

UNIVERSITÉ DE MONTRÉAL

# Rapport : Devoir1

Pierre Gérard  
Mathieu Bouchard

**IFT3395-6390 Fondements de l'apprentissage machine**  
Pascal Vincent, Alexandre de Brébisson et César Laurent

Année académique 2015 - 2016

# Table des matières

<b>1</b>	<b>Partie Théorique : Estimation de densité</b>	<b>2</b>
1.1	Densité avec des fenêtres de Parzen a noyau Gaussien isotropique . . . . .	2
1.2	Densité paramétrique avec Gaussienne isotropique . . . . .	2
1.3	Densité paramétrique avec Gaussienne diagonale . . . . .	2
1.4	Choix du modèle . . . . .	3
<b>2</b>	<b>Partie Théorique : Classifieurs de Bayes</b>	<b>3</b>
2.1	Classifieur de Bayes obtenu avec des densités paramétriques Gaussiennes diagonales . . .	3
2.2	Classifieur de Bayes obtenu avec des fenêtres de Parzen à noyau Gaussien isotropique . .	3
2.3	Classifieur de Bayes obtenu avec des fenêtres de Parzen à noyau Gaussien isotropique . .	4

# 1 Partie Théorique : Estimation de densité

## 1.1 Densité avec des fenêtres de Parzen a noyau Gaussien isotropique

a)  $\mathcal{N}_{x^{(i)}, \sigma^2} = \frac{1}{(2\pi)^{d/2} \sigma^d} e^{-\frac{1}{2} \frac{d(x^{(i)}, x)^2}{\sigma^2}}$

b) Si on choisit de prendre tous les points du voisinages  $p(x) = \frac{1}{n} \sum_{i=1}^n \mathcal{N}_{x^{(i)}, \sigma^2}$

c) Hyper-paramètre :  $\sigma \in \mathbb{R}$

Paramètre : On mémorise l'ensemble des données d'entraînement.  $(n \times \mathbb{R}^d)$

On considère cette méthode comme non-paramétrique car le nombre de paramètre varie avec la taille de l'ensemble de données.

d) 3 propriété a prouvé : integral, integral=1, toujours positive (ou presque)

distribuer  
pour virer  
la somme ?

Le petit d

## 1.2 Densité paramétrique avec Gaussienne isotropique

a)  $p(x) = \frac{1}{(2\pi)^{d/2} \sigma^d} e^{-\frac{\|x-\mu\|^2}{2\sigma^2}}$

$\mu$  est le est le vecteur des moyennes, de dimension d.  $\Sigma$  est la matrice de covariance, de dimension 1.

b) Dans le cas de l'estimation par fenêtre de Parzen, il nous faut effectuer une moyenne sur n densités Gaussiennes et, pour chacune d'elles, le point  $x^{(i)}$  sert de moyenne  $\mathcal{N}_{x^{(i)}, \sigma^2}$  avec  $\sigma$  qui est un hyperparamètre. Dans le cas présent, on utilise une seule densité Gaussienne pour laquelle il faudra apprendre les valeurs de moyenne  $\mu$  et d'écart-type  $\sigma$ .

c)

Paramètres :

–  $\mu$  est le est le vecteur des moyennes, de dimension d.

–  $\Sigma$  est la matrice de covariance, de dimension 1.

Pas d'hyper-paramètre.

## 1.3 Densité paramétrique avec Gaussienne diagonale

a)  $p(x) = \frac{1}{(2\pi)^{d/2} \sqrt{|\Sigma|}} e^{-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)}$

–  $\mu$  est le est le vecteur des moyennes, de dimension d.

–  $\Sigma$  est la matrice de covariance, de dimension d.

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 & 0 & \dots & 0 \\ 0 & \sigma_2 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & \sigma_n \end{bmatrix}$$

b) Indépendance des composantes d'un vecteur aléatoire suivant une distribution Gaussienne diagonale :

Soit  $x \in \mathbb{R}^d$  un vecteur aléatoire qui suit une distribution Gaussienne diagonale  $\mathcal{N}_{\mu, \sigma(x)}$

Posons  $x - \mu = k$  et  $\frac{1}{(2\pi)^{d/2} \sqrt{|\Sigma|}} = s$

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 & 0 & \dots & 0 \\ 0 & \sigma_2 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & \sigma_n \end{bmatrix} \quad \Sigma^{-1} = \begin{bmatrix} \frac{1}{\sigma_1} & 0 & 0 & \dots & 0 \\ 0 & \frac{1}{\sigma_2} & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & \frac{1}{\sigma_n} \end{bmatrix}$$

$$p(x) = \mathcal{N}_{\mu, \Sigma}(x) = \frac{1}{(2\pi)^{d/2} \sqrt{|\Sigma|}} e^{-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)} = s \times e^{-\frac{1}{2} k^T \Sigma^{-1} k}$$

$$p(x) = s \times e^{-\frac{1}{2} \left( \frac{k_1}{\sigma_1^2} \frac{k_2}{\sigma_2^2} \dots \frac{k_d}{\sigma_d^2} \right) k} = s \times e^{-\frac{1}{2} \sum_{i=1}^d \frac{k_i^2}{\sigma_i^2}}$$

$p(x) = s \prod_{i=1}^d e^{-\frac{1}{2} \frac{k_i^2}{\sigma_i^2}} = s \prod_{i=1}^d e^{-\frac{1}{2} (x^i - \mu)^T \Sigma^{-1} (x^i - \mu)}$  où  $x^i$  est un vecteur de taille d avec la  $i^{eme}$  valeur de x en position i et des 0 partout ailleurs.

$$p(x) = \frac{1}{(2\pi)^{d/2} \sqrt{|\Sigma|}} \prod_{i=1}^d e^{-\frac{1}{2}(x^i - \mu)^T \Sigma^{-1} (x^i - \mu)}$$

$$p(x) = \prod_{i=1}^d p(x^i)$$

La densité estimée pour à partir d'un vecteur aléatoire  $x \in R$  qui suit une distribution gaussienne diagonale est égale au produit des densités estimées des composantes de ce vecteur aléatoire. On peut en conclure que ces composantes sont des variables aléatoires indépendantes.

c)

gd c

d)

gd d

## 1.4 Choix du modèle

La première tâche à effectuer pour choisir le modèle le plus approprié est de sélectionner un ensemble de tous les modèles que nous pourrions effectuer. Ici, il s'agit des approches présentées aux sections 1.1, 1.2 et 1.3. Ensuite, pour chacun des modèles sélectionnés, on procède ainsi :

1. Soit  $A$ , le modèle courant.
2. Si le modèle contient des hyper-paramètres, répéter la procédure pour toutes les valeurs acceptables que ces hyper-paramètres peuvent prendre. Appelons cette configuration  $\lambda$ . Rappelons que l'approche des fenêtres de Parzen à noyau isotropique contient un hyper- paramètre (l'écart-type  $\sigma$ ) tandis que les approches 2 et 3 n'en n'ont aucun ( $\lambda = \emptyset$ ).
3. Mélanger les éléments de l'ensemble  $D$  s'ils étaient ordonnés.
4. Séparer l'ensemble d'entraînement  $D$  en deux sous-ensembles  $D_{train}$  et  $D_{valid}$ . L'ensemble d'entraînement devrait être plus grand que l'ensemble de validation.  $x \in \mathbb{R}$  qui suit une distribution Gaussienne
5. Entraîner  $A$  avec la configuration d'hyper-paramètres choisis sur l'ensemble d'entraînement. On obtient alors :  $\hat{f}(A_\lambda) = A_\lambda(D_{train})$
6. Evaluer le résultat à l'aide de l'ensemble de validation.  $e_{(A_\lambda)} = \hat{R}(\hat{f}_{(A_\lambda)}, D_{valid})$

Une fois que ce travail a été effectuer pour tous les algorithmes et, pour chacun, sur chaque configurations d'hyper-paramètres  $\lambda$  jugés admissibles, on regarde les valeurs de  $e_{(A_\lambda)}$  obtenues. Le modèle  $A_\lambda$  qui retourne une erreur de validation minimale  $e_{A_\lambda}$  sera alors sélectionner. Nous aurions également pu séparer  $D$  en 3 sous-ensembles plutôt que 2 et ainsi ajouter un ensemble de test. Une fois l'algorithme optimal trouvé à l'aide de l'ensemble d'entraînement et de l'ensemble de validation, nous aurions pu calculer un estimé non-biaisé de la performance de généralisation de l'algorithme choisit en calculant l'erreur empirique de validation sur l'ensemble de test ; soit un ensemble de données dont le valeurs n'ont jamais été utilisées pour sélectionner l'algorithme dans les étapes précédantes. Cette étape nous permet de nous assurer (si nous avons assez de données pour pouvoir effectuer un test concluant) que l'algorithme sélectionné est capable de généraliser sur de nouvelles données et qu'il n'est pas juste efficace sur les données de  $D_{train}$  et  $D_{valid}$ .

## 2 Partie Théorique : Classifieurs de Bayes

### 2.1 Classifieur de Bayes obtenu avec des densités paramétriques Gaussiennes diagonales

a)

b)

c)

### 2.2 Classifieur de Bayes obtenu avec des fenêtres de Parzen à noyau Gaussien isotropique

a)

b)

### 2.3 Classifieur de Bayes obtenu avec des fenêtres de Parzen à noyau Gaussien isotropique

- a)
- b)