

CIÊNCIA DE DADOS APLICADA A ANÁLISE ESPORTIVA UTILIZANDO PYTHON AVANÇADO

## MACHINE LEARNING: REGRESSÃO

DIEGO RODRIGUES DSC

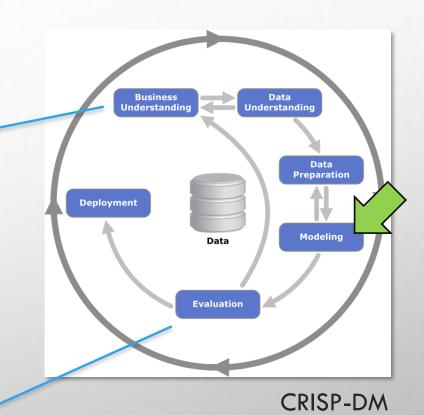
INFNET

#### AGENDA

- PARTE 1 : TEORIA
  - CONCEITOS
  - CASE : CONSUMO DE COMBUSTÍVEL DE NAVIO
  - CASE : PREVENDO NÚMERO DE GOLS EM UMA PARTIDA DE FUTEBOL

### CRONOGRAMA

NÚMERO	ÁREA	AULA	TRABALHOS
1	Intro	Introdução a Disciplina e Organização do Ambiente	
2	Dados	Coleta de Dados e Sensoriamento	
3		Variáveis Aleatórias	Grupos
4		Análise Exploratória	
5		Estatísticas para Ranqueamento	
6	Estatística	Ranqueamento Estatístico : ELO	
7		Ranqueamento Estatístico : Glicko	
8		Ranqueamento Estatístico : TrueSkill	
9		Ranqueamento Estatístico : XELO	Base de Dados
10		Modelos de Aprendizado de Máquina	N /
11		Machine Learning: Classificação	~
12	ML	Machine Learning: Regressão	
13		Machine Learning: Agrupamento	Pesquisa
14		Machine Learning: Visão Computacional	
15		Aplicações & Artigos: Esportes Independentes	Modelo
16	Espertes	Aplicações & Artigos: Esportes de Objeto	
17	Esportes	Aplicações & Artigos: Esportes de Combate	
18		Aplicações & Artigos : Betting	
19		Workshop	



#### **AMBIENTE PYTHON**

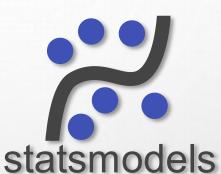


4. Variáveis Aleatórias



5. Visualização

6. Estimação e Inferência



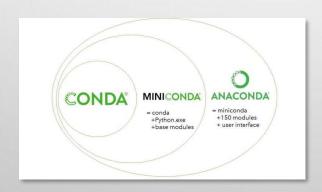


7. Machine Learning





1. Editor de Código



2. Gestor de Ambiente



3. Ambiente Python do Projeto

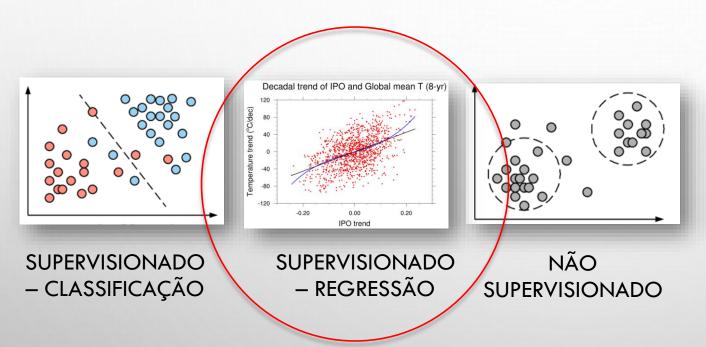


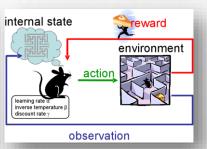
3. Notebook Dinâmico



## CONCEITOS

## PARADIGMAS DE MODELAGEM ESTATÍSTICA







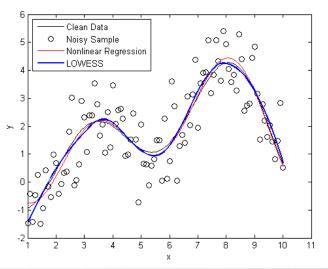


**GENERATIVO** 

# Regressão

O objetivo da regressão é modelar as relações funcionais entre dois conjuntos de variáveis.

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

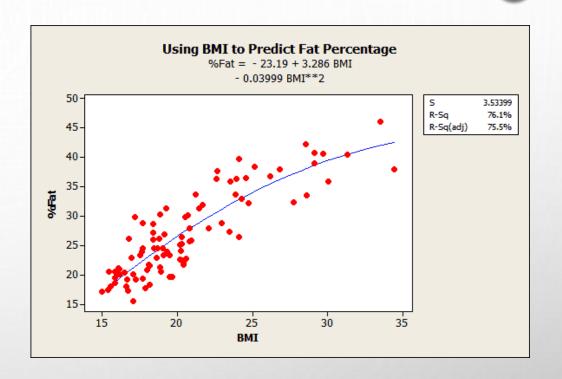


As vezes quando o mundo não é linear & gaussiano...

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).

## Regressão Linear: Modelo Matemático

#### Formulation [edit]

Given a data set  $\{y_i, x_{i1}, \dots, x_{ip}\}_{i=1}^n$  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

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

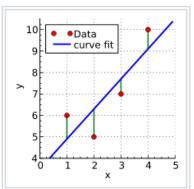
where  $^{\mathsf{T}}$  denotes the transpose, so that  $\mathbf{x}_i^{\mathsf{T}}\boldsymbol{\beta}$  is the inner product between vectors  $\mathbf{x}_i$  and  $\boldsymbol{\beta}$ .

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

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

where

$$\mathbf{y} = egin{bmatrix} egin{aligned} \mathbf{y} &= egin{bmatrix} egin{aligned} y_2 \ dots \ y_n \end{bmatrix}, \ \mathbf{X} &= egin{bmatrix} \mathbf{x}_1^{\mathsf{T}} \ \mathbf{x}_2^{\mathsf{T}} \ dots \ \mathbf{x}_n^{\mathsf{T}} \end{bmatrix} = egin{bmatrix} 1 & x_{11} & \cdots & x_{1p} \ 1 & x_{21} & \cdots & x_{2p} \ dots & dots & \ddots & dots \ 1 & x_{n1} & \cdots & x_{np} \end{bmatrix}, \ eta &= egin{bmatrix} eta_0 \ eta_1 \ eta_2 \ dots \ eta_p \end{bmatrix}, & oldsymbol{arepsilon} & oldsymbol{arepsilon} & eta \ egin{bmatrix} arepsilon_1 \ arepsilon_2 \ dots \ eta_n \end{bmatrix}. \ \ egin{bmatrix} eta &= egin{bmatrix} arepsilon_1 \ eta_2 \ dots \ eta_n \end{bmatrix}. \end{aligned}$$



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).

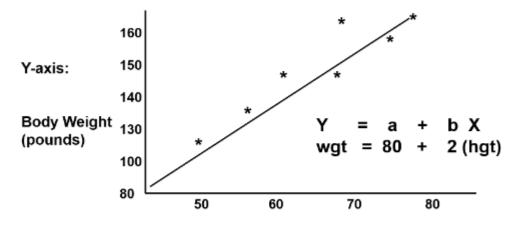
$$y = \sum_{i}^{V} \beta_{i} x_{i} + \varepsilon$$



## 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.



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

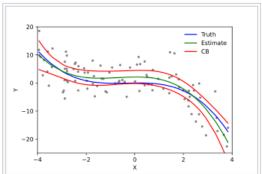
#### Assumptions [edit]

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



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 | 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.

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.)

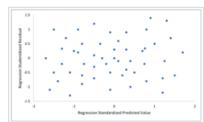
- Additivity: f(x + y) = f(x) + f(y).
- Homogeneity of degree 1:  $f(\alpha x) = \alpha f(x)$  for all  $\alpha$ .

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

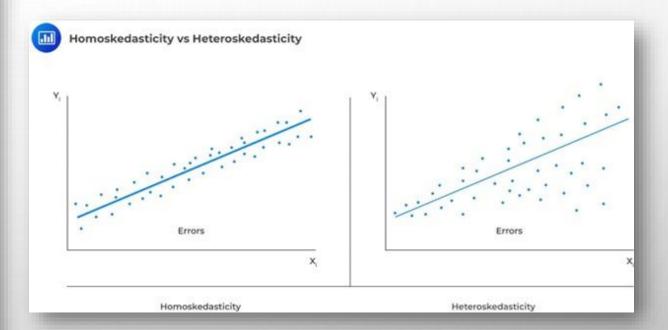
spread as one moves left to right on the plot). A plot of the absolute or squared residuals versus the predicted values (or each predictor) can also be examined for a trend or curvature. Formal tests can also be used; see Heteroscedasticity. The presence of heteroscedasticity will result in an overall "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 lognormal 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).



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.

inflation factor). In the case of perfect multicollinearity, the parameter vector  $\boldsymbol{\beta}$  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.



## 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.,  $\mathbf{x}_i = [x_{i1},x_{i2},\ldots,x_{ip}]^{\mathrm{T}}$ . 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,

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

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  $\mathbf{x}_i^{\mathrm{T}}$  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  $(\beta)$ .

#### Matrix/vector formulation [edit]

Consider an overdetermined system

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

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

$$X\beta = y$$
,

where

$$\mathbf{X} = egin{bmatrix} X_{11} & X_{12} & \cdots & X_{1p} \ X_{21} & X_{22} & \cdots & X_{2p} \ dots & dots & \ddots & dots \ X_{n1} & X_{n2} & \cdots & X_{np} \end{bmatrix}, \qquad oldsymbol{eta} = egin{bmatrix} eta_1 \ eta_2 \ dots \ eta_p \end{bmatrix}, \qquad \mathbf{y} = egin{bmatrix} y_1 \ y_2 \ dots \ y_n \end{bmatrix}.$$

(Note: for a linear model as above, not all elements in  ${\bf X}$  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  ${\bf p}$  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  $\beta$  which fit the equations "best", in the sense of solving the quadratic minimization problem

$$\hat{oldsymbol{eta}} = rg \min_{oldsymbol{eta}} S(oldsymbol{eta}),$$

where the objective function  $oldsymbol{S}$  is given by

$$S(oldsymbol{eta}) = \sum_{i=1}^n \left| y_i - \sum_{j=1}^p X_{ij} eta_j 
ight|^2 = \left\| \mathbf{y} - \mathbf{X} oldsymbol{eta} 
ight\|^2.$$

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*:

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

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

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

or

$$\hat{oldsymbol{eta}} = oldsymbol{eta} + \left(\mathbf{X}^{\mathrm{T}}\mathbf{X}
ight)^{-1}\mathbf{X}^{\mathrm{T}}oldsymbol{arepsilon}.$$



## 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]

$$E[\varepsilon \mid X] = 0.$$

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

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]

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

Usually, it is also assumed that the regressors have finite moments up to at least the second moment. Then the matrix  $Q_{xx} = E[X^TX/n]$  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:<sup>[18]</sup>

$$\operatorname{Var}[\,arepsilon\mid X\,] = \sigma^2 I_n,$$

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:  $E[s_i^2|X] = \sigma^2$ , 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:  $\mathbb{E}\left[\varepsilon_{i}\varepsilon_{j}|X\right]=0$  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:<sup>[19]</sup>

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

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.<sup>[20]</sup>

## Validação : Statsmodels

#### OLS Regression Results

Dep. Variable:	Ticker	R-squared:	0.782
Model:	OLS	Adj. R-squared:	0.776
Method:	Least Squares	F-statistic:	128.4
Date:	Sun, 11 Sep 2022	Prob (F-statistic):	3.25e-79
Time:	15:28:31	Log-Likelihood:	-844.20
No. Observations:	259	AIC:	1704.
Df Residuals:	251	BIC:	1733.

Df Model: 7 Covariance Type: nonrobust

	coef	std err	t	P> t	[0.025	0.975]
const	-45.0907	27.593	-1.634	0.103	99.434	9.253
UK FTSE	0.0227	0.002	9.162	0.000	0.018	0.028
ASX200 (yesterday)	0.0358	0.003	12.808	0.000	0.039	0.041
Japan Nikkei (y'day)	-0.0028	0.001	-4.576	0.000	-0.004	-0.002
AUD TWI 4pm	-2.1353	0.413	-5.175	0.000	-2.948	-1.323
Iron Ore futures (\$US/t)	-0.3426	0.034	-10.080	0.000	-0.409	-8,276
Uranium, weekly (\$US/lb)	0.5427	0.092	5.879	0.000	0.361	0.725
Copper (\$US/t)	0.0041	0.001	6.149	0.000	0.003	0.005
=======================================					======	

Kurtosis:	2.719	Cond. No.	2.16e+06
Skew:	0.001	Prob(JB):	0.653
Prob(Omnibus):	0.674	Jarque-Bera (JB):	0.853
Omnibus:	0.789	Durbin-Watson:	0.618

#### Warnings:

- [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [2] The condition number is large, 2.16e+06. This might indicate that there are strong multicollinearity or other numerical problems.

#### Coeficiente de Determinação R<sup>2</sup>

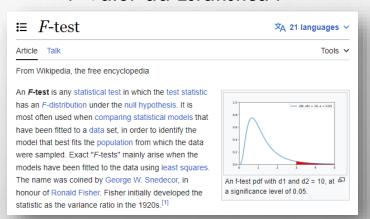
$$R^2 = 1 - rac{RSS}{TSS}$$

 $R^2$  = coefficient of determination

RSS = sum of squares of residuals

TSS = total sum of squares

#### P Valor da Estatística F

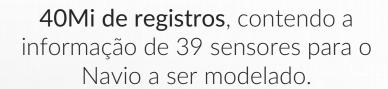


#### P Valor dos Coeficientes

Número de Condicionamento

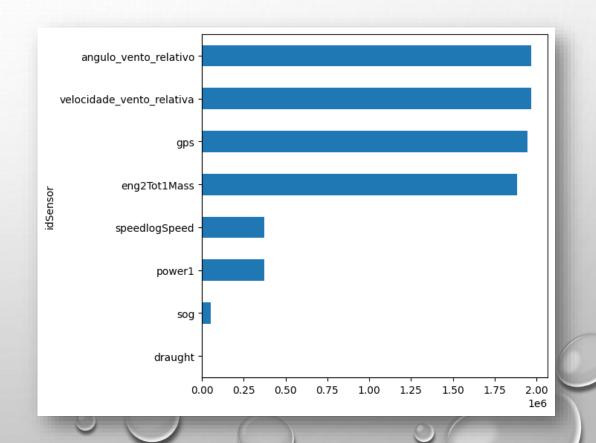


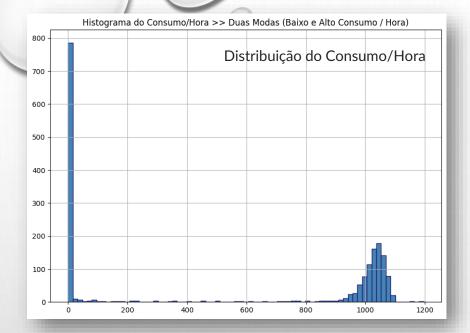
## EXEMPLO II – CONSUMO DE COMBUSTÍVEL

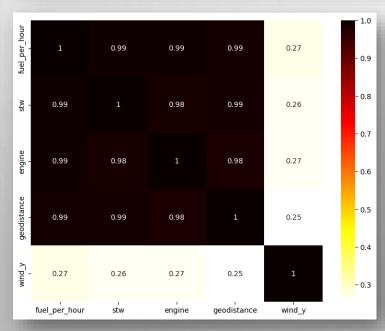


	idSensor	date	value	day
0	cog	1635724807952	19.5	2021-11-01
1	sog	1635724807952	0.2	2021-11-01
2	angulo_vento_relativo	1635724807972	226.7	2021-11-01
3	velocidade_vento_relativa	1635724807972	2.4	2021-11-01
4	bussola	1635724807999	104.1	2021-11-0
***			***	3.7
14545	angulo_vento_relativo	1668793146825	83.3	2022-11-18
14546	velocidade_vento_relativa	1668793146825	6.5	2022-11-18
14547	heading	1668793156581	333.1	2022-11-18
14548	speed	1668793157080	0	2022-11-18
14549	cog	1668793157080	282.4	2022-11-18

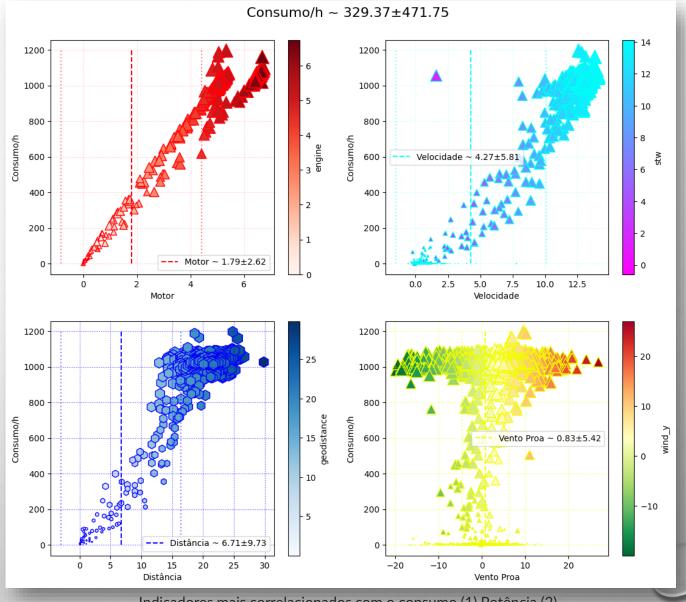
Fundidas em 5Mi de linhas em milissegundos, 2990 linhas completas, em horas, com Consumo, GPS, Velocidade pela Água, Vento e Potência.



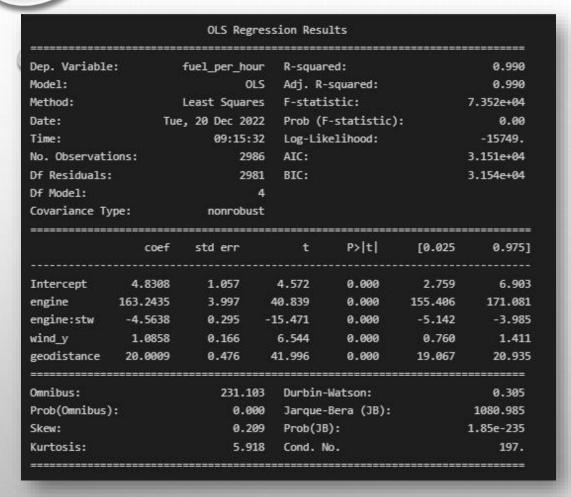




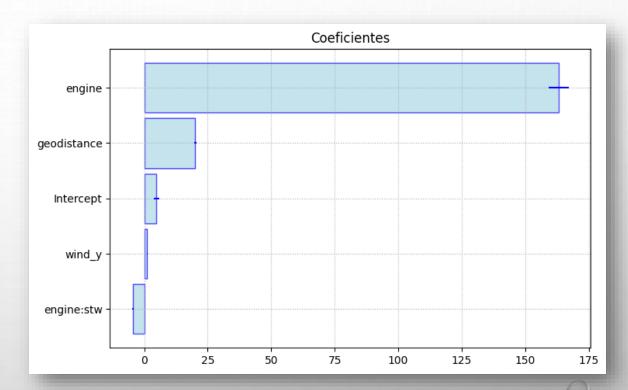
Correlação entre o Consumo e os indicadores do motor – Speed Through Water, Potência do Motor e Distância Geodésica



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



Resultados da Regressão



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

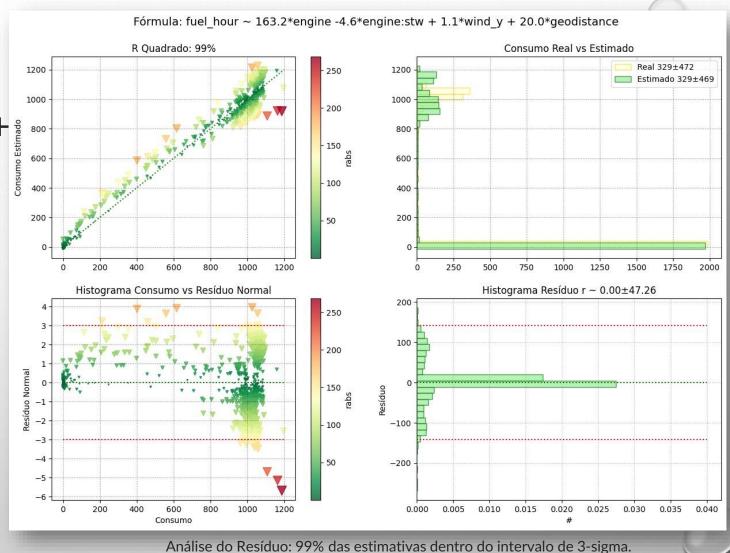
## Fuel\_hour ~ 163.2\*engine - 4.6\*engine:stw + 1.1\*wind\_y + 20.0\*geodistance

		fuel_per_hour	y_est
day	hour		
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

		fuel_per_hour	y_est
day	hour		
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

10 Melhores Estimativas > 500 fuel\_per\_hour



## DESAFIO: LEITURA DO ARTIGO E IMPLEMENTAÇÃO DA REGRESSÃO

PRÓXIMA AULA LEITURA: AGRUPAMENTO