

Projet d'étude de Statistiques

Maxime Baba, Alexandre Demarquet, Félix de Brandois, Tristan Gay

2024-01-24

Contents

1	Introduction	2
2	Analyse descriptive des données	2
2.1	Analyse unidimensionnelle	2
2.2	Analyse multidimensionnelle	3
3	Classification des EPCI	6
3.1	Clustering	6
3.2	Analyse discriminante linéaire	6
4	EMS	8
4.1	Modèle linéaire	8
4.2	Modèle linéaire généralisé	9
5	Conclusion	9

List of Figures

1	Boxplot des variables nox_kg,co_kg,so2_kg	2
2	Histogramme de la variable co_kg en brute, scale et scale(log())	2
3	Corrélation entre les variables	3
4	ACP des variables quantitatives	5
5	MCA avec découpage des données en 3, 4 et 5 intervalles	6

1 Introduction

Le but de ce projet est d'étudier différents polluants mesurés par de nombreux EPCI d'Occitanie. Nous disposons du jeu de données suivant : `Data-projetmodIA-2324.csv`.

Dans la suite de ce rapport, on utilise les notations suivantes :

- a
- b
- c

2 Analyse descriptive des données

On commence par interpréter les éléments jeu de données.

2.1 Analyse unidimensionnelle

On s'intéresse dans un premier temps aux variables quantitatives du jeu de données (et en particulier aux émissions de polluants).

La figure 1 présente une visualisation de quelques variables quantitatives brutes.

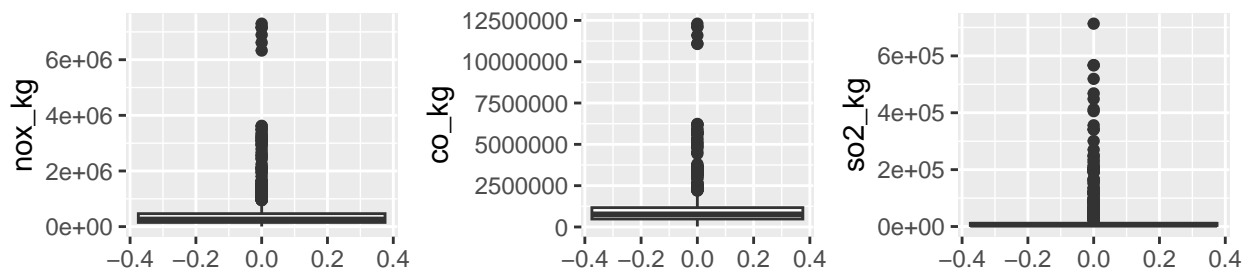


Figure 1: Boxplot des variables `nox_kg`, `co_kg`, `so2_kg`

On observe que ... (pas scale...)

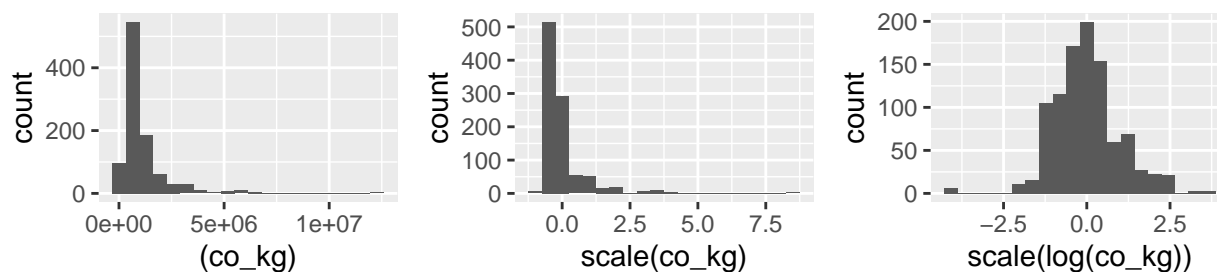


Figure 2: Histogramme de la variable `co_kg` en brute, `scale` et `scale(log())`

On effectue une transformation des données car d'après les boxplots de la figure 1 on remarque une variance énorme de certaines données comme `co_kg`. En examinant l'histogramme des données quantitatives, on observe une distribution fortement asymétrique. On peut donc appliquer une log-transformation pour normaliser la distribution des données.

Certaines variables ont pour unité la tonne et d'autre le kg on peut donc scale les données. On peut visualiser l'interet de ces transformations grâce à la figure 2 avec la variable `co_kg`.

Par la suite, on manipule les variables quantitatives transformées `scale(log())`.

On étudie ensuite la corrélation entre les variables quantitatives.

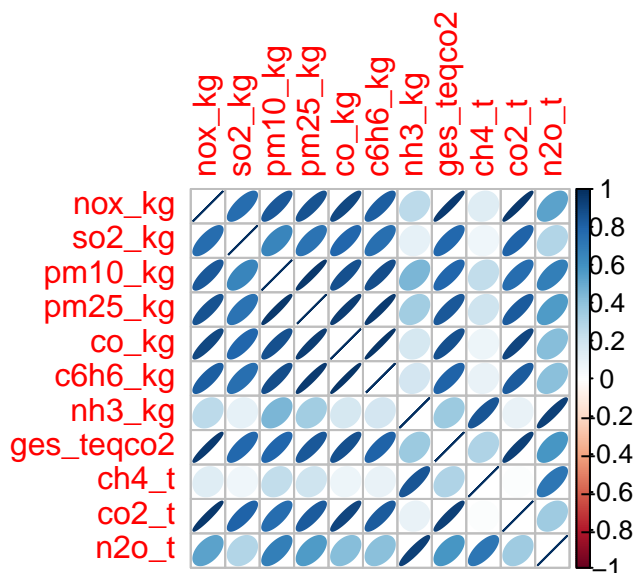


Figure 3: Corrélation entre les variables

L'analyse de la figure 3 nous permet d'identifier rapidement les relations significatives entre nos variables. Les ellipses fortement allongées suggèrent une corrélation plus forte, tandis que les ellipses plus circulaires indiquent une corrélation plus faible.

2.2 Analyse multidimensionnelle

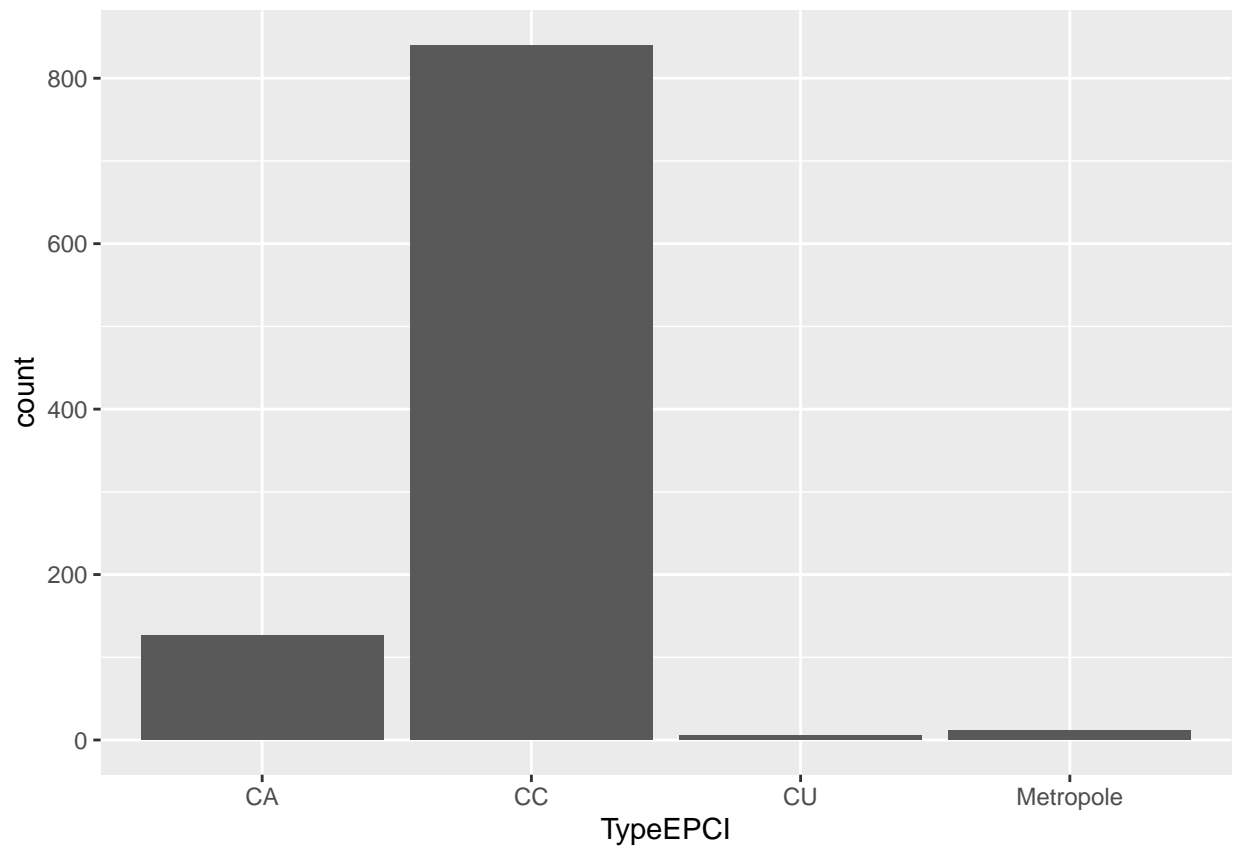
Dans le jeu de données nous avons aussi des variables qualitatives comme le code `epci`, le `lib_epci` ou des infos sur les départements.

```
data_quali=data[,c("code_epci","lib_epci","annee_inv","TypeEPCI","nomdepart")]
table(data_quali[,c("nomdepart")])
```

```
##
##                Ardèche,Gard                Ariège
##                6                48
##                Aude                Aude,Haute-Garonne,Tarn
##                48                6
##                Aude,Pyrénées-Orientales                Aveyron
##                6                102
##                Aveyron,Lot                Aveyron,Lozère
```

```
##          12          6
##          Gard          Gard,Hérault
##          78          6
##          Gard,Lozère          Gard,Vaucluse
##          6          6
##          Gers          Gers,Haute-Garonne
##          84          6
##          Gers,Landes Gers,Lot-et-Garonne,Tarn-et-Garonne
##          6          6
##          Haute-Garonne          Haute-Garonne,Tarn
##          96          6
##          Hautes-Pyrénées Hautes-Pyrénées,Pyrénées-Atlantiques
##          48          12
##          Hérault          Hérault,Tarn
##          90          6
##          Lot          Lozère
##          48          54
##          Pyrénées-Orientales          Tarn
##          66          72
##          Tarn-et-Garonne          Tarn,Tarn-et-Garonne
##          48          6
```

```
ggplot(data=data_quali)+geom_bar(aes(x = TypeEPCI))
```



2.2.1 Analyse en Composantes Principales (ACP) des variables quantitatives

On visualise les individus à partir des émissions de polluants.

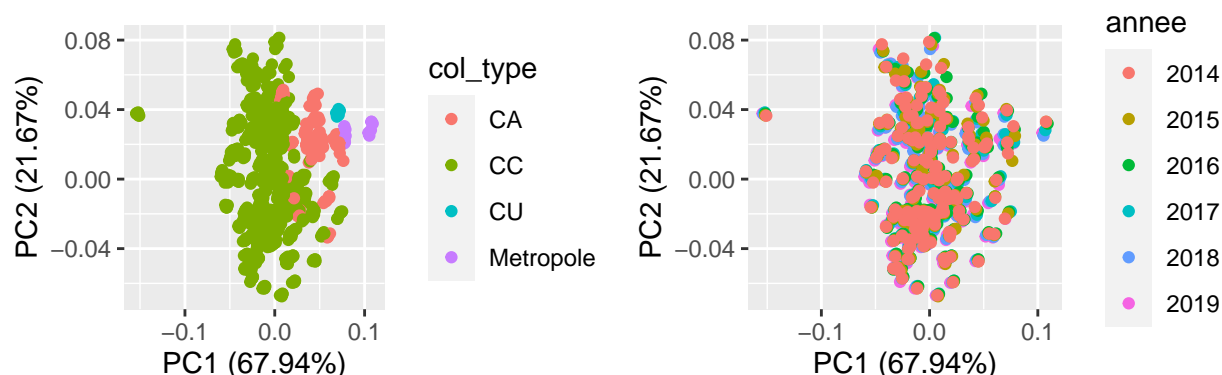


Figure 4: ACP des variables quantitatives

On observe sur la figure 4 que ...

2.2.2 Réduction de dimension (MCA)

Dans cette partie, on cherche à effectuer une réduction de dimension pour les polluants et du type EPCI. Nous allons donc utiliser une MCA (Multiple Correspondance Analysis). Les polluants sont des variables quantitatives nous avons donc besoin de discrétiser ces variables. Nous allons former un nombre fini d'intervalles qui formeront les modalités des nouvelles variables qualitatives. Nous allons aussi retirer les valeurs aberrantes c'est-à-dire en-dehors des quantiles (voir boxplot) : En effet, la MCA est sensible aux valeurs extrêmes car elle vise à maximiser la variance des données. Les outliers, en raison de leur nature inhabituelle, peuvent influencer significativement la variance et ainsi biaiser les résultats de l'analyse.

Les données quantitatives sont enrichies en incluant la colonne avec la variable qualitative, puis les données quantitatives sont transformées en données qualitatives afin de réaliser une Analyse en Composantes Principales (MCA) à l'aide de FactoMineR.

Ensuite, nous appliquons l'Analyse en Composantes Principales à l'aide de la bibliothèque factoMineR, en variant les intervalles de découpage des données quantitatives en données qualitatives.

```
## Warning: ggrepel: 34 unlabeled data points (too many overlaps). Consider  
## increasing max.overlaps
```

```
## Warning: ggrepel: 46 unlabeled data points (too many overlaps). Consider  
## increasing max.overlaps
```

L'analyse des résultats de la MCA révèle une structure significative lorsque les variables sont regroupées selon un découpage en trois intervalles. Dans ce scénario, les variables partageant le même découpage d'intervalles présentent un regroupement cohérent, suggérant une association claire entre ces catégories.

Les deux premiers axes principaux de l'Analyse en Composantes Principales (MCA) capturent un pourcentage significatif de la variance totale, avec des valeurs respectives de 27% et 17%. Ces résultats indiquent que ces axes fournissent une représentation robuste des relations entre les variables, soulignant des patterns structurés dans les données.

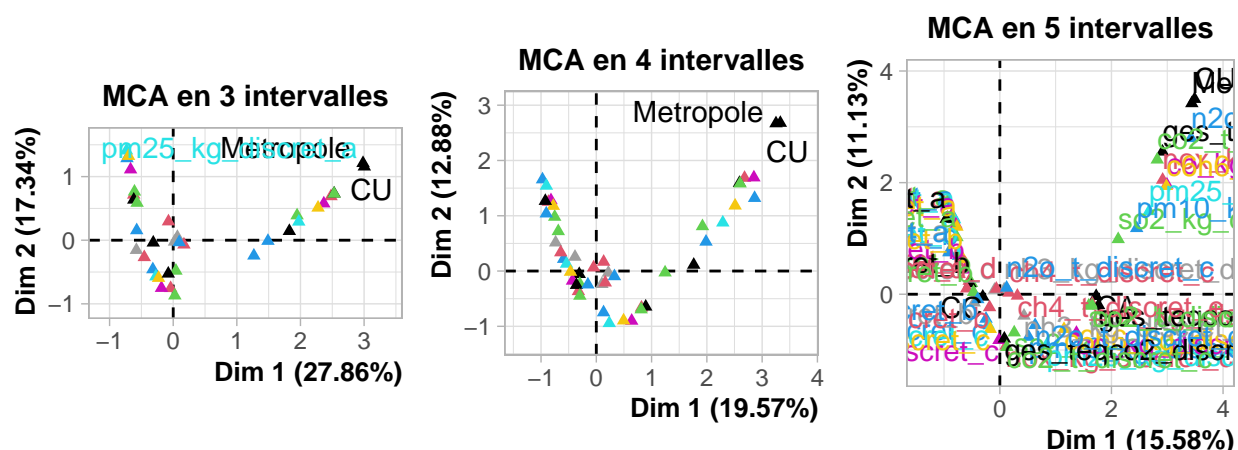


Figure 5: MCA avec découpage des données en 3, 4 et 5 intervalles

Cependant, lorsqu'on effectue un découpage en un plus grand nombre d'intervalles, les pourcentages associés aux axes principaux diminuent, suggérant une dispersion accrue des données. Cela peut être interprété comme une indication que le découpage en trois intervalles offre une simplification pertinente, condensant l'information tout en préservant la structure sous-jacente, tandis qu'un découpage plus fin pourrait introduire du bruit ou de la complexité excessive.

En résumé, l'analyse suggère que le découpage en trois intervalles optimise la représentation des variables, offrant une compréhension significative des relations dans les données, tandis qu'un découpage plus fin pourrait conduire à une perte de clarté et à une dilution de l'information utile.

3 Classification des EPCI

On cherche à classer les EPCI en fonction de leurs émissions de polluants. On utilise pour cela différentes méthodes de classification.

3.1 Clustering

On met en place différents algorithmes de clustering :

Blabla sur les méthodes de clustering

3.2 Analyse discriminante linéaire

Explication sur la méthode de l'analyse discriminante linéaire

3.2.1 LDA sur le dépassement d'émissions de méthane de 1000 tonnes par an

```
# On sépare les données en train et test (70% train, 30% test)
taille_train=round(0.7*nrow(data_mel2))
d_train=data_mel2[1:taille_train,]
d_test=data_mel2[taille_train:nrow(data_mel2),]

# On applique la LDA
lda_model <- lda(ch4_t ~ .,data=d_train)
```

```
## Warning in lda.default(x, grouping, ...): les variables sont colinéaires
```

```
# On colorie les points en fonction du dépassement ou non de $1000$ tonnes par an
color2 <- data_lda2$ch4_t ;
color2[color2=="TRUE"] <- "black";
color2[color2=="FALSE"] <- "red"

# Afficher les résultats de la LDA
#print(lda_model)
```

Commentaire sur les résultats obtenus.

3.2.2 LDA sur le type d'EPCI

```
## Call:
## lda(TypeEPCI ~ ., data = d_train)
##
## Prior probabilities of groups:
##          CA          CC          CU  Metropole
## 0.127721335 0.851959361 0.008708273 0.011611030
##
## Group means:
##          nox_kg    so2_kg    pm10_kg    pm25_kg    co_kg    c6h6_kg
## CA          1.325364  1.1611561  1.0012335  1.1355026  1.3289946  1.2178029
## CC         -0.220940 -0.1737318 -0.1796963 -0.2091062 -0.2369778 -0.2253901
## CU          2.230647  1.7879132  1.9975620  2.2075201  2.4962003  2.4671038
## Metropole   2.937184  2.2400560  2.4146530  2.6588233  2.9154855  2.6660273
##          nh3_kg ges_teqco2    ch4_t    co2_t    n2o_t
## CA          0.0526407206  1.2995476 -0.095701039  1.4070196  0.38221266
## CC         -0.0005481848 -0.2243859 -0.007160466 -0.2304081 -0.06621661
## CU         -0.5323667351  2.3399442 -0.097655937  2.2593576  0.58235390
## Metropole   0.0007450938  3.0051703  0.539421145  2.9011474  1.09396625
##
## Coefficients of linear discriminants:
##          LD1          LD2          LD3
## nox_kg      0.6524815 -1.54740044  1.93472949
## so2_kg      0.2469881 -0.08494908  0.13266331
## pm10_kg     -0.2034815  0.98645084 -0.71869623
## pm25_kg     0.5791775 -0.13097609  2.71491721
## co_kg      -1.5971732  0.75738057  4.29939222
```

```
## c6h6_kg      0.3523484 -0.97759416 -5.46199964
## nh3_kg       2.4583801 -5.60529147 -0.04955041
## ges_teqco2 -0.7300010  0.29650270 -2.71039308
## ch4_t       -0.7056128  1.84084802  1.30422054
## co2_t       -0.2364000 -0.59601327 -0.17929090
## n2o_t       -1.7782893  4.23073821 -0.68767128
##
## Proportion of trace:
##      LD1      LD2      LD3
## 0.8673 0.1225 0.0102
```

4 EMS

4.1 Modèle linéaire

4.1.1 Modèle d'ANOVA

On explique le gaz à effet de serre en fonction des variables Type et années.

On utilise un modèle d'ANOVA à deux facteurs avec interaction :

$$Y_{ij} = \mu + \alpha_i + \beta_j + \gamma_{ij} + \epsilon_{ij}$$

EXPLIQUER LA SIGNIFICATION DES TERMES DU MODELE

```
##
## Call:
## lm(formula = ges_teqco2 ~ TypeEPCI * annee_inv, data = dlog)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.2367 -0.4233 -0.0383  0.3863  2.8469
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    -18.397236   80.407295  -0.229    0.819
## TypeEPCICC       4.064479   86.227217   0.047    0.962
## TypeEPCICU       7.818578  377.143642   0.021    0.983
## TypeEPCIMetropole -19.949436  272.674402  -0.073    0.942
## annee_inv        0.009761   0.039875   0.245    0.807
## TypeEPCICC:annee_inv -0.002779   0.042761  -0.065    0.948
## TypeEPCICU:annee_inv -0.003354   0.187029  -0.018    0.986
## TypeEPCIMetropole:annee_inv  0.010806   0.135222   0.080    0.936
##
## Residual standard error: 0.7644 on 976 degrees of freedom
## Multiple R-squared:  0.4198, Adjusted R-squared:  0.4157
## F-statistic: 100.9 on 7 and 976 DF,  p-value: < 2.2e-16
```

-> Commentaire sur la valeur de R^2 obtenue.

On essaie de simplifier le modèle en enlevant les interactions avec un test de sous-modèle :

$$\begin{aligned}\mathcal{H}_0 : & Y_{ij} = \mu + \alpha_i + \beta_j + \epsilon_{ij} \\ \mathcal{H}_1 : & Y_{ij} = \mu + \alpha_i + \beta_j + \gamma_{ij} + \epsilon_{ij}\end{aligned}$$

On obtient une p-value de $1 > 0.05$.

On ne rejette pas l'hypothèse de nullité des interactions.

On garde donc le modèle suivant :

$$Y_{ij} = \mu + \alpha_i + \beta_j + \epsilon_{ij}$$

On essaie de simplifier le modèle en enlevant les variables non significatives (on fait 2 tests de sous-modèle) :

$$\begin{aligned}\mathcal{H}_0 : & Y_{ij} = \mu + \alpha_i + \epsilon_{ij} \\ \mathcal{H}_1 : & Y_{ij} = \mu + \alpha_i + \beta_j + \epsilon_{ij}\end{aligned}$$

et

$$\begin{aligned}\mathcal{H}_0 : & Y_{ij} = \mu + \beta_j + \epsilon_{ij} \\ \mathcal{H}_1 : & Y_{ij} = \mu + \alpha_i + \beta_j + \epsilon_{ij}\end{aligned}$$

Pour le modèle dépendant uniquement du type d'EPCI, on obtient une p-value de $0.599 > 0.05$.

On peut donc enlever l'année dans le modèle :

$$Y_{ij} = \mu + \alpha_i + \epsilon_{ij}$$

On essaie à nouveau de simplifier le modèle en enlevant les variables non significatives.

On obtient cette fois une p-value de $0 < 0.05$.

On ne peut donc pas enlever le type d'EPCI dans le modèle.

On vérifie finalement la cohérence du modèle retenu :

On obtient une p-value de $0.99 > 0.05$ donc le modèle est cohérent. On garde donc le modèle :

4.1.2 Régression linéaire

4.1.3 ANCOVA

4.2 Modèle linéaire généralisé

5 Conclusion