

Rappel: Formule de Bayes
et Proba totale

Bayes:

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

Proba totale: Si A_1, A_2, \dots disjoints, et $\bigcup_{i \geq 1} A_i = \Omega$

$$P(B) = \sum_{i \geq 1} P(B|A_i) P(A_i)$$

↳

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

$$P(B) = P(B|A) P(A) + P(B|A^c) \cdot P(A^c)$$

Exemple: Oubli de la fréquence
de base (Expl. 2.4.6)

Situation: Maladie dans une population

p = proportion de gens infectés

On a un teste pour vérifier si qqn est infecté ou non.

Teste rate de temps en temps:

- Test dit "infecté" alors que l'individu est sain: \rightarrow proba $\frac{3}{1000}$

- Test dit "sain" alors que l'individu est infecté: \rightarrow proba $\frac{1}{1000}$

Modélisation:

$$\Omega = \{I_+, I_-\} \times \{T_+, T_-\}$$

\uparrow ind. infecté
 \downarrow ind. sain
 \uparrow teste pos.
 \downarrow teste nég.

Q1: $P(\text{test rate}) = ?$

On connaît: $* P(\text{ind infecté}) = P(\{(I_+, T_-), (I_+, T_+)\}) = p$

$* P(\text{ind sain}) = 1 - p$

$* P(\text{test rate} | \text{Ind sain}) = \frac{9}{1000}$

$* P(\text{test rate} | \text{Ind inf.}) = \frac{1}{1000}$

proba tot.

$$P(\text{test rate}) = \underbrace{P(\text{test rate} | \text{Ind sain})}_{= 9/1000} \underbrace{P(\text{ind sain})}_{= 1-p} + \underbrace{P(\text{test rate} | \text{Ind inf.})}_{= 1/1000} \underbrace{P(\text{ind. inf.})}_{= p}$$

$$= \frac{9}{1000} (1-p) + \frac{1}{1000} p \stackrel{p \leq 1}{\leq} \frac{10}{1000} = \frac{1}{100}$$

$1-p \leq 1$

\hookrightarrow test est fiable à $99/100$ indépendamment

↳ test est fiable à $\approx 100\%$ indépendamment de la valeur de p .

Q2: Si on teste un individu positif, quelle est la proba qu'il soit effectivement infecté ?

$$\mathbb{P}(\underbrace{\{(I_+, T_-), (I_+, T_+)\}}_{\text{infecté}} \mid \underbrace{\{(I_-, T_+), (I_+, T_+)\}}_{\text{test positif}}) = ?$$

$$= 1 - \mathbb{P}(\text{test nég} \mid \text{inf}) = 1 - \frac{1}{1000}$$

$$\mathbb{P}(\text{infecté} \mid \text{test pos.}) = \frac{\mathbb{P}(\text{test pos} \mid \text{inf.}) \mathbb{P}(\text{inf.})}{\mathbb{P}(\text{test pos} \mid \text{inf.}) \mathbb{P}(\text{inf.}) + \mathbb{P}(\text{test pos} \mid \text{sain}) \mathbb{P}(\text{sain})}$$

$$= \frac{\frac{999}{1000} \cdot p}{\frac{999}{1000} \cdot p + \frac{9}{1000} (1-p)} = \frac{999p}{999p + 9}$$

⚠ cette proba est petite si p petit.

Exemple 2.4.7 (Paradoxe de Simpson)

	All		Women		Men	
	Applicants	Admitted	Applicants	Admitted	Applicants	Admitted
Total	12763	41%	4321	35%	8442	44%

Admission à Berkeley.

Dep.	All		Women		Men	
	Applicants	Admitted	Applicants	Admitted	Applicants	Admitted
A	933	64%	108	82%	825	62%
B	585	63%	25	68%	560	63%
C	918	35%	593	34%	325	37%
D	792	34%	375	35%	417	33%
E	584	25%	393	24%	191	28%
F	714	6%	341	7%	373	6%
Total	4526	39%	1835	30%	2691	45%

Dep.	All		Women		Men	
	Applicants	Admitted	Applicants	Admitted	Applicants	Admitted
A	933	64%	108	82%	825	62%
C	918	35%	593	34%	325	37%
Total	1851	39%	701	38%	1150	62%

55

$$\mathbb{P}(\text{Admis} | H) = \mathbb{P}(\text{Admis} | H, \text{depA}) \mathbb{P}(\text{depA} | H) + \mathbb{P}(\text{Admis} | H, \text{depC}) \mathbb{P}(\text{depC} | H)$$

$$= \frac{62}{100} \cdot \frac{825}{1150} + \frac{37}{100} \cdot \frac{325}{1150}$$

$$\mathbb{P}(\text{Admis} | F) = \frac{82}{100} \cdot \frac{108}{701} + \frac{34}{100} \cdot \frac{593}{701}$$

Le "paradoxe" veut dire que les hommes postulent plus au département qui accepte plus.

Corrélation (section 2.5)

Variance: Soit X une variable aléatoire est

$$\text{Var}(X) := \mathbb{E}[(X - \mathbb{E}(X))^2] = \mathbb{E}[X^2] - \mathbb{E}[X]^2$$

↳ mesure à quel point X est typiquement

↳ mesure à quel point X est typiquement loin de $E[X]$.

Écart-type: (deviation standard)

$$\sigma_X = \sqrt{\text{Var}(X)} \rightarrow \text{"typiquement"}, X \in [E(X) - \sigma_X, E(X) + \sigma_X].$$

Remarque: On peut définir alternativement

$$\text{Var}(X) = \inf_{\mu \in \mathbb{R}} E[(X - \mu)^2]$$

$$E[X] = \arg \uparrow$$

"Prv": * $E[(X - \mu)^2] = E[X^2] + \mu^2 - 2\mu E[X]$

* On cherche les minima en μ :

$$2\mu - 2E[X] = 0 \Leftrightarrow \mu = E[X].$$

Covariance: Soient X, Y deux variables aléatoires,

$$\text{Cov}(X, Y) := E[XY] - E[X]E[Y].$$

Remarque 1: Si X et Y sont indépendantes,

$$\begin{aligned} \text{Cov}(X, Y) &= E[XY] - E[X]E[Y] \\ &\stackrel{\text{indép}}{=} E[X]E[Y] - E[X]E[Y] = 0 \end{aligned}$$

$$\stackrel{\text{indép}}{=} E[X]E[Y] - E[X]E[Y] = 0.$$

Beaucoup de résultats sont "simples" à montrer sous l'hyp. "X et Y sont indép.", l'hyp "Cov(X, Y) = 0" est souvent une manière naturelle d'affaiblir l'hyp. d'indép.

Remarque 2:

La covariance apparaît naturellement quand on calcule la variance d'une somme:

$$\begin{aligned} \text{Var}(X+Y) &= E[(X+Y)^2] - E[X+Y]^2 \\ &= E[X^2] + E[Y^2] + 2E[XY] \\ &\quad - E[X]^2 - E[Y]^2 - 2E[X]E[Y] \\ &= \text{Var}(X) + \text{Var}(Y) + 2\text{Cov}(X, Y) \end{aligned}$$

↳ Cov(X, Y) est une correction à "l'additivité" pour la variance.

Propriétés: (lemme 2.5.2 + def)

- 1) $\text{Cov}(X, X) = \text{Var}(X)$
- 2) $\text{Cov}(X, Y) = \text{Cov}(Y, X)$
- 3) $\text{Cov}(aX, bY) = a \cdot b \cdot \text{Cov}(X, Y) \quad \forall a, b \in \mathbb{R}$
- 4) $\text{Cov}(X, Y_1 + Y_2) = \text{Cov}(X, Y_1) + \text{Cov}(X, Y_2)$

↳ $(X, Y) \mapsto \text{Cov}(X, Y)$ est linéaire en chaque

$\hookrightarrow (X, Y) \mapsto \text{Cov}(X, Y)$ est linéaire en chaque coordonnée.

Remarque: Indép \Rightarrow Cov = 0

Mais Cov = 0 $\not\Rightarrow$ indép.

Expl: Soient X, Z deux v.a. indép. telles que

- $\mathbb{P}(X=10) = \mathbb{P}(X=0) = 1/2$
- $\mathbb{P}(Z=-1) = \mathbb{P}(Z=+1) = 1/2$

On définit $Y = X \cdot Z$.

$$\begin{aligned} \text{On a } \text{Cov}(X, Y) &= E[XY] - E[X]E[Y] \\ &\stackrel{\text{def } Y}{=} E[X^2 Z] - E[X]E[X \cdot Z] \\ &\stackrel{X, Z \text{ indép}}{=} E[X^2] \underbrace{E[Z]}_{=0} - E[X]^2 \underbrace{E[Z]}_{=0} = 0 \end{aligned}$$

Mais X et Y ne sont pas indépendantes:

$$\begin{aligned} \mathbb{P}(X=0, Y=0) &= \mathbb{P}(X=0, XZ=0) \\ &\stackrel{Z \neq 0 \text{ p.s.}}{=} \mathbb{P}(X=0) = 1/2 \end{aligned}$$

$$\begin{aligned} \mathbb{P}(X=0) \mathbb{P}(Y=0) &= \mathbb{P}(X=0) \mathbb{P}(XZ=0) \\ &\stackrel{Z \neq 0 \text{ p.s.}}{=} \mathbb{P}(X=0)^2 = 1/4 \neq 1/2 \end{aligned}$$

On a un cas particulier:

lemme (2.5.1) c. v. v. . . .

Lemme (2.5.1) Si X et Y sont des variables de Bernoulli ($\mathbb{P}(X \in \{0,1\}) = 1$
 $\mathbb{P}(Y \in \{0,1\}) = 1$)

Alors X et Y sont indép si et seulement si $\text{Cov}(X, Y) = 0$.

Indication: $\forall c \in \mathbb{R}, \text{Cov}(X, c) = 0$.

Coefficient de corrélation: Soient X, Y deux v.c.

$$\rho_{XY} = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \text{Var}(Y)}} = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} \in [-1, 1].$$

↑ plus tard

ρ_{XY} est une manière indép. des unités de mesurer la corrélation:

Si deux quantités sont exprimées en X_1, Y_1 m. ou en X_2, Y_2 m.m.

$$\begin{aligned} \hookrightarrow \text{Cov}(X_1, Y_1) &= \text{Cov}\left(\frac{X_2}{100000}, \frac{Y_2}{100000}\right) \\ &= \frac{1}{100000} \text{Cov}(X_2, Y_2) \end{aligned}$$

mais $\rho_{X_1 Y_1} = \rho_{X_2 Y_2} \quad \nabla$