

Pour information

- Page web du cours <http://www.di.ens.fr/~fbach/courses/fall2010/>

1.1 Rappels de proba

1.1.1 Notations

Soit $\{X_1, \dots, X_n\}$ un ensemble de variables aléatoires, de distribution :

$$\mathbb{P}(X_1 = x_1, \dots, X_p = x_p) = p(x_1, \dots, x_p)$$

On note (de manière ambiguë) $\mathbb{P}(X_i = x_i) = p(x_i)$.

Pour $A, B \in \{1, \dots, p\}$, on note $X_A = (X_i)_{i \in A}$ et $X_B = (X_j)_{j \in B}$.

1.1.2 Quelques définitions / formules

- Loi marginale : $p(x_A) = \sum_{x_{A^c}} p(x_A, x_{A^c})$
- Loi conditionnelle : $p(x_A | x_B) = \frac{p(x_A, x_B)}{p(x_B)}$ si $p(x_B) \neq 0$
- Formule de Bayes : $p(x_A | x_B) = \frac{p(x_B | x_A) p(x_A)}{p(x_B)}$

1.1.3 Espérances

- Espérance de X_i : $\mathbb{E}[X_i] = \sum_{x_i} x_i \cdot p(x_i)$
- Espérance de $f(X_i)$, pour f mesurable : $\mathbb{E}[f(X_i)] = \sum_{x_i} f(x_i) \cdot p(x_i)$
- Variance :

$$\begin{aligned} \text{Var}(X_i) &= \mathbb{E}[(X_i - \mathbb{E}[X_i])^2] \\ &= \mathbb{E}[X_i^2] - \mathbb{E}[X_i]^2 \end{aligned}$$

- Covariance :

$$\begin{aligned} \text{Cov}(X_i, X_j) &= \mathbb{E}[(X_i - \mathbb{E}[X_i])(X_j - \mathbb{E}[X_j])] \\ &= \mathbb{E}[X_i X_j] - \mathbb{E}[X_i] \mathbb{E}[X_j] \end{aligned}$$

1.1.4 Indépendance

- Indépendance de 2 variables :
 X_i, X_j indépendantes $\Leftrightarrow p(x_i, x_j) = p(x_i)p(x_j) \quad \forall x_i, x_j$.
 On note $X_i \perp X_j$.
- Indépendance de n variables :
 X_1, \dots, X_n indépendantes $\Leftrightarrow p(x_1, \dots, x_n) = p(x_1) \cdots p(x_n) \quad \forall x_1, \dots, x_n$.
- Attention : indépendance de n variables \Rightarrow indépendance 2 à 2,
 mais *la réciproque est fautive* !

1.1.5 Espérances et indépendance conditionnelles

En considérant la distribution (conditionnelle) $p_{x_j}(x_i) = p(x_i|x_j) = \frac{p(x_i, x_j)}{p(x_j)}$ au lieu de $p(x_i)$, on obtient les notions d'espérances conditionnelles.

- Espérance conditionnelle de X_i : $\mathbb{E}[X_i|X_j] = \sum_{x_i} x_i \cdot p(x_i|x_j)$
- Variance conditionnelle :

$$\begin{aligned} \text{Var}(X_i) &= \mathbb{E}[(X_i - \mathbb{E}[X_i|X_j])^2 | X_j] \\ &= \mathbb{E}[X_i^2 | X_j] - \mathbb{E}[X_i | X_j]^2 \end{aligned}$$

- Indépendance conditionnelle :
 X_i, X_k indépendantes conditionnellement à $X_j \Leftrightarrow p(x_i, x_k|x_j) = p(x_i|x_j)p(x_k|x_j) \quad \forall x_i, x_j, x_k$
 On note $X_i \perp X_k | X_j$.

Exemples et contres-exemples : on réalise deux lancers d'une pièce de monnaie. Soient les variables aléatoire X_1, X_2 valant 1 si le lancer donne pile, 0 sinon, et $X_3 = XOR(X_1, X_2)$.

On montre alors que :

- les $(X_i)_{i=1,2,3}$ sont indépendantes 2 à 2 ;
- les $(X_i)_{i=1,2,3}$ ne sont pas indépendantes ;
- deux des $(X_i)_{i=1,2,3}$ ne sont pas indépendantes sachant la troisième.

1.2 Modèles à un noeud

1.2.1 Modèle ?

Notations

X_1, \dots, X_n désignent des variables aléatoires IID (Indépendantes et Identiquement Distribuées) ;

x_1, \dots, x_n sont des observations de X ;

$\{x_1, \dots, x_n\}$ est l'échantillon observé.

Définition

Un modèle est un ensemble de distributions avec des paramètres : $\{p_\theta(x), \theta \in \Theta\}$, avec généralement $\Theta = \mathbb{R}^p$.

Par exemple : pour une gaussienne à une dimension, les paramètres sont la moyenne et la variance : $\theta = (\mu, \sigma^2) \in \mathbb{R} \times \mathbb{R}_+^*$.

L'objectif est de déterminer (d'inférer) les paramètres θ du modèle à partir d'observations (de réalisations) de ce dernier.

Approche fréquentiste

Les paramètres optimaux θ^* sont définis comme extrênum d'une fonction de contraste (ou fonction de coût / de perte).

$\hat{\theta} = T(X_1, \dots, X_n)$ étant un estimateur, le but est de faire tendre $\hat{\theta}$ vers θ^* .

Approche bayésienne

On fait un *a priori* $p(\theta)$ sur les paramètres que peut prendre le modèle.

Les observations sont des probabilités, connaissant les paramètres : $p(x|\theta)$.

La probabilité *a posteriori* du modèle est alors $p(\theta|x) = \frac{p(x|\theta)p(\theta)}{p(x)} \propto p(x|\theta)p(\theta)$.

Maximum de vraisemblance (Fisher)

La vraisemblance (*likelihood*) d'un modèle est la quantité $L(\theta) = p_\theta(x_1, \dots, x_n)$, vue comme fonction de θ .

On utilise souvent la log-vraisemblance, $l(\theta) = \log L(\theta)$.

L'estimateur du maximum de vraisemblance est le maximiseur de la vraisemblance :

$$\hat{\theta} = \arg \max_{\theta \in \Theta} L(\theta)$$

1.2.2 Modèles élémentaires**Loi de Bernoulli****Définition**

$X \in \{0, 1\}$ est une variable aléatoire et $\theta \in [0, 1]$ un paramètre.

X suit une loi de Bernoulli, $X \sim \text{Ber}(\theta)$, si :

$$\begin{cases} p(X = 1) &= \theta \\ p(X = 0) &= 1 - \theta \end{cases}$$

Une autre manière de l'écrire :

$$p(x) = \theta^x (1 - \theta)^{(1-x)} \tag{1.1}$$

Propriétés

- $\mathbb{E}[X] = \theta$
- $\text{Var}(X) = \theta(1 - \theta)$
- Si $X_1, \dots, X_n \sim \text{Ber}(\theta)$ iid, alors $Z = \sum_i X_i \sim \text{Binom}(n, \theta)$

Maximum de vraisemblance

Soient X_1, \dots, X_n iid de loi $\text{Ber}(\theta)$.

Vraisemblance :

$$L(\theta) = p_\theta(x_1, \dots, x_n) = \prod_{i=1}^n p_\theta(x_i)$$

Log-vraisemblance, en utilisant l'expression (??) :

$$\begin{aligned} l(\theta) &= \sum_{i=1}^n \log p_\theta(x_i) \\ &= \sum_{i=1}^n (x_i \log \theta + (1 - x_i) \log(1 - \theta)) \\ &= \left(\sum_{i=1}^n x_i \right) \log \theta + \left(n - \sum_{i=1}^n x_i \right) \log(1 - \theta) \end{aligned}$$

Cette expression est dérivable :

$$\frac{\partial l}{\partial \theta}(\theta) = \frac{1}{\theta} \sum_{i=1}^n x_i - \frac{1}{1 - \theta} \left(n - \sum_{i=1}^n x_i \right)$$

Cette log-vraisemblance étant concave, son maximum est atteint quand la dérivée s'annule, d'où :

$$\hat{\theta} = \frac{1}{n} \sum_{i=1}^n x_i$$

Loi multinomiale**Définition**

$X \in \{1, \dots, q\}$ est une variable aléatoire.

Les paramètres sont $\Pi = (\pi_1, \dots, \pi_q) \in \mathbb{R}_+^q$, qui vérifient :

$$\sum_k \pi_k = 1 \tag{1.2}$$

La distribution d'une loi multinomiale $X \sim \text{Multi}(1, \pi_1, \dots, \pi_q)$ est : $p_\Pi(k) = \pi_k \quad \forall k$.

Une autre manière de l'écrire :

$$p_{\Pi}(x) = \prod_{k=1}^q \pi_k^{\delta(x,k)} \quad (1.3)$$

Pour $X_i \sim \text{Multi}(1, \pi_1, \dots, \pi_q)$, on définit $\Delta_i = (\delta_{i,1}, \dots, \delta_{i,q})$ avec $\delta_{i,k} = \delta(X_i, k) \quad \forall k$.

Propriétés

Si $X_1, \dots, X_n \sim \text{Multi}(1, \pi_1, \dots, \pi_q)$ iid, alors :

$$Z' = \sum_{i=1}^n \Delta_i \sim \text{Multi}(n, \pi_1, \dots, \pi_q)$$

Et :

$$p(Z' = (n_1, \dots, n_q)) = \frac{n!}{n_1! \dots n_q!} \prod_{k=1}^q \pi_k^{n_k}$$

Maximum de vraisemblance

Soient X_1, \dots, X_n iid de loi $\text{Multi}(1, \pi_1, \dots, \pi_q)$ (et donc $\Delta_1, \dots, \Delta_n$ sont iid).

Vraisemblance :

$$L(\Pi) = p_{\Pi}(x_1, \dots, x_n) = \prod_{i=1}^n p_{\Pi}(x_i)$$

Log-vraisemblance, utilisant l'expression (??) :

$$\begin{aligned} l(\Pi) &= \sum_{i=1}^n \log p_{\Pi}(x_i) \\ &= \sum_{i=1}^n \sum_{k=1}^q (\delta_{i,k} \log \pi_k) \\ &= \sum_{k=1}^q n_k \log \pi_k \end{aligned}$$

Où on a posé $n_k = \sum_i \delta_{i,k}$ le nombre d'observations ayant pour valeur k .

Pour prendre en compte la contrainte (??), considérons le lagrangien :

$$\mathcal{L}(\Pi, \lambda) = - \underbrace{\sum_{k=1}^q n_k \log \pi_k}_{\text{log-likelihood}} + \lambda \left(\sum_{k=1}^q \pi_k - 1 \right)$$

L'objectif est de maximiser la log-vraisemblance sous la contrainte (??), ce qui revient à chercher :

$$\min_{\Pi \in \mathbb{R}^q} \left[\max_{\lambda \in \mathbb{R}_+} \mathcal{L}(\Pi, \lambda) \right]$$

\mathcal{L} étant convexe pour Π , ce problème est équivalent à :

$$\max_{\lambda \in \mathbb{R}_+} \left[\min_{\Pi \in \mathbb{R}^q} \mathcal{L}(\Pi, \lambda) \right]$$

Dérivons alors \mathcal{L} par rapport aux composantes de Π :

$$\frac{\partial \mathcal{L}}{\partial \pi_k}(\Pi, \lambda) = -\frac{n_k}{\pi_k} + \lambda$$

Le minimum de \mathcal{L} par rapport à Π est atteint en annulant cette expression $\forall k$, d'où :
 $n_k = \lambda \pi_k$.

En sommant cette relation pour $k = 1 \dots q$ et en utilisant la contrainte (??), il vient que
 $n = \lambda$, et la relation précédente donne donc :

$$\hat{\pi}_k = \frac{n_k}{n}$$

Loi gaussienne (1D)

Définition

X est une variable aléatoire réelle, et $\theta = (\mu, \sigma^2) \in \mathbb{R} \times \mathbb{R}_+^*$ sont les paramètres.
 X suit une loi gaussienne, $X \sim \mathcal{N}(\mu, \sigma^2)$, signifie :

$$p_{\mu, \sigma^2}(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2} \frac{(x - \mu)^2}{\sigma^2}\right)$$

Maximum de vraisemblance

Soient $X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$ iid.

Vraisemblance :

$$L(\mu, \sigma^2) = p_{\mu, \sigma^2}(x_1, \dots, x_n) = \prod_{i=1}^n p_{\mu, \sigma^2}(x_i)$$

Log-vraisemblance :

$$\begin{aligned} l(\mu, \sigma^2) &= \sum_{i=1}^n \log p_{\mu, \sigma^2}(x_i) \\ &= -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \end{aligned}$$

Dérivées, à annuler pour obtenir le maximum :

$$\begin{aligned}\frac{\partial l}{\partial \mu}(\mu, \sigma^2) &= \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) \\ &= \frac{1}{\sigma^2} \left(\sum_{i=1}^n x_i - n\mu \right) \\ \frac{\partial l}{\partial \sigma^2}(\mu, \sigma^2) &= -\frac{n}{2} \frac{2\pi}{2\pi\sigma^2} + \frac{1}{2} \frac{1}{\sigma^4} \sum_{i=1}^n (x_i - \mu)^2 \\ &= -\frac{1}{2\sigma^4} \left(n\sigma^2 - \sum_{i=1}^n (x_i - \mu)^2 \right)\end{aligned}$$

Ce qui nous donne :

$$\begin{aligned}\hat{\mu} &= \frac{1}{n} \sum_{i=1}^n x_i \\ \hat{\sigma}^2 &= \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2\end{aligned}$$

Loi gaussienne (kD)

Définition

X est une variable aléatoire à valeur dans \mathbb{R}^d .

$\theta = (\mu, \Sigma) \in \mathbb{R}^d \times \mathbb{R}^{d \times d}$ sont les paramètres (avec Σ définie positive).

X suit une loi gaussienne, $X \sim \mathcal{N}(\mu, \Sigma)$, signifie :

$$p_{\mu, \Sigma}(x) = \frac{1}{\sqrt{(2\pi)^d \det \Sigma}} \exp\left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)\right)$$

Maximum de vraisemblance

Soient $X_1, \dots, X_n \sim \mathcal{N}(\mu, \Sigma)$ iid.

Log-vraisemblance :

$$\begin{aligned}-l(\mu, \Sigma) &= -\sum_{i=1}^n \log p_{\mu, \Sigma}(x_i) \\ &= \frac{nd}{2} \log(2\pi) - \frac{n}{2} \log(\det \Sigma^{-1}) + \frac{1}{2} \sum_{i=1}^n (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)\end{aligned}$$

Par convexité, le minimum de $-l(\mu, \Sigma)$ est le point en lequel son gradient s'annule.

Gradient par rapport à μ :

$$\begin{aligned} -\nabla_{\mu} l(\mu, \Sigma) &= \sum_{i=1}^n \Sigma^{-1} (x_i - \mu) \\ &= \Sigma^{-1} \left(\sum_{i=1}^n x_i - n\mu \right) \end{aligned}$$

Gradient par rapport à Σ^{-1} :

$$-\nabla_{\Sigma^{-1}} l(\mu, \Sigma) = -\frac{n}{2} (\Sigma^{-1})^{-1} + \frac{1}{2} \sum_{i=1}^n (x_i - \mu) (x_i - \mu)^T$$

D'où :

$$\begin{aligned} \hat{\mu} &= \frac{1}{n} \sum_{i=1}^n x_i \\ \hat{\Sigma} &= \frac{1}{n} \sum_{i=1}^n (x_i - \mu) (x_i - \mu)^T \end{aligned}$$