



SCIENCE

MACHINE LEARNING
E DATA MINING



BIG DATA SCIENCE

ESTATÍSTICA APLICADA – PARTE 2

ANÁLISE DE REGRESSÃO

DIÓGENES JUSTO



REVISÃO DA LISTA DE EXERCÍCIOS





REVISÃO DO CONTEÚDO DA PRIMEIRA AULA

MÉDIA, DESVIO PADRÃO SÃO MÉTRICAS RICAS?

```
mean(anscombe$x1)
```

```
...
```

```
mean(anscombe$x4)
```

```
mean(anscombe$y1)
```

```
...
```

```
sd(anscombe$x1)
```

```
...
```

```
cor(anscombe$x1, anscombe$y1)
```

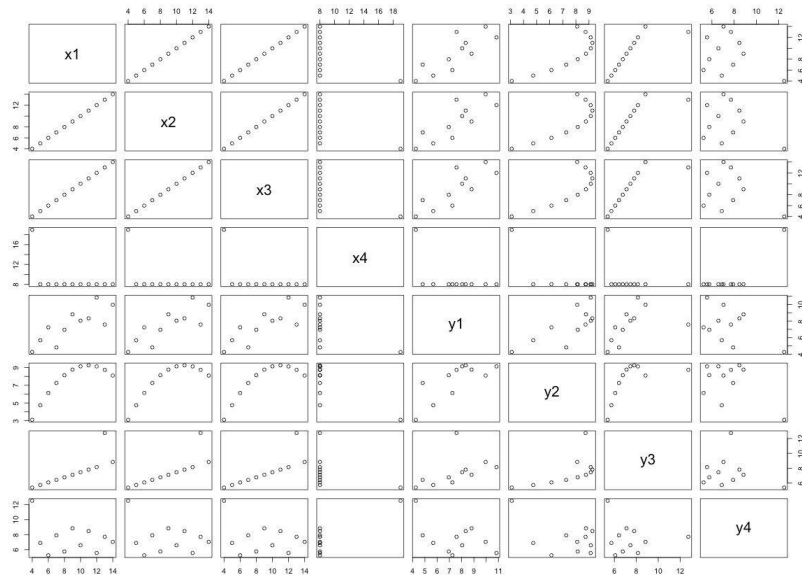
```
...
```

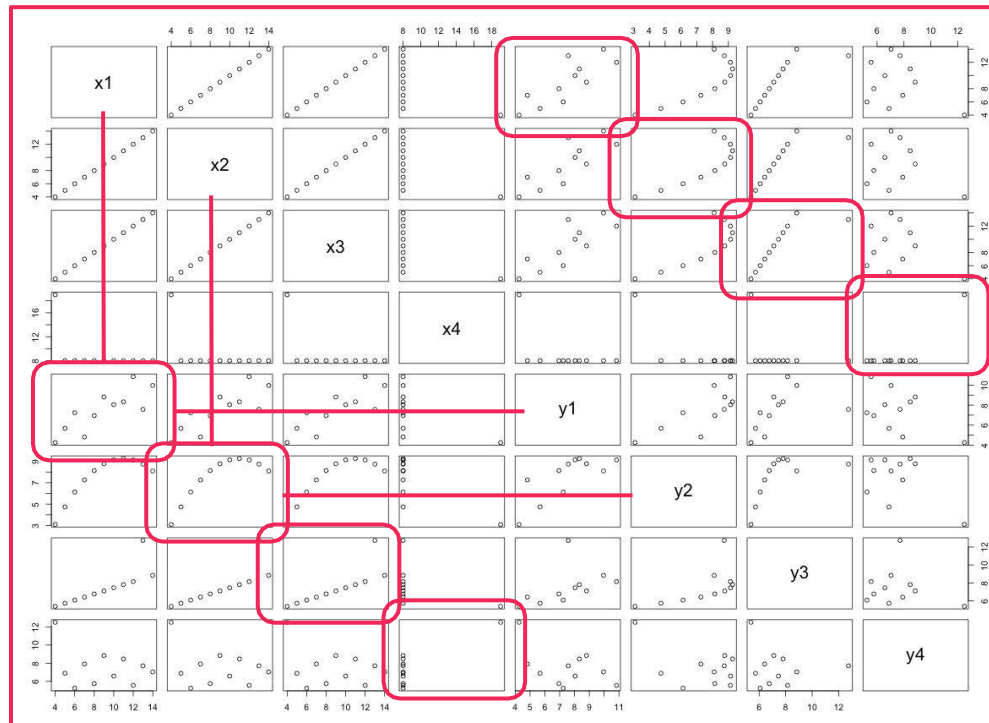
MÉDIA, DESVIO PADRÃO SÃO MÉTRICAS RICAS?

```
par(mfrow=c(2,2))  
plot(anscombe$x1, anscombe$y1)  
plot(anscombe$x2, anscombe$y2)  
plot(anscombe$x3, anscombe$y3)  
plot(anscombe$x4, anscombe$y4)  
par(mfrow=c(1,1))
```

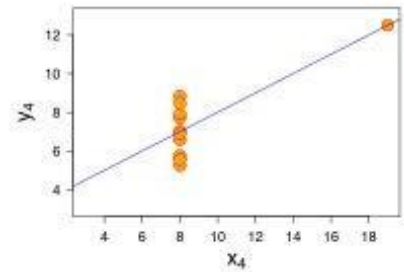
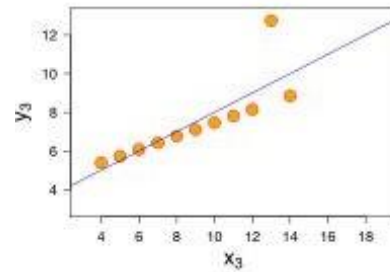
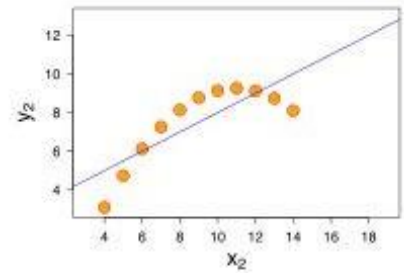
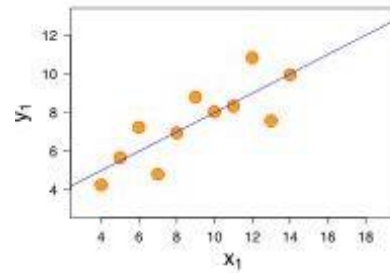
MÉDIA, DESVIO PADRÃO SÃO MÉTRICAS RICAS?

`pairs(anscombe)`





QUARTETO DE ANSCOMBE



https://en.wikipedia.org/wiki/Anscombe's_quartet

LISTA DE EXERCÍCIOS

Média, Mediana, Moda

Score Z

Variância e Desvio Padrão

COVARIÂNCIA

CORRELAÇÃO

- Ex.1:

R – `mtcars`, `mpg` x `wt` (peso do carro x consumo)

`cor(mtcars$mpg, mtcars$wt)`

`plot(mtcars$mpg, mtcars$wt)`

CORRELAÇÃO

Ex. 2: Excel, calcule a correlação linear do conjunto abaixo.

i	y	x
1	1	5
2	2	4
3	3	3
4	4	2
5	5	1

CORRELAÇÃO

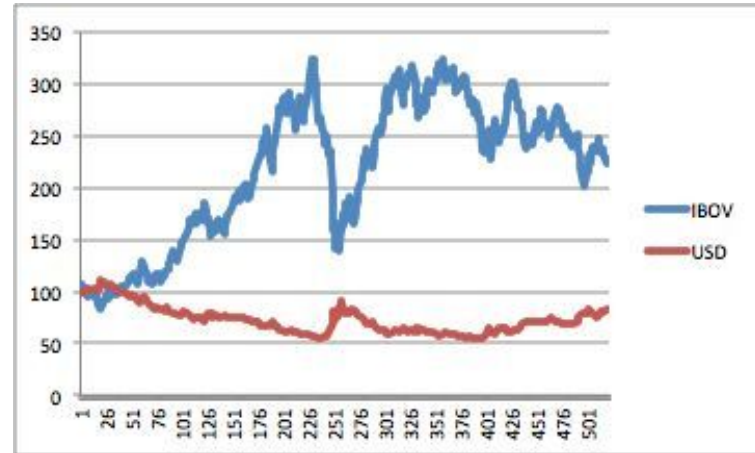
Ex. 3: Excel, calcule a correlação linear do índice iBovespa x dólar

```
=GoogleFinance("IBOV","price", DATE(2004, 1, 1),  
DATE(2013, 12, 31), "WEEKLY")
```

```
=GoogleFinance("USDBRL","price", DATE (2003, 12, 27), DATE(2013,  
12, =31),"WEEKLY")
```

Coloque (1) ambas séries na mesma base, faça (2) gráficos de linha, (3) linha em escalas diferentes, (4) dispersão e calcule a (5) correlação linear.

Qual a relação entre o índice
iBovespa e a cotação USD x REAL?

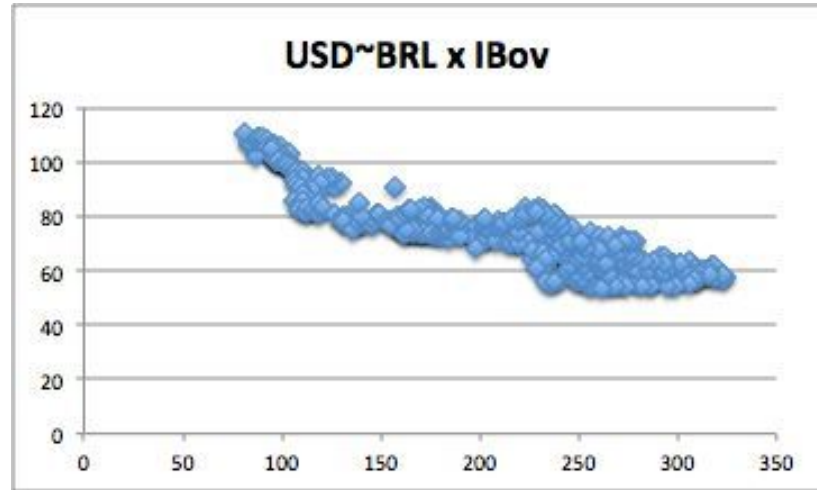


Qual a relação entre o índice
iBovespa e a cotação USD x REAL?



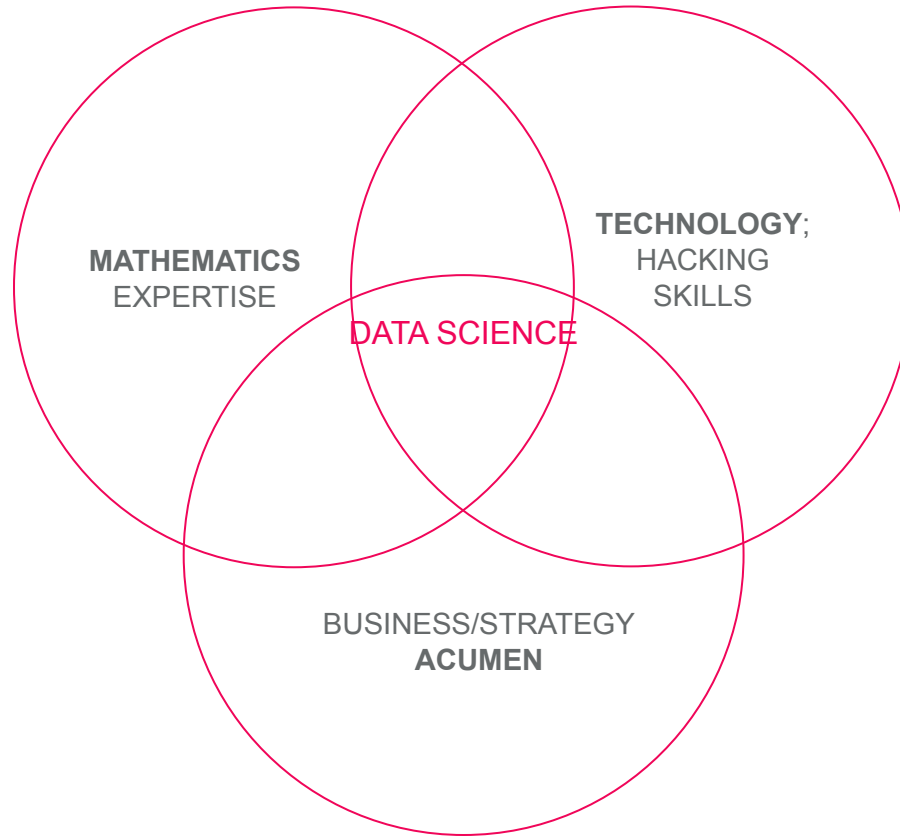
Qual a relação entre o índice
iBovespa e a cotação USD x REAL?

$\rho \sim -0,9$





DATA SCIENCE



STEPS IN DATA ANALYSIS

- Define the **QUESTION**
- Define the **IDEAL DATA SET**
- Determine **WHAT DATA** you can access
- **OBTAIN** the data
- **CLEAN** the data
- Exploratory **DATA ANALYSIS**
- Statistical **PREDICTION/MODELING**
- **INTERPRET** Results
- **CHALLENGE** results
- Synthesize/**WRITE UP** results
- Create reproducible **CODE**

Prever o PIB

Mensal / Anual, etc, como explicar?

PIB (público),

Tráfego (ABCR) PIB (Site Itau),

Tráfego (Site ABCR)

Gráficos, Estatística descritiva

Análise de regressão

OBTENDO DADOS

google -> pib itau mensal

<https://www.italy.com.br/italyba-pt/analises-economicas/nossas-series-economicas/pib-mensal-italy-unibanco>

google -> tráfico abcr download

<http://www.abcr.org.br/Download.ashx?arquivo=IndiceABCR.xls>



ANÁLISE EXPLORATÓRIA



- Verificar se os DADOS ESTÃO COMPLETOS
- GRÁFICOS
- ESTATÍSTICA descritiva
- CORRELAÇÃO É ALTA?



REGRESSÃO LINEAR

REGRESSÃO LINEAR

Chamaremos de modelo de regressão linear uma aplicação, com equação no seguinte formato:

$$y = \alpha \cdot x + \beta$$

A reta será uma aproximação, portanto, para cada ponto há uma diferença entre o ponto real e o ponto expressado por esta equação, Incluiremos, portanto, este fator de erro.

$$y_i = \alpha \cdot x_i + \beta + \epsilon_i$$

PARA O MODELO:

$$y_i = \alpha \cdot x_i + \beta + \varepsilon_i$$

Obs 1: alfa será o mesmo para todos valores. Utilizamos a notação indicial, pois posteriormente incluiremos outros fatores alfa.

Obs 2: épsilon tem índice i pois cada valor é independente, sendo correto incluir também o índice em y e x, portanto.

Obs 3: alfa, em geometria, é chamado de coeficiente angular e beta, coeficiente linear ou intercepto.

Obs 4: a variável x é dita independente e y dependente (depende de x). Ou ainda x exógena (exo=fora do modelo) e y endógena (endo=explicada dentro do modelo).



•

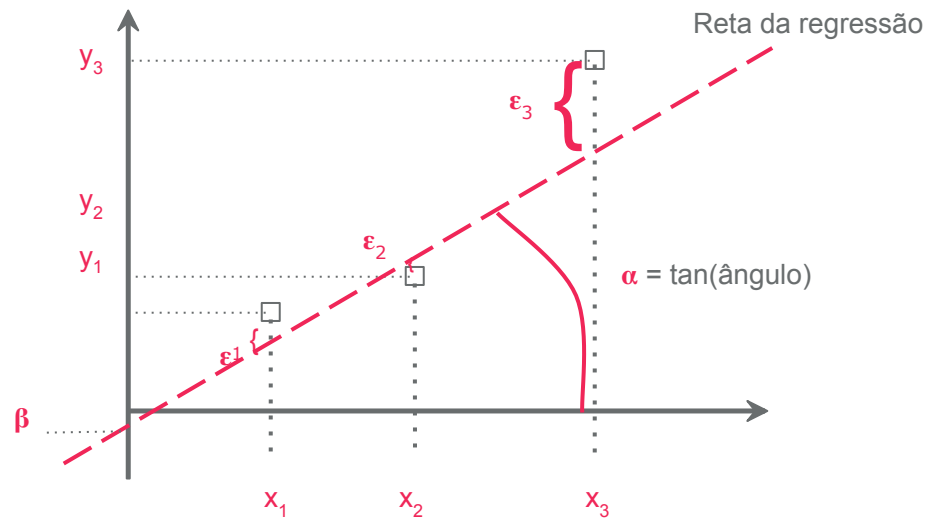


•



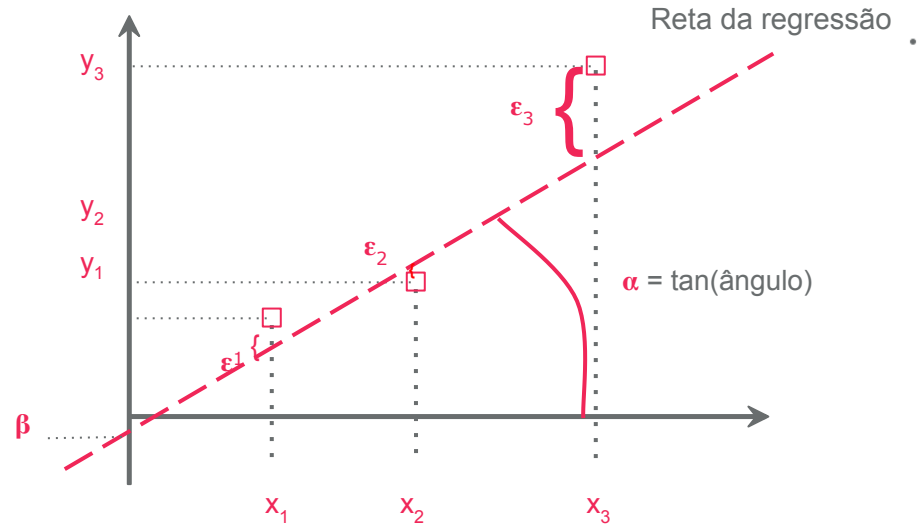
PARA O MODELO:

$$y_i = \alpha \cdot x_i + \beta + \varepsilon_i$$



PERG.: MAS COMO DETERMINAR
OS COEFICIENTES ALFA E BETA, A
PARTIR DOS DADOS?

Res.: Através de um método
matemático chamado mínimos
quadrados, onde procura-se minimizar
os erros (epsilon),
daí o nome do método.





MÍNIMOS QUADRADOS



- DADO SEU RIGOR MATEMÁTICO acadêmico, a explicação do método foge do escopo da disciplina em nível de pós-graduação. Para maiores referências, sugere-se consultar:
- <http://livrommq.blogspot.com.br/>
- Dentro do escopo da disciplina, utilizaremos o EXCEL como ferramenta de cálculo (que se utiliza do método em suas funções):
 - Suplemento ANÁLISE DE DADOS, Opção Regressão Linear
 - Funções: para beta, INTERCEPT (INTERCEPÇÃO) e; para alfa, SLOPE (INCLINAÇÃO)



PARA O MODELO:

$$y_i = a.x_i + b + e_i$$

Representação em formato de equação:

$$y_1 = a.x_1 + b + e_1$$

$$y_2 = a.x_2 + b + e_2$$

$$y_n = a.x_n + b + e_n$$



PARA O MODELO:

$$y_i = a \cdot x_i + b + e_i$$

Representação em formato tabular:

i	y_i	a	x_i	b	e_i
1	y_1	a	x_1	b	e_1
2	y_2	a	x_2	b	e_2
-	-		-		-
n	y_n	a	x_n	b	e_n

Obs: como os coeficientes alfa e beta são constantes, em geral os omitimos no formato tabular.

PARA O MODELO:

$$y_i = a.x_i + b + s_i$$

Representação em formato tabular:

i	y _i	x _i	s _i
1	y ₁	x ₁	s ₁
2	y ₂	x ₂	s ₂
-	-	-	-
n	y _n	x _n	s _n

REGRESSÃO LINEAR MULTIPLA

Expandimos o modelo inserindo variáveis independentes adicionais:

$$y_i = a_1 \cdot x_{1i} + a_2 \cdot x_{2i} + \beta + s_i$$

Ou ainda, generalizando para k variáveis:

$$y_i = a_1 \cdot x_{1i} + a_2 \cdot x_{2i} + \dots + a_k \cdot x_{ki} + \beta + s_i$$

Em um modelo múltiplo (ou multivariado) são tratadas duas ou mais variáveis independentes, explicativas, ou ainda, exógenas.



MODELO AUTOREGRESSIVO



Um modelo AR utiliza a própria "história" de uma variável para criar uma "segunda dimensão", dessa forma;

$$y_i = a_1 \cdot y_{i-1} + a_2 \cdot y_{i-2} + p + s_i$$

No exemplo acima, criamos um modelo AR-2 (com duas variáveis). Abaixo, generalizamos para AR-K (modelo autoregressivo de ordem K):

$$y_i = a_1 \cdot y_{i-1} + a_2 \cdot y_{i-2} + \dots + a_k \cdot y_{i-k} + p + s_i$$

Obs: aparentemente o modelo pode parecer complexo, mas veremos que sua aplicação é bem simples, na verdade, utilizando a representação tabular.



CONSTRUÇÃO 1/2

Expandiremos o modelo abaixo, para um AR-3:

$$y_i = a \cdot y_{i-1} + b + s_i$$

Representação em formato tabular:

i	y_i	y_{i-1}	y_{i-2}	y_{i-3}	s_i
1	y_1				s_1
2	y_2	y_1			s_2
3	y_3	y_2	y_1		s_3
4	y_4	y_3	y_2	y_1	s_4
n	y_n	y_{n-1}	y_{n-2}	y_{n-3}	s_n

CONSTRUÇÃO 2/2

Repare que a coluna y_{i-1} nada mais é do que a coluna y_i , deslocada uma linha para baixo. E y_{i-2} segue o mesmo comportamento, porém com deslocamento de duas linhas:

i	y_i	y_{i-1}	y_{i-2}	y_{i-3}	s_i
1	y_1				s_1
2	y_2	y_1			s_2
3	y_3	y_2	y_1		s_3
4	y_4	y_3	y_2	y_1	s_4
n	y_n	y_{n-1}	y_{n-2}	y_{n-3}	s_n

VARIÁVEIS DUMMY

- Uma variável dummy (também chamada indicador ou VARIÁVEL DE CATEGORIA OU BOOLEANA) é aplicada utilizando valor 0 ou 1 para indicar o que se propõe.
- É aplicável para:
 - Categorizar OUTLIERS;
 - Criar CATEGORIAS (fumante/não fumante, masc./fem., etc)
 - Sinalizar SAZONALIDADES em séries temporais
 - ENTRE OUTRAS categorizações

VARIÁVEL DUMMY PARA OUTLIER

Basta criar uma nova variável e, quando foi identificado um outlier (ver teste z), assinalar com 1.

i	y_i	x_i	d_i	s_i
1	y_1	x_1	0	s_1
2	y_2	x_2	1	s_2
3	y_3	x_3	0	s_3
n	y_n	x_n	0	s_n

Obs: no exemplo acima, somente x_2 foi identificado como outlier.

MODELO REGRESSIVO COM VARIÁVEL DUMMY

Expandiremos o modelo abaixo:

$$y_i = a_1 \cdot x_{1i} + a_2 \cdot x_{2i} + \beta + s_i$$

Que ficará agora:

$$y_i = a_1 \cdot x_i + a_2 \cdot d_{1i} + \beta + s_i$$

Onde $d_{1i} = \begin{cases} 0, & \text{se } x_i \text{ não for outlier} \\ 1, & \text{se } x_i \text{ for outlier} \end{cases}$

Obs: utilizamos aqui a notação indicial d_1 propositalmente, pois é comum utilizarmos várias variáveis *dummies*

VARIÁVEL DUMMY PARA SAZONALIDADES

Para captar o efeito sazonal, incluiremos uma variável dummy para cada evento sazonal.

Para o caso de identificarmos, por exemplo, que a cada mês do ano há uma sazonalidade (isto é, janeiros são similares ao longo dos anos, o mesmo ocorrendo para fevereiro e todos demais meses) as dummies serão:

$$\begin{aligned} D_{jan} &= \begin{cases} 1 & \text{se_for_mês_janeiro} \\ 0 & \text{caso_contrário} \end{cases} \\ D_{fev} &= \begin{cases} 1 & \text{se_for_mês_fevereiro} \\ 0 & \text{caso_contrário} \end{cases} \\ &\vdots \\ D_{nov} &= \begin{cases} 1 & \text{se_for_mês_novembro} \\ 0 & \text{caso_contrário} \end{cases} \end{aligned}$$

MODELO REGRESSIVO COM DUMMY SAZONAL

Nosso modelo será expandido para:

$$y_i = a_1 \cdot x_i + a_2 \cdot d_{1i} + a_2 \cdot d_{2i} + \dots + a_{12} \cdot d_{12i} + p + s_i$$

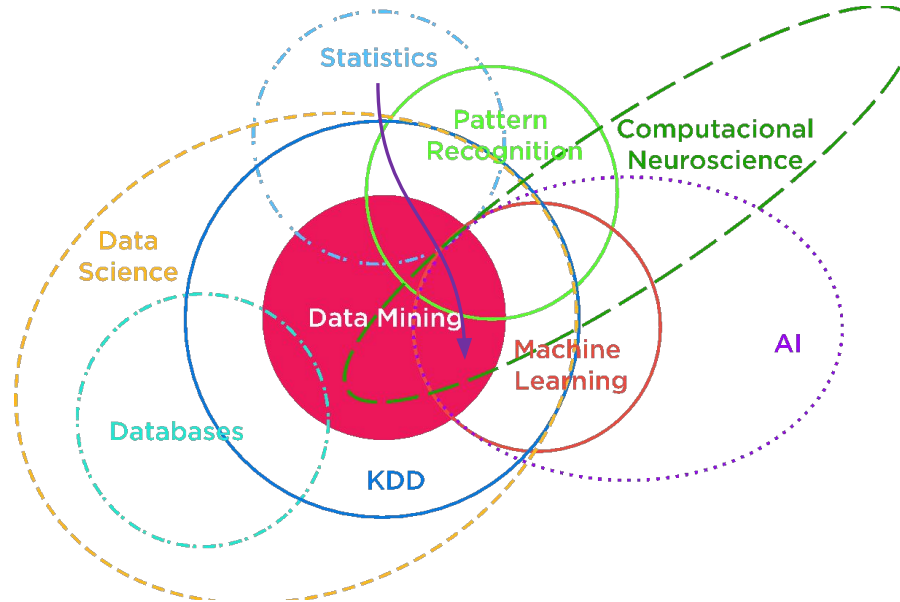
Pode parecer complexo a primeira vista, mas veremos no formato tabular que é simples.

MODELO REGRESSIVO COM DUMMY SAZONAL

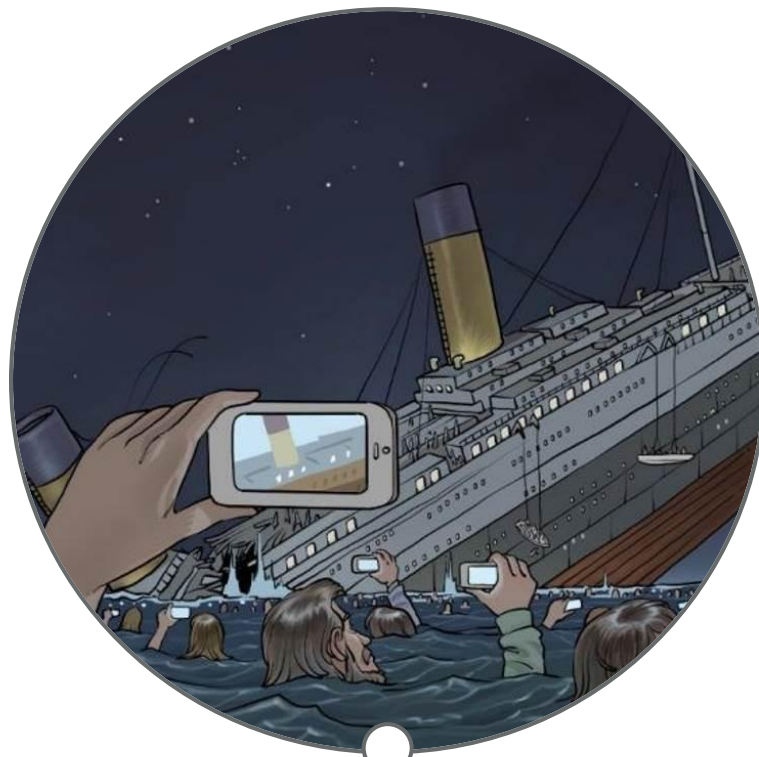
Basta criar uma coluna para cada evento sazonal

i	y_i	x_i	d_i	d_i	...	d_i	s_i
1	y_1	x_1	1	0		0	s_1
2	y_2	x_2	0	1		0	s_2
12	y_{12}	x_{12}	0	0		1	s_{12}
13	y_{13}	x_{13}	1	0		0	s_{13}
14	y_{14}	x_{14}	0	1		0	s_{14}
n	y_n	x_n	d_{1n}	d_{2n}		d_{12n}	s_n

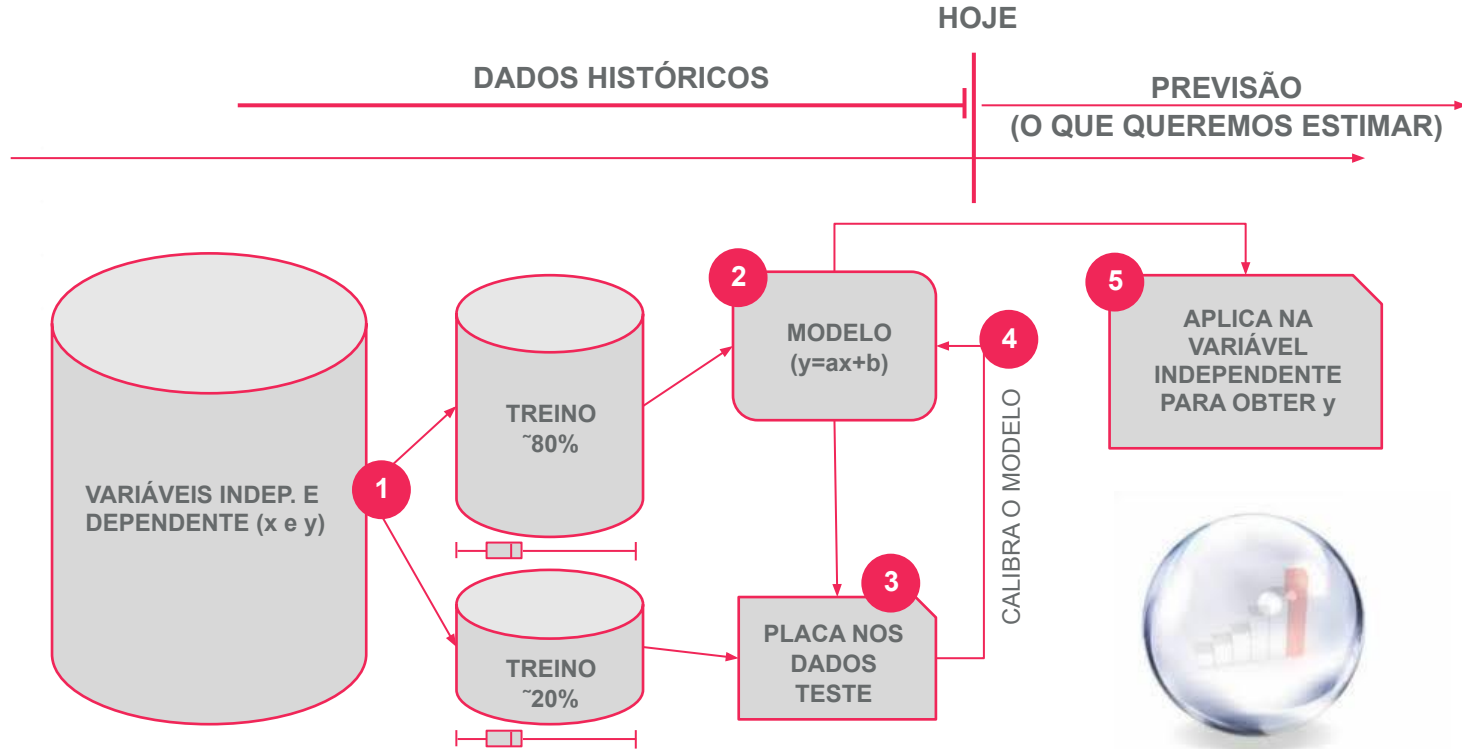
Por onde andamos no nosso "PASSEIO"?



DESAFIO



QUAIS DADOS USAR NA ESTIMAÇÃO DO MODELO?



Obrigado



ProfDiogenes.Justo@fiap.com.br



Copyright © 2018 | Diógenes Justo

Todos os direitos reservados. Reprodução ou divulgação total ou parcial deste documento, é expressamente proibido sem consentimento formal, por escrito, do professor/autor.