

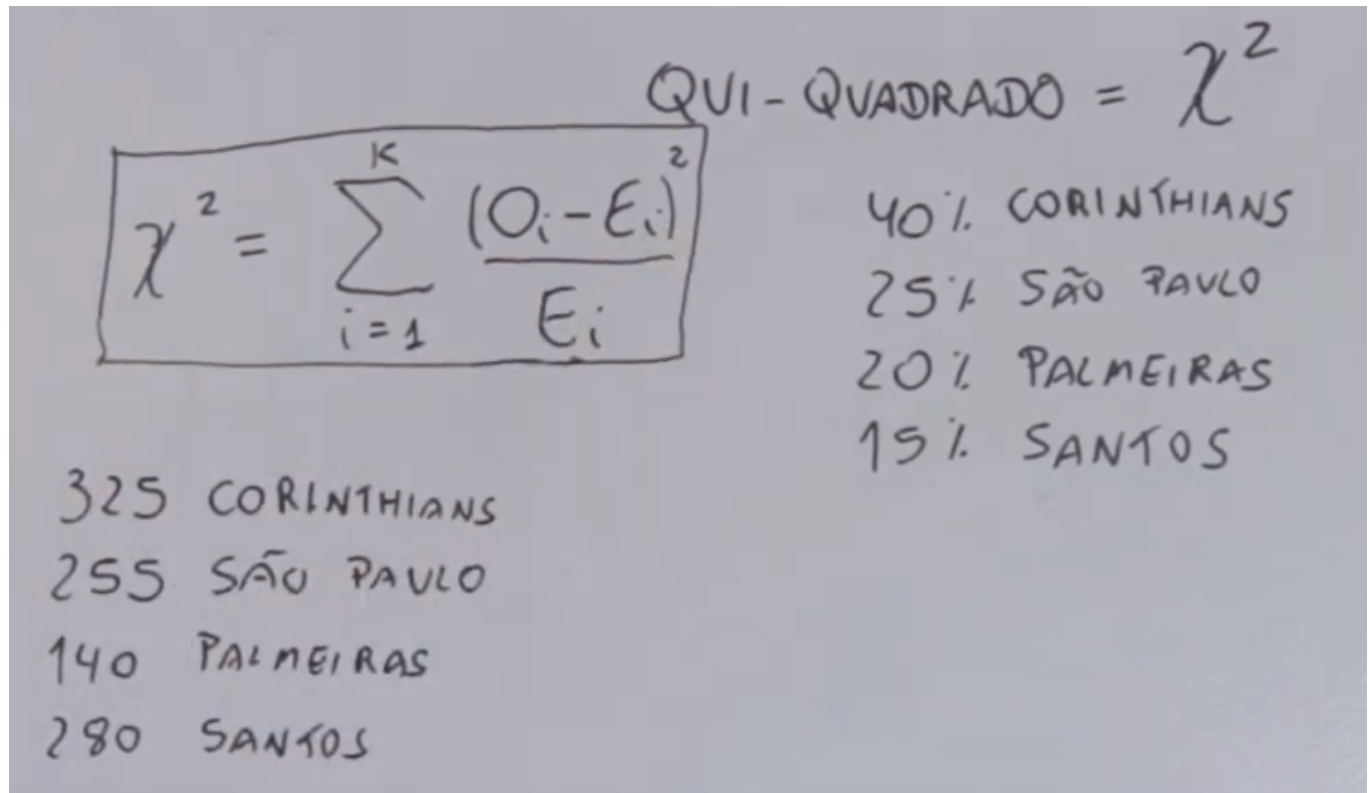
# Qui-Quadrado

Essa técnica visa selecionar as melhores variáveis para o modelo, ou seja, é possível definir as 10 mais correlacionadas com a variável target para que a predição seja feita da melhor forma.

## Exemplo de uso:

Dado um campeonato de futebol brasileiro, alguns times se saíram super bem, um ganhando libertadores, e outros 2 sendo rebaixados. Um estatístico resolveu medir a variação da torcida, se após esse campeonato, tiveram mudanças significativas. Para isso foi utilizado o qui-quadrado.

Após uma pesquisa de 1000 torcedores, temos a seguinte relação:



As porcentagens acima são o resultado esperado para os torcedores de cada clube, ou seja, 40% da torcida é corinthiana, 25% são paulina, e assim por diante. Já os debaixo, foram extraídos da pesquisa de 1000 torcedores, sendo que de 1000, 32,5% são corinthianos, 25,5% são paulinos, e assim por diante.

Aplicando o conceito de qui-quadrado, temos:

$$\chi^2 = \frac{(325-400)^2}{400} + \frac{(255-250)^2}{250} + \frac{(440-200)^2}{200} + \frac{(280-150)^2}{150}$$
$$\chi^2 \approx 145$$

O valor extraído era de 325, mas como o esperado era 40% de 1000, temos 400 torcedores, dividido por 400. Isso é feito para cada clube e somado no final.

Quanto maior a variação do que aconteceu, maior o valor do qui-quadrado.

## **Aplicando o conceito em um dataset:**

QUI - QUADRADO =  $\chi^2$

$$\chi^2 = \sum_{i=1}^K \frac{(O_i - E_i)^2}{E_i}$$

$X_A$	$X_B$	$X_C$	$Y$
12	2	30	1
15	11	6	1
16	8	90	1
5	3	20	0
4	14	5	0
2	5	70	0

$\bar{X}_A = 9 \quad \bar{X}_B = 7,17 \quad \bar{X}_C = 36,83$

$$\chi^2_{X_A} = \frac{(12 + 15 + 16 - 9 \cdot 3)^2}{9 \cdot 3} + \frac{(5 + 4 + 2 - 9 \cdot 3)^2}{9 \cdot 3} = 18,96$$

$$\chi^2_{X_B} = \frac{(2 + 11 + 8 - 7,17 \cdot 3)^2}{7,17 \cdot 3} + \frac{(3 + 14 + 5 - 7,17 \cdot 3)^2}{7,17 \cdot 3} = 0,023$$

$$\chi^2_{X_C} = 4,348$$

Primeiro deve-se calcular a média de cada variável, como mostrado abaixo do dataset.

Depois para cada variável, para calcular o Qui-quadrado, devemos considerar as duas classes da variável target e somar ou seja  $C1 + C2$ . O cálculo de  $C1$  é a soma de cada valor de  $X_a$  por exemplo subtraindo a média vezes a quantidade de target 1 no dataset, ficando:

$$X_{xa} = (12 + 15 + 16 - 9 \cdot 3) / 9 \cdot 3$$

9 é a média, e temos 3 valores 1 observados.

Agora considerando as 3 variáveis já calculadas, podemos definir as 2 melhores, no caso os valores de qui-quadrado de  $X_a$  e  $X_c$  são os maiores, portanto podemos desconsiderar  $X_b$  que não influencia muito na variável target.

O conceito por trás da variação está em analisar a influência dos valores nas classes, por exemplo, podemos ver que quando o valor da variável target é 1, os valores de  $X_a$  são maiores do que quando a target é 0, ou seja, é possível observar a variação.