

# PAO : machine learning et clicks prediction

Octave QUEFFELEC et Sandratra RASENDRASOA

11/06/2017

## Contents

<b>I</b>	<b>Introduction</b>	<b>3</b>
<b>II</b>	<b>Partie théorique</b>	<b>4</b>
1	Supervised learning	4
2	Regression Lineaire	5
3	Regression Logistique	6
4	Fonction de coût	7
5	Perceptron	8
5.1	Descente de gradient . . . . .	9
<b>III</b>	<b>Implémentation</b>	<b>10</b>
<b>6</b>	<b>Présentation du jeu de données</b>	<b>10</b>
6.1	Présentation générale . . . . .	10
6.2	Zoom sur la représentation réelle des attributs . . . . .	12
<b>7</b>	<b>Présentation des différents fichiers/classes</b>	<b>13</b>
<b>8</b>	<b>Warm-up</b>	<b>13</b>
8.1	Taux de clic moyen . . . . .	13
8.2	Nombre de tokens unique dans les données d'apprentissage . . . .	14
<b>9</b>	<b>Implémentation du perceptron</b>	<b>14</b>
9.1	Matlab . . . . .	14
9.2	Python . . . . .	15

<b>10 Descente de gradient stochastique (SGD)</b>	<b>17</b>
<b>IV Evaluation du modèle</b>	<b>19</b>
10.1 Matrice de confusion . . . . .	19
10.1.1 Apprentissage simple . . . . .	19
10.1.2 Régularisation . . . . .	20
10.1.3 LogRegression . . . . .	20
10.1.4 Tirage aléatoire . . . . .	20
10.2 Courbe d'erreur moyenne . . . . .	21
<b>11 Interprétation des résultats</b>	<b>21</b>
<b>V Répartition des tâches</b>	<b>22</b>
<b>VI Conclusion</b>	<b>22</b>
<b>VII Bibliographie</b>	<b>23</b>

## List of Figures

1	Apprentissage supervisé VS non-supervisé . . . . .	5
2	Prix des maisons dans l'état d'Oregon . . . . .	5
3	Modélisation du jeu de données . . . . .	6
4	Fonction logistique associée à $D(x)$ . . . . .	7
5	Perceptron . . . . .	8
6	Algorithme Perceptron . . . . .	9
7	Descente de gradient . . . . .	10
8	Exemple de session . . . . .	11
9	Système OR . . . . .	15
10	Système AND . . . . .	16
11	2 classes séparables . . . . .	16
12	plusieurs solutions pour 2 classes séparables . . . . .	16
13	Erreur moyenne en fonction du nombre de pas . . . . .	21
14	The road to becoming a data scientist . . . . .	23

## List of Algorithms

1	taux de clic moyen . . . . .	13
2	Nombre de tokens unique . . . . .	14
3	perceptron sous matlab . . . . .	15
4	perceptron en python . . . . .	15

5	Classe Weights . . . . .	17
6	simple train . . . . .	18
9	logreg . . . . .	19
7	regularized train . . . . .	19
8	sigmoïde . . . . .	19

## Part I

# Introduction

“We are drowning in information and starving for knowledge.”  
Rutherford D. Roger<sup>1</sup>

January 2009 Hal Varian, Google’s Chief Economist, tells the McKinsey Quarterly: “I keep saying the sexy job in the next ten years will be statisticians. People think I’m joking, but who would’ve guessed that computer engineers would’ve been the sexy job of the 1990s? The ability to take data—to be able to understand it, to process it, to extract value from it, to visualize it, to communicate it—that’s going to be a hugely important skill in the next decades...”<sup>2</sup>

Avec l’avènement des technologies de l’information, le Data Science (ou science des données) est une discipline ayant pour objectif de répondre au besoin d’analyser des immenses ensembles de données pour en extraire des informations -potentiellement- utiles. Les applications sont nombreuses, que ce soit dans le secteur privé (la recommandation de séries par Netflix, les moteurs de recherche, l’analyse financière), le domaine médical (l’aide aux diagnostics) et bien d’autres encore.

Le but de notre PAO est de s’initier aux méthodes et algorithmes issus du Machine Learning. Pour cela, nous avons à notre disposition un jeu de données issu de la base de données en ligne Kaggle, qui a été créé pour répondre à la problématique suivante : déterminer un modèle pour estimer le taux de clicks(ou CTR) qui est une unité de mesure importante pour évaluer les publicités en ligne. Nous avons utilisé comme base un projet de l’université de Washington, en réalisant la partie pratique du projet et en ayant une réunion toutes les deux semaines avec nos responsables pédagogiques M.Gasso et M.Gauzere pour aborder ces différents points et les parties théoriques s’y référant.

L’ensemble du projet a été réalisé en langage Python qui est un des langages proposé dans le sujet de l’université de Washington. Nous avons notamment

---

<sup>1</sup>T.Hastie, R. Tibshirani, J. Friedman, The Elements of Statistical Learning: Data Mining, Inference, and Prediction (Springer Verlag)

<sup>2</sup>Gil Press, “A Very Short History Of Data Science”,Forbes

utilisé les librairies Numpy et Scipy pour nos algorithmes. Nous avons aussi dans un premier temps réalisé quelques algorithmes sous Matlab pour se familiariser avec les jeux de données. Pour faciliter notre travail et permettre un suivi régulier, nous avons utilisé le gestionnaire de source de la plateforme MonProjet.

Dans ce rapport, nous allons vous présenter tout d'abord les éléments théoriques nécessaires pour ce projet, puis l'ensemble des travaux pratiques que nous avons réalisé ainsi que les résultats qui y sont liés. Nous terminerons par une interprétation de ces résultats et une conclusion globale.

## Part II

# Partie théorique

### 1 Supervised learning

Il existe plusieurs algorithmes d'apprentissage automatique(ou machine learning) qui se différencient par le mode d'apprentissage qu'ils utilisent. Dans le cadre de cette initiation, nous ne verrons qu'un mode d'apprentissage, mais il nous semble important de mentionner que d'autres techniques existent pour répondre aux différents besoins en machine learning.

L'apprentissage supervisé (ou supervised learning) est la technique d'apprentissage la plus commune pour répondre à des problèmes de machine learning. L'objectif est le suivant :

A partir des données  $(x_i, y_i) \in X \times Y, i = \dots, N$ , estimer les dépendances entre X et Y.

On parle ici d'apprentissage supervisé (en opposition à l'apprentissage non-supervisé) car les  $y_i$  (en pratique, il s'agira des cas déjà traités et validés) permettent de guider le processus d'estimation. Un exemple concret pour différencier l'apprentissage supervisé de l'apprentissage non-supervisé est présenté dans la figure 1.

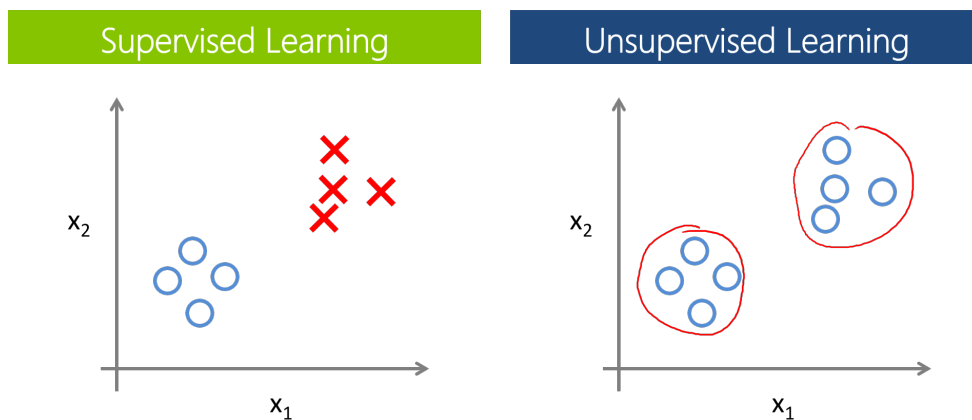


Figure 1: Apprentissage supervisé VS non-supervisé

En apprentissage supervisé, on sait d'avance qu'il y a exactement deux catégories (cercle bleu et croix rouge), tandis qu'en apprentissage non supervisé, on essaie de partitionner et classer les données dans des groupes homogènes (en anglais, on parle aussi de clustering). Dans le cadre de notre projet, nous travaillons en apprentissage supervisé, nous possédons ainsi plusieurs jeux de données correspondants à  $X$  et  $Y$ , que nous préciserons dans la partie pratique.

## 2 Regression Lineaire

La regression lineaire est une des méthodes d'apprentissage supervisé la plus simple. Une base de données d'apprentissage est un ensemble de couples entrée-sortie  $(x_i, y_i)$ ,  $x \in X, y \in Y, i = 1..n$  qui suit une loi inconnue. Un apprentissage est supervisé lorsque l'on connaît les sorties  $y_i$  et qu'elles nous permettent de guider le processus d'estimation. Le but d'un algorithme d'apprentissage supervisé est donc de généraliser pour des entrées inconnues ce qu'il a pu « apprendre » grâce aux données déjà traitées, ceci de façon « raisonnable ». On dit que la fonction de prédiction apprise doit avoir de bonnes garanties en généralisation. Le but d'un problème de regression est de prédire grâce à cette fonction de prédiction la sortie  $y$  d'une entrée  $x$  dont on ne connaît pas par avance la véritable sortie.

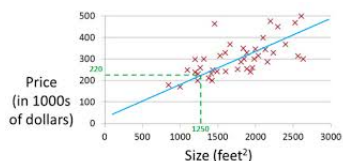


Figure 2: Prix des maisons dans l'état d'Oregon

Voici un graphe avec le prix des maisons dans l'état d'Oregon en fonction de leur superficie en mètre carré. Le dataSet nous est donné. Le problème posé est le suivant : pour une superficie donnée, quel serait le prix le plus adéquat de la maison ?

Pour cela, il nous faut modéliser ce jeu de données. On pourrait par exemple tracer une droite en minimisant chaque distance entre la droite et un point des données pour représenter le modèle au mieux.

Pour ce faire, on dispose d'un training set, un ensemble de couple  $(x_i, y_i)$ .

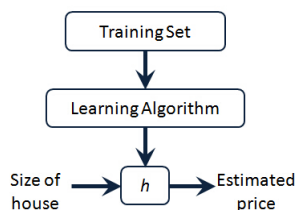


Figure 3: Modélisation du jeu de données

A l'aide du training set, l'algorithme d'apprentissage doit déterminer une fonction  $h$ . Cette fonction doit prédire un  $y$  pour chaque entrée  $x$  passée en entrée de la fonction.

Généralement, cette fonction est définie dans l'absolu comme :

$$h(x) = \langle w, x \rangle, x \in \mathbb{R}^n, w \in \mathbb{R}^n$$

Dans notre cas de régression linéaire, on aurait quelque chose de la forme :

$$h(x) = w_0 + w_1 * x, x \in \mathbb{R}^2, w \in \mathbb{R}^2$$

En somme, il "suffit" que notre algo détermine correctement les coefficients  $w_0$  et  $w_1$ .

### 3 Régression Logistique

La régression logistique est un cas particulier de la régression linéaire, c'est un modèle de régression binomiale. Ainsi, la variable de sortie prends 2 modalités :  $y = 0, 1$

On cherche à prédire une probabilité d'appartenance à une classe ou non.

*Règle de décision de Bayes :*

$$D(x) = C1 \text{ si } P(C1|X)/P(C2|X) > 1, C2 \text{ sinon}$$

Avec  $P(C1|X) = \exp(wTx)/(1 + \exp(wTx))$ . La fonction logistique associée donne donc :

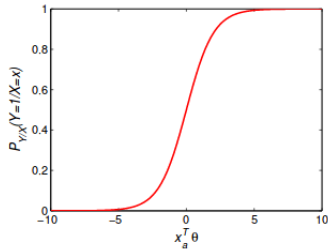


Figure 4: Fonction logistique associée à  $D(x)$

## 4 Fonction de coût

En apprentissage supervisé, l'objectif est de trouver une fonction  $f : X \rightarrow Y$  qui permet d'estimer la valeur  $y$  associée à  $x$ . On peut alors se poser la question suivante : la fonction  $f$  que nous avons construite est-elle optimale pour résoudre notre problème ?

On introduit alors la notion de coût  $L(Y, f(X))$ , dont l'objectif est d'évaluer la pertinence de la prédiction réalisée par  $f$ , et de pénaliser les erreurs. Dans l'idéal, on chercherait donc une fonction  $f$  qui prédit au mieux  $y$ ; en d'autres termes, une fonction  $f$  qui minimise l'erreur entre la vérité et la valeur estimée. D'où la fonction  $f$  telle que :

$$R(f) = E_{X,Y}[L(Y, f(X))]$$

$R(f)$  où  $R$  est appelé le risque moyen ou erreur de généralisation.

Un exemple de fonction de coût et de risque moyen associé est le coût quadratique (ou moindres carrés) :

$$L(Y, f(X)) = (Y - f(X))^2$$

$$R(f) = E[(Y - f(X))^2] = \int (y - f(x))^2 p(x, y) dx dy$$

La fonction de coût est associée à l'ensemble  $Y$  dont on souhaite prédire les valeurs, et on peut distinguer deux types de problèmes en fonction de cet ensemble  $Y$ .

1. Si l'ensemble  $Y$  est un sous-ensemble de  $R^d$  (donc continu), on traitera le problème comme étant un problème dit de régression, où les moindres carrés sont la fonction de coût usuelle.
2. Dans le cadre de notre projet,  $Y$  est un ensemble discret non-ordonné (clic ou non clic, 0 ou 1 dans le jeu de données), il s'agit d'un problème dit de classification. Ici, on cherchera plutôt à approcher la fonction de coût suivante :  $\theta(-yf(x))$  où  $\theta$  est la fonction échelon.

## 5 Perceptron

Le perceptron est un classifieur lineaire de la forme :

$$f(x) = \begin{cases} 1 & \text{si } \langle w, x \rangle + b > 0 \\ 0 & \text{sinon} \end{cases}, w \in \mathbb{R}^n$$

$w$  le vecteur des poids,  $b \in \mathbb{R}$ , le biais,  $\langle w, x \rangle$ , le produit scalaire entre le vecteur de poids  $w$  et le vecteur d'entrée  $x$ .  $\langle w, x \rangle + b = 0$  est alors l'équation de l'hyperplan affine qui sépare l'espace en 2 classes. Le but est donc d'entraîner notre vecteur de poids  $w$  afin de connaître par la suite notre fonction  $f$  qui nous permettra de prédire la classe d'un vecteur d'entrée  $x$  donné.

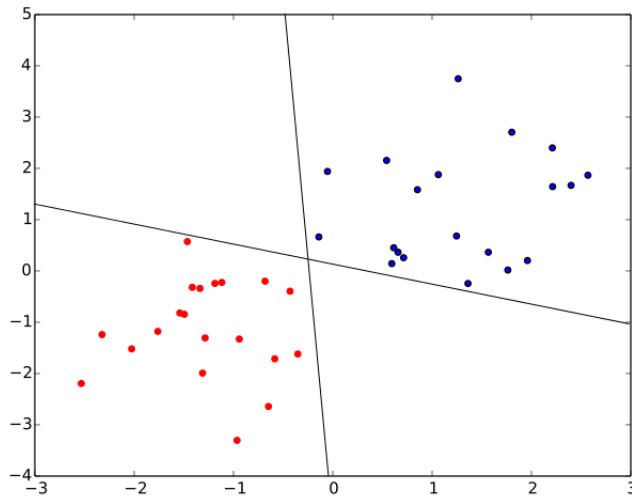


Figure 5: Perceptron

L'équation séparatrice des 2 classes n'est pas unique. Ici, 2 droites sont quasi-perpendiculaires l'une à l'autre mais le perceptron n'a aucun moyen d'en privilégier l'une au profit de l'autre.

Voici l'algorithme du perceptron, qui nous permet de mettre à jour à chaque iteration les poids  $w$  :



**Algorithme d'apprentissage du Perceptron**  
**Entrée :**  $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_l, y_l)\}$ , un échantillon complété linéairement séparable de  $\mathbb{R}^{n+1} \times \{-1, 1\}$   
 $\mathbf{w}_0 = \mathbf{0} \in \mathbb{R}^{n+1}, k = 0$   
**Répéter**  
    **Pour**  $i = 1$  à  $l$   
        **Si**  $y_i \langle \mathbf{w}_k, \mathbf{x}_i \rangle \leq 0$  **alors**  
             $\mathbf{w}_{k+1} = \mathbf{w}_k + y_i \mathbf{x}_i$   
             $k = k + 1$   
    **FinPour**  
**Jusqu'à ce qu'il n'y ait plus d'erreurs**  
**Sortie :**  $\mathbf{w}_k$

Figure 6: Algorithme Perceptron

Le perceptron peut être vu comme un type de réseau de neurones simplifié. Nous avons donc voulu pour nous entraîner, implémenter en matlab puis en python l'algorithme.

## 5.1 Descente de gradient

De manière générale, la descente de gradient est un algorithme cherchant à minimiser une fonction  $J(\omega_0, \omega_1)$ . Il pourrait donc dans notre cas s'agir de la fonction de coût que nous utilisons en régression logistique. Le principe général est le suivant : Soit une fonction

$$J(\omega_0, \omega_1)$$

on veut

$$\min_{\omega_0, \omega_1} (J(\omega_0, \omega_1))$$

Étapes:

1. on débute à  $\omega_0, \omega_1$
2. on modifie les paramètres  $\omega_0, \omega_1$  pour minimiser notre fonction et arriver à un minimum.

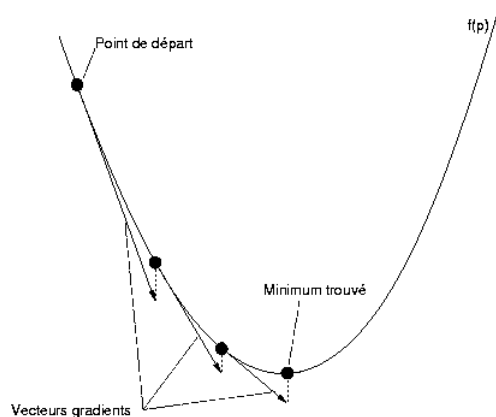


Figure 7: Descente de gradient

## Part III

# Implémentation

Le CTR (ou taux de clic) est une mesure cohérente pour évaluer la popularité d'une publicité. En pratique, notre objectif sera d'entraîner notre modèle pour espérer prédire le CTR sur un jeu de données d'environ un million d'exemples.

## 6 Présentation du jeu de données

### 6.1 Présentation générale

Les données que nous manipulons sont issues de la compétition 2012 KDD Cup Track 2. Nous avons à notre disposition 3 fichiers CSV : un fichier pour l'apprentissage "train.txt" et deux fichiers pour la validation "test.txt" et "test\_label.txt". Pour mieux visualiser les données, prenons l'exemple d'un utilisateur ayant réalisé une requête sur moteur de recherche et nous allons observer sur quelles publicités l'utilisateur a cliqué :

1. Nicolas réalise une recherche sur le "très utilisé" moteur de recherche Bing, où il a écrit la requête "perruque".
2. En plus des résultats, Bing affiche plusieurs publicités contenant des images et du texte (description, titre, etc ).
3. Nicolas clique alors sur la première publicité.

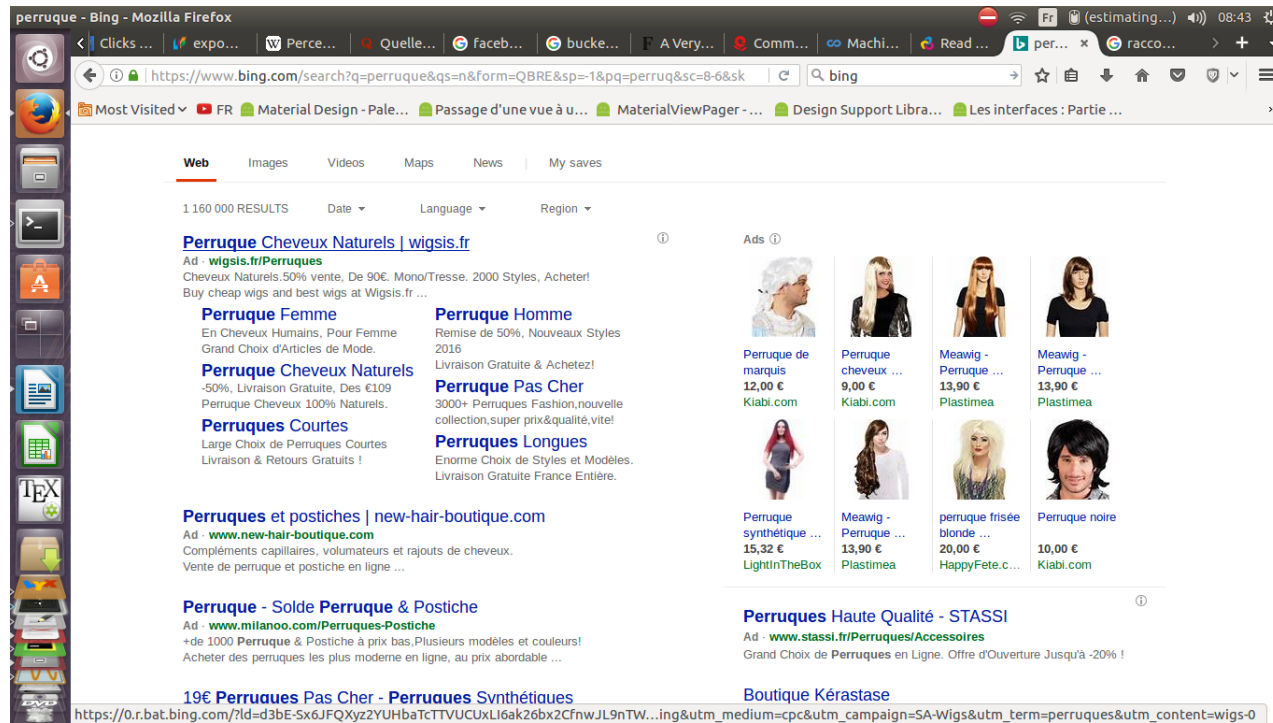


Figure 8: Exemple de session

Ces 3 étapes correspondent à une session. A la fin de cette session, Bing a recueilli un certain nombre d'enregistrement :

*Clicked = 1|Depth = 10|Position = 1|Nicolas|TextAd1*  
*Clicked = 0|Depth = 10|Position = 2|Nicolas|TextAd2*  
*Clicked = 0|Depth = 10|Position = 3|Nicolas|TextAd3*  
*etc...*

En plus de ces d'enregistrements, Bing possède d'autres informations sur l'utilisateur. Le format complet d'une ligne du jeu de données est présenté ici :

Clicked | Depth | Position | Userid | Gender | Age | Text Tokens  
of Ad

Les caractéristiques sont délibérément laissées en anglais afin de rester cohérent avec le code présenté ci-après.

Nous allons maintenant présenter chacune des caractéristiques :

**Clicked:** 0 ou 1, indique si l'utilisateur a cliqué ou non sur la publicité (0 non clic et 1 clic).

**Depth:** prend une valeur dans  $\{1, 2, \dots\}$ , indique le nombre de publicités affichées dans chaque session.

**Position:** prend une valeur dans  $\{1, 2, \dots, Depth\}$ , indique le rang/position de la publicité parmi les publicités affichées.

**Userid:** l'identifiant d'un utilisateur.

**Gender:** le genre  $\{-1, 0, 1\}$  d'un utilisateur : -1 homme, 1 femme et 0 inconnu.

**Age:** prend une valeur dans  $\{0, 1, 2, 3, 4, 5, 6\}$ , indique la tranche d'âge d'un utilisateur : 0 si inconnu, 1 si entre (0, 12], 2 si entre (12, 18], 3 si entre (18, 24], 4 si entre (24, 30], 5 si entre (30, 40] et 6 si au-delà de 40.

**TextTokensofAd:** une liste séparée par des virgules des "clés" de chaque mot d'une publicité. Par exemple "9,30,151" correspond aux clés des mots "mot#9", "mot#30", "mot#151" du dictionnaire utilisé (la correspondance avec des mots réels n'est pas réalisée pour des soucis de confidentialité).

Voici donc un échantillon complet d'une ligne du jeu de données "train.txt" :

0|2|1|490234|1|3|0,133,1374,164,17005,1891,208,211,246,2527,3917,459,5705,625,850,8726,888,906,93

Comme on peut le constater, l'ensemble des attributs correspond ici à des scalaires pour faciliter le traitement. La seule différence entre le fichier d'apprentissage et de validation est que la première caractéristique du fichier de validation "test.txt" est séparée du reste du fichier et enregistrée dans le fichier "test\_label.txt".

## 6.2 Zoom sur la représentation réelle des attributs

De manière théorique, il est facile représenter le vecteur des attributs :

$$x^t = [x_1^t, \dots, x_d^t]$$

où chaque attribut est un élément du vecteur  $x^t$ .

En pratique, la construction de ce vecteur nécessite plusieurs prérequis :

1. L'attribut "Userid" ne sera pas considéré comme une caractéristique du vecteur; il ne sera donc pas pris en compte.
2. Les nombres présents dans la liste des tokens étant assignés de manière arbitraire, il est inutile de les considérer dans notre modèle. Il est conseillé de visualiser la liste des tokens  $L = l_1, l_2, \dots$  comme étant la représentation compacte d'un vecteur creux où  $b(i) = 1 \forall i \in L$ . Ainsi les valeurs non nulles de ce vecteur correspondent aux clés des mots utilisés pour décrire une publicité.
3. Le reste des attributs sera utilisé en tant que tel.

## 7 Présentation des différents fichiers/classes

L'énoncé à mis à disposition un package avec différentes classes python afin de gérer assez facilement l'accès aux données.

**BasicAnalysis.py** Permet de faire le warm up, on y calcule entre autre les unique users, unique tokens, ctr...

**DataInstance.py** La classe DataInstance sectionne une instance des données en plusieurs attributs qui sont les différentes features : clicked, depth, position, id, gender, age, tokens... Cela facilite par la suite l'accès aux features.

**DataSet.py** La classe DataSet permet de charger un jeu de données et de naviguer entre les différentes instances à l'aide de fonction comme : hasNext(), nextInstance(), reset()...

**DummyLoader.py** Sert à vérifier si les données chargent correctement. Il print un nombre défini d'instances.

**LogisticRegression.py** Permet d'entraîner les poids grâce au training set et de prédire sur le testing set.

## 8 Warm-up

Notre premier traitement sur les données correspondra à plusieurs opérations simples, pour s'assurer qu'on peut correctement accéder et manipuler les données. En somme, les notions utilisés ici sont la manipulation des fichiers CSV en Python et l'utilisation basique des fonctionnalités de la librairie Numpy. Nous profitons des opérations déjà fournis avec le template de l'université de Washington pour accéder aux attributs qui nous intéressent.

### 8.1 Taux de clic moyen

---

**Algorithm 1** taux de clic moyen

---

```
def average_ctr(self, dataset):
    temp = 0
    x = []
    while dataset.hasNext():
        temp = dataset.nextInstance().clicked
        x.append(temp)
    return numpy.sum(x)/dataset.size
```

---

Résultat : *Average CTR = 3.36552848438 %*

## 8.2 Nombre de tokens unique dans les données d'apprentissage

---

**Algorithm 2** Nombre de tokens unique

---

```
def uniq_tokens(self, dataset):
    X = []
    instance = dataset.nextInstance()
    while dataset.hasNext():
        y = instance.tokens
        # tokens uniques de la ligne i
        x = np.unique(y)
        X.append(x)
        instance = dataset.nextInstance()
    X = np.array(X)
    # tokens uniques de X = [x1' x2' ... xi' ... xn']
    X = np.unique(X)
    return X
```

---

Résultat : *there are 24 unique tokens*

$X_{unique} = [1\ 3\ 26\ 144\ 149\ 178\ 255\ 440\ 466\ 582\ 685\ 771\ 1181$   
 $1691\ 1766\ 2854\ 5838\ 5977\ 6800\ 8511\ 8833\ 9930\ 10948\ 15645]$

## 9 Implémentation du perceptron

### 9.1 Matlab

Nous avons implementé différents systèmes assez simple.  
Pour chaque système, nous avons gardé la même notation :

1.  $l$ : le nombre d'entrée
2.  $n$ : la dimension des entrées
3.  $X = [x_1, x_2, \dots, x_l]$ , on concatène les  $l$  données dans  $X$
4.  $x_i = [a_1, a_2, \dots, a_n, 1]$  on concatène le biais à la fin de chaque vecteur d'entrée
5.  $y = [1, -1, 1, \dots] \in \mathbb{R}^l$

---

**Algorithm 3** perceptron sous matlab

---

```
w1=zeros(n+1,1);
k=0;
K=67;
for compteur=1:K
    indice=randperm(l);
    for i=1:l i=indice(i);
        if(y(i)*(x(i,:)*w1)<=0)
            w1=w1+y(i)*x(i,:);
            k=k+1;
        end
    end
end
end
```

---

## 9.2 Python

Nous avons ensuite implementé le perceptron en python, afin de nous familiariser avec les librairies numpy, math, et matplotlib. Les résultats obtenus sont identiques à ceux en matlab.

---

**Algorithm 4** perceptron en python

---

```
l = np.size(x[:,1])
for count in range(1,k):
    indice = np.random.permutation(range(l))
    for i in range(1,l):
        i = indice[i]
        if(np.dot(y[i],np.inner(x[i,:],w))<=0):
            w += np.dot(y[i],x[i,:])
```

---

Résultats pour les différents systèmes :

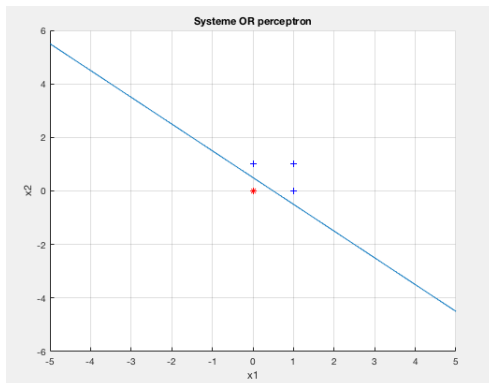


Figure 9: Système OR

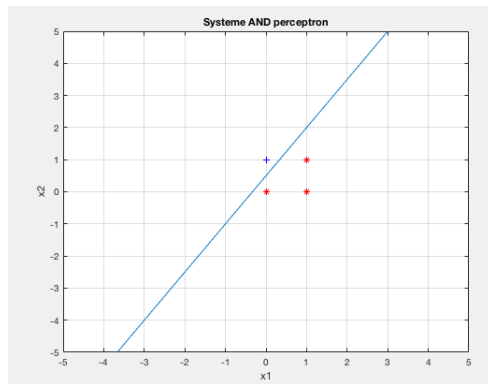


Figure 10: Système AND

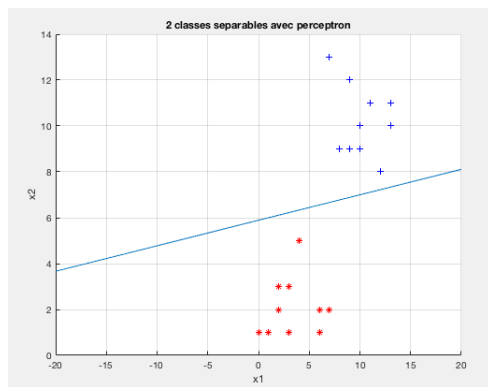


Figure 11: 2 classes séparables

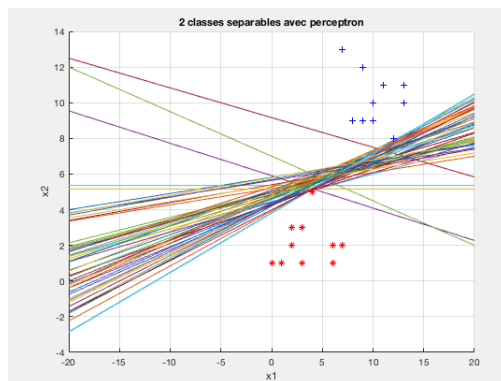


Figure 12: plusieurs solutions pour 2 classes séparables



## 10 Descente de gradient stochastique (SGD)

L'objectif de la descente de gradient stochastique est de réaliser une estimation du "vrai" gradient en ne mettant à jour que ce qui nous intéresse. Le principe est le suivant :

1. Choisir un vecteur initial de paramètres  $\omega$ , et un taux d'apprentissage  $\eta$ .
2. Répéter jusqu'à ce qu'un minimum approché (assez précisément) soit obtenu :
3. Mélanger aléatoirement les échantillons de l'ensemble d'apprentissage.
4. Pour  $i = 1, 2, \dots, n$  faire :
5.  $\omega = \omega - \eta \nabla Q_i(w)$

Dans le projet, le vecteur de paramètres  $\omega$  est une classe nommée `Weights`, possédant qui possèdent plusieurs paramètres :

---

**Algorithm 5** Classe `Weights`

---

```
# This class represents the weights in the logistic regression model.
class Weights:
    def __init__(self):
        self.w0 = self.w_age = self.w_gender = self.w_depth = self.w_position = 0
        # token feature weights
        self.w_tokens = np.zeros(shape=(1070659,1))
        # to keep track of the access timestamp of feature weights.
        # use this to do delayed regularization. self.access_time = {}
    def __str__(self):
        formatter = "{0:.2f}"
        string = ""
        string += "Intercept: " + formatter.format(self.w0) + "\n"
        string += "Depth: " + formatter.format(self.w_depth) + "\n"
        string += "Position: " + formatter.format(self.w_position) + "\n"
        string += "Gender: " + formatter.format(self.w_gender) + "\n"
        string += "Age: " + formatter.format(self.w_age) + "\n"
        string += "Tokens: " + str(self.w_tokens) + "\n"
        return string
```

---

- La mise à jour des paramètres  $\omega$  est réalisé dans la fonction `train(self, dataset, lambduh, step, avg_loss)`. Cette fonction retourne les paramètres mis à jour ainsi que la courbe d'erreur cumulative.

---

**Algorithm 6** simple train

---

```
def train(self, dataset, lambduh, step, avg_loss):
    N = dataset.size
    weights= Weights()
    n_epoch = 1
    N00=0
    N01=0
    N10=0
    N11=0
    count = 0
    nbStep = 100
    T = np.linspace(0,N,N/nbStep)
    for epoch in range(n_epoch):
        while (dataset.hasNext()):
            instance = dataset.nextInstance()
            prediction = self.predict(weights, instance)
            error = instance.clicked - prediction
            if(error==0):
                if(prediction==0):
                    N00+=1
                else:
                    N11+=1
            else:
                if(prediction==0):
                    N10+=1
                else:
                    N01+=1
        # if (error!=0):
        weights.w0 = weights.w0 + step * error
        weights.w_age = weights.w_age + step * error * instance.age
        weights.w_gender = weights.w_gender + step * error * instance.gender
        weights.w_depth = weights.w_depth + step * error * instance.depth
        weights.w_position = weights.w_position + step * error * instance.position
        for indice in instance.tokens:
            weights.w_tokens[indice]= weights.w_tokens[indice]+step*error
    # record the average loss for each step 100
    avg_loss[0] = (1 / 2) * (error * error)
    j = count % nbStep
    if (j == 0 and count / nbStep != 0):
        avg_loss[int(count / nbStep)] = (1 / (2 * count)) * (error * error) +
    avg_loss[int(count / nbStep) - 1]
    count += 1
    print("train DONE")
    return weights,N00,N10,N01,N11,T,avg_loss
```

---

---

**Algorithm 9** logreg

---

```
def regularised_train(self, dataset, lambduh, step, avg_loss):  
    #—même chose que le regularized train—  
    prediction = self.sigmoid(self.compute_weight_feature_product(weights, instance))  
    #—même chose que le regularized train—
```

---

- Régularisation

---

**Algorithm 7** regularized train

---

```
def regularised_train(self, dataset, lambduh, step, avg_loss):  
    #—même chose que le simple train—  
    weights.w0 = (1-step*lambduh/N)*weights.w0 + step * error  
    weights.w_age = (1-step*lambduh/N)*weights.w_age + step * error * instance.age  
    weights.w_gender = (1-step*lambduh/N)*weights.w_gender + step * error * instance.gender  
    weights.w_depth = (1-step*lambduh/N)*weights.w_depth + step * error * instance.depth  
    weights.w_position = (1-step*lambduh/N)*weights.w_position + step * error * instance.position  
    #—même chose que le simple train—
```

---

- logreg

---

**Algorithm 8** sigmoïde

---

```
def sigmoid(self, z):  
    return 1 / (1 + math.exp(-1 * z))
```

---

## Part IV

# Evaluation du modèle

L'évaluation est réalisée pour  $\lambda = \{0.01\}$  pour la régularisation et le logreg.

### 10.1 Matrice de confusion

#### 10.1.1 Apprentissage simple

Prédiction $\hat{y}$ / Vérité $y$	Positifs	Négatifs
Positifs	45978	1906
Négatifs	1916	200

- Ratio de reussite pour le training 0.92356
- Average ctr for training 0.04212
- Average ctr for training predicted 0.04232

#### 10.1.2 Régularisation

Prédiction $\hat{y}$ / Vérité $y$	Positifs	Négatifs
Positifs	45989	1898
Négatifs	1905	208

- Ratio de reussite pour le training 0.92394
- Average ctr for training 0.04212
- Average ctr for training predicted 0.04226

#### 10.1.3 LogRegression

Prédiction $\hat{y}$ / Vérité $y$	Positifs	Négatifs
Positifs	45989	1898
Négatifs	1905	208

- Ratio de reussite pour le training 0.92394
- Average ctr for training 0.04212
- Average ctr for training predicted 0.04226

#### 10.1.4 Tirage aléatoire

Dans ce dernier cas, la probabilité d'obtenir la classe 1 est la même que d'obtenir la classe 0.

Prédiction $\hat{y}$ / Vérité $y$	Positifs	Négatifs
Positifs	2809	2165
Négatifs	2150	2876

- Ratio de reussite pour le training 0.5685
- Average ctr for training 0.5041
- Average ctr for training predicted 0.5026

## 10.2 Courbe d'erreur moyenne

On cherche à tracer la courbe d'erreur moyenne  $L$  en fonction du nombre de pas  $T$ , où

$$L(T) = \frac{1}{T} \sum_{t=1}^T (\hat{y}^t - y^t)^2$$

Dans le graphique suivant, on prend le cas où la probabilité d'obtenir un clic/non-clic est équiprobable (cf 10.1.4).

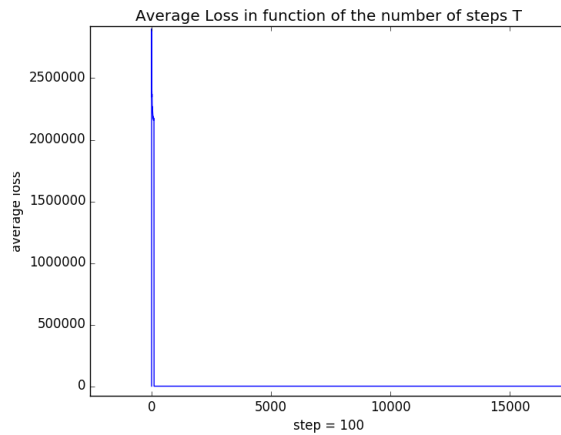


Figure 13: Erreur moyenne en fonction du nombre de pas

## 11 Interprétation des résultats

Nous nous sommes heurtés à plusieurs problèmes pendant la mise en pratique des différentes méthodes.

- La gestion des données nous a été problématique sur certains aspects. Nous avons décidé d'utiliser les classes python de l'énoncé afin d'accéder plus simplement aux données. Cependant, l'accès se fait instance par instance donc il nous faut stocker les entrées  $t-1$  si nous voulons les utiliser avec celle de  $t$ . De plus, nous ne pouvions pas accéder à une instance particulière  $i$ , l'accès se faisant en séquentiel.
- Les tokens sont pour chaque instance une liste de taille indéfinie d'index. Ainsi pour gérer la mise à jour des tokens, nous avons décidé d'utiliser des SparseVector afin de ne pas avoir à créer de vecteurs trop coûteux en dimensions (le nombre de tokens au total s'élève à plus de 1 million) et de pouvoir gérer cette taille de tokens qui varie pour chaque instance. Cette méthode a été vite oubliée car elle n'était pas très efficace de par sa

manipulation et de son temps à l'exécution. Nous avons opté finalement de créer un seul vecteur pour le poids des tokens, de taille le nombre total des tokens, en utilisant la librairie numpy.

- Dans l'échantillon du training, les classes n'étaient pas représentées de manière équitables. En effet, la Classe 0 c'est à dire unclicke est présente dans 97% de l'échantillon total et la classe 1 c'est à dire clicke seulement 3%. Nous avons d'abord décidé d'effectuer un tirage avec remise des données... Mais cela s'est révélé trop coûteux pour nous donc nous avons tiré les données séquentiellement.
- Dans les 2 cas, les résultats sont ambigus. En effet, le taux de réussite est excessivement élevé mais le nombre de Classe 1 prédit par rapport aux Classe 1 total est très faible. Cela est dû à la surprésence de la classe 0 par rapport à la classe 1. L'algorithme considère alors la classe 1 comme une irrégularité dans les données.
- Après la dernière séance, nous avons donc décidé de tirer aléatoirement les données en "piochant" à tour de rôle un coup la classe 1 puis la Classe 0. Les résultats sont beaucoup plus parlants.

## Part V

# Répartition des tâches

Tâches	Octave	Sandratra
taux de clic moyen		X
unique tokens	X	
unique users	X	
Descente de gradient stochastique	X	X
Régularisation	X	X
Evaluation du modèle	X	X
Rédaction du rapport	X	X

## Part VI

# Conclusion

Nous avons fait ce PAO dans l'optique qu'il nous serve d'introduction à la matière DataMining de 4ème année. Nous n'avions aucune connaissance préalable en machine learning ni des librairies s'y référant.

Nous nous sommes familiarisés avec les librairies scientifiques tels que numpy et scipy et avons pu répondre concrètement à un problème d'apprentissage supervisé. En pratique, l'une des problématiques nécessitant le plus d'attention est

*So, what's next ?*



Figure 14: The road to becoming a data scientist

discuter et comprendre les enjeux derrière ce vaste domaine.

s'y intéresser, nous ne doutons pas un instant du brillant avenir de ce domaine.

## Part VII

# Bibliographie

1. Département ASI, Fouille de données : <https://moodle.insa-rouen.fr/course/view.php?id=92>
2. NumPy for Matlab users : <https://docs.scipy.org/doc/numpy-dev/user/numpy-for-matlab-users.html#numpy-for-matlab-users>
3. T. Hastie, R. Tibshirani, J. Friedman, The Elements of Statistical Learning: Data Mining, Inference, and Prediction (Springer Verlag) : <http://statweb.stanford.edu/~tibs/ElemS>
4. University of Washington, Machine Learning for Big Data par Emily Fox : <http://courses.cs.washington.edu/courses/cse547/15sp/>
5. Cours sur le perceptron : <http://pageperso.lif.univ-mrs.fr/~francois.denis/IAAM1/transparentes-Cours3.pdf>
6. Coursera, Machine learning par Andrew NG : <https://fr.coursera.org/learn/machine-learning>
7. Wikipédia : [https://fr.wikipedia.org/wiki/Apprentissage\\_supervise](https://fr.wikipedia.org/wiki/Apprentissage_supervise)