



Insper

Machine Learning

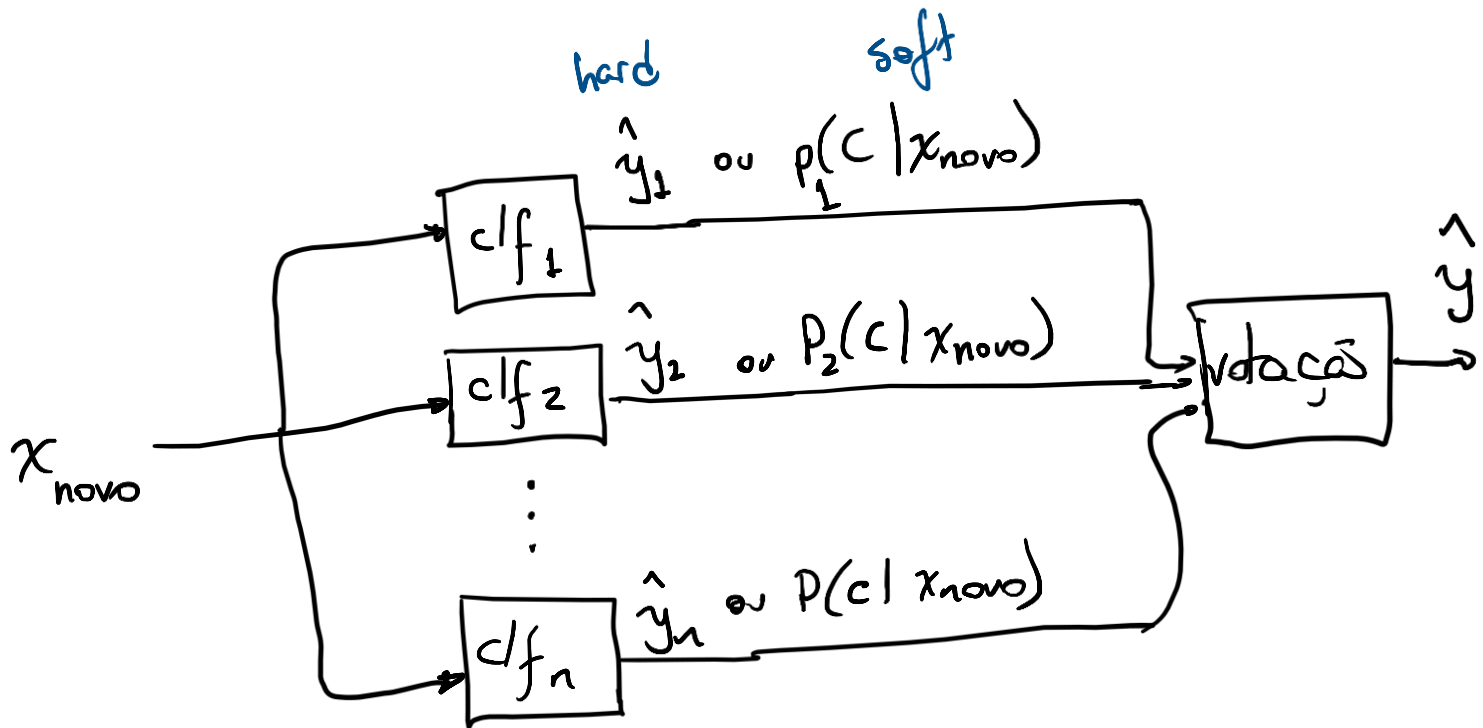
Ensemble Models

Fábio Ayres <fabioja@insper.edu.br>

- Wisdom of the crowds
 - múltiplos modelos "fracos" em modelo "forte"
- Requer:
 - prob. de acerto $> 50\%$ p/cada modelo fraco
 - num. grande de modelos p/cancelar o erro
 - independência
 - ↳ efeito manada não existe
 - ↳ $P(A|B) = P(A)$
 - ↳ Algoritmos diferentes, mesmos dados
 - ↳ Mesmo algoritmo, dados diferentes.

Voting Classifier

3



Bootstrap

4

- distribuição desconhecida F (população)

- Goleto m amostras i.i.d.

$(x_1, x_2, \dots, x_m) \leftarrow$ conjunto de observações

- Quero estimar $\bar{F} = E_F[X]$

$$\bar{F} = \frac{1}{n} \sum_{i=1}^m x_i \quad (\text{media amostral})$$

Considere o seguinte:

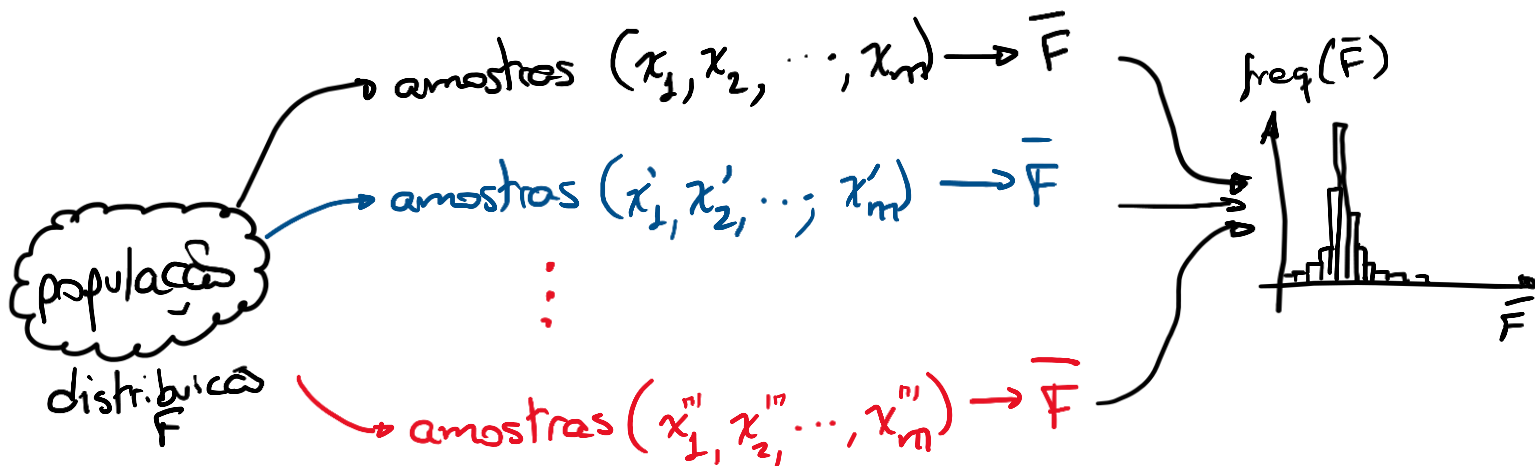
- cada x_i foi amostrado de F

- Se eu repetir a amostragem, obtenho outros x_i
 \Rightarrow obtenho outros \bar{F}

$\Rightarrow \bar{F}$ é variável aleatória \Rightarrow tem uma distribuição!

Histograma de $\bar{F} = \frac{1}{m} \sum_{i=1}^m X_i$ onde X_i v.a. com distribuição F

em outras palavras:



Exemplo: Será que o pãozinho da padaria pesa em média mais que 50g? (6)

=> teste de hipótese!

1) $H_0: \mu \leq 50$
 $H_A: \mu > 50$

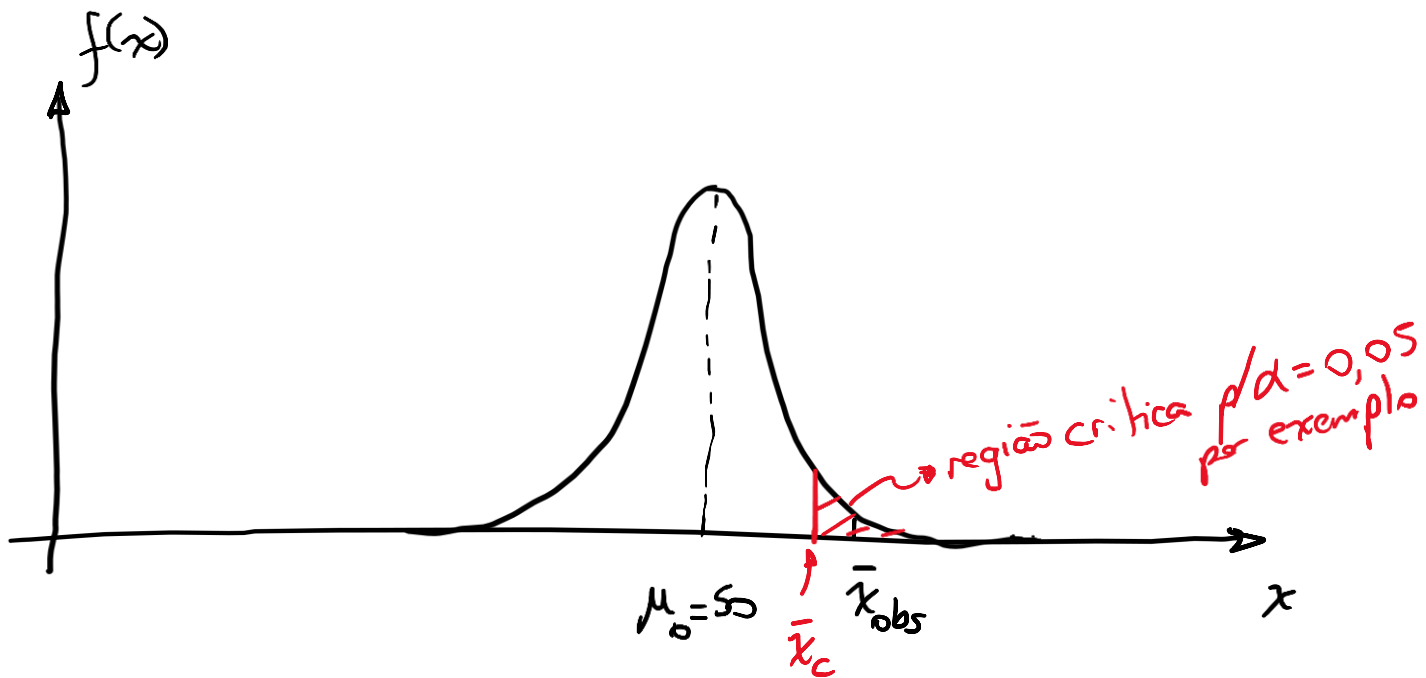
← média populacional (desconhecida)

requer populações
com distr.
normal!

2) Estatística de teste: estatística t-student

3) Critério de rejeição de H_0 :
↳ escolher um α
↳ Região crítica
↳ valor-p

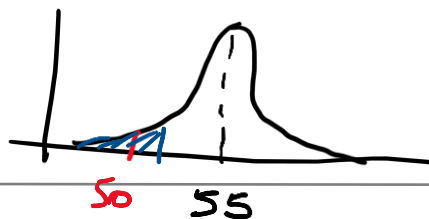
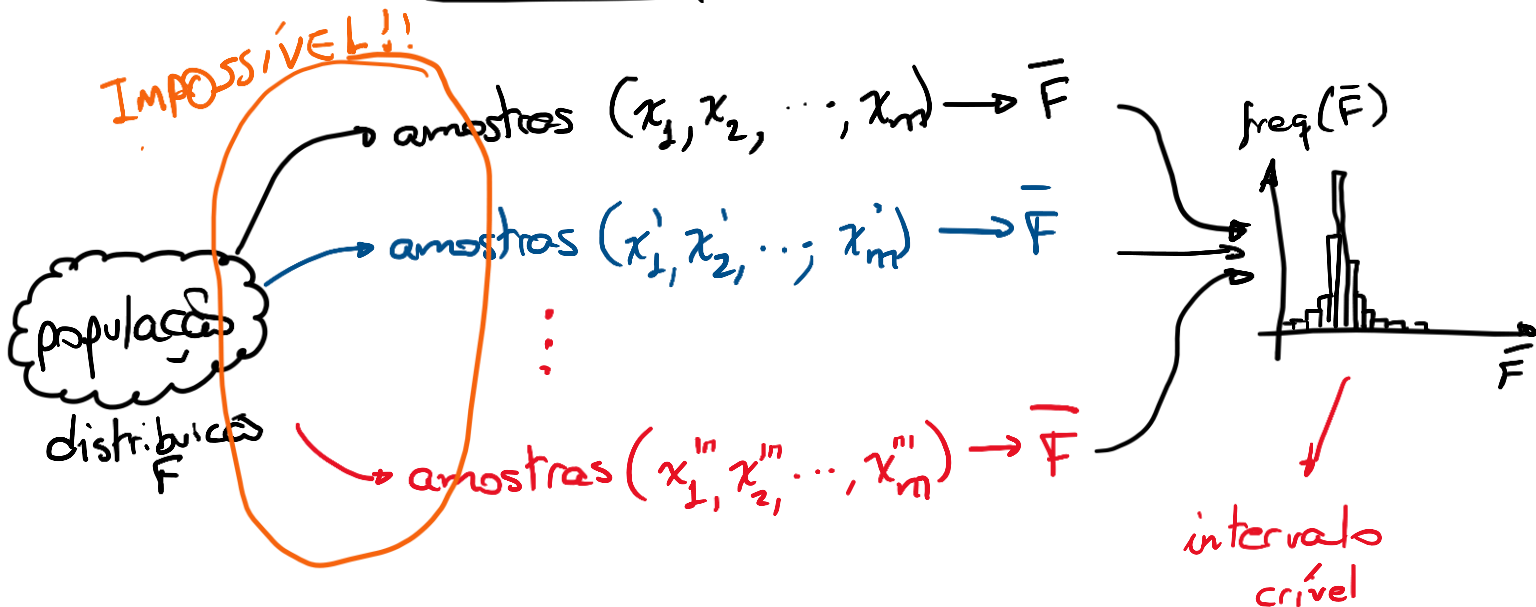
4) Calcula \bar{X} → verifica se H_0 foi rejeitada



É se não quero supor $x_i \sim \text{Normal}(\mu, \sigma^2)$?
Ou nenhuma outra?

8

\Rightarrow métodos não-paramétricos!

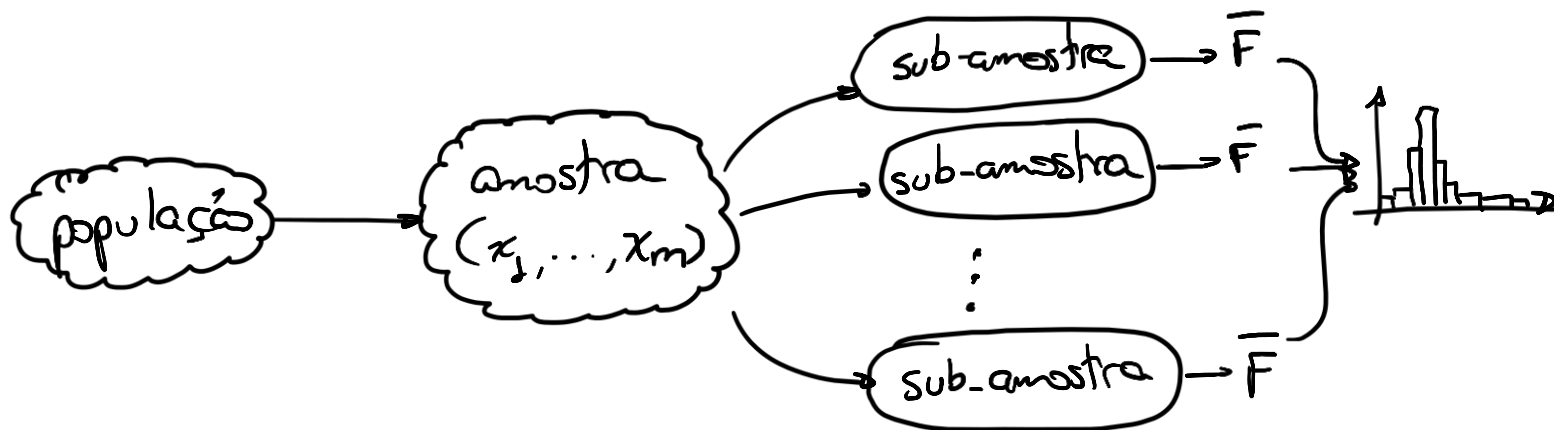


E se eu tenho apenas 1 conjunto de observações?

9

BOOTSTRAP!

- Gerar conjuntos de observações a partir de um único conjunto através de amostragem com repetições

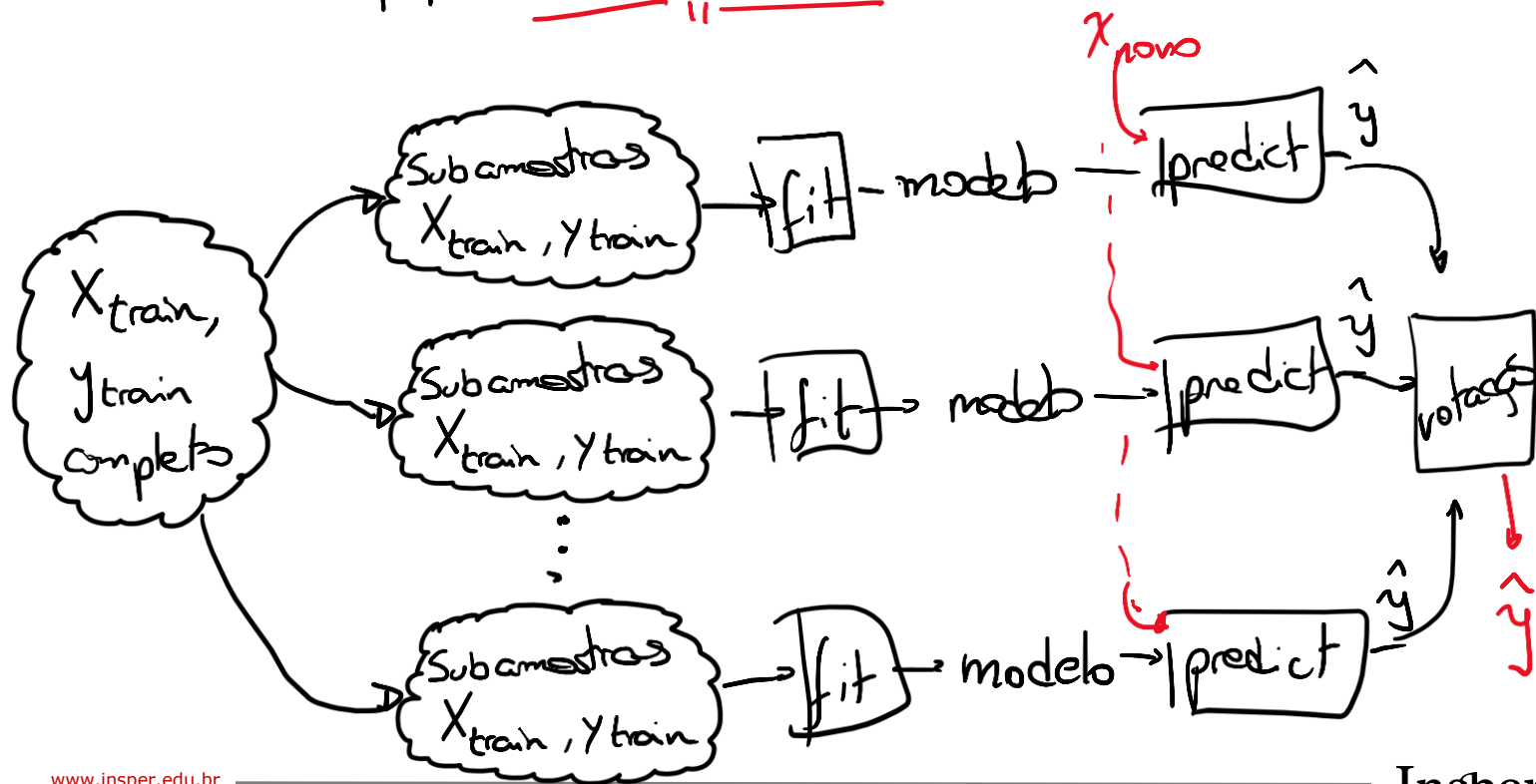


"to pull oneself up by one's own bootstraps"
se erguer sozinho a partir do nada

Bradley Efron (1979)

Bootstrap: permite obter amostras de uma estatística de interesse a partir de um conjunto de observações, sem super nada sobre a distribuição populacional

10



Bagging : Bootstrap aggregating (com repetições)

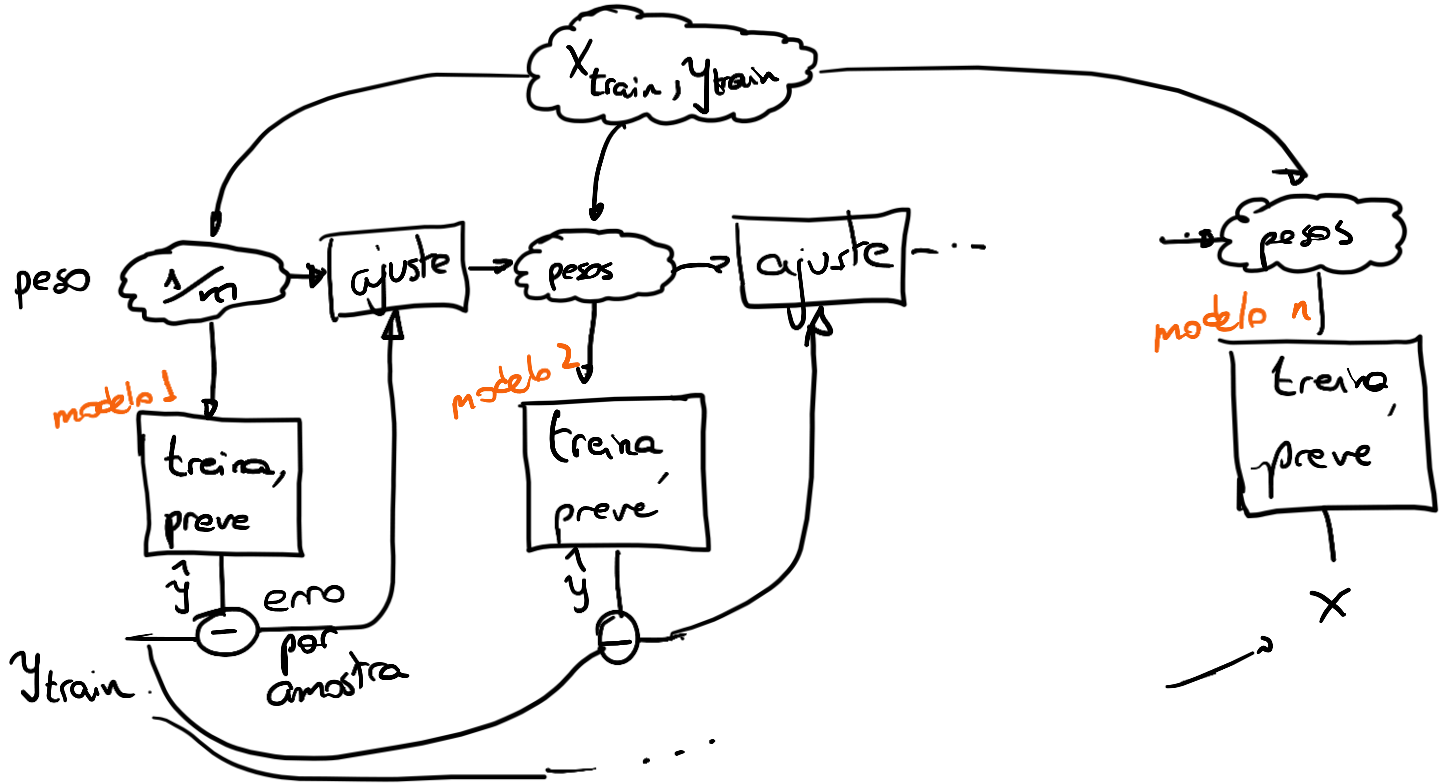
↓
Pasting : o mesmo, mas sem repetições

Boosting: aprendizado sequencial de uma coleção de modelos

12

AdaBoost

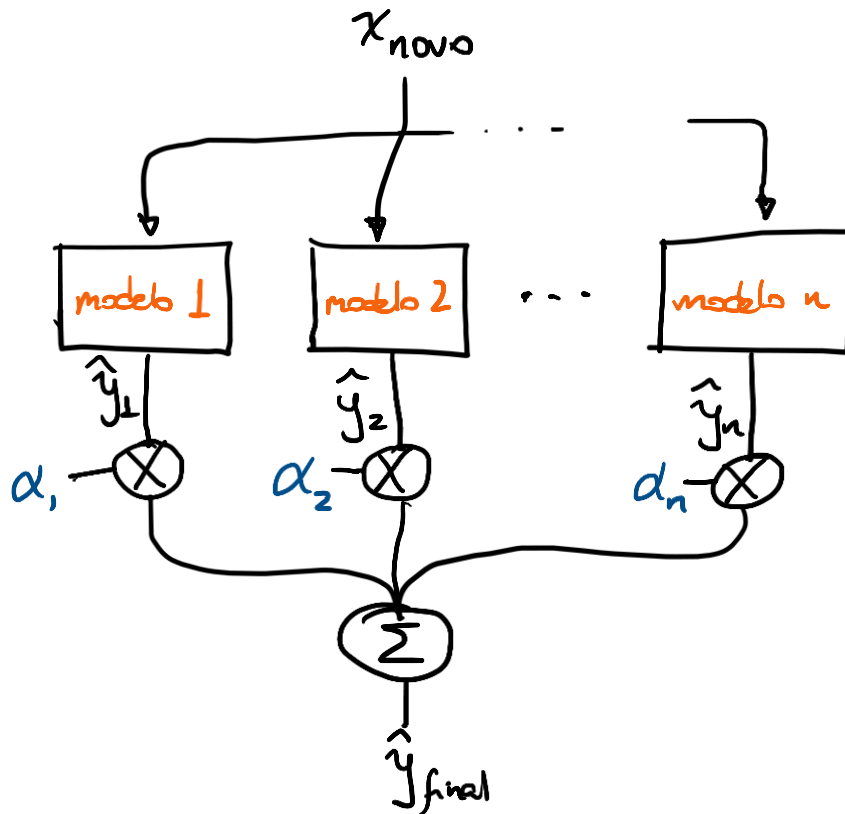
treinamento



AdaBoost

predições

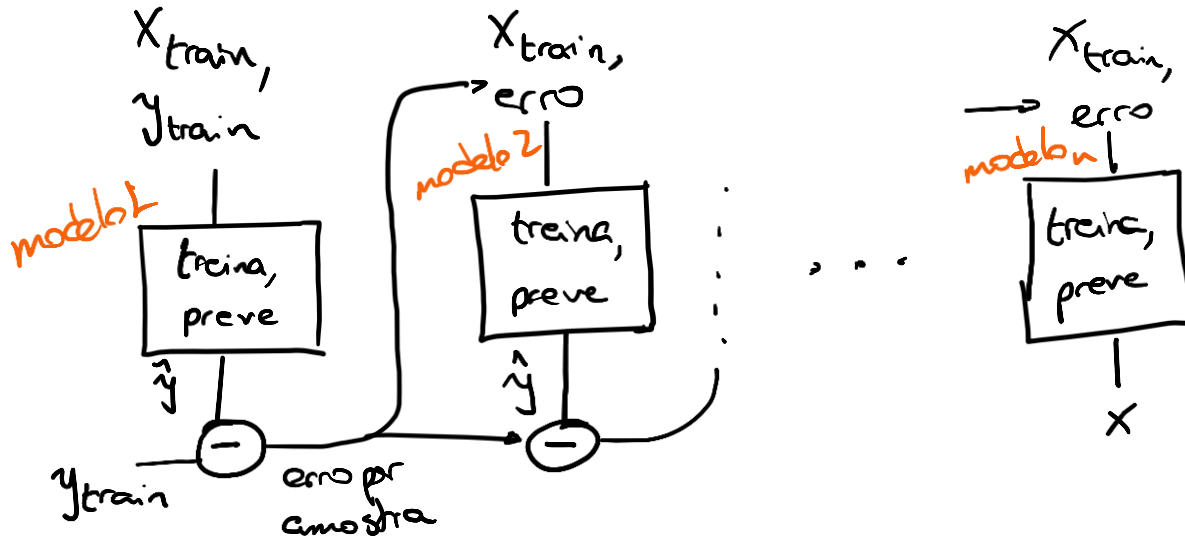
13



Gradient Boosting:

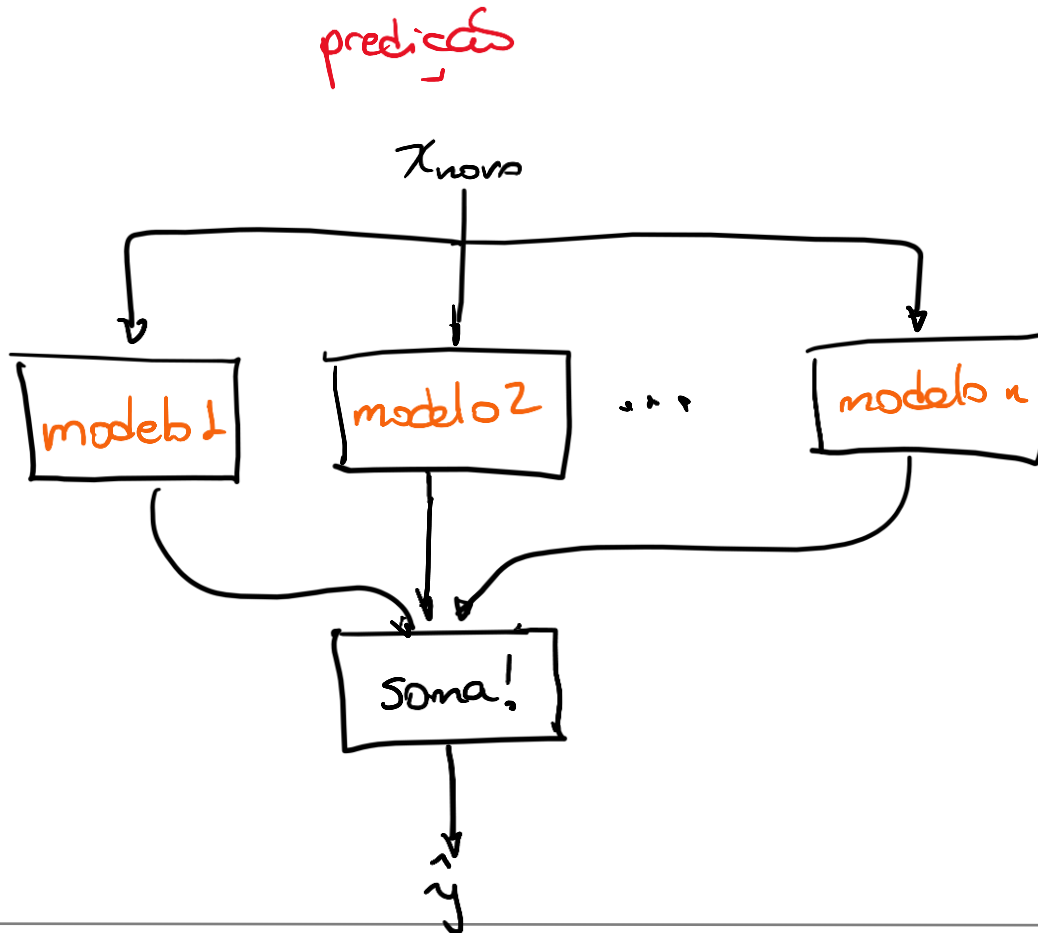
14

treinamento



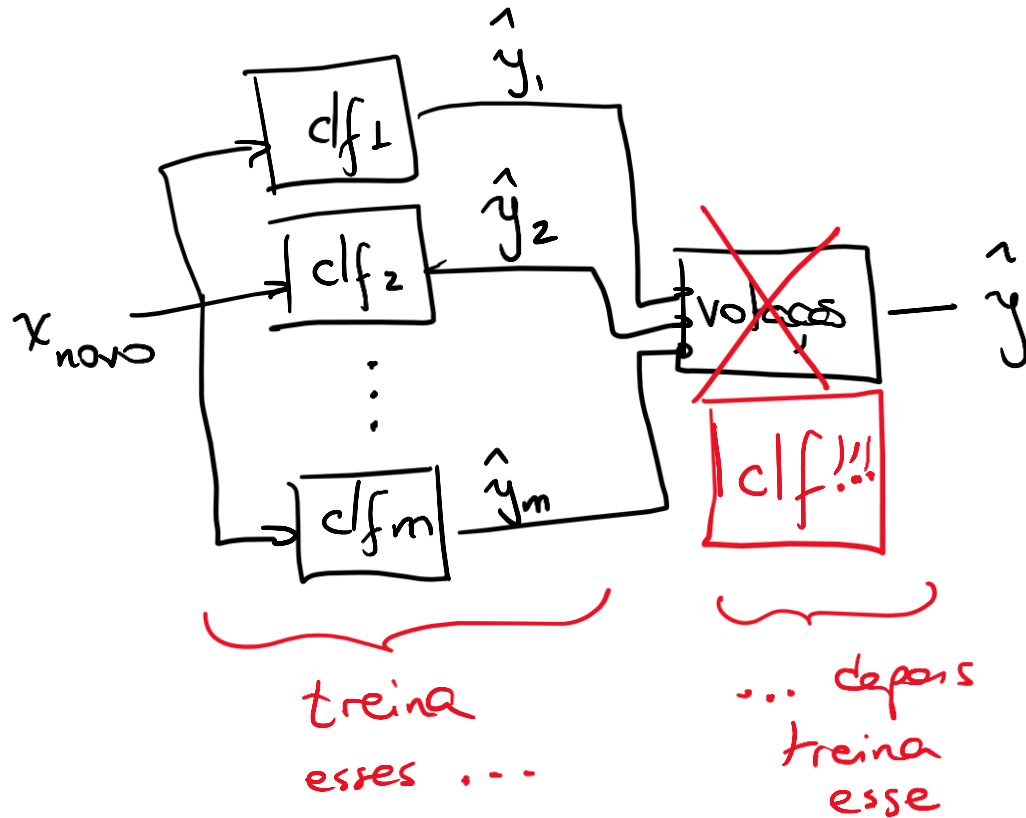
Gradient Boosting

15



Stacking

16



The background of the slide is composed of several concentric, partial arcs in red and grey, creating a dynamic, circular pattern. The arcs vary in thickness and are positioned at different radii from the center, which is where the word 'Insper' is located. The overall effect is a modern, geometric design.

Insper