

\mathcal{M} : classe de modelos

A : algoritmo: amostra \rightarrow modelo

Perguntas:

Ha' modelos idênticos em $\hat{t}_n^1, \dots, \hat{t}_n^M$?
"próximos"

\rightarrow Isto depende de uma distância definida em \mathcal{M}

Cadeias estocásticas com memória de alcance variável
1983 - J. Rissanen
"A universal system for data compression"
Minimal Description Length

$k \geq 1$

$$\mathcal{M}_k(A) = \left\{ p: A^k \times A \rightarrow [0,1] : \forall a_{-k} \in A^k, \text{ vale } \sum_{b \in A} p(b|a_{-k}) = 1 \right\}$$

$p \in \mathcal{M}_k(A)$ e $(X_n)_n$ foi gerada por p
↑
cadeia de Markov de alcance k .

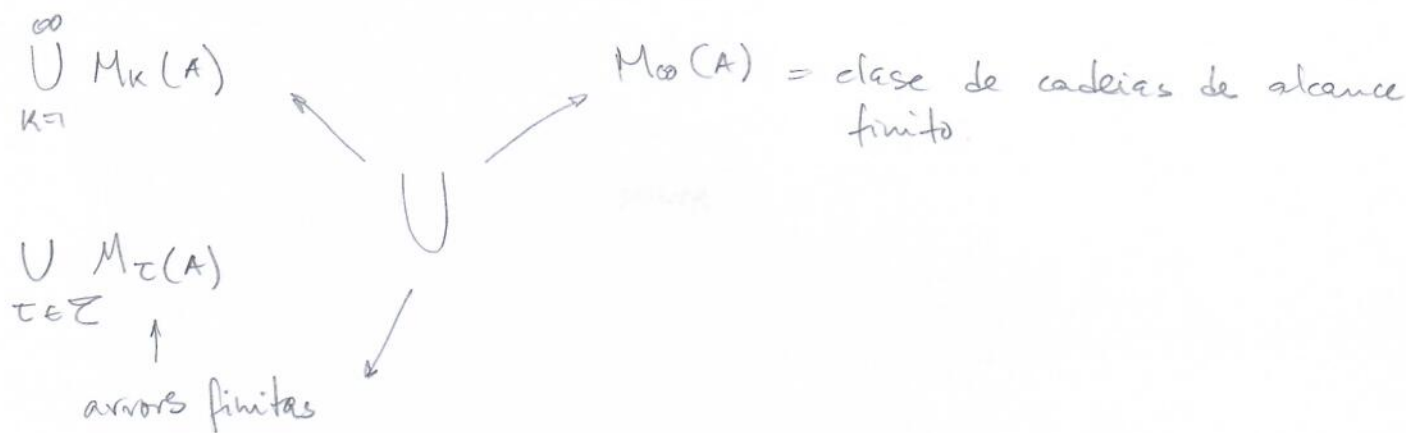
τ árvore de contextos de altura k .

$$\mathcal{M}_\tau(A) = \left\{ p: \tau \times A \rightarrow [0,1], \forall w \in \tau, \sum_{b \in A} p(b|w) = 1 \right\}$$

Se altura de $\tau = k$

$$M_\tau(A) \subset M_k(A)$$

$\tau \rightarrow$ "captura" informações sobre a dependência temporais da cadeia



Seja $p \in M_k(A)$, gero amostra x_1, x_2, \dots, x_n usando p

- Perguntas:
- 1) Conhecendo k , estimar p
 - 2) Se k e p for desconhecidos, seleccionar \hat{k}_n y \hat{p}_n que "melhor" descrebam a amostra.

Seja $p \in M_k(A)$ com τ finito, gero a amostra x_1, \dots, x_n usando (τ, p) .

- 3) Se τ e p foram desconhecidos, seleccionar $\hat{\tau}_n$ e \hat{p}_n que "melhor" descrevam a amostra.

Perguntas:

- 1) Dada a amostra x_1, \dots, x_n qual é o ^{máximo} alcance k que podemos identificar?

Resposta: Só conseguimos identificar alcances k tais que $|A|^k \ll n$

$$e^{k \ln |A|} < e^{\ln n}$$

$$\Rightarrow k < \frac{\ln n}{\ln |A|}$$

T. Shannon Mc. Millan Breese
Lema de Kac

2. Como seleccionar k (ou τ) com alcance $< \frac{\ln n}{\ln |A|}$?

(3)

Estimador de máxima verosimilhança em $M_k(A)$ ou $M_\tau(A)$

amostra: $X_{-k}^n = X_{-k}, X_{-k+1}, \dots, X_n$

$$k = \frac{\ln n}{\ln |A|} \quad \text{No:}n(a^{-k}b) = \sum_{t=0}^n \mathbb{1}_{\{X_{t-k}^{t-1} = a^{-k}, X_t = b\}}$$

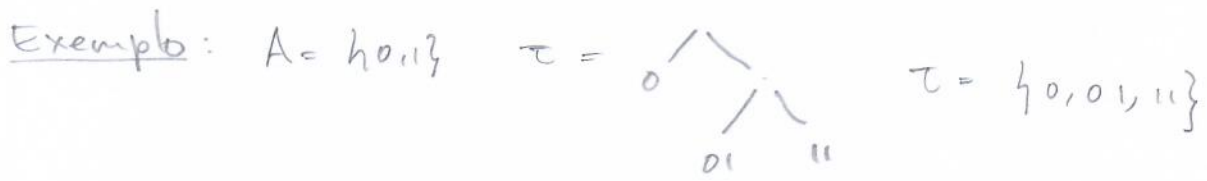
$$\hat{P}_n^{[k]}(b|a^{-k}) = \frac{\text{No:}n(a^{-k}b)}{\sum_{z \in A} \text{No:}n(a^{-k}z)} = \frac{\text{No:}n(a^{-k}b)}{\text{No:}n-1(a^{-k})}$$

Para caso: τ árvore finita

$$\hat{P}_n^\tau(b|w) = \frac{\text{No:}n(wb)}{\text{No:}n-1(w)}, \quad \forall w \in \tau$$

onde: $\text{No:}n(wb) = \sum_{t=0}^n \mathbb{1}_{\{X_{t-\ell(w)}^{t-1} = w, X_t = b\}}$, onde

$\ell(w)$ = comprimento de w .



$$\text{No:}n((01)0) = \sum_{t=0}^n \mathbb{1}_{\{X_{t-2}^{t-1} = (0,1), X_t = 0\}}$$

Vimos que se a amostra foi gerada por $p \in M_k(A)$ (ou por $M_\tau(A)$)

então $\hat{P}_n^{[k]}(b|a^{-k}) \xrightarrow{n \rightarrow \infty} p(b|a^{-k})$

$(\hat{P}_n^\tau(b|w) \xrightarrow{n \rightarrow \infty} p(b|w))$

Leia los grandes numeros.

Se $p \in M_{1k}(A)$ e (X_n) foi gerada por p

$$P(X_n = b \mid X_{n-k} = a^{-k}, \underbrace{X_{n-(k+1)} = z}_{\text{informação inútil}})$$

||

$$p(a \mid a^{-k}) \text{ qualquer seja } z$$

$$\hat{p}_n^{[k+1]}(b \mid a^{-k} z) = \frac{N_{0:n}(z a^{-k} \dots a^{-1} b)}{N_{0:n-1}(z a^{-k} \dots a^{-1})} \rightarrow p(b \mid a^{-k})$$

$\underbrace{\hspace{10em}}_{\text{últimos } k \text{ valores}}$ \uparrow valor $k+1$ passos atrás \uparrow isto converge a porque não depende de z

Então $\hat{p}_n^{[k+1]}(b \mid a^{-k} z) \stackrel{!?!?!}{\approx} \hat{p}_n^{[k]}(b \mid a^{-k})$ se n for grande

!?!?! o que quer dizer estatisticamente igual

Dado $\epsilon > 0$, ϵ pequeno

$$P\left(\left|\hat{p}_n^{[k]}(b \mid a^{-k}) - p(b \mid a^{-k})\right| > \epsilon\right) \leq \delta(\epsilon, n)$$

como calcular $\delta(\epsilon, n)$?

\uparrow
pequeno. sob a hipótese nula:
A cadeia foi gerada por p .

Teorema Limite Central

$$\sqrt{n} \left[\frac{1}{n} \sum_{i=1}^n \xi_i - p \right] \xrightarrow{\text{Distribuição}} N(0, p(1-p))$$

Ley dos grandes números com que velocidade?

ξ_1, ξ_2, \dots iid $\xi_i \in \{0,1\}$

$$\frac{1}{n} \sum_{i=1}^n \xi_i \xrightarrow{n \rightarrow \infty} P(\xi_i = 1) = p$$

\swarrow com que velocidade?

$$\mathbb{P} \left(\hat{p}(b|w) - p(b|w) > \epsilon \right)$$

$$= \mathbb{P} \left(\frac{N_{0:n}(wb)}{N_{0:n-1}(w)} - p(b|w) > \epsilon \right)$$

$$\left[\begin{array}{l} N_{0:n}(wb) = N(wb) \\ p(b|w) = p \end{array} \right]$$

$$= \mathbb{P} \left(N(wb) - pN(w) > N(w)\epsilon \right)$$

$$M_n = N_{0:n}(wb) - N_{0:n-1}(w) p(b|w)$$

↑ Isto é um MARTINGALE!! Doob (anos 50)

$$\mathbb{E} \left(M_n \mid X_{-k}^{n-1} \right) = M_{n-1}$$

↑
esperança condicional

$$\rightarrow \mathbb{P} \left(M_n > N(w)\epsilon \right)$$

se fosse um número e não uma v.a poderíamos usar a Des. de Hoeffding para martingais.

$$\mathbb{P} \left(M_n > N(w)\epsilon ; N(w) < m \right) \leq \mathbb{P} \left(N(w) < m \right)$$

caso pessimista

$$+ \mathbb{P} \left(M_n > N(w)\epsilon ; N(w) \geq m \right) \leq \mathbb{P} \left(M_n > m\epsilon ; N(w) \geq m \right)$$

ótimo, w aparece mais de m vezes

Juntando temos: evento altamente provável

$$\mathbb{P} \left(\hat{p}(b|w) - p(b|w) > \epsilon \right) \leq \mathbb{P} \left(N(w) < m \right) + \mathbb{P} \left(M_n > m\epsilon \right)$$

Aqui posso usar Hoeffding

$(X_1, Y_1), \dots, (X_n, Y_n)$ iid

$X_n \in \mathcal{X}$

$Y_n \in \{0, 1\}$

$\mathcal{G}: \mathcal{X} \rightarrow \{0, 1\}$

↑
Classe de funções ("modelos")
candidatas a classificador

Objetivo: Encontrar o classificador que a cada amostra X associa Y a partir de amostra $(X_1, Y_1), \dots, (X_n, Y_n)$

Risco do classificador $g \in \mathcal{G}$

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

Risco mínimo: $R^* = \inf \{ R(f) : f: \mathcal{X} \rightarrow \{0, 1\} \}$

f mensurável | $\{x: f(x)=1\}$ é um evento conjunto mensurável

Classificador de Bayes:

$$f^*(x) = \mathbb{1}_{\{\eta(x) \geq 1/2\}} \quad \text{onde} \quad \eta(x) = \mathbb{P}(Y=1 | X=x)$$

↑
função de regressão

Teorema:

$$R(f^*) = R^*$$

Demonstração ver página

Risco empírico:

→ podemos calcular a partir da amostra

$$R_n(g) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{g(X_i) \neq Y_i\}}$$

↑
variável aleatória

Pergunta central: Quão longe $R_n(g)$ está de $R(g)$?

$R(g)$ se não conhecemos \mathbb{P} não podemos calcular,

$$P(R_n(g) - R(g) > \varepsilon) \leq \underbrace{e^{-2n\varepsilon^2}}_{\delta} \quad (7)$$

\Downarrow

$$P(|R_n(g) - R(g)| > \varepsilon) \leq 2e^{-2n\varepsilon^2}$$

$$\delta = e^{-2n\varepsilon^2} \Rightarrow \ln \delta = -2n\varepsilon^2 \Rightarrow \varepsilon^2 = -\frac{\ln \delta}{2n} \Rightarrow \varepsilon = \sqrt{\frac{\ln 1/\delta}{2n}}$$

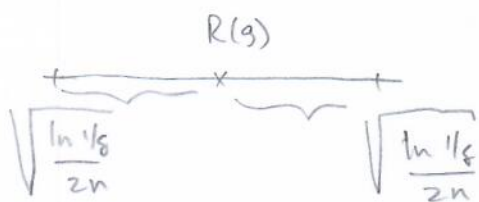
$$P\left(R_n(g) - R(g) > \sqrt{\frac{\ln 1/\delta}{2n}}\right) \leq \delta$$

$1/\delta \nearrow +\infty$ quando $\delta \searrow 0$ precisão de estimação \Rightarrow cresce quando $n \uparrow \infty$
 decresce quando $\delta \searrow 0$.

Idem

$$P\left(|R_n(g) - R(g)| > \sqrt{\frac{\ln 1/\delta}{2n}}\right) \leq \delta$$

ótimo, porque essa majoração não depende da função g
 (só usamos os fatos: g assume os valores 0 e 1 e $(x_1, y_1), \dots, (x_n, y_n)$ iid)



$$R(g) \leq R_n(g) + \sqrt{\frac{\ln 1/\delta}{2n}} \quad \text{com probabilidade } \geq 1-\delta$$

$C_g =$ conjunto ruim

$$C_g = \left\{ |R_n(g) - R(g)| > \sqrt{\frac{\ln 1/\delta}{2n}} \right\}$$

$$P(C_g) \leq \delta$$

Problema: se usamos uma outra função \tilde{g} , teremos um outro conjunto ruim $C_{\tilde{g}}$

Vamos supor que $\mathcal{G} = \{g_1, g_2\}$. Quero encontrar uma stima
tiva que seja boa para ambas (8)

$$\mathbb{P} \left(\exists g \in \mathcal{G} : |R_n(g) - R(g)| > \sqrt{\frac{\ln 1/\delta}{2n}} \right)$$

$$= \mathbb{P} \left(\bigcup_{g \in \mathcal{G}} C_g \right) \leq \sum_{g \in \mathcal{G}} \mathbb{P}(C_g) = \mathbb{P}(C_{g_1}) + \mathbb{P}(C_{g_2}) = 2\delta$$

Se $\mathcal{G} = \{g_1, \dots, g_N\}$

$$\mathbb{P} \left(\bigcup_{g \in \mathcal{G}} \left\{ |R_n(g) - R(g)| > \sqrt{\frac{\ln 1/\delta}{2n}} \right\} \right) \leq \sum_{i=1}^N \mathbb{P}(C_{g_i}) = N \cdot \delta$$

Vamos refazer desde o comenzo:

$$\mathbb{P}(R_n(f) - R(f) > \epsilon) \leq e^{-2n\epsilon^2}$$

$$\{R_n(f) - R(f) > \epsilon\} = D_f$$

$$\mathbb{P}(\exists f \in \mathcal{G} : R_n(f) - R(f) > \epsilon) \leq \sum_{f \in \mathcal{G}} \mathbb{P}(R_n(f) - R(f) > \epsilon)$$

Se está usando

$$\mathbb{P} \left(\bigcup_{i=1}^N D_i \right) \leq \sum_{i=1}^N \mathbb{P}(D_i)$$

$$= |\mathcal{G}| e^{-2n\epsilon^2} = N e^{-2n\epsilon^2}$$

$$\text{Se } \mathcal{G} = \{g_1, \dots, g_N\}$$

↑ tamanho da amostra
↑ # de funções na base

$$\delta = N \cdot e^{-2n\epsilon^2}$$

$$\Rightarrow \epsilon = \sqrt{\frac{\ln N + \ln 1/\delta}{2n}}$$

$$\mathbb{P} \left(\bigcup_{f \in \mathcal{G}} R_n(f) - R(f) > \sqrt{\frac{\ln N + \ln 1/\delta}{2n}} \right) \leq \delta$$

Ou seja

9

$$\mathbb{P} \left(R_n(g) < R(f) + \sqrt{\frac{\ln N + \ln 1/\delta}{2n}}, \forall f \in \mathcal{G} \right) \geq 1 - \delta$$

O numerador $\nearrow +\infty$, para compensar o tamanho da amostra $n \gg \ln N + \ln 1/\delta$

Quero que:

$$\sup_{g \in \mathcal{G}} |R_n(g) - R(g)| \leq \varepsilon \quad \text{e tipicamente } |\mathcal{G}| \text{ e' infinito!!!}$$

Como enfrentar isso? Teoria de VC.

Demostração que classificador de Bayes realiza o risco de Bayes

$$\eta(x) = \mathbb{P}(Y=1 | X=x)$$

$$f^*(x) = \mathbb{1}_{\{\eta(x) \geq 1/2\}}$$

Teorema: $R(f^*) = R^*$

onde $R^* = \inf \{ R(g) : g : \mathcal{X} \rightarrow \{0,1\} \}$ g mensurável

Temos que demonstrar que:

$$R(g) - R(f^*) \geq 0 \quad \text{para toda função } g \text{ mensurável.}$$

Ora:

$$R(g) - R(f^*) = \mathbb{P}(g(X) \neq Y) - \mathbb{P}(f^*(X) \neq Y)$$

$$\mathbb{P}(g(X) \neq Y) = \int_{\mathcal{X}} \mathbb{P}_x(dx) \mathbb{P}(g(x) \neq Y | X=x)$$

$$\mathbb{P}(f^*(X) \neq Y) = \int_{\mathcal{X}} \mathbb{P}_x(dx) \mathbb{P}(f^*(x) \neq Y | X=x)$$

ou seja

$$R(g) - R(f^*) = \int_{\mathcal{X}} \mathbb{P}_x(dx) \left[\mathbb{P}(g(X) \neq Y | X=x) - \mathbb{P}(f^*(X) \neq Y | X=x) \right]$$

Basta demonstrarmos que esta diferença é sempre ≥ 0 .

(10)

$$\begin{aligned} \mathbb{P}(g(X) \neq Y | X=x) &= 1 - \mathbb{P}(g(X) = Y | X=x) \\ &= 1 - \left[\mathbb{P}(g(x)=1, Y=1 | X=x) + \mathbb{P}(g(x)=0, Y=0 | X=x) \right] \\ &= 1 - \left[\mathbb{1}_{\{g(x)=1\}} \underbrace{\mathbb{P}(Y=1 | X=x)}_{\eta(x)} + \mathbb{1}_{\{g(x)=0\}} \underbrace{\mathbb{P}(Y=0 | X=x)}_{1-\eta(x)} \right] \\ &= 1 - \left[\mathbb{1}_{\{g(x)=1\}} \eta(x) + \mathbb{1}_{\{g(x)=0\}} (1-\eta(x)) \right] \end{aligned}$$

Idem:

$$\mathbb{P}(f^*(X) \neq Y | X=x) = 1 - \left[\mathbb{1}_{\{f^*(x)=1\}} \eta(x) + \mathbb{1}_{\{f^*(x)=0\}} (1-\eta(x)) \right]$$

$$\begin{aligned} &\mathbb{P}(g(X) \neq Y | X=x) - \mathbb{P}(f^*(X) \neq Y | X=x) \\ &= (2\eta(x) - 1) \left[\mathbb{1}_{\{f^*(x)=1\}} - \mathbb{1}_{\{g(x)=1\}} \right] \end{aligned}$$

Analisar casos: $\eta(x) > 1/2$ e $\eta(x) < 1/2$, para ambos casos da ≥ 0 .