



Table des matières

R	emer	rciements	vii
A	vant-	Propos	ix
Ι	In	troduction au modèle linéaire	1
1	La	régression linéaire simple	3
	1.1	Introduction	3
		1.1.1 Un exemple : la pollution de l'air	3
		1.1.2 Un second exemple : la hauteur des arbres	5
	1.2	Modélisation mathématique	7
		1.2.1 Choix du critère de qualité et distance à la droite	7
		1.2.2 Choix des fonctions à utiliser	9
	1.3	Modélisation statistique	10
	1.4	Estimateurs des moindres carrés	11
		1.4.1 Calcul des estimateurs de β_j , quelques propriétés	11
		1.4.2 Résidus et variance résiduelle	15
		1.4.3 Prévision	15
	1.5	Interprétations géométriques	16
		1.5.1 Représentation des individus	16
		1.5.2 Représentation des variables	17
	1.6	Inférence statistique	19
	1.7	Exemples	22
	1.8	Exercices	28
2	La	régression linéaire multiple	31
	2.1	Introduction	31
	2.2	Modélisation	32
	2.3	Estimateurs des moindres carrés	34
	-	2.3.1 Calcul de $\hat{\beta}$	35
		2.3.2 Interprétation	37
		2.3.3 Quelques propriétés statistiques	38
		2.3.4 Résidus et variance résiduelle	40







"regression" — 2024/12/31 — 17:10 — page xii — #6



xii Régression avec Python

		2.3.5 Prévision	41	
	2.4	Interprétation géométrique	42	
	2.5	Exemples	43	
	2.6	Exercices	46	
3	Val	idation du modèle	51	
	3.1	Analyse des résidus	52	
		3.1.1 Les différents résidus	52	
		3.1.2 Ajustement individuel au modèle, valeur aberrante	53	
		3.1.3 Analyse de la normalité	54	
		3.1.4 Analyse de l'homoscédasticité	55	
		3.1.5 Analyse de la structure des résidus	56	
	3.2	Analyse de la matrice de projection	59	
	3.3	Autres mesures diagnostiques	60	
	3.4	Effet d'une variable explicative	63	
		3.4.1 Ajustement au modèle	63	
		3.4.2 Régression partielle : impact d'une variable	64	
		3.4.3 Résidus partiels et résidus partiels augmentés	65	
	3.5	Exemple: la concentration en ozone	67	
	3.6	Exercices	71	
4	4 Extensions : non-inversibilité et (ou) erreurs corrélées			
	4.1	Régression ridge	73	
		4.1.1 Une solution historique	74	
		4.1.2 Minimisation des MCO pénalisés	75	
		4.1.3 Equivalence avec une contrainte sur la norme des coefficients		
	4.0.1	4.1.4 Propriétés statistiques de l'estimateur ridge β_{ridge}	76	
	4.21	Erreurs corrélées : moindres carrés généralisés	78	
		4.2.1 Erreurs hétéroscédastiques	79	
		4.2.2 Estimateur des moindres carrés généralisés	82	
	4.0	4.2.3 Matrice Ω inconnue	84	
	4.3	Exercices	85	
5	Rég	ression polynomiale et régression spline	89	
	5.1	Régression polynomiale	89	
	5.2	Régression spline	93	
		5.2.1 Introduction	93	
		5.2.2 Spline de régression	94	
	5.3	Spline de lissage	98	
	5.4	Exercices	101	
II	Ir	nférence	103	
6		rence dans le modèle gaussien	105	
	6.1	Estimateurs du maximum de vraisemblance	105	







"regression" — 2024/12/31 — 17:10 — page xiii — #7



		Т	able des matières	xiii
	6.2	Nouvelles propriétés statistiques		
	6.3	Intervalles et régions de confiance		
	6.4	Prévision		
	6.5	Les tests d'hypothèses		
		6.5.1 Introduction		
	0.0	6.5.2 Test entre modèles emboîtés		
	6.6	Applications		
	6.7	Exercices		
	6.8	Notes		
		6.8.1 Intervalle de confiance : bootstrap		
		6.8.2 Test de Fisher pour une hypothèse linéair		
		6.8.3 Propriétés asymptotiques		127
7	Var	riables qualitatives : ANCOVA et ANOVA		131
	7.1	Introduction		
	7.2	Analyse de la covariance		
		7.2.1 Introduction : exemple des eucalyptus .		
		7.2.2 Modélisation du problème		
		7.2.3 Hypothèse gaussienne		
		7.2.4 Exemple: la concentration en ozone		
		7.2.5 Exemple: la hauteur des eucalyptus		142
	7.3	Analyse de la variance à 1 facteur		
		7.3.1 Introduction		144
		7.3.2 Modélisation du problème		145
		7.3.3 Interprétation des contraintes		147
		7.3.4 Estimation des paramètres		147
		7.3.5 Hypothèse gaussienne et test d'influence	du facteur	148
		7.3.6 Exemple: la concentration en ozone		150
		7.3.7 Une décomposition directe de la variance		154
	7.4	Analyse de la variance à 2 facteurs		155
		7.4.1 Introduction		155
		7.4.2 Modélisation du problème		156
		7.4.3 Estimation des paramètres		158
		7.4.4 Analyse graphique de l'interaction		159
		7.4.5 Hypothèse gaussienne et test de l'interact	ion	160
		7.4.6 Exemple: la concentration en ozone		163
	7.5	Exercices		164
	7.6	Note : identifiabilité et contrastes		167
п	I 1	Réduction de dimension		171
8		pix de variables		173
	8.1	Introduction		
	8.2	Notations		175











xiv Régression avec Pyth	non
--------------------------	-----

	0.9		176
	8.3		176
			176
			178
			179
			181
	8.4	Critères classiques de choix de modèles	183
			184
		8.4.2 Le \mathbb{R}^2	185
		8.4.3 Le \mathbb{R}^2 ajusté	186
		8.4.4 Le C_p de Mallows	187
			189
			191
	8.5		193
			193
			193
	8.6		195
	0.0	•	195
			196
	8.7	-	197
	8.8		199
	0.0	1,000 . Op 00 blads de selection	100
	-	The state of the s	വ
9	Rég	ularisation des moindres carrés : Ridge, Lasso et elastic-net 2	
9	9.1	Introduction	203
9	_	Introduction	
9	9.1	Introduction	203
9	9.1 9.2	Introduction	203 206
9	9.1 9.2 9.3	Introduction	203 206 207
9	9.1 9.2 9.3	Introduction	203 206 207 207
9	9.1 9.2 9.3		203 206 207 207 211
9	9.1 9.2 9.3	$ \begin{array}{cccccccccccccccccccccccccccccccccccc$	203 206 207 207 211 212
9	9.1 9.2 9.3 9.4	$ \begin{array}{cccccccccccccccccccccccccccccccccccc$	203 206 207 207 211 212
9	9.1 9.2 9.3 9.4	Introduction 2 Problème du centrage-réduction des variables 2 Ridge, lasso et elastic net 2 Propriétés des régressions Ridge et lasso 2 9.4.1 Interprétation géométrique 2 9.4.2 Simplification quand les X sont orthogonaux 2 9.4.3 Choix de λ par validation croisée 2 Régressions avec le package scikitlearn 2 9.5.1 Estimation des paramètres 2	203 206 207 207 211 212 218
9	9.1 9.2 9.3 9.4	Introduction 2 Problème du centrage-réduction des variables 2 Ridge, lasso et elastic net 2 Propriétés des régressions Ridge et lasso 2 9.4.1 Interprétation géométrique 2 9.4.2 Simplification quand les X sont orthogonaux 2 9.4.3 Choix de λ par validation croisée 2 Régressions avec le package scikitlearn 2 9.5.1 Estimation des paramètres 2 9.5.2 Chemin de régularisation 2	203 206 207 207 221 221 221 218 218
9	9.1 9.2 9.3 9.4	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	203 206 207 207 211 212 218 218 218 2218
9	9.1 9.2 9.3 9.4	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	203 206 207 207 211 212 218 218 218 218 222 222 222 222
9	9.1 9.2 9.3 9.4 9.5	Introduction2Problème du centrage-réduction des variables2Ridge, lasso et elastic net2Propriétés des régressions Ridge et lasso29.4.1 Interprétation géométrique29.4.2 Simplification quand les X sont orthogonaux29.4.3 Choix de λ par validation croisée2Régressions avec le package scikitlearn29.5.1 Estimation des paramètres29.5.2 Chemin de régularisation29.5.3 Choix du paramètre de régularisation α 29.5.4 Mise en pratique2Intégration de variables qualitatives2	203 206 207 207 207 211 212 218 218 218 222 222 222 222 222
9	9.1 9.2 9.3 9.4 9.5	Introduction2Problème du centrage-réduction des variables2Ridge, lasso et elastic net2Propriétés des régressions Ridge et lasso29.4.1 Interprétation géométrique29.4.2 Simplification quand les X sont orthogonaux29.4.3 Choix de λ par validation croisée2Régressions avec le package scikitlearn29.5.1 Estimation des paramètres29.5.2 Chemin de régularisation29.5.3 Choix du paramètre de régularisation α 29.5.4 Mise en pratique2Intégration de variables qualitatives2Exercices2	203 206 207 207 211 212 218 218 218 222 222 222 222 222
9	9.1 9.2 9.3 9.4 9.5	Introduction2Problème du centrage-réduction des variables2Ridge, lasso et elastic net2Propriétés des régressions Ridge et lasso29.4.1 Interprétation géométrique29.4.2 Simplification quand les X sont orthogonaux29.4.3 Choix de λ par validation croisée2Régressions avec le package scikitlearn29.5.1 Estimation des paramètres29.5.2 Chemin de régularisation29.5.3 Choix du paramètre de régularisation α 29.5.4 Mise en pratique2Intégration de variables qualitatives2	203 206 207 207 211 212 218 218 218 222 222 222 222 222
	9.1 9.2 9.3 9.4 9.5 9.6 9.7 9.8 Rég	Introduction 2 Problème du centrage-réduction des variables 2 Ridge, lasso et elastic net 2 Propriétés des régressions Ridge et lasso 2 9.4.1 Interprétation géométrique 2 9.4.2 Simplification quand les X sont orthogonaux 2 9.4.3 Choix de λ par validation croisée 2 Régressions avec le package scikitlearn 2 9.5.1 Estimation des paramètres 2 9.5.2 Chemin de régularisation 2 9.5.3 Choix du paramètre de régularisation α 2 9.5.4 Mise en pratique 2 Intégration de variables qualitatives 2 Exercices 2 Note : lars et lasso 2 ression sur composantes : PCR et PLS 2	203 206 207 207 211 212 218 218 218 222 222 222 222 222
	9.1 9.2 9.3 9.4 9.5 9.6 9.7 9.8 Rég	Introduction 2 Problème du centrage-réduction des variables 2 Ridge, lasso et elastic net 2 Propriétés des régressions Ridge et lasso 2 9.4.1 Interprétation géométrique 2 9.4.2 Simplification quand les X sont orthogonaux 2 9.4.3 Choix de λ par validation croisée 2 Régressions avec le package scikitlearn 2 9.5.1 Estimation des paramètres 2 9.5.2 Chemin de régularisation 2 9.5.3 Choix du paramètre de régularisation α 2 9.5.4 Mise en pratique 2 Intégration de variables qualitatives 2 Exercices 2 Note : lars et lasso 2 ression sur composantes : PCR et PLS 2 Régression sur composantes principales (PCR) 2	203 206 2206 2207 2213 2218 2218 2218 2218 2228 2222 2222
	9.1 9.2 9.3 9.4 9.5 9.6 9.7 9.8 Rég	Introduction 2 Problème du centrage-réduction des variables 2 Ridge, lasso et elastic net 2 Propriétés des régressions Ridge et lasso 2 9.4.1 Interprétation géométrique 2 9.4.2 Simplification quand les X sont orthogonaux 2 9.4.3 Choix de λ par validation croisée 2 Régressions avec le package scikitlearn 2 9.5.1 Estimation des paramètres 2 9.5.2 Chemin de régularisation 2 9.5.3 Choix du paramètre de régularisation α 2 9.5.4 Mise en pratique 2 Intégration de variables qualitatives 2 Exercices 2 Note : lars et lasso 2 ression sur composantes : PCR et PLS 2 Régression sur composantes principales (PCR) 2 10.1.1 Changement de base 2	203 206 207 207 211 212 218 218 218 222 222 222 222 222
	9.1 9.2 9.3 9.4 9.5 9.6 9.7 9.8 Rég	Introduction 2 Problème du centrage-réduction des variables 2 Ridge, lasso et elastic net 2 Propriétés des régressions Ridge et lasso 2 9.4.1 Interprétation géométrique 2 9.4.2 Simplification quand les X sont orthogonaux 2 9.4.3 Choix de λ par validation croisée 2 Régressions avec le package scikitlearn 2 9.5.1 Estimation des paramètres 2 9.5.2 Chemin de régularisation 2 9.5.3 Choix du paramètre de régularisation α 2 9.5.4 Mise en pratique 2 Intégration de variables qualitatives 2 Exercices 2 Note : lars et lasso 2 ression sur composantes : PCR et PLS 2 Régression sur composantes principales (PCR) 2 10.1.1 Changement de base 2	203 206 2206 2207 2211 2212 2218 2218 2218 2218 2222 2222 2222 2222 2222 2222 2222 2222
	9.1 9.2 9.3 9.4 9.5 9.6 9.7 9.8 Rég	Introduction 2 Problème du centrage-réduction des variables 2 Ridge, lasso et elastic net 2 Propriétés des régressions Ridge et lasso 2 9.4.1 Interprétation géométrique 2 9.4.2 Simplification quand les X sont orthogonaux 2 9.4.3 Choix de λ par validation croisée 2 Régressions avec le package scikitlearn 2 9.5.1 Estimation des paramètres 2 9.5.2 Chemin de régularisation 2 9.5.3 Choix du paramètre de régularisation α 2 9.5.4 Mise en pratique 2 Intégration de variables qualitatives 2 Exercices 2 Note : lars et lasso 2 ression sur composantes : PCR et PLS 2 Régression sur composantes principales (PCR) 2 10.1.1 Changement de base 2 10.1.2 Estimateurs des MCO 2	203 206 207 207 211 212 218 218 218 2218 222 222 222 22
	9.1 9.2 9.3 9.4 9.5 9.6 9.7 9.8 Rég	Introduction 2 Problème du centrage-réduction des variables 2 Ridge, lasso et elastic net 2 Propriétés des régressions Ridge et lasso 2 9.4.1 Interprétation géométrique 2 9.4.2 Simplification quand les X sont orthogonaux 2 9.4.3 Choix de λ par validation croisée 2 Régressions avec le package scikitlearn 2 9.5.1 Estimation des paramètres 2 9.5.2 Chemin de régularisation 2 9.5.3 Choix du paramètre de régularisation α 2 9.5.4 Mise en pratique 2 Intégration de variables qualitatives 2 Exercices 2 Note : lars et lasso 2 ression sur composantes : PCR et PLS 2 Régression sur composantes principales (PCR) 2 10.1.1 Changement de base 2 10.1.2 Estimateurs des MCO 2 10.1.3 Choix de composantes/variables 2	203 206 207 207 211 212 218 218 218 2218 222 222 222 22
	9.1 9.2 9.3 9.4 9.5 9.6 9.7 9.8 Rég	Introduction 2 Problème du centrage-réduction des variables 2 Ridge, lasso et elastic net 2 Propriétés des régressions Ridge et lasso 2 9.4.1 Interprétation géométrique 2 9.4.2 Simplification quand les X sont orthogonaux 2 9.4.3 Choix de λ par validation croisée 2 Régressions avec le package scikitlearn 2 9.5.1 Estimation des paramètres 2 9.5.2 Chemin de régularisation 2 9.5.3 Choix du paramètre de régularisation α 2 9.5.4 Mise en pratique 2 Intégration de variables qualitatives 2 Exercices 2 Note : lars et lasso 2 ression sur composantes : PCR et PLS 2 Régression sur composantes principales (PCR) 2 10.1.1 Changement de base 2 10.1.2 Estimateurs des MCO 2 10.1.3 Choix de composantes/variables 2 10.1.4 Retour aux données d'origine 2 2 2	203 206 207 2207 2212 2218 2218 2218 2218 2222 2222 222





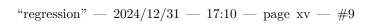






	Table des matières	X
10	2. Dágmagaign agus maindrea cannág mantialg (DIC)	240
10.	2 Régression aux moindres carrés partiels (PLS)	$\frac{240}{242}$
	10.2.1 Algorithmes PLS	242
	10.2.2 Choix de composantes/variables	243
	10.2.4 La régression PLS en pratique	$\frac{243}{245}$
10	3 Exercices	246
	4 Notes	248
10.	10.4.1 ACP et changement de base	248
	10.4.1 ACI et changement de base $\cdot \cdot \cdot$	249
11 Co	mparaison des différentes méthodes, étude de cas réels	253
11.	1 Erreur de prévision et validation croisée	253
11.	2 Analyse de l'ozone	257
	11.2.1 Préliminaires	257
	11.2.2 Méthodes et comparaison	258
11.	3 Transformation des variables : feature engineering	261
	11.3.1 Modèle de prévision avec interactions	261
	11.3.2 Modèle de prévision avec des polynômes	262
	11.3.3 Modèle de prévision avec des splines	262
	11.3.4 Modèle de prévision avec des splines et de l'interaction	263
	11.3.5 Conclusion	264
IV	Le modèle linéaire généralisé	265
12 Ré	gression logistique	267
	1 Présentation du modèle	267
12.	12.1.1 Exemple introductif	267
	12.1.2 Modélisation statistique	268
	12.1.3 Variables explicatives qualitatives, interactions	271
12.	2 Estimation	273
	12.2.1 La vraisemblance	274
	12.2.2 Calcul des estimateurs : l'algorithme IRLS	275
	12.2.3 Propriétés asymptotiques de l'EMV	277
12.	3 Intervalles de confiance et tests	278
	12.3.1 IC et tests sur les paramètres du modèle	278
	12.3.2 Test sur un sous-ensemble de paramètres	280
	12.3.3 Prévision	282
12.	4 Adéquation du modèle	284
	12.4.1 Le modèle saturé	285
	12.4.2 Tests d'adéquation de la déviance et de Pearson	287
	12.4.3 Analyse des résidus	289
12	5 Choix de variables	293
12.	12.5.1 Tests entre modèles emboîtés	$\frac{233}{294}$
	12.5.2 Procédures automatiques	294
		~ -











xvi	Régression	avec	Python

13 Régression de Poisson 30 13.1 Le modèle linéaire généralisé (GLM) 36 13.2 Exemple : modélisation du nombre de visites 36 13.3 Régression Log-linéaire 36 13.3.1 Le modèle 36 13.3.2 Estimation 31 13.3.3 Tests et intervalles de confiance 31 13.3.4 Choix de variables 31 13.4 Exercices 31 14 Régularisation de la vraisemblance 32 14.1 Régressions ridge, lasso et elastic-net 32 14.2 Choix du paramètre de régularisation λ 32 14.3 Group-lasso 32 14.4 Exercices 33 15 Comparaison en classification supervisée 33 15.1 Prévision en classification supervisée 33 15.2 Performance d'une règle 33 15.2.1 Erreur de classification et accuracy 33 15.2.2 Sensibilité (recall) et taux de faux pésitifs 33 15.2.3 Spécificité et taux de faux positifs 33 15.3.1 Courbe ROC 34 15.3.2 Courbe lift 34 15.4.1 Respect des proportions initiales 34 15.4.2 Maximisation d'un coût moyen 34			
13.1 Le modèle linéaire généralisé (GLM) 30 13.2 Exemple : modélisation du nombre de visites 30 13.3 Régression Log-linéaire 30 13.3.1 Le modèle 30 13.3.2 Estimation 31 13.3.3 Tests et intervalles de confiance 31 13.3.4 Choix de variables 31 14.7 Exercices 31 14 Régularisation de la vraisemblance 32 14.1 Régressions ridge, lasso et elastic-net 32 14.2 Choix du paramètre de régularisation λ 32 14.3 Group-lasso 32 14.4 Exercices 33 15 Comparaison en classification supervisée 33 15.1 Prévision en classification supervisée 33 15.2 Performance d'une règle 33 15.2.1 Erreur de classification et accuracy 33 15.2.2 Sensibilité (recall) et taux de faux négatifs 33 15.2.3 Spécificité et taux de faux positifs 33 15.3.1 Courbe ROC 34 15.3.2 Courbe lift 34 15.4 Choix du seuil 34 15.4.2 Maximisation d'indices ad hoc 34 15.4.3 Maximisation d'un coût moyen 34		12.6	Exercices
13.1 Le modèle linéaire généralisé (GLM) 30 13.2 Exemple : modélisation du nombre de visites 30 13.3 Régression Log-linéaire 30 13.3.1 Le modèle 30 13.3.2 Estimation 31 13.3.3 Tests et intervalles de confiance 31 13.3.4 Choix de variables 31 14.7 Exercices 31 14 Régularisation de la vraisemblance 32 14.1 Régressions ridge, lasso et elastic-net 32 14.2 Choix du paramètre de régularisation λ 32 14.3 Group-lasso 32 14.4 Exercices 33 15 Comparaison en classification supervisée 33 15.1 Prévision en classification supervisée 33 15.2 Performance d'une règle 33 15.2.1 Erreur de classification et accuracy 33 15.2.2 Sensibilité (recall) et taux de faux négatifs 33 15.2.3 Spécificité et taux de faux positifs 33 15.3.1 Courbe ROC 34 15.3.2 Courbe lift 34 15.4 Choix du seuil 34 15.4.2 Maximisation d'indices ad hoc 34 15.4.3 Maximisation d'un coût moyen 34	13	Rég	ression de Poisson 30
13.2 Exemple : modélisation du nombre de visites 36 13.3 Régression Log-linéaire 30 13.3.1 Le modèle 30 13.3.2 Estimation 31 13.3.3 Tests et intervalles de confiance 31 13.3.4 Choix de variables 31 13.4 Exercices 31 14 Régularisation de la vraisemblance 32 14.1 Régressions ridge, lasso et elastic-net 32 14.2 Choix du paramètre de régularisation λ 32 14.3 Group-lasso 32 14.4 Exercices 33 15 Comparaison en classification supervisée 33 15.1 Prévision en classification supervisée 33 15.2 Performance d'une règle 33 15.2.1 Erreur de classification et accuracy 33 15.2.2 Sensibilité (recall) et taux de faux négatifs 33 15.2.2 Sensibilité (recall) et taux de faux négatifs 33 15.2.3 Spécificité et taux de faux positifs 33 15.2.4 Mesure sur les tables de contingence 34 15.3.1 Courbe ROC 34 15.3.2 Courbe lift 34 15.4.2 Maximisation d'un coût moyen 34 15.4.2 Maximisation d'un coût moyen<	10	_	
13.3 Régression Log-linéaire 36 13.3.1 Le modèle 36 13.3.2 Estimation 31 13.3.3 Tests et intervalles de confiance 31 13.3.4 Choix de variables 31 13.4 Exercices 31 14 Régularisation de la vraisemblance 32 14.1 Régressions ridge, lasso et elastic-net 32 14.2 Choix du paramètre de régularisation λ 32 14.3 Group-lasso 32 14.4 Exercices 33 15 Comparaison en classification supervisée 33 15.1 Prévision en classification supervisée 33 15.2 Performance d'une règle 33 15.2.1 Erreur de classification et accuracy 33 15.2.2 Sensibilité (recall) et taux de faux négatifs 33 15.2.2 Sensibilité (recall) et taux de faux négatifs 33 15.2.4 Mesure sur les tables de contingence 34 15.3.1 Courbe ROC 34 15.3.2 Courbe lift 34 15.4.1 Respect des proportions initiales 34 15.4.2 Maximisation d'indices ad hoc 34 15.4.3 Maximisation d'un coût moyen 34 15.5.1 Les données 34			
13.3.1 Le modèle 36 13.3.2 Estimation 31 13.3.3 Tests et intervalles de confiance 31 13.3.4 Choix de variables 31 13.3.4 Exercices 31 14 Régularisation de la vraisemblance 32 14.1 Régressions ridge, lasso et elastic-net 32 14.2 Choix du paramètre de régularisation λ 32 14.3 Group-lasso 32 14.4 Exercices 33 15 Comparaison en classification supervisée 33 15.1 Prévision en classification supervisée 33 15.2 Performance d'une règle 33 15.2.1 Erreur de classification et accuracy 33 15.2.2 Sensibilité (recall) et taux de faux négatifs 33 15.2.3 Spécificité et taux de faux positifs 33 15.3 Performance d'un score 34 15.3.1 Courbe ROC 34 15.3.2 Courbe lift 34 15.4.1 Respect des proportions initiales 34 15.4.2 Maximisation d'un coût moyen 34 15.5.1 Les données 34 15.5.2 Méthodes et comparaison 34 15.6 Transformation des variables : feature engineering 35			1
13.3.2 Estimation 31 13.3.3 Tests et intervalles de confiance 31 13.3.4 Choix de variables 31 13.4 Exercices 31 14 Régularisation de la vraisemblance 32 14.1 Régressions ridge, lasso et elastic-net 32 14.2 Choix du paramètre de régularisation λ 32 14.3 Group-lasso 32 14.4 Exercices 33 15 Comparaison en classification supervisée 33 15.1 Prévision en classification supervisée 33 15.2 Performance d'une règle 33 15.2.1 Erreur de classification et accuracy 33 15.2.2 Sensibilité (recall) et taux de faux négatifs 33 15.2.2 Spécificité et taux de faux positifs 33 15.2.3 Spécificité et taux de faux positifs 33 15.3 Performance d'un score 34 15.3.1 Courbe ROC 34 15.3.2 Courbe lift 34 15.4.1 Respect des proportions initiales 34 15.4.2 Maximisation d'indices ad hoc 34 15.5.1 Les données 34 15.5.2 Méthodes et comparaison 34 15.6.1 Modèle de prévision avec interactions 35 </td <td></td> <td>10.0</td> <td></td>		10.0	
13.3.3 Tests et intervalles de confiance 31 13.3.4 Choix de variables 31 13.4 Exercices 31 14 Régularisation de la vraisemblance 32 14.1 Régressions ridge, lasso et elastic-net 32 14.2 Choix du paramètre de régularisation λ 32 14.3 Group-lasso 32 14.4 Exercices 33 15 Comparaison en classification supervisée 33 15.1 Prévision en classification supervisée 33 15.2 Performance d'une règle 33 15.2.1 Erreur de classification et accuracy 33 15.2.2 Sensibilité (recall) et taux de faux négatifs 33 15.2.3 Spécificité et taux de faux positifs 33 15.2.4 Mesure sur les tables de contingence 34 15.3.1 Courbe ROC 34 15.3.2 Courbe lift 34 15.4.2 Maximisation d'indices ad hoc 34 15.4.3 Maximisation d'un coût moyen 34 15.5.1 Les données 34 15.5.2 Méthodes et comparaison 34 15.6.1 Modèle de prévision avec interactions 35 15.6.2 Modèle de prévision avec des polynômes 35 15.7 Exercices			
$13.3.4 \text{ Choix de variables} \qquad 31$ $13.4 \text{ Exercices} \qquad 31$ $14 \text{ Régularisation de la vraisemblance} \qquad 32$ $14.1 \text{ Régressions ridge, lasso et elastic-net} \qquad 32$ $14.2 \text{ Choix du paramètre de régularisation } \lambda \qquad 32$ $14.3 \text{ Group-lasso} \qquad 32$ $14.4 \text{ Exercices} \qquad 33$ $15 \text{ Comparaison en classification supervisée} \qquad 33$ $15.1 \text{ Prévision en classification supervisée} \qquad 33$ $15.2 \text{ Performance d'une règle} \qquad 33$ $15.2.2 \text{ Sensibilité } (recall) \text{ et aux de faux négatifs} \qquad 33$ $15.2.3 \text{ Spécificité et taux de faux positifs} \qquad 33$ $15.2.4 \text{ Mesure sur les tables de contingence} \qquad 34$ $15.3 \text{ Performance d'un score} \qquad 34$ $15.3.1 \text{ Courbe ROC} \qquad 34$ $15.3.2 \text{ Courbe lift} \qquad 34$ $15.4.1 \text{ Respect des proportions initiales} \qquad 34$ $15.4.2 \text{ Maximisation d'indices ad hoc} \qquad 34$ $15.4.3 \text{ Maximisation d'un coût moyen} \qquad 34$ $15.5 \text{ Analyse des données chd} \qquad 34$ $15.5.1 \text{ Les données} \qquad 34$ $15.5.2 \text{ Méthodes et comparaison} \qquad 34$ $15.6.1 \text{ Modèle de prévision avec interactions} \qquad 35$ $15.6.2 \text{ Modèle de prévision avec des polynômes} \qquad 35$ $15.6.2 \text{ Modèle de prévision avec des polynômes} \qquad 35$ $15.6.2 \text{ Nodèle de prévision avec des polynômes} \qquad 35$ $16.1 \text{ Données déséquilibrées} \qquad 35$ $16.1 \text{ Données déséquilibrées} \text{ et modèle logistique} \qquad 36$			
$13.4 \ \text{Exercices} \qquad \qquad$			
14 Régularisation de la vraisemblance 32 14.1 Régressions ridge, lasso et elastic-net 32 14.2 Choix du paramètre de régularisation λ 32 14.3 Group-lasso 32 14.4 Exercices 33 15 Comparaison en classification supervisée 33 15.1 Prévision en classification supervisée 33 15.2 Performance d'une règle 33 15.2.1 Erreur de classification et accuracy 33 15.2.2 Sensibilité (recall) et taux de faux négatifs 33 15.2.3 Spécificité et taux de faux positifs 33 15.2.4 Mesure sur les tables de contingence 34 15.3 Performance d'un score 34 15.3.2 Courbe lift 34 15.4 Choix du seuil 34 15.4.1 Respect des proportions initiales 34 15.4.2 Maximisation d'indices ad hoc 34 15.5 Analyse des données chd 34 15.5.2 Méthodes et comparaison 34 15.6.1 Les données 34 15.6.2 Modèle de prévision avec interactions 35 15.6.2 Modèle de prévision avec des polynômes 35 15.6.2 Modèle de prévision avec des polynômes 35		19.4	
$14.1 \ \text{Régressions ridge, lasso et elastic-net} \qquad 32$ $14.2 \ \text{Choix du paramètre de régularisation} \lambda \qquad 32$ $14.3 \ \text{Group-lasso} \qquad 32$ $14.4 \ \text{Exercices} \qquad 33$ $15.4 \ \text{Exercices} \qquad 33$ $15.1 \ \text{Prévision en classification supervisée} \qquad 33$ $15.2 \ \text{Performance d'une règle} \qquad 33$ $15.2.1 \ \text{Erreur de classification et } accuracy \qquad 33$ $15.2.2 \ \text{Sensibilité } (recall) \ \text{et aux de faux négatifs} \qquad 33$ $15.2.3 \ \text{Spécificité et taux de faux positifs} \qquad 33$ $15.2.4 \ \text{Mesure sur les tables de contingence} \qquad 34$ $15.3 \ \text{Performance d'un score} \qquad 34$ $15.3.1 \ \text{Courbe ROC} \qquad 34$ $15.3.2 \ \text{Courbe lift} \qquad 34$ $15.4.1 \ \text{Respect des proportions initiales} \qquad 34$ $15.4.2 \ \text{Maximisation d'indices ad hoc} \qquad 34$ $15.4.3 \ \text{Maximisation d'indices ad hoc} \qquad 34$ $15.4.3 \ \text{Maximisation d'indices ad hoc} \qquad 34$ $15.5.1 \ \text{Les données} \ \text{des données chd} \qquad 34$ $15.5.1 \ \text{Les données} \ \text{des données chd} \qquad 34$ $15.6 \ \text{Transformation des variables : feature engineering} \qquad 35$ $15.6.1 \ \text{Modèle de prévision avec interactions} \qquad 35$ $15.6.2 \ \text{Modèle de prévision avec des polynômes} \qquad 35$ $15.7 \ \text{Exercices} \qquad 35$ $16 \ \text{Données déséquilibrées} \qquad 35$ $16.1 \ \text{Données déséquilibrées} \qquad 35$ $16.1.1 \ \text{Un exemple} \qquad 35$ $16.1.2 \ \text{Rééquilibrage pour le modèle logistique} \qquad 36$		13.4	Exercices
14.2 Choix du paramètre de régularisation λ 3214.3 Group-lasso3214.4 Exercices3315 Comparaison en classification supervisée3315.1 Prévision en classification supervisée3315.2 Performance d'une règle3315.2.1 Erreur de classification et accuracy3315.2.2 Sensibilité (recall) et taux de faux négatifs3315.2.3 Spécificité et taux de faux positifs3315.2.4 Mesure sur les tables de contingence3415.3 Performance d'un score3415.3.1 Courbe ROC3415.3.2 Courbe lift3415.4 Choix du seuil3415.4.2 Maximisation d'indices ad hoc3415.4.3 Maximisation d'un coût moyen3415.5 Analyse des données3415.5.1 Les données3415.5.2 Méthodes et comparaison3415.6 Transformation des variables : feature engineering3515.6.2 Modèle de prévision avec interactions3515.6.2 Modèle de prévision avec des polynômes3515.7 Exercices3516 Données déséquilibrées3516.1.1 Un exemple3516.1.2 Rééquilibrage pour le modèle logistique363636	14	_	
14.3 Group-lasso 32 14.4 Exercices 33 15 Comparaison en classification supervisée 33 15.1 Prévision en classification supervisée 33 15.2 Performance d'une règle 33 15.2.1 Erreur de classification et accuracy 33 15.2.2 Sensibilité (recall) et taux de faux négatifs 33 15.2.3 Spécificité et taux de faux positifs 33 15.2.4 Mesure sur les tables de contingence 34 15.3 Performance d'un score 34 15.3.1 Courbe ROC 34 15.3.2 Courbe lift 34 15.4.1 Respect des proportions initiales 34 15.4.2 Maximisation d'indices ad hoc 34 15.4.3 Maximisation d'un coût moyen 34 15.5.1 Les données 34 15.5.2 Méthodes et comparaison 34 15.6 Transformation des variables : feature engineering 35 15.6.2 Modèle de prévision avec interactions 35 15.6.2 Modèle de prévision avec des polynômes 35 15.7 Exercices 35 16 Données déséquilibrées 35 16.1.2 Rééquilibrage pour le modèle logistique 36 36			
14.4 Exercices 33 15 Comparaison en classification supervisée 33 15.1 Prévision en classification supervisée 33 15.2 Performance d'une règle 33 15.2.1 Erreur de classification et accuracy 33 15.2.2 Sensibilité (recall) et taux de faux négatifs 33 15.2.3 Spécificité et taux de faux positifs 33 15.2.4 Mesure sur les tables de contingence 34 15.3 Performance d'un score 34 15.3.1 Courbe ROC 34 15.3.2 Courbe lift 34 15.4 Choix du seuil 34 15.4.2 Maximisation d'indices ad hoc 34 15.4.3 Maximisation d'indices ad hoc 34 15.5 Analyse des données chd 34 15.5.1 Les données 34 15.5.2 Méthodes et comparaison 34 15.6.1 Modèle de prévision avec interactions 35 15.6.2 Modèle de prévision avec des polynômes 35 15.7 Exercices 35 16 Données déséquilibrées 35 16.1.1 Un exemple 35 16.1.2 Rééquilibrage pour le modèle logistique 36			
15 Comparaison en classification supervisée 33 15.1 Prévision en classification supervisée 33 15.2 Performance d'une règle 33 15.2.1 Erreur de classification et accuracy 33 15.2.2 Sensibilité (recall) et taux de faux négatifs 33 15.2.3 Spécificité et taux de faux positifs 33 15.2.4 Mesure sur les tables de contingence 34 15.3 Performance d'un score 34 15.3.1 Courbe ROC 34 15.3.2 Courbe lift 34 15.4 Choix du seuil 34 15.4.1 Respect des proportions initiales 34 15.4.2 Maximisation d'indices ad hoc 34 15.5 Analyse des données chd 34 15.5.1 Les données 34 15.5.2 Méthodes et comparaison 34 15.6 Transformation des variables : feature engineering 35 15.6.2 Modèle de prévision avec interactions 35 15.6.2 Modèle de prévision avec des polynômes 35 15.7 Exercices 35 16 Données déséquilibrées 35 16.1.1 Un exemple 35 16.1.2 Rééquilibrage pour le modèle logistique 36		14.3	Group-lasso
15.1 Prévision en classification supervisée 33 15.2 Performance d'une règle 33 15.2.1 Erreur de classification et accuracy 33 15.2.2 Sensibilité (recall) et taux de faux négatifs 33 15.2.3 Spécificité et taux de faux positifs 33 15.2.4 Mesure sur les tables de contingence 34 15.3 Performance d'un score 34 15.3.1 Courbe ROC 34 15.3.2 Courbe lift 34 15.4 Choix du seuil 34 15.4.1 Respect des proportions initiales 34 15.4.2 Maximisation d'indices ad hoc 34 15.5 Analyse des données chd 34 15.5.1 Les données 34 15.5.2 Méthodes et comparaison 34 15.6 Transformation des variables : feature engineering 35 15.6.1 Modèle de prévision avec interactions 35 15.6.2 Modèle de prévision avec des polynômes 35 15.7 Exercices 35 16 Données déséquilibrées 35 16.1.1 Un exemple 35 16.1.2 Rééquilibrage pour le modèle logistique 36		14.4	Exercices
15.2 Performance d'une règle 33 15.2.1 Erreur de classification et accuracy 33 15.2.2 Sensibilité (recall) et taux de faux négatifs 33 15.2.3 Spécificité et taux de faux positifs 33 15.2.4 Mesure sur les tables de contingence 34 15.3 Performance d'un score 34 15.3.1 Courbe ROC 34 15.3.2 Courbe lift 34 15.4 Choix du seuil 34 15.4.1 Respect des proportions initiales 34 15.4.2 Maximisation d'indices ad hoc 34 15.5 Analyse des données chd 34 15.5.1 Les données 34 15.5.2 Méthodes et comparaison 34 15.6.1 Modèle de prévision avec interactions 35 15.6.2 Modèle de prévision avec des polynômes 35 15.7 Exercices 35 16 Données déséquilibrées 35 16.1 Données déséquilibrées et modèle logistique 35 16.1.2 Rééquilibrage pour le modèle logistique 36	15	Con	paraison en classification supervisée 33
15.2 Performance d'une règle 33 15.2.1 Erreur de classification et accuracy 33 15.2.2 Sensibilité (recall) et taux de faux négatifs 33 15.2.3 Spécificité et taux de faux positifs 33 15.2.4 Mesure sur les tables de contingence 34 15.3 Performance d'un score 34 15.3.1 Courbe ROC 34 15.3.2 Courbe lift 34 15.4 Choix du seuil 34 15.4.1 Respect des proportions initiales 34 15.4.2 Maximisation d'indices ad hoc 34 15.5 Analyse des données chd 34 15.5.1 Les données 34 15.5.2 Méthodes et comparaison 34 15.6.1 Modèle de prévision avec interactions 35 15.6.2 Modèle de prévision avec des polynômes 35 15.7 Exercices 35 16 Données déséquilibrées 35 16.1 Données déséquilibrées et modèle logistique 35 16.1.2 Rééquilibrage pour le modèle logistique 36		15.1	Prévision en classification supervisée
15.2.1 Erreur de classification et accuracy 33 15.2.2 Sensibilité (recall) et taux de faux négatifs 33 15.2.3 Spécificité et taux de faux positifs 33 15.2.4 Mesure sur les tables de contingence 34 15.3 Performance d'un score 34 15.3.1 Courbe ROC 34 15.3.2 Courbe lift 34 15.4 Choix du seuil 34 15.4.1 Respect des proportions initiales 34 15.4.2 Maximisation d'indices ad hoc 34 15.5 Analyse des données chd 34 15.5 Les données 34 15.5.1 Les données 34 15.5.2 Méthodes et comparaison 34 15.6 Transformation des variables : feature engineering 35 15.6.1 Modèle de prévision avec interactions 35 15.7 Exercices 35 16 Données déséquilibrées 35 16.1 Données déséquilibrées et modèle logistique 35 16.1.2 Rééquilibrage pour le modèle logistique 36			
15.2.2 Sensibilité (recall) et taux de faux négatifs 33 15.2.3 Spécificité et taux de faux positifs 33 15.2.4 Mesure sur les tables de contingence 34 15.3 Performance d'un score 34 15.3.1 Courbe ROC 34 15.3.2 Courbe lift 34 15.4 Choix du seuil 34 15.4.1 Respect des proportions initiales 34 15.4.2 Maximisation d'indices ad hoc 34 15.5 Analyse des données chd 34 15.5.1 Les données 34 15.5.2 Méthodes et comparaison 34 15.6 Transformation des variables : feature engineering 35 15.6.1 Modèle de prévision avec interactions 35 15.6.2 Modèle de prévision avec des polynômes 35 15.7 Exercices 35 16 Données déséquilibrées 35 16.1 Données déséquilibrées et modèle logistique 35 16.1.2 Rééquilibrage pour le modèle logistique 36			
15.2.3 Spécificité et taux de faux positifs 33 15.2.4 Mesure sur les tables de contingence 34 15.3 Performance d'un score 34 15.3.1 Courbe ROC 34 15.3.2 Courbe lift 34 15.4 Choix du seuil 34 15.4.1 Respect des proportions initiales 34 15.4.2 Maximisation d'indices ad hoc 34 15.4.3 Maximisation d'un coût moyen 34 15.5 Analyse des données chd 34 15.5.1 Les données 34 15.5.2 Méthodes et comparaison 34 15.6.1 Modèle de prévision avec interactions 35 15.6.2 Modèle de prévision avec des polynômes 35 15.7 Exercices 35 16 Données déséquilibrées 35 16.1 Données déséquilibrées et modèle logistique 35 16.1.2 Rééquilibrage pour le modèle logistique 36 16.1.2 Rééquilibrage pour le modèle logistique 36			
15.2.4 Mesure sur les tables de contingence 34 15.3 Performance d'un score 34 15.3.1 Courbe ROC 34 15.3.2 Courbe lift 34 15.4 Choix du seuil 34 15.4.1 Respect des proportions initiales 34 15.4.2 Maximisation d'indices ad hoc 34 15.4.3 Maximisation d'un coût moyen 34 15.5 Analyse des données chd 34 15.5.1 Les données 34 15.5.2 Méthodes et comparaison 34 15.6 Transformation des variables : feature engineering 35 15.6.1 Modèle de prévision avec interactions 35 15.7 Exercices 35 16 Données déséquilibrées 35 16.1 Données déséquilibrées et modèle logistique 35 16.1.2 Rééquilibrage pour le modèle logistique 36 16.1.2 Rééquilibrage pour le modèle logistique 36			· /
15.3 Performance d'un score 34 15.3.1 Courbe ROC 34 15.3.2 Courbe lift 34 15.4 Choix du seuil 34 15.4.1 Respect des proportions initiales 34 15.4.2 Maximisation d'indices ad hoc 34 15.4.3 Maximisation d'un coût moyen 34 15.5 Analyse des données chd 34 15.5.1 Les données 34 15.5.2 Méthodes et comparaison 34 15.6 Transformation des variables : feature engineering 35 15.6.1 Modèle de prévision avec interactions 35 15.7 Exercices 35 16 Données déséquilibrées 35 16.1 Données déséquilibrées et modèle logistique 35 16.1.1 Un exemple 35 16.1.2 Rééquilibrage pour le modèle logistique 36 36 36			
15.3.1 Courbe ROC 34 15.3.2 Courbe lift 34 15.4 Choix du seuil 34 15.4.1 Respect des proportions initiales 34 15.4.2 Maximisation d'indices ad hoc 34 15.4.3 Maximisation d'un coût moyen 34 15.5 Analyse des données chd 34 15.5.1 Les données 34 15.5.2 Méthodes et comparaison 34 15.6 Transformation des variables : feature engineering 35 15.6.1 Modèle de prévision avec interactions 35 15.6.2 Modèle de prévision avec des polynômes 35 15.7 Exercices 35 16 Données déséquilibrées 35 16.1 Données déséquilibrées et modèle logistique 35 16.1.2 Rééquilibrage pour le modèle logistique 36		15.3	
15.3.2 Courbe lift 34 15.4 Choix du seuil 34 15.4.1 Respect des proportions initiales 34 15.4.2 Maximisation d'indices ad hoc 34 15.4.3 Maximisation d'un coût moyen 34 15.5 Analyse des données chd 34 15.5.1 Les données 34 15.5.2 Méthodes et comparaison 34 15.6 Transformation des variables : feature engineering 35 15.6.1 Modèle de prévision avec interactions 35 15.6.2 Modèle de prévision avec des polynômes 35 15.7 Exercices 35 16 Données déséquilibrées 35 16.1 Données déséquilibrées et modèle logistique 35 16.1.2 Rééquilibrage pour le modèle logistique 36		10.0	
15.4 Choix du seuil 34 15.4.1 Respect des proportions initiales 34 15.4.2 Maximisation d'indices ad hoc 34 15.4.3 Maximisation d'un coût moyen 34 15.5 Analyse des données chd 34 15.5.1 Les données 34 15.5.2 Méthodes et comparaison 34 15.6 Transformation des variables : feature engineering 35 15.6.1 Modèle de prévision avec interactions 35 15.6.2 Modèle de prévision avec des polynômes 35 15.7 Exercices 35 16 Données déséquilibrées 35 16.1 Données déséquilibrées et modèle logistique 35 16.1.2 Rééquilibrage pour le modèle logistique 36			
15.4.1 Respect des proportions initiales 34 15.4.2 Maximisation d'indices ad hoc 34 15.4.3 Maximisation d'un coût moyen 34 15.5 Analyse des données chd 34 15.5.1 Les données 34 15.5.2 Méthodes et comparaison 34 15.6 Transformation des variables : feature engineering 35 15.6.1 Modèle de prévision avec interactions 35 15.6.2 Modèle de prévision avec des polynômes 35 15.7 Exercices 35 16 Données déséquilibrées 35 16.1 Données déséquilibrées et modèle logistique 35 16.1.2 Rééquilibrage pour le modèle logistique 36		15.4	
15.4.2 Maximisation d'indices ad hoc 34 15.4.3 Maximisation d'un coût moyen 34 15.5 Analyse des données chd 34 15.5.1 Les données 34 15.5.2 Méthodes et comparaison 34 15.6 Transformation des variables : feature engineering 35 15.6.1 Modèle de prévision avec interactions 35 15.6.2 Modèle de prévision avec des polynômes 35 15.7 Exercices 35 16 Données déséquilibrées 35 16.1 Données déséquilibrées et modèle logistique 35 16.1.2 Rééquilibrage pour le modèle logistique 36		10.1	
15.4.3 Maximisation d'un coût moyen 34 15.5 Analyse des données chd 34 15.5.1 Les données 34 15.5.2 Méthodes et comparaison 34 15.6 Transformation des variables : feature engineering 35 15.6.1 Modèle de prévision avec interactions 35 15.6.2 Modèle de prévision avec des polynômes 35 15.7 Exercices 35 16 Données déséquilibrées 35 16.1.1 Un exemple 35 16.1.2 Rééquilibrage pour le modèle logistique 36			1 1 1
15.5 Analyse des données chd 34 15.5.1 Les données 34 15.5.2 Méthodes et comparaison 34 15.6 Transformation des variables : feature engineering 35 15.6.1 Modèle de prévision avec interactions 35 15.6.2 Modèle de prévision avec des polynômes 35 15.7 Exercices 35 16 Données déséquilibrées 35 16.1 Données déséquilibrées et modèle logistique 35 16.1.1 Un exemple 35 16.1.2 Rééquilibrage pour le modèle logistique 36			
15.5.1 Les données 34 15.5.2 Méthodes et comparaison 34 15.6 Transformation des variables : feature engineering 35 15.6.1 Modèle de prévision avec interactions 35 15.6.2 Modèle de prévision avec des polynômes 35 15.7 Exercices 35 16 Données déséquilibrées 35 16.1 Données déséquilibrées et modèle logistique 35 16.1.1 Un exemple 35 16.1.2 Rééquilibrage pour le modèle logistique 36		15.5	· · · · · · · · · · · · · · · · · · ·
15.5.2 Méthodes et comparaison 34 15.6 Transformation des variables : feature engineering 35 15.6.1 Modèle de prévision avec interactions 35 15.6.2 Modèle de prévision avec des polynômes 35 15.7 Exercices 35 16 Données déséquilibrées 35 16.1 Données déséquilibrées et modèle logistique 35 16.1.1 Un exemple 35 16.1.2 Rééquilibrage pour le modèle logistique 36		10.0	
15.6 Transformation des variables : feature engineering 35 15.6.1 Modèle de prévision avec interactions 35 15.6.2 Modèle de prévision avec des polynômes 35 15.7 Exercices 35 16 Données déséquilibrées 35 16.1 Données déséquilibrées et modèle logistique 35 16.1.1 Un exemple 35 16.1.2 Rééquilibrage pour le modèle logistique 36			
15.6.1 Modèle de prévision avec interactions 35 15.6.2 Modèle de prévision avec des polynômes 35 15.7 Exercices 35 16 Données déséquilibrées 35 16.1 Données déséquilibrées et modèle logistique 35 16.1.1 Un exemple 35 16.1.2 Rééquilibrage pour le modèle logistique 36			
15.6.2 Modèle de prévision avec des polynômes 35 15.7 Exercices 35 16 Données déséquilibrées 35 16.1 Données déséquilibrées et modèle logistique 35 16.1.1 Un exemple 35 16.1.2 Rééquilibrage pour le modèle logistique 36			
15.7 Exercices 35 16 Données déséquilibrées 35 16.1 Données déséquilibrées et modèle logistique 35 16.1.1 Un exemple 35 16.1.2 Rééquilibrage pour le modèle logistique 36			
16 Données déséquilibrées 35 16.1 Données déséquilibrées et modèle logistique 35 16.1.1 Un exemple 35 16.1.2 Rééquilibrage pour le modèle logistique 36			
16.1 Données déséquilibrées et modèle logistique 35 16.1.1 Un exemple 35 16.1.2 Rééquilibrage pour le modèle logistique 36		15.7	Exercices
16.1.1 Un exemple	16	Don	nées déséquilibrées 35
16.1.1 Un exemple		16.1	Données déséquilibrées et modèle logistique
			16.1.2 Rééquilibrage pour le modèle logistique
10.1.5 Exemples de schema de reequinorage			16.1.3 Exemples de schéma de rééquilibrage







"regression" — 2024/12/31 — 17:10 — page xvii — #11



			Table des matières	xvii
		Stratégies pour données déséquilibrées 16.2.1 Quelques méthodes de rééquilibrage 16.2.2 Critères pour données déséquilibrées Choisir un algorithme de rééquilibrage		366 367 372 375 376
	16.4	16.3.2 Application aux données d'images publ Exercices		$377 \\ 381$
A	A.2	Rappels d'algèbre		
Bibliographie		raphie		393
In	\mathbf{dex}			393
No	Notations			397
Fo	Fonctions et packages python			397
Fo	Fonctions et packages R		397	



