



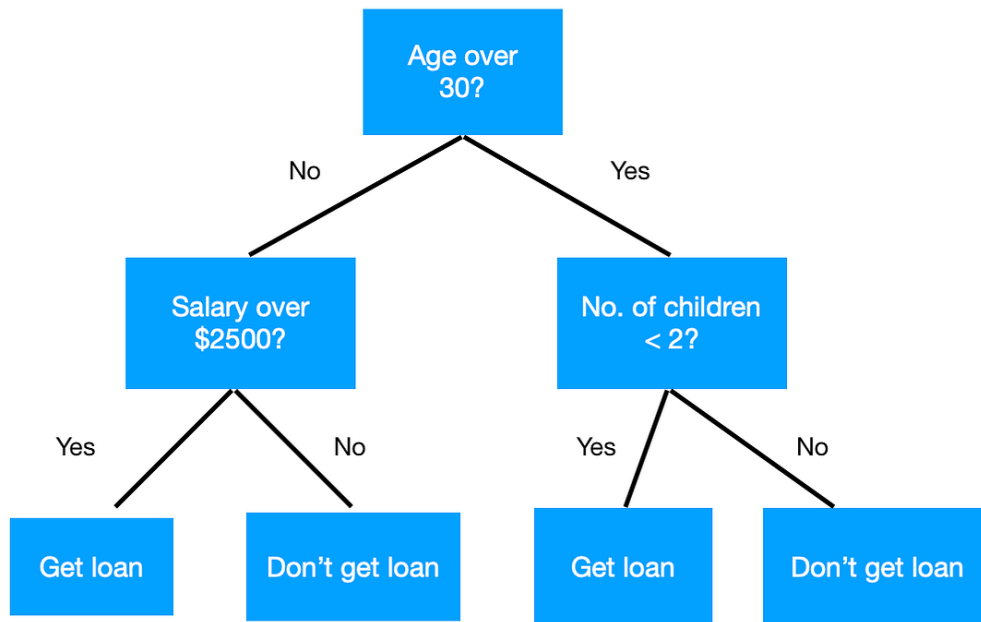
Fundamentos de Machine Learning

Comitês de Modelos

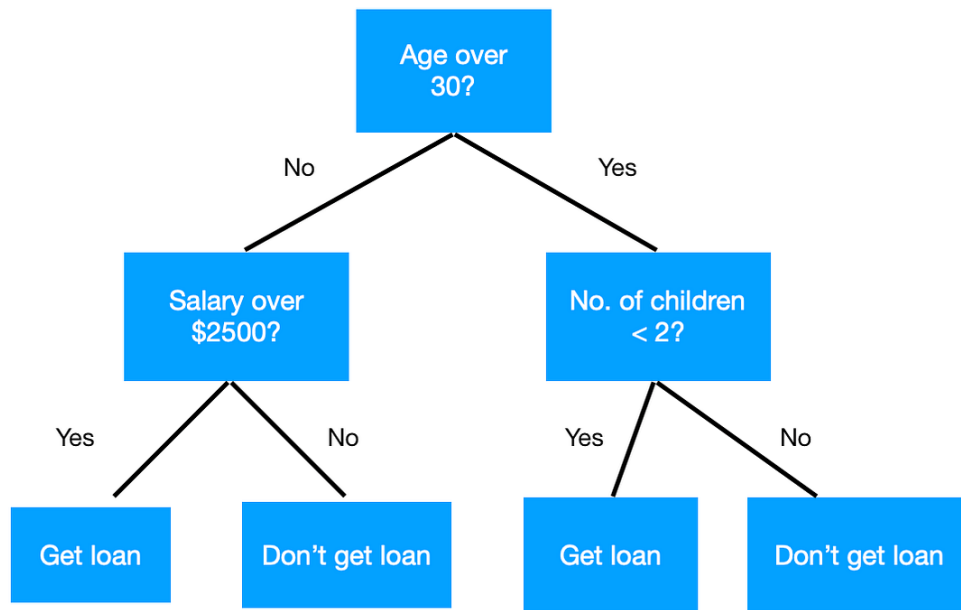
Aula 1.1 Medindo erros
Prof. Pedro Calais



Lembra da árvore de decisão?



- Fácil de construir
- Fácil de usar
- Fácil de interpretar



Problema: árvores de decisão são instáveis

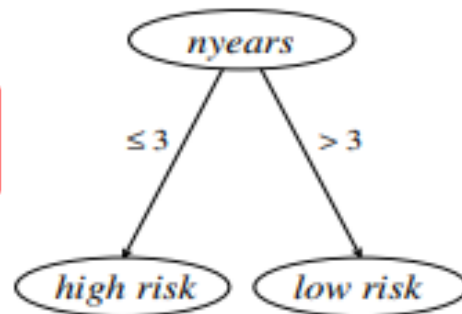
Instability of Decision Tree Classification Algorithms

Ruey-Hsia Li
Lightspeed Semiconductor
209 N. Fair Oaks Avenue
Sunnyvale, CA 94085
rli@lightspeed.com

Geneva G. Belford
Department of Computer Science
University of Illinois at Urbana-Champaign
Urbana, IL 61801
belford@cs.uiuc.edu

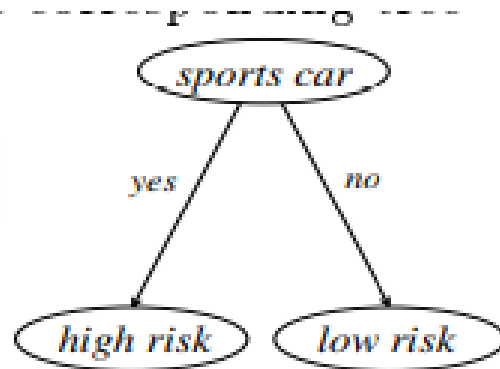
Um pequeno conjunto de dados

nyears	sports car	risk
1	yes	high
2	no	high
2	yes	high
4	no	low
8	no	low



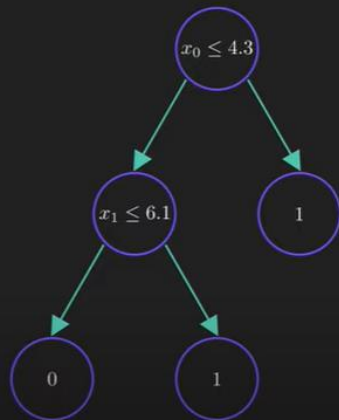
E se o registro #2 for diferente?

nyears	sports car	risk
1	yes	high
4	yes	high
2	yes	high
4	no	low
8	no	low



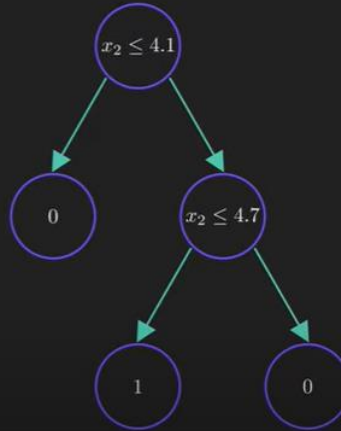
Mais um exemplo

<i>id</i>	x_0	x_1	x_2	x_3	x_4	y
0	4.3	4.9	4.1	4.7	5.5	0
1	3.9	6.1	5.9	5.5	5.9	0
2	2.7	4.8	4.1	5.0	5.6	0
3	6.6	4.4	4.5	3.9	5.9	1
4	6.5	2.9	4.7	4.6	6.1	1
5	2.7	6.7	4.2	5.3	4.8	1



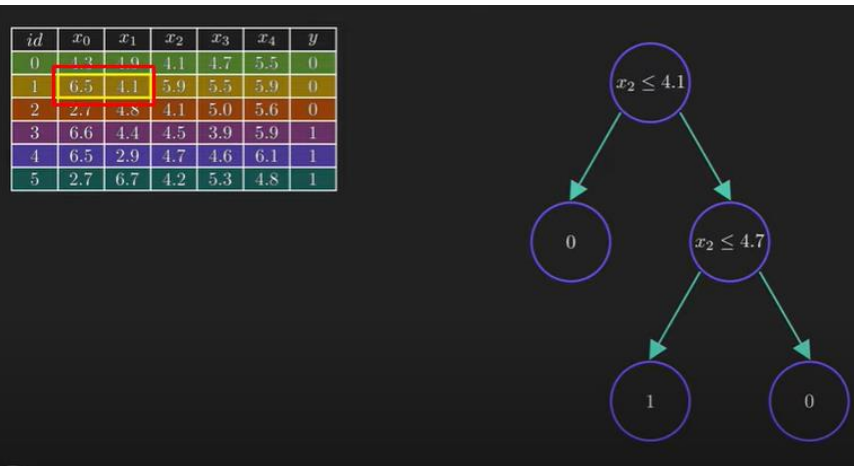
Introduzindo uma pequena mudança

id	x_0	x_1	x_2	x_3	x_4	y
0	1.3	1.0	4.1	4.7	5.5	0
1	6.5	4.1	5.9	5.5	5.9	0
2	2.1	4.8	4.1	5.0	5.6	0
3	6.6	4.4	4.5	3.9	5.9	1
4	6.5	2.9	4.7	4.6	6.1	1
5	2.7	6.7	4.2	5.3	4.8	1



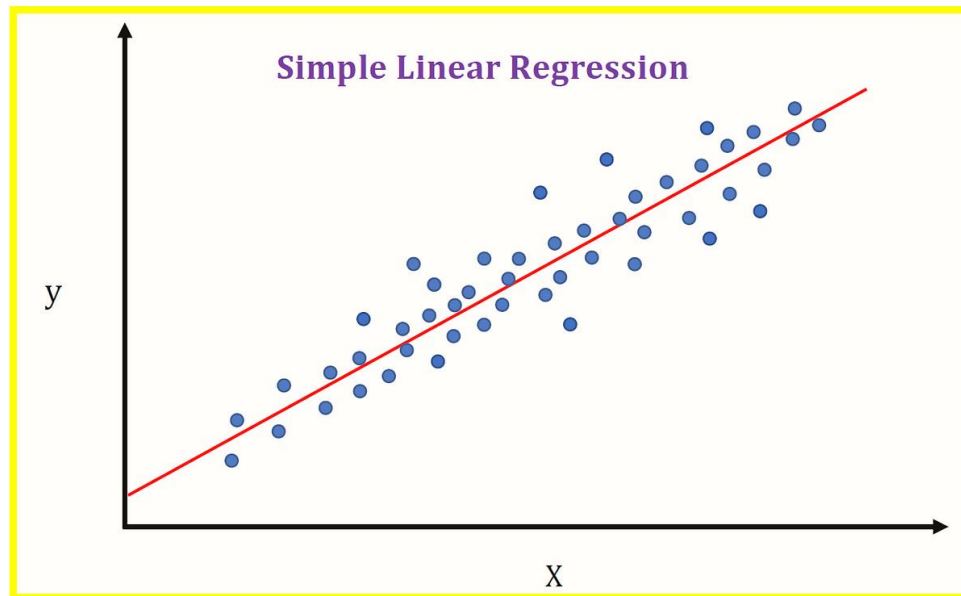
Árvores de decisão têm alta variância

- Pequenas mudanças nos dados ou nos hiperparâmetros
- → Modelos muito diferentes



Contraponto

- Regressão linear tende a ter baixa variância

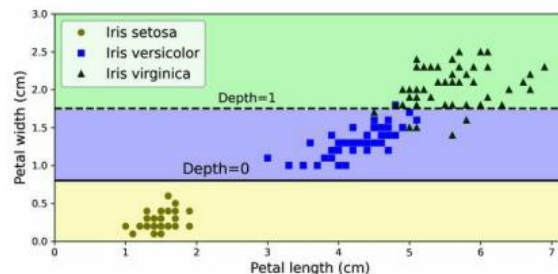


Árvores de decisão têm alta variância



Decision Trees Have a High Variance

More generally, the main issue with decision trees is that they have quite a **high variance**: small changes to the hyperparameters or to the data may produce very different models. In fact, since the training algorithm used by Scikit-Learn is stochastic—it randomly selects the set of features to evaluate at each node—even retraining the same decision tree on the exact same data may produce a very different model, such as the one represented in **Figure 6-9** (unless you set the `random_state` hyperparameter). As you can see, it looks very different from the previous decision tree (**Figure 6-2**).



Por que alta variância é um problema?

- Imprevisibilidade
- Aprendendo o ruído
- Treina o modelo várias vezes, resultados muito diferentes
- Difícil validar
- Modelo está se ajustando em excesso?

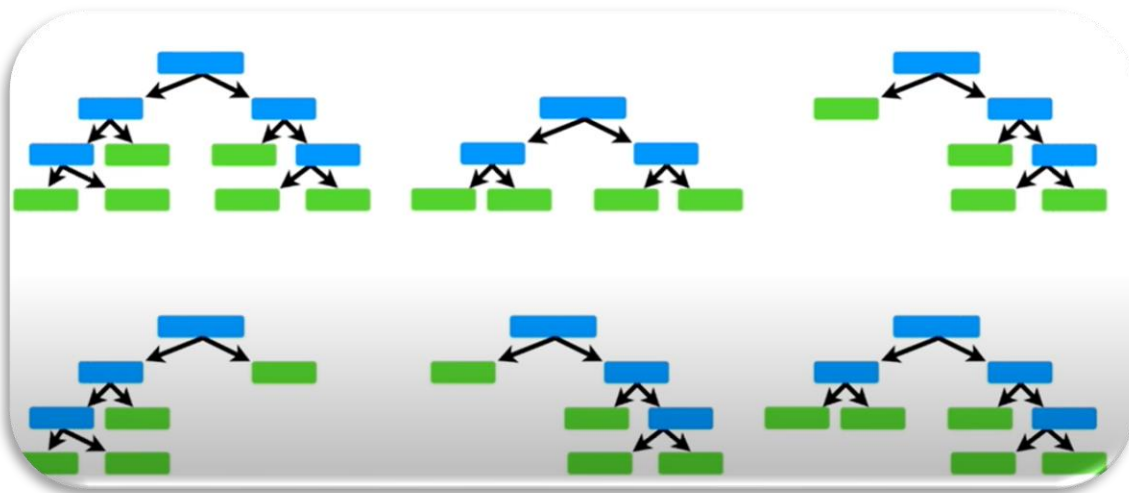


Como reduzir a variância de árvores de decisão?



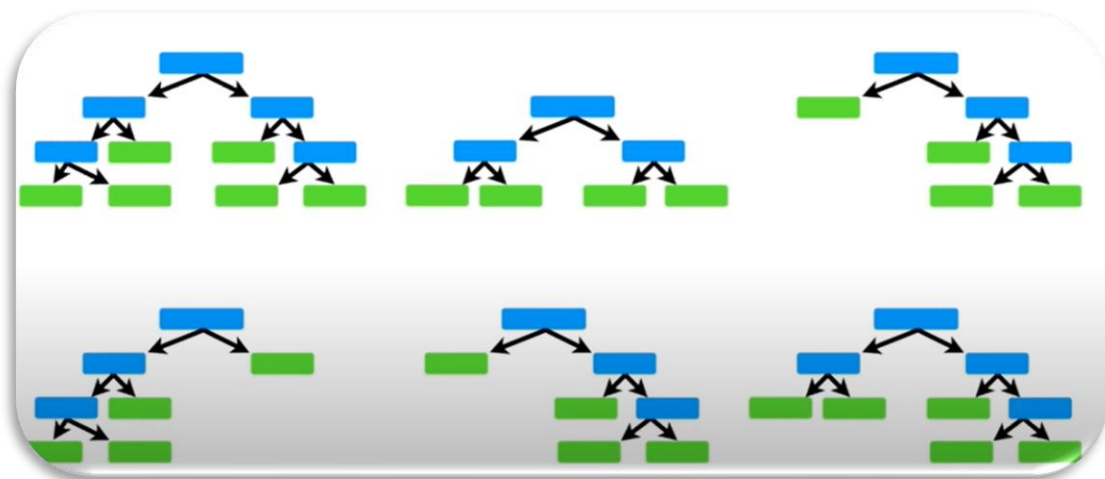
Como reduzir a variância de árvores de decisão?

- Gere várias árvores
- Tire a média das previsões



As árvores precisam ser diferentes, certo?

- Selecione os registros aleatoriamente
- Selecione colunas aleatoriamente



Aleatorize os registros

id	x_0	x_1	x_2	x_3	x_4	y
0	4.3	4.9	4.1	4.7	5.5	0
1	3.9	6.1	5.9	5.5	5.9	0
2	2.7	4.8	4.1	5.0	5.6	0
3	6.6	4.4	4.5	3.9	5.9	1
4	6.5	2.9	4.7	4.6	6.1	1
5	2.7	6.7	4.2	5.3	4.8	1

id
2
0
2
4
5
5

id
2
1
3
1
4
4

id
4
1
3
0
0
2

id
3
3
2
5
1
2

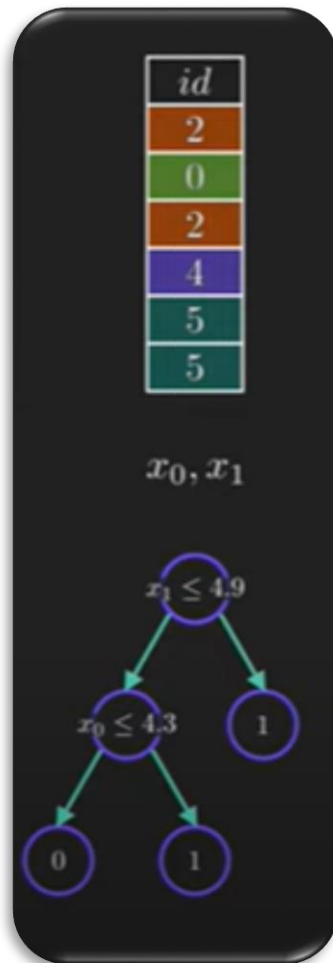


Aleatorize os atributos

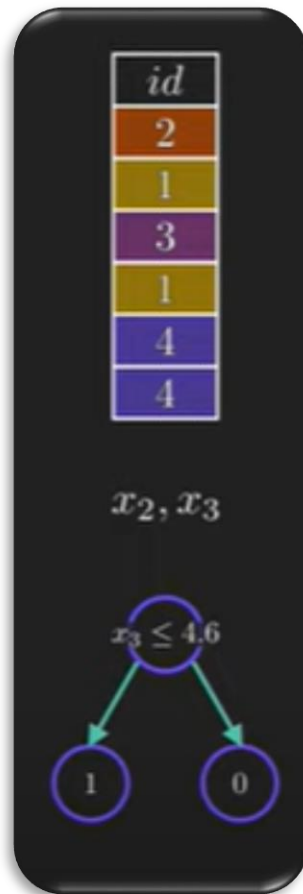
<table><tr><th><i>id</i></th></tr><tr><td>2</td></tr><tr><td>0</td></tr><tr><td>2</td></tr><tr><td>4</td></tr><tr><td>5</td></tr><tr><td>5</td></tr></table> x_0, x_1	<i>id</i>	2	0	2	4	5	5	<table><tr><th><i>id</i></th></tr><tr><td>2</td></tr><tr><td>1</td></tr><tr><td>3</td></tr><tr><td>1</td></tr><tr><td>4</td></tr><tr><td>4</td></tr></table> x_2, x_3	<i>id</i>	2	1	3	1	4	4	<table><tr><th><i>id</i></th></tr><tr><td>4</td></tr><tr><td>1</td></tr><tr><td>3</td></tr><tr><td>0</td></tr><tr><td>0</td></tr><tr><td>2</td></tr></table> x_2, x_4	<i>id</i>	4	1	3	0	0	2	<table><tr><th><i>id</i></th></tr><tr><td>3</td></tr><tr><td>3</td></tr><tr><td>2</td></tr><tr><td>5</td></tr><tr><td>1</td></tr><tr><td>2</td></tr></table> x_1, x_3	<i>id</i>	3	3	2	5	1	2
<i>id</i>																															
2																															
0																															
2																															
4																															
5																															
5																															
<i>id</i>																															
2																															
1																															
3																															
1																															
4																															
4																															
<i>id</i>																															
4																															
1																															
3																															
0																															
0																															
2																															
<i>id</i>																															
3																															
3																															
2																															
5																															
1																															
2																															



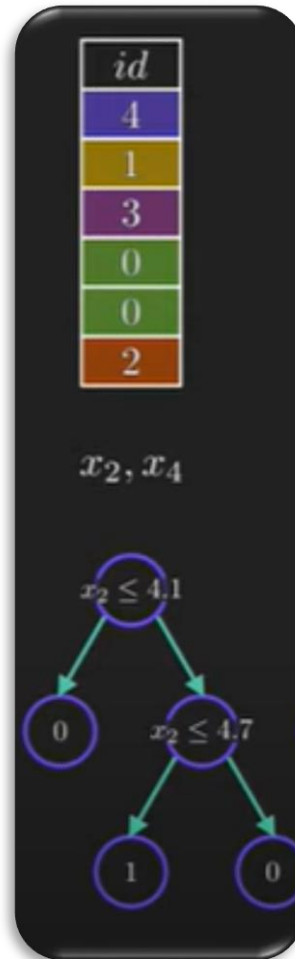
Árvore # 1



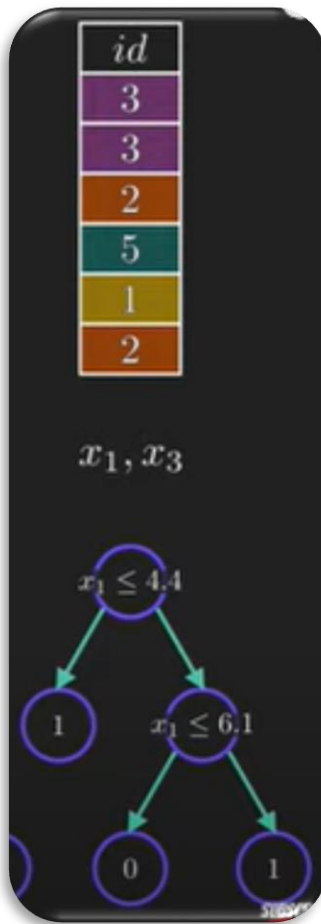
Árvore # 2



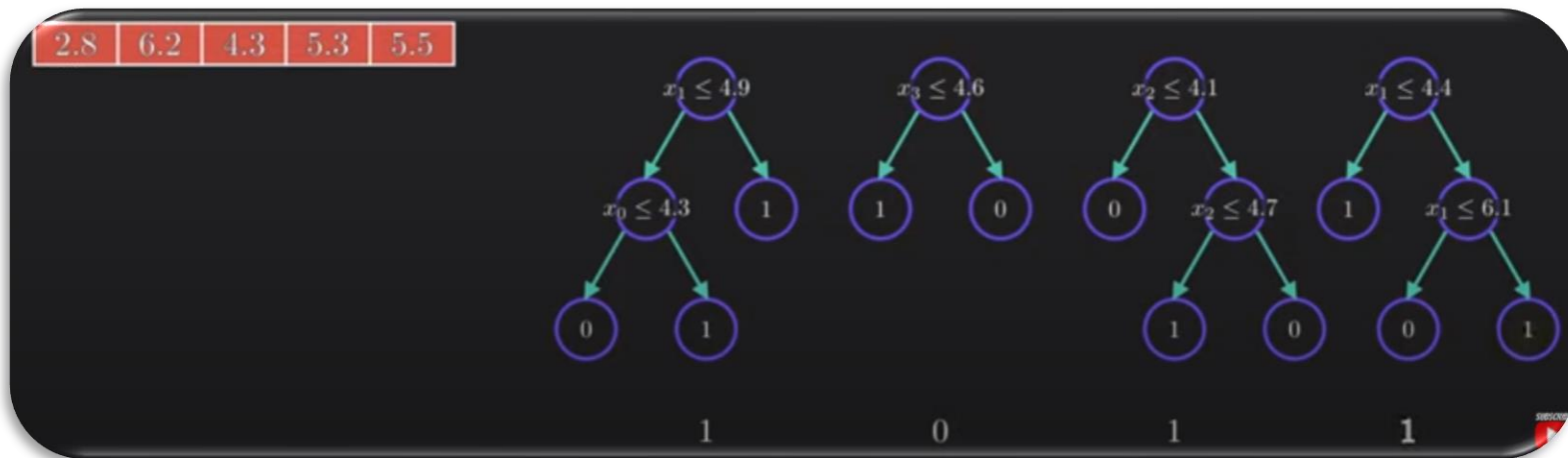
Árvore # 3



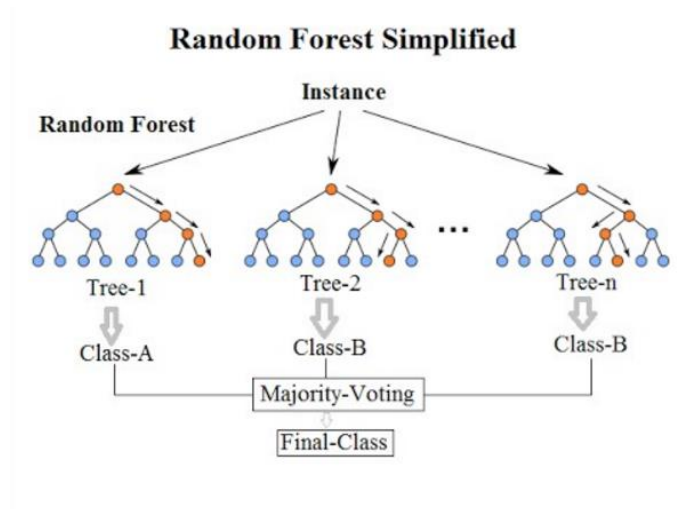
Árvore # 4



Como fazer uma previsão?

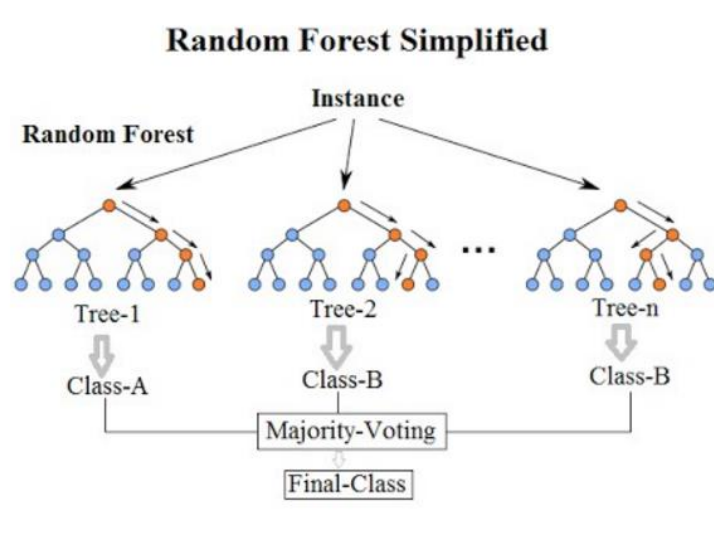


Como fazer uma previsão?



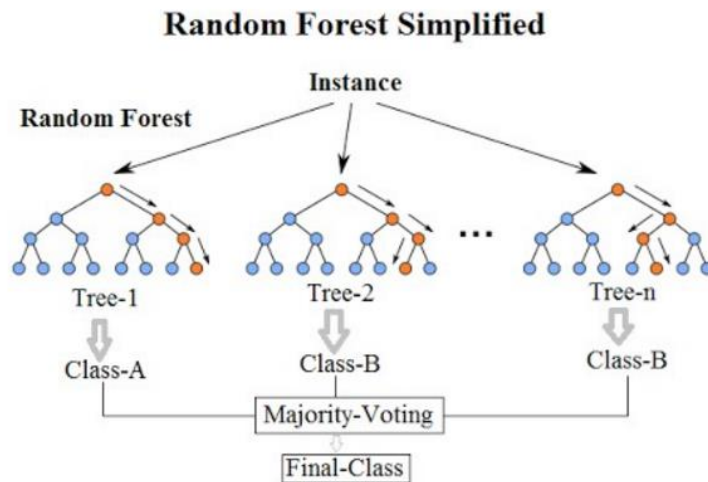
Por que a floresta tende a diminuir a variância?

- Várias árvores diminuem a sensibilidade a registros e colunas específicas
- Menos *overfitting*!



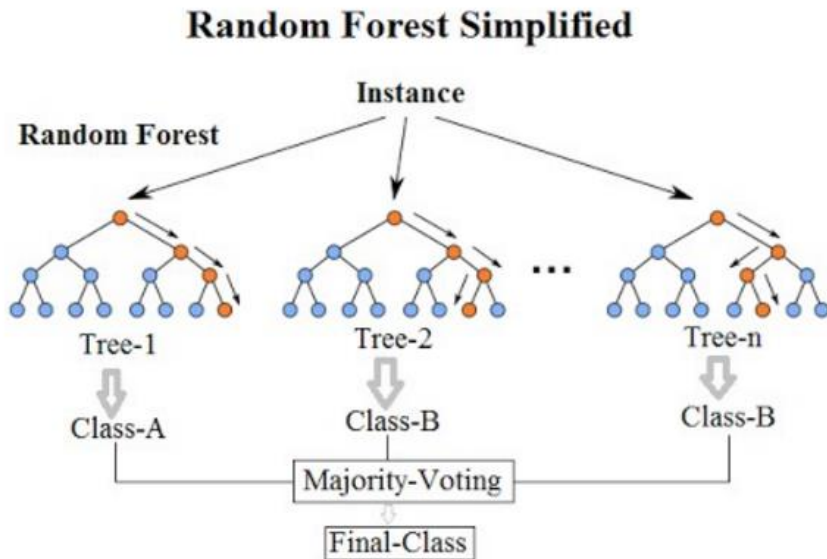
Por que a floresta tende a diminuir a variância?

Fraquezas individuais de cada árvore tendem a se anular



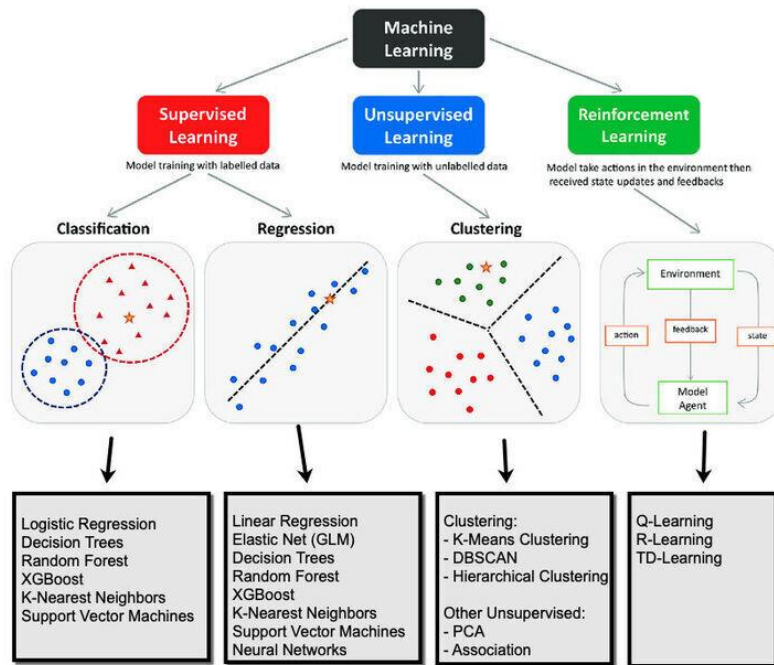
Conclusão

Bem-vindo às *random forests* :-)

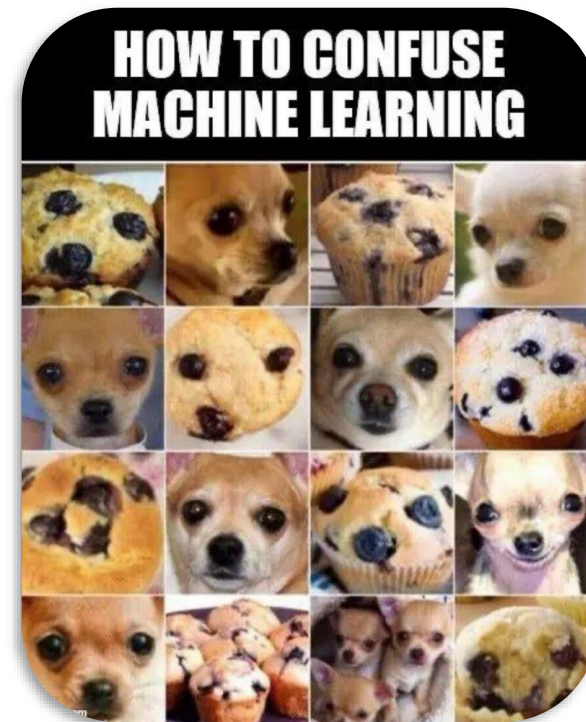


Os tipos de aprendizado de máquina

3 Types of Machine Learning (Every Data Scientist Should Know)



Como a máquina se confunde?



Modelos e fronteiras de decisão

Classification Model *Decision Boundaries*

DataInterview.com

How do ML classifiers form *decision boundaries*?

