



# Métricas para os Modelos

Avaliando a qualidade da modelagem

# Métricas

## Regressão

- Mean Absolute Error (MAE)
- Mean Squared Error (MSE)
- Root Mean Squared Error (RMSE)
- $R^2$  (R-Squared)

## Classificação

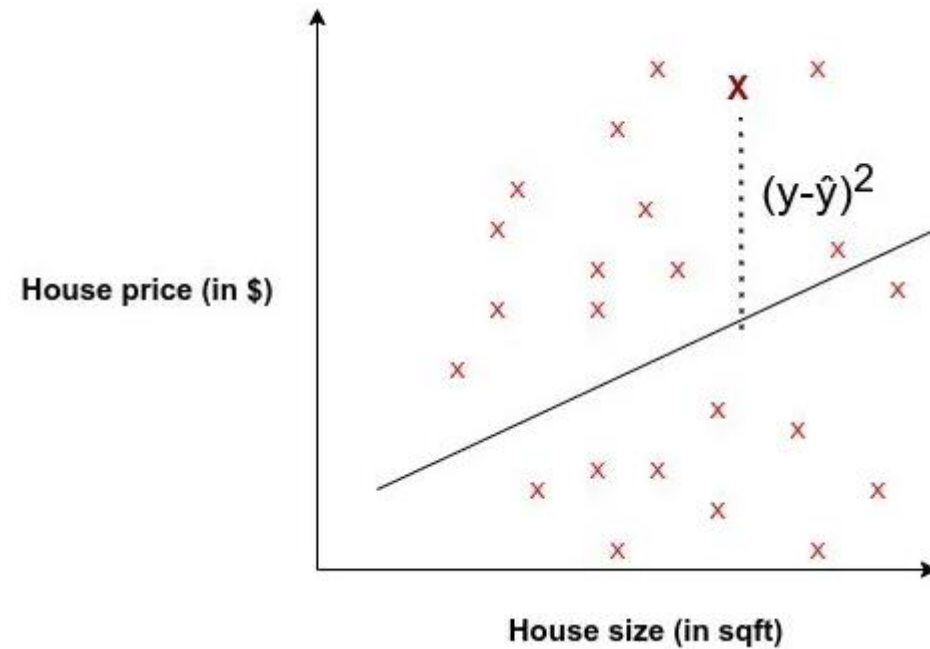
- Accuracy
- Confusion Matrix (not a metric but fundamental to others)
- Precision and Recall
- F1-score
- AU-ROC

## Clusterização

- Rand index
- Mutual Information based scores
- Homogeneity, completeness and V-measure
- Fowlkes-Mallows scores
- Silhouette Coefficient
- Calinski-Harabasz Index
- Contingency Matrix
- Pair Confusion Matrix

# Mean Squared Error

$$MSE = \frac{1}{N} \sum_{j=1}^N (y_j - \hat{y}_j)^2$$



# Índice de exatidão (acurácia)

A porcentagem do quanto o método acertou dentro do conjunto total de amostras na base de dados

$$Exatidão = \frac{\sum_{i=1}^m A_{nn}}{\sum_{i=1}^m \sum_{j=1}^m A_{ij}}$$

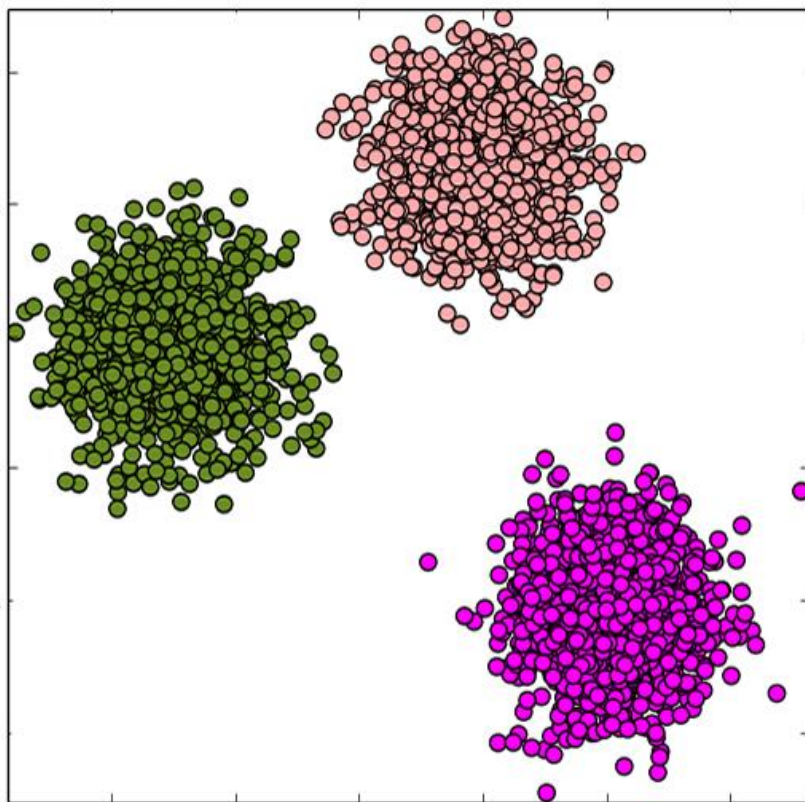
# Matriz de Confusão

Classes	1	2	3	4	Amostras incluídas (comissionadas)	Total de amostras na linha	Erro de comissão (%)
1	33	4	2	1	7	40	17,5
2	2	35	3	0	5	40	12,5
3	0	3	34	3	6	40	15
4	0	1	2	37	3	40	7,5

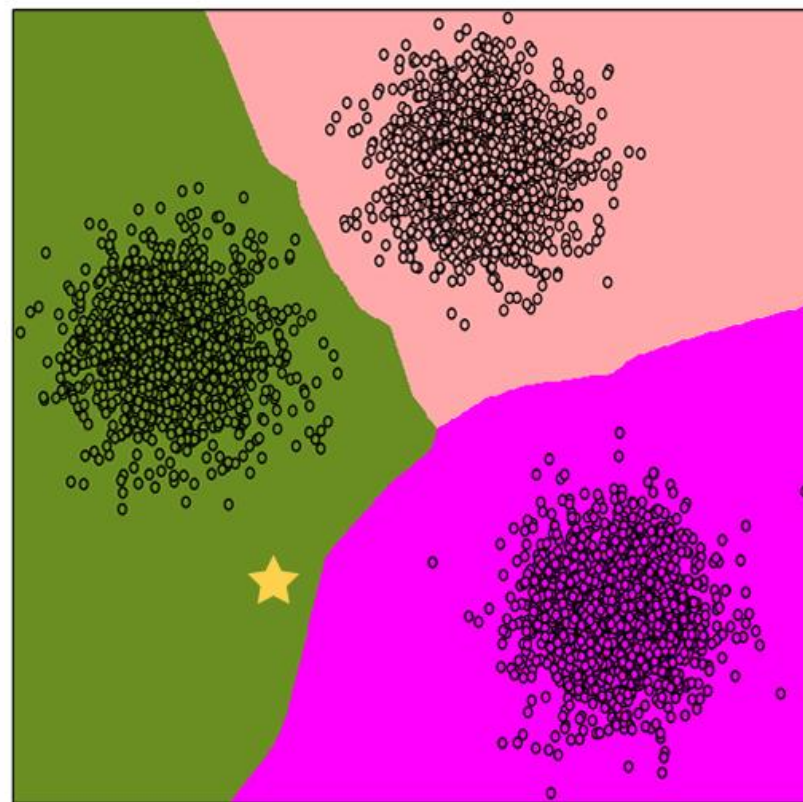
Amostras omitidas	2	8	7	4
Total de amostras na coluna	35	43	41	41
Erro de omissão (%)	5,7	19	17	10

# Fronteiras de decisão

Conjunto de dados com três classes e 4000 atributos



Fronteira de decisão do conjunto de dados



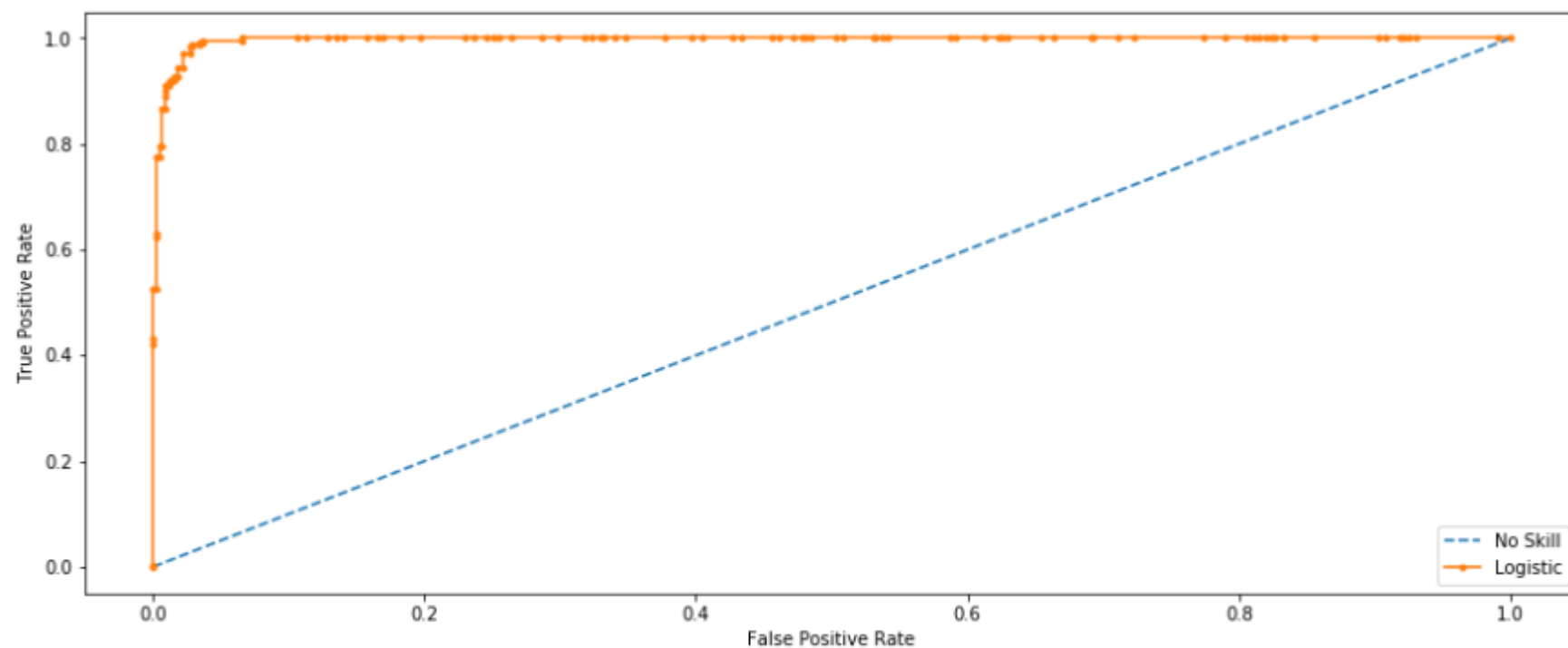
# F1-score

$$F_1 = \frac{2}{\frac{1}{precision} + \frac{1}{recall}}$$

$$R = \frac{TP}{TP+FN} = \frac{\text{Cancer patients correctly identified}}{\text{Cancer patients correctly identified} + \text{incorrectly labelled non-cancer patients as cancerous}}$$

$$P = \frac{TP}{TP+FP} = \frac{\text{Cancer patients correctly identified}}{\text{Cancer patients correctly identified} + \text{incorrectly labelled cancer patients as non-cancerous}}$$

# AU-ROC



*No Skill: ROC AUC=0.500*

*Logistic: ROC AUC=0.996*

$$\text{TPR} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

$$\text{FPR} = \frac{\text{FP}}{\text{FP} + \text{TN}}$$



# Índice Rand Corrigido

Um mais usados para validação externa de agrupamentos e determina a similaridade de partições.

varia no intervalo  $[-1,1]$ , onde quanto mais próximo de 1 mais similares são as partições comparadas



# Semantix<sup>®</sup>

All about data

[hader.azzini@semantix.com.br](mailto:hader.azzini@semantix.com.br)

[www.semantix.com.br](http://www.semantix.com.br)