

Algorithmique des données

Régression

Charlotte Pelletier

MCF Univ. Bretagne Sud – IRISA Vannes

Basé sur le cours de Chloé Friguet (MCF UBS/IRISA).

19 février 2020

Introduction

Le modèle logistique

- Loi de Bernoulli

- Fonction logistique

Algorithme de descente du gradient pour la régression logistique

- Rappel

- Pour la régression logistique

- Optimisation

Application

Rappel

- Modélisation de la relation entre la variable à expliquer y et les variables explicatives X^1, X^2, \dots, X^d
 - y est une variable à expliquer **qualitative** (binaire)
 - X^i : variables explicatives quantitatives ou qualitatives

→ Le **modèle linéaire** n'est pas adapté

Rappel

- Modélisation de la relation entre la variable à expliquer y et les variables explicatives X^1, X^2, \dots, X^d
 - y est une variable à expliquer **qualitative** (binaire)
 - X^i : variables explicatives quantitatives ou qualitatives
- Le **modèle linéaire** n'est **pas adapté** car y n'est pas **quantitative** !

Exemple B

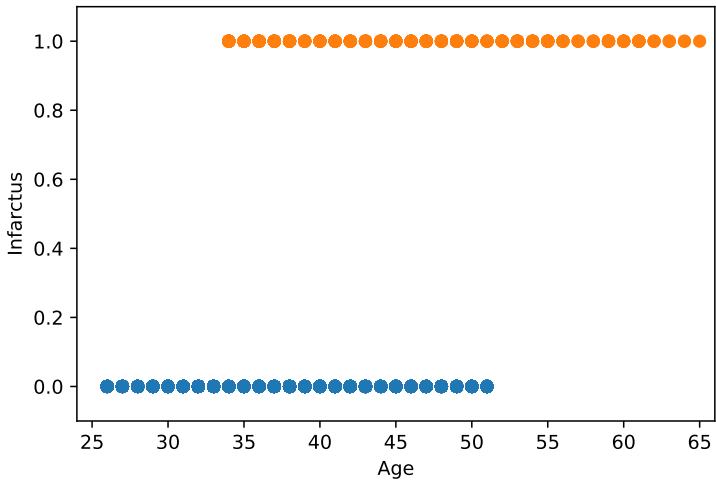
On souhaite évaluer l'existence d'un risque plus élevé de survenue d'un infarctus du myocarde chez les femmes à partir d'une étude cas-témoins. Les facteurs d'exposition recueillis sont : la prise de contraceptif, l'âge, le poids, la taille, la consommation de tabac, l'hypertension artérielle, les antécédents familiaux de maladies cardio-vasculaires, etc.

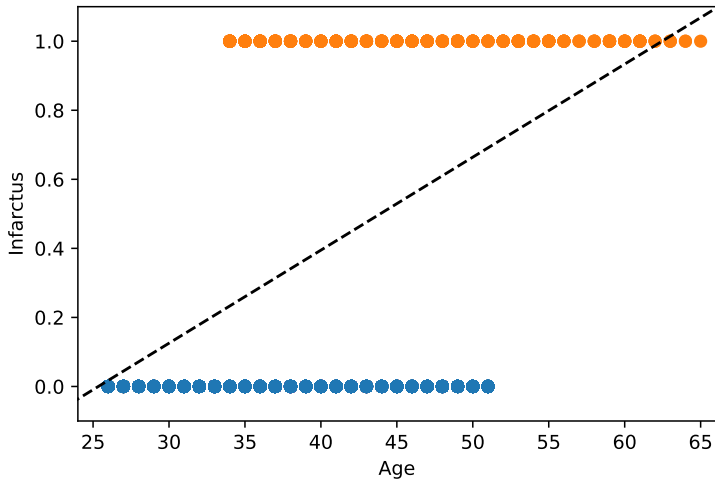
Extrait des données :

Infarctus	Contraceptif oral	Age	Poids	...
non	non	47	48	...
non	non	35	53	...
non	oui	62	41	...
oui	oui	47	45	...
oui	oui	63
oui	non	45	69	...
oui	non	90	84	...
.
.
.

Données Institut de Santé Publique, d'Epidémiologie et de Développement (ISPED), Bordeaux

Source : <http://www.biostatisticien.eu/springerR/jeuxDonnees5.html>





Introduction

Le modèle logistique

Loi de Bernoulli

Fonction logistique

Algorithme de descente du gradient pour la régression logistique

Rappel

Pour la régression logistique

Optimisation

Application

On s'adapte

- on ne va plus modéliser Y par une relation linéaire
- on s'intéresse aux probabilités : $\mathbf{P}(Y = 0|X = x)$ et $\mathbf{P}(Y = 1|X = x)$
- la connaissance de $\mathbf{P}(Y = 1|X = x)$ implique la connaissance de $\mathbf{P}(Y = 0|X = x)$ car

On s'adapte

- on ne va plus modéliser Y par une relation linéaire
- on s'intéresse aux probabilités : $\mathbf{P}(Y = 0|X = x)$ et $\mathbf{P}(Y = 1|X = x)$
- la connaissance de $\mathbf{P}(Y = 1|X = x)$ implique la connaissance de $\mathbf{P}(Y = 0|X = x)$ car $\mathbf{P}(Y = 0|X = x) = 1 - \mathbf{P}(Y = 1|X = x)$

Posons

$$\mathbf{P}(Y = 1|X = x) = \pi(x)$$

et

$$\mathbf{P}(Y = 0|X = x) = 1 - \pi(x),$$

avec $\pi(x) \in [0, 1]$.

On peut donc modéliser Y par

$$Y|X = x \sim \mathcal{B}(\pi(x))$$

(loi de Bernoulli)

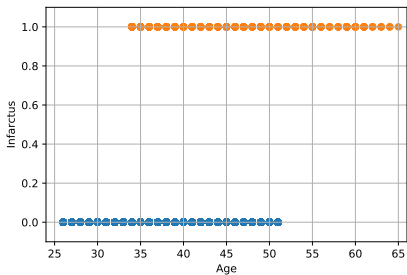
Propriétés (loi de Bernoulli)

$$\mathbf{P}(Y = y|X = x) = \pi(x)^y (1 - \pi(x))^{1-y}$$

$$\mathbb{E}(Y|X = x) = \pi(x) \quad \text{Var}(Y|X = x) = \pi(x)(1 - \pi(x))$$

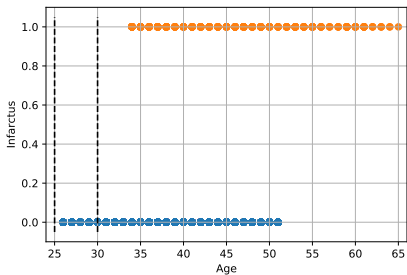
- Estimation de $\pi(x)$: $\hat{\pi}(x)$ = proportion de $Y = 1$ pour X mis en classes

classes	n_c	#0	#1	$\hat{\pi}(x)$



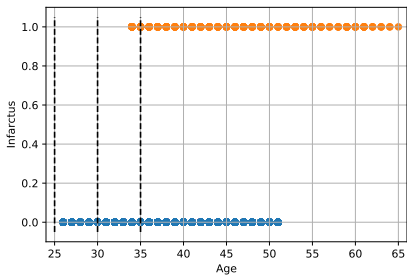
- Estimation de $\pi(x)$: $\hat{\pi}(x)$ = proportion de $Y = 1$ pour X mis en classes

classes	n_c	#0	#1	$\hat{\pi}(x)$
[25,30]	454	454	0	0.00



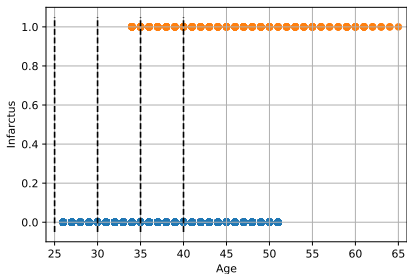
- Estimation de $\pi(x)$: $\hat{\pi}(x)$ = proportion de $Y = 1$ pour X mis en classes

classes	n_c	#0	#1	$\hat{\pi}(x)$
[25,30]	454	454	0	0.00
(30,35]	958	913	45	0.04



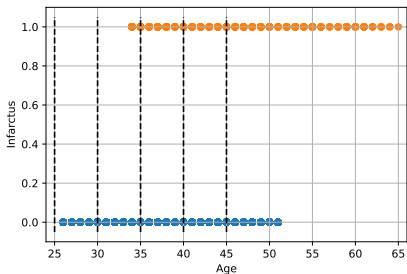
- Estimation de $\pi(x)$: $\hat{\pi}(x)$ = proportion de $Y = 1$ pour X mis en classes

classes	n_c	#0	#1	$\hat{\pi}(x)$
[25,30]	454	454	0	0.00
(30,35]	958	913	45	0.04
(35,40]	780	709	71	0.10



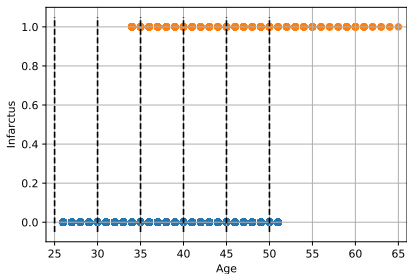
- Estimation de $\pi(x)$: $\hat{\pi}(x)$ = proportion de $Y = 1$ pour X mis en classes

classes	n_c	#0	#1	$\hat{\pi}(x)$
[25,30]	454	454	0	0.00
(30,35]	958	913	45	0.04
(35,40]	780	709	71	0.10
(40,45]	608	505	103	0.17



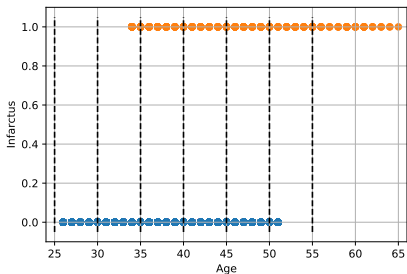
- Estimation de $\pi(x) : \hat{\pi}(x) = \text{proportion de } Y = 1 \text{ pour } X \text{ mis en classes}$

classes	n_c	#0	#1	$\hat{\pi}(x)$
[25,30]	454	454	0	0.00
(30,35]	958	913	45	0.04
(35,40]	780	709	71	0.10
(40,45]	608	505	103	0.17
(45,50]	554	355	199	0.36



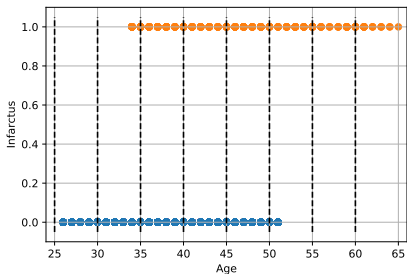
- Estimation de $\pi(x) : \hat{\pi}(x) = \text{proportion de } Y = 1 \text{ pour } X \text{ mis en classes}$

classes	n_c	#0	#1	$\hat{\pi}(x)$
[25,30]	454	454	0	0.00
(30,35]	958	913	45	0.04
(35,40]	780	709	71	0.10
(40,45]	608	505	103	0.17
(45,50]	554	355	199	0.36
(50,55]	134	30	104	0.78



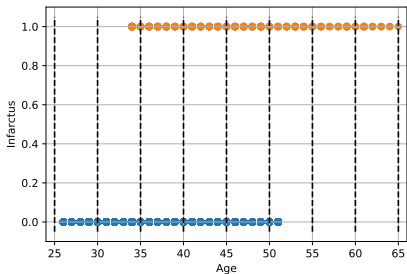
- Estimation de $\pi(x) : \hat{\pi}(x) = \text{proportion de } Y = 1 \text{ pour } X \text{ mis en classes}$

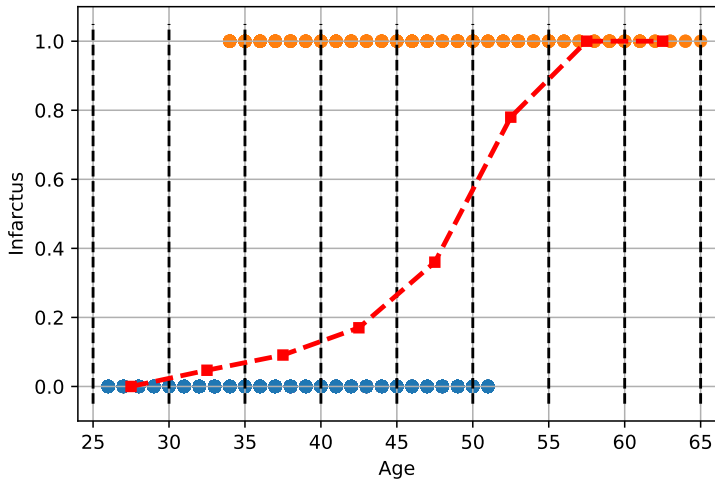
classes	n_c	#0	#1	$\hat{\pi}(x)$
[25,30]	454	454	0	0.00
(30,35]	958	913	45	0.04
(35,40]	780	709	71	0.10
(40,45]	608	505	103	0.17
(45,50]	554	355	199	0.36
(50,55]	134	30	104	0.78
(55,60]	33	0	33	1.00



- Estimation de $\pi(x)$: $\hat{\pi}(x)$ = proportion de $Y = 1$ pour X mis en classes

classes	n_c	#0	#1	$\hat{\pi}(x)$
[25,30]	454	454	0	0.00
(30,35]	958	913	45	0.04
(35,40]	780	709	71	0.10
(40,45]	608	505	103	0.17
(45,50]	554	355	199	0.36
(50,55]	134	30	104	0.78
(55,60]	33	0	33	1.00
(60,65]	17	0	17	1.00





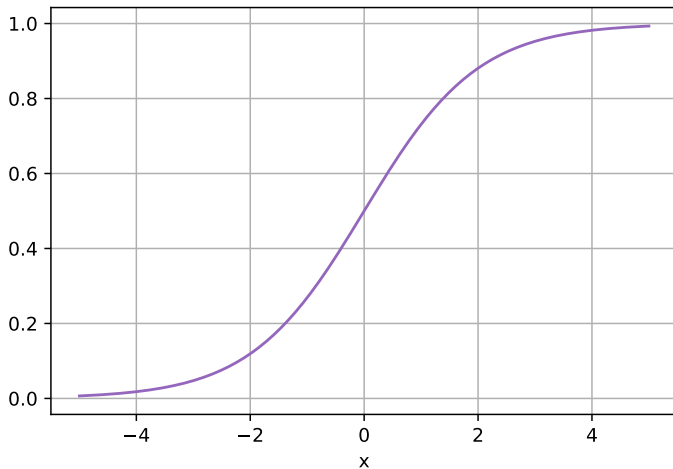
- $\pi(x)$ n'est pas linéaire en x !
- On a plutôt envie de modéliser $\pi(x)$ par une fonction "en S" \Rightarrow une possibilité est la fonction logistique

Fonction logistique

$$s : \mathbb{R} \rightarrow]0, 1[$$
$$x \rightarrow s(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{1 + e^x}$$

Cas particulier de la fonction sigmoïde

Fonction logistique



Allure de la courbe de la fonction (plus générale)

$$s : \mathbb{R} \rightarrow]0, 1[$$

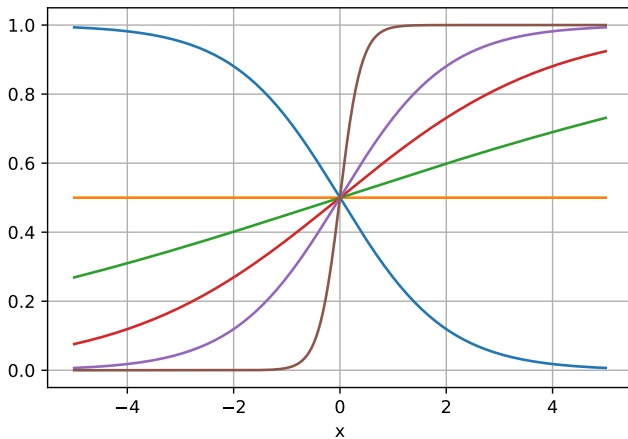
$$x \mapsto s(x) = \frac{1}{1 + e^{-x\beta}} = \frac{e^x}{1 + e^{x\beta}}$$

pour différentes valeurs de β

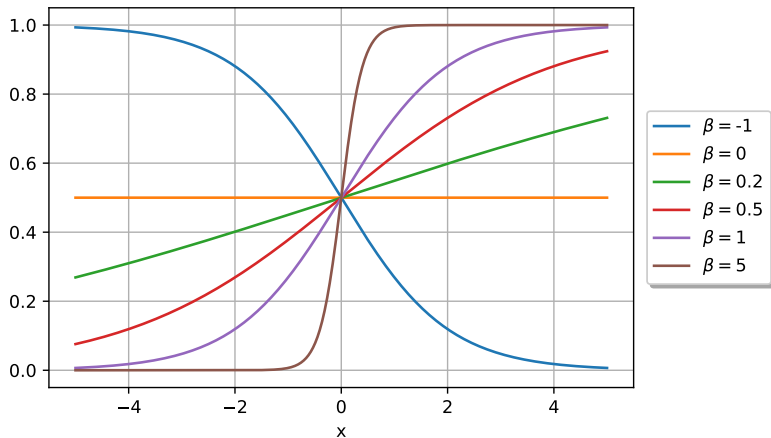
- $\beta = 0 \Rightarrow$ fonction constante
- β "petit" \Rightarrow large plage de valeurs de x pour lesquelles la fonction se situe aux alentours de 0.5 \Rightarrow discrimination difficile
- β "grand" \Rightarrow petite plage de valeurs de x pour lesquelles la fonction se situe aux alentours de 0.5 \Rightarrow discrimination plus facile

⚠ Dépend de l'échelle de x !

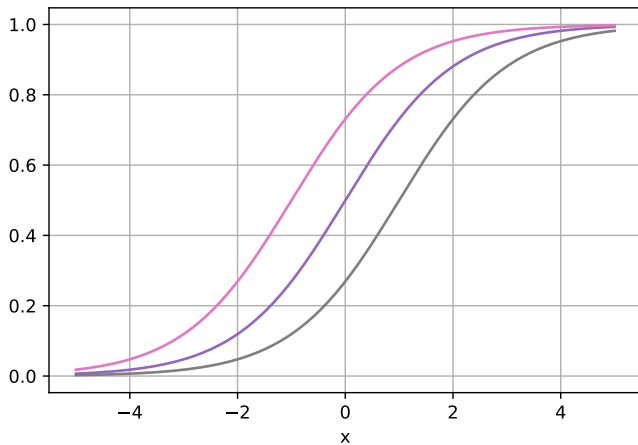
Fonction logistique



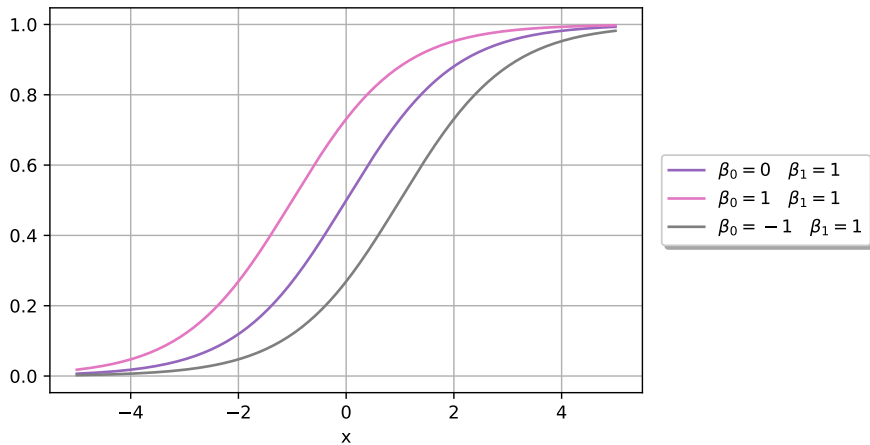
Fonction logistique



Fonction logistique



Fonction logistique



Modèle logistique

$$\begin{aligned}\pi(\mathbf{x}) &= \frac{e^{\tilde{\mathbf{x}}\boldsymbol{\beta}}}{1 + e^{\tilde{\mathbf{x}}\boldsymbol{\beta}}} = \frac{e^{\beta_0 + \beta_1 x^1 + \dots + \beta_d x^d}}{1 + e^{\beta_0 + \beta_1 x^1 + \dots + \beta_d x^d}} \\ &\quad \Updownarrow \\ \text{logit}(\pi(\mathbf{x})) &= \log\left(\frac{\pi(\mathbf{x})}{1 - \pi(\mathbf{x})}\right) = \tilde{\mathbf{x}}\boldsymbol{\beta} = \beta_0 + \beta_1 x^1 + \dots + \beta_d x^d\end{aligned}$$

- Transformation g qui permet d'avoir un lien linéaire entre $g(\pi(\mathbf{x}))$ et \mathbf{x} :
fonction de lien

Fonction logit

$$\begin{aligned}\text{logit} &:]0, 1[\rightarrow \mathbb{R} \\ p &\rightarrow \text{logit}(p) = \log\left(\frac{p}{1-p}\right)\end{aligned}$$

- On pose donc :

$$f_{\beta}(\mathbf{x}_i) = \frac{e^{\mathbf{x}_i' \beta}}{1 + e^{\mathbf{x}_i' \beta}} = \mathbb{P}(Y = 1 | X = \mathbf{x}_i)$$

- On cherche β tel que $f_{\beta}(\mathbf{x}_i)$ est proche de y_i pour toutes les données d'apprentissage $\{\mathbf{x}_i, y_i\}_{i=1}^m$ avec $\mathbf{x}_i \in \mathbb{R}^d$

- Vraisemblance du modèle :

$$L(\beta, y) = \prod_{i=1}^m \mathbf{P}(Y = y_i | X = \mathbf{x}_i) = \prod_{i=1}^m f_{\beta}(\mathbf{x}_i)^{y_i} \times (1 - f_{\beta}(\mathbf{x}_i))^{1-y_i}$$

- Vraisemblance du modèle :

$$L(\beta, y) = \prod_{i=1}^m \mathbf{P}(Y = y_i | X = \mathbf{x}_i) = \prod_{i=1}^m f_{\beta}(\mathbf{x}_i)^{y_i} \times (1 - f_{\beta}(\mathbf{x}_i))^{1-y_i}$$

- Log-vraisemblance

$$\mathcal{L}(\beta, y) = \sum_{i=1}^m \left[y_i \log(f_{\beta}(\mathbf{x}_i)) + (1 - y_i) \log(1 - f_{\beta}(\mathbf{x}_i)) \right]$$

- Vraisemblance du modèle :

$$L(\beta, y) = \prod_{i=1}^m \mathbf{P}(Y = y_i | X = \mathbf{x}_i) = \prod_{i=1}^m f_{\beta}(\mathbf{x}_i)^{y_i} \times (1 - f_{\beta}(\mathbf{x}_i))^{1-y_i}$$

- Log-vraisemblance

$$\mathcal{L}(\beta, y) = \sum_{i=1}^m \left[y_i \log(f_{\beta}(\mathbf{x}_i)) + (1 - y_i) \log(1 - f_{\beta}(\mathbf{x}_i)) \right]$$

- Maximisation de la (log-)vraisemblance
⇒ annuler les dérivées de $L(\beta, y)$
 - **pas de solution analytique explicite**
 - besoin d'**algorithmes d'optimisation itératifs**

- **Objectif** : trouver le meilleur β pour minimiser le coût (quadratique) **global** des erreurs :

$$\underset{\beta}{\operatorname{argmin}} \left(J(\beta) \right)$$

avec

$$J(\beta) = -\frac{1}{m} \sum_{i=1}^m \left[y_i \log(f_{\beta}(\mathbf{x}_i)) + (1 - y_i) \log(1 - f_{\beta}(\mathbf{x}_i)) \right]$$

Introduction

Le modèle logistique

Loi de Bernoulli

Fonction logistique

Algorithme de descente du gradient pour la régression logistique

Rappel

Pour la régression logistique

Optimisation

Application

- Objectif : trouver le minimum d'une fonction-coût
- Principe : algorithme itératif
 1. initialisation : $\beta^{(0)}$
 2. à chaque étape k , modifier $\beta^{(k-1)}$ pour faire diminuer $\mathbf{J}(\beta^{(k)})$
 3. arrêt lorsque le minimum est atteint

Itération k de l'algorithme de descente du gradient

Pour le paramètre β_j

$$\beta_j^{(k)} := \beta_j^{(k-1)} - \alpha \frac{\partial}{\partial \beta_j} \mathbf{J}(\beta^{(k-1)})$$

avec :

- $\frac{\partial}{\partial \beta_j}$:
- α :

- Objectif : trouver le minimum d'une fonction-coût
- Principe : algorithme itératif
 1. initialisation : $\beta^{(0)}$
 2. à chaque étape k , modifier $\beta^{(k-1)}$ pour faire diminuer $\mathbf{J}(\beta^{(k)})$
 3. arrêt lorsque le minimum est atteint

Itération k de l'algorithme de descente du gradient

Pour le paramètre β_j

$$\beta_j^{(k)} := \beta_j^{(k-1)} - \alpha \frac{\partial}{\partial \beta_j} \mathbf{J}(\beta^{(k-1)})$$

avec :

- $\frac{\partial}{\partial \beta_j}$: dérivée partielle
- α : pas d'apprentissage (à fixer par l'utilisateur)

- Itération k de l'algorithme de descente du gradient?

Itération k de l'algorithme de descente du gradient - régression logistique

$$\beta_j^{(k)} := \beta_j^{(k-1)} - \alpha \frac{1}{m} \sum_{i=1}^m \left(f_{\beta^{(k-1)}}(\mathbf{x}_i) - y_i \right) x_i^j$$

Remarque : $\forall i, x_i^0 = 1$

- Itération k de l'algorithme de descente du gradient?
- Même formule que pour la régression linéaire...

Itération k de l'algorithme de descente du gradient - régression logistique

$$\beta_j^{(k)} := \beta_j^{(k-1)} - \alpha \frac{1}{m} \sum_{i=1}^m \left(f_{\beta^{(k-1)}}(\mathbf{x}_i) - y_i \right) x_i^j$$

Remarque : $\forall i, x_i^0 = 1$

Il existe des versions plus avancées de l'algorithme de descente du gradient :

- gradient conjugué, BFGS¹, L-BFGS (à mémoire limitée)
 - pas besoin de fixer α (fait automatiquement dans l'algorithme)
 - plus rapide (moins d'itérations) que l'algorithme de descente de gradient
 - **mais** plus complexe à implémenter (ces variantes sont disponibles dans les principaux logiciels/langages)

1. Broyden-Fletcher-Goldfarb-Shanno

Introduction

Le modèle logistique

Loi de Bernoulli

Fonction logistique

Algorithme de descente du gradient pour la régression logistique

Rappel

Pour la régression logistique

Optimisation

Application

Application

Les données portent sur une étude clinique sur le cancer du sein menée à l'Université du Wisconsin. Il s'agit de prévoir le statut de la tumeur (maligne-1 ou bénigne-0) à partir de caractéristiques de cellules prélevées chez les patientes

- radius (mean of distances from center to points on the perimeter)
- texture (standard deviation of gray-scale values)
- perimeter
- area
- ...

Objectif : étudier le lien entre la probabilité d'avoir une tumeur maligne en fonction de certaines caractéristiques

The data was originally published by W.N. Street, W.H. Wolberg and O.L. Mangasarian. Nuclear feature extraction for breast tumor diagnosis. IS&T/SPIE 1993 International Symposium on Electronic Imaging : Science and Technology, volume 1905, pages 861-870, San Jose, CA, 1993.

- Importer les données
- Construire les fonctions nécessaires : f , coût, gradient puis implémenter l'algorithme de descente de gradient
 - Cas de la régression logistique simple
 - Etendre à la régression logistique multiple
 - Retourner les coefficients du modèle et les valeurs de la fonction-coût pour toutes les itérations
- **Faire varier les paramètres de l'algorithme (initialisation, pas) et commenter.**
- **Comparer avec les solutions d'optimisation natives de Python (pas de solution exacte ici)**