $e_i^2$ , but it also includes a penalization term on the absolute values of the regression coefficients

$$\hat{\beta} = \arg \min_{\beta} \sum_{i=1}^N e_i^2 + \alpha \sum_{j=1}^p |\beta_j| \rightarrow \hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \dots + \hat{\beta}_p X_p.$$

- The penalization term uses the  $l_1$  norm of the coefficients vector. A similar penalization, but employing the  $l_2$  norm would be the case of **ridge regression**.
- The value  $\alpha$  controls the importance of the penalization. High values would force some coefficients to become null.
- A penalization term combining  $l_1$  and  $l_2$  norms by means of a convex combination is called **Elastic Net**





**Mondragon  
Unibertsitatea**

Escuela Politécnica  
Superior

# **Validation and model selection**

# Explained variance

- Let  $y$  and  $\hat{y}$  be real and predicted values respectively
- The **explained variance** is:

$$\text{ExplainedVariance}(y, \hat{y}) = 1 - \frac{\text{Var}(y - \hat{y})}{\text{Var}(y)}$$

- The best possible value, without error, would be 1 (maximum)
- In *Python*, we would use **explained\_variance\_score** function, available in the part of scikit-learn concerning *scores* (metrics).
- Concretely, it is imported with: **from sklearn.metrics import explained\_variance\_score**
- It is calculated using the prediction by: **explained\_variance\_score(y\_true, y\_pred)**

# Mean absolute error

- Let  $y$  and  $\hat{y}$  be real and predicted values respectively
- The **mean absolute error** is:

$$MAE(y, \hat{y}) = \frac{1}{N} \sum_{i=0}^{N-1} |y_i - \hat{y}_i|$$

- The best possible value, without error, would be 0 (minimum)
- In *Python*, we would use **mean\_absolute\_error** function, available in the part of scikit-learn concerning *scores* (metrics).
- Concretely, it is imported with: **from sklearn.metrics import mean\_absolute\_error**
- It is calculated using the prediction by: **mean\_absolute\_error(y\_true, y\_pred)**

# Mean squared error

- Let  $y$  and  $\hat{y}$  be real and predicted values respectively
- The **mean squared error** is:

$$MSE(y, \hat{y}) = \frac{1}{N} \sum_{i=0}^{N-1} (y_i - \hat{y}_i)^2$$

- The best possible value, without error, would be 0 (minimum)
- It is a quadratic measure. If we want an error value in a similar range as MAE, we must take its square root (denoted by RMSE, from *root mean squared error*).
- Regarding its minimization, it is irrelevant as the maxima and minima of a function  $f(x)$  are located at the same points as the ones of  $g(x) = \sqrt{f(x)}$
- In *Python*, we would use **mean\_squared\_error** function, available in the part of scikit-learn concerning scores (metrics).
- Concretely, it is imported with: **from sklearn.metrics import mean\_squared\_error**
- It is calculated using the prediction by: **mean\_squared\_error(y\_true, y\_pred)**

# Mean squared logarithmic error

- Let  $y$  and  $\hat{y}$  be real and predicted values respectively
- The **mean squared logarithmic error** is:

$$MLSE(y, \hat{y}) = \frac{1}{N} \sum_{i=0}^{N-1} (\ln(1 + y_i) - \ln(1 + \hat{y}_i))^2$$

- The best possible value, without error, would be 0 (minimum)
- Due to the shape of the logarithmic function, errors by default are penalized harder than errors by excess
- In *Python*, we would use **mean\_squared\_log\_error** function, available in the part of scikit-learn concerning scores (metrics).
- Concretely, it is imported with: **from sklearn.metrics import mean\_squared\_log\_error**
- It is calculated using the prediction by: **mean\_squared\_log\_error(y\_true, y\_pred)**

# Coefficient of determination $R^2$

- Let  $y$  and  $\hat{y}$  be real and predicted values respectively, and let  $\bar{y}$  be the average of real values
- The **coefficient of determination** is:

$$R^2(y, \hat{y}) = 1 - \frac{\sum_{i=0}^{N-1} (y_i - \hat{y}_i)^2}{\sum_{i=0}^{N-1} (y_i - \bar{y})^2}$$

- The best possible value, without error, would be 1 (maximum)
- It could take values arbitrarily negative
- In case of predicting always the mean, we would obtain a value 0
- In *Python*, we would use **r2\_score** function, available in the part of scikit-learn concerning *scores* (metrics).
- Concretely, it is imported with: **from sklearn.metrics import r2\_score**
- It is calculated using the prediction by: **r2\_score(y\_true, y\_pred)**



**Mondragon  
Unibertsitatea**

Escuela Politécnica  
Superior

Eskerrik asko  
Muchas gracias  
Thank you

**Carlos Cernuda**

ccernuda@mondragon.edu

MGEP

Goiru, 2

20500 Arrasate – Mondragon

Tlf. 662420414