

Aprendizagem Estatística
Machine learning
Learning from data

Seleção Estatística de modelos

Bibliografia

- Hastie, Tibshirani, Friedman. The elements of Statistical Learning.
- Notas de Rafael.

Objetivo : Queremos encontrar padrões em dados.

O que é um padrão ?

Exemplo 1 : Um dos períodos do sono é REM (rapid eye movement)

Os registros REM parecem os registros feitos durante vigília

Sidarta Ribeiro (UFRRN) tem a seguinte conjectura :

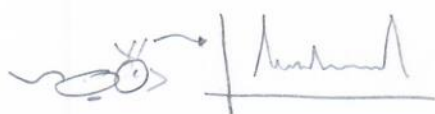
O cérebro "reverbera" as experiências de vigília durante o sono REM e assim aprende, constitui memória.

D. Brillinger, famoso estatístico dos anos 70

Tive a ideia de discretizar uma sinal a sequencias de 0's e 1's



ACORDADO



SONO

?

Exemplo 2 : (Seq 1)

Seq 1 : 0 1 1 0 1 0 0 0 1 0 1 0 0 1

Seq 2 : 0 1 0 0 1 0 0 0 1

tem um padrão comum ?

Ambas sequências têm proporção de números 1 ≈ 1/2 (talvez isso é um padrão)

$\hat{p}_n \xrightarrow{n \rightarrow \infty} 1/2$ Ley dos grandes números

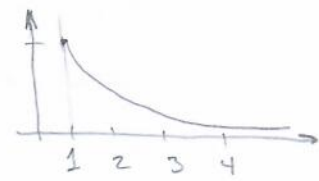
O que é um padrão ???

Exemplo 3: (Seq 2)

1 0000 ... 0 1111 ... 100 ... 01

C: comprimento é aleatório e tem uma distribuição geométrica de parametro p (1/2 < p < 1)

$P(C=k) = p^{k-1} (1-p)$
↑
k = 1, 2, 3, ...

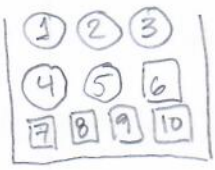


$\sum_{k=1}^{\infty} p^{k-1} (1-p) = \sum_{j=0}^{\infty} p^j (1-p)$
 $= (1-p) \sum_{j=0}^{\infty} p^j$ (projeção geométrica)
 $= (1-p) \frac{1}{1-p}$
 $= 1$

00000 11111 ... 1 000 ... 01 ... 1
C1 C1 C2 C2

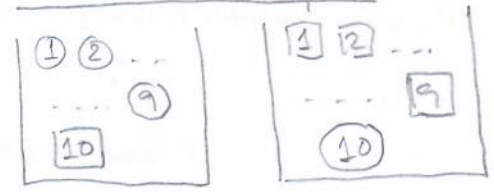
$P(C_i^0 = k) = p^{k-1} (1-p)$
 $P(C_i^1 = k) = p^{k-1} (1-p)$

Gerador Seq 1



- 1. Escolho uma bola:
- 2. Si a bola é brameca (O) coloco 0.
Si a bola é amarela (□) coloco 1.

Gerador Seq 2



- 1. Escolho uma das 2 urnas ao acaso com prob. 1/2, 1/2
- 2. Escolho ao acaso uma bola da urna sorteada. Se a bola for da cor semelhante renomeo na urna

e faço novo sorteio

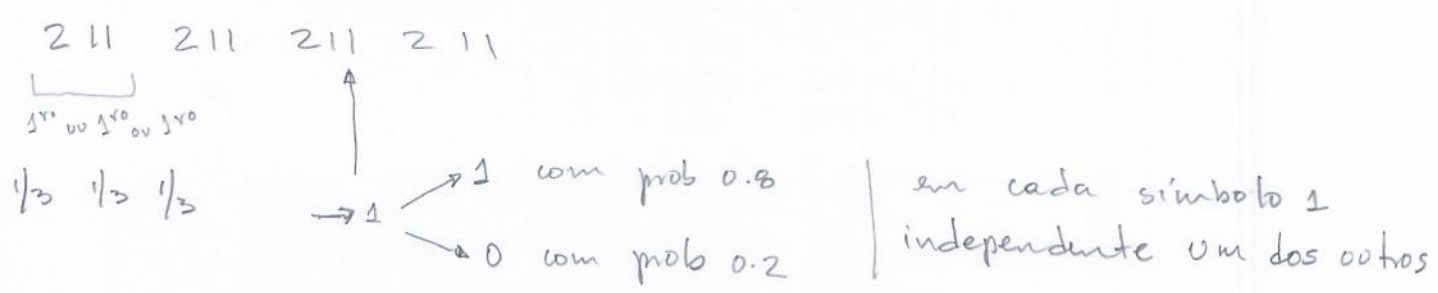
3. Se a bola for da cor minoritária, mudo de urna

4. Depois de cada sorteio reponho a bola na urna.

Exercício (para casa)

Escrever pseudo-códigos para implementar algoritmos gerando as 2 experiências de urna; mais a seguinte experiência

Jogo do Goleiro:



Geração

1. Escolho começar com 2
ou com 1º símbolo 1 com probabilidade 1/3, 1/3, 1/3
ou com 2º símbolo 1
2. A partir do símbolo escolhido concateno
3. Atualizo cada símbolo 1 de maneira iid da seguinte maneira:
1 permanece 1 com prob 0.8
1 é alterado para 0 com prob. 0.2

M. Gromov (geômetra)

Padrão: (Definição provisória)

Conjunto coerente de regularidade estatística.

Resumo: O que fizemos na 1ª parte.

- Discutimos como "classificar" seqüências de 0's e 1's.
- Tentamos fazer isso utilizando proporções relativas.
- No final atribuímos um algoritmo de geração para cada seqüência.

Sequência $x_1, x_2, \dots, x_n = x$, $x_i \in \{0, 1\}$

(4)

↓
modelo - [algoritmo de
geração

Problema de classificação
sem supervisão

Classificação com supervisão:

\mathcal{S} : conjunto

$\xi_1, \xi_2, \dots, \xi_n$ são elementos de \mathcal{S}

f ↙
 y_1, y_2, \dots, y_n , $y_n = f(\xi_n)$, $y_n \in \{0, 1\}$

(ξ_n, y_n) são iid

Hipótese simplificadora: $y_n =$ função determinista de ξ_n

Dado amostra $(\xi_1, y_1), \dots, (\xi_n, y_n)$ iid.

Sei que $y_n =$ função (ξ_n) e sei também que essa função desconhecida pertence ao conjunto \mathcal{H} de funções

$\mathcal{H} =$ conjunto de "hipóteses" ou de "modelos" que tenho à minha disposição

Problema: Escolher $\hat{f}_n \in \mathcal{H}$ que "melhor se ajuste" aos dados

Erro do ajuste (Risco):

Tomo $h \in \mathcal{H}$,

$R(h) =$ probabilidade de $h(\xi) \neq y$

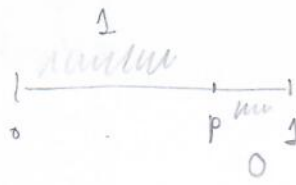
$= \mathbb{P}(h(\xi) \neq y)$

↳ não conheço \mathbb{P} .

Obs:

$f = \arg \min \{ R(h) : h \in \mathcal{H} \}$, $R(f) = 0$ por hipótese

Exemplo:



(5)

$$\mathcal{S} = [0, 1]$$

$\xi_1, \xi_2, \dots, \xi_n$ são sorteios independentes feitos com uma distribuição uniforme em $[0, 1]$.

$$y_n = \mathbb{1}_{\{\xi_n \leq p\}} \quad \text{onde } 0 < p < 1 \text{ fixado.}$$

$$\mathcal{H} = \left\{ h : [0, 1] \rightarrow \{0, 1\} : h \text{ é da forma } h(u) = \mathbb{1}_{\{u \leq q\}} \text{ para algum } q \in [0, 1] \right\}$$

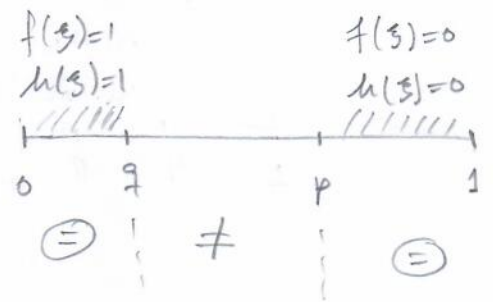
\mathbb{P} = distribuição uniforme em $[0, 1]$

$$R(h) = \mathbb{P}(h(\xi) \neq f(\xi)) \quad \text{onde } h(\xi) = \mathbb{1}_{\{\xi \leq q\}}$$

$$\mathbb{P}([q, p]) = p - q$$

Vamos supor que $q < p$

$$f(\xi) = \begin{cases} 1 & \text{se } \xi \leq p \\ 0 & \text{se } \xi > p \end{cases}$$



Em geral $R(\mathbb{1}_{\{\xi \leq q\}}) = |p - q|$, e $\min h(h) = 0$ ocorre quando $p = q$.

Escolho $h \in \mathcal{H}$, isto é h da forma $\mathbb{1}_{[0, q[}$ para algum $q \in [0, 1]$

$$\hat{R}_n(h) = \frac{1}{n} \sum_{m=1}^n \mathbb{1}_{\{h(\xi_m) \neq y_m\}} \quad \text{onde } (\xi_1, y_1), \dots, (\xi_n, y_n)$$

é a amostra de que dispomos

$$\hat{R}_n(h) \xrightarrow{n \rightarrow \infty} R(h) \quad \left| \text{Ley dos grandes números} \right.$$

Fixo δ ,

$$P(|\hat{R}_n(h) - R(h)| > \delta) \leq \epsilon(n) \downarrow 0$$

gostaria de ter a menor majoração possível.

$$\hat{R}_n(h) - R(h)$$

$$\frac{1}{n} \sum_{m=1}^n \mathbb{1}_{\{h(\xi_m) \neq y_m\}}$$

$$P(h(\xi) \neq y) \text{ onde } y = f(\xi)$$

$z_m \rightarrow$ var. aleatórias iid

$$E(\mathbb{1}_{\{h(\xi) \neq y\}})$$

Obs.:

$$A \subset [0, 1]$$

$$\mathbb{1}_A = \begin{cases} 1 & \text{se } u \in A \\ 0 & \text{se } u \notin A \end{cases}$$

$$P(A) = E(\mathbb{1}_A)$$

$$E(\mathbb{1}_A) = 1 \cdot P(u: u \in A) + 0 \cdot P(u: u \notin A) = P(A)$$

$$\hat{R}_n(h) = \frac{1}{n} \sum_{m=1}^n z_m, \quad R(h) = E(z)$$

$$\hat{R}_n(h) - R(h) = \frac{1}{n} \sum_{m=1}^n z_m - E(z)$$

media empírica (conheço)

media teórica (não conheço)

Lei dos grandes números diz que sob certas hipóteses a media empírica converge para a media teórica.

$$\mathbb{P}(|\hat{R}_n(n) - R(n)| > \delta) = \mathbb{P}\left(\left|\frac{1}{n} \sum_{m=1}^n z_m - \mathbb{E}(z)\right| > \delta\right) \quad (7)$$

$$= \mathbb{P}\left(\left|\sum_{m=1}^n z_m - n\mathbb{E}(z)\right| > n\delta\right)$$

$$\leq \frac{1}{n} \frac{\text{var}(z)}{\delta^2} \quad \text{Desigualdade de Chebyshev}$$

"
 $\mathbb{E}(n)$

esto pode ser majorado por $1/4$

Na verdade podemos obter:

$$\mathbb{E}(n) = e^{-2n\delta^2}$$

Desigualdade de Hoeffding
 \sim (1960).