

## TP 3 : Classification binaire avec approches probabilistes

### Introduction

Dans notre TP on peut considérer ham et spam en valeur numérique 0 et 1.

**VP (Vrais Positifs)** : Le nombre de prédictions correctes pour les échantillons positifs.

**FN (Faux Négatifs)** : Le nombre d'échantillons positifs que le modèle n'a pas correctement identifiés.

**FP (Faux positifs)** : Nombre d'échantillons négatifs 0(ham) classés à tort comme positifs 1(spam).

### Exercice 1

#### Question 5 :

Écriture mathématique	Définitions
<p><b>Accuracy:</b></p> $Accuracy = \frac{\text{Vrais Positifs (VP)} + \text{Vrai négatifs(VN)}}{\text{Total}}$	
<p><b>Recall:</b></p> $Rappel(recall) = \frac{\text{Vrais Positifs (VP)}}{\text{Vrais positifs(VP)} + \text{Faux Négatifs(FN)}}$	Un score recall près de 1 meilleur est la détection car elle serait à 100% de précision. Le recall est particulièrement important lorsqu'on manque d'échantillons, comme dans les diagnostics médicaux ou la détection de fraude.
<p><b>F1_Score:</b></p> $F1 = 2 * \frac{\text{Vrais Positifs (VP)}}{\text{Vrais positifs(VP)} + \text{Faux Positifs (FP)} + \text{Faux Négatifs(FN)}}$	Le F1-score atteint sa valeur maximale de 1 (100%) lorsque la précision et le rappel sont parfaits. Un F1-score faible indique soit une mauvaise précision, soit un faible rappel, ou les deux
<p><b>Macro Avg:</b></p> $macro\ avg = \frac{\sum_{i=1}^N Mi}{N}$	C'est une moyenne arithmétique des métrique (précision, rappel, F1_score) N = le nombre de classe, dans notre cas cela correspond à v1 et v2 Mi = la métrique (précision, rapport ou F1_score) pour la classe i (1 et 2 dans notre cas)

<p>Weighted avg :</p> $weighted\ avg = \frac{\sum_{i=1}^N Si.Mi}{\sum_{i=1}^N Si}$	<p>C'est une moyenne pondérée par le support des métrique (précision, rappel, F1_score)</p> <p>N = le nombre de classe, dans notre cas cela correspond à v1 et v2</p> <p>Mi = la métrique (précision, rapport ou F1_score) pour la classe i (1 et 2 dans notre cas)</p> <p>Si = Support (nombre de d'échantillon) de la classe i (v1 ou v2)</p>
--	---

En supplément on peut utiliser la métrique **hamming\_loss** qui permet de calculer la proportion de class mal classifiés, en prenant en compte chaque classe individuellement.

Question 6 :

L'intérêt d'un modèle génératif par rapport aux autres modèles dans notre cas est qu'il est le plus adapté pour les données textuelles vectorisées (CountVectorizer) en matrice, est Bayes et le plus simple et rapide lorsque l'on a beaucoup de caractères comme nous avons ici.

Exercice 2

Question 6 :

Si l'on regarde en le rapport de classification du modèle complément Naïve Bayes dans la Figure 1 : tableau répertoriant les rapports de classification par modèle, nous avons :

La colonne réalité :	Ligne de prediction :
<p>Represente les classes du csv réel</p> <p>Negative : 0 (classe de ham )</p> <p>Positive : 1 (classe spam)</p>	<p>Represente les prediction du modele</p> <p>Negative : 0 (le modele predit le nombre de classe de ham )</p> <p>Positive : 1 (le modele predit le nombre de classe spam)</p>

A partir de toutes ces valeurs on peut calculer nos différentes métriques : Accuracy, Precision, recall, f1\_score, heureusement il existe des fonctions dans scikit-learn pour faire les différents calculs.

Exercice 3:

Exercice 4 : Comparaison des performances

Question 1 : table de comparaison

<

Figure 1 : tableau répertoriant les rapports de classification par modèle

Question 3: Comparaison

Accuracy	Analyses
Generative Bayes	Précision élevée pour les spams (0.95) et les non-spams.

<b>Complement Bayes</b>	Précision plus faible pour les spams (0.80), mais meilleure que la Régression Logistique.
<b>Regression logistique</b>	Précision de 0.80 pour les spams.
<b>Choix :</b> Naïves Bayes est le meilleur en termes de précisions pour minimiser les faux positif (erreurs les ham)	

Recall	Analyses
<b>Generative Bayes</b>	Meilleur rappel pour les spams (0.92), ce qui signifie qu'il manque peu de spams
<b>Complement Bayes</b>	Meilleur rappel (0.94), légèrement supérieur au Naïve Bayes standard
<b>Regression logistique</b>	Appel plus faible (0.85) pour les spams, ce qui signifie qu'il manque davantage de spams par rapport au 2 autres modèles
<b>Choix :</b> Complément naïves bayes est légèrement le 1 <sup>er</sup> en termes de recall	

F1_Score	Analyses
<b>Generative Bayes</b>	Très bon score F1-score pour les spams (0.93) et les non-spams (0.99), elle est bien équilibrée
<b>Complement Bayes</b>	F1-score plus faible pour les spams (0.86), malgré un bon call, malheureusement l'accuracy était plus basse
<b>Regression logistique</b>	Compromis moyen en spam et ham.
<b>Choix :</b> Naïves Bayes est le plus équilibre en termes de F1_score	

Discussion :

Analyse temps d'entraînement :

<b>Generative Bayes</b>	0.196854829788208 secondes
<b>Complement Bayes</b>	0.19664859771728516 secondes
<b>Regression logistique</b>	0.7287561893463135 secondes

En termes de temps d'entraînement on a un très léger avantage du modèle compléments naïves Bayes

Analyse complexité :

<b>Generative Bayes</b>	O(n)
<b>Complement Bayes</b>	
<b>Regression logistique</b>	

Pertinence de chaque approche :

Naïves Bayes est la plus idéale pour les données de type textuelle, lorsque les classes sont équilibrées.

Complément Naïves Bayes : Plus adapté pour les classes déséquilibrées ce qui peut être plus utile pour la détection d'anomalies.

Régression logistique : utile pour une meilleure flexibilité dans la séparation des classes, peut donc être utilisée dans des csv plus brouillons.

**Conclusion :**

Pour des problèmes de type textuels simples et rapide. Pour un dataset déséquilibré il vaut mieux choisir Complément naïves bayes.

Pour des données complexes il vaut mieux utiliser la régression logistique car plus robuste au bruit comme nous l'avons vu lors des précédents TP.