

Classificação $\left\{ \begin{array}{l} \text{Sem supervisão} \rightarrow \text{sequências de 0's e 1's geradas} \\ \text{por uma fonte probabilística} \\ \text{com supervisão} \rightarrow \text{último mes Deep Learning.} \end{array} \right.$

O que era Ξ_n passa a ser X_n

Espaço X de "entradas"
(conjunto) "sinais"

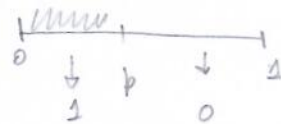
Y conjunto de "saídas"
"etiquetas"
"categorias"

Exemplo 1ª aula:

$$X = [0, 1]$$

$$Y = \{0, 1\} \quad \text{classificação binária}$$

$$X \in [0, 1] \quad Y = f(X) = \mathbb{1}_{\{X \leq p\}}$$



$(X, Y) \in X \times Y$ são aleatórias
Não conhecemos a distribuição conjunta

\mathcal{F} = conjunto de funções de X em Y
com certas restrições

objetivo: encontrar $f^* \in \mathcal{F}$ que minimize o "risco" da classificação

$$f \in \mathcal{F}$$

$$R(f) = \mathbb{P}(f(X) \neq Y) \quad \text{risco}$$

No exemplo da última aula X é aleatório e Y é uma função determinística de X .

$$P(X=x, Y=1) = P(X=x, f(X)=1)$$

$$= P(X=x, f(x)=1) = P(X=x) \mathbb{1}_{\{f(x)=1\}}$$



$$\tilde{Y} = 1 \text{ se } X \leq p$$

$$Y = \mathbb{1}_{\{X \leq p\}} + z$$



z
↑ ruído da linha.

Hipótese: Ruído z é independente da entrada X

$$P(X \leq p, Y=1) = P(X \leq p) P(z=1)$$

sem ruído na
classificação seria

$$= p \cdot P(z=1) = p(1-\epsilon)$$

Vamos supor que $P(z=1) = 1-\epsilon$ onde ϵ é pequeno ($\epsilon < 1/2$)

Escolho \mathcal{F} como sendo a classe de funções da forma:

$$g: [0,1] \rightarrow \{0,1\}$$

$$g(x) = \mathbb{1}_{\{x \leq q\}}$$

← parâmetro

$$P(Y \neq g(x)) = ?$$

Na aula anterior $\epsilon=0$, $P(z=1)=1$

$$P(g(x) \neq Y) = |q-p|$$

$$\inf_{g \in \mathcal{F}} P(g(x) \neq Y) = 0$$

↑ é atingido por $f^*(x) = \mathbb{1}_{\{x \leq p\}}$

se $\epsilon \in (0, 1/2)$ quero calcular

$$P(g(x) \neq Y) \text{ quando } Y = \mathbb{1}_{\{X \leq p\}} \cdot z$$

$$g \in \mathcal{F}$$

3

Definimos:

$$R(g) = \mathbb{P}(g(x) \neq Y)$$

$$R^* = \inf_{g \in \mathcal{F}} R(g) \quad \text{Risco de Bayes}$$

f^* = classificador de Bayes satisfaz $R(f^*) = R^*$

Texto: Bousquet, Bouchieron, Lugosi
Introduction to Statistical Learning theory

Função de regressão:

$$x \in \mathcal{X} \longmapsto \eta(x) \in \mathcal{Y} = \{0,1\}$$

↑
qualquer

$$\eta(x) = \mathbb{1} \left\{ \mathbb{P}(Y=1 | X=x) > \frac{1}{2} \right\}$$

Teorema

$\eta(x)$ = classificador de Bayes

Exemplo:

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

$X \in [0,1]$ e tem distribuição uniforme

isto é $\mathbb{P}(X \leq r) = r$

$Z \in \{0,1\}$ e $\mathbb{P}(Z=1) = 1-\varepsilon$, onde $0 \leq \varepsilon \leq 1$

Vamos supor que X, Z são independentes

Defino $Y = \mathbb{1} \{X \leq p\} \cdot Z$

Exercício: calcular

$$\eta(x) = \mathbb{1} \left\{ \mathbb{P}(Y=1 | X=x) > \frac{1}{2} \right\}$$

Na prática NÃO conhecemos a lei conjunta de $X \in Y$ e portanto não temos como calcular $R^* = \inf_{g \in \mathcal{F}} \mathbb{P}(g(X) \neq Y)$ (9)

nem $\eta(x) = \mathbb{1}_{\{\mathbb{P}(Y=1|X=x) > 1/2\}}$

Porém temos amostras $(x_1, y_1), \dots, (x_n, y_n)$ com n sorteios independentes de (X, Y)

Dado $g \in \mathcal{F}$, calculo $\hat{R}_n(g) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{g(x_i) \neq y_i\}}$

Vimos na última aula no caso $\epsilon=0$

$$\hat{p}_n = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{y_i=1\}} \xrightarrow{n \rightarrow \infty} p$$

parâmetro correto

Aplicação da
 Ley dos grandes
 números.

Então se chamamos $\hat{f}_n(x) = \mathbb{1}_{\{x \leq \hat{p}_n\}}$ então

$$\hat{f}_n \xrightarrow{n \rightarrow \infty} f(x) = \mathbb{1}_{\{x \leq p\}}$$

$$\mathbb{P}\left(\left|\hat{R}_n(g) - R(g)\right| > \epsilon\right) = \mathbb{P}\left(\left|\frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{g(x_i) \neq y_i\}} - \underbrace{\mathbb{P}(g(X) \neq Y)}_{\mathbb{E}(\mathbb{1}_{\{g(X) \neq Y\}})}\right| > \epsilon\right)$$

$$= \mathbb{P}\left(\left|\frac{1}{n} \sum_{i=1}^n z_i - \mathbb{E}(z)\right| > \epsilon\right) \leq 2e^{-2n\epsilon^2}$$

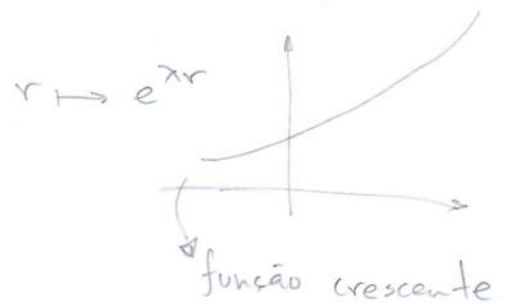
\uparrow v.a iid $z_i \in \{0,1\}$ \uparrow Hoeffding's.

$$\mathbb{P}\left(\left|\frac{1}{n} \sum_{i=1}^n z_i - \mathbb{E}(z)\right| > \epsilon\right) = \mathbb{P}\left(\left|\sum_{i=1}^n z_i - n\mathbb{E}(z)\right| > n\epsilon\right)$$

$$= \mathbb{P}\left(\left\{\sum_{i=1}^n z_i > n[\mathbb{E}(z) + \epsilon]\right\} \cup \left\{\sum_{i=1}^n z_i < n[\mathbb{E}(z) - \epsilon]\right\}\right)$$

$$= \mathbb{P} \left(\sum_{i=1}^n z_i > n [\mathbb{E}(z) + \varepsilon] \right) + \mathbb{P} \left(\sum_{i=1}^n z_i < n [\mathbb{E}(z) - \varepsilon] \right)$$

$$= \mathbb{P} \left(e^{\lambda \sum_{i=1}^n z_i} > e^{\lambda n [\mathbb{E}(z) + \varepsilon]} \right)$$



→ método de Cramer-Chernoff ~ 1940

Desigualdade de Markov ~ 1905

v.a W positiva

$$\mathbb{P}(W \geq a) \leq \frac{\mathbb{E}(W)}{a}$$

$$W = e^{\lambda \sum_{i=1}^n z_i}, \quad a = e^{\lambda n (\mathbb{E}(z) + \varepsilon)}$$

→ usando Markov segue que:

$$\mathbb{P} \left(e^{\lambda \sum_{i=1}^n z_i} > e^{\lambda n (\mathbb{E}(z) + \varepsilon)} \right) \leq \frac{\mathbb{E} \left(e^{\lambda \sum_{i=1}^n z_i} \right)}{e^{\lambda n (\mathbb{E}(z) + \varepsilon)}}$$

$$= \frac{\prod_{i=1}^n \mathbb{E} \left(e^{\lambda z_i} \right)}{e^{\lambda n (\mathbb{E}(z) + \varepsilon)}} = \frac{\mathbb{E} \left(e^{\lambda z} \right)^n}{e^{\lambda n (\mathbb{E}(z) + \varepsilon)}}$$

isso vale para todo $\lambda > 0$

Final da conta: achar "o melhor" λ para essa majoração

Hoeffding diz que:

$$\mathbb{P} \left(|\hat{R}_n(g) - R(g)| > \varepsilon \right) \leq \underbrace{2e^{-2n\varepsilon^2}}_{\delta}$$

$$\text{Eu quero } \delta = 2e^{-2n\varepsilon^2}$$

$$\frac{\delta}{2} = e^{-2n\varepsilon^2}$$

$$\ln \frac{\delta}{2} = -2n\varepsilon^2, \quad -\frac{1}{2n} \ln \frac{\delta}{2} = \varepsilon^2, \quad \sqrt{-\frac{1}{2n} \ln \frac{\delta}{2}} = \varepsilon$$

Se queremos que a probabilidade de erro seja δ , tomamos $\epsilon = \sqrt{\frac{\ln 2/\delta}{2n}}$

Isto acontece

$R(\beta) \leq \hat{R}_n(\beta) + \sqrt{\frac{\ln 2/\delta}{2n}}$ com probabilidade $\geq 1-\delta$

Jogo do Goleiro:

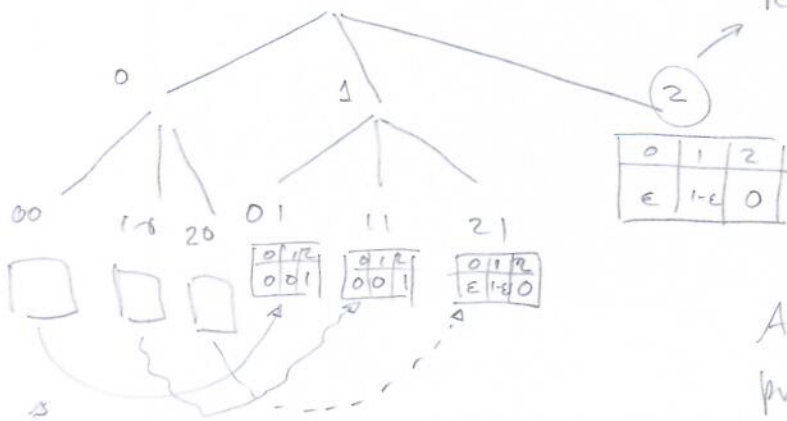
Batedor: Começa com a sequência determinística

2 1 1 2 1 1 2 1 1 ...

cada símbolo 1 pode ser ou transformado em 0 com probabilidade ϵ ou mantido com prob $1-\epsilon$.

Cada uma das escolhas é feita independentemente das anteriores.

$X_n = ?$



Representa o conjunto das sequências que acabam com o símbolo 2.

Temos

Árvore, uma família de probabilidades de transição indexada pelas folhas da árvore.

Algoritmo:

- 1. Como começar?
- 2. Como escolher próximo passo.

Árvore τ define uma partição no conjunto de todas as sequências de símbolos passados.

O que é Árvore: 1ª resposta: Grafo sem ciclos com etiqueta. Grafo orientado (por laço de hereditário 1, 1)

usamos

Alfabeto $A = \{0, 1, 2\}$

Vamos representar τ por suas folhas

$$\tau = \{2, 01, 11, 21, 00, 10, 20\}$$

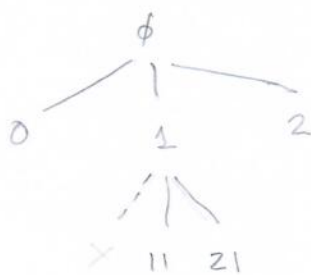
Outro exemplo, $\tau' = \{2, 12, 0, 1\}$
incompatível



Nenhum elemento de τ pode ser sufixo próprio de outro.

Candidato a árvore

$$\tilde{\tau} = \{0, 11, 21, 2\}$$



Problema: E se a sequência terminar com 01?

Ela pode ser usada só para algoritmos que nunca geram pares 01.

Exemplo (exercício)

1. Tomo a sequência periódica

$$2101 \ 2101 \ 2101 \ \dots$$

2. Transformo símbolos 1 em zero com prob ϵ e mantenho 1 com prob. $1-\epsilon$. Faço isso de maneira iid

τ árvore que define uma partição de todos os passados que podem aparecer no arquivo.

Represento τ por as folhas, $\text{CONTEXT} = \text{FOLHA}$

$$(\alpha_{-m}, \dots, \alpha_{-2}, \alpha_{-1}) = \alpha_{-m}^{-1}$$

$$\dots \alpha_{-m}, \alpha_{-(m-1)}, \dots, \alpha_{-1} = \alpha_{-\infty}^{-1}$$

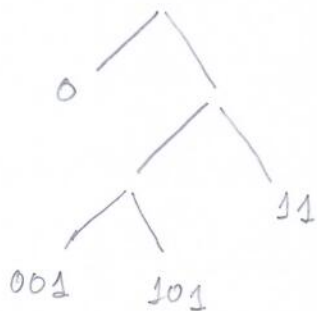
$$c_{\tau} : \alpha_{-\infty}^{-1} \mapsto c_{\tau}(\alpha_{-\infty}^{-1}) \text{ que é o único sufixo de } \alpha_{-\infty}^{-1} \text{ que}$$

pertence a Σ .

Goleiro

Escolhe o próximo símbolo usando a probabilidade de transição associada ao contexto que termina no símbolo anterior.

Como começar: Começo com um contexto de comprimento máximo.



$$l(x_{-k}^{-1}) = k$$

↑ comprimento da sequência

Altura da árvore

$$\max \{ l(w) : w \in \Sigma \} = h$$

↑
comprimento máximo

• Associamos uma árvore a uma sequência
• x_0, x_1, \dots, x_n