

Examen final

Durée 2h - Sans document
Vendredi 27 octobre 2023

Exercice 1 (4 pts)

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

Les données considérées sont les résultats d'une analyse chimique de vins cultivés dans la même région d'Italie par trois cultivateurs différents. Treize mesures ont été effectuées pour différents constituants présents dans les trois types de vin.

```
import numpy as np
import sklearn.linear_model as lm
from sklearn import preprocessing
from sklearn import datasets
dataset = datasets.load_wine()
X = dataset.data
X.shape
(178, 13)

y = dataset.target
np.unique(y)
array([0, 1, 2])

X = X[y != 0]
y = y[y != 0]
X.shape
(119, 13)

np.unique(y)
array([1, 2])

scaler = preprocessing.StandardScaler().fit(X)
X = scaler.transform(X)
X.mean(axis=0)
array([ 1.41810000e-15,  4.73943947e-16,  6.67253367e-15,  3.69452368e-16,
       1.20293598e-16, -2.21811365e-16, -4.09569670e-16, -6.69865657e-16,
      -2.00120033e-16, -4.31027763e-16,  7.57563946e-16, -9.19899078e-16,
       2.78721957e-17])

X.var(axis=0)
array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1.])
```

```

model1 = lm.LogisticRegression()
model1.fit(X,y)
model1.penalty
'12'
model1.C
1.0
model1.intercept_
array([-1.93238543])
model1.coef_
array([[ 0.67061806,  0.54717719,  0.63550791,  0.24868629,  0.16614675,
       -0.20010196, -1.34821171, -0.09664867, -0.56443121,  1.26159206,
      -1.16294376, -0.99888229,  0.33789776]])

from sklearn.model_selection import cross_val_score
from sklearn.model_selection import RepeatedKFold
Sel = RepeatedKFold(n_splits=5, n_repeats=100)
scores = cross_val_score(lm.LogisticRegression(), X, y, cv=Sel)
scores.shape
(500,)
scores.mean()
0.9775869565217391

from sklearn.model_selection import GridSearchCV
parameters = {'C':np.logspace(-4, 4, 50), 'penalty':['l1','l2']}
model2 = GridSearchCV(lm.LogisticRegression(), parameters, cv=Sel)

model2.fit(X,y)
model2.best_params_
{'C': 0.02811768697974228, 'penalty': 'l2'}

model2.best_score_
0.9889710144927536

```

Exercice 2 (6 pts)

Décrire et commenter les résultats obtenus à l'aide du code Python ci-dessous. Les données considérées sont décrites en annexe.

```

import os
import numpy as np

import pandas as pd
from sklearn.linear_model import LogisticRegression
from sklearn import preprocessing
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from math import sqrt
from sklearn.metrics import classification_report

```

```

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

data = pd.read_csv('Breast-Cancer-Wisconsin.csv')
y = data['diagnosis']
X = data.drop(['id','diagnosis','Unnamed: 32'], axis = 1)
X.head()

  radius_mean  texture_mean  ...  symmetry_worst  fractal_dimension_worst
0      17.99        10.38  ...          0.4601            0.11890
1      20.57        17.77  ...          0.2750            0.08902
2      19.69        21.25  ...          0.3613            0.08758
3      11.42        20.38  ...          0.6638            0.17300
4      20.29        14.34  ...          0.2364            0.07678
[5 rows x 30 columns]

pd.value_counts(y, sort = True).sort_index()

B    357
M    212
Name: diagnosis, dtype: int64

scaler = preprocessing.StandardScaler().fit(X)
X = scaler.transform(X)

X_train, X_test, y_train, y_test = train_test_split(X,y,test_size = 0.2,
random_state = 2023)

pd.value_counts(y_train, sort = True).sort_index()

B    286
M    169
Name: diagnosis, dtype: int64

pd.value_counts(y_test, sort = True).sort_index()

B    71
M    43
Name: diagnosis, dtype: int64

forest = RandomForestClassifier(n_estimators=500, max_features=int(sqrt(30)),
oob_score=True,n_jobs=-1)
forest = forest.fit(X_train,y_train)
y_chap = forest.predict(X_test)
print(classification_report(y_test,y_chap))

      precision    recall   f1-score   support
B        0.95     0.99     0.97       71
M        0.97     0.91     0.94       43

accuracy                           0.96      114
macro avg       0.96     0.95     0.95      114

```

```

weighted avg      0.96      0.96      0.96      114

forest.oob_score_
0.9604395604395605

forest.oob_decision_function_
array([[0.        , 1.        ],
       [1.        , 0.        ],
       [0.05747126, 0.94252874],
       [0.        , 1.        ],
       [1.        , 0.        ],
       [0.99456522, 0.00543478],
       ...
       [1.        , 0.        ],
       [0.00534759, 0.99465241],
       [0.96907216, 0.03092784],
       ...
       [1.        , 0.        ],
       [1.        , 0.        ],
       [0.98809524, 0.01190476],
       [0.99450549, 0.00549451],
       [0.9893617 , 0.0106383 ],
       [0.12429379, 0.87570621],
       [0.56666667, 0.43333333],
       [0.        , 1.        ],
       [1.        , 0.        ],
       [0.87434555, 0.12565445],
       [1.        , 0.        ],
       [0.34104046, 0.65895954],
       [1.        , 0.        ],
       [1.        , 0.        ],
       [0.        , 1.        ],
       [0.9902439 , 0.0097561 ],
       [1.        , 0.        ],
       [1.        , 0.        ],
       [0.98969072, 0.01030928],
       [0.97354497, 0.02645503],
       [0.95854922, 0.04145078]])

```

Exercice 3 (8 pts)

On considère un n -échantillons $(y_1, x_1), \dots, (y_n, x_n)$ du couple (y, x) où y est une variable aléatoire à valeurs dans $\{0, 1\}$ et x est une variable fixée prenant également les valeurs $\{0, 1\}$. On suppose qu'il existe une variable latente y^* telle que

$$f(y^*|x) = \frac{\exp(-(y^* - \beta_0 - \beta_1 x))}{(1 + \exp(-(y^* - \beta_0 - \beta_1 x)))^2}$$

et

$$y = \begin{cases} 1 & \text{si } y^* \leq 0 \\ 0 & \text{si } y^* > 0 \end{cases}.$$

1 (3 pts) Montrer qu'il s'agit d'un modèle linéaire généralisée et l'expliciter.

Indication : la primitive de $\frac{\exp(-(y - \mu))}{(1 + \exp(-(y - \mu)))^2}$ est $\frac{1}{1 + \exp(-(y - \mu))}$.

2 (2 pts) Est-ce le lien canonique qui a été choisi ?

3 (3 pts) Nous disposons de 100 observations (y_i, x_i) décrites dans la table de contingence suivante

		$y = 0$	$y = 1$
		30	20
$x = 0$	30	20	
$x = 1$	20	30	

Donner la log-vraisemblance en fonction de β_0 et β_1 .

Exercice 4 (2 pts)

Décrire la méthode de construction d'arbre de régression.

Annexe

Breast cancer data (Données sur le cancer du sein)

About Dataset

Features are computed from a digitized image of a fine needle aspirate (FNA) of a breast mass. They describe characteristics of the cell nuclei present in the image. In the 3-dimensional space is that described in:

[K. P. Bennett and O. L. Mangasarian: "Robust Linear Programming Discrimination of Two Linearly Inseparable Sets", Optimization Methods and Software 1, 1992, 23-34].

Attribute Information:

- 1) ID number
- 2) Diagnosis (M = malignant, B = benign)

3-32)

Ten real-valued features are computed for each cell nucleus:

- a) radius (mean of distances from center to points on the perimeter)
- b) texture (standard deviation of gray-scale values)
- c) perimeter
- d) area
- e) smoothness (local variation in radius lengths)
- f) compactness ($\text{perimeter}^2 / \text{area} - 1.0$)

- g) concavity (severity of concave portions of the contour)
- h) concave points (number of concave portions of the contour)
- i) symmetry
- j) fractal dimension ("coastline approximation" - 1)

The mean, standard error and "worst" or largest (mean of the three largest values) of these features were computed for each image, resulting in 30 features. For instance, field 3 is Mean Radius, field 13 is Radius SE, field 23 is Worst Radius.

All feature values are recoded with four significant digits.

Missing attribute values: none

Class distribution: 357 benign, 212 malignant

Correction Examens Finel

(A)

27 octobre 2023

HAX 912X

Exercise 1

4 points seront tirés au hasard

1) Importation des données
avec un filtre (après suppression
de certains individus)

(APT)

119 individus

13 variables contiennent une
variable à valeur négative
Toutes les variables

des 12 variables continues
peut contenir valeurs

2] Théorie de la régression logistique ②
NEC permet de faire un test de
MIS en oeuvre simple et ① pt
contient la vérification C = 1
Tous les coefficients estimés sont
non nuls -

3] Estimation du taux de
bonne classification par
validation croisée à 5 ① pt
ensembles répartis 100 fois
résultat # 0,98 % (sur 500
des 500 estimations)

4] Choisir le meilleur modèle de
régression entre tous (l1)
et nivage (l2) pour 50 valeurs
de C entre e^{-4} et e^4 uniformément
réparties sur l'échelle log -

Reilles marlak

③

Moyenne $\text{NFC} < +0,03$

Tous les bacs classés moins 0,99%
Nombre d'échantillon fausse /
nombre d'analyse utilisée) pour donner
- le tout correctement -

Exercice 2

3 points doivent être mentionnés

- 1) Les données - les cellules d'une image d'un trame sont dans une grille 212x212 et il y a 569 cellules, il y a 569 cellules, il y a 569 cellules.
- 2) Les séances - Objectif : construire une fonction graphique prenant en entrée une cellule de la grille et renvoyant un résultat entre 0 et 100.
- 3) Les résultats - 30 méthodes pour faire

les 30 derniers mois

4

sont centrés sur les

hommes non séparés en deux
échantillons : un échantillon

2 pts

d'apprentissage contenant 455

individus et un échantillon d'essai
contenant 116 individus.

2) On classifie par fonction

et intérêt en 1 échantillon

d'apprentissage. Il contient

500 verbs et il devra servir

la meilleure manière de choisir

versi. un nouveau tiré

aléatoirement de 30 mots -

Tous les hommes classifiés sur l'échantillon

ont fait 96% -

3) Avec les probabilités électriques,
 on peut déterminer une estimation
 statistique en utilisant une loi de
 probabilité d'un classifieur. Pour un
 film, il suffit d'utiliser
 l'estimation out-of-bag sur
 l'échantillon d'apprentissage -
 Dans ce cas, si tous les films
 classifiés sont estimés à
 96% - alors c'est l'estimation
 obtenue sur l'échantillon utilisé.

Evidence 3

1) $y \sim \mathcal{B}(p(1))$

La loi du Bernoulli effectue une
 loi exponentielle simple.

$$f(y|p) = p^y (1-p)^{1-y} \sum_{y \in \{0,1\}} S(dy) \quad (6)$$

$$f(y|p) = \exp \left(y \log \left(\frac{p}{1-p} \right) + \log(1-p) \right) S(dy)$$

$$\theta = \log \left(\frac{p}{1-p} \right) \Rightarrow p = \frac{e^\theta}{1+e^\theta}$$

$$f(\theta) = -\log \left(\frac{1}{1+e^\theta} \right) = \log(1+e^\theta)$$

Na willens, $p(x) = P(y=1)$

$$p(x) = P(y^* \leq 0)$$

$$p(x) = \int_0^\infty \frac{\exp(-|y^* - \beta_0 - \beta_1 x|)}{(1 + \exp(-|y^* - \beta_0 - \beta_1 x|))^2} dy^*$$

$$\Leftrightarrow p(x) = \left[\frac{1}{1 + \exp(-|y^* - \beta_0 - \beta_1 x|)} \right]_0^\infty$$

$$\Leftrightarrow p(x) = \frac{1}{1 + e^{\beta_0 + \beta_1 x}} = \left[\frac{e^{-\beta_0 - \beta_1 x}}{1 + e^{-\beta_0 - \beta_1 x}} \right]$$

(7)

$$E(y) = f'(0) = \frac{e^0}{1+e^0} = p$$

$$E(y) = \frac{1}{1+e^{B_0+B_1K}} = f^{-1}(B_0+B_1K)$$

Il s'agit bien d'un modélis linéaire
générique -

2) New neurons

$$\theta = \log \left(\frac{1}{1+e^{B_0+B_1K}} \times \frac{1+e^{B_0+B_1K}}{e^{B_0+B_1K}} \right)$$

$$\theta = -B_0 - B_1 K + B_0 + B_1 K$$

équivalent si on considère les
paramètres $\tilde{B}_0 = -B_0$ et $\tilde{B}_1 = -B_1$

$$\text{alors } \theta = \tilde{B}_0 + \tilde{B}_1 K$$

Dans formellement - on est pas à la
même que pour la théorie mais
peut $(\tilde{B}_0, \tilde{B}_1)$ soit être le -

Entwurf mathematisches C'at min. ⑧

Entwurf der Wahrscheinlichkeit, Cat mini.

$$3) \angle(\beta_0, \beta_1) = - \sum_{i=1}^{100} y_i \log(1 + e^{\beta_0 + \beta_1 x_i})$$

$$+ \sum_{i=1}^{100} (1-y_i) \log\left(\frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}}\right)$$

$$\Leftrightarrow \angle(\beta_0, \beta_1) = - \sum_{i=1}^{100} y_i (\beta_0 + \beta_1 x_i)$$

$$+ \sum_{i=1}^{100} \log\left(\frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}}\right)$$

$$\Leftrightarrow \angle(\beta_0, \beta_1) = - 50\beta_0 - 30\beta_1$$

$$+ 100\beta_0 + 50\beta_1 - \sum_{i=1}^{100} \log(1 + e^{\beta_0 + \beta_1 x_i})$$

$$\Leftrightarrow \angle(\beta_0, \beta_1) = 50\beta_0 + 20\beta_1 - 50 \log(1 + e^{\beta_0})$$
$$- 50 \log(1 + e^{\beta_0 + \beta_1})$$

Eriogonum

(9)

Can contain up to 100 species
now determined in the area
lacks a median angle
(Wright, 1911) indicating the number
of members - which is often
at most 2 or 3 members
which allows for few angles
about it feasible to group
as parts of one -