

Hiver 2019

Méthodes d'apprentissage  
**IFT603**

Formulation probabiliste

Par  
Pierre-Marc Jodoin  
et  
Hugo Larochelle

1

## Variable aléatoire

- La théorie des probabilités est l'outil idéal pour formaliser nos **hypothèses et incertitudes** par rapport à nos données
- On va traiter nos données comme des **variables aléatoires**
  - la valeur d'une variable aléatoire est incertaine (avant de l'observer)
  - la loi de probabilité de la variable aléatoire caractérise notre incertitude par rapport à sa valeur

2

## Variable aléatoire

- Soient  $X$  et  $Y$  des variables aléatoires **discrètes**
  - $X$  peut prendre comme valeurs  $x_1, \dots, x_M$
  - $Y$  peut prendre comme valeurs  $y_1, \dots, y_M$
- La **probabilité jointe** qu'on observe  $X=x_i$  et  $Y=y_j$  est notée

$$P(X = x_i, Y = y_j)$$

et se lit comme la « probabilité d'observer à **la fois**  $x_i$  **et**  $y_j$  ».

- Note:

$$P(X = x_i, Y = y_j) = P(Y = y_j, X = x_i)$$

3

## Probabilité marginale

Une **probabilité marginale** est lorsqu'on ne s'intéresse pas à toutes les variables aléatoire qu'on a défini

Exemple : la probabilité marginale d'observer  $X=x_i$

$$P(X = x_i) = \sum_{j=1}^N P(X = x_i, Y = y_j)$$

4

## Probabilité conditionnelle

Une **probabilité conditionnelle** est lorsqu'on s'intéresse la valeur d'une variable aléatoire «étant donnée» une valeur assignée à d'autres variables

$$P(X = x_i | Y = y_j)$$

Se lit : la probabilité que  $X = x_j$  étant donné que  $Y = y_i$

<https://www.npr.org/2016/11/14/501737150/rural-voters-played-a-big-part-in-helping-trump-defeat-clinton>

5

## Probabilité conditionnelle

**Exemple, élections américaines 2016**

$$P(\text{Voter républicain}) = 46.1\%$$

**VS**

$$\left\{ \begin{array}{l} P(\text{Voter républicain} | \text{Zone urbaine}) = 35\% \\ P(\text{Voter républicain} | \text{Zone rurale}) = 62\% \\ P(\text{Voter républicain} | \text{Banlieu}) = 50\% \end{array} \right.$$

<https://www.npr.org/2016/11/14/501737150/rural-voters-played-a-big-part-in-helping-trump-defeat-clinton>

6

## Produit des probabilités

$x_i$  et  $y_j$  ont disparu,  
seulement pour  
simplifier la notation

Une probabilité jointe peut toujours être décomposée par le produit d'une probabilité conditionnelle et marginale

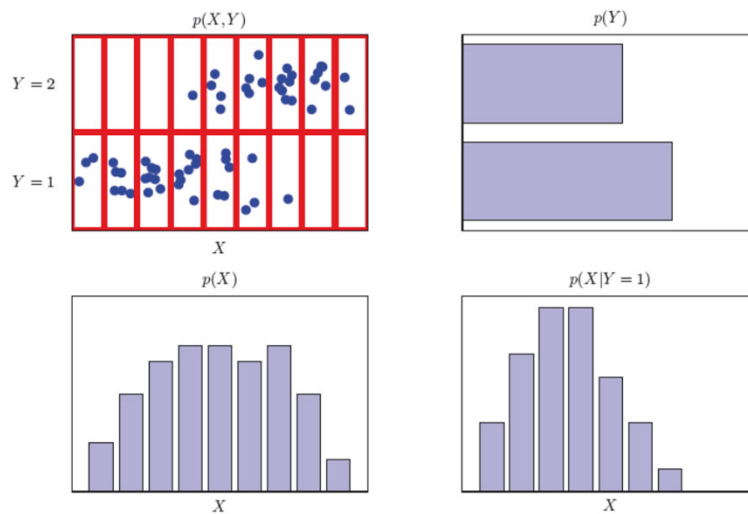
$$P(X, Y) = P(X | Y)P(Y)$$

En mots :

la probabilité d'observer  $X=x_i$  ET  $Y=y_j$ , c'est la probabilité d'observer  $Y=y_j$  multipliée par la probabilité d'observer  $X=x_i$  **étant donné que**  $Y=y_j$

7

## Probabilités jointes, marginale et conditionnelles



Crédit : Bishop

8

## Bayes

La règle de Bayes permet d'inverser l'ordre de la conditionnelle

$$P(Y | X) = \frac{P(X | Y)P(Y)}{P(X)}$$

$p(Y)$  est appelée loi de probabilité *a priori* (*prior*)

$p(Y | X)$  est appelée loi de probabilité *a posteriori* (*posterior*)

9

## Indépendance

Deux variables aléatoires  $X$  et  $Y$  sont indépendantes si

➤  $P(X, Y) = P(X)P(Y)$  ou

➤  $P(X | Y) = P(X)$  ou

➤  $P(Y | X) = P(Y)$

➤ Observer la valeur d'une variable ne nous apprend rien sur la valeur de l'autre

10

## Variable aléatoire continue

Soit  $X$  une **variable aléatoire continue**

- $X$  peut prendre un nombre infini de valeurs possibles (e.g.  $\mathbb{R}$ )
- $X$  est associée à une **fonction de densité** de probabilité  $p(x)$

la probabilité que  $X$  appartienne à un intervalle  $(a,b)$  est

$$p(x \in (a,b)) = \int_a^b p(x) dx$$

11

## Variables aléatoires continues

Soit  $X$  une **variable aléatoire continue**

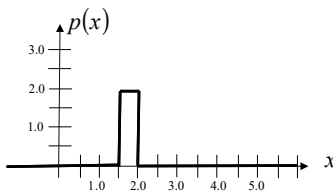
- la fonction de densité doit satisfaire

$$p(x) \geq 0$$

$$\int_{-\infty}^{\infty} p(x) dx = 1$$

à noter que, contrairement aux probabilités d'une variable discrète, la fonction de densité peut être  $> 1$ .

Exemple



$$p(x) = \begin{cases} 2 & \text{is } x \in [1.5, 2.0] \\ 0 & \text{sinon} \end{cases}$$

12

## Variables aléatoires continues

Soit  $X$  une **variable aléatoire continue**

la **fonction de répartition**  $P(z)$  (*cumulative distribution function*) donne la probabilité que  $X$  appartienne à l'intervalle  $(-\infty, z)$

$$P(z) = \int_{-\infty}^z p(x) dx$$

13

## Variables aléatoires continues

Soient  $X$  et  $Y$  deux **variables aléatoires continues**

➤ elles sont associées à une **fonction de densité jointe**  $p(x,y)$  telle que :

$$p(x \in [a, b], y \in [c, d]) = \int_a^b \int_c^d p(x, y) dx dy$$

14

## Variables aléatoires continues

Soient  $X$  et  $Y$  deux **variables aléatoires continues**

➤ La **fonction de densité marginale** s'obtient en intégrant l'autre variable

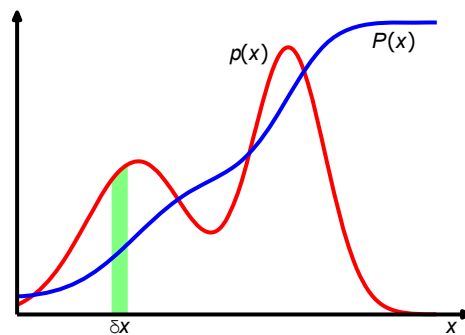
$$p(x) = \int p(x, y) dy$$

➤ La **fonction de densité conditionnelle** s'obtient comme auparavant

$$p(y | x) = \frac{p(x, y)}{p(x)}$$

15

## Variables aléatoires continues



16



## Expérance mathématique

L'**espérance** d'une **variable  $X$**  est la moyenne qu'on obtient si on répète un grand nombre de fois une expérience

$$E[X] = \sum_x xp(x) \quad (\text{cas discret})$$

$$E[X] = \int xp(x)dx \quad (\text{cas continu})$$

17

## Expérance mathématique

L'**espérance** d'une **fonction  $f(x)$**  est la moyenne qu'on obtient si on génère un grand nombre de valeurs pour cette fonction

$$E[f] = \sum_x f(x)p(x) \quad (\text{cas discret})$$

$$E[f] = \int f(x)p(x)dx \quad (\text{cas continu})$$

18

## Variance

- La **variance** d'une **variable**  $X$  est

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

- La **variance** d'une **fonction**  $f(x)$  est

$$\text{var}[f] = \mathbb{E}[(f(x) - \mathbb{E}[f(x)])^2]$$

La variance mesure à quel point les valeurs varient autour de l'espérance

19

## Propriétés de l'espérance et de la variance

**Transformation linéaire** de l'espérance

$$\begin{aligned} \mathbb{E}[ax + by] &= \sum_x \sum_y (ax + by)p(x, y) && \text{a, b sont réels} \\ &= a\mathbb{E}[x] + b\mathbb{E}[y] && \text{Si x, y indépendants} \end{aligned}$$

**Transformation linéaire** de la variance

$$\text{var}[ax + by] = a^2 \text{var}[x] + b^2 \text{var}[y] \quad \text{Si x, y indépendants}$$

20

## Espérance et variance conditionnelles

L'espérance et la variance se généralise au cas **conditionnel** :

$$E[x | y] = \sum x p(x | y)$$
$$E[f(x) | y] = \sum_x f(x) p(x | y)$$

$$\text{var}[x | y] = E\left[(x - E[x | y])^2\right]$$
$$\text{var}[f(x) | y] = E\left[(f(x) - E[f(x) | y])^2\right]$$

21

## Covariance

La covariance entre 2 variables aléatoires X et Y

$$\text{cov}[x, y] = E_{xy}[(x - E_x[x])(y - E_y[y])]$$
$$= E_{xy}[xy] - E_x[x]E_y[y]$$

mesure à quel point on peut prédire X à partir de Y (linéairement), et vice-versa  
si X et Y sont indépendantes, alors la covariance est 0

22

## Variables aléatoires multidimensionnels

Une variable aléatoire peut être un vecteur

L'espérance d'un vecteur est le vecteur des espérances

$$E[\vec{x}] = (E[x_1], \dots, E[x_D])^T$$

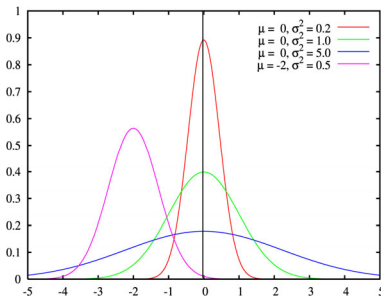
Et la covariance de deux vecteurs est

$$\text{cov}[\vec{x}, \vec{y}] = E_{\vec{x}\vec{y}}[\vec{x}\vec{y}^T] - E_{\vec{x}}[\vec{x}]E_{\vec{y}}[\vec{y}]$$

23

## Loi de probabilité gaussienne

$$N(x | \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



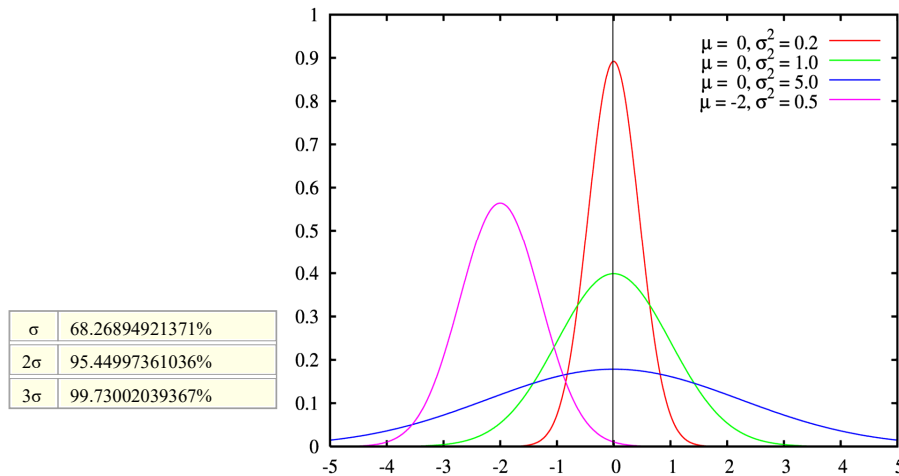
$$\text{Moyenne : } E[x] = \int_{-\infty}^{\infty} N(x; \mu, \sigma) x dx = \mu$$

$$\text{Variance : } \text{var}[x] = \int_{-\infty}^{\infty} N(x; \mu, \sigma) (x - \mu)^2 dx = \sigma^2$$

$$\text{Écart type : } \sqrt{\text{var}[x]} = \sigma$$

24

## Loi de probabilité gaussienne



25

## Gaussienne multivariée

$$N(\vec{x} | \vec{\mu}, \Sigma) = \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2} (\vec{x} - \vec{\mu})^T \Sigma^{-1} (\vec{x} - \vec{\mu})\right\}$$

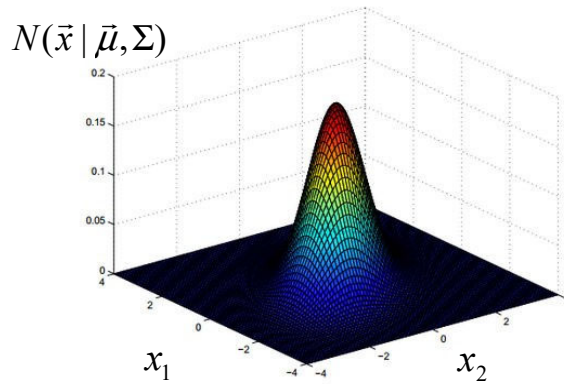
$$\text{Moyenne : } E[\vec{x}] = \vec{\mu}$$

$$\text{Variance : } \text{cov}[\vec{x}] = \Sigma$$

26

## Gaussienne multivariée

Exemple :  $\vec{x} = (x_1, x_2)$

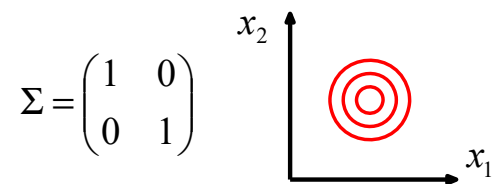


27

## Gaussienne multivariée

Exemple :  $\vec{x} = (x_1, x_2)$

Courbes de niveaux de  $N(\vec{x} | \vec{\mu}, \Sigma)$

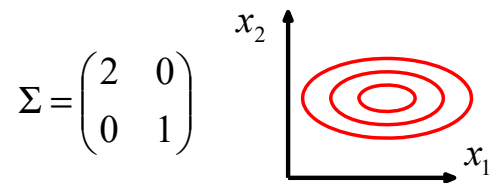


28

## Gaussienne multivariée

Exemple :  $\vec{x} = (x_1, x_2)$

Courbes de niveaux de  $N(\vec{x} | \vec{\mu}, \Sigma)$

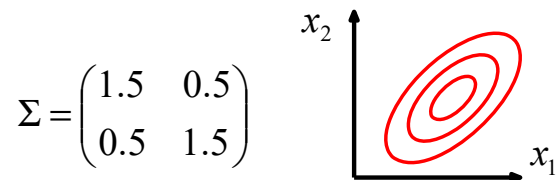


29

## Gaussienne multivariée

Exemple :  $\vec{x} = (x_1, x_2)$

Courbes de niveaux de  $N(\vec{x} | \vec{\mu}, \Sigma)$



30

## Gaussienne multivariée

Une **combinaison linéaire** de variables aléatoires gaussiennes est également gaussienne

- Exemple
  - soit  $x$  une variable gaussienne de moyenne  $\mu_1$  et variance  $\sigma_1^2$
  - soit  $y$  une variable gaussienne de moyenne  $\mu_2$  et variance  $\sigma_2^2$
  - alors  $ax + by$  suit une loi gaussienne de moyenne  $a\mu_1 + b\mu_2$  et variance  $a^2\sigma_1^2 + b^2\sigma_2^2$  ( $x$  et  $y$  sont indépendantes)

31

# Régression:

Maximum de vraisemblance

vs

Maximum a posteriori

32

32



# Introduction au tableau

Maximum de vraisemblance

vs

Maximum a posteriori

33

33

## Régression 1D

Retournons à notre exemple de régression

➤ **Données**

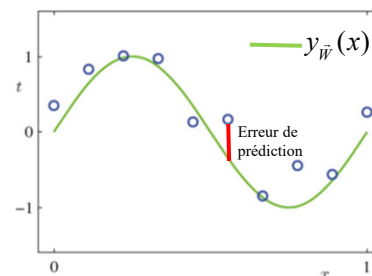
- ✓ entrée : scalaire  $x$
- ✓ cible : scalaire  $t$

➤ **Ensemble d'entraînement  $D$  contient:**

- ✓  $X = (x_1, \dots, x_N)^T$
- ✓  $T = (t_1, \dots, t_N)^T$

➤ **Objectif :**

- ✓ Faire une prédiction  $\hat{t}$  pour chaque nouvelle entrée  $\hat{x}$



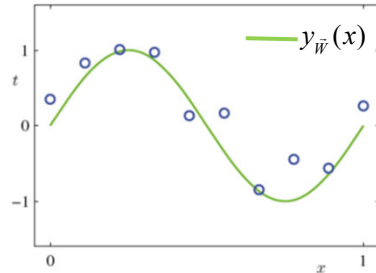
34

34

# Régression polynomiale

➤ On va supposer que nos données suivent une **forme polynomiale**

$$y_{\vec{w}}(x) = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M$$
$$= \sum_{i=0}^M w_i x^i$$



➤  $y_{\vec{w}}(x)$  est notre **modèle**

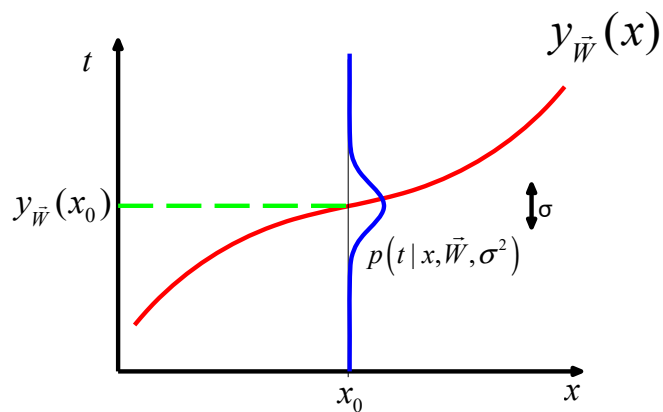
- ✓ Représente nos hypothèses sur le problème à résoudre
- ✓ Un modèle a toujours des paramètres qu'on doit trouver (ici  $\vec{w}$ )

35

35

# Régression polynomiale

Loi gaussienne conditionnelle

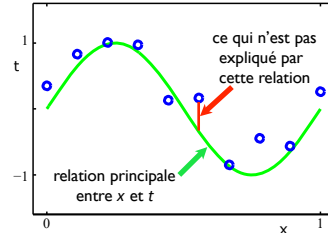


36

# Régression polynomiale

Loi gaussienne conditionnelle

On va formuler l'hypothèse que chaque donnée  $t$  a été générée selon une gaussienne de moyenne  $y_{\vec{w}}(x)$  et de variance  $\sigma^2$



$$p(T | X, \vec{W}, \sigma^2) = \prod_{n=1}^N p(t_n | x_n, \vec{W}, \sigma^2)$$

$$= \prod_{n=1}^N N(t_n | y_{\vec{w}}(x_n), \sigma^2)$$

**Variables indépendantes et identiquement distribuées (i.i.d)**

37

# Régression polynomiale

Maximum de vraisemblance  
(cf. notes au tableau)

Cherche les meilleurs paramètres  $\vec{W}$  en **maximisant la vraisemblance**

**Connue** → **Inconnue**

$$W = \arg \max_{\vec{w}} p(T | X, \vec{W}, \sigma^2)$$

$$= \arg \max_{\vec{w}} \prod_{n=1}^N N(t_n | x_n, \vec{W}, \sigma^2) \Rightarrow \text{Données gaussiens et (i.i.d)}$$

$$= \arg \max_{\vec{w}} \ln \left[ \prod_{n=1}^N N(t_n | y_{\vec{w}}(x_n), \sigma^2) \right]$$

$$= \arg \max_{\vec{w}} \sum_{n=1}^N \ln N(t_n | y_{\vec{w}}(x_n), \sigma^2)$$

38

# Régression polynomiale

Maximum de vraisemblance  
(cf. notes au tableau)

Cherche les meilleurs paramètres  $W$  en **maximisant la vraisemblance**

$$\begin{aligned} W &= \arg \max_W \sum_{n=1}^N \ln \mathcal{N}(t_n | y_W(x_n), \sigma^2) \\ &= \arg \max_W \sum_{n=1}^N \ln \left[ \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(t_n - y_W(x_n))^2}{2\sigma^2}} \right] \\ &= \arg \max_W \sum_{n=1}^N -\frac{(t_n - y_W(x_n))^2}{2\cancel{\sigma^2}} + \ln \left[ \frac{1}{\cancel{\sqrt{2\pi}\sigma}} \right] \\ &= \arg \min_W \frac{1}{2} \sum_{n=1}^N (t_n - y_W(x_n))^2 \end{aligned}$$

Indépendant de  $W$

39

# Régression polynomiale

Maximum de vraisemblance

Cherche les meilleurs paramètres  $W$  en **maximisant la vraisemblance**

$$W = \arg \min_W \frac{1}{2} \sum_{n=1}^N (t_n - y_W(x_n))^2$$

40

# Régression polynomiale

Maximum *a posteriori* (MAP)  
(cf. notes au tableau)

Cherche les meilleurs paramètres  $W$  en **maximisant la probabilité *a posteriori***

**Inconnue** →      ← **Connues**

$$\begin{aligned} W &= \arg \max_W p(W | X, T, \sigma^2) \\ &= \arg \max_W \frac{p(T | X, W, \sigma^2) p(W)}{P(X, T, \sigma^2)} \quad \Rightarrow \text{Par le théorème de Bayes} \\ &= \arg \max_W p(T | X, W, \sigma^2) p(W) \quad \leftarrow \text{Constante par rapport à } W \\ &= \arg \max_W \prod_{n=1}^N N(t_n | y_W(x_n), \sigma^2) p(W) \quad \Rightarrow \text{Données gaussiennes et (i.i.d)} \end{aligned}$$

41

# Régression polynomiale

Maximum *a posteriori* (MAP)  
(cf. notes au tableau)

Cherche les meilleurs paramètres  $W$  en **maximisant la probabilité *a posteriori***

$$W = \arg \max_W \prod_{n=1}^N N(t_n | y_W(x_n), \sigma^2) p(W)$$

Si on présuppose que les paramètres  $W$  suivent une distribution gaussienne centrée à zéro avec une matrice de variance  $\Sigma$

$$W = \arg \max_W \prod_{n=1}^N N(t_n | y_W(x_n), \sigma^2) N(W | 0, \Sigma)$$

$$W = \arg \max_W \ln \left[ \prod_{n=1}^N N(t_n | y_W(x_n), \sigma^2) N(W | 0, \Sigma) \right]$$

42

# Régression polynomiale

Maximum *a posteriori* (MAP)  
(cf. notes au tableau)

Cherche les meilleurs paramètres  $W$  en **maximisant la probabilité *a posteriori***

$$\begin{aligned}
 W &= \arg \max_W \ln \left[ \prod_{n=1}^N N(t_n | y_W(x_n), \sigma^2) N(W | 0, \Sigma) \right] \\
 &= \arg \max_W \sum_{n=1}^N \ln [N(t_n | y_W(x_n), \sigma^2) N(W | 0, \Sigma)] \\
 &= \arg \max_W \sum_{n=1}^N \ln \left[ \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(t_n - y_W(x_n))^2}{2\sigma^2}} \right] + \ln \left[ \frac{1}{(2\pi)^{1/M} |\Sigma|} e^{-\frac{W^T \Sigma^{-1} W}{2}} \right] \\
 &= \arg \max_W \sum_{n=1}^N -\frac{(t_n - y_W(x_n))^2}{2\sigma^2} - \frac{W^T \Sigma^{-1} W}{2} + \ln \left[ \frac{1}{\sqrt{2\pi}\sigma} \right] + \ln \left[ \frac{1}{(2\pi)^{1/M} |\Sigma|} \right]
 \end{aligned}$$

43

# Régression polynomiale

Maximum *a posteriori* (MAP)  
(cf. notes au tableau)

Cherche les meilleurs paramètres  $W$  en **maximisant la probabilité *a posteriori***

$$W = \arg \max_W \sum_{n=1}^N -\frac{(t_n - y_W(x_n))^2}{2\sigma^2} - \frac{W^T \Sigma^{-1} W}{2} + \ln \left[ \frac{1}{\sqrt{2\pi}\sigma} \right] + \ln \left[ \frac{1}{(2\pi)^{1/M} |\Sigma|} \right]$$

Constante par rapport à  $W$

De plus, comme on ne connaît généralement pas  $\Sigma$ , on suppose qu'elle est isotropique

$$\Sigma = \begin{pmatrix} \alpha^2 & 0 & \dots & 0 \\ 0 & \alpha^2 & & \vdots \\ \vdots & & \ddots & \\ 0 & 0 & \dots & \alpha^2 \end{pmatrix} = \alpha^2 I$$

44

# Régression polynomiale

Maximum *a posteriori* (MAP)  
(cf. notes au tableau)

Cherchez les meilleurs paramètres  $W$  en **maximisant la probabilité *a posteriori***

$$\begin{aligned} W &= \arg \max_W \sum_{n=1}^N -\frac{(t_n - y_W(x_n))^2}{2\sigma^2} - \frac{W^T \Sigma^{-1} W}{2} \\ &= \arg \max_W \sum_{n=1}^N -\frac{(t_n - y_W(x_n))^2}{2\sigma^2} - \frac{W^T W}{2\alpha^2} \\ &= \arg \min_W \sum_{n=1}^N \frac{(t_n - y_W(x_n))^2}{2} + \lambda \frac{W^T W}{2} \quad \text{où } \lambda = \frac{\sigma^2}{\alpha^2} \end{aligned}$$

45

# Régression polynomiale

Maximum *a posteriori* (MAP)

Cherchez les meilleurs paramètres  $W$  en **maximisant la probabilité *a posteriori***

$$W = \arg \min_W \sum_{n=1}^N \frac{(t_n - y_W(x_n))^2}{2} + \lambda \frac{W^T W}{2}$$

**NOTE**

Formule également connue sous le nom de  
« **régression de Ridge** »

Voir sklearn pour une implémentation simple

[scikit-learn.org/stable/modules/generated/sklearn.linear\\_model.Ridge.html](https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.Ridge.html)

46

# Théorie de l'information

47

47

# Théorie de l'information

- Les probabilités sont également utiles pour **quantifier l'information** présente dans des données  
exemple : quel est le nombre minimum de bits nécessaire pour encoder un message ?
- Cette question est intimement liée à la **probabilité d'observer ce message**  
plus le message est «surprenant» (improbable), plus on aura besoin de bits

48

48



# Théorie de l'information

Codage de Huffman :

- façon optimale d'encoder des symboles indépendants de façon binaire
- plus un symbole est «fréquent» (probable), plus son code sera court

'a'	'b'	'c'	'd'
0.4	0.05	0.2	0.35

49

# Théorie de l'information

Codage de Huffman :

- façon optimale d'encoder des symboles indépendants de façon binaire
- plus un symbole est «fréquent» (probable), plus son code sera court

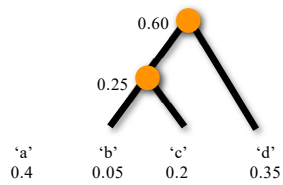
	0.25		
	└─┘		
'a'	'b'	'c'	'd'
0.4	0.05	0.2	0.35

50

# Théorie de l'information

Codage de Huffman :

- façon optimale d'encoder des symboles indépendants de façon binaire
- plus un symbole est «fréquent» (probable), plus son code sera court

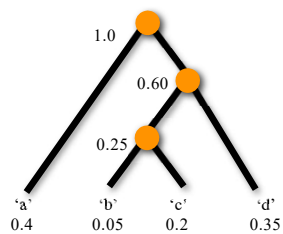


51

# Théorie de l'information

Codage de Huffman :

- façon optimale d'encoder des symboles indépendants de façon binaire
- plus un symbole est «fréquent» (probable), plus son code sera court

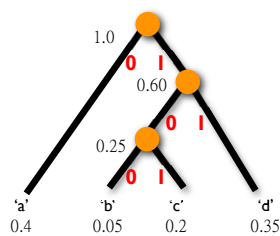


52

# Théorie de l'information

Codage de Huffman :

- façon optimale d'encoder des symboles indépendants de façon binaire
- plus un symbole est «fréquent» (probable), plus son code sera court



Symbole	Code	Prob
'a'	0	40%
'b'	100	5%
'c'	101	20%
'd'	11	35%

53

# Entropie

Symbole	Code	Prob
'a'	0	40%
'b'	100	5%
'c'	101	20%
'd'	11	35%

- Soit  $p(x)$  la probabilité d'observer le symbole  $x$

la taille moyenne du code d'un symbole est

$$0.4 \times 1 + 0.05 \times 3 + 0.2 \times 3 + 0.35 \times 2 = \mathbf{1.85 \text{ (bits)}}$$

- **Entropie :**

$$H[x] = -\sum_x p(x) \log_2 p(x) \approx \mathbf{1.739 \text{ (bits)}}$$

Claude Shannon a démontré qu'il est impossible de compresser l'information dans un plus petit code moyen **sans perte d'information**

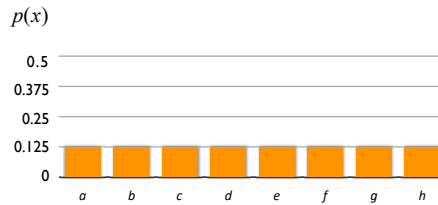
$-\log_2 p(x)$  est l'information contenue par  $x$

54

# Entropie

Plus  $p(x)$  est proche d'une **loi uniforme**, plus l'**entropie est élevée**

exemple :  $x \in \{a, b, c, d, e, f, g, h\}$



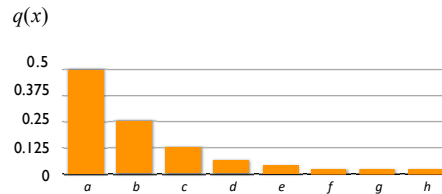
$$\begin{aligned}
 H[x] &= -\sum_x p(x) \log_2 p(x) \\
 &= -\left(\frac{1}{8} \log_2 \left(\frac{1}{8}\right)\right) - \left(\frac{1}{8} \log_2 \left(\frac{1}{8}\right)\right) - \left(\frac{1}{8} \log_2 \left(\frac{1}{8}\right)\right) - \left(\frac{1}{8} \log_2 \left(\frac{1}{8}\right)\right) - \left(\frac{1}{8} \log_2 \left(\frac{1}{8}\right)\right) - \left(\frac{1}{8} \log_2 \left(\frac{1}{8}\right)\right) - \left(\frac{1}{8} \log_2 \left(\frac{1}{8}\right)\right) - \left(\frac{1}{8} \log_2 \left(\frac{1}{8}\right)\right) \\
 &= -8 \left(\frac{1}{8} \log_2 \left(\frac{1}{8}\right)\right) \\
 &= 3 \text{ bits}
 \end{aligned}$$

55

# Entropie

Plus  $q(x)$  s'éloigne d'une **loi uniforme**, plus l'**entropie est faible**

exemple :  $x \in \{a, b, c, d, e, f, g, h\}$



$$\begin{aligned}
 H[x] &= -\sum_x q(x) \log_2 q(x) \\
 &= -\left(\frac{1}{2} \log_2 \left(\frac{1}{2}\right)\right) - \left(\frac{1}{4} \log_2 \left(\frac{1}{4}\right)\right) - \left(\frac{1}{8} \log_2 \left(\frac{1}{8}\right)\right) - \left(\frac{1}{16} \log_2 \left(\frac{1}{16}\right)\right) - \left(\frac{1}{32} \log_2 \left(\frac{1}{32}\right)\right) - 3 \left(\frac{1}{64} \log_2 \left(\frac{1}{64}\right)\right) \\
 &= 2.06 \text{ bits}
 \end{aligned}$$

56

# Entropie

L'entropie se généralise aux variables continues

$$H[x] = - \int_{-\infty}^{\infty} p(x) \log_2 p(x) dx$$

57

## Entropie relative et divergence de Kullback-Leibler

- Si on ne connaît pas  $p(x)$ , on va vouloir l'estimer
- Si  $q(x)$  est notre estimation, on définit la **divergence de Kullback-Leibler** (K-L) comme suit :

$$\begin{aligned} KL(p(x) \parallel q(x)) &= - \sum_x p(x) \log_2 q(x) - \left( - \sum_x p(x) \log_2 p(x) \right) \\ &= \sum_x p(x) \log_2 \frac{p(x)}{q(x)} \end{aligned}$$

➤ correspond au **nombre de bits additionnels** par rapport à ce qui serait optimal

58

## Entropie jointe

L'entropie est une fonction d'une loi de probabilité

- elle reflète l'**incertitude représentée par la loi**
- si  $p(x) = 1$  pour une seule valeur de  $x$ , l'entropie est 0

On peut généraliser l'entropie à **plusieurs variables**

$$H[x, y] = - \sum_x \sum_y p(x, y) \log_2 p(x, y)$$

59

## Entropie conditionnelle

L'entropie conditionnelle quantifie l'**information additionnelle** qu'apporte une **nouvelle observation  $y$**

$$H[x | y] = - \sum_x \sum_y p(x, y) \log_2 p(x | y)$$

On peut démontrer que

$$H[x, y] = H[y | x] + H[x]$$

60

## Information mutuelle

- Mesure à quel point deux variables sont indépendantes

$$\begin{aligned} I(x, y) &= KL(p(x, y) \parallel p(x)p(y)) \\ &= \sum_{x,y} p(x, y) \log_2 \frac{p(x, y)}{p(x)p(y)} \end{aligned}$$

- On appelle cette mesure l'**information mutuelle**