# CRIANDO MODELOS SIMPLES DE MACHINE LEARNING III

DIEGO RODRIGUES DSC

INFNET

# MODEL LIFECYCLE : CRIANDO MODELOS SIMPLES DE MACHINE LEARNING III

**PARTE 1 : TEORIA**

- BUSINESS UNDERSTANDING
  - REGRESSÃO

- MODELING
  - REGRESSÃO LINEAR ORDINÁRIA

- EVALUATION
  - $R^2$
  - ESTATÍSTICAS DO MODELO
  - ANÁLISE GRÁFICA DO RESULTADO

- MODELING ++
  - MODELO LOG LINEAR
  - REGULARIZAÇÃO
    - REGRESSÃO LINEAR PONDERADA
    - RIDGE, LASSO, ELASTIC NET

**PARTE 2 : PRÁTICA**

- NOTEBOOK REGRESSÃO IRIS

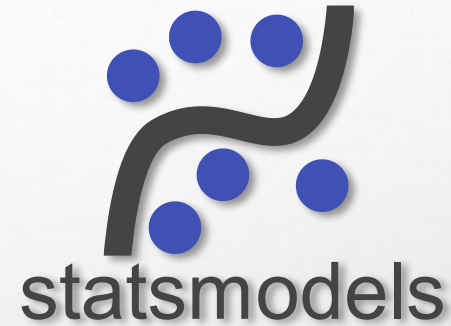Produzir Ação

# CICLO DE VIDA DO MODELO

Baseado em Dados

Cross Industry Standard Process for Data Mining - IBM



**1) Requerimentos e Análise de Negócio**

Entendimento do problema decisório, dados relacionados & revisão bibliográfica.

**2) Preparação dos Dados**

Entendimento das fontes de dados, dos tipos e elaboração da representação.

**3) Modelagem**

Análise Exploratória, Seleção de atributos e treinamento.

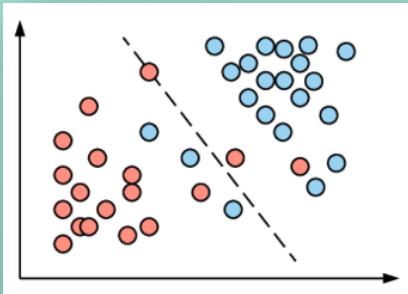**4) Avaliação**

Seleção do melhor modelo.

**5) Liberação**

Liberação do modelo no ambiente de produção.

# BUSINESS UNDERSTANDING

# REGRESSÃO

O objetivo da regressão é

**modelar as relações funcionais**

entre dois conjuntos de variáveis.



As vezes quando o mundo não é linear & gaussiano...

As variáveis que representam as causas são chamadas de **variáveis independentes**, e as variáveis cujo objetivo é prever, são chamadas **variáveis dependentes**.

Então, uma **regressão** é um modelo utilizado para prever **uma ou mais variáveis dependentes**, baseado em causas, ou variáveis independentes.

# MODELOS DE REGRESSÃO

1) **Regressão Linear**

2) Regressão Não-Linear

3) Processos Gaussianos

4) Máquina de Vetores Suporte

5) Redes Neurais



Algoritmos de regressão geralmente são modelados combinando uma **parte determinística e uma parte aleatória.** Os parâmetros correspondente à parte determinística são encontrados utilizando estimadores como máxima verossimilhança ou máximo a posteriori (MAP).

# MODELOS DE REGRESSÃO

$$Y = F(X) + \varepsilon$$

Parte Determinística    Parte Estocástica

$$Y = \alpha^T x + \varepsilon$$

$$Y = X\alpha + \varepsilon$$

LOGISTIC REGRESSION

$$Y = \frac{1}{1 + e^{\alpha^t x + \varepsilon}}$$

Input layer    Hidden layer    Output layer

$$Y = \varphi(x) + \varepsilon$$

Regressor Comparison

Minimal regularization — Lots of regularization

# MODELING

# Regressão Linear : Modelo Matemático

## Formulation [ edit ]

Given a data set $\{y_i, x_{i1}, \ldots, x_{ip}\}_{i=1}^{n}$ of $n$ statistical units, a linear regression model assumes that the relationship between the dependent variable $y$ and the vector of regressors $\mathbf{x}$ is linear. This relationship is modeled through a *disturbance term* or *error variable* $\varepsilon$ — an unobserved random variable that adds "noise" to the linear relationship between the dependent variable and regressors. Thus the model takes the form

$$y_i = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip} + \varepsilon_i = \mathbf{x}_i^\mathsf{T}\boldsymbol{\beta} + \varepsilon_i, \qquad i = 1, \ldots, n,$$

where $\mathsf{T}$ denotes the transpose, so that $\mathbf{x}_i^\mathsf{T}\boldsymbol{\beta}$ is the inner product between vectors $\mathbf{x}_i$ and $\boldsymbol{\beta}$.

Often these $n$ equations are stacked together and written in matrix notation as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon},$$

where

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix},$$

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^\mathsf{T} \\ \mathbf{x}_2^\mathsf{T} \\ \vdots \\ \mathbf{x}_n^\mathsf{T} \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & \cdots & x_{1p} \\ 1 & x_{21} & \cdots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \cdots & x_{np} \end{bmatrix},$$

$$\boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{bmatrix}, \quad \boldsymbol{\varepsilon} = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}.$$

In linear regression, the observations (**red**) are assumed to be the result of random deviations (**green**) from an underlying relationship (**blue**) between a dependent variable ($y$) and an independent variable ($x$).

$$y = \sum_i^V \beta_i x_i + \varepsilon$$

# Exemplo I: Altura e Peso

## Simple Linear Regression

Regression analysis makes use of mathematical models to describe relationships. For example, suppose that height was the only determinant of body weight. If we were to plot height (the independent or 'predictor' variable) as a function of body weight (the dependent or 'outcome' variable), we might see a very linear relationship, as illustrated below.



$$Y = a + b\,X$$
$$wgt = 80 + 2\,(hgt)$$

X-axis: Height (inches)

We could also describe this relationship with the equation for a line, Y = a + b(x), where 'a' is the Y-intercept and 'b' is the slope of the line. We could use the equation to predict weight if we knew an individual's height. In this example, if an individual was 70 inches tall, we would predict his weight to be:

$$Weight = 80 + 2 \times (70) = 220 \text{ lbs.}$$

In this simple linear regression, we are examining the impact of one independent variable on the outcome. If height were the only determinant of body weight, we would expect that the points for individual subjects would lie close to the line. However, if there were other factors (independent variables) that influenced body weight besides height (e.g., age, calorie intake, and exercise level), we might expect that the points for individual subjects would be more loosely scattered around the line, since we are only taking height into account.

# Premissas I

## Assumptions [ edit ]

*See also: Ordinary least squares § Assumptions*

Standard linear regression models with standard estimation techniques make a number of assumptions about the predictor variables, the response variables and their relationship. Numerous extensions have been developed that allow each of these assumptions to be relaxed (i.e. reduced to a weaker form), and in some cases eliminated entirely. Generally these extensions make the estimation procedure more complex and time-consuming, and may also require more data in order to produce an equally precise model.
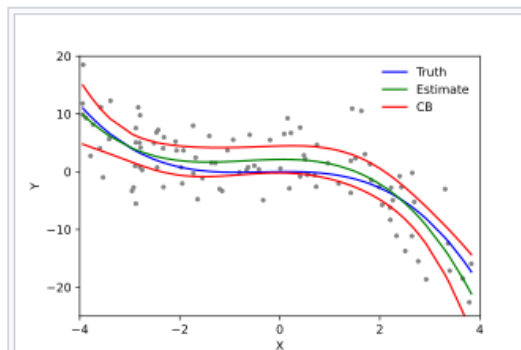
The following are the major assumptions made by standard linear regression models with standard estimation techniques (e.g. ordinary least squares):

- **Weak exogeneity**. This essentially means that the predictor variables $x$ can be treated as fixed values, rather than random variables. This means, for example, that the predictor variables are assumed to be error-free—that is, not contaminated with measurement errors. Although this assumption is not realistic in many settings, dropping it leads to significantly more difficult errors-in-variables models.
- **Linearity**. This means that the mean of the response variable is a linear combination of the parameters (regression coefficients) and the predictor variables. Note that this assumption is much less restrictive than it may at first seem. Because the predictor variables are treated as fixed values (see above), linearity is really only a restriction on the parameters. The predictor variables themselves can be arbitrarily transformed, and in fact multiple copies of the same underlying predictor variable can be added, each one transformed differently. This technique is used, for example, in polynomial regression, which uses linear regression to fit the response variable as an arbitrary polynomial function (up to a given degree) of a predictor variable. With this much flexibility, models such as polynomial regression often have "too much power", in that they tend to overfit the data. As a result, some kind of regularization must typically be used to prevent unreasonable solutions coming out of the estimation process. Common examples are ridge regression and lasso regression. Bayesian linear regression can also be used, which by its nature is more or less immune to the problem of overfitting. (In fact, ridge regression and lasso regression can both be viewed as special cases of Bayesian linear regression, with particular types of prior distributions placed on the regression coefficients.)



Example of a cubic polynomial regression, which is a type of linear regression. Although *polynomial regression* fits a nonlinear model to the data, as a statistical estimation problem it is linear, in the sense that the regression function $E(y \mid x)$ is linear in the unknown parameters that are estimated from the data. For this reason, polynomial regression is considered to be a special case of multiple linear regression.

- Additivity: $f(x + y) = f(x) + f(y)$.
- Homogeneity of degree 1: $f(\alpha x) = \alpha\, f(x)$ for all $\alpha$.

# Premissas II

- **Constant variance** (a.k.a. **homoscedasticity**). This means that the variance of the errors does not depend on the values of the predictor variables. Thus the variability of the responses for given fixed values of the predictors is the same regardless of how large or small the responses are. This is often not the case, as a variable whose mean is large will typically have a greater variance than one whose mean is small. For example, a person whose income is predicted to be $100,000 may easily have an actual income of $80,000 or $120,000—i.e., a **standard deviation** of around $20,000—while another person with a predicted income of $10,000 is unlikely to have the same $20,000 standard deviation, since that would imply their actual income could vary anywhere between −$10,000 and $30,000. (In fact, as this shows, in many cases—often the same cases where the assumption of normally distributed errors fails—the variance or standard deviation should be predicted to be proportional to the mean, rather than constant.) The absence of homoscedasticity is called **heteroscedasticity**. In order to check this assumption, a plot of residuals versus predicted values (or the values of each individual predictor) can be examined for a "fanning effect" (i.e., increasing or decreasing vertical spread as one moves left to right on the plot). A plot of the absolute or squared residuals versus the predicted values (or each predictor) can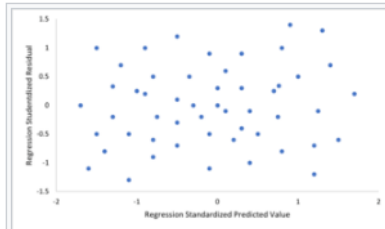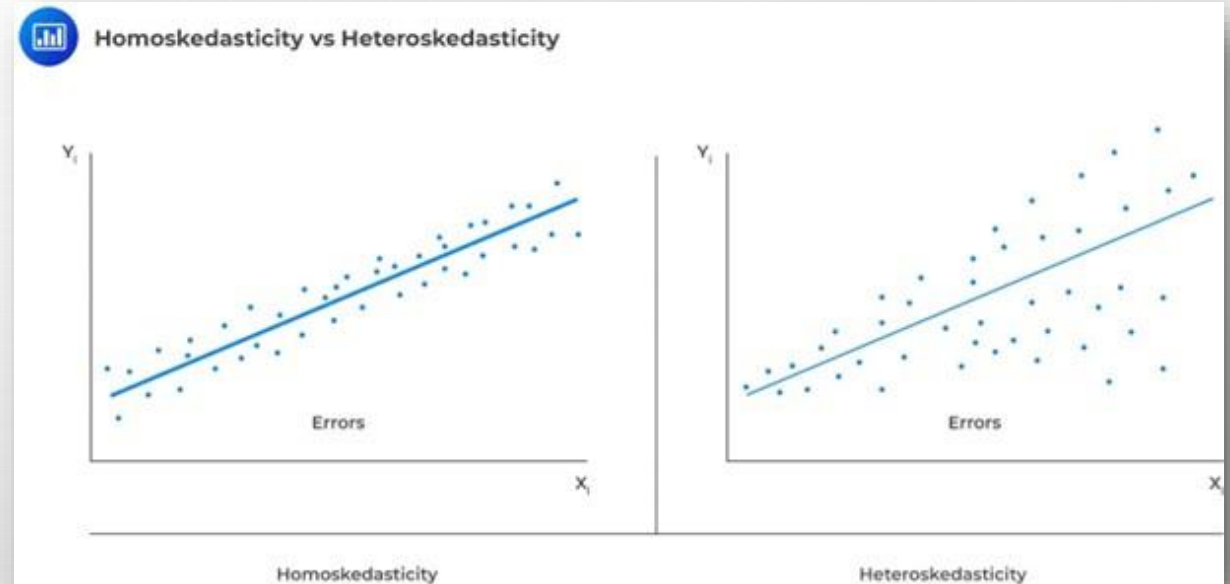 also be examined for a trend or curvature. Formal tests can also be used; see **Heteroscedasticity**. The presence of heteroscedasticity will result in an overall "average" estimate of variance being used instead of one that takes into account the true variance structure. This leads to less precise (but in the case of **ordinary least squares**, not biased) parameter estimates and biased standard errors, resulting in misleading tests and interval estimates. The **mean squared error** for the model will also be wrong. Various estimation techniques including **weighted least squares** and the use of **heteroscedasticity-consistent standard errors** can handle heteroscedasticity in a quite general way. **Bayesian linear regression** techniques can also be used when the variance is assumed to be a function of the mean. It is also possible in some cases to fix the problem by applying a transformation to the response variable (e.g., fitting the **logarithm** of the response variable using a linear regression model, which implies that the response variable itself has a **log-normal distribution** rather than a **normal distribution**).

- **Independence of errors**. This assumes that the errors of the response variables are uncorrelated with each other. (Actual **statistical independence** is a stronger condition than mere lack of correlation and is often not needed, although it can be exploited if it is known to hold.) Some methods such as **generalized least squares** are capable of handling correlated errors, although they typically require significantly more data unless some sort of **regularization** is used to bias the model towards assuming uncorrelated errors. **Bayesian linear regression** is a general way of handling this issue.

- **Lack of perfect multicollinearity** in the predictors. For standard **least squares** estimation methods, the design matrix $X$ must have full **column rank** $p$; otherwise perfect **multicollinearity** exists in the predictor variables, meaning a linear relationship exists between two or more predictor variables. This can be caused by accidentally duplicating a variable in the data, using a linear transformation of a variable along with the original (e.g., the same temperature measurements expressed in Fahrenheit and Celsius), or including a linear combination of multiple variables in the model, such as their mean. It can also happen if there is too little data available compared to the number of parameters to be estimated (e.g., fewer data points than regression coefficients). Near violations of this assumption, where predictors are highly but not perfectly correlated, can reduce the precision of parameter estimates (see **Variance inflation factor**). In the case of perfect multicollinearity, the parameter vector $\beta$ will be **non-identifiable**—it has no unique solution. In such a case, only some of the parameters can be identified (i.e., their values can only be estimated within some linear subspace of the full parameter space $R^p$). See **partial least squares regression**. Methods for fitting linear models with multicollinearity have been developed,[5][6][7][8] some of which require additional assumptions such as "effect sparsity"—that a large fraction of the effects are exactly zero. Note that the more computationally expensive iterated algorithms for parameter estimation, such as those used in **generalized linear models**, do not suffer from this problem.



Visualization of heteroscedasticity in a scatter plot against 100 random fitted values using Matlab



To check for violations of the assumptions of linearity, constant variance, and independence of errors within a linear regression model, the residuals are typically plotted against the predicted values (or each of the individual predictors). An apparently random scatter of points about the horizontal midline at 0 is ideal, but cannot rule out certain kinds of violations such as **autocorrelation** in the errors or their correlation with one or more covariates.



Homoskedasticity vs Heteroskedasticity

# Encontrando os Coeficientes : Mínimos Quadrados Ordinários

Pseudo-inversa de Moore-Penrose

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

## Linear model [edit]

*Main article: Linear regression model*

Suppose the data consists of $n$ observations $\{\mathbf{x}_i, y_i\}_{i=1}^n$. Each observation $i$ includes a scalar response $y_i$ and a column vector $\mathbf{x}_i$ of $p$ parameters (regressors), i.e., $\mathbf{x}_i = [x_{i1}, x_{i2}, \ldots, x_{ip}]^T$. In a linear regression model, the response variable, $y_i$, is a linear function of the regressors:

$$y_i = \beta_1\, x_{i1} + \beta_2\, x_{i2} + \cdots + \beta_p\, x_{ip} + \varepsilon_i,$$

or in vector form,

$$y_i = \mathbf{x}_i^T \boldsymbol{\beta} + \varepsilon_i,$$

where $\mathbf{x}_i$, as introduced previously, is a column vector of the $i$-th observation of all the explanatory variables; $\boldsymbol{\beta}$ is a $p \times 1$ vector of unknown parameters; and the scalar $\varepsilon_i$ represents unobserved random variables (errors) of the $i$-th observation. $\varepsilon_i$ accounts for the influences upon the responses $y_i$ from sources other than the explanatory variables $\mathbf{x}_i$. This model can also be written in matrix notation as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon},$$

where $\mathbf{y}$ and $\boldsymbol{\varepsilon}$ are $n \times 1$ vectors of the response variables and the errors of the $n$ observations, and $\mathbf{X}$ is an $n \times p$ matrix of regressors, also sometimes called the design matrix, whose row $i$ is $\mathbf{x}_i^T$ and contains the $i$-th observations on all the explanatory variables.

Typically, a constant term is included in the set of regressors $\mathbf{X}$, say, by taking $x_{i1} = 1$ for all $i = 1, \ldots, n$. The coefficient $\beta_1$ corresponding to this regressor is called the *intercept*. Without the intercept, the fitted line is forced to cross the origin when $x_i = \vec{0}$.

Regressors do not have to be independent: there can be any desired relationship between the regressors (so long as it is not a linear relationship). For instance, we might suspect the response depends linearly both on a value and its square; in which case we would include one regressor whose value is just the square of another regressor. In that case, the model would be *quadratic* in the second regressor, but none-the-less is still considered a *linear* model because the model *is* still linear in the parameters ($\boldsymbol{\beta}$).

## Matrix/vector formulation [edit]

Consider an overdetermined system

$$\sum_{j=1}^p x_{ij}\beta_j = y_i, \ (i = 1, 2, \ldots, n),$$

of $n$ linear equations in $p$ unknown coefficients, $\beta_1, \beta_2, \ldots, \beta_p$, with $n > p$. This can be written in matrix form as

$$\mathbf{X}\boldsymbol{\beta} = \mathbf{y},$$

where

$$\mathbf{X} = \begin{bmatrix} X_{11} & X_{12} & \cdots & X_{1p} \\ X_{21} & X_{22} & \cdots & X_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ X_{n1} & X_{n2} & \cdots & X_{np} \end{bmatrix}, \qquad \boldsymbol{\beta} = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{bmatrix}, \qquad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}.$$

(Note: for a linear model as above, not all elements in $\mathbf{X}$ contains information on the data points. The first column is populated with ones, $X_{i1} = 1$. Only the other columns contain actual data. So here $p$ is equal to the number of regressors plus one).

Such a system usually has no exact solution, so the goal is instead to find the coefficients $\boldsymbol{\beta}$ which fit the equations "best", in the sense of solving the quadratic minimization problem

$$\hat{\boldsymbol{\beta}} = \underset{\boldsymbol{\beta}}{\arg\min}\, S(\boldsymbol{\beta}),$$

where the objective function $S$ is given by

$$S(\boldsymbol{\beta}) = \sum_{i=1}^n \left| y_i - \sum_{j=1}^p X_{ij}\beta_j \right|^2 = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2.$$

A justification for choosing this criterion is given in Properties below. This minimization problem has a unique solution, provided that the $p$ columns of the matrix $\mathbf{X}$ are linearly independent, given by solving the so-called *normal equations*:

$$(\mathbf{X}^T \mathbf{X})\, \hat{\boldsymbol{\beta}} = \mathbf{X}^T \mathbf{y}.$$

The matrix $\mathbf{X}^T \mathbf{X}$ is known as the *normal matrix* or Gram matrix and the matrix $\mathbf{X}^T \mathbf{y}$ is known as the moment matrix of regressand by regressors.[2] Finally, $\hat{\boldsymbol{\beta}}$ is the coefficient vector of the least-squares hyperplane, expressed as

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}.$$

or

$$\hat{\boldsymbol{\beta}} = \boldsymbol{\beta} + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\varepsilon}.$$

# Mínimos Quadrados Ordinários: Premissas

## Classical linear regression model   [ edit ]

The classical model focuses on the "finite sample" estimation and inference, meaning that the number of observations $n$ is fixed. This contrasts with the other approaches, which study the asymptotic behavior of OLS, and in which the number of observations is allowed to grow to infinity.

- **Correct specification**. The linear functional form must coincide with the form of the actual data-generating process.
- **Strict exogeneity**. The errors in the regression should have conditional mean zero:[16]

$$\mathrm{E}[\,\varepsilon \mid X\,] = 0.$$

The immediate consequence of the exogeneity assumption is that the errors have mean zero: $\mathrm{E}[\varepsilon] = 0$ (for the law of total expectation), and that the regressors are uncorrelated with the errors: $\mathrm{E}[X^\mathrm{T}\varepsilon] = 0$.

The exogeneity assumption is critical for the OLS theory. If it holds then the regressor variables are called *exogenous*. If it doesn't, then those regressors that are correlated with the error term are called *endogenous*,[17] and the OLS estimator becomes biased. In such case the method of instrumental variables may be used to carry out inference.

- **No linear dependence**. The regressors in $X$ must all be linearly independent. Mathematically, this means that the matrix $X$ must have full column rank almost surely:[18]

$$\Pr\big[\ \mathrm{rank}(X) = p\ \big] = 1.$$

Usually, it is also assumed that the regressors have finite moments up to at least the second moment. Then the matrix $Q_{xx} = \mathrm{E}[X^\mathrm{T}X/n]$ is finite and positive semi-definite.

When this assumption is violated the regressors are called linearly dependent or perfectly multicollinear. In such case the value of the regression coefficient $\beta$ cannot be learned, although prediction of $y$ values is still possible for new values of the regressors that lie in the same linearly dependent subspace.

- **Spherical errors**:[18]

$$\mathrm{Var}[\,\varepsilon \mid X\,] = \sigma^2 I_n,$$

where $I_n$ is the identity matrix in dimension $n$, and $\sigma^2$ is a parameter which determines the variance of each observation. This $\sigma^2$ is considered a nuisance parameter in the model, although usually it is also estimated. If this assumption is violated then the OLS estimates are still valid, but no longer efficient.

It is customary to split this assumption into two parts:

- **Homoscedasticity**: $\mathrm{E}[\,\varepsilon_i^2 \mid X\,] = \sigma^2$, which means that the error term has the same variance $\sigma^2$ in each observation. When this requirement is violated this is called heteroscedasticity, in such case a more efficient estimator would be weighted least squares. If the errors have infinite variance then the OLS estimates will also have infinite variance (although by the law of large numbers they will nonetheless tend toward the true values so long as the errors have zero mean). In this case, robust estimation techniques are recommended.
  - **No autocorrelation**: the errors are uncorrelated between observations: $\mathrm{E}[\,\varepsilon_i\varepsilon_j \mid X\,] = 0$ for $i \neq j$. This assumption may be violated in the context of time series data, panel data, cluster samples, hierarchical data, repeated measures data, longitudinal data, and other data with dependencies. In such cases generalized least squares provides a better alternative than the OLS. Another expression for autocorrelation is *serial correlation*.
- **Normality**. It is sometimes additionally assumed that the errors have normal distribution conditional on the regressors:[19]

$$\varepsilon \mid X \sim \mathcal{N}(0, \sigma^2 I_n).$$

This assumption is not needed for the validity of the OLS method, although certain additional finite-sample properties can be established in case when it does (especially in the area of hypotheses testing). Also when the errors are normal, the OLS estimator is equivalent to the maximum likelihood estimator (MLE), and therefore it is asymptotically efficient in the class of all regular estimators. Importantly, the normality assumption applies only to the error terms; contrary to a popular misconception, the response (dependent) variable is not required to be normally distributed.[20]
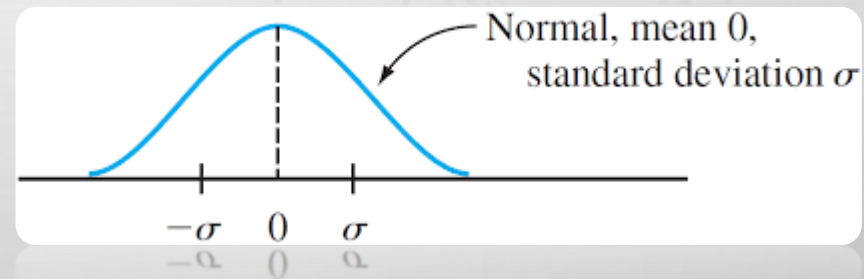
# EVALUATION

# FIGURAS DE MÉRITO - REGRESSÃO

- R QUADRADO

$$R^2 = 1 - \frac{SS_{RES}}{SS_{TOT}} = 1 - \frac{\sum_i (y_i - \hat{y}_i)^2}{\sum_i (y_i - \overline{y})^2}$$

- RESÍDUO NORMAL DE MÉDIA

ZERO E VARIÂNCIA CONSTANTE



Normal, mean 0, standard deviation $\sigma$

$-\sigma$   0   $\sigma$

# VALIDAÇÃO : STATSMODELS

Coeficiente de Determinação $R^2$

P Valor da Estatística F

P Valor dos Coeficientes

Número de Condicionamento

```
                          OLS Regression Results
==============================================================================
Dep. Variable:          petal_length   R-squared:                       0.968
Model:                            OLS   Adj. R-squared:                  0.967
Method:                 Least Squares   F-statistic:                     1473.
Date:                Mon, 05 May 2025   Prob (F-statistic):           6.98e-109
Time:                        16:30:17   Log-Likelihood:                 -39.408
No. Observations:                 150   AIC:                             86.82
Df Residuals:                     146   BIC:                             98.86
Df Model:                           3
Covariance Type:            nonrobust
==============================================================================
                 coef    std err          t      P>|t|      [0.025      0.975]
------------------------------------------------------------------------------
Intercept      -0.2627      0.297     -0.883      0.379      -0.850       0.325
petal_width     1.4468      0.068     21.399      0.000       1.313       1.580
sepal_length    0.7291      0.058     12.502      0.000       0.614       0.844
sepal_width    -0.6460      0.068     -9.431      0.000      -0.781      -0.511
==============================================================================
Omnibus:                        2.520   Durbin-Watson:                   1.783
Prob(Omnibus):                  0.284   Jarque-Bera (JB):                2.391
Skew:                           0.073   Prob(JB):                        0.303
Kurtosis:                       3.601   Cond. No.                        79.3
==============================================================================

Notes:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
```
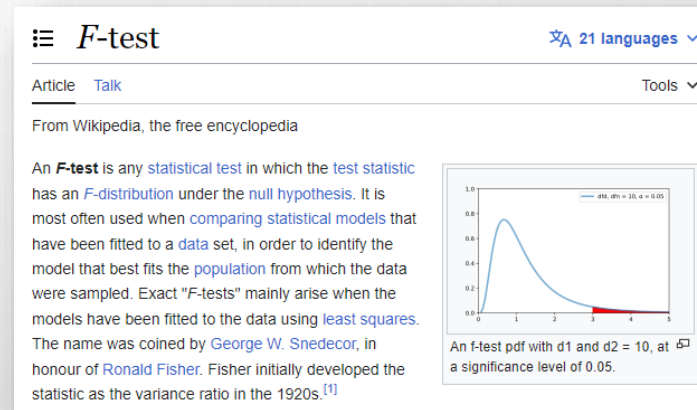
F-test

21 languages

Article  Talk                                Tools

From Wikipedia, the free encyclopedia

An **F-test** is any statistical test in which the test statistic has an F-distribution under the null hypothesis. It is most often used when comparing statistical models that have been fitted to a data set, in order to identify the model that best fits the population from which the data were sampled. Exact "F-tests" mainly arise when the models have been fitted to the data using least squares. The name was coined by George W. Snedecor, in honour of Ronald Fisher. Fisher initially developed the statistic as the variance ratio in the 1920s.[1]

An f-test pdf with d1 and d2 = 10, at a significance level of 0.05.
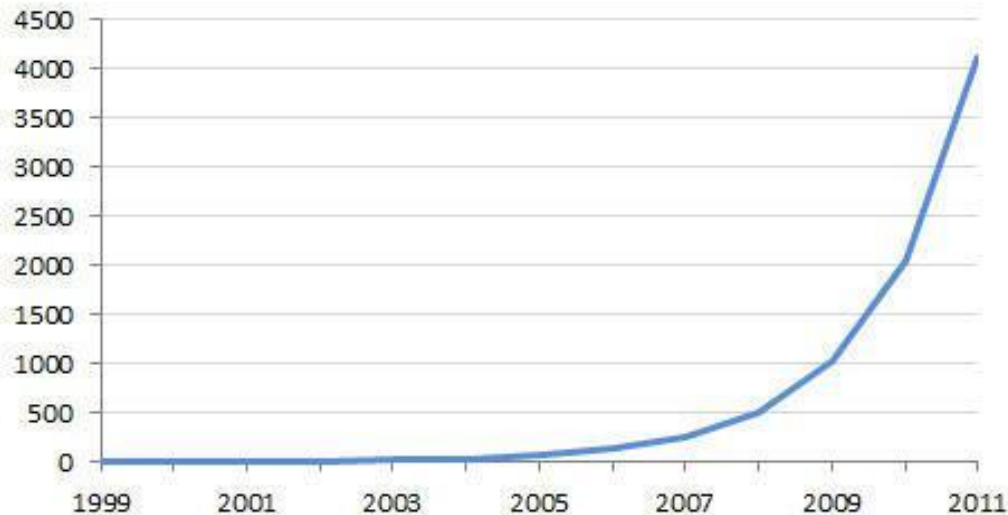
# VALIDAÇÃO : GRÁFICOS DE APOIO

# MODELING++

# MODELO LOG-LINEAR

A **log-linear model** is a mathematical model that takes the form of a function whose logarithm equals a linear combination of the parameters of the model, which makes it possible to apply (possibly multivariate) linear regression. That is, it has the general form
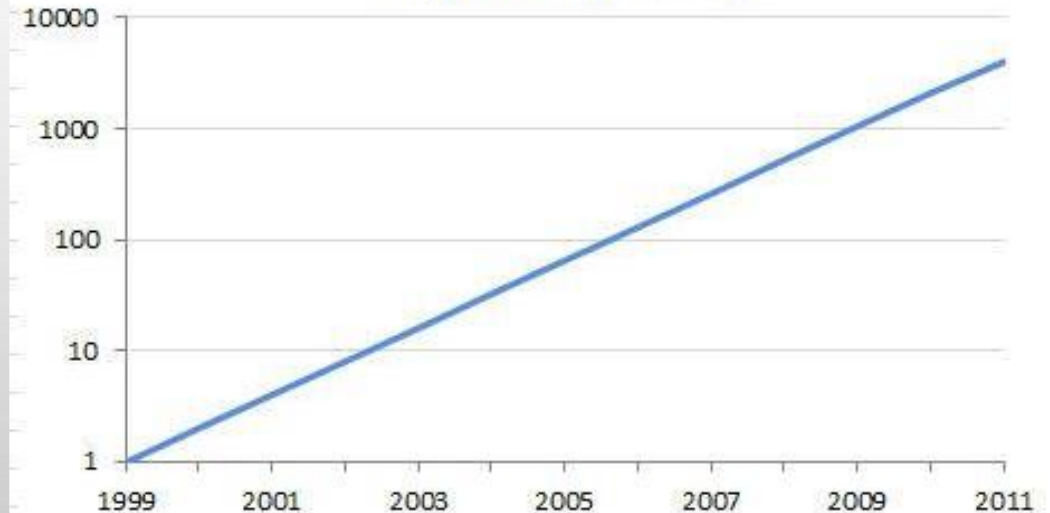
$$\exp\left(c + \sum_i w_i f_i(X)\right),$$

in which the $f_i(X)$ are quantities that are functions of the variable $X$, in general a vector of values, while $c$ and the $w_i$ stand for the model parameters.

# REGULARIZAÇÃO

In mathematics, statistics, finance,[1] and computer science, particularly in machine learning and inverse problems, **regularization** is a process that changes the result answer to be "simpler". It is often used to obtain results for ill-posed problems or to prevent overfitting.[2]
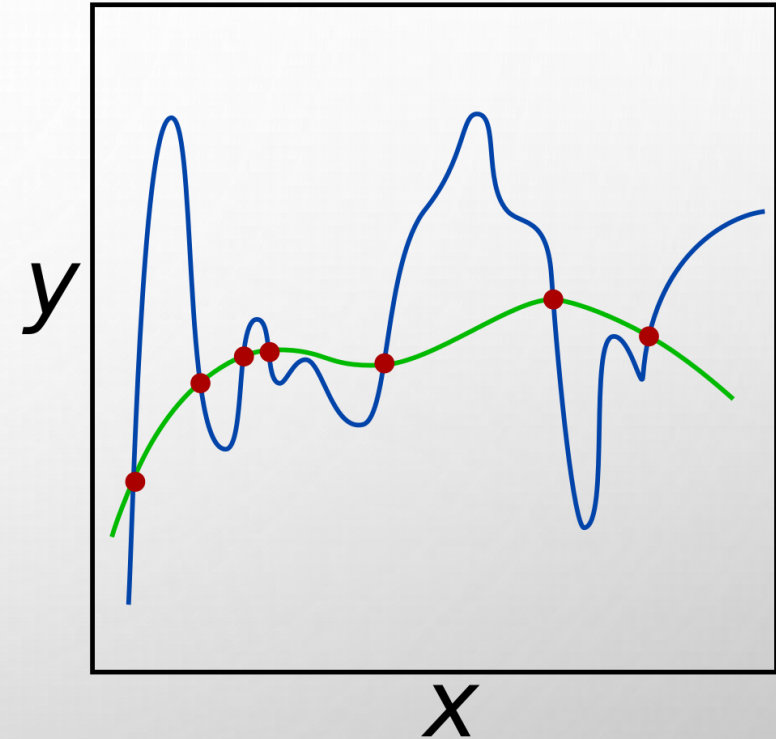
Although regularization procedures can be divided in many ways, the following delineation is particularly helpful:

- **Explicit regularization** is regularization whenever one explicitly adds a term to the optimization problem. These terms could be priors, penalties, or constraints. Explicit regularization is commonly employed with ill-posed optimization problems. The regularization term, or penalty, imposes a cost on the optimization function to make the optimal solution unique.
- **Implicit regularization** is all other forms of regularization. This includes, for example, early stopping, using a robust loss function, and discarding outliers. Implicit regularization is essentially ubiquitous in modern machine learning approaches, including stochastic gradient descent for training deep neural networks, and ensemble methods (such as random forests and gradient boosted trees).

In explicit regularization, independent of the problem or model, there is always a data term, that corresponds to a likelihood of the measurement and a regularization term that corresponds to a prior. By combining both using Bayesian statistics, one can compute a posterior, that includes both information sources and therefore stabilizes the estimation process. By trading off both objectives, one chooses to be more addictive to the data or to enforce generalization (to prevent overfitting). There is a whole research branch dealing with all possible regularizations. In practice, one usually tries a specific regularization and then figures out the probability density that corresponds to that regularization to justify the choice. It can also be physically motivated by common sense or intuition.

In machine learning, the data term corresponds to the training data and the regularization is either the choice of the model or modifications to the algorithm. It is always intended to reduce the generalization error, i.e. the error score with the trained model on the evaluation set and not the training data.[3]

One of the earliest uses of regularization is Tikhonov regularization (ridge regression), related to the method of least squares.

# MODELO MÍNIMOS QUADRADOS PONDERADO

**Weighted least squares** (WLS), also known as **weighted linear regression**,[1][2] is a generalization of ordinary least squares and linear regression in which knowledge of the unequal variance of observations (*heteroscedasticity*) is incorporated into the regression. WLS is also a specialization of generalized least squares, when all the off-diagonal entries of the covariance matrix of the errors, are null.

## Formulation [ edit ]

The fit of a model to a data point is measured by its residual, $r_i$, defined as the difference between a measured value of the dependent variable, $y_i$ and the value predicted by the model, $f(x_i, \boldsymbol{\beta})$:

$$r_i(\boldsymbol{\beta}) = y_i - f(x_i, \boldsymbol{\beta}).$$

If the errors are uncorrelated and have equal variance, then the function

$$S(\boldsymbol{\beta}) = \sum_i r_i(\boldsymbol{\beta})^2,$$

is minimised at $\hat{\boldsymbol{\beta}}$, such that $\frac{\partial S}{\partial \beta_j}(\hat{\boldsymbol{\beta}}) = 0$.

The Gauss–Markov theorem shows that, when this is so, $\hat{\boldsymbol{\beta}}$ is a best linear unbiased estimator (BLUE). If, however, the measurements are uncorrelated but have different uncertainties, a modified approach might be adopted. Aitken showed that when a weighted sum of squared residuals is minimized, $\hat{\boldsymbol{\beta}}$ is the BLUE if each weight is equal to the reciprocal of the variance of the measurement

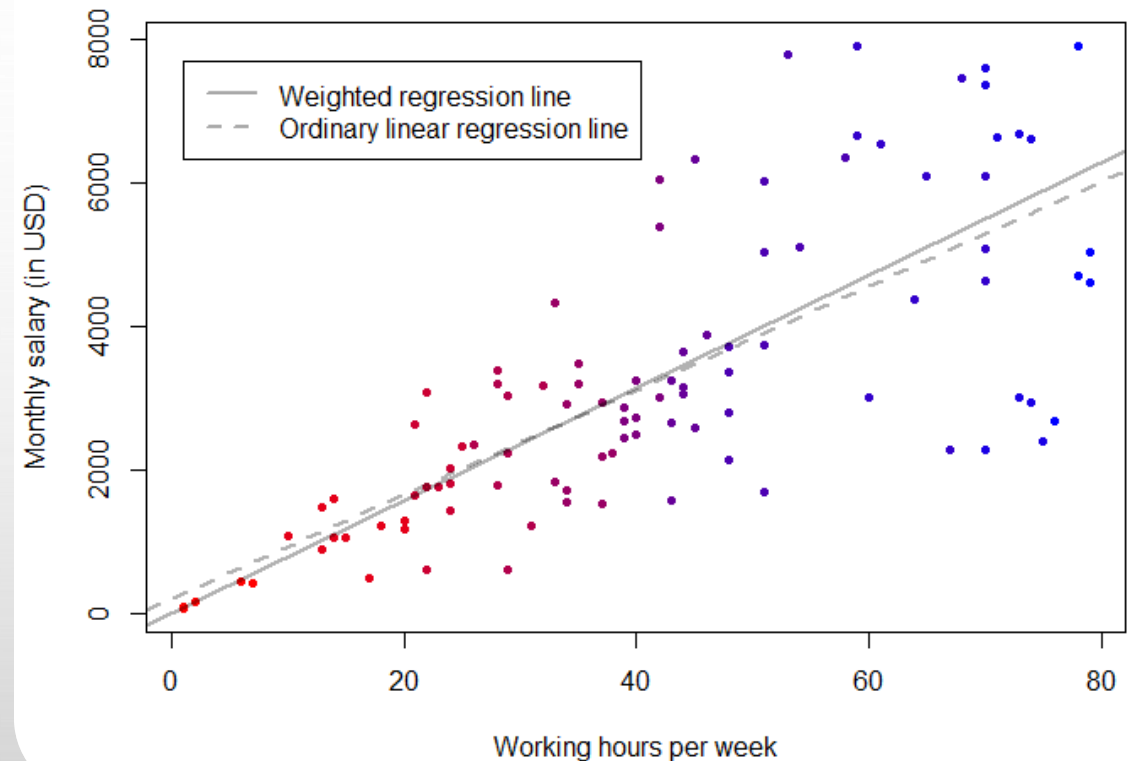$$S = \sum_{i=1}^{n} W_{ii} r_i^2, \qquad W_{ii} = \frac{1}{\sigma_i^2}$$

The gradient equations for this sum of squares are

$$-2 \sum_i W_{ii} \frac{\partial f(x_i, \boldsymbol{\beta})}{\partial \beta_j} r_i = 0, \quad j = 1, \ldots, m$$

which, in a linear least squares system give the modified normal equations,

$$\sum_{i=1}^{n} \sum_{k=1}^{m} X_{ij} W_{ii} X_{ik} \hat{\beta}_k = \sum_{i=1}^{n} X_{ij} W_{ii} y_i, \quad j = 1, \ldots, m.$$



**Weighted Regression vs Ordinary Linear Regression**

Legend: — Weighted regression line; - - - Ordinary linear regression line

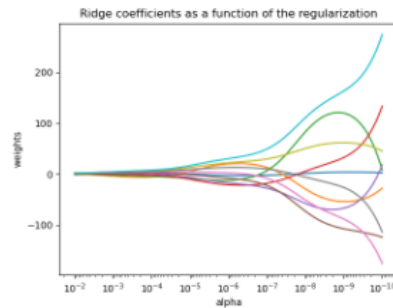X-axis: Working hours per week. Y-axis: Monthly salary (in USD)

# MODELO RIDGE

## 1.1.2.1. Regression
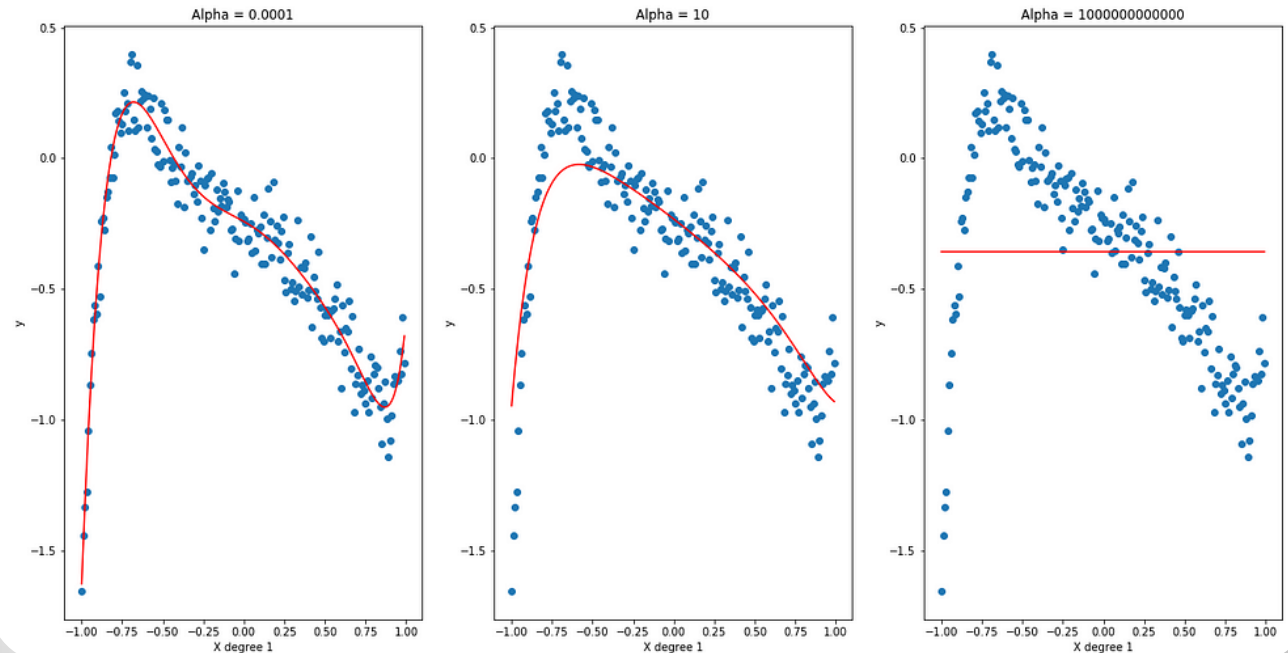
Ridge regression addresses some of the problems of Ordinary Least Squares by imposing a penalty on the size of the coefficients. The ridge coefficients minimize a penalized residual sum of squares:

$$\min_{w} ||Xw - y||_2^2 + \alpha ||w||_2^2$$

The complexity parameter $\alpha \geq 0$ controls the amount of shrinkage: the larger the value of $\alpha$, the greater the amount of shrinkage and thus the coefficients become more robust to collinearity.



Ridge coefficients as a function of the regularization



Ridge Regression model fits for different tuning parameters alpha

# MODELO LASSO

## 1.1.3. Lasso

The `Lasso` is a linear model that estimates sparse coefficients. It is useful in some contexts due to its tendency to prefer solutions with fewer non-zero coefficients, effectively reducing the number of features upon which the given solution is dependent. For this reason, Lasso and its variants are fundamental to the field of compressed sensing. Under certain conditions, it can recover the exact set of non-zero coefficients (see Compressive sensing: tomography reconstruction with L1 prior (Lasso)).

Mathematically, it consists of a linear model with an added regularization term. The objective function to minimize is:

$$\min_{w} \frac{1}{2n_{\text{samples}}} ||Xw - y||_2^2 + \alpha ||w||_1$$

The lasso estimate thus solves the minimization of the least-squares penalty with $\alpha ||w||_1$ added, where $\alpha$ is a constant and $||w||_1$ is the $\ell_1$-norm of the coefficient vector.

# MODELO ELASTIC-NET
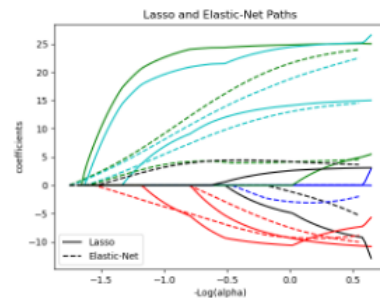
## 1.1.5. Elastic-Net

`ElasticNet` is a linear regression model trained with both $\ell_1$ and $\ell_2$-norm regularization of the coefficients. This combination allows for learning a sparse model where few of the weights are non-zero like `Lasso`, while still maintaining the regularization properties of `Ridge`. We control the convex combination of $\ell_1$ and $\ell_2$ using the `l1_ratio` parameter.

Elastic-net is useful when there are multiple features that are correlated with one another. Lasso is likely to pick one of these at random, while elastic-net is likely to pick both.

A practical advantage of trading-off between Lasso and Ridge is that it allows Elastic-Net to inherit some of Ridge's stability under rotation.

The objective function to minimize is in this case

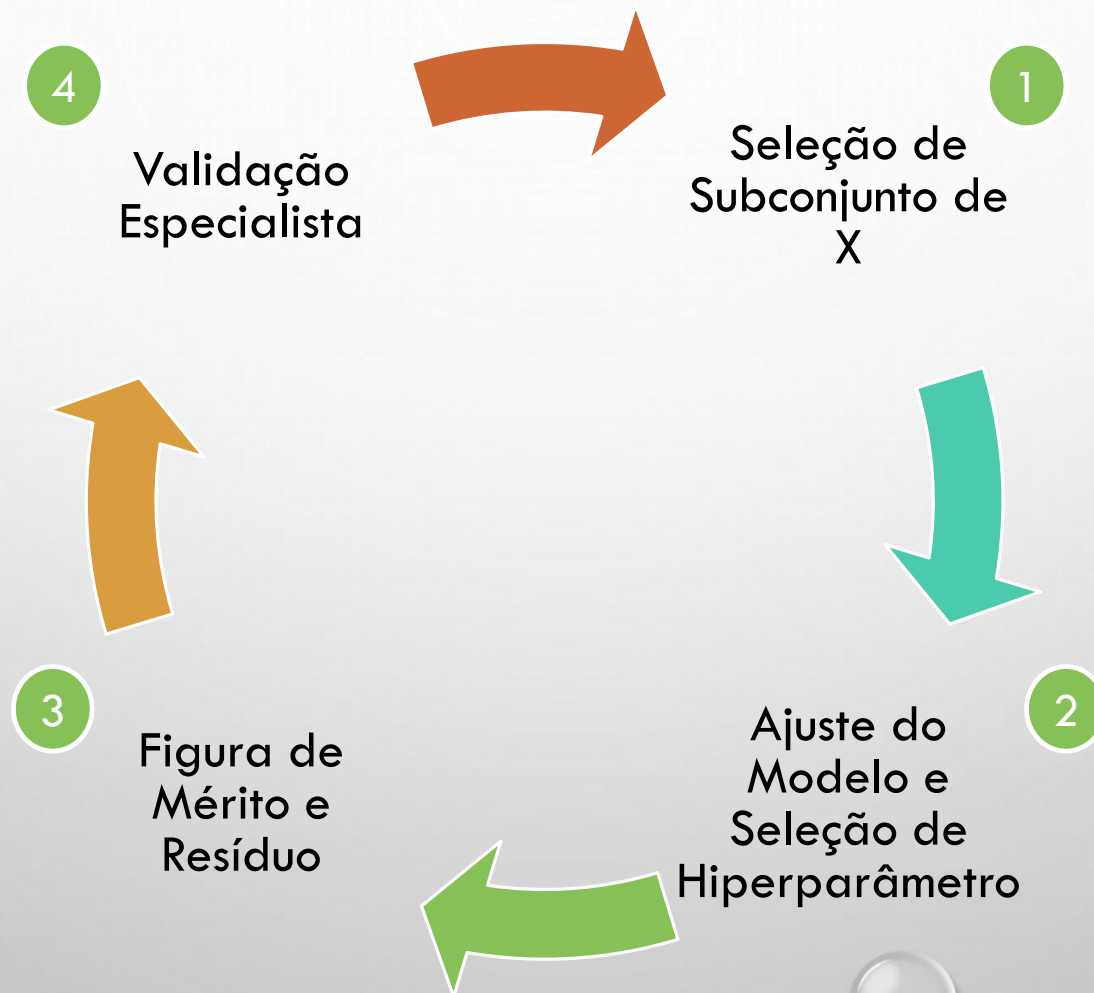$$\min_{w} \frac{1}{2n_{\text{samples}}}||Xw - y||_2^2 + \alpha\rho||w||_1 + \frac{\alpha(1-\rho)}{2}||w||_2^2$$



The class `ElasticNetCV` can be used to set the parameters `alpha` ($\alpha$) and `l1_ratio` ($\rho$) by cross-validation.

# CASE : CONSUMO DE COMBUSTÍVEL

# META-HEURÍSTICA DE TREINAMENTO

**4** Validação Especialista

**1** Seleção de Subconjunto de X

**2** Ajuste do Modelo e Seleção de Hiperparâmetro

**3** Figura de Mérito e Resíduo

40Mi de registros, contendo a informação de 39 sensores para o Navio a ser modelado.

Fundidas em 5Mi de linhas em milissegundos, 2990 linhas completas, em horas, com **Consumo, GPS, Velocidade pela Água, Vento e Potência.**

Histograma do Consumo/Hora >> Duas Modas (Baixo e Alto Consumo / Hora)

Distribuição do Consumo/Hora

Consumo/h ~ 329.37±471.75

Motor ~ 1.79±2.62

Velocidade ~ 4.27±5.81

Distância ~ 6.71±9.73

Vento Proa ~ 0.83±5.42

Correlação entre o Consumo e os indicadores do motor – Speed Through Water, Potência do Motor e Distância Geodésica

Indicadores mais correlacionados com o consumo (1) Potência (2) Velocidade sobre a Água (3) Distância (4) Vento Proa

```
                    OLS Regression Results
==============================================================================
Dep. Variable:          fuel_per_hour   R-squared:                      0.990
Model:                            OLS   Adj. R-squared:                 0.990
Method:                 Least Squares   F-statistic:                7.352e+04
Date:                Tue, 20 Dec 2022   Prob (F-statistic):              0.00
Time:                        09:15:32   Log-Likelihood:                -15749.
No. Observations:                2986   AIC:                        3.151e+04
Df Residuals:                    2981   BIC:                        3.154e+04
Df Model:                           4
Covariance Type:            nonrobust
==============================================================================
                  coef    std err          t      P>|t|      [0.025      0.975]
------------------------------------------------------------------------------
Intercept       4.8308      1.057      4.572      0.000       2.759       6.903
engine        163.2435      3.997     40.839      0.000     155.406     171.081
engine:stw     -4.5638      0.295    -15.471      0.000      -5.142      -3.985
wind_y          1.0858      0.166      6.544      0.000       0.760       1.411
geodistance    20.0009      0.476     41.996      0.000      19.067      20.935
==============================================================================
Omnibus:                      231.103   Durbin-Watson:                  0.305
Prob(Omnibus):                  0.000   Jarque-Bera (JB):            1080.985
Skew:                           0.209   Prob(JB):                   1.85e-235
Kurtosis:                       5.918   Cond. No.                        197.
==============================================================================
```
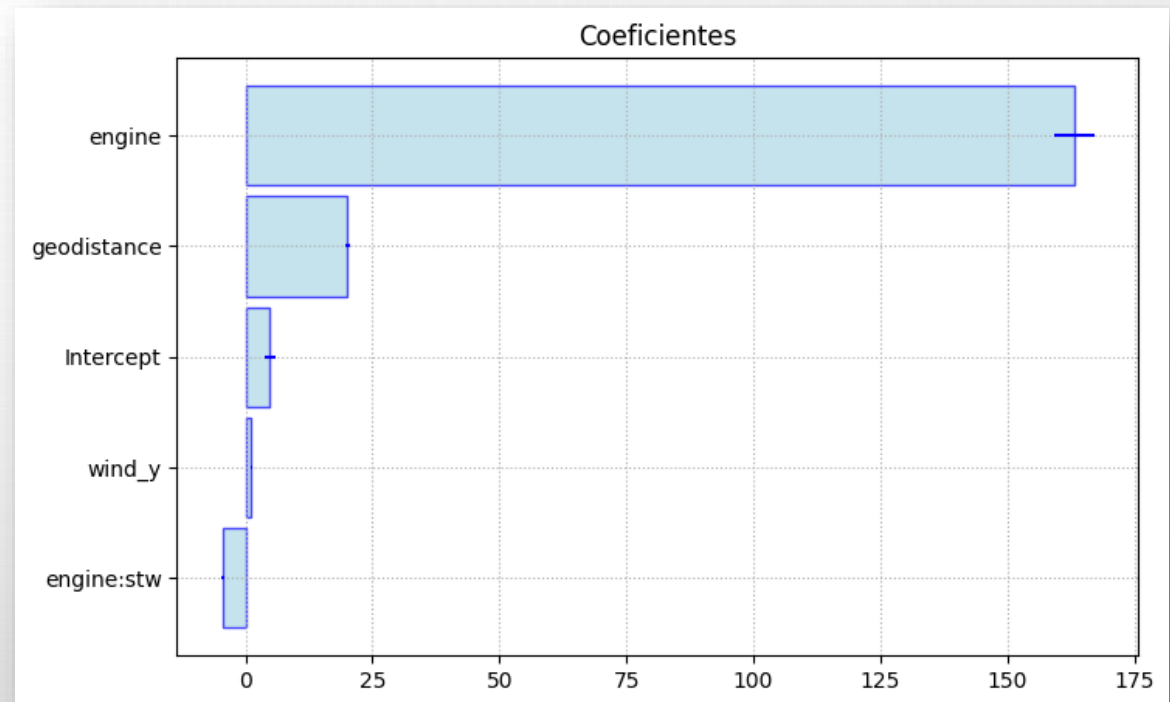
Resultados da Regressão



Coeficientes :
a influência de cada variável no consumo/hora.

# Fuel_hour ~ 163.2*engine - 4.6*engine:stw + 1.1*wind_y + 20.0*geodistance



Fórmula: fuel_hour ~ 163.2*engine -4.6*engine:stw + 1.1*wind_y + 20.0*geodistance

**10 Melhores Estimativas**

| day | hour | fuel_per_hour | y_est |
|---|---|---|---|
| 2022-07-23 | 594 | 6.0 | 5.998875 |
| 2022-09-10 | 1764 | 25.0 | 25.018087 |
| 2022-08-02 | 828 | 8.0 | 7.970102 |
| 2022-09-22 | 2052 | 0.0 | -0.033011 |
| 2022-10-07 | 2415 | 0.0 | 0.034579 |
| 2022-09-24 | 2106 | 0.0 | 0.035541 |
| 2022-08-06 | 935 | 7.0 | 6.957175 |
| 2022-11-02 | 3044 | 0.0 | 0.061800 |
| 2022-10-11 | 2496 | 0.0 | -0.072006 |
| 2022-09-15 | 1877 | 0.0 | 0.078089 |

**10 Melhores Estimativas > 500 fuel_per_hour**

| day | hour | fuel_per_hour | y_est |
|---|---|---|---|
| 2022-08-15 | 1148 | 1055.0 | 1054.644730 |
| 2022-09-27 | 2164 | 1036.0 | 1035.530274 |
| 2022-08-10 | 1025 | 1043.0 | 1043.645077 |
| 2022-10-13 | 2557 | 1050.0 | 1049.257795 |
| 2022-07-17 | 452 | 1081.0 | 1081.805467 |
|  | 450 | 1063.0 | 1062.108331 |
|  | 455 | 1071.0 | 1071.957986 |
| 2022-08-08 | 977 | 1083.0 | 1084.356100 |
| 2022-10-21 | 2755 | 928.0 | 926.483196 |
| 2022-10-08 | 2437 | 992.0 | 990.373949 |

Análise do Resíduo: 99% das estimativas dentro do intervalo de 3-sigma.

# COMPARAR O DESEMPENHO DE MODELOS DE MACHINE LEARNING UTILIZANDO VALIDAÇÃO CRUZADA I