



Lisbon School  
of Economics  
& Management  
Universidade de Lisboa



Carlos J. Costa

# REGRESSIONS

# Regression

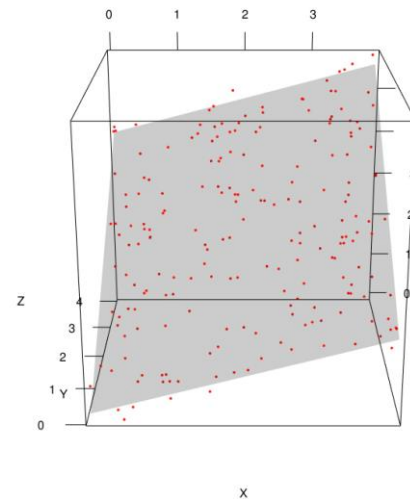
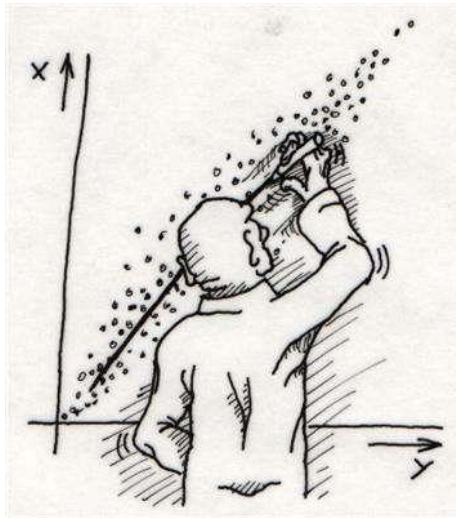
- Statistical processes for estimating the relationships among variables.
- Dependent variable, outcome variable, target
- Independent variables, predictor, covariates, or features

# Regression

- simple regression/multivariate regression

$$Y_i = \beta_0 + \beta_1 X_i + e_i$$

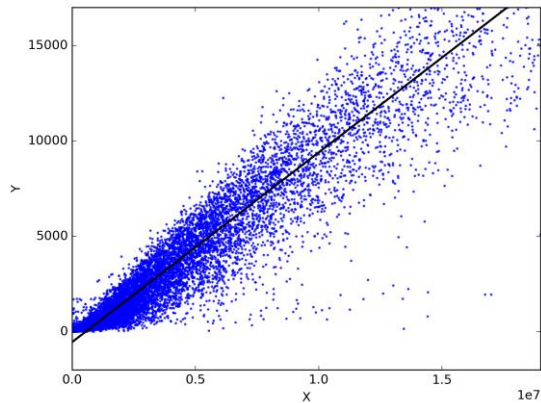
$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + e_i.$$



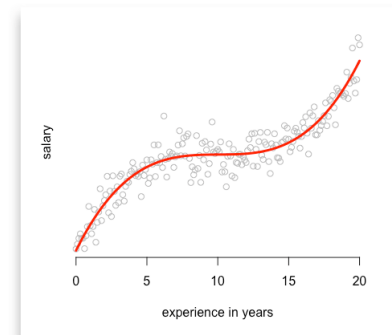
# Regression

- .Linear/non linear

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad i = 1, \dots, n.$$



$$y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \varepsilon_i, \quad i = 1, \dots, n.$$



# Regression



### CURVE-FITTING METHODS AND THE MESSAGES THEY SEND

<p><b>LINEAR</b></p> <p>"HEY, I DID A REGRESSION."</p>	<p><b>QUADRATIC</b></p> <p>"I WANTED A CURVED LINE, SO I MADE ONE WITH MATH."</p>	<p><b>LOGARITHMIC</b></p> <p>"LOOK, IT'S TAPERING OFF!"</p>
<p><b>EXPONENTIAL</b></p> <p>"LOOK, IT'S GROWING UNCONTROLLABLY!"</p>	<p><b>LOESS</b></p> <p>"I'M SOPHISTICATED, NOT LIKE THOSE BUMBLING POLYNOMIAL PEOPLE."</p>	<p><b>LINEAR, NO SLOPE</b></p> <p>"I'M MAKING A SCATTER PLOT BUT I DON'T WANT TO."</p>
<p><b>LOGISTIC</b></p> <p>"I NEED TO CONNECT THESE TWO LINES, BUT MY FIRST IDEA DIDN'T HAVE ENOUGH MATH."</p>	<p><b>CONFIDENCE INTERVAL</b></p> <p>"LISTEN, SCIENCE IS HARD. BUT I'M A SERIOUS PERSON DOING MY BEST."</p>	<p><b>PIECEWISE</b></p> <p>"I HAVE A THEORY, AND THIS IS THE ONLY DATA I COULD FIND."</p>
<p><b>CONNECTING LINES</b></p> <p>"I CLICKED 'SMOOTH LINES' IN EXCEL."</p>	<p><b>AD-HOC FILTER</b></p> <p>"I HAD AN IDEA FOR HOW TO CLEAN UP THE DATA. WHAT DO YOU THINK?"</p>	<p><b>HOUSE OF CARDS</b></p> <p>"AS YOU CAN SEE, THIS MODEL SMOOTHLY FITS THE— WAIT NO NO DON'T EXTEND IT AAAAAA!!!"</p>

# Regression

- OLS

$$\sum_{i=1}^M (y_i - \hat{y}_i)^2 = \sum_{i=1}^M \left( y_i - \sum_{j=0}^p w_j \times x_{ij} \right)^2$$

- Ridge

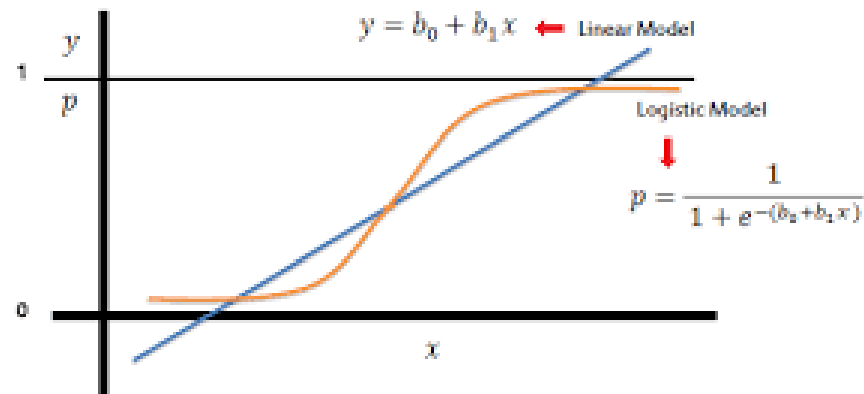
$$\sum_{i=1}^M (y_i - \hat{y}_i)^2 = \sum_{i=1}^M \left( y_i - \sum_{j=0}^p w_j \times x_{ij} \right)^2 + \lambda \sum_{j=0}^p w_j^2$$

- Lasso

$$\sum_{i=1}^M (y_i - \hat{y}_i)^2 = \sum_{i=1}^M \left( y_i - \sum_{j=0}^p w_j \times x_{ij} \right)^2 + \lambda \sum_{j=0}^p |w_j|$$

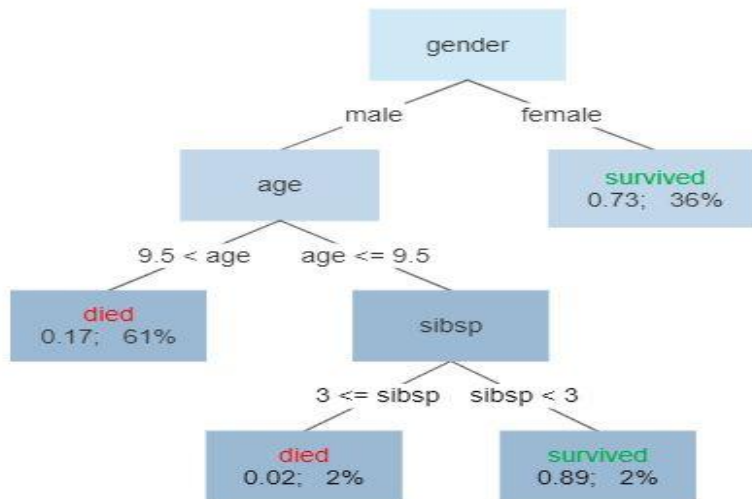
# Logistics Regression

- Logistic regression is a statistical model that in its basic form uses a logistic function to model a binary dependent variable, although many more complex extensions exist



# Decision Tree

## Survival of passengers on the Titanic



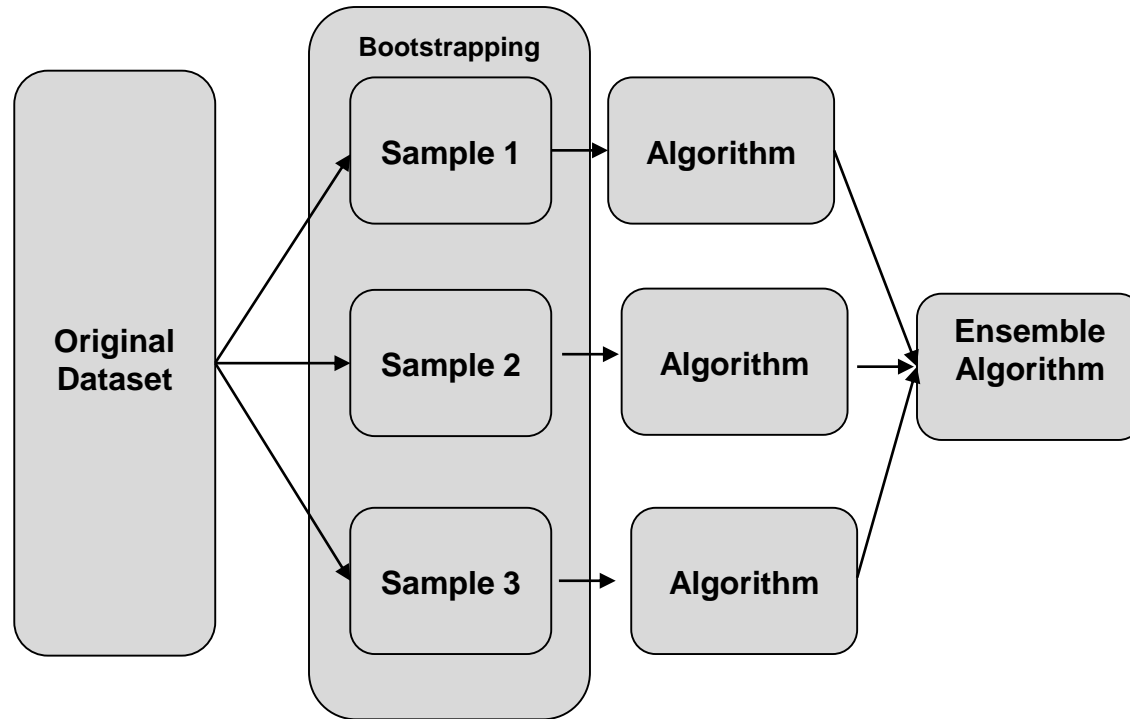
- Classification or Regression
- breaks down a data set into smaller and smaller subsets
- final result is a tree with decision nodes and leaf nodes.



# Ensemble

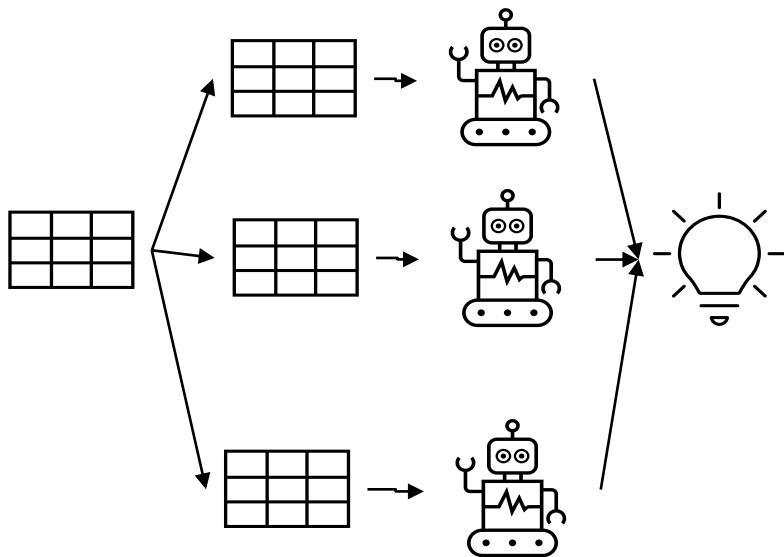
- Ensemble is a Machine Learning concept in which the idea is to train multiple models using the same learning algorithm.
- classification, regression and other tasks
- multitude of decision trees at training time
- outputting the class that is the mode of the classes (classification) or mean prediction (regression) of the individual trees

# Ensemble

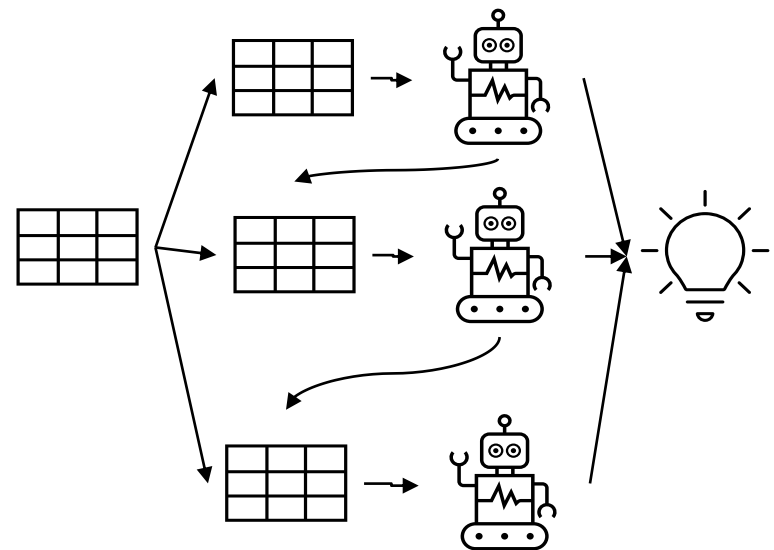


# Bagging vs. Boosting

Bagging



Boosting



# Bagging vs. Boosting

- classification, regression and other tasks
- multitude of decision trees at training time
- outputting the class that is the mode of the classes (classification) or mean prediction (regression) of the individual trees.

# Bagging

- **create multiple bootstrap samples**
- **fit a weak learner**
- **aggregate** -> “average” their
- outputs is an ensemble model with less variance than its components..

$$s_L(\cdot) = \frac{1}{L} \sum_{l=1}^L w_l(\cdot) \quad (\text{simple average, for regression problem})$$

$$s_L(\cdot) = \arg \max_k [\text{card}(l | w_l(\cdot) = k)] \quad (\text{simple majority vote, for classification problem})$$

# Boosting

- in fitting sequentially multiple weak learners in a very adaptative way
- each new model focus its efforts on the most difficult observations to fit up to now
- at the end of the process, is obtained a strong learner with lower bias
- Boosting can also have the effect of reducing variance

$$(c_l, w_l(\cdot)) = \arg \min_{c, w(\cdot)} E(s_{l-1}(\cdot) + c \times w(\cdot)) = \arg \min_{c, w(\cdot)} \sum_{n=1}^N e(y_n, s_{l-1}(x_n) + c \times w(x_n))$$

# Random Forest

- is a bagging method where deep trees, fitted on bootstrap samples, are combined to produce an output with lower variance
- classification, regression and other tasks
- multitude of decision trees at training time
- outputting the class that is the mode of the classes (classification) or mean prediction (regression) of the individual trees.

# Python Libraries

