

Analyse prédictive et factorielle sous SCILAB



Pierrick Michel
Wahiba Azzoug
EL Mehdi EL ALAMI KAABOUC
Kafil ELkhadir

Master 2 Statistiques et Informatique pour la science des données 2017-2018

Plan

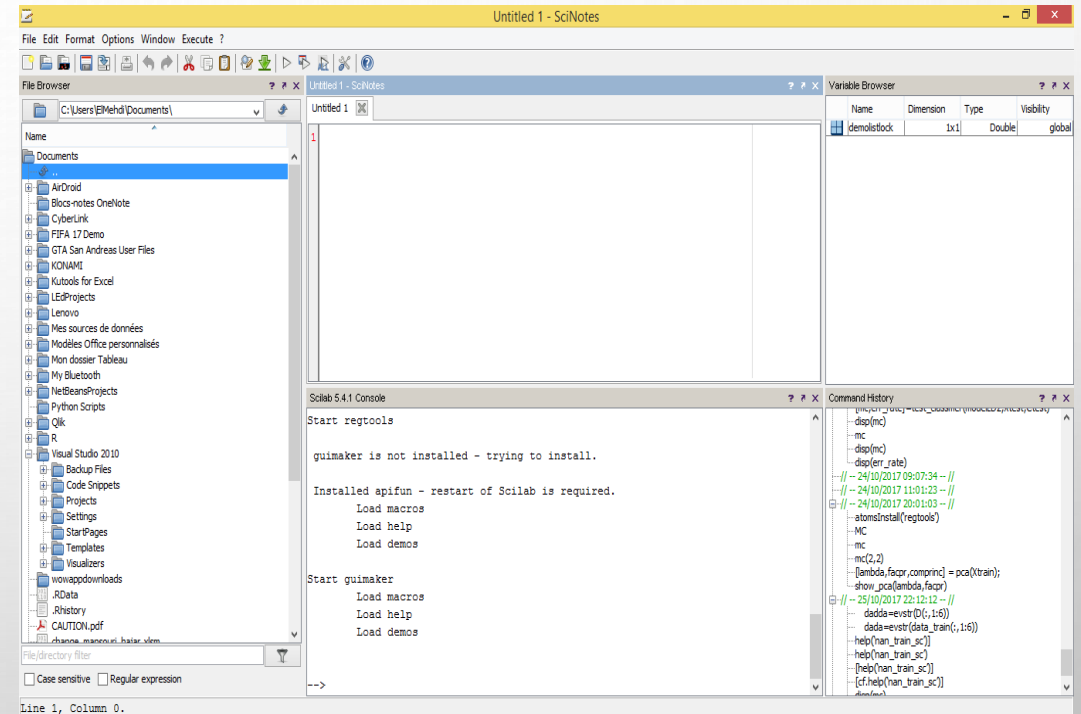
- 1 Introduction
- 2 Analyse statistique avec Scilab
- 3 Régression linéaire simple et multiple
- 4 Analyse Discriminante
- 5 **Analyse en composantes principales (Exemple)**

1. Introduction

Dans le cadre des ateliers techniques de formation aux outils, nous avons voulu réaliser des tests statistiques à l'aide du logiciel Scilab.

Ce logiciel ressemble beaucoup à Rstudio en termes d'interface et de capacités à réaliser des analyses statistiques

Utilisation de la version 5.4.1 plutôt que 6 car packages non mis à jour.



1. Introduction

Téléchargement en ligne du package utile pour nos analyses,
installation du package:

```
1 | atomsInstall("C:\Users\Desktop\M2-SISE\Scilab\nan_1.3.4-1.bin.x64.windows.zip")
```

Les fonctions du package

Data Correlation and Covariance

```
"nan_covm"  
"nan_ecovm"  
"nan_decovm"  
"nan_xcovf"  
"nan_conv2nan"  
"nan_cor"  
"nan_cov"  
"nan_corrcoef"  
"nan_corrcof"  
"nan_rankcorr"  
"nan_partcorrcoef"  
"nan_tiedrank"
```

hypothesis Tests

```
"nan_ttest"  
"nan_ttest2"
```

Classification

```
"nan_train_sc"  
"nan_test_sc"  
"nan_classify"  
"nan_xval"  
"nan_kappa"  
"nan_train_lda_sparse"  
"nan_fss"  
"nan_cat2bin"  
"nan_row_col_deletion"  
"nan_mahal"  
"nan_rocplot"  
"nan_confusionmat"  
"nan_partest"
```

cluster Analysis

```
"nan_kmeans"  
"nan_pdist"  
"nan_linkage"  
"nan_squareform"  
"nan_cluster"  
nan_crosstab"
```

Il existe peu de package de
statistique sous Scilab,

- Nan
- Libvsm
- Regtools

Les packages ne sont pas remis
à jour pour la dernière versions
de Scilab.

2. Analyse statistiques avec Scilab



Les analyses que nous avons choisi de traiter sont :

- Régression linéaire simple
- Régression linéaire multiple
- Analyse discriminante

Il existe des packages sous Scilab pour réaliser ces analyses.

Le but étant d'apprendre à coder avec ce logiciel, seuls les packages «nan» et «libsvm» sont nécessaires.

3. Régression linéaire simple et multiple

Y est une variable quantitative qui est expliquée, par une ou plusieurs variables quantitatives $X_j(1, \dots, p)$

Notre objectif est :

- estimer les paramètres du modèle par la méthode des moindres carrés

$$\hat{a} = (X'X)^{-1}X'Y$$

- évaluer globalement la pertinence du modèle

$$F_{calc} = \frac{\frac{R^2}{p}}{\frac{1-R^2}{n-p-1}},$$

Jeu de données :

Parmi l'un des polluants dans l'air nous retrouvons le NO₂. Selon la concentration d'O₃ dans l'air et des températures journalières nous allons prédire la concentration de ce gaz.

```
1 function [b] = inverse(X)
2
3 *****Pré-traitement-
4 *****Calculer l'inverse
5 b = inv(X'*X)
6
7 endfunction
8
9 ////-DATA-
10 //Variables-
11 X1=strtod(data_ozone(1:61,2:2))
12 X2=strtod(data_ozone(1:61,3:3))
13 X3=strtod(data_ozone(1:61,4:4))
14 X4=strtod(data_ozone(1:61,5:5))
15 X5=strtod(data_ozone(1:61,6:6))
16
17 //Création du matrice
18 X=[ones(61,1) X1 X2 X3 X4 X5]//X6-X7]
```

4. L'analyse discriminante

→ Est Une extension de la régression dans le cas où la variable à expliquer est qualitative.

L'analyse discriminante sous SCILAB :

- (1) charger l'échantillon d'apprentissage
- (2) le coder en implémentant une fonction

- (3) Construction du modèle prédictif avec la Toolbox « Nan » - une librairie de l'apprentissage automatique.

```
datatype: "classifier:statistical:ld2"  
Labels: [1,2]  
MD: hypermat  
NN: hypermat  
weights: [7x2 constant]
```

Evaluation du modèle sur un échantillon test

- (1) charger l'échantillon test
- (2) le coder en respectant le schéma utilisé pour le fichier d'apprentissage
- (3) produire la prédiction en appliquant le modèle sur l'échantillon test,
- (4) créer une fonction pour la matrice de confusion en confrontant les valeurs prédites par le modèle avec celles observées dans l'échantillon test
- (5) déduire de la matrice de confusion le taux d'erreur.

- (1) `Data_train=csvRead("C:\Users\EIMehdi\Desktop\SCILAB\data.txt","t",".", "string",[],[],[2 1 150 9])`
- (2) `function [descripteurs, cible]=recodage(base)`
..
`descripteurs: evstr(base(:,1:8)) //transformer les colonnes 1 à 8 en numérique`
`y=base(:,9:9) //Récupération de la variable cible`

- (3) **Il faut avoir installé et chargé la Toolbox « Nan » pour pouvoir poursuivre.**
`atomsInstall('nom de la toolbox') // pour les installer`
`atomsLoad('nom de la toolbox') // pour les charger`
`modelLD2=nan_train_sc(XTrain,CTrain,'LD2') //apprentissage - méthode LD2`
`disp(modelLD2)`

```
//chargement des données - format chaîne  
Data_test=csvRead("C:\Users\EIMehdi\Desktop\SCILAB\data.txt","t",".", "string",[],[],[2 1 60 9])  
//le coder en utilisant la fonction recodage  
[XTest,CTest]=recodage(Data_test)  
disp(tabul(data_test)) // Pour afficher le nombre d'observation dans chaque classe.
```

- (3) et (4) et (5)
`function [MC, ERR_RATE]=test_classifier(classifier, descripteurs, cible) //Fonction pour l'évaluation d'un classifieur`
`pred=nan_test_sc(classifier,descripteurs) //prédiction`
`//matrice de confusion : cible observée vs. prédiction`
`MC=nan_confusionmat(cible,pred.classlabel)`
`//taux d'erreur`
`ERR_RATE = 1.0 - sum(diag(MC))/sum(MC)`
`endfunction`
`//évaluation`
`[mc,err_rate] = test_classifier(modelLD2,XTest,CTest)`

4. L'analyse discriminante

→ **L'objectif** : Prédire la classe d'un nouvel objet décrit par la valeur de ces attributs.

Exemple :

Donner une prévision d'avoir un cancer de prostate ou pas ,d'une personne de 64 ans ayant un niveau de sérum de 0.4 et qui a des résultats d'une radiographie positive et sa taille grande ainsi leur résultat d'une biopsie moins sérieux .

Résultats de l'analyse discriminante sous SCILAB : `disp(modelD2.weights)` → Pour accéder aux champs de l'objets `weights`.

```
datatype: "classifier:statistical:ld2"  
Labels: [1,2]  
MD: hypermat  
NN: hypermat  
weights: [7x2 constant]  
  
0.5330020 - 0.5330020  
- 0.0075632 0.0075632  
- 0.0060413 0.0060413  
0.4370939 - 0.4370939  
0.4058681 - 0.4058681  
0.7191982 - 0.7191982
```

$$Y = a_0 + a_1 * \text{age} + a_2 * \text{nv_serum} + a_3 * \text{radio} + a_4 * \text{taille} + a_5 * \text{biop}$$

5. Analyse en composantes principales

Exemple :

```
a=rand(100,10,'n'); [lambda,facpr,comprinc] = pca(a);  
show_pca(lambda,facpr)
```

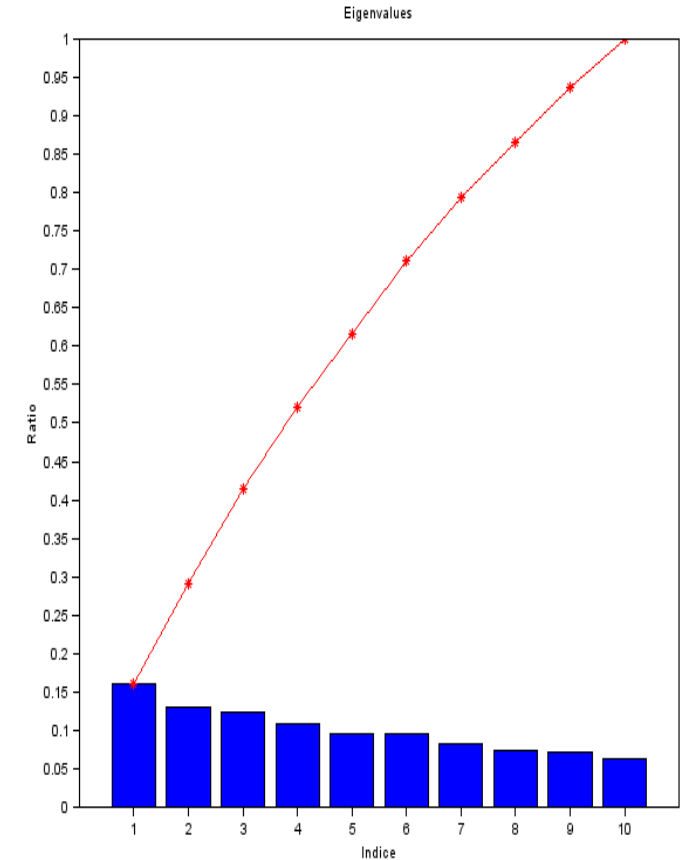
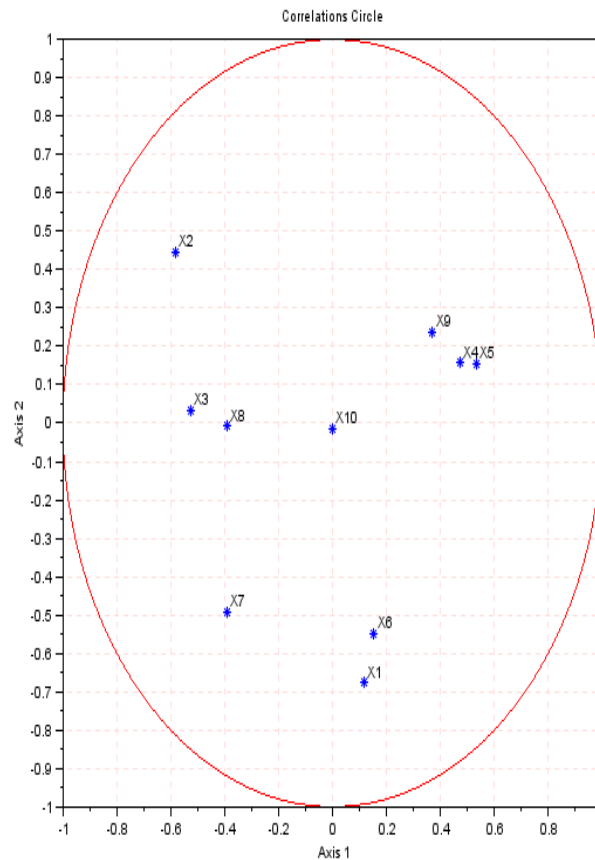
Lambda : est une matrice numérique $p \times 2$. Dans la première colonne, nous trouvons les valeurs propres de V , où V est la matrice de corrélation $p \times p$ et dans la seconde colonne, les rapports de la valeur propre correspondante sur la somme des valeurs propres.

Facpr : sont les principaux facteurs: vecteurs propres de V . Chaque colonne est un vecteur propre du dual de \mathbb{R}^p .

comprinc : sont les principaux composants. Chaque colonne ($c_i = X_{u_i}$) de cette matrice $n \times n$ est la projection M -orthogonale des individus sur l'axe principal. Chacune de ces colonnes est une combinaison linéaire des variables x_1, \dots, x_p avec la variance maximale sous condition $u_i^T M^{-1} u_i = 1$

show_pca - Visualisation des résultats de l'analyse des composants principaux

princomp - Analyse en composantes principales



The background features a low-poly, geometric design. On the left side, there is a large, intricate cluster of triangles in various shades of blue, ranging from light sky blue to deep navy blue. These triangles are arranged in a way that creates a sense of depth and crystalline structure. Towards the bottom center, there is a smaller, more compact cluster of triangles in shades of teal and green. The rest of the background is a plain, light gray with very subtle, larger-scale geometric patterns that echo the low-poly style of the main clusters.

**MERCI POUR VOTRE
ATTENTION !**