

Analyse de données
M1 Statistique et économétrie - 2013-2014
V. Monbet
Exploration des données avec SAS

L'objectif de ce TD est de revoir comment mener une analyse descriptive simple sous SAS pour un ensemble d'observations de variables aléatoires quantitatives,

1 Les données

Les données (fichier `depart.dat` sur la page du cours http://perso.univ-rennes1.fr/valerie.monbet/Cours_AD/AD.html) proviennent du Groupe d'Etude et de Reflexion Inter-régional (GERI). Elles portent sur 4 grands thèmes : la démographie, l'emploi, la fiscalité directe locale, la criminalité. Les indicateurs sont mesurés sur l'ensemble des départements français métropolitains ainsi que la Corse pendant l'année 1990, ils sont, pour la plupart, des taux calculés relativement à la population totale du département concerné. On observe les variables suivantes :

- numéro de département,
- code du département,
- code de la région,
- URBR indicateur de concentration de la population mesurant le caractère urbain ou rural du département,
- TXCR taux de croissance de la population sur la période intercensitaire 1983-1990,
- JEUN part des 0-19 ans dans la population totale,
- AGE part des plus de 65 ans dans la population totale,
- FE90 taux de fécondité (pour 1000) égal au nombre de naissances rapporté au nombre de femmes fécondes (15 à 49 ans) en moyenne triennale,
- ETRA part des étrangers dans la population totale,
- CHOM taux de chômage,
- CRIM taux de criminalité : nombre de délits par habitant,
- FISC produit, en francs constants 1990 et par habitant des quatre taxes locales (professionnelle, habitation, foncier bâti, foncier non bâti).

On observe également les parts de chaque profession en catégorie socioprofessionnelle (PCS) dans la population active occupée du département :

- AGRI : agriculteurs,
- ARTI : artisans,
- CADR : cadres supérieurs,
- EMPL : employés,
- OUVR : ouvriers,
- PROF : professions intermédiaires.

2 Exploration avec le logiciel SAS

2.1 Lecture des données Tous

Utiliser, par exemple, les instructions suivantes pour lire le fichier de données.

```
libname TPexplor '~/AnalyseDonnees/TP1' ;
data TPexplor.depart ;
infile '~/AnalyseDonnees/TP1/depart.dat' ;
input num $ depart $ region $ txcr etra urbr jeun age chom agri
      arti cadr empl ouvr prof fisc crim fe90 ;
run;
```

2.2 Analyse interactive des données

En utilisant l'outil d'analyse interactive les données (ou SAS Insight)

- Visualiser les distributions des différentes variables (histogrammes, boîtes à moustaches).
- Tracer des nuages de points des variables deux à deux (scatter plot).
- Etudier, rapidement, la relation linéaire entre la variable criminalité et les autres (fit).

Choisir dans le menu déroulant **Solution** la ligne **Analyse Interactive des données** puis sélectionner les variables à étudier et utiliser ensuite le menu **Analyze**)

2.3 PROC UNIVARIATE et PROC CORR

Répondre de nouveau aux questions précédentes à l'aide des procédures PROC UNIVARIATE et PROC CORR.

Dans la procédure UNIVARIATE l'option KERNEL permet d'obtenir le graphe de l'estimateur à noyau, l'option K= permet de choisir le noyau et l'option C= la largeur de fenêtre standardisée. Pour C=, on peut choisir une valeur à la main ou choisir C=MISE. Dans le premier cas, on trace l'estimation pour plusieurs valeurs de C et on retient celle qui nous semble la plus raisonnable. Dans le second cas, la largeur de fenêtre est choisie telle que celle ci minimise le critère AMISE (approximate mean integrated square error) pour une loi de Gauss ayant la moyenne et la variance estimées dans l'échantillon.

Tester les deux approches et commenter.

Rq : le plus souvent on combine ces deux approches, la seconde donnant une valeur initiale cohérente pour la première.

```
proc univariate data = TPexplor.depart plots ;
var crim empl ouvr fe90 jeun ;
run ;
```

```
proc univariate data = TPexplor.depart plots ;
histogram crim empl ouvr fe90 jeun ;
run ;
```

```
proc univariate data = TPexplor.depart plots ;  
  histogram crim / KERNEL (K=NORMAL C=.3);  
run;
```

On peut aussi utiliser cette fonction pour ajuster des modèles de loi. Les paramètres sont alors estimés par maximum de vraisemblance.

```
proc univariate data = TPexplor.depart plots ;  
  histogram crim / GAMMA (THETA=EST ALPHA=EST SIGMA=EST);  
run;
```

```
ods graphics on;  
title 'Données des crimes';  
proc corr data=TPexplor.depart plots/*=matrix(histogram)*/;  
  var crim empl ouvr fe90 jeun ;  
run;  
ods graphics off;
```