



INSTITUT
DE LA
communication

Université Lumière Lyon 2

Librairie MIMOSA

Mixed Input Multinomial Optimization for Statistical Analysis

*Régression logistique multinomiale sur données mixtes,
par descente de gradient*

Par

Linh Nhi Le Dinh
Antoine Oruezabala
Béranger Thomas

Rapport présenté dans le cadre du
Master 2 SISE - Promotion 2024/2025

Table des matières

1	Théorie	1
1.1	La régression logistique multinomiale	1
1.1.1	Principes Généraux	1
1.1.2	Modélisation mathématique[makisreglog]	2
1.2	La descente de gradient	2
1.2.1	Fonction de Coût	3
1.3	L'algorithme de descente de gradient	4
1.3.1	Modélisation	4
1.3.2	Test	4
2	Implémentation Pratique en R	5
2.1	Structure du programme	5
2.1.1	main.r	5
2.1.2	utils.r	5
2.1.3	preprocessing.r	6
2.1.4	LogisticRegression.r	6
2.1.5	Répertoire docs	7
2.2	Vérification des données	7
2.2.1	Format des objets passés en argument	7
2.2.2	Format des données	7
2.2.3	Cohérence des données	7
2.3	Pré-traitement des données	7
2.3.1	Variable cible y	7
2.3.2	Variables explicatives X	7
3	Éléments introductifs	9
3.1	La régression logistique	9
3.1.1	Notations	9
3.1.2	Algorithmes	10
3.1.3	Champs d'application	11
4	Lexique	12

Résumé

Ce rapport détaille l'implémentation de la régression logistique multinomiale pour des données mixtes, en langage R, et en utilisant la descente de gradient. Nous couvrons tout d'abord quelques fondements théoriques, puis l'approche pratique et enfin discutons des résultats obtenus. En fin d'ouvrage, une bibliographie donnera quelques pistes au lecteur curieux d'approfondir certains aspects, ou désireux d'apprendre par le biais d'une pédagogie quelque fois différente de celle de ce rapport. Un lexique vient également compléter nos écrits et donne accès à des définitions et formules mathématiques de façon plus condensés que dans le corps du texte.

Ce rapport a été rédigé en novembre 2024 dans le cadre d'un projet du master 2 SISE.

Chapitre 1

Théorie

1.1 La régression logistique multinomiale

1.1.1 Principes Généraux

La régression logistique est un modèle statistique permettant d'étudier les relations entre un ensemble de variables X_i et une variable qualitative y . Elle permet de prédire la probabilité qu'un événement arrive (valeur de 1) ou non (valeur de 0) à partir de l'optimisation des coefficients de régression.

On distingue trois types de régression logistique :

1. La régression logistique binaire. Elle vise à prédire la probabilité qu'un événement binaire se produise (oui/non, 1/0) en fonction de variables explicatives. Elle utilise une fonction logistique (ou sigmoïde) pour transformer la combinaison linéaire des variables explicatives en une probabilité comprise entre 0 et 1.
2. La régression logistique **multinomiale** est une généralisation de la régression logistique binaire. Elle permet d'estimer la probabilité d'appartenance à chacune des catégories en fonction de variables explicatives quantitatives et/ou qualitatives. Contrairement à la régression logistique binaire qui modélise une seule probabilité p , la régression logistique multinomiale estime simultanément les probabilités d'appartenance à toutes les catégories (p_1, p_2, \dots, p_K) , avec la contrainte que leur somme soit égale à 1.
3. La régression logistique ordinale. Ce type de modèle de régression logistique est utilisé lorsque la variable de réponse a trois modalités possibles ou plus et que ces dernières ont un ordre défini.

Ce document ne traitera sauf mention contraire que la régression logistique multinomiale nominale, c'est-à-dire que l'on ne suppose aucun ordre particulier entre les modalités de la variable à prédire.

1.1.2 Modélisation mathématique[makisreglog]

Formulation du problème

Soit :

- Y une variable catégorielle cible prenant K valeurs distinctes : $\{1, 2, \dots, K\}$,
- X est un vecteur des p variables explicatives pour une observation donnée,
- β_k est un vecteur des coefficients associés à la classe k ,
- $P(Y = k \mid X = \mathbf{x})$ est la probabilité conditionnelle pour chaque classe $k \in \{1, 2, \dots, K\}$.

Colinéarité

Pour éviter les problèmes de colinéarités, on choisit une classe de référence que l'on exclut de la modélisation.

Les probabilités des $k - 1$ classes se calculent avec :

$$P_k = \frac{e^{C_k}}{1 + \sum_{k=1}^{K-1} e^{C_k}} \quad (1.1)$$

Avec C_k le logit pour la modalité k de Y ; et comme indiqué précédemment $k \in \{1, 2, \dots, K - 1\}$.

La probabilité pour la classe de référence K peut être déduite après coup :

$$P_K = 1 - \sum_{k=1}^{K-1} P_k \quad (1.2)$$

1.2 La descente de gradient

L'estimation des coefficients β peut se faire par descente de gradient.

La descente de gradient est un algorithme d'optimisation itératif, qui permet de trouver le minimum d'une fonction.

L'idée est la suivante :

- on part d'un point choisi aléatoirement,
- depuis ce point nous cherchons le gradient (la pente),
- et nous descendons d'un pas dans la direction opposé au gradient,
- nous recommençons le processus à partir du nouveau point atteint.

Les critères pour s'arrêter peuvent être :

- le nombre d'itérations ; on définit un maximum au début de l'algorithme, que l'on ne dépasse pas,
- la différence de "niveau" entre deux points ; si elle passe en dessous d'un seuil très faible, l'approximation est jugée suffisamment bonne,
- une combinaison de ces deux critères, qui a l'avantage de s'arrêter le plus tôt possible.

Vocabulaire :

- la taille du pas est dénommé taux d'apprentissage (en anglais learning rate),

— la différence minimale souhaitée entre deux points consécutifs est la tolérance,

—
—

Mathématiquement, la pente est définie par le gradient. La taille du pas est dénommé taux d'apprentissage.

Il existe trois types d'algorithmes d'apprentissage par descente de gradient : la descente de gradient simple (sur tous les points), la descente de gradient stochastique et la descente de gradient par mini-lots.

1.2.1 Fonction de Coût

La fonction de coût utilisée pour la régression logistique multinomiale est la cross-entropie catégorielle. Elle évalue la divergence entre les probabilités prédites par le modèle et les probabilités réelles des classes. Elle est définie comme suit :

$$J(\beta) = -\frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K y_{i,k} \log (P(Y = k | X_i))$$

où :

- N est le nombre total d'observations. - K est le nombre de classes. - $y_{i,k}$ est une variable indicatrice qui vaut 1 si l'observation i appartient à la classe k , et 0 sinon. - $P(Y = k | X_i)$ est la probabilité prédite que l'observation i appartienne à la classe k .

Minimiser cette fonction de coût permet d'ajuster les paramètres β du modèle pour améliorer la précision des prédictions. La descente de gradient est une méthode d'optimisation utilisée pour trouver les valeurs optimales de β qui minimisent $J(\beta)$.

La fonction *SoftMax* est essentielle dans ce modèle car elle transforme les scores linéaires en probabilités. Elle garantit que les sorties du modèle sont comprises entre 0 et 1 et que la somme des probabilités pour toutes les classes est égale à 1 pour chaque observation. La fonction *SoftMax* est définie par :

$$P(Y = k | X) = \frac{\exp(X\beta_k)}{\sum_{j=1}^K \exp(X\beta_j)}$$

où K est le nombre total de classes. Cette transformation est basée sur l'hypothèse que les logits peuvent être exponentiés et normalisés pour obtenir des probabilités, ce qui est cohérent avec le principe de maximum d'entropie.

L'utilisation de la fonction *SoftMax* implique également que les classes sont mutuellement exclusives, ce qui signifie qu'une observation ne peut appartenir qu'à une seule classe. Cela correspond à la nature des problèmes de classification multinomiale où le but est de prédire une seule catégorie parmi plusieurs possibles.

1.3 L'algorithme de descente de gradient

1.3.1 Modélisation

1.3.2 Test

Chapitre 2

Implémentation Pratique en R

2.1 Structure du programme

L'architecture du programme est la suivante :

```
projet_racine/  
├── datasets/  
├── docs/  
│   ├── Rapport.pdf  
│   └── Rapport.tex  
├── utils/  
│   ├── LogisticRegression.r  
│   ├── preprocessing.r  
│   └── utils.r  
└── main.r
```

Les fichiers sont décrits ci-après, avec du pseudo-code pour certaines fonctions.

2.1.1 main.r

Le fichier `main.r` contient le programme principal, qui instancie l'objet `LogisticRegression`, qui contient les fonctions suivantes :

- `fit` : pour modéliser
- `transform` : pour
- `probas`

2.1.2 utils.r

Le fichier `utils.r` contient une fonction `check_type`, qui vérifie si l'objet passé en premier argument est du même type que le type passé en second argument.

Si ce n'est pas le cas, le programme s'arrête.

2.1.3 preprocessing.r

Le fichier `preprocessing.r` contient plusieurs fonctions de prétraitement :

- `imputer_numeriques` : réalise une imputation par la moyenne des valeurs manquantes dans le jeu de données passé en argument.
- `imputer_categorielles` : réalise une imputation par le mode des valeurs manquantes dans le jeu de données passé en argument.
- `split_train_test` : réalise une séparation du jeu de données passé en argument

2.1.4 LogisticRegression.r

Le fichier `LogisticRegression.r` contient toutes les fonctions nécessaires à la réalisation de la régression logistique :

Fonction « sigmoid »

Applique la [fonction sigmoïde](#) à la valeur passée en argument et renvoie le résultat.

Fonction « softmax »

Applique la [fonction softmax](#) à la valeur passée en argument et renvoie le résultat.

Fonction « cout »

Fonction de coût pour la régression logistique, avec cross-entropy :

```
cout <- function(X, y, beta) {  
  n <- nrow(X)  
  scores <- X %*% beta  
  probs <- exp(scores) / rowSums(exp(scores))  
  cost <- -mean(sum(y * log(probs)))  
  return(cost)  
}
```

Fonction « hypothese »

Fonction « descente_gradient »

Applique une descente de gradient pour la régression logistique :

```
descente_gradient <- function(X, y, beta) {  
  n <- nrow(X)  
  scores <- X %*% beta  
  probs <- exp(scores) / rowSums(exp(scores))  
  gradient <- -t(X) %*% (y - probs) / n  
  return(gradient)  
}
```

2.1.5 Répertoire docs

Il contient :

- ce rapport au format pdf : `Rapport.pdf`,
- et au format Latex : `Rapport.tex`

2.2 Vérification des données

Afin de s'assurer de l'intégrité des données, nous procédons à plusieurs vérifications lors de l'instanciation de l'objet *LogisticRegression*.

2.2.1 Format des objets passés en argument

Les deux vérifications suivantes sont réalisées à l'aide de la fonction *check_type* du fichier *utils.r* :

1. les données explicatives X doivent être dans un `data.frame`
2. la cible y doit être un vecteur ou un caractère

2.2.2 Format des données

1. Les données de la cible y doivent être de type catégorielle (caractère ou `factor`).

2.2.3 Cohérence des données

1. La cible y doit être présente dans les données X ,
2. il faut au moins deux modalités pour la variable cible y .

2.3 Pré-traitement des données

2.3.1 Variable cible y

Si le format des données de la variable cible ne correspondaient pas aux pré-requis (variable catégorielle et avec au moins deux modalités), le programme l'aura signalé lors de l'instanciation de l'objet *LogisticRegression*.

Il n'y a pas d'autre contrainte à respecter pour la variable cible y , mais une opération de suppression des lignes sans la variable cible est menée ici. Le but étant bien évidemment de ne pas charger à modéliser sur des données manquantes pour la cible, ce qui n'aurait pas de sens puisque nous sommes en apprentissage supervisé.

2.3.2 Variables explicatives X

Les variables explicatives peuvent être mixtes, numériques et/ou catégorielles. Dans les deux cas, des pré-traitements sont nécessaires.

Variables numériques

Afin de pouvoir comparer des données aux unités quelques fois différentes, on centre et réduit les données numériques. Les coefficients seront alors normalisés ce qui permettra leur comparaison directe.

Variables catégorielles

Chaque variable catégorielle subit deux étapes pour les préparer à la descente de gradient :

1. un *codage disjonctif* (one-hot encoding) pour produire des valeurs numériques,
2. la suppression d'une des modalités.

Cette dernière étape permet d'éviter la colinéarité, c'est-à-dire la redondance d'information entre variables (une variable désignant ici une des colonnes produites par le one-hot encoding ; c'est-à-dire une des modalités de la variable encodée).

Chapitre 3

Éléments introductifs

3.1 La régression logistique

Quelle soit binaire (la variable cible possède deux modalités) ou multinomiale (la variable cible possède plus de deux modalités), le principe de la régression logistique est de maximiser la vraisemblance, ou pour le dire autrement, de minimiser la déviance.

Quel outil pour mesurer ce paramètre ? La log-vraisemblance.

3.1.1 Notations

Soit une population Ω de n individus, définie par J variables explicatives notées $\{X_1, \dots, X_J\}$, et une variable cible Y possédant K valeurs :

Ω	Cible	X_1	...	X_J
1	Y_1			
...	...			
ω	Y_ω	$X_1(\omega)$...	$X_J(\omega)$
...	...			
n	Y_K			

Dans le cas binaire, la probabilité d'un individu ω d'être positif à priori se note p par commodité, pour simplifier $p(\omega)$, lui-même une notation simplifiée de $P[Y(\omega) = +]$.

Toujours dans le cas binaire, la probabilité d'un individu ω d'être positif à posteriori, c'est-à-dire la probabilité que l'on modélisera en apprentissage supervisée, se note π par commodité, pour simplifier $\pi(\omega)$, lui-même une notation simplifiée de $P[Y(\omega) = +/X(\omega)]$.

La fonction LOGIT pour cet individu ω est :

$$\text{LOGIT}(\omega) = \ln \left[\frac{\pi(\omega)}{1 - \pi(\omega)} \right] = a_0 + a_1 X_1(\omega) + \dots + a_J X_J(\omega) \quad (3.1)$$

Soit en écriture matricielle :

$$\ln \left[\frac{\pi(\omega)}{1 - \pi(\omega)} \right] = X(\omega) \times a \quad (3.2)$$

Avec a_0, \dots, a_J les paramètres que l'on souhaite estimer.

3.1.2 Algorithmes

Un algorithme utilisé pour optimiser la log-vraisemblance est celui de Newton-Raphson. Il s'agit d'une approche itérative, qui utilise la dérivée de la fonction considérée pour approcher une solution.

Il a toutefois pour défaut d'utiliser la matrice hessienne, qui peut être coûteuse en mémoire et temps de calcul.

Nous expliquerons et utiliserons un autre algorithme ici, la descente de gradient.

Il s'agit aussi d'un algorithme itératif, dont l'objectif est de trouver le minimum d'une fonction $f(x)$ en partant d'un point arbitraire x_0 . On se « déplace » pour cela dans la direction opposée au gradient, c'est-à-dire que l'on « descend » le long de la pente.

Mathématiquement, on part d'une valeur arbitraire x_0 , puis pour trouver x_1 on utilise la formule suivante :

$$x_{n+1} = x_n - \alpha \nabla f(x_n) \quad (3.3)$$

Avec :

- x_n la position actuelle,
- α le taux d'apprentissage,
- $\nabla f(x_n)$ le gradient de la fonction au point x_n

Les étapes sont donc les suivantes :

- choisir un point x_0 arbitraire pour commencer les calculs,
- calculer le gradient à ce point,
- mettre à jour la position en se déplaçant dans la direction opposée,
- vérifier que l'on a atteint un critère d'arrêt :
 - nombre d'itérations maximum atteint,
 - et/ou précision souhaitée (différence entre deux itérations) atteinte,

3.1.3 Champs d'application

La régression logistique n'émet pas d'hypothèse directement sur les distributions des probabilités $P(X/Y = +)$ et $P(X/Y = -)$, mais uniquement sur leur rapport.

Elle permet donc, en théorie, une application plus large, et est nommé semi-paramétrique, pour la différencier des modèles qui supposent une loi donnée sur la distribution des probabilités.

Chapitre 4

Lexique

CENTRAGE ET RÉDUCTION ■ L'opération de centrage et réduction transforme les données pour avoir une moyenne de 0 et un écart-type de 1 :

$$x_{cr} = \frac{x_i - \bar{x}}{\sigma} \quad (4.1)$$

Centrer et réduire les données ne change pas les distributions, seulement l'échelle des données. Le centrage déplace le centre de la distribution au point zéro de l'axe des abscisses, et la réduction modifie l'échelle pour que l'écart-type soit égal à 1.

Il rend les variables indépendantes de leur unité de mesure initiale et permet donc de comparer des variables mesurées dans des unités ou sur des échelles différentes.

FONCTION SIGMOÏDE ■ La fonction sigmoïde ou logistique s'applique aux nombres réels et est utilisée pour transformer des valeurs réelles en probabilités. Elle produit une courbe en forme de S, avec des valeurs comprises entre 0 et 1.

Elle est l'inverse de la fonction LOGIT. Sa formule est :

$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad (4.2)$$

Où :

- $\sigma(x)$ représente la fonction sigmoïde
- e est la base du logarithme naturel (constante d'Euler)
- x est la variable d'entrée

Les deux fonctions logistique et LOGIT étant inverses l'une de l'autre, cela signifie que :

- $\text{logit}(\sigma(x)) = x$
- $\sigma(\text{logit}(p)) = p$

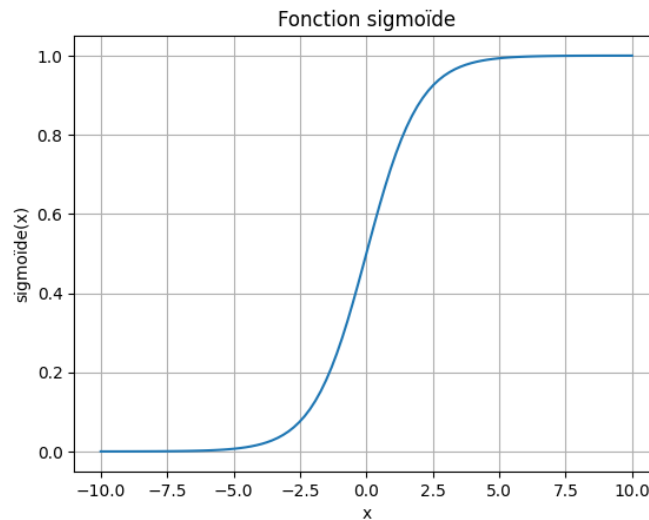


FIGURE 4.1 – La fonction sigmoïde

FONCTION SOFTMAX ■ La fonction softmax est une fonction d'activation qui transforme un vecteur de nombres réels en une distribution de probabilités.

MODALITÉ DE RÉFÉRENCE ■ Modalité de référence par rapport à laquelle le LOGIT est exprimé.

TAUX D'APPRENTISSAGE ■ Ou learning Rate : il fixe la n^{e} grandeur z du pas de chaque itération de la descente de gradient.

DESCENTE DE GRADIENT ■ La descente de gradient est un algorithme d'optimisation couramment utilisé pour entraîner les modèles de machine learning et les réseaux neuronaux. Ce type d'algorithme entraîne les modèles de machine learning par réduction des erreurs entre les résultats prédits et les résultats réels. ■ Le point de départ n'est qu'un point arbitraire qui nous permet d'évaluer les performances. À partir de ce point de départ, nous allons trouver la dérivée (ou la pente) et, à partir de là, nous pourrions utiliser une ligne tangente pour observer l'inclinaison de la pente. La pente renseigne sur les mises à jour des paramètres, c'est-à-dire les poids et les biais. La pente au point de départ est plus forte, mais au fur et à mesure que de nouveaux paramètres sont générés, elle devrait progressivement diminuer jusqu'à atteindre le point le plus bas de la courbe, dénommé point de convergence. ■ Comme pour trouver la ligne de meilleur ajustement dans la régression linéaire, l'objectif de la descente de gradient est de minimiser la fonction de coût, ou l'erreur entre y prédit et y réel. Pour ce faire, deux points de données sont nécessaires : une orientation et un taux d'apprentissage. Ces facteurs déterminent les calculs de dérivée partielle des itérations futures, ce qui lui permet d'atteindre progressivement le minimum local ou global (c'est-à-dire le point de convergence). La fonction de perte dans la descente de gradient agit spécifiquement comme un baromètre, évaluant sa précision à chaque itération des mises à jour de paramètres. Jusqu'à ce que la fonction soit proche ou égale à zéro, le modèle continue à ajuster ses paramètres pour obtenir l'erreur la plus faible possible. ■ Il existe trois types d'algorithmes d'apprentissage par descente de gradient :

la descente de gradient par lots, la descente de gradient stochastique et la descente de gradient par mini-lots.

Table des figures

4.1 La fonction sigmoïde	13
------------------------------------	----