

Rappel:

Th. Si $(X_n)_{n \geq 1}$ source sans mémoire (iid) de la p d'alphabet fini \mathcal{X} ,

alors \exists un code préfixe $C: \mathcal{X}^* \rightarrow \{0,1\}^*$ et de longueur $L: \{0,1\}^* \rightarrow \mathbb{N}$ tq.

$$H(X) \leq \mathbb{E}[L(C(X))] < H(X) + 1$$

Exercice Wrong code.

Supposons que l'on connaisse pas la vraie loi p; mais seulement q.
 Supposons que l'on choisisse un code C_q dont la longueur

vérifie $L(C_q(x)) = \left\lceil \log \frac{1}{q(x)} \right\rceil$. Notons $H(p) = H(X)$ (calculée avec vraie loi de X)

Mq

$$H(p) + D(p||q) \leq \mathbb{E}[L(C_q(X))] < H(p) + D(p||q) + 1$$

Preuve: $\mathbb{E}[L(C_q(X))] = \sum_{x \in \mathcal{X}} p(x) \left\lceil \log \frac{1}{q(x)} \right\rceil < \sum_{x \in \mathcal{X}} p(x) \left(\log \frac{p(x)}{q(x)} + \frac{1}{p(x)} + 1 \right)$
 $= D(p||q) + H(p) + 1$

de m pour la preuve de la borne inf

Donc ~~intéressant~~ croire que la vraie loi est q(x) alors qu'en fait elle est p(x)
 génère 1 accroissement de $D(p||q)$ de la longueur moyenne du code.

Question: quel accroissement de la longueur moyenne du code si l'on ne connaît pas la loi?

Objectif: si p inconnue, trouver un ~~meilleur~~ code qui rende la longueur de description du message la + petite possible (v. loi processus dans 1 famille)

Notation

X processus iid d'alphabet $\mathcal{X} = \{a_1, a_2, \dots, a_{|\mathcal{X}|}\}$ ^{indice de alphabet}

p loi de X

x^n 1 réalisation de X de longueur $n =$ message à coder
 $= (x_1, \dots, x_i, \dots, x_n)$ ^{indice du vecteur}

q_{x^n} ou $q =$ type de x^n i.e. loi empirique de x^n
 \leftarrow si pas d'ambiguïté, on supprime x^n pour ne pas surcharger les notations.

$\mathcal{Z}_n =$ ens des types de vecteurs de longueur n .

$\tilde{p} =$ une loi arbitraire.

$T(\tilde{p}) =$ classe de type $\tilde{p} =$ ens des vecteurs de longueur n et de type \tilde{p}
 $= \{x^n : q_{x^n} = \tilde{p}\}.$

$$H(\tilde{p}) = - \sum_{a \in \mathcal{X}} \tilde{p}(a) \log_2 \tilde{p}(a).$$

Résumé

$$|\mathcal{Z}_n| \leq (n+1)^{|\mathcal{X}|}$$

Méthode des types

Stratégie 2 étapes

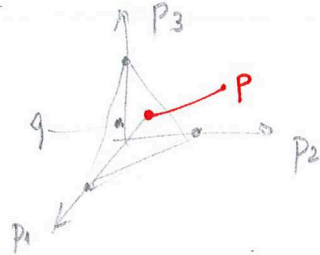
- 1) estimer la loi à partir de x^n ; et transmettre cette estimation
- 2) coder x^n selon cette estimation

Déf type du vecteur x^n est sa loi empirique

$$\forall j \in \llbracket 1, |\mathcal{X}| \rrbracket, q_{x^n}(a_j) = \frac{|\{i: x_i = a_j\}|}{n} = \frac{n_j}{n} = q(a_j)$$

Rappel le simplexe est $\{(p_1, \dots, p_j, \dots, p_{|\mathcal{X}|}) : \forall j \in \llbracket 1, |\mathcal{X}| \rrbracket, p_j \geq 0 \text{ et } \sum_{j=1}^{|\mathcal{X}|} p_j = 1\}$

Ex $|\mathcal{X}| = 3$



Soit $\mathcal{L}_m =$ ens des types de vecteurs de taille n

(CT Th 11.1.1) Montrer que

a) $|\mathcal{L}_m| \leq (n+1)^{|\mathcal{X}|}$

b) $|\mathcal{L}_m| \leq (n+1)^{|\mathcal{X}|-1}$

c) $|\mathcal{L}_m| = \binom{n+|\mathcal{X}|-1}{|\mathcal{X}|-1} \approx n^{|\mathcal{X}|-1}$

indice : ^{c'est le pb.} $|\mathcal{X}|$ candidats, n votants; $|\mathcal{L}_m| =$ nbr de résultats de votes différents

Pour le résultat: considérer le pb.

$|\mathcal{X}|$ boules de couleurs $\neq (a_1, \dots, a_{|\mathcal{X}|})$

n tirages avec remise (au réservoir de boules ∞)

\neq résultats de tirages \neq

Pour la suite a) est suffisant.

a) $m q_{|X|} = (m_1, \dots, m_j, \dots, m_{|X|})$ est 1 vecteur de taille $|X|$

avec $0 \leq m_j \leq m$, i.e. chaque élément du vecteur peut prendre $m+1$ valeurs.
(c'est seul + borne sup, car on n'a pas pris en compte les contraintes) \square

b) On sait que $\sum_{j=1}^{|X|} m_j = m$; donc $m_{|X|}$ est déterminé par les autres valeurs $m_j, 1 \leq j < |X|$

donc il y a au plus $(m+1)^{|X|-1}$ types possibles \square

c) $m_1 = \#$ fois où boule de couleur 1 est sortie

$m_2 =$

\vdots

$m_{|X|} =$

$n = \sum m_j =$ nbr de tirages

Il y a correspondance entre:

• un choix de $m_1, m_2, \dots, m_{|X|}$, $m_j \geq 0$ et $\sum m_j = n$

• une manière de placer $|X|-1$ boules séparatrices parmi $n+|X|-1$



nbr de choix possibles \equiv nbr de manières de placer les boules séparatrices \dots \square

Par formule de Stirling: $n! \approx \sqrt{2\pi n} \left(\frac{n}{e}\right)^n$ $n \rightarrow +\infty$

$$\log \binom{n+|X|-1}{|X|-1} \approx (|X|-1) \log(n+|X|-1) + \text{cte}$$

$$\approx (|X|-1) \log n \quad n \rightarrow +\infty$$

11.1.2 (CT)

Th. Soit x^n processus iid de loi $p(x)$. La probabilité de x^n dépend uniquement de son type et satisfait

$$P(x^n) = 2^{n(-D(q_n \| p) - H(q_n))}$$

Corollaire Si le type de x^n est p , alors

$$P(X^n = x^n) = 2^{-nH(p)}$$

Preuve

$$\begin{aligned} P(X^n = x^n) &= \prod_{i=1}^n p(x_i) = \prod_{a \in \mathcal{X}} p(a)^{\# \text{occurrences de } a \text{ dans } x^n} \\ &= \prod_{a \in \mathcal{X}} 2^{nq(a) \log p(a)} = \prod_{a \in \mathcal{X}} 2^{nq(a) \log p(a) + nq(a) \log q(a) - nq(a) \log q(a)} \\ &= 2^{n \sum_{a \in \mathcal{X}} (-q(a) \log q(a) + q(a) \log p(a) + q(a) \log q(a))} \\ &= 2^{n(-D(q \| p) - H(q))} \quad \square \end{aligned}$$

Th 11.1.3 (CT). Taille de la classe de type q . Quelle que soit le type q de la séquence x^n .

$$|T(q)| \leq 2^{nH(q)}$$

Preuve: $|T(q)| =$ nb de manières d'arranger $nq(a_1), nq(a_2), \dots, nq(a_{|\mathcal{X}|})$ de 1 vecteur de taille n éléments, éléments bas

$$= \text{coefficient multinomial} = \binom{n}{nq(a_1), nq(a_2), \dots, nq(a_{|\mathcal{X}|})}$$

difficile à manipuler, d'où calcul de borne sup.

$$1 \geq \sum_{x^n \in T(q)} q(x^n) = \sum_{x^n} 2^{-nH(q)} = |T(q)| 2^{-nH(q)}$$

⚠ proba évolue avec le type q . c'est différent de $P(X^n = x^n)$ évalué avec vraie loi

Donc $|T(q)| \leq 2^{nH(q)}$

□

Th 11.1.4 CT.

probabilité de la classe de type q , (selon la loi \tilde{p}) ~~Th. 11.1.4~~

* Soit un type $q \in \mathcal{A}_n$. Soit 1 loi arbitraire \tilde{p} .

la proba de la classe de type q selon la loi \tilde{p} satisfait

$$\tilde{p}(T(q)) \leq 2^{-nD(q||\tilde{p})}$$

$$\tilde{p}(T(q)) = \sum_{x^n \in T(q)} \tilde{p}(x^n) = \sum_{x^n \in T(q)} 2^{-n(D(q||\tilde{p}) + H(q))}$$

11.1.2

$$\stackrel{11.1.3}{=} \underbrace{|T(q)|}_{\leq 2^{nH(q)}} 2^{-n(D(q||\tilde{p}) + H(q))} \leq 2^{-nD(q||\tilde{p})} \quad \square$$

En particulier $p(T(q)) \leq 2^{-nD(q||p)}$

Résumé

- 11.1.1 $|\mathcal{A}_n| \leq (m+1)^{|X|}$
- 11.1.2 $p(x^n) = 2^{-n(D(q||p) + H(q))}$
- 11.1.3 $|T(q)| \leq 2^{nH(q)}$
- 11.1.4 $p(T(q)) \leq 2^{-nD(q||p)}$

Définition ~~Soit~~ Un code à longueur fixe de ~~longueur~~ R pour le processus X_1, X_2, \dots de vraie distribution inconnue p , consiste en

fonction codage $f_n: \mathcal{X}^n \rightarrow \{1, \dots, 2^{nR}\}$

— de décodage $g_n: \{1, \dots, 2^{nR}\} \rightarrow \mathcal{X}^n$

- R est appelé débit (il faut R bits pour coder 1 symbole X_i)
- la proba d'erreur pour ce code (calculée avec la vraie loi p) est

$$P_e^{(n)} = \mathbb{P}_p(X^n : g_n(f_n(X^n)) \neq X^n)$$

Évalué avec p .

Définition Un code de débit R est dit universel si f_n et g_n ne dépendent pas de p et

$$\lim_{n \rightarrow \infty} P_e^{(n)} \rightarrow 0$$

Théorème de codage universel. Il existe une suite de codes universels de paramètres $(2^{nR}, n)$ tq pour toute source de loi $p(x)$ tq $H(p) < R$, alors $P_e^n \rightarrow 0$ $n \rightarrow \infty$.

Interprétation: même taux de compression asymptotique $\left\{ \begin{array}{l} \text{avec} \\ \text{sans} \end{array} \right.$ connaissance de la loi!!!

Preuve Soit R le débit de ce code. Soit

$$R_n = R - |\mathcal{X}| \frac{\log(n+1)}{n}$$

$$f_n: \mathcal{X}^n \rightarrow \{1, \dots, 2^{nR}\}$$

$x^n \mapsto$

← identification
logs q fournis
 2_{2n}

← identification vecteur de type q

longueur = $n R_n$

logues $|\mathcal{X}| \log(n+1)$

CONCLUSION

- même taux de compression asymptotique avec/sans connaissance loi.

"normal"

$$\text{codes loi} \approx |X| \log n$$

$$\text{codes vecteur} \approx H(p) \log n$$

- à taille finie le coût supplémentaire est

(admis)

$$\frac{\# \text{ degré de liberté loi}}{2} \log n$$

$$\text{dit} \quad \frac{|X|-1}{2} \log n$$

$$\text{Markov order 1} \quad \frac{|X|-1}{2} |X| \log n$$

- l'algorithme en 2 étapes est optimal.

- ce principe est utilisé en ML pour faire sélection modèle
(model selection : MDL minimum description length)

Conclusion

le principe du codage universel CV est utilisé en ML, statistique et s'offre MDL =

appel stat : à partir d'une observation d'un v.o. que dire sur la loi
[codage sans perte : ... message, trouver la description la + courte]

le principe MDL : choisir "la bonne loi" est la loi qui comprime bien les données

CV est aussi utilisé en ML :

Classification

x_1, \dots, x_n

n points

à séparer

en K classes (linéaires) mais combien de classes ?



algo : init K classes (ex centres) (C_k random)

iterate : - association à centre

$$\forall k \quad C_k = \{x_i \text{ s.t. } |x_i - c_k| < |x_i - c_{k'}|, \forall k' \neq k\}$$

- calculer centroide

$c_k = \text{centre de } C_k$

$$c_k = \arg \min_{c \in C_k} \sum_{x \in C_k} \|x - c\|^2$$

Chose K tq $\forall x_i, \quad x_i \text{ codage } (x_i) : \overbrace{\quad}^k \quad \overbrace{\quad}^{x_i - c_k}$

$$\min_K \sum_{1 \leq i \leq n} L(\text{codage}_K(x_i)) + L(\text{codage}_K(x_i - c_k))$$

$Q \rightarrow P$
m-juni

Cette technique est 1 parmi d'autres (pour les sources iid)

LZ est 1 code à longueur variable pour les sources ergodiques, dont la caractéristique de la loi est compléte.

Algo Lempel Ziv 78-78 = repose sur analyse syntaxique (parsing) du message à coder.

parsing d'un chaîne de binaire (binary string) = est une segmentation de la chaîne en phrases séparées par des virgules

distinct parsing = parsing où toutes les phrases sont \neq

LZ parsing = distinct parsing où chaque phrase est la + petite phrase qui n'est pas déjà apparue

$x_1, x_n = 0, 1, 01, 11, 00, 10, 000, 101, 001, 110, 1010, 111$
 (étiquettes 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11)
 $c(x_1, \dots, x_n) = \#$ phrases dans le LZ parsing de la chaîne x_1, \dots, x_n

Algo pour coder x_1, \dots, x_n

1. faire LZ parsing.
2. compter $c(x_1, \dots, x_n)$ et coder $c(x_1, \dots, x_n)$ avec m bits.
3. étiqueter chaque phrase avec 1 mot de longueur m bits.
4. La k ième phrase est du type $y_k = \begin{cases} y_{mk} 0 & \text{ou } m_k = \text{étiquette de } y_{mk} \\ y_{mk} 1 \end{cases}$

Coder cette k ième phrase (y_k) sous forme $\begin{cases} (m_k, 0) \\ (m_k, 1) \end{cases}$

~~est~~ nécessaire. $\log(c(x_1, \dots, x_n) + 1) + 1$
 (explique gain si $y_{mk} = \emptyset$ ou 1)

Ex $c(x_1, \dots, x_n) = 11$
 $c(\dots) + 1 = 12 \rightarrow 4$ bits

$LZ(x_1, \dots, x_n) = \underbrace{0000, 0000}_0, \underbrace{1}_1, \underbrace{0001, 1}_1, \underbrace{0010, 1}_2, \underbrace{0001, 0}_1, \underbrace{0010, 0}_2, \underbrace{0101, 0}_5$

En pratique $\left\{ \begin{array}{l} \cdot \text{ si mémoire finie, les phrases anciennes sont érasées par les nouvelles} \\ \cdot \text{ variante célèbre de Storer Szymanski: on ne code que si codage est + court que la phrase d'origine} \end{array} \right.$
 en ajoutant 1 bit de flag
 0 si codage direct
 1 si codage avec position