

# SY19 - TP10

Hugo Martin, Louis Wahart, Zhifan Huang

May 13, 2023

## 1 Phonemes

Le jeu de données "phonemes" est tiré d'un ensemble de 4509 observations de 512 trames de discours, ayant été adaptées pour faire de la reconnaissance de la parole en log-périodogrammes composés de 256 variables. Parmi toutes ces données, seules 2250 ont été gardées pour entraînement. Ces données sont associées au son émis par 50 hommes lors de la prononciation d'un phonème, élément sonore du langage parlé. Le problème posé est un problème de classification sur la variable  $y$ , qui représente le phonème prononcé parmi "aa", "ao", "dcl", "iy", "sh". Ce problème correspond donc à un problème de classification non binaire mais multiclasse.

### 1.1 Analyse exploratoire

Lorsqu'on regarde les données sans prendre en compte la classe associée aux observations, on se rend compte que la moyenne de toutes les variables se situe très proche de 0, et la variance très proche de 1: on peut donc supposer que les variables ont été centrées et réduites. En regardant les estimations des densités de chacune des variables, on voit que toutes les variables semblent être principalement uni-modales, et semblent presque suivre des gaussiennes.

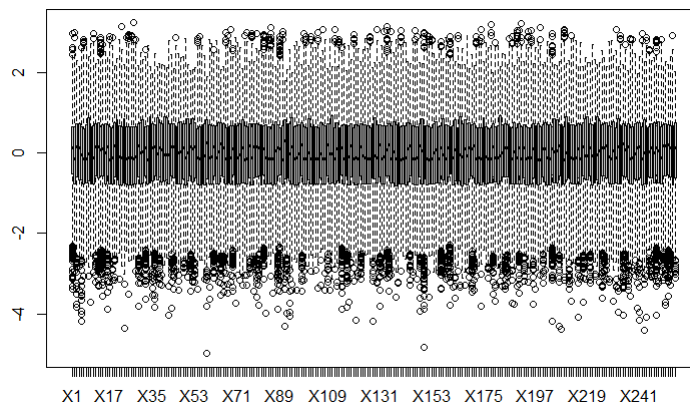


Figure 1: Boxplots de l'ensemble des variables

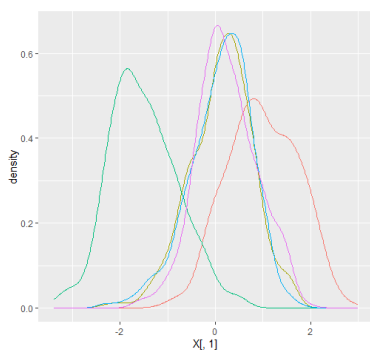


Figure 2: Estimation des densités pour X1

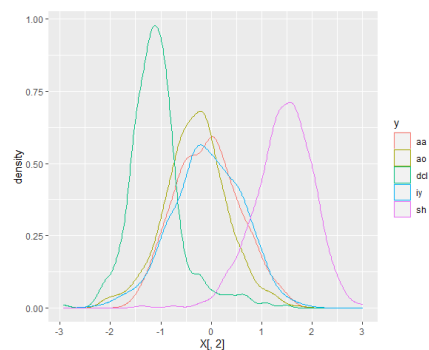


Figure 3: Estimation des densités pour X2

En ce qui concerne la variable de classe à prédire, nous avons remarqué que parmi les 5 phonèmes à prédire, deux d'entre eux se ressemblaient tout particulièrement, et que les 3 autres phonèmes étaient bien plus naturellement distinguables. En effet, les phonèmes "aa" et "ao" se ressemblent beaucoup lors de la prononciation orale, ce qui implique de fortes similitudes dans les données, tandis que les phonèmes "dcl", "iy", et "sh" sont très différents en termes de prononciation, et cela se ressent dans les données.

De plus, du fait de l'apparence de la distribution des variables, nous nous sommes demandé si il existait une simple relation linéaire entre plusieurs variables. Par conséquent, nous avons établi une matrice de corrélation, qui nous a permis d'observer que de nombreuses variables étaient très corrélées linéairement, et qu'il était possible de visualiser des blocs de variables corrélées pour les classes "aa" et "ao". Par exemple, on voit que parmi les premières variables, beaucoup semblent être très corrélées, ce qui a du sens car les deux phonèmes commencent par le même son, mais il est aussi intéressant de remarquer que les dernières variables ont tout de même un coefficient de corrélation élevé, ce qui laisse penser que le son "o" du phonème "ao" est un peu similaire au son "a" du phonème "aa". Nous avons utilisé le coefficient de corrélation de Pearson, qui permet justement de détecter la présence ou l'absence d'une relation linéaire entre des variables, ce qui était ce que nous voulions observer.

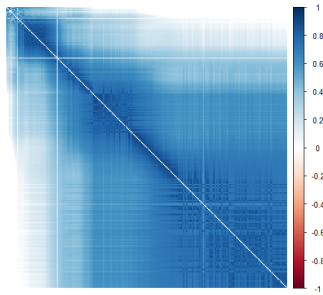


Figure 4: Graphe de corrélation pour toutes les classes (X1-X256)

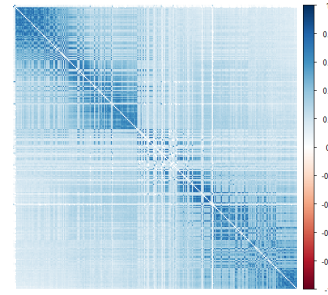


Figure 5: Graphe de corrélation pour les classes "aa" et "ao" (X1-X256)

## 1.2 Choix de modèle

Du fait de la proportion non négligeable de variables corrélées, nous avons gardé en tête que certaines des variables étaient possiblement redondantes, et qu'il serait peut-être judicieux de ne pas prendre l'ensemble des prédicteurs en compte pour notre modèle. En effet, notre réflexion était qu'un nombre de variables aussi important pouvait parfois rendre le modèle plus complexe qu'il avait besoin d'être, et que réduire le nombre de periodogrammes pris en compte dans le modèle pouvait nous permettre d'obtenir une meilleure précision.

De par la nature relativement uniforme des données et du grand nombre de prédicteurs, nous avons commencé par chercher sur des modèles qui profiteraient de la nature des données. Le fait que celles-ci semblent parfois suivre des pseudo-gaussiennes nous donnerait presque envie de tester un classifieur Bayésien naïf, mais le nombre de variables conséquent ne s'y prête évidemment pas. Ainsi, nous avons commencé par des modèles d'apprentissage basiques, ou basés sur une frontière de décision linéaire, comme un modèle des  $k$  plus proches voisins, ou une analyse discriminante linéaire (le nombre de variable fait que le jeu de données ne se prête pas à l'analyse discriminante quadratique). Puis, après des résultats corrects avec la LDA, nous avons tout de même envie de faire mieux, et nous avons alors essayé des modèles un peu plus complexes, comme les SVM ou les modèles basés sur des forêts aléatoires. Ces modèles fonctionnaient bien, mais obtenaient des taux d'erreurs assez similaires à ceux d'une LDA. Afin d'améliorer ces résultats, nous avons utilisé les hypothèses mentionnées précédemment concernant les corrélations entre les variables. Nous avons donc essayé de ne plus utiliser les variables "doublons", c'est-à-dire les variables parmi les couples de variables dont le coefficient de Pearson étaient supérieur à une certaine valeur, que nous avons déterminé par validation croisée, au même titre que le kernel du SVM, son hyperparamètre sigma (0.0168503129656732) ou encore la valeur *mtry* de la random forest.

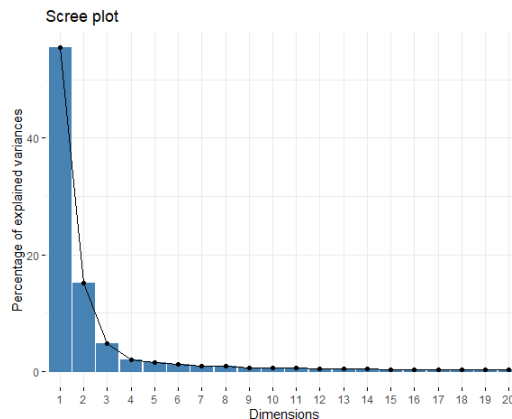
## 1.3 Résultats et améliorations possibles

Comme on peut l'observer sur le tableau 1, les résultats du SVM et de la random forest sont légèrement meilleurs après avoir écarté les variables redondantes. Nous avons retenu le SVM de Weston-Watkins car il est réputé pour obtenir de très bons résultats dans des problèmes multi-classes, et correspondait également au type de SVM qui nous permettait d'avoir les meilleurs taux d'erreur parmi tous les types de SVM disponibles en R.

**Table 1: Résultats récapitulatif des différentes méthodes utilisées, sur les données d’entraînement (taux d’erreurs obtenus avec une validation croisée)**

Sans sélection de variables	Taux d’erreur
LDA	0.081
SVM (C support vector, cost C = 1, linear kernel)	0.105
SVM (Weston, Watkins native multi-class, radial basis kernel)	0.078
Random Forest (500 arbres, mtry=129)	0.078
Avec sélection de variables	
LDA	0.085
SVM (C support vector, cost C = 1, linear kernel)	0.103
SVM (Weston, Watkins native multi-class, radial basis kernel)	0.076
Random Forest (500 trees, mtry=104)	0.080

Cependant, ces résultats sont grandement impactés par la difficulté du classifieur à distinguer le phonème "aa" du phonème "ao". En effet, quand on sépare les deux problèmes (distinguer "aa" de "ao" et distinguer les 3 autres phonèmes), on se rend compte que la grande majorité des cas où le classifieur se trompe de prédiction correspond au problème de distinction entre les deux phonèmes mentionnés précédemment. Pour faire face à ce problème, plusieurs options s’offrent à nous. Nous pourrions combiner plusieurs classifieurs et choisir la classe en fonction de la probabilité moyenne, calculée en combinant les probabilités des différents modèles, ou bien choisir uniquement la deuxième moitié des variables dans le cas des deux phonèmes difficiles à séparer, étant donné que les différences ont toutes les raisons de se trouver à la fin du son. Nous pourrions également séparer le problème en deux étapes: la première consisterait à entraîner notre modèle le plus efficace sur toutes les données, en prenant en compte toutes les classes. Ensuite, il s’agit de trouver un modèle permettant de très bien séparer les deux phonèmes que nous avons du mal à distinguer et de l’entraîner sur les données dont la prédiction était soit "aa" soit "ao", puis de faire une nouvelle prédiction de ces données, cette fois plus spécifiquement. Cependant, la mise en place de cette solution et des autres tentatives n’ont pas eu les résultats escomptés, et nous n’avons jamais vraiment réussi à déroger aux 20 pourcents environ d’erreur de prédiction pour les deux premiers phonèmes. Une de nos hypothèses à ce sujet est que si la mise en place de ces dispositifs ne nous permet pas de nous écarter d’un certain taux d’erreur non-négligeable (environ 20 pourcents), alors il nous faudrait peut-être changer de base, et essayer d’effectuer une ACP. En effet, peut-être qu’en effectuant un changement de base, on pourrait faire ressortir les features les plus importantes pour faire une meilleure distinction entre les classes, en se libérant du facteur de corrélation entre les variables.



**Figure 6: Pourcentages de variance expliquée par les axes tirés de l’analyse en composantes principales**

Après avoir effectué une PCA sur les données, on observe le fameux "coude" sur la figure 6 lorsqu'on affiche la variance expliquée par chacun des axes, ce qui confirme l'utilité d'une PCA, et nous permet ainsi de n'utiliser qu'une petite partie des axes. En utilisant les 40 premiers axes, nous avons pu obtenir des résultats légèrement meilleurs, mais sans changement radical (avec un SVM de Weston Watkins à noyau radial, par validation croisée: 0.072).

float

## 2 Robotics

The Robotics dataset contains data related to the kinematics of a robot arm. It has eight features and one response 'y'. And, y is a continuous number, so the learning task is a regression problem. Because these predictors describe the movement of a robot arm, they may be related to six degrees of freedom. So the non-learner model may be better than the learner model. In addition, the eight variables don't have names, so the model does not need to have a good comprehensibility. But, expectation maximization algorithm can't be used, because we don't have any former knowledge based on these feature

### 2.1 Analyzing Data

In the boxplot of variables, the eight features have almost the same range, they may have the same distribution. Furthermore, their mean is close to zero, and their variance is about 1. So they have been normalized. Also their covariances are around 0.01, so they are not linearly correlated. The Figure 7 shows the distribution of X2, whose middle part is nearly uniform, with nothing special.

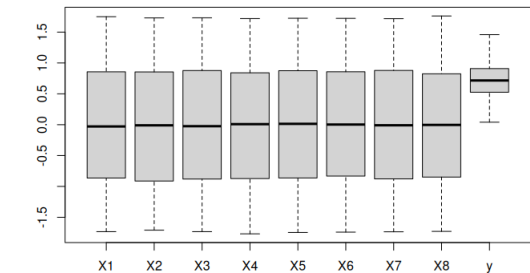


Figure 7: Boxplot of variables

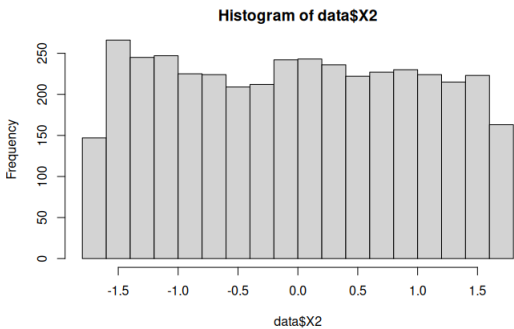


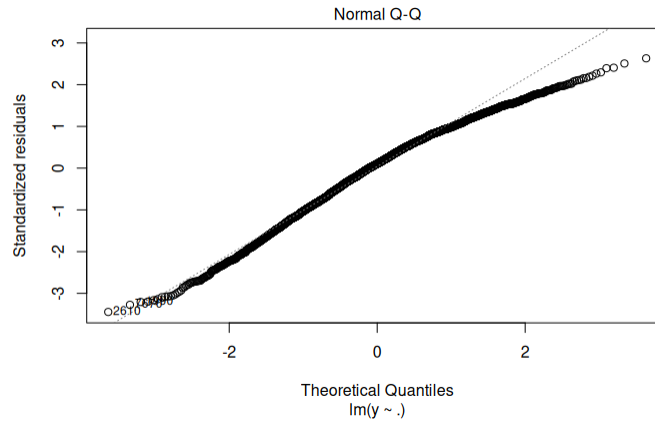
Figure 8: Hist of variables X2

	PC1	PC2	PC3	PC4	PC5	PC6	PC7	PC8
Standard deviation	1.0331	1.0192	1.0130	1.0116	0.9948	0.9876	0.9795	0.9741
Proportion of Variance	0.1329	0.1294	0.1278	0.1275	0.1232	0.1215	0.1195	0.1182
Cumulative Proportion	0.1329	0.2623	0.3901	0.5176	0.6408	0.7623	0.8818	1.0000

This table shows the importance of robotic's features, as generated by the command `prcomp`. Each component has a similar variance and almost has its own unique role. So Dimensionality reduction becomes a difficult task. So methods for decreasing dimensions like subset selection and shrinkage are not attractive anymore. And consider it's density  $4000^{1/8} = 2.82$ . For any point, there are not enough adjacent points to predict the response  $y$ . k-means will lead to the curse of dimensionality.

## 2.2 Linear Model

Because of the simplicity of linear regression, we can try it out and deepen our understanding of the data. However, as previously guessed, the linear model may not be suitable for this dataset. According to the summary of the linear model, each feature's p-value is extremely little, which corresponds to the previous component analysis. However, the adjusted R-squared is 0.4015. This means the model failed to predict the response  $y$ . 50% variance of  $y$  is not explained by this model. The MSE on test data is 0.04 The photo shows the Q-Q plot of the residual.



**Figure 9: Linear regression residual'qqplot**

When the residuals are greater than zero, they do not satisfy the normal distribution, which means the assumption of epsilon failed. Since the value of  $y$  is greater than 0, I also tried to take the log of  $y$ , but it did not improve the residuals' problem.

I then tried using nature spline to see if the degree of each feature would improve the results. By using the following formula. But that also failed.

$$y \sim ns(X1, df=3) + ns(X2, df=3) + ns(X3, df=3) + ns(X4, df=3) + ns(X5, df=3) + ns(X6, df=3) + ns(X7, df=3) + ns(X8, df=3)$$

Finally, I correlate multiple variables by the formula  $y \sim X1 * X2 * X3 * X4 * X5 * X6 * X7 * X8$ . It is exciting to see that the results have improved significantly, the R-squared increase to 0.7, and the test MSE drop to 0.02. So these variables are not independent of each other.

## 2.3 Non-linear Model

### 2.3.1 Random Forest

Since there are only 8 variables, the value of `mtry` for the random forest is 3 regardless of whether the default  $n/3$  or  $\log_2(n)$ . And the test MSE is not special, 0.021. By increasing the `ntree` parameter, the MSE little optimized from 0.021 to 0.020. And the result of bagging is similar to the random forest.

### 2.3.2 Support Vector Machine

As said before, linear models are not very good, so kernel tricks are important to nonlinear transform the input feature space to a new feature space.

rbfdot	polydot	vanilladot	tanhdot	laplacedot	besseldot	anovadot
0.01271	0.04295	0.04295	662.1834	0.0156	0.2601	22.2761'

The table above shows the cross-validation result of different kernels. The rbfdot and laplacedot have the best result. Their MSEs are already better than previous models. Besides, polydot use polynomial of features, it's MSE closed to the spline models.

We then focused on their different parameters by using cross-validation.

The grid of cross-validation about cost and sigma with kernel laplacedot						
C sigma	0.03125	0.0625	0.125	0.25	0.5	1
32	0.008338947	0.008306725	0.008107607	0.008436020	0.008921340	0.01079038
64	0.008339537	0.008212678	0.008330160	0.008416510	0.008672814	0.01099427
128	0.008256747	0.008220027	0.008400844	0.008296268	0.008876666	0.01081118
256	0.008145284	0.008213203	0.008224357	0.008380859	0.008837824	0.01105220
512	0.008309967	0.008223594	0.008313709	0.008319953	0.008829467	0.01085927
1024	0.008257692	0.008198213	0.008409451	0.008416197	0.008759666	0.01093048

The grid of cross-validation about cost and sigma with kernel rbfdot							
C sigma	0.03125	0.0625	0.125	0.25	0.5	1	2
0.5	0.02685705	0.019234443	0.011482773	0.008134370	0.010568028	0.02630474	0.05855071
1	0.02476976	0.016385046	0.008990384	0.007034462	0.009253777	0.02068916	0.05256697
2	0.02198867	0.013724261	0.007574106	0.006673018	0.008719502	0.01997557	0.04971728
4	0.01941255	0.011470439	0.006813427	0.006739043	0.009127632	0.02004326	0.04973362
8	0.01739689	0.009400735	0.006567940	0.007250834	0.008966554	0.02020160	0.04970631
16	0.01569961	0.008281557	0.006474320	0.007836754	0.008923694	0.01992549	0.04976910
32	0.01319922	0.007324950	0.006771780	0.007980618	0.009191220	0.02011846	0.04951920
64	0.01098693	0.007087883	0.007666046	0.008211932	0.009254396	0.02005148	0.04946180

These two tables show how the two kernels behave under different parameters. The parameter interval is exponential in two to cover more ranges. For laplacedot, its performance verify on the cost and sigma with no evidence of tendency. When the sigma is smaller than 0.25, the MSE seems to change randomly based on the cost. For rbfdot, it has better performance than laplacedot. Furthermore, it has a clear minimum MSE interval. MSE is minimized when sigma around 0.25, cost around 16.

Because epsilon defines a margin of tolerance where no penalty is given to errors and error is determined concerning cost. So I do another grid search with epsilon and cost.

The grid of cross-validation about cost and epsilon with kernel rbfdot						
C epsilon	0.0125	0.025	0.05	0.1	0.2	0.4
6	0.006668170	0.006806161	0.006589806	0.006670819	0.006711102	0.008297290
8	0.006685697	0.006597636	0.006554867	0.006409632	0.006748865	0.008125630
10	0.006718945	0.006549523	0.006505093	0.006609048	0.006706116	0.007997654
12	0.006757601	0.006577276	0.006493934	0.006553620	0.006755135	0.007979025
14	0.007013213	0.006752403	0.006716760	0.006607136	0.006664013	0.008067838
16	0.006782465	0.006810728	0.006617948	0.006580512	0.006762107	0.008226855

The result indicates that epsilon = 0.1 is the best choice. Finally, the following command is used to train SVM model.

```
ksvm(y~., data=data, kernel='rbfdot', C=8, epsilon=0.1, kpar=list(sigma =0.125))
```

### 2.3.3 Neural Network

The cross-validation of hidden units number in nnet					
29	36	43	50	57	64
0.006773181	0.006802221	0.006909978	0.006457816	0.006649864	0.006742320
71	78	85	92	99	
0.006483055	0.006189515	0.006319409	0.005934737	0.005896900	

The neural network has a good performance in predicting the response. It's performance increases with the number of hidden neurons.

## 2.4 Conclusion

For this dataset, the nonlinear model is better than the linear model. However, by adjusting the linear regression formula, like adding cross terms, the model can be improved, even comparable to the random forest. SVM and neural networks perform the best, probably because the features record the space coordinates of the robot arm. While these two models can extract information like trajectory shapes (circles, ellipses) from them, so they perform better.

## 3 Communities

Ce jeu de données contient 1000 observations de 122 variables socio-économiques et policières recueillies dans autant de communautés aux États-Unis dans les années 90. Nous devons construire un modèle permettant de prédire le nombre de crimes violent pour 100,000 personnes d'une part et d'autre part déterminer quelles variables recueillies influent le plus sur cette valeur.

### 3.1 Prétraitement

On peut tout d'abord noter que l'on disposait également de 5 variables dites non-prédictives. Parmi ces variables on a le numéro de l'État, ce qui pourrait avoir une valeur prédictive, puisque certains États pourraient avoir un taux de crimes violents plus importants, mais cela n'apporterait que peu d'informations par rapport à l'analyse socio-économique des variables associées à ce taux. On ne prendra donc pas en compte cette variable. Il n'est ici pas nécessaire de standardiser les variables car d'après le descriptif du jeu de données elles ont toutes déjà été normalisées et sont donc sur la même échelle, comme on peut le voir sur la figure 10. Par ailleurs on remarque que deux variables, "PolicePerPop" et "LemasSwornFT" ont une corrélation de 1, elles ont toujours la même valeur. On supprime donc la deuxième variable du jeu de données. Nous séparons le jeu de données en un ensemble d'entraînement et un ensemble de test selon un découpage 80%/20%.

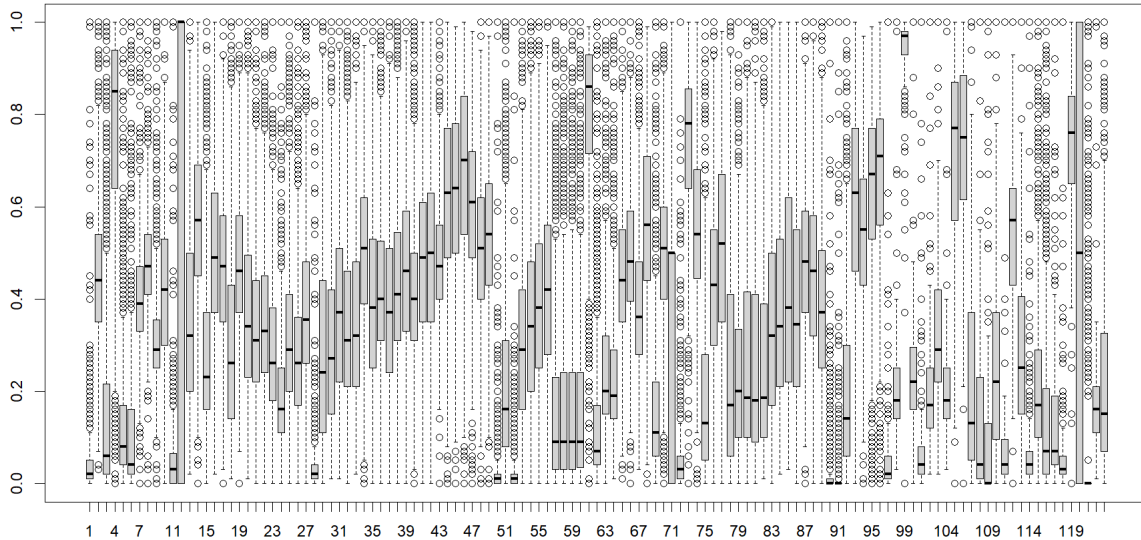


Figure 10: Boxplots de l'ensemble des variables

### 3.2 Gestion et analyse des données manquantes

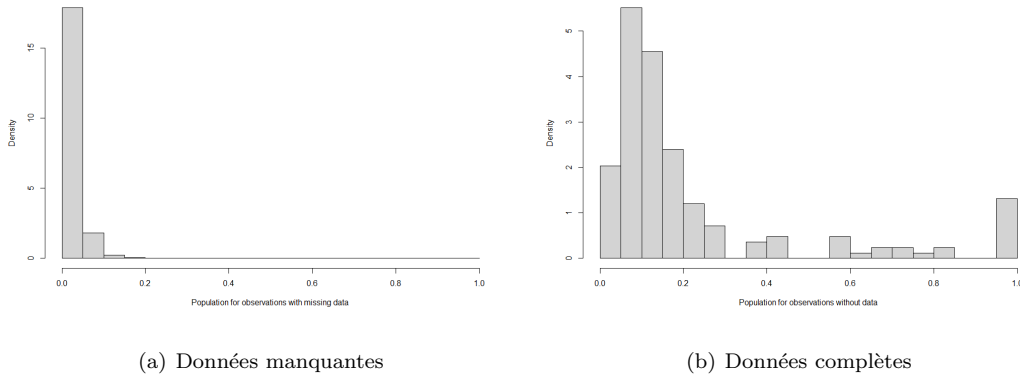
Il manque la donnée "OtherPerCap" pour une des communautés sans qu'il n'y ait de raison donnée, nous supposons donc que cette absence est complètement aléatoire. Nous décidons ainsi de remplacer cette donnée par une valeur représentative de cette variable. Nous réalisons un test de normalité de Shapiro-Wilk pour déterminer si

la moyenne empirique est une statistique représentative de la distribution pour imputer cette donnée. Nous trouvons une p-value inférieure à  $2.2 * 10^{-16}$ , on peut rejeter l'hypothèse nulle de normalité avec une grande confiance. On décide donc d'utiliser la médiane plutôt que la moyenne pour cette variable.

Le reste des données manquantes touche les variables liées à l'état des forces de police dans les communautés étudiées. Pour ces variables, nous envisageons trois approches différentes :

1. Suppression des variables aux données incomplètes
2. Suppression des communautés aux données incomplètes
3. Imputation des données manquantes

La première approche a un problème évident : on ne pourra pas déterminer l'impact des variables supprimées sur la variable cible. La seconde approche pose deux problèmes. D'abord les observations concernées par les données manquantes sont très nombreuses : 833 sur 1000. Ensuite, les données ne sont pas manquantes de manière aléatoire. En effet d'après la description du jeu de données, les données liées aux forces de police sont issues de l'étude LEMAS qui ne concernait que les départements de police comptant plus de 100 officiers, ainsi qu'une sélection aléatoire de plus petits départements. Les données manquantes correspondent donc toutes à des départements de police comptant moins de 100 officiers (petites communautés, communautés "tranquilles", communautés en manque de moyen, etc), ce qui peut introduire un biais dans l'analyse et conduire à de moins bons résultats sur les données tests, où l'on devra appliquer le modèle à des communautés ayant toute taille de départements de police. De plus, même en supprimant les lignes concernées il sera nécessaire de trouver une méthode d'imputation pour les variables manquantes des données de test. La dernière méthode envisagée est l'imputation des données manquantes. Comme précédemment on pourrait remplacer les données manquantes par la moyenne ou la médiane des données. Ici cela ne nous semble pas pertinent puisque comme expliqué précédemment les données manquantes ne sont pas complètement aléatoires mais liées au nombre de policiers dans le département de police de la communauté. Ainsi la distribution de certaines variables au sein des observations incomplètes est parfois très différente de celles observée pour les données complètes, comme on peut le voir sur la figure 11 pour l'attribut "population" par exemple. On peut supposer que cela soit également le cas pour certaines des données manquantes. En faisant la médiane ou la moyenne, on inclurait donc dans le calcul des forces de police de petits départements des données de départements bien plus grands. Nous pouvons tenter de construire un modèle de prédiction pour chaque attribut en fonction des autres attributs (complets) présents dans le jeu de données.



**Figure 11: Distribution de l'attribut "population" parmi les données manquantes et les données complètes**

### 3.3 Imputation des données de police manquantes

Pour chaque variables manquantes nous construisons un modèle de prédiction à partir des autres features. Nous décidons de nous restreindre à des modèles linéaires obtenus par sélection itérative des variables en "forward" basée sur le  $R^2$ , puis par une comparaison des meilleurs modèles à  $n$  variables par validation croisée ( $k = 5$  répétée 3 fois).



### 3.4 Modèles de prédiction

Nous voulons obtenir un modèle interprétable, nous pouvons donc par exemple construire un modèle linéaire pour prédire la variable cible. Le coefficient associé à une variable  $X_i$  représentera l'augmentation répercutée sur la variable cible si l'on augmente la variable  $X_i$  d'une unité. On apprend tout d'abord un modèle linéaire utilisant toutes les variables avec une régularisation sur les paramètres de régularisation élastique  $\alpha$  et  $\lambda$ , en incluant le cas d'une absence de régularisation ( $\alpha = 0$  et  $\lambda = 0$ ). Nous entraînons trois modèles de ce type sur les trois jeux de données décrits précédemment. On teste ensuite un modèle obtenu de la même manière que les modèles de la section précédente (sélection itérative). Enfin, on teste un modèle gardant les 50 premières variables sélectionnées itérativement tout en appliquant une régularisation élastique. L'erreur quadratique moyenne sur le jeu de test (MSE) et le  $R^2$  pour chacun de ces modèles sont rapportés en table ??.

Nous entraînons également des modèles de forêt aléatoire sur chacun des jeux de données, qui ne seront pas interprétables mais qui pourront peut-être fournir des prédictions plus précises. Pour ce faire, on effectue dans le cas du jeu de données imputé une validation croisée ( $k = 10$ , 5 répétitions) avec différentes valeurs de *mtry*, le nombre de variables considérées à chaque division. On fait une première validation croisée avec des valeurs assez éloignées les unes des autres sur l'espace de la variable : 3, 10, 20, 50, 80, 100. On obtient les meilleurs résultats pour les valeurs 10 (0.01922) et 20 (0.01923). On effectue ensuite une seconde validation croisée d'affinement avec les valeurs de 8 à 30 par pas de 2. On obtient le meilleur résultat avec la valeur 24 (0.01789). On obtient une valeur de 0.01948 sur le jeu de données de test. On répète les mêmes opérations pour les deux autres jeux de données, mais on obtient de moins bons résultats.

### 3.5 Analyse de l'influence des variables sur le taux de crimes violents

Avant de s'intéresser aux variables les plus influentes/importantes des modèles précédents, il convient de rappeler que d'une part les modèles linéaires ne prennent pas en compte l'interaction des variables entre elles, et d'autre part ces modèles n'établissent pas des liens de causalité mais des associations entre les variables. Ainsi il est parfois possible de prédire efficacement une variable  $Y$  avec une variable  $X_1$  n'ayant pas de lien causal (dans un sens ou l'autre)  $Y$  mais ayant un "parent commun"  $X_2$  ayant un effet causal sur  $X_1$  et  $Y$ . On emploiera donc les termes "effet", "importance", "influence" dans cette section dans le sens d'intensité de l'association apprise par les modèles, et non d'une quelconque causalité supposée.

Les modèles linéaires développés précédemment sont tous interprétables, on peut savoir dans quelle mesure chaque variable a participé à la prédiction. Les coefficients étant sur la même échelle de 0 à 1, on peut calculer la moyenne et l'écart-type des coefficients obtenus avec chaque modèle développé pour les jeux de données 1 et 3 (présentant chacun toutes les variables) pour obtenir une valeur permettant de comparer l'influence des coefficients entre eux. On trouve les valeurs présentées en figure 12.

On remarque que l'ensemble des modèles sont globalement unanimes en ce qui concerne le sens de l'influence des variables les plus importantes. Par exemple pour tous les modèles "PctKids2Par" (pourcentage d'enfants vivant en famille avec deux parents) semble avoir un effet négatif pour le taux de crimes violents, ce qui paraît intuitif. On note cependant pour ce coefficient et pour nombre d'autres une variance assez importante entre les modèles pour le jeu de données 2, par exemple pour "PctTeens2Par" qui passe de 0 à 0.98. Cela est probablement dû au manque de données disponibles pour ce jeu de données, puisque les coefficients des modèles développés sur les autres jeux de données présentent des variances moins importantes.

Les variables liées aux forces de police, présentes dans les jeux de données 2 et 3, semblent avoir chacune un peu d'influence sur la variable cible.

Parmi les variables qui semblent être les plus importantes selon ces modèles, on peut citer "PctKids2Par", "racepctblack", "NumStreet", ainsi que des variables mesurant le capital économique des communautés (données liées au loyer, au revenu, aux emplois majoritaires, etc suivant les modèles), toujours présentes dans les dix variables les plus importantes de chaque modèle. Les variables du jeu de données présentent un grand nombre de corrélations, il est donc possible que chaque modèle choisisse plus ou moins aléatoirement une variable représentative de la situation économique, et extraient le nécessaire d'information économique de celle-ci sans avoir besoin du reste des informations du même type.

Le problème des modèles régularisés est que l'on a pas de p-value associée à chaque coefficient, on a donc pas d'assurance que le coefficient soit différent de 0 de manière statistiquement significative. On a cependant accès aux p-values des coefficients des variables des modèles linéaires aux variables sélectionnées itérativement, que l'on peut lire figure 13.

On retrouve dans les trois modèles les coefficients statistiquement significativement différents de 0 (à 95%) associés aux variables suivantes : "NumStreet" (nombre de sans-abris, positivement), "MedOwnCostPctIncNoMtg"

	coeffs	elastic	net	coeffs	iterative	selection	coeffs	it.	sel.	+ elastic	net	mean	std
PctKids2Par	-0.300650231			-0.31819991			-0.355955034			-0.324935057	1.597354e-03		
racepctblack	0.20645678			0.36156693			0.109062051			0.225758219	3.242729e-02		
PctNotHSGrad	0.136396356			0.18481573			0.137016807			0.152742966	1.543186e-03		
NumStreets	0.098112767			0.12525736			0.086604655			0.112324926	2.458117e-03		
HouseVacant	0.069510608			0.124825510			0.118272213			0.104436442	1.836973e-02		
RentLow	0.000000000			-0.28867777			0.000000000			-0.096215923	5.555652e-02		
MalePctDivorce	0.093407642			0.12395623			0.061159047			0.092840972	1.972225e-03		
MedRent	0.000000000			0.23228509			0.000000000			0.077428362	3.597091e-02		
MedOwnCostPctIncNoMtg	-0.072768852			-0.13229918			-0.014626442			-0.07321493	6.923758e-03		
PctIllegal	0.112826714			0.000000000			0.102497274			0.07192813	7.783215e-03		
PctNotSpeakEngWell	0.000000000			-0.20715887			0.000000000			-0.069052957	2.860986e-02		
PctRecImmig8	0.000000000			0.18682710			0.000000000			0.062275701	2.326958e-02		
PctHouseOccup	-0.060215161			-0.10083788			-0.024707349			-0.061920129	2.902289e-03		
PctNotHSGrad	0.016059469			0.13504905			0.000000000			0.050369506	1.088489e-02		
PctPopUnderPov	0.000000000			-0.14348586			0.000000000			-0.047828620	1.372546e-02		
PctWorkMomYoungKids	-0.038693969			-0.09710473			-0.005948897			-0.047249531	4.264398e-04		
racePctHisp	0.010484915			0.12586035			0.000000000			0.045448421	9.754084e-03		
PctImmigRec10	0.000000000			-0.11000422			0.000000000			-0.036668073	8.067286e-03		
MedRentPctHouseInc	0.011090623			0.08317142			0.000000000			0.031420680	4.078709e-03		
racePctWhite	-0.029366725			0.04425087			-0.105903235			-0.030340362	1.127485e-02		
PctPop15to19	-0.4422115			-0.150455093			0.000000000			0.102722216	6.979989e-02		
PctVacMore6Mos	0.000000000			-0.07420500			0.000000000			-0.024735001	3.670922e-03		
pctwFarmSelf	0.000000000			0.06938543			0.000000000			0.023128476	3.209558e-03		
numbUrban	0.046132934			0.000000000			0.019943648			0.022025527	1.070625e-03		
MedVHouseBuilt	0.000000000			-0.05651399			0.000000000			-0.018837996	2.129221e-03		
PctForeignBorn	0.045646047			0.000000000			0.000000000			0.016094902	1.313310e-02		
NumInShelters	0.036033753			0.000000000			0.000000000			0.012011251	8.656209e-04		
agePct12to29	0.000000000			0.03074501			0.000000000			0.010248336	6.301703e-04		
PctImmigRecent	-0.022739582			0.000000000			0.000000000			-0.007579861	3.447257e-04		
pctwSecSec	0.022159562			0.000000000			0.000000000			0.007398654	3.284405e-04		
PerKerenttoCchous	0.000926740			0.000000000			0.000000000			0.00120200	5.444082e-03		
AsianPerCap	0.006722474			0.000000000			0.000000000			0.002240825	3.012777e-05		
PctSameCity85	0.003687108			0.000000000			0.000000000			0.001229036	9.063179e-06		
MedNumR	-0.003032770			0.000000000			0.000000000			-0.001010923	6.131797e-06		

(a) Jeu de données 1

	coeffs	elastic	net	coeffs	iterative	selection	coeffs	it.	sel.	+ elastic	net	mean	std
PctKids2Par	-0.35232163			-0.26697935			-0.002138000			-0.073313209	6.278006e-03		
PctTeen2Par	0.000000000			0.86799070			0.000000000			0.329330234	6.507504e-01		
PctNotHSGrad	0.000000000			-0.94773395			0.000000000			-0.115911941	5.988023e-01		
PctForeignBorn	0.000000000			0.53582173			0.000000000			0.218504489	2.608361e-01		
PctPopMenu	-0.063903959			-0.40617838			-0.154375173			-0.208152043	6.291389e-02		
HouseVacant	-0.4422115			-0.069510608			-0.006915400			-0.150455093	1.307718e-02		
Houseless3B8	0.077600408			0.30247216			0.068731227			0.149608001	3.509004e-02		
PctUnempLOY	0.028643228			0.11017182			0.101687803			0.147500512	4.250856e-02		
racepctblack	0.04667242			0.13612278			0.062103926			0.145163085	4.935461e-02		
PctLess9thGrade	0.000000000			0.40646062			0.000000000			0.135468729	1.101402e-01		
PerCapInc	-0.19994178			0.00943086			-0.00943086			-0.116475438	7.496213e-02		
NumStreet	0.093779921			0.10094164			0.142235077			0.118191935	1.040153e-03		
PctRecentImmig	-0.00245335			-0.10224535			0.000000000			-0.100748456	6.090150e-02		
LenasapctPopInc	-0.00193436			-0.20981527			-0.07330818			-0.096423042	2.114486e-02		
PctUrban	0.000000000			0.26341811			0.000000000			0.093464152	4.594313e-02		
PctBornsamState	0.000000000			-0.21366709			-0.060813076			-0.091160729	2.435284e-02		
PctHouseOwns	0.000000000			0.25633547			0.000000000			0.085451570	4.380025e-02		
LandArea	0.000000000			0.15992770			0.095966085			0.081997828	1.164572e-02		
racePctWhite	-0.107903398			-0.07788883			-0.006906167			-0.007074659	4.258416e-02		
PctPersDenseHous	0.026918422			-0.27458081			0.00902342			-0.0792540145	5.737357e-02		
PctPopMenu	0.000000000			0.23050842			0.000000000			0.076769471	3.563131e-02		
PctPoliceAsian	0.04783769			0.05510476			0.000000000			0.065001334	1.580186e-03		
LenasapctPopInc	0.05895077			0.04274891			0.04859797			0.061931957	8.665216e-02		
NumIndurpSetz	0.000000000			0.15510702			0.029145410			0.081520448	1.360708e-02		
PctPolicePop	0.000000000			0.16684428			0.000000000			0.055564150	1.852749e-02		
LenasapctPopInc	0.000000000			0.13883676			0.001380662			0.0467391405	1.272391e-01		
PctFam2Par	-0.108754466			0.000000000			0.000000000			-0.032511621	7.889556e-03		
PctHouseOccup	-0.076494093			0.10783690			0.000000000			0.035945317	7.752511e-03		
PctYoungKids2Par	0.000000000			0.000000000			0.000000000			-0.025165444	3.921127e-03		
PctTeen2Par	0.000000000			0.000000000			0.000000000			0.223030491	2.912018e-03		
PctHouseOwns	0.000000000			0.000000000			-0.08835138			-0.019611727	2.307756e-03		
PctImmigRecent	0.000000000			0.000000000			-0.05813735			-0.017604783	1.891276e-02		
PctwInc	0.000000000			0.000000000			0.000000000			-0.014962524	1.127697e-03		
PctPoliceAsian	0.000000000			0.000000000			0.000000000			0.013167179	1.040336e-03		
PctPolicePop	0.000000000			0.000000000			0.000000000			0.010260194	7.182146e-02		
PctPolicePop	0.000000000			0.000000000			0.000000000			0.0094807716	5.993102e-04		
PctPolicePop	0.000000000			0.000000000			0.000000000			0.009384303	5.238996e-04		
PctPolicePop	0.000000000			0.000000000			0.000000000			0.008851224	4.812398e-04		
PctPolicePop	0.000000000			0.000000000			0.000000000			0.007380266	3.780736e-04		
PctPolicePop	0.000000000			0.000000000			0.000000000			0.007322466	3.409836e-04		
PctPolicePop	0.000000000			0.000000000			0.000000000			0.006713295	2.704006e-04		
PctPolicePop	0.000000000			0.000000000			0.000000000			0.0064952721	2.531314e-04		
PctPolicePop	0.000000000			0.000000000			0.000000000			0.004871342	1.102077e-04		
PctPolicePop	0.000000000			0.000000000			0.000000000			0.002676367	3.955347e-05		
NumIndurpSetz	0.002401681			0.000000000			0.000000000			0.000012269	3.511787e-06		
MalePctDivorce	0.093407642			0.12395623			0.061159047			0.092840972	1.972225e-03		
MedRent	0.000000000			0.23228509			0.000000000			0.077428362	3.597091e-02		
MedOwnCostPctIncNoMtg	-0.072768852			-0.13229918			-0.014626442			-0.07321493	6.923758e-03		
PctIllegal	0.112826714			0.000000000			0.102497274			0.07192813	7.783215e-03		
PctNotSpeakEngWell	0.000000000			-0.20715887			0.000000000			-0.069052957	2.860986e-02		
PctRecImmig8	0.000000000			0.18682710			0.000000000			0.062275701	2.326958e-02		
PctHouseOccup	-0.060215161			-0.10083788			-0.024707349			-0.061920129	2.902289e-03		
PctNotHSGrad	0.016059469			0.13504905			0.000000000			0.050369506	1.088489e-02		
PctPopUnderPov	0.000000000			-0.14348586			0.000000000			-0.047828620	1.372546e-02		
PctWorkMomYoungKids	-0.038693969			-0.09710473			-0.005948897			-0.047249531	4.264398e-04		
racePctHisp	0.010484915			0.12586035			0.000000000			0.045448421	9.754084e-03		
PctImmigRec10	0.000000000			-0.11000422			0.000000000			-0.036668073	8.067286e-03		
MedRentPctHouseInc	0.011090623			0.08317142			0.000000000			0.031420680	4.078709e-03		
racePctWhite	-0.029366725			0.04425087			-0.105903235			-0.030340362	1.127485e-02		
PctPop15to19	-0.4422115			-0.069510608			0.000000000			0.102722216	6.979989e-02		
PctVacMore6Mos	0.000000000			-0.07420500			0.000000000			-0.024735001	3.670922e-03		

	Estimate	Pr(> t )		Estimate	Pr(> t )		Estimate	Pr(> t )
racepctblack	0.31281969	7.707916e-07	pctWInvInc	-2.48918782	4.457991e-09	racepctblack	0.36704306	6.095796e-07
MalePctDivorce	0.27480483	2.196469e-06	PctNothSGrad	-1.81112586	1.217309e-08	MalePctDivorce	0.66941858	1.740247e-05
NumStreet	0.24138179	6.473620e-05	NumStreet	0.40292570	2.129761e-07	PctWorkMom	-0.13034518	3.151508e-04
PctWorkMom	-0.15590875	1.474232e-04	PctBornSameState	-1.13229332	2.203827e-07	PctIlleg	0.21095313	9.297773e-04
PctForeignBorn	0.23298201	1.620926e-04	MedOwnCostPctInc	-0.73189027	4.520012e-07	MedOwnCostPctIncNoMtg	-0.10128210	1.042250e-03
PctEmploy	0.30701275	2.542964e-04	MedRent	4.17851026	6.412365e-06	TotalPctDiv	-0.56302153	1.128115e-03
pctWwage	-0.31052828	2.895374e-04	pctWPubAsst	0.99256379	1.614498e-05	PctLargHouseFam	-0.23549706	1.227578e-03
MedOwnCostPctIncNoMtg	-0.11074995	3.186656e-04	PctOccupMgmtProf	1.41138095	2.686698e-05	PctPersDenseHous	0.29055260	1.557314e-03
PctHousOccup	-0.11607010	4.209519e-04	LemasWtFieldPerPop	2.54451458	3.525119e-05	whitePerCap	-0.66442240	2.726547e-03
whitePerCap	-0.73643942	7.282029e-04	PctForeignBorn	1.57777600	5.227244e-05	PctPolicAsian	0.13375936	4.462702e-03
RentLowQ	-0.28040199	9.142691e-04	PctImmigrRec5	-2.55259353	8.178226e-05	PolicPerPop	-0.03815597	6.830964e-03
MalePctNevMarr	0.18330726	1.331937e-03	PolicPerPop	-2.13030943	1.194449e-04	PolicReqPerOffic	0.10561911	7.443131e-03
PctPersDenseHous	0.25284568	4.402140e-03	PctRecImmig10	-2.84453763	1.365557e-04	perCapInc	0.61739612	8.094139e-03
PctIlleg	0.17706825	5.296439e-03	PctSameState85	0.70297061	1.640443e-04	PctPolicHisp	0.10310956	1.323897e-02
perCapInc	0.64767388	5.985130e-03	racePctAsian	-0.42864970	1.679086e-04	PctPolicBlack	0.10155110	1.711592e-02
HousVacant	0.11791783	1.354201e-02	PolicCars	-0.47063098	2.489026e-04	PctNotSpeakEnglWell	-0.18507565	1.893450e-02
PctLargHouseFam	-0.17895486	1.957155e-02	IndianPerCap	0.42168616	1.627711e-03	PctRecImmig5	0.19178816	2.309826e-02
MedRent	0.21441193	2.281982e-02	blackPerCap	-0.71439300	1.986876e-03	PctHousOccup	-0.08330648	2.464276e-02
PctHousOwnOcc	0.12837453	3.994922e-02	AsianPerCap	-0.42157170	3.254307e-03	householdsSize	0.14967889	2.527427e-02
			MedNumBR	-0.18847571	3.558794e-03	NumStreet	0.15225780	3.565732e-02
			MedYrHousBuilt	-0.33892953	4.314384e-03	PctKids2Par	-0.19473892	4.149383e-02
			PctImmigrRec10	1.39855213	6.730074e-03	MedYrHousBuilt	-0.05952732	4.653226e-02
			PctEmploy	1.03740053	6.910475e-03			
			PctFam2Par	1.75461044	7.042760e-03			
			PctIlleg	0.55326000	7.784692e-03			
			RentLowQ	-1.16921722	8.783941e-03			
			LemasPctPoliconPatr	-0.29313752	1.013110e-02			
			whitePerCap	-0.66687450	1.186371e-02			
			medIncome	-1.43668182	2.464411e-02			
			(Intercept)	1.33256437	2.497551e-02			
			PctRecImmig5	1.08001715	2.616070e-02			
			pctWwage	-0.71244608	3.485095e-02			
			PersPerRentOccHous	-1.08775476	3.678513e-02			
			PctPolicAsian	0.15968419	4.508829e-02			
			PctHousOwnOcc	1.93489830	4.664628e-02			

(a) Jeu de données 1

(b) Jeu de données 2

(c) Jeu de données 3

**Figure 13: Liste des coefficients statistiquement significatifs (à 95%) et leur p-value associée calculés par les modèles linéaires aux variables sélectionnées itérativement**

étant arrivé il y a moins de 5, 8 et 10 ans), etc

Chacune de ces variables générales influe négativement sur le taux de crimes violents. Ces variables générales peuvent également s'influencer entre elles (il y a par exemple des corrélations entre la situation économique et la race des habitants). Enfin il convient de rappeler qu'il faudrait des analyse sociologiques théoriques comme empiriques plus poussées pour tenter d'établir des liens de causalité.