

Résultats de recherche sur la classification supervisée

M. Deker Sylvain¹

Mme. Courdy-Bahsoun Clémence²

¹ M1 IGAI Université Paul Sabatier

² M1 IGAI Université Paul Sabatier

sylvain.deker@univ-tlse3.fr

clemence.courdy-bahsoun@univ-tlse3.fr

Résumé

Cet article vise à comparer les avantages et les désavantages des classifications supervisées par bayésienne avec gaussienne et par la méthode des k-plus proches voisins, en s'appuyant sur les résultats obtenus sur un même jeu de données.

Mots Clef

Apprentissage, classification bayésienne, classification k-ppv.

1 Introduction

La classification bayésienne avec gaussienne et l'approche par la méthode des k-ppv sont des classifications supervisées dont l'objectif est de répartir les données par classe avec un taux d'erreur minimisé. Les tests sont menés sur un jeu de données contenant une centaine d'occurrences de la prononciation d'une dizaine de voyelle au format csv.

A l'aide des travaux réalisés précédemment et en adaptant un peu le code nous avons pu implémenter le modèle de Bayes. L'implémentation pour une approche par les k-ppv, le sujet nous imposait l'utilisation de *Scikitlearn*.

Dans les tests menés les données fournies ont été réparties de telle sorte à avoir 80% de données d'apprentissage et 20% de données de tests. Dans les démarches d'apprentissage le choix d'un taux de données d'apprentissage élevé favorise l'obtention de meilleurs résultats, ici une meilleure prédiction de la répartition par classe.

2 Classification bayésienne avec Gaussienne

La classification bayésienne avec Gaussienne est une minimisation de la vraisemblance qui vise à prédire la classe d'appartenance la plus probable d'une donnée. Néanmoins avec des données réelles il est difficile d'obtenir des résultats satisfaisants car les données ne sont pas simples. Ainsi pour améliorer le modèle on tient compte de la probabilité a priori en appliquant la loi de Bayes.

Dans les tps précédant, nous avons déjà implémenté en python l'algorithme de prédiction permettant de classer les

données fournies. Ici nous sommes repartis de ce code qui était basé sur une simplification de la vraisemblance en log-vraisemblance. C'est-à-dire que la vraisemblance qui est une expression exponentielle a été simplifiée par l'application du logarithme népérien afin d'obtenir un calcul un peu plus optimisé. De plus toutes les constantes ne dépendant pas des données ont été supprimées, n'apportant pas une précision indispensable dans la reconnaissance des données.

3 Classification k-NN

La méthode des K plus proches voisins (K-NN) a pour but de classer des points cibles (classe méconnue) en fonction de leurs distances par rapport à des points constituant un échantillon d'apprentissage (c'est-à-dire dont la classe est connue a priori). Le calcul de la distance peut dépendre de l'expérience menée, Scikitlearn se base sur un calcul euclidien de la distance entre les données de l'échantillon, placées dans un repère après un traitement par une Analyse en Composantes Principales, et les centroïdes de la base d'apprentissage. Scikit learn est une bibliothèque libre pour python dont certaines fonctionnalités ont été utilisées pour implémenter la classification par k-NN. En effet toutes les fonctions appelées sur la variable n de la fonction k_NN sont des fonctions appelées pour la classe KNeighborsClassifier, une classe de la bibliothèque permettant d'effectuer les traitements nécessaires à une classification par plus proche voisin.

Algorithm 1 k_NN

Require: train_data, train_labels, test_data, test_labels

Ensure: m : matrice de confusion

p : precision

n = neighbors.KNeighborsClassifier(n_neighbors, weights = 'distance')

n.fit(train_data, train_labels) %apprentissage

y = n.predict(test_data) mconf = matrice-Conf(test_labels, y) %calcul matrice de confusion, fonction rédigée

return (mconf[0], mconf[1])

4 Evaluation par validation croisée

La validation croisée est une méthode d'estimation de la fiabilité d'un modèle fondée sur la technique d'échantillonnage. Les ensembles de données d'apprentissage sont les fichiers de test data2.csv, data3.csv et data12.csv sur lequel on peut entraîner le modèle. Nous appelons "learning ratio" la proportion des données d'un fichier dédiée à l'apprentissage, la partie restante sert d'échantillon de validation. Cette méthode de test s'appelle "testset validation" ou "holdout method".

5 Résultats

Nous rappelons que le "learning ratio" est la proportion des données d'un fichier dédiée à l'apprentissage, la partie restante servant d'échantillon de validation. Pour évaluer la performance d'un algorithme, nous disposons de deux métriques, le premier étant le temps nécessaire au processus d'apprentissage et de validation, et le second est la précision avec laquelle les résultats sont fiables. Une série de 50 tests ont été effectués sur les fichiers tests data2.csv, data3.csv et data12.csv pour les learning ratio suivant : 0.2, 0.4, 0.6 et 0.8. Les résultats ci-dessous présentent les moyennes sur ces 50 tests.

Pour reproduire les résultats observés il suffit de faire varier les paramètres dans le programme (cad le learning rate, nombre de test, affichage etc...). Par défaut la fonction testmain est configuré pour tester 50 fois l'opération suivante : calcul de la durée et de la précision des algorithmes KNN et Bayes pour un learning ratio fixé. En redirigeant la sortie standard vers un fichier (ex : python3 main.py > resultat.csv) il suffit ensuite de copier/coller correctement les résultats dans le tableur resultat_test.ods pour faire apparaître les résultats. Les FIGURE 1, FIGURE 2, FIGURE 3 et FIGURE 4 ont été réalisés avec le logiciel GNUPlot avec les commandes suivantes :

5.1 Matrice de confusion

Afin d'analyser et de pouvoir comparer les résultats obtenus par les deux méthodes étudiées ici, on calcule pour chaque échantillon la matrice de confusion, permettant de déduire un taux d'erreur. Cette matrice est construite à l'aide d'une classe présente dans la bibliothèque Scikit-learn et en appliquant les calculs nécessaires pour calculer le taux d'erreur.

5.2 Jeu de test

Les jeux de test fournis montrent des classes distinctes les unes des autres pour les fichiers data2.csv (voir FIGURE 1) et pour le fichier data3.csv (voir FIGURE 2, FIGURE 3, FIGURE 4)

5.3 Observations

La FIGURE 5 nous montre que plus la proportion des données liées à l'apprentissage est basse (cad plus la pro-

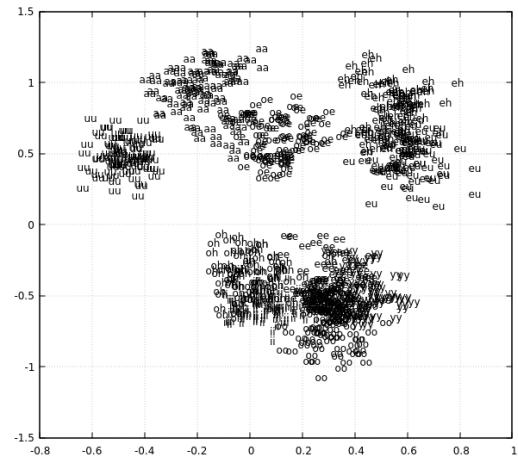


FIGURE 1 – Jeu de test data2.csv

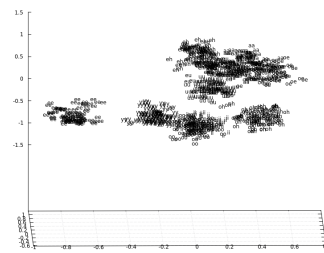


FIGURE 2 – Jeu de test data3.csv vue 1

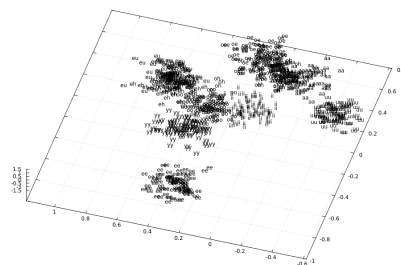


FIGURE 3 – Jeu de test data3.csv vue 2

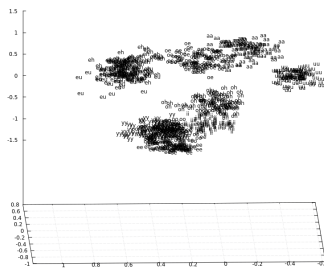


FIGURE 4 – Jeu de test data3.csv vue 3

portion lié à la validation est haute) plus le temps d'exécution global de l'algorithme est long (cad temps d'apprentissage + temps de validation). La restitution des données pour la validation demande plus de temps de calcul que pour l'apprentissage. Cette différence de durée est observé aussi avec la méthode KNN (VOIR FIGURE 6) mais de manière moins drastique. Le pic correspondant à un learning ratio de 0.8 sur le jeu de test data12.csv de la FIGURE 6 reste inexpliqué.

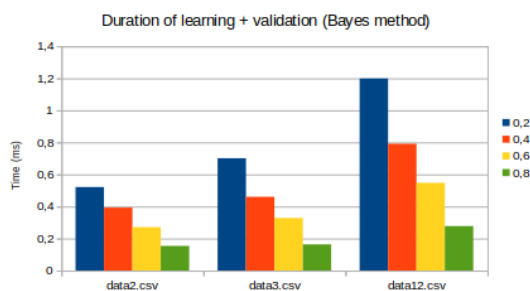


FIGURE 5 – Mesure du temps de l'algorithme Bayes sur les différents jeux de test pour chaque learning rate observé

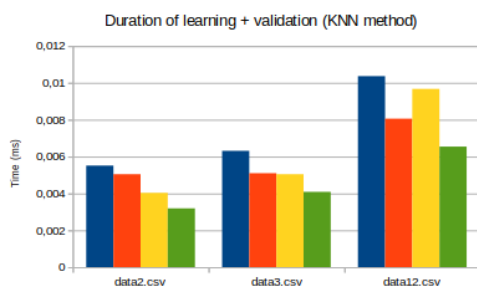


FIGURE 6 – Mesure du temps de l'algorithme KNN sur les différents jeux de test pour chaque learning rate observé

Les FIGURE 7 et FIGURE 8 montre que sur ces jeux de test les méthode KNN et Bayes sont équivalentes, en effet

la différence de précision entre ces deux méthode ne sont pas significative. Des données chiffrés sont disponible dans les FIGURE 9, FIGURE 10, FIGURE 11 et FIGURE 12.

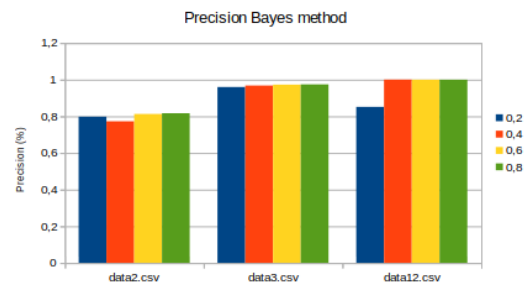


FIGURE 7 – Mesure du la précision de l'algorithme Bayes sur les différents jeux de test pour chaque learning rate observé

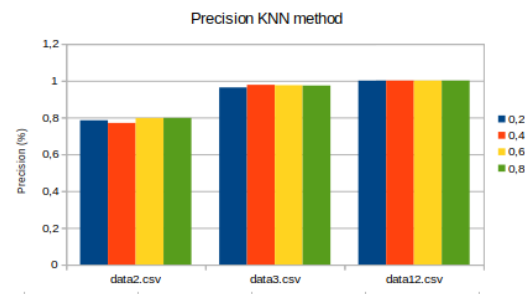


FIGURE 8 – Mesure du la précision de l'algorithme KNN sur les différents jeux de test pour chaque learning rate observé

Learning ratio = 0.2.

20 % des données ont été utilisé pour l'apprentissage de facons aléatoire pour chaque jeu de test (fichier.csv). Les données pour chaque jeu de test est une moyenne de 50 tests réalisés dans le mêmes conditions. Dans ce contexte nous observons que l'algorithme KNN est environ 40 fois plus rapide de celui de Bayes. La précision reste stable, KNN est très légèrement moins précis que Bayes pour les deux premiers jeux de test (voir FIGURE 9).

| Files test | Bayes | | KNN | | Evolution from Bayes to KNN | |
|------------|-----------------|-----------|-----------------|-----------|-----------------------------|-----------|
| | Times | Precision | Times | Precision | Time Acceleration | Precision |
| data2.csv | 0.1536511166992 | 0.8168 | 0.0031902637753 | 0.7916 | 48.226789918047 | 10.0189 |
| data3.csv | 0.1633808898925 | 0.9769 | 0.0060544984981 | 0.9743 | 40.900541652741 | 10.0006 |
| data12.csv | 0.2776164579391 | 0.9997 | 0.0065380048752 | 1 | 42.461953338842 | 10.0003 |

FIGURE 9 – Resultat pour 20% des données dédiées à l'apprentissage

Learning ratio = 0.4.

40 % des données ont été utilisé pour l'apprentissage de facons aléatoire pour chaque jeu de test (fichier.csv). Les données pour chaque jeu de test est une moyenne de 50

tests réalisés dans le mêmes conditions. Dans ce contexte nous observons que l’algorithme KNN est environ 60 fois plus rapide de celui de Bayes. La précision reste stable, KNN est très légèrement plus précis que Bayes pour les deux derniers jeux de test (voir FIGURE 10).

| Files test | Bayes | | KNN | | Evolution from Bayes to KNN | |
|------------|-----------------|-----------|-----------------|-----------|-----------------------------|-----------------|
| | Times | Precision | Times | Precision | Time Acceleration | Delta Precision |
| data2.csv | 0.778641374588 | 0.81256 | 0.0040290403388 | 0.7911 | 87.172684450365 | 0.01546 |
| data3.csv | 0.3260823006821 | 0.9727 | 0.005098731232 | 0.9749 | 85.005384984512 | 0.0022 |
| data12.csv | 0.5476139545441 | 0.9995 | 0.0096665334702 | 0 | 86.650500019978 | 0.0005 |

FIGURE 10 – Resultat pour 40% des données dédiées à l’apprentissage

Learning ratio = 0.6.

60 % des données ont été utilisé pour l’apprentissage de facons aléatoire pour chaque jeu de test (fichier.csv). Les données pour chaque jeu de test est une moyenne de 50 tests réalisés dans le mêmes conditions. Dans ce context nous observons que l’algorithme KNN est environ 80 fois plus rapide de celui de Bayes. La précision reste stable, KNN est très légèrement plus précis que Bayes pour les deux derniers jeux de test (voir FIGURE 11).

| Files test | Bayes | | KNN | | Evolution from Bayes to KNN | |
|------------|------------------|------------------|-----------------|------------------|-----------------------------|------------------|
| | Times | Precision | Times | Precision | Time Acceleration | Delta Precision |
| data2.csv | 0.35231762920703 | 0.80803333333333 | 0.0050431013107 | 0.79543333333333 | 77.701837519797 | 0.0126 |
| data3.csv | 0.4601988983154 | 0.97066666666667 | 0.005856089952 | 0.973 | 78.584675163727 | 0.00233333333333 |
| data12.csv | 0.790848543377 | 0.9995 | 0.0108829483745 | 0 | 72.812255278951 | 0.0006 |

FIGURE 11 – Resultat pour 60% des données dédiées à l’apprentissage

Learning ratio = 0.8.

80 % des données ont été utilisé pour l’apprentissage de facons aléatoire pour chaque jeu de test (fichier.csv). Les données pour chaque jeu de test est une moyenne de 50 tests réalisés dans le mêmes conditions. Dans ce contexte nous observons que l’algorithme KNN est environ 110 fois plus rapide de celui de Bayes. La précision reste stable, KNN est très légèrement plus précis que Bayes pour les deux derniers jeux de test (voir FIGURE 12).

| Files test | Bayes | | KNN | | Evolution from Bayes to KNN | |
|------------|-----------------|-----------|-----------------|-----------|-----------------------------|-----------------|
| | Times | Precision | Times | Precision | Time Acceleration | Delta Precision |
| data2.csv | 0.5208087568809 | 0.79715 | 0.005508310463 | 0.7845 | 84.600324910451 | 0.01325 |
| data3.csv | 0.7010438108444 | 0.9596 | 0.0063081979752 | 0.96305 | 111.1321828523910 | 0.00345 |
| data12.csv | 0.158473110199 | 0.8511 | 0.0103621292114 | 0 | 115.65896214433 | 0.1489 |

FIGURE 12 – Resultat pour 80% des données dédiées à l’apprentissage

6 Conclusion

Les tests menés sur trois jeux de données différents, et sur une cinquantaine d’exécution nous permet de mettre en evidence la corrélation entre la taille de la base d’apprentissage et le temps d’exécution des algorithmes de classification supervisées étudiés est important. En effet on remarque, et ce de manière plus probante sur l’application de la méthode de Bayes, que le temps d’exécution tant à diminuer pour des bases d’apprentissage fournies généreusement. Les tests menés montre également que pour un

même ratio de partage des données, l’algorithme K-NN est beaucoup plus performant en temps d’exécution que l’algorithme de Bayes avec gaussienne. Une explication possible à ce résultat et que les fonctions utilisé par K-NN appartenant à une bibliothèque propre à python, sont des fonctions optimisés au maximum. A cela peu s’ajouter un manque d’optimisation lors de la rédaction des fonctions pour la méthode par bayésienne, en exploitant pas suffisamment les optimisations prévu dans le langage python par un manque de connaissance. On peut seulement déduire que dans le contexte fournit il est plus interessant d’opter pour la méthode des k plus proches voisins . Ce raisonnement est encouragé par les résultats obtenus lors des tests de précision qui se valent pour les deux méthodes. D’après nos résultats on peut en conclure que pour de grands jeux de données il est plus interessant d’opter pour la méthode des k-NN, et eventuellement favoriser la méthode par approche bayésienne pour des petits jeux de données. Cependant il s’agit de résultats expérimentaux. De plus il existe de nombreuses autres méthodes de classification supervisées qui n’ont pas été étudiés ici et qui auront probablement elles aussi leurs propre avantage en fonction du contexte, comme par exemple les resaux de neurones.