

Massih-Reza Amini • Renaud Blanch • Marianne Clausel  
Jean-Baptiste Durand • Eric Gaussier • Jérôme Malick  
Christophe Picard • Vivien Quéma • Georges Quénot

# DATA SCIENCE

Cours et exercices



EYROLLES

# Table des matières

Table des figures .....	IX
Liste des algorithmes .....	XIII
CHAPITRE 1	
Introduction .....	1
CHAPITRE 2	
Prétraitement des données .....	5
2.1 Prétraitement des données textuelles .. . . . .	5
2.1.1 Segmentation . . . . .	6
2.1.2 Normalisation et filtrage . . . . .	7
2.1.3 Filtrage . . . . .	10
2.2 Prétraitement des données image ou vidéo . . . . .	11
2.2.1 Représentations globales . . . . .	12
2.2.2 Représentations locales . . . . .	16
2.2.3 Agrégation de représentations locales . . . . .	17
2.2.4 Représentations apprises . . . . .	19
2.3 Exercices .. . . . .	19
CHAPITRE 3	
Gestion de données large-échelle et systèmes distribués .....	21
3.1 Les limites des systèmes traditionnels de gestion de données .. . . . .	21

3.1.1	Les besoins liés au traitement de grandes masses de données . . . . .	22
3.1.2	Limites des architectures incrémentales . . . . .	23
<b>3.2</b>	<b>L'architecture Lambda pour le traitement de grandes masses de données . . . . .</b>	<b>24</b>
<b>3.3</b>	<b>La couche batch : traitement de données par lots . . . . .</b>	<b>27</b>
3.3.1	Caractéristiques du jeu de données principal . . . . .	28
3.3.2	Stockage du jeu de données principal . . . . .	29
3.3.3	Traitement de données par lots . . . . .	30
<b>3.4</b>	<b>La couche service : stockage et requêtes sur les vues batch . . . . .</b>	<b>33</b>
3.4.1	Remarque préliminaire sur le stockage des vues batch . . . . .	33
3.4.2	Stockage des vues batch . . . . .	35
<b>3.5</b>	<b>La couche vitesse : traitement de flux de données, stockage et requêtes sur les vues temps réel . . . . .</b>	<b>36</b>
3.5.1	Traitemet de flux de données . . . . .	37
3.5.2	Stockage des vues temps réel . . . . .	38
<b>3.6</b>	<b>Conclusion . . . . .</b>	<b>38</b>

## CHAPITRE 4

### Calcul haute performance . . . . . **41**

<b>4.1</b>	<b>Introduction . . . . .</b>	<b>41</b>
4.1.1	Motivations . . . . .	42
4.1.2	Hiérarchies du parallélisme . . . . .	44
4.1.3	Classification des plateformes . . . . .	47
4.1.4	Coûts de communication . . . . .	48
<b>4.2</b>	<b>Principes de conception des algorithmes . . . . .</b>	<b>50</b>
4.2.1	Techniques de décomposition . . . . .	51
4.2.2	Caractéristiques des tâches et des interactions . . . . .	53
4.2.3	Équilibrage des ressources . . . . .	54
4.2.4	Modèles d'algorithmes parallèles . . . . .	56

<b>4.3 Modèles analytiques . . . . .</b>	57
4.3.1 Métriques de performances . . . . .	58
4.3.2 Passage à l'échelle des systèmes parallèles . . . . .	59
4.3.3 Effet de la granularité . . . . .	61
4.3.4 Notion d'iso-efficacité . . . . .	62
<b>4.4 Conclusion . . . . .</b>	62
 CHAPITRE 5	
<b>Optimisation pour l'analyse de données . . . . .</b>	65
<b>5.1 Apprentissage et optimisation . . . . .</b>	66
<b>5.2 Introduction à l'optimisation . . . . .</b>	72
5.2.1 Problèmes d'optimisation . . . . .	72
5.2.2 Analyse convexe pour impatients . . . . .	74
5.2.3 Algorithmes d'optimisation . . . . .	77
<b>5.3 Algorithmes en science des données . . . . .</b>	82
5.3.1 Algorithmes incrémentaux . . . . .	83
5.3.2 Algorithmes distribués . . . . .	87
5.3.3 Au-delà de ce chapitre . . . . .	91
<b>5.4 Exercices . . . . .</b>	91
 CHAPITRE 6	
<b>Décomposition matricielle/tensorielle . . . . .</b>	95
<b>6.1 Motivations . . . . .</b>	95
<b>6.2 La SVD . . . . .</b>	96
6.2.1 Quelques rappels d'algèbre linéaire . . . . .	96
6.2.2 Approximation de rang faible . . . . .	97
6.2.3 SVD et analyse en composantes principales . . . . .	99
6.2.4 Algorithme pour déterminer la SVD d'une matrice . . . . .	100
<b>6.3 Décomposition en matrices non négatives . . . . .</b>	104
6.3.1 Algorithme de Seung et Lee . . . . .	105
6.3.2 Algorithme des moindres carrés alternés . . . . .	107
6.3.3 Comparaison de la SVD et de la NMF . . . . .	107
<b>6.4 Décomposition tensorielle . . . . .</b>	108
6.4.1 Décomposition canonique polyadique . . . . .	109

<b>6.5 Conclusion</b>	111
<b>6.6 Exercices</b>	111

## CHAPITRE 7

<b>Modèles génératifs</b>	115
---------------------------	-----

<b>7.1 Motivations</b>	115
7.1.1 Modèles graphiques	117
7.1.2 Modèles à variables latentes	117
<b>7.2 Introduction à la statistique bayésienne</b>	
7.2.1 Généralités	120
7.2.2 Algorithmes génériques pour l'inférence bayésienne	123
<b>7.3 Inférence dans les modèles à variables latentes</b>	127
7.3.1 Modèles probabilistes graphiques	127
7.3.2 Mélanges	129
7.3.3 Analyse en composantes principales probabiliste	132
7.3.4 Chaînes de Markov cachées	134
7.3.5 Modèles hiérarchiques à variables latentes	136
<b>7.4 Références</b>	138
<b>7.5 Exercices</b>	140

## CHAPITRE 8

<b>Modèles discriminants</b>	145
------------------------------	-----

<b>8.1 Approches supervisées</b>	146
8.1.1 Modèles binaires	147
8.1.2 Modèles multi-classes	161
<b>8.2 Approches semi-supervisées</b>	164
8.2.1 Modèles graphiques	165
8.2.2 Modèles de mélange	171
8.2.3 Modèles discriminants	171
<b>8.3 Exercices</b>	172

<b>CHAPITRE 9</b>	
<b>Deep learning .....</b>	<b>177</b>
<b>9.1 Neurone formel .....</b>	<b>178</b>
<b>9.2 Réseaux simples .....</b>	<b>179</b>
9.2.1 Perceptron .....	180
9.2.2 ADALINE .....	183
9.2.3 Perceptrons multicouches (PMC) .....	185
<b>9.3 Réseaux à propagation avant .....</b>	<b>186</b>
9.3.1 Composition de fonctions .....	186
9.3.2 Fonction-objectif et descente de gradient stochastique par mini-lots .....	187
9.3.3 Calcul des gradients par rétro-propagation de l'erreur .....	188
9.3.4 Architecture modulaire .....	190
9.3.5 Réseaux quelconques .....	195
9.3.6 Différentiation automatique .....	196
<b>9.4 Réseaux convolutifs .....</b>	<b>197</b>
9.4.1 Couche de convolution .....	197
9.4.2 Changements de résolution .....	199
9.4.3 Passage à des couches complètement connectées .....	200
9.4.4 Un exemple : AlexNet .....	200
<b>9.5 Optimisations supplémentaires .....</b>	<b>202</b>
9.5.1 Traitement par mini-lots .....	202
9.5.2 Moment .....	202
9.5.3 Fonctions d'activation .....	203
9.5.4 Dropout .....	204
9.5.5 Normalisation de lots .....	204
9.5.6 Augmentation de données .....	205
<b>9.6 Réseaux pour la catégorisation d'images .....</b>	<b>206</b>
<b>9.7 Exercices .....</b>	<b>207</b>
<b>CHAPITRE 10</b>	
<b>Visualisation interactive d'information.....</b>	<b>211</b>
<b>10.1 Introduction .....</b>	<b>211</b>
<b>10.2 Des données au graphique .....</b>	<b>214</b>
10.2.1 Les données .....	214

10.2.2 L'image . . . . .	216
10.2.3 Encodage visuel . . . . .	224
<b>10.3 Encodages avancés . . . . .</b>	<b>227</b>
10.3.1 Utilisation multiple des variables graphiques . . . . .	227
10.3.2 Encodage des liens entre individus . . . . .	230
<b>10.4 Pour aller plus loin . . . . .</b>	<b>234</b>
<b>10.5 Exercices . . . . .</b>	<b>235</b>
<b>Bibliographie . . . . .</b>	<b>239</b>
<b>Index . . . . .</b>	<b>251</b>