



# Table des matières

Remerciements				
A	vant-	Propos	ix	
I Introduction au modèle linéaire				
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 $\beta_j$ , quelques propriétés	11	
		1.4.2 Résidus et variance résiduelle	15	
	1 5	1.4.3 Prévision	15	
	1.5	Interprétations géométriques	16 16	
		1.5.1 Représentation des individus	10	
	1.6	1.5.2 Représentation des variables	19	
	$1.0 \\ 1.7$	Exemples	$\frac{19}{22}$	
	1.8	Exercices	28	
	1.0	Exercices	20	
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" — 2025/2/11 — 17:35 — page xii — #6



#### xii Régression avec Python

		2.3.5 Prévision	41		
	2.4	Interprétation géométrique	42		
	$\frac{2.1}{2.5}$	Exemples	43		
	2.6	Exercices	46		
			10		
3					
	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		
	0.0	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		
	2.5	3.4.3 Résidus partiels et résidus partiels augmentés	65		
	3.5	Exemple: la concentration en ozone	67		
	3.6	Exercices	71		
4	Ext	ensions : non-inversibilité et (ou) erreurs corrélées	<b>73</b>		
	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			
		4.1.4 Propriétés statistiques de l'estimateur ridge $\beta_{\text{ridge}}$	76		
	$4.2\mathrm{I}$	Erreurs corrélées : moindres carrés généralisés	78		
		4.2.1 Erreurs hétéroscédastiques	79		
		4.2.2 Estimateur des moindres carrés généralisés	81		
		4.2.3 Matrice $\Omega$ inconnue	84		
	4.3	Exercices	85		
5	Rég	gression polynomiale et régression spline	87		
	5.1	Régression polynomiale	87		
	5.2	Régression spline	91		
		5.2.1 Introduction	91		
		5.2.2 Spline de régression	92		
	5.3	Spline de lissage	96		
	5.4	Exercices	99		
II	Ir	nférence	101		
6	Infé	érence dans le modèle gaussien	103		
_	6.1	Estimateurs du maximum de vraisemblance	103		







### "regression" — 2025/2/11 — 17:35 — page xiii — #7



_			Table des	matières	Xi
	6.2		s propriétés statistiques		
	6.3	Intervalle	es et régions de confiance		
	6.4	Prévision			
	6.5		d'hypothèses		
			ntroduction		
			'est entre modèles emboîtés		
	6.6		ions		
	6.7		3		
	6.8				
			ntervalle de confiance : bootstrap		
			est de Fisher pour une hypothèse linéaire quelco		
		6.8.3 P	ropriétés asymptotiques		. 125
7	Var	iables qu	alitatives : ANCOVA et ANOVA		129
	7.1	Introduc	tion		. 129
	7.2		de la covariance		
			ntroduction: exemple des eucalyptus		
			Iodélisation du problème		
			[ypothèse gaussienne		
			xemple: la concentration en ozone		
			xemple: la hauteur des eucalyptus		
	7.3		de la variance à 1 facteur		
			ntroduction		
			Iodélisation du problème		
			nterprétation des contraintes		
			stimation des paramètres		
			Typothèse gaussienne et test d'influence du facte		
			xemple: la concentration en ozone		
			Ine décomposition directe de la variance		
	7.4		de la variance à 2 facteurs		
			ntroduction		
			fodélisation du problème		
			stimation des paramètres		
			nalyse graphique de l'interaction		
			Sypothèse gaussienne et test de l'interaction		. 158
			xemple: la concentration en ozone		
	7.5	Exercices			
	7.6	Note: id	entifiabilité et contrastes		. 165
II	I 1	Réducti	on de dimension		167
8		oix de va			169
	8.1		tion		. 169
	8.2	Notation	8		. 171











### ${\rm xiv} \qquad {\rm R\'egression \ avec \ Python}$

	8.3	Choix incorrect de variables : conséquences	172
		8.3.1 Biais des estimateurs	172
		8.3.2 Variance des estimateurs	174
		8.3.3 Erreur quadratique moyenne	17
			17
	8.4	Critères classiques de choix de modèles	179
			180
			18
			182
			183
			18
			18
	8.5		189
			189
			189
	8.6		19:
			19:
			192
	8.7		19:
	8.8		19
	0.0	2. oct of ot state de selection 1	
9	Rég	ularisation des moindres carrés : Ridge, Lasso et elastic-net 1	99
	_		
	9.1	Introduction	199
	9.1 9.2		199 202
	-	Problème du centrage-réduction des variables	
	9.2	Problème du centrage-réduction des variables	202
	9.2	Problème du centrage-réduction des variables	202 203
	9.2	Problème du centrage-réduction des variables	202 203 207
	9.2	Problème du centrage-réduction des variables	202 203 203 208
	9.2 9.3	Problème du centrage-réduction des variables	202 203 203 208 213
	9.2 9.3	Problème du centrage-réduction des variables 2 Propriétés des régressions Ridge et lasso 2 9.3.1 Interprétation géométrique 2 9.3.2 Simplification quand les $X$ sont orthogonaux 2 9.3.3 Choix de $\lambda$ par validation croisée 2 Régularisation avec le module <b>scikitlearn</b> 2 9.4.1 Estimation des paramètres 2	202 203 208 208 213 213
	9.2 9.3	Problème du centrage-réduction des variables 2 Propriétés des régressions Ridge et lasso 2 9.3.1 Interprétation géométrique 2 9.3.2 Simplification quand les $X$ sont orthogonaux 2 9.3.3 Choix de $\lambda$ par validation croisée 2 Régularisation avec le module <b>scikitlearn</b> 2 9.4.1 Estimation des paramètres 2 9.4.2 Chemin de régularisation 2	202 203 203 208 213 213
	9.2 9.3	Problème du centrage-réduction des variables 2 Propriétés des régressions Ridge et lasso 2 9.3.1 Interprétation géométrique 2 9.3.2 Simplification quand les $X$ sont orthogonaux 2 9.3.3 Choix de $\lambda$ par validation croisée 2 Régularisation avec le module <b>scikitlearn</b> 2 9.4.1 Estimation des paramètres 2 9.4.2 Chemin de régularisation 2 9.4.3 Choix du paramètre de régularisation $\alpha$ 2	202 203 208 212 213 214 214 216
	9.2 9.3	Problème du centrage-réduction des variables 2 Propriétés des régressions Ridge et lasso 2 9.3.1 Interprétation géométrique 2 9.3.2 Simplification quand les $X$ sont orthogonaux 2 9.3.3 Choix de $\lambda$ par validation croisée 2 Régularisation avec le module <b>scikitlearn</b> 2 9.4.1 Estimation des paramètres 2 9.4.2 Chemin de régularisation 2 9.4.3 Choix du paramètre de régularisation $\alpha$ 2 9.4.4 Mise en pratique 2	202 203 207 208 212 213 214 215 216 218
	9.2 9.3 9.4	Problème du centrage-réduction des variables 2 Propriétés des régressions Ridge et lasso 2 9.3.1 Interprétation géométrique 2 9.3.2 Simplification quand les $X$ sont orthogonaux 2 9.3.3 Choix de $\lambda$ par validation croisée 2 Régularisation avec le module <b>scikitlearn</b> 2 9.4.1 Estimation des paramètres 2 9.4.2 Chemin de régularisation 2 9.4.3 Choix du paramètre de régularisation $\alpha$ 2 9.4.4 Mise en pratique 2 Intégration de variables qualitatives 2	202 203 208 212 213 214 214 216
	9.2 9.3 9.4	Problème du centrage-réduction des variables 2 Propriétés des régressions Ridge et lasso 2 9.3.1 Interprétation géométrique 2 9.3.2 Simplification quand les $X$ sont orthogonaux 2 9.3.3 Choix de $\lambda$ par validation croisée 2 Régularisation avec le module <b>scikitlearn</b> 2 9.4.1 Estimation des paramètres 2 9.4.2 Chemin de régularisation 2 9.4.3 Choix du paramètre de régularisation 2 9.4.4 Mise en pratique 2 1.1 Exercices 2 2.2 Exercices 2 2.3 Exercices 2 2.5 Exercices 2 2.5 Exercices 2 2.6 Exercices 2 2.7 Exercices 2 2.7 Exercices 2 2.7 Exercices 2 2.8 Exercices 2 2.8 Exercices 2 2.9 Exercices 2 2.0 Ex	202 203 203 203 213 214 214 218 218 218 2218
	9.2 9.3 9.4 9.5 9.6	Problème du centrage-réduction des variables 2 Propriétés des régressions Ridge et lasso 2 9.3.1 Interprétation géométrique 2 9.3.2 Simplification quand les $X$ sont orthogonaux 2 9.3.3 Choix de $\lambda$ par validation croisée 2 Régularisation avec le module <b>scikitlearn</b> 2 9.4.1 Estimation des paramètres 2 9.4.2 Chemin de régularisation 2 9.4.3 Choix du paramètre de régularisation 2 9.4.4 Mise en pratique 2 1.1 Exercices 2 2.2 Exercices 2 2.3 Exercices 2 2.5 Exercices 2 2.5 Exercices 2 2.6 Exercices 2 2.7 Exercices 2 2.7 Exercices 2 2.7 Exercices 2 2.8 Exercices 2 2.8 Exercices 2 2.8 Exercices 2 2.9 Exercices 2 2.9 Exercices 2 2.0 Ex	202 203 203 203 214 214 214 216 218 218
10	9.2 9.3 9.4 9.5 9.6 9.7	Problème du centrage-réduction des variables 2 Propriétés des régressions Ridge et lasso 2 9.3.1 Interprétation géométrique 2 9.3.2 Simplification quand les $X$ sont orthogonaux 2 9.3.3 Choix de $\lambda$ par validation croisée 2 Régularisation avec le module <b>scikitlearn</b> 2 9.4.1 Estimation des paramètres 2 9.4.2 Chemin de régularisation 2 9.4.3 Choix du paramètre de régularisation 2 9.4.4 Mise en pratique 2 Intégration de variables qualitatives 2 Exercices 2 Note : lars et lasso 2	202 203 203 203 213 214 214 218 218 218 2218
10	9.2 9.3 9.4 9.5 9.6 9.7 <b>Rég</b>	Problème du centrage-réduction des variables 2 Propriétés des régressions Ridge et lasso 2 9.3.1 Interprétation géométrique 2 9.3.2 Simplification quand les $X$ sont orthogonaux 2 9.3.3 Choix de $\lambda$ par validation croisée 2 Régularisation avec le module <b>scikitlearn</b> 2 9.4.1 Estimation des paramètres 2 9.4.2 Chemin de régularisation 2 9.4.3 Choix du paramètre de régularisation 2 9.4.4 Mise en pratique 2 Intégration de variables qualitatives 2 Exercices 2 Note : lars et lasso 2  ression sur composantes : PCR et PLS 2	202 203 203 203 213 214 214 218 2218 2218 2223
10	9.2 9.3 9.4 9.5 9.6 9.7 <b>Rég</b>	Problème du centrage-réduction des variables 2 Propriétés des régressions Ridge et lasso 2 9.3.1 Interprétation géométrique 2 9.3.2 Simplification quand les $X$ sont orthogonaux 2 9.3.3 Choix de $\lambda$ par validation croisée 2 Régularisation avec le module scikitlearn 2 9.4.1 Estimation des paramètres 2 9.4.2 Chemin de régularisation 2 9.4.3 Choix du paramètre de régularisation 2 9.4.4 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	202 203 203 203 213 214 214 214 218 218 2218 2218 2218 2218
10	9.2 9.3 9.4 9.5 9.6 9.7 <b>Rég</b>	Problème du centrage-réduction des variables 2 Propriétés des régressions Ridge et lasso 2 9.3.1 Interprétation géométrique 2 9.3.2 Simplification quand les $X$ sont orthogonaux 2 9.3.3 Choix de $\lambda$ par validation croisée 2 Régularisation avec le module <b>scikitlearn</b> 2 9.4.1 Estimation des paramètres 2 9.4.2 Chemin de régularisation 2 9.4.3 Choix du paramètre de régularisation 2 9.4.4 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	202 203 203 203 213 214 214 214 215 2218 2218 2223 2223 2230
10	9.2 9.3 9.4 9.5 9.6 9.7 <b>Rég</b>	Problème du centrage-réduction des variables 2 Propriétés des régressions Ridge et lasso 2 9.3.1 Interprétation géométrique 2 9.3.2 Simplification quand les $X$ sont orthogonaux 2 9.3.3 Choix de $\lambda$ par validation croisée 2 Régularisation avec le module <b>scikitlearn</b> 2 9.4.1 Estimation des paramètres 2 9.4.2 Chemin de régularisation 2 9.4.3 Choix du paramètre de régularisation 2 9.4.4 Mise en pratique 2 Intégration de variables qualitatives 2 Exercices 2 Note : lars et lasso 2  Pression 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	202 203 203 203 213 214 214 2218 2218 2223 2223 2230
10	9.2 9.3 9.4 9.5 9.6 9.7 <b>Rég</b>	Problème du centrage-réduction des variables 2 Propriétés des régressions Ridge et lasso 2 9.3.1 Interprétation géométrique 2 9.3.2 Simplification quand les $X$ sont orthogonaux 2 9.3.3 Choix de $\lambda$ par validation croisée 2 Régularisation avec le module <b>scikitlearn</b> 2 9.4.1 Estimation des paramètres 2 9.4.2 Chemin de régularisation 2 9.4.3 Choix du paramètre de régularisation 2 9.4.4 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	202 203 203 208 213 214 214 218 2218 2218 2223 2230 2230 2230
10	9.2 9.3 9.4 9.5 9.6 9.7 <b>Rég</b>	Problème du centrage-réduction des variables 2 Propriétés des régressions Ridge et lasso 2 9.3.1 Interprétation géométrique 2 9.3.2 Simplification quand les $X$ sont orthogonaux 2 9.3.3 Choix de $\lambda$ par validation croisée 2 Régularisation avec le module <b>scikitlearn</b> 2 9.4.1 Estimation des paramètres 2 9.4.2 Chemin de régularisation 2 9.4.3 Choix du paramètre de régularisation 2 9.4.4 Mise en pratique 2 Intégration de variables qualitatives 2 Exercices 2 Note : lars et lasso 2  Pression 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	202 203 203 203 213 214 214 218 2218 222 233 2233 2233 2233 2233





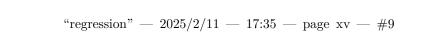




	Table des matières	X
	10.2.1 Algorithmes PLS	240
	10.2.2 Choix de composantes/variables	240
	10.2.3 Retour aux données d'origine	241
	10.2.4 La régression PLS en pratique	243
10.3	Exercices	244
10.4	Notes	246
	10.4.1 ACP et changement de base	246
	10.4.2 Colinéarité parfaite : $ X'X  = 0 \dots \dots \dots \dots$	247
11 Con	nparaison des différentes méthodes, étude de cas réels	<b>251</b>
	Erreur de prévision et validation croisée	
	Analyse de l'ozone	
	11.2.1 Préliminaires	
	11.2.2 Méthodes et comparaison	256
11.3	Transformation des variables : feature engineering	
	11.3.1 Modèle de prévision avec interactions	
	11.3.2 Modèle de prévision avec des polynômes	
	11.3.3 Modèle de prévision avec des splines	
	11.3.4 Modèle de prévision avec interactions et splines	261
	11.3.5 Conclusion	262
IV I	Le modèle linéaire généralisé	263
12 Rég	ression logistique	265
12 Rég	ression logistique Présentation du modèle	<b>265</b> 265
12 Rég	ression logistique Présentation du modèle	265 265 265
12 Rég	Présentation du modèle	265 265 265 266
<b>12</b> Rég 12.1	Présentation du modèle  12.1.1 Exemple introductif  12.1.2 Modélisation statistique  12.1.3 Variables explicatives qualitatives, interactions	265 265 266 266 269
<b>12</b> Rég 12.1	Présentation du modèle	265 265 265 266 269 271
<b>12</b> Rég 12.1	Présentation du modèle  12.1.1 Exemple introductif  12.1.2 Modélisation statistique  12.1.3 Variables explicatives qualitatives, interactions Estimation  12.2.1 La vraisemblance	265 265 265 266 269 271 271
<b>12</b> Rég 12.1	Présentation du modèle	265 265 265 266 269 271 271 273
12 Rég 12.1 12.2	Présentation du modèle  12.1.1 Exemple introductif  12.1.2 Modélisation statistique  12.1.3 Variables explicatives qualitatives, interactions Estimation  12.2.1 La vraisemblance	265 265 265 266 269 271 271 273 274
12 Rég 12.1 12.2	Présentation du modèle	265 265 265 266 269 271 271 273 274 275
12 Rég 12.1 12.2	Présentation du modèle  12.1.1 Exemple introductif  12.1.2 Modélisation statistique  12.1.3 Variables explicatives qualitatives, interactions Estimation  12.2.1 La vraisemblance  12.2.2 Calcul des estimateurs : l'algorithme IRLS  12.2.3 Propriétés asymptotiques de l'EMV Intervalles de confiance et tests  12.3.1 IC et tests sur les paramètres du modèle	265 265 266 269 271 271 273 274 275 276
12 Rég 12.1 12.2	Présentation du modèle  12.1.1 Exemple introductif  12.1.2 Modélisation statistique  12.1.3 Variables explicatives qualitatives, interactions Estimation  12.2.1 La vraisemblance  12.2.2 Calcul des estimateurs : l'algorithme IRLS  12.2.3 Propriétés asymptotiques de l'EMV Intervalles de confiance et tests	265 265 266 269 271 271 273 274 275 276 277
12 Rég 12.1 12.2	Présentation du modèle  12.1.1 Exemple introductif  12.1.2 Modélisation statistique  12.1.3 Variables explicatives qualitatives, interactions Estimation  12.2.1 La vraisemblance  12.2.2 Calcul des estimateurs : l'algorithme IRLS  12.2.3 Propriétés asymptotiques de l'EMV Intervalles de confiance et tests  12.3.1 IC et tests sur les paramètres du modèle  12.3.2 Test sur un sous-ensemble de paramètres  12.3.3 Prévision	265 265 266 269 271 271 273 274 275 276 277 280
12 Rég 12.1 12.2	Présentation du modèle  12.1.1 Exemple introductif  12.1.2 Modélisation statistique  12.1.3 Variables explicatives qualitatives, interactions Estimation  12.2.1 La vraisemblance  12.2.2 Calcul des estimateurs: l'algorithme IRLS  12.2.3 Propriétés asymptotiques de l'EMV  Intervalles de confiance et tests  12.3.1 IC et tests sur les paramètres du modèle  12.3.2 Test sur un sous-ensemble de paramètres  12.3.3 Prévision  Adéquation du modèle	265 265 266 269 271 271 273 274 275 276 277 280 282
12 Rég 12.1 12.2	Présentation du modèle  12.1.1 Exemple introductif  12.1.2 Modélisation statistique  12.1.3 Variables explicatives qualitatives, interactions Estimation  12.2.1 La vraisemblance  12.2.2 Calcul des estimateurs : l'algorithme IRLS  12.2.3 Propriétés asymptotiques de l'EMV Intervalles de confiance et tests  12.3.1 IC et tests sur les paramètres du modèle  12.3.2 Test sur un sous-ensemble de paramètres  12.3.3 Prévision	265 265 266 269 271 273 274 275 276 277 280 282 283
12 Rég 12.1 12.2	Présentation du modèle  12.1.1 Exemple introductif  12.1.2 Modélisation statistique  12.1.3 Variables explicatives qualitatives, interactions Estimation  12.2.1 La vraisemblance  12.2.2 Calcul des estimateurs : l'algorithme IRLS  12.2.3 Propriétés asymptotiques de l'EMV  Intervalles de confiance et tests  12.3.1 IC et tests sur les paramètres du modèle  12.3.2 Test sur un sous-ensemble de paramètres  12.3.3 Prévision  Adéquation du modèle  12.4.1 Le modèle saturé  12.4.2 Tests d'adéquation de la déviance et de Pearson	265 265 266 269 271 271 273 274 275 276 277 280 282 283 285
12 Rég 12.1 12.2 12.3	Présentation du modèle  12.1.1 Exemple introductif  12.1.2 Modélisation statistique  12.1.3 Variables explicatives qualitatives, interactions Estimation  12.2.1 La vraisemblance  12.2.2 Calcul des estimateurs : l'algorithme IRLS  12.2.3 Propriétés asymptotiques de l'EMV  Intervalles de confiance et tests  12.3.1 IC et tests sur les paramètres du modèle  12.3.2 Test sur un sous-ensemble de paramètres  12.3.3 Prévision  Adéquation du modèle  12.4.1 Le modèle saturé	265 265 266 269 271 271 273 274 275 276 277 280 282 283 285 288
12 Rég 12.1 12.2 12.3	Présentation du modèle  12.1.1 Exemple introductif  12.1.2 Modélisation statistique  12.1.3 Variables explicatives qualitatives, interactions Estimation  12.2.1 La vraisemblance  12.2.2 Calcul des estimateurs : l'algorithme IRLS  12.2.3 Propriétés asymptotiques de l'EMV  Intervalles de confiance et tests  12.3.1 IC et tests sur les paramètres du modèle  12.3.2 Test sur un sous-ensemble de paramètres  12.3.3 Prévision  Adéquation du modèle  12.4.1 Le modèle saturé  12.4.2 Tests d'adéquation de la déviance et de Pearson  12.4.3 Analyse des résidus	265 265 266 269 271 271 273 274 275 276 277 280 282 283 285 288 292
12 Rég 12.1 12.2 12.3	Présentation du modèle  12.1.1 Exemple introductif  12.1.2 Modélisation statistique  12.1.3 Variables explicatives qualitatives, interactions Estimation  12.2.1 La vraisemblance  12.2.2 Calcul des estimateurs : l'algorithme IRLS  12.2.3 Propriétés asymptotiques de l'EMV  Intervalles de confiance et tests  12.3.1 IC et tests sur les paramètres du modèle  12.3.2 Test sur un sous-ensemble de paramètres  12.3.3 Prévision  Adéquation du modèle  12.4.1 Le modèle saturé  12.4.2 Tests d'adéquation de la déviance et de Pearson  12.4.3 Analyse des résidus  Choix de variables	265 265 266 269 271 271 273 274 275 276 277 280 282 283 285 288 292











### ${\it xvi}$ Régression avec Python

13 Rég	ression de Poisson	301
13.1	Le modèle linéaire généralisé (GLM)	30
	Exemple: modélisation du nombre de visites	304
	Régression Log-linéaire	30'
10.0	13.3.1 Le modèle	30'
		308
	13.3.2 Estimation	
	13.3.3 Tests et intervalles de confiance	309
	13.3.4 Choix de variables	313
13.4	Exercices	314
14 Rég	ularisation de la vraisemblance	319
14.1	Régressions ridge, lasso et elastic-net	319
	Choix du paramètre de régularisation $\lambda$	$32^{2}$
	Group-lasso	32'
	Exercices	329
14.4	LACICIOCS	02.
15 Con	nparaison en classification supervisée	331
15.1	Prévision en classification supervisée	33
	Performance d'une règle	333
	15.2.1 Erreur de classification et accuracy	336
	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	338
15.9		
10.5	Performance d'un score	339
	15.3.1 Courbe ROC	339
	15.3.2 Courbe lift	34
15.4	Choix du seuil	342
	15.4.1 Respect des proportions initiales	342
	15.4.2 Maximisation d'indices ad hoc	342
	15.4.3 Maximisation d'un coût moyen	343
15.5	Analyse des données chd	$34^{2}$
	15.5.1 Les données	34
	15.5.2 Méthodes et comparaison	$34^{4}$
15.6	Transformation des variables : feature engineering	351
10.0	15.6.1 Modèle de prévision avec interactions	352
	15.6.2 Modèle de prévision avec des polynômes	
15.7	Exercices	
10.1	DACTORES	995
16 Don	mées déséquilibrées	357
	Données déséquilibrées et modèle logistique	357
	16.1.1 Un exemple	357
	16.1.2 Rééquilibrage pour le modèle logistique	359
	16.1.3 Exemples de schéma de rééquilibrage	360
16.9	Stratégies pour données déséquilibrées	365
10.2	16.2.1 Quelques méthodes de rééquilibrage	365
	10.2.1 Queiques methodes de reequinitage	506







## "regression" — 2025/2/11 — 17:35 — page xvii — #11



		r	Table des matières	xvi	
		16.2.2 Critères pour données déséquilibrées . Choisir un algorithme de rééquilibrage 16.3.1 Rééquilibrage et validation croisée 16.3.2 Application aux données d'images publication	citaires	373 374 375	
	A.1 A.2	Rappels d'algèbre		384	
Bibliographie				391	
Index					
Notations					
For	Fonctions et modules python				



