



Classificação

Identificando a qual categoria um objeto pertence

Classificação

Métodos

- Logistic Regression
- Stochastic Gradient Descent
- Naïve Bayes
- K-Nearest Neighbors
- Decision Tree
- Random Forest
- Support Vector Machine

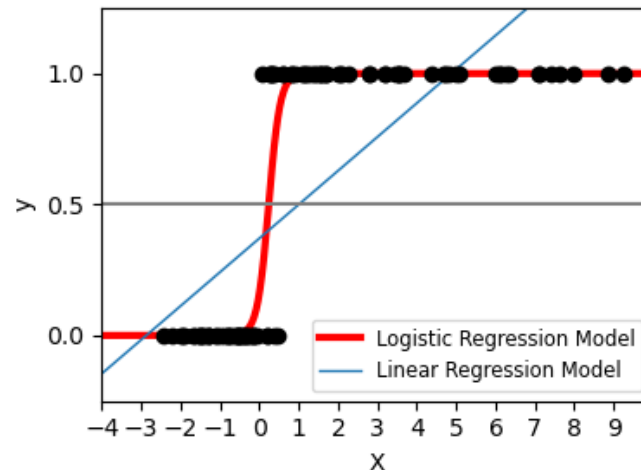
Regressão Logística

A regressão logística, apesar do nome, é um modelo linear para classificação, em vez de regressão. Uma curva em “S”

$\text{odds} = p / (1-p)$ = probability of event occurrence / probability of not event occurrence

$\ln(\text{odds}) = \ln(p/(1-p))$

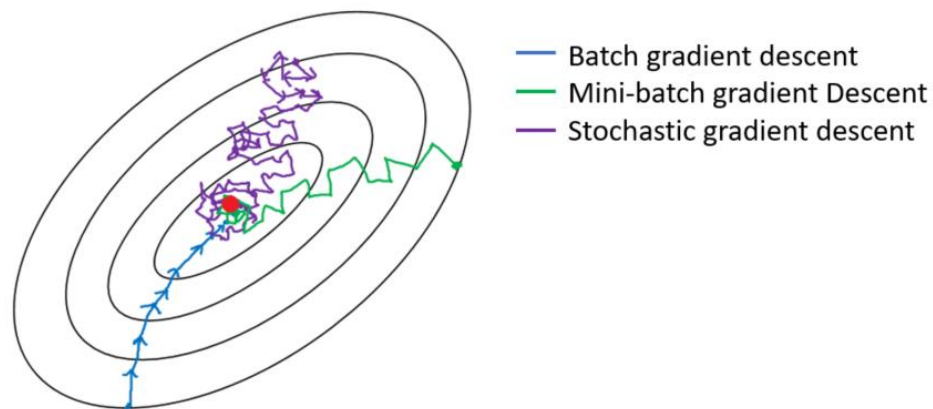
$\text{logit}(p) = \ln(p/(1-p)) = b_0 + b_1X_1 + b_2X_2 + b_3X_3 + \dots + b_kX_k$



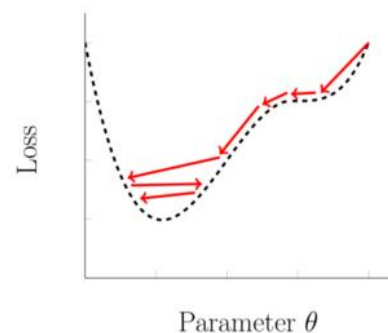
Stochastic Gradient Descent

A rigor, SGD é apenas uma técnica de otimização e não corresponde a uma família específica de modelos de aprendizado de máquina. É apenas uma forma de treinar um modelo.

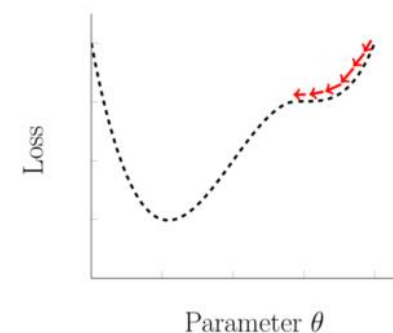
É uma abordagem simples, mas muito eficiente para ajustar classificadores lineares e regressores em funções de perda convexa, como SVM (linear) e regressão logística



High Learning Rate



Low Learning Rate



Naive Bayes

Baseado no teorema de Bayes e assume a independência entre as features.

S: Spam
H: Ham (not spam)
B: 'Buy'
C: 'Cheap'

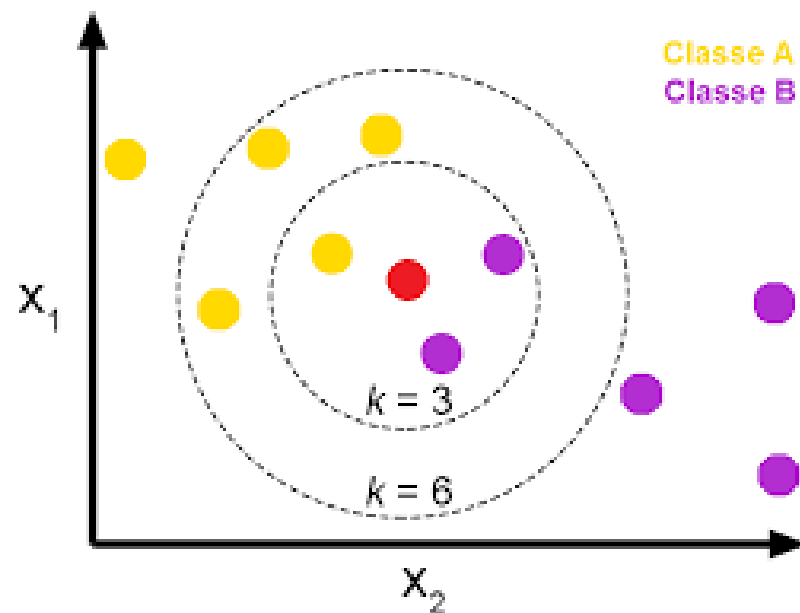
Naive Bayes

$$P(S | B \cap C) = \frac{P(B|S)P(C|S)P(S)}{P(B|S)P(C|S)P(S) + P(B|H)P(C|H)P(H)}$$

$$\begin{aligned} P(\text{spam if "Buy" \& "Cheap"}) &= \frac{\frac{20}{25} \cdot \frac{15}{25} \cdot \frac{25}{100}}{\frac{20}{25} \cdot \frac{15}{25} \cdot \frac{25}{100} + \frac{5}{75} \cdot \frac{10}{75} \cdot \frac{75}{100}} \\ &= 94.737\% \end{aligned}$$

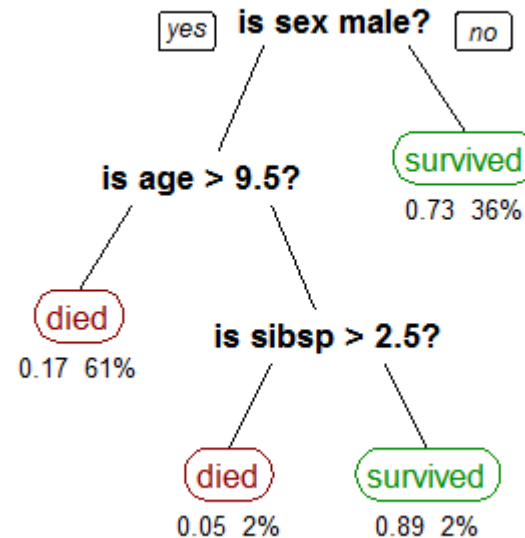
K-Nearest Neighbors

KNN funciona encontrando as distâncias entre uma amostra e os K exemplos mais próximos, em seguida, vota para o rótulo mais frequente (no caso de classificação) ou calcula a média dos rótulos (em o caso de regressão)



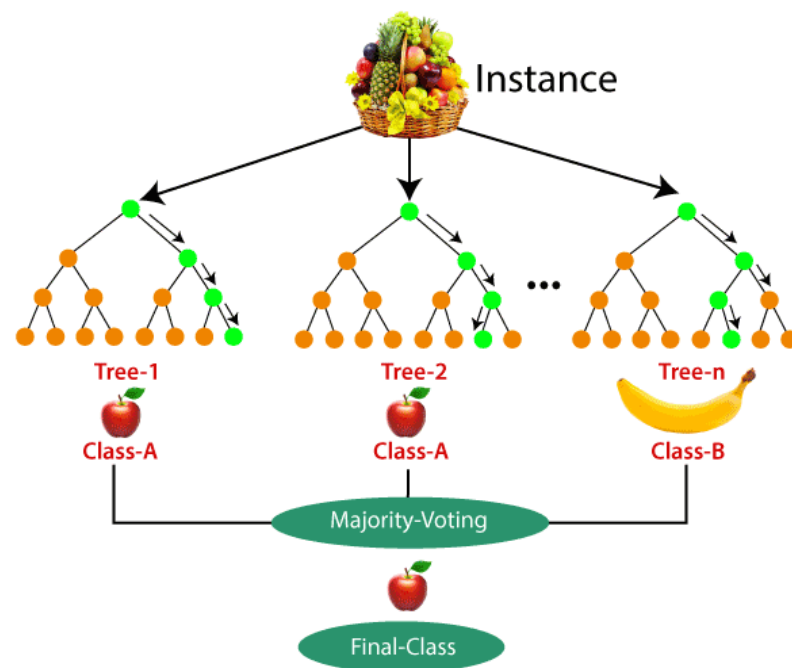
Decision Tree

Todas os features são consideradas e diferentes pontos de divisão são experimentados e testados usando uma função de custo. A divisão com o melhor custo (ou menor custo) é selecionada.



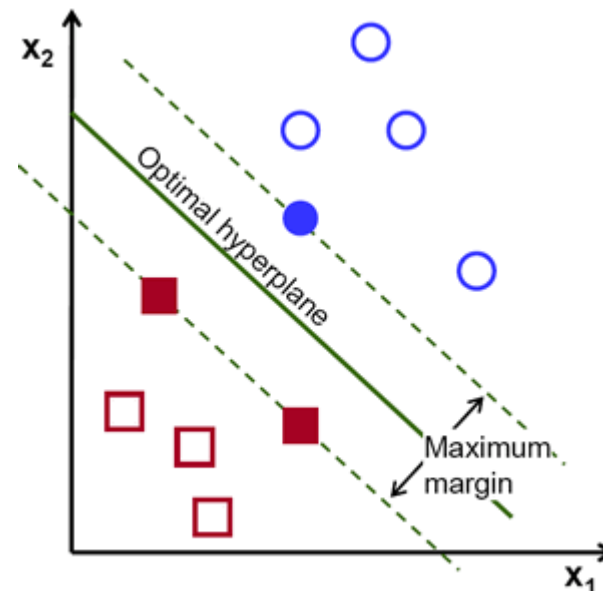
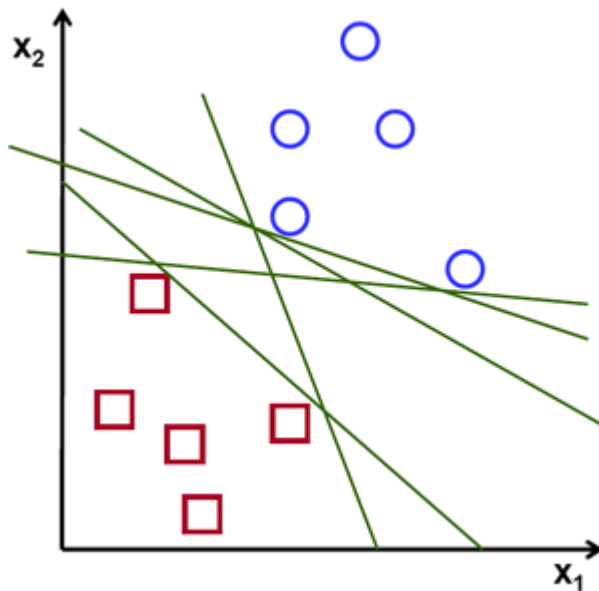
Random Forest

Random Forest é composta por várias Decision Trees. A união é feita por Bagging.



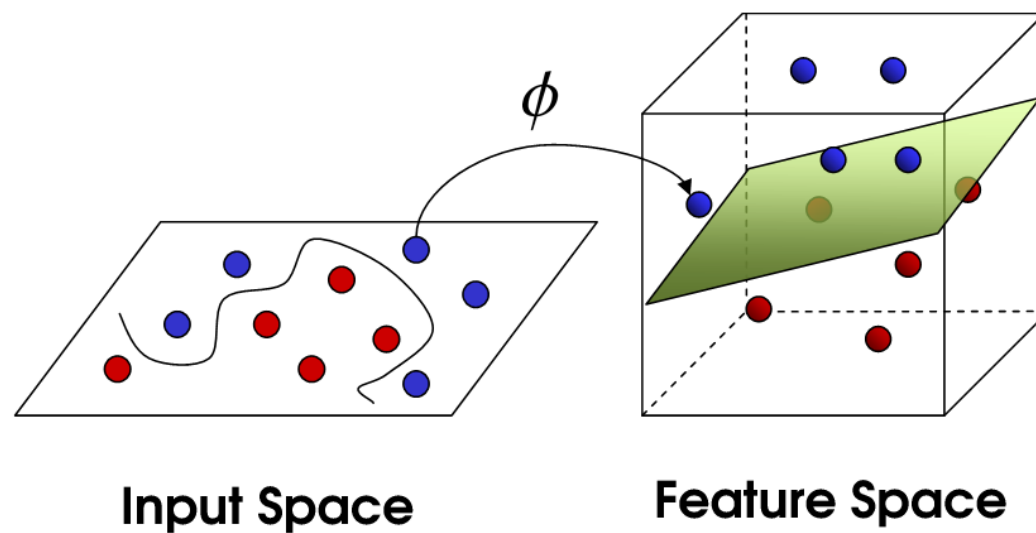
Support Vector Machine

O SVM tenta aumentar a margem de distância entre amostra de diferentes grupos



Support Vector Machine

O truque do kernel é uma transformação que permite a aplicação de um plano em dados que originalmente não poderia ser separados desta forma



Support Vector Machine

O truque do kernel é uma transformação que permite a aplicação de um plano em dados que originalmente não poderia ser separados desta forma

