

Examen de contrôle continu

Durée 1h - Sans document
Vendredi 20 octobre 2023

Exercice 1 (10 pts)

Décrire et commenter les résultats obtenus à l'aide du code Python ci-dessous.

```
import os
import pandas as pd
import numpy as np

from sklearn import preprocessing
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import RepeatedKFold
from sklearn.model_selection import GridSearchCV

path = os.getcwd()
os.chdir('/Users/marin/TEACHING/2324/M2-GLM-HAX912X/TP')

data = pd.read_csv("creditcard.csv")
data.shape

(284807, 31)

count_classes = pd.value_counts(data['Class'], sort = True).sort_index()
count_classes

0    284315
1     492

X = data.loc[:, data.columns != 'Class']
y = data['Class']
scaler = preprocessing.StandardScaler().fit(X)
X = scaler.transform(X)
X_train, X_test, y_train, y_test = train_test_split(X,y,test_size = 0.3,
random_state = 1974)

print("Number transactions train dataset: ", len(X_train))
print("Number transactions test dataset: ", len(X_test))
print("Total number of transactions: ", len(X_train)+len(X_test))
```

```

Number transactions train dataset: 199364
Number transactions test dataset: 85443
Total number of transactions: 284807

```

```

logit = LogisticRegression()
logit.fit(X_train,y_train)
y_chap = logit.predict(X_test)
table = pd.crosstab(y_test,y_chap)
table

```

```

      0      1
0  85272  19
1   48    104

```

```

Crossvalid = RepeatedKfold(n_splits=2, n_repeats=2)
parameters = {'C':np.linspace(0.1,2.1,num=5)}
logit = GridSearchCV(LogisticRegression(), parameters,
                    cv=Crossvalid,scoring='neg_log_loss')
logit.fit(X_train,y_train)
print(logit.best_params_)

```

```
{'C': 0.1}
```

```

y_chap = logit.predict(X_test)
table=pd.crosstab(y_test,y_chap)
table

```

```

      0      1
0  85273  18
1   49    103

```

Exercice 2 (10 pts)

Nous considérons n variables aléatoires indépendantes y_1, \dots, y_n telles que y_i est distribué suivant une loi de Pareto de paramètre $(1, k_i)$ où $k_i > 0$. La densité de la loi de Pareto de paramètre $(1, k)$ est telle que

$$f(y; k) = \frac{k}{y^{k+1}} \mathbf{1}_{[1, +\infty[}(y).$$

Pour tout $i = 1, \dots, n$, nous supposons que $\log(k_i) = \beta_0 + \beta_1 x_i$ où $x_i \in \mathbb{R}$ est supposé connu, β_0 et β_1 sont des paramètres inconnus.

1 (4 pts) Montrer qu'il s'agit d'un modèle linéaire généralisé. La fonction de lien canonique a-t-elle été utilisée?

2 (6 pts) Donner la log-vraisemblance et les équations de vraisemblance. Pouvons-nous calculer les expressions analytique des estimateurs du maximum de vraisemblance de β_0 et β_1 ?

Correction E-women part 1

20 octobre 2023

HAX912X

(1)

Exercice 1

6 points répartis sur les questions

1] Importation données

284 307 individus

31 variables

(1 pt)

2]

Variable `primin` class
global déséquilibre entre
les modalités 0 et 1

(1 pt)

3]

`y` variable `primin` à
expliquer class

`X` matrice des variables
explicatives

les variables contenues dans `X`
sont centres réduites.

(1 pt)

4] Les données sont séparées dans
un échantillon d'apprentissage
contenant 199 344 individus
et un échantillon de test
contenant 85 443 individus (2)
(1 pt)

5] Le tableau de confusion
est mis en œuvre par défaut
le matériau de confusion sur
l'échantillon de test montre
surtout 67 mal classés parmi
85 443 sur l'ensemble
jeu (#8 par 10 000)
sur l'ensemble - beaucoup plus
fort parmi la classe 2 =
48 sur 152 (#32 par 100)

(3 pts)

6] Problème de régression logistique
régularisée avec le paramètre
de régularisation choisi
par validation croisée à
2 ensembles répété 2 fois
sur 5 valeurs de C régulièrement
espacés entre 0,1 et 2,1 et
pour le critère de la log-vraisemblance
négative du modèle de régression

(3pts) logistique. Le paramètre de
régularisation choisi est $C = 0,2$
c'est un peu inhabituel car on
trouve dans la littérature de recherche
les résultats sur l'échantillon de
test sont très similaires à ceux
obtenus sans validation de C

Exercice 2

(4)

$$\Delta) f(y; k) = \frac{k}{y^{k+1}} \mathbb{1}(y) \mathbb{1}(dy)$$

$$f(y; k) = \exp\left\{-k \log(y) + \log(k)\right\} \frac{1}{y} \mathbb{1}(dy)$$

Changement de variable $v(dy)$

$$u = \log(y) \Rightarrow y = e^u$$

$$du = \left(\frac{1}{y}\right) dy$$

mesure m de l'axe réel

$$f(u; k) = k e^{-\frac{1}{2}u} \mathbb{1}(u) \mathbb{1}(du)$$

$$u, v \text{ sur } (k)$$

$$\mathbb{P}_k(du) = \exp(-ku + \log(k)) \mathbb{1}(u) \mathbb{1}(du)$$

$$\theta = -k$$

$$h(\theta) = -\log(1 - \theta)$$

$v(dy)$

Estimons,

$$\log(k) = \beta_0 + \beta_1 K \Rightarrow k = e^{\beta_0 + \beta_1 K} \quad (5)$$

$$E(u) = \frac{1}{k} = e^{-\beta_0 - \beta_1 K} = f^{-1}(\beta_0 + \beta_1 K)$$

Pratiquement le changement de variable $u = \ln(y)$, se fait d'un modèle linéaire généralisé appelé la fonction de lien

$$f(z) = -\log(z) \quad (f^{-1}(z) = e^{-z})$$

$$f(z) \neq f'(z) = -\frac{1}{z}$$

ce n'est pas le bon choix que l'on a pu utiliser.

$$\begin{aligned} \underline{\Delta} \quad L(\beta_0, \beta_1) &= \sum_{i=1}^M [\beta_0 + \beta_1 x_i] \\ &- \sum_{i=1}^M e^{\beta_0 + \beta_1 x_i} \log(y_i) \end{aligned}$$

$$\frac{dL}{d\beta_0}(\beta_0, \beta_1) = M - \sum_{i=1}^M \log(y_i) e^{(\beta_0 + \beta_1 x_i)}$$

6

$$\frac{dL}{d\beta_1}(\beta_0, \beta_1) = \sum_{i=1}^M x_i - \sum_{i=1}^M x_i \log(y_i) e^{(\beta_0 + \beta_1 x_i)}$$

Nous ne pouvons pas calculer les expressions explicites des estimateurs du maximum de vraisemblance.