

iris-classification

August 19, 2020

1 0. Introdução

Trabalho Clustering:

Aluno: Gabriel Luiz

Disciplina: Tópico em Aprendizado de Máquina

Objetivos :

- Escolha dois datasets rotulados.
- Realize a análise estatística, visualização e pré-processamento dos dados.
- Realize os experimentos criando duas bases de teste distintas:
 - – considerando todos os atributos do dataset ;
 - – selecionando alguns atributos e descartando outros;
- Aplique três métodos de clustering distintos nas duas bases acima.
- Para cada dataset , em cada uma das bases, analise os resultados segundo medidas de qualidade de clustering , usando índices de validação interna (SSW, SSB, silhueta, Calinski-Harabasz, Dunn e Davis-Bouldin) e externa (pureza, entropia, acurácia, F-measure , ARI, NMI).
- Proponha uma maneira adicional de comparar os resultados obtidos além das medidas acima.
- Compare e interprete os resultados dos dois experimentos em cada dataset

1.1 0.1 Dependências

Para realização da tarefa foram utilizados as seguintes bibliotecas:

```
[96]: from datetime import datetime
import numpy as np
import pandas as pd
from sklearn import preprocessing
from sklearn.cluster import *
import seaborn as sns
from sklearn.preprocessing import StandardScaler
import matplotlib.pyplot as plt
```

```

from sklearn.feature_selection import SelectKBest
from sklearn.feature_selection import chi2
from sklearn.model_selection import train_test_split

# KFold
from sklearn.model_selection import KFold
import random

# Classificadores
from sklearn.neighbors import KNeighborsClassifier
from sklearn.naive_bayes import GaussianNB
from sklearn.tree import DecisionTreeClassifier
from sklearn.svm import SVC

#Metricas
from sklearn.metrics import accuracy_score
from sklearn.metrics import recall_score
from sklearn.metrics import precision_score
from sklearn.metrics import f1_score
from sklearn.metrics import cohen_kappa_score
from sklearn.metrics import roc_auc_score
from sklearn.metrics import balanced_accuracy_score

from sklearn.datasets import load_iris

```

2 1. Dados

Para realização das tarefas envolvidas neste relatório utilizou-se o arquivo **dim128.csv** que contém dados não descritos, onde foram feitos para a realização de clustering que se encontram no site: <http://cs.uef.fi/sipu/datasets/>

2.1 1.1 Carregamento do arquivo

```
[97]: (data, label) = load_iris(return_X_y=True, as_frame = True)
```

```
[98]: scaler = preprocessing.StandardScaler()
data_scaler = scaler.fit_transform(X = data)

data_results = np.array(label)

## 3.2 Selecionando atributos do dataset
```

```
[99]: data_reduzida = pd.DataFrame(SelectKBest(chi2, k=2).fit_transform(data, label))
data_reduzida.shape
```

```
data_reduzida = data_reduzida.to_numpy()

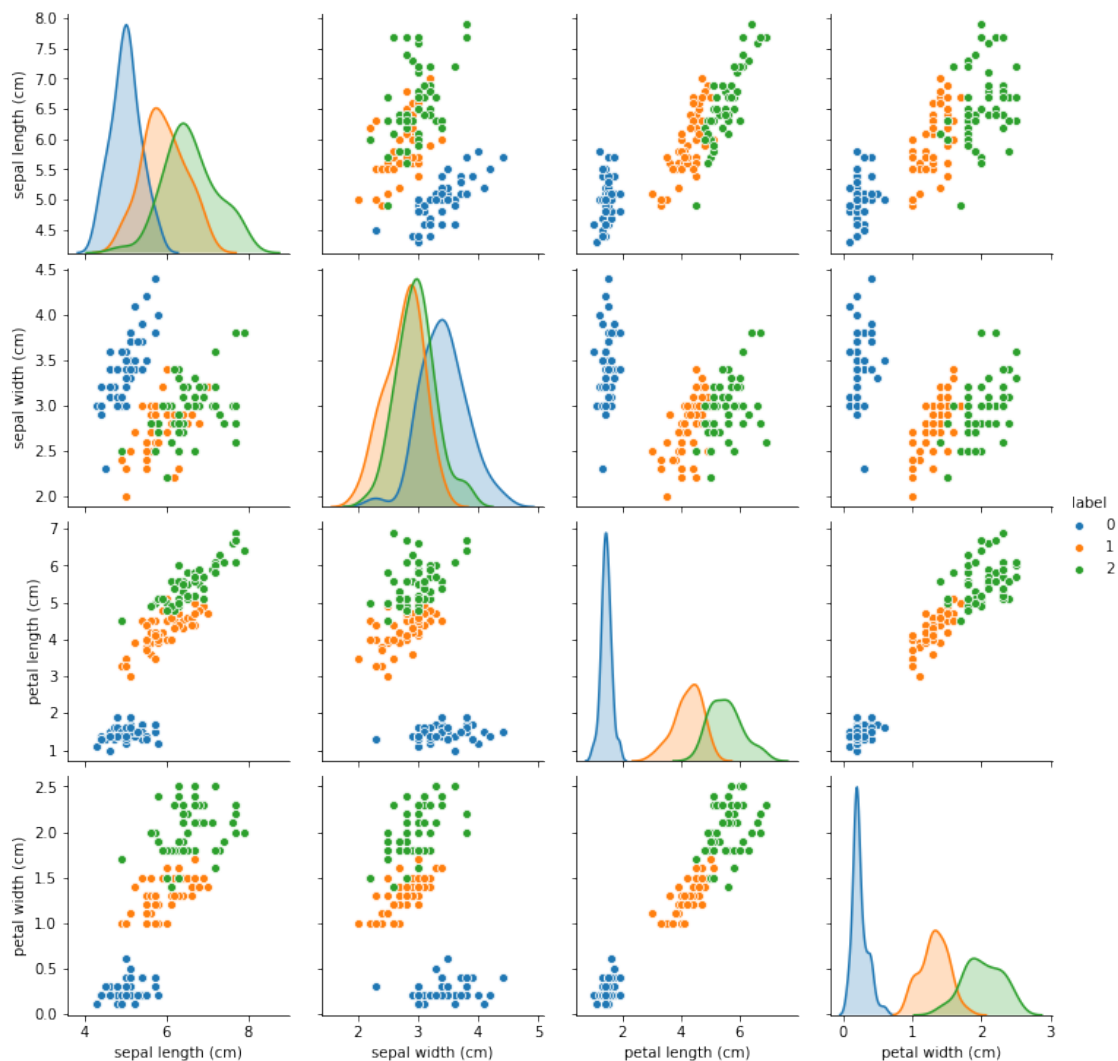
data_scaler2 = scaler.fit_transform(X = data_reduzida)
```

```
[100]: # data_scaled2 = pd.DataFrame(data_scaler2)
# data_scaled2.head()
```

```
[126]: df = data
df = df.assign(label = label)
```

```
[127]: sns.pairplot(df, diag_kind="kde",hue='label')
```

```
[127]: <seaborn.axisgrid.PairGrid at 0x7ff65902ae80>
```



2.2 Classificando

2.3 Funções necessárias

```
[101]: def calcula_metricas(metricas, y_test, y_predict):  
    metricas['acc'] += (accuracy_score(y_test, y_predict))  
    metricas['recall'] += (recall_score(y_test, y_predict, average='micro'))  
    metricas['precision'] += (precision_score(y_test, y_predict,   
→average='macro'))  
    metricas['f1'] += f1_score(y_test, y_predict, average='weighted')  
    # metricas['roc'] += roc_auc_score(y_test, y_predict)  
    metricas['kappa'] += cohen_kappa_score(y_test, y_predict)  
    metricas['balanced_acc'] += balanced_accuracy_score(y_test, y_predict)
```

```
[102]: def save_metricas(name, metricas):  
    f = open(name, 'w')  
    f.write('Acuária: ' + str(metricas['acc']) + '\n')  
    f.write('Recall: ' + str(metricas['recall']) + '\n')  
    f.write('Precisão: ' + str(metricas['precision']) + '\n')  
    f.write('F-Measure: ' + str(metricas['f1']) + '\n')  
    # f.write('Curva Roc: ' + str(metricas['roc']) + '\n')  
    f.write('Indice Kappa: ' + str(metricas['kappa']) + '\n')  
    f.write('Acuária Balanceada: ' + str(metricas['balanced_acc']) + '\n')  
    f.close()
```

```
[103]: def show_metricas(metricas):  
    print('Acuária:', metricas['acc'])  
    print('Recall:', metricas['recall'])  
    print('Precisão:', metricas['precision'])  
    print('F-Measure:', metricas['f1'])  
    # print('Curva Roc:', metricas['roc'])  
    print('Indice Kappa:', metricas['kappa'])  
    print('Acuária Balanceada:', metricas['balanced_acc'])
```

```
[104]: def write_metricas(name_file, metricas, metodo):  
    f = open(name_file, "a")  
    f.write(metodo + ',')  
    f.write(str(metricas['acc']) + ',')  
    f.write(str(metricas['recall']) + ',')  
    f.write(str(metricas['precision']) + ',')  
    f.write(str(metricas['f1']) + ',')  
    # f.write(str(round(metricas['roc'],4)) + ';')  
    f.write(str(metricas['kappa']) + ',')  
    f.write(str(metricas['balanced_acc']) + '\n')  
    f.close()
```

2.4 Aplicando KNN com K-fold

2.5 DataFrame Cru

```
[105]: formato = 'Cru'

       folds_value = 16
```

```
[106]: # TODO change split function
kf = KFold(n_splits=2, shuffle=True, random_state=random.randint(0, 10))
data_kfold = kf.split(data)

train = []
test = []

for train_index, test_index in data_kfold:
    train.append(train_index)
    test.append(test_index)
```

```
[107]: name_file = 'metricas.csv'

# Roc;
f = open(name_file, "w")
f.write('Acurácia,Recall,Precisão,F1,Kappa,Acurácia Balanceada\n')
f.close()
```

2.6 Aplicando KNN com K-fold

```
[108]: # 'roc': 0,
metodo = 'KNN'
metricas = {'acc': 0, 'recall': 0, 'precision': 0, 'f1': 0, 'kappa': 0,
            ↪ 'balanced_acc': 0}

x_train, x_test, y_train, y_test = train_test_split(data_scaler, data_results)

neigh = KNeighborsClassifier(n_neighbors=10)
neigh.fit(x_train, y_train)

y_predict = neigh.predict(x_test)

calcula_metricas(metricas, y_test, y_predict)

show_metricas(metricas)
write_metricas(name_file, metricas, metodo)
```

Acuária: 1.0
Recall: 1.0
Precisão: 1.0
F-Measure: 1.0
Indice Kappa: 1.0
Acuária Balanceada: 1.0

```
[109]: # # 'roc': 0,
# metodo = 'KNN'
# metricas = {'acc': 0, 'recall': 0, 'precision': 0, 'f1': 0, 'kappa': 0,
#           ↪ 'balanced_acc': 0}
#
# for train_index, test_index in zip(train, test):
#     x_train, x_test = data_scaler[train_index], data_scaler[test_index]
#     y_train, y_test = data_results[train_index], data_results[test_index]
#
#     neigh = KNeighborsClassifier(n_neighbors=100)
#     neigh.fit(x_train, y_train)
#
#     y_predict = neigh.predict(x_test)
#
#     calcula_metricas(metricas, y_test, y_predict)
#
# for metrica, value in metricas.items():
#     metricas[metrica] = value/10
#
# show_metricas(metricas)
# write_metricas(name_file, metricas, metodo)
```

2.7 Aplicando GaussianNB com K-fold

```
[110]: metodo = 'Gauss'
metricas = {'acc': 0, 'recall': 0, 'precision': 0, 'f1': 0, 'roc': 0, 'kappa': 0,
           ↪ 'balanced_acc': 0}

x_train, x_test, y_train, y_test = train_test_split(data_scaler, data_results)

gauss = GaussianNB()
gauss.fit(x_train, y_train)

y_predict = gauss.predict(x_test)

calcula_metricas(metricas, y_test, y_predict)

show_metricas(metricas)
```

```
write_metricas(name_file, metricas, metodo)
```

Acuária: 0.9473684210526315
Recall: 0.9473684210526315
Precisão: 0.9458874458874459
F-Measure: 0.9473684210526315
Indice Kappa: 0.9206680584551148
Acuária Balanceada: 0.9458874458874459

```
[111]: # metodo = 'Gauss'
# metricas = {'acc': 0, 'recall': 0, 'precision': 0, 'f1': 0, 'roc': 0, 'kappa':
→ 0, 'balanced_acc': 0}
#
# for train_index, test_index in zip(train, test):
#     x_train, x_test = data_scaler[train_index], data_scaler[test_index]
#     y_train, y_test = data_results[train_index], data_results[test_index]
#
#     gauss = GaussianNB()
#     gauss.fit(x_train, y_train)
#
#     y_predict = gauss.predict(x_test)
#
#     calcula_metricas(metricas, y_test, y_predict)
#
# for metrica, value in metricas.items():
#     metricas[metrica] = value/10
#
# show_metricas(metricas)
# write_metricas(name_file, metricas, metodo)
```

2.8 Aplicando DecisionTreeClassifier com K-fold

```
[112]: metodo = 'Tree'
metricas = {'acc': 0, 'recall': 0, 'precision': 0, 'f1': 0, 'roc': 0, 'kappa': 0,
→ 'balanced_acc': 0}

x_train, x_test, y_train, y_test = train_test_split(data_scaler, data_results)

tree = DecisionTreeClassifier()
tree.fit(x_train, y_train)

y_predict = tree.predict(x_test)

calcula_metricas(metricas, y_test, y_predict)

show_metricas(metricas)
```

```
write_metricas(name_file, metricas, metodo)
```

Acuária: 1.0
Recall: 1.0
Precisão: 1.0
F-Measure: 1.0
Índice Kappa: 1.0
Acuária Balanceada: 1.0

```
[113]: # metodo = 'Tree'
# metricas = {'acc': 0, 'recall': 0, 'precision': 0, 'f1': 0, 'roc': 0, 'kappa':
→ 0, 'balanced_acc': 0}
#
# for train_index, test_index in zip(train, test):
#     x_train, x_test = data_scaler[train_index], data_scaler[test_index]
#     y_train, y_test = data_results[train_index], data_results[test_index]
#
#     tree = DecisionTreeClassifier()
#     tree.fit(x_train, y_train)
#
#     y_predict = tree.predict(x_test)
#
#     calcula_metricas(metricas, y_test, y_predict)
#
# for metrica, value in metricas.items():
#     metricas[metrica] = value/10
#
# show_metricas(metricas)
# write_metricas(name_file, metricas, metodo)
```

2.9 Aplicando SVM com K-fold

```
[114]: metodo = 'SVM'
metricas = {'acc': 0, 'recall': 0, 'precision': 0, 'f1': 0, 'roc': 0, 'kappa': 0,
→ 'balanced_acc': 0}

x_train, x_test, y_train, y_test = train_test_split(data_scaler, data_results)

svm = SVC()
svm.fit(x_train, y_train)

y_predict = svm.predict(x_test)

calcula_metricas(metricas, y_test, y_predict)

show_metricas(metricas)
```



```
write_metricas(name_file, metricas, metodo)
```

Acuária: 0.9473684210526315
Recall: 0.9473684210526315
Precisão: 0.9555555555555556
F-Measure: 0.9468557758031441
Índice Kappa: 0.9208333333333334
Acuária Balanceada: 0.9444444444444445

```
[115]: # metodo = 'SVM'
# metricas = {'acc': 0, 'recall': 0, 'precision': 0, 'f1': 0, 'roc': 0, 'kappa':
→ 0, 'balanced_acc': 0}
#
# for train_index, test_index in zip(train, test):
#     x_train, x_test = data_scaler[train_index], data_scaler[test_index]
#     y_train, y_test = data_results[train_index], data_results[test_index]
#
#     svm = SVC()
#     svm.fit(x_train, y_train)
#
#     y_predict = svm.predict(x_test)
#
#     calcula_metricas(metricas, y_test, y_predict)
#
# for metrica, value in metricas.items():
#     metricas[metrica] = value/10
#
# show_metricas(metricas)
# write_metricas(name_file, metricas, metodo)
```

2.10 DataFrame Selecionado

2.11 Aplicando

```
[116]: kf = KFold(n_splits=2, shuffle=True, random_state=random.randint(0, 10))
data_kfold = kf.split(data_scaler2)

train = []
test = []

for train_index, test_index in data_kfold:
    train.append(train_index)
    test.append(test_index)
```

2.12 Aplicando KNN com K-fold

```
[117]: metodo = 'KNNSELECT'
metricas = {'acc': 0, 'recall': 0, 'precision': 0, 'f1': 0, 'roc': 0, 'kappa': 0, 'balanced_acc': 0}

x_train, x_test, y_train, y_test = train_test_split(data_scaler, data_results)

neigh = KNeighborsClassifier(n_neighbors=10)
neigh.fit(x_train, y_train)

y_predict = neigh.predict(x_test)

calcula_metricas(metricas, y_test, y_predict)

show_metricas(metricas)
write_metricas(name_file, metricas, metodo)
```

Acuária: 1.0
Recall: 1.0
Precisão: 1.0
F-Measure: 1.0
Indice Kappa: 1.0
Acuária Balanceada: 1.0

```
[118]: # metodo = 'KNNSELECT'
# metricas = {'acc': 0, 'recall': 0, 'precision': 0, 'f1': 0, 'roc': 0, 'kappa': 0, 'balanced_acc': 0}
#
# for train_index, test_index in zip(train, test):
#     x_train, x_test = data_scaler2[train_index], data_scaler2[test_index]
#     y_train, y_test = data_results[train_index], data_results[test_index]
#
#     neigh = KNeighborsClassifier(n_neighbors=20)
#     neigh.fit(x_train, y_train)
#
#     y_predict = neigh.predict(x_test)
#
#     calcula_metricas(metricas, y_test, y_predict)
#
# for metrica, value in metricas.items():
#     metricas[metrica] = value/10
#
# show_metricas(metricas)
# write_metricas(name_file, metricas, metodo)
```

2.13 Aplicando GaussianNB com K-fold

```
[119]: metodo = 'GaussSELECT'
metricas = {'acc': 0, 'recall': 0, 'precision': 0, 'f1': 0, 'roc': 0, 'kappa': 0, 'balanced_acc': 0}

x_train, x_test, y_train, y_test = train_test_split(data_scaler, data_results)

gauss = GaussianNB()
gauss.fit(x_train, y_train)

y_predict = gauss.predict(x_test)

calcula_metricas(metricas, y_test, y_predict)

show_metricas(metricas)
write_metricas(name_file, metricas, metodo)
```

Acuária: 0.9473684210526315
Recall: 0.9473684210526315
Precisão: 0.9440559440559442
F-Measure: 0.9473684210526315
Indice Kappa: 0.9206680584551148
Acuária Balanceada: 0.9440559440559442

```
[120]: # metodo = 'GaussSELECT'
# metricas = {'acc': 0, 'recall': 0, 'precision': 0, 'f1': 0, 'roc': 0, 'kappa': 0, 'balanced_acc': 0}
#
# for train_index, test_index in zip(train, test):
#     x_train, x_test = data_scaler2[train_index], data_scaler2[test_index]
#     y_train, y_test = data_results[train_index], data_results[test_index]
#
#     gauss = GaussianNB()
#     gauss.fit(x_train, y_train)
#
#     y_predict = gauss.predict(x_test)
#
#     calcula_metricas(metricas, y_test, y_predict)
#
# for metrica, value in metricas.items():
#     metricas[metrica] = value/10
#
# show_metricas(metricas)
# write_metricas(name_file, metricas, metodo)
```

2.14 Aplicando DecisionTreeClassifier com K-fold

```
[121]: metodo = 'TreeSELECT'
metricas = {'acc': 0, 'recall': 0, 'precision': 0, 'f1': 0, 'roc': 0, 'kappa': 0, 'balanced_acc': 0}

x_train, x_test, y_train, y_test = train_test_split(data_scaler, data_results)

tree = DecisionTreeClassifier()
tree.fit(x_train, y_train)

y_predict = tree.predict(x_test)

calcula_metricas(metricas, y_test, y_predict)

show_metricas(metricas)
write_metricas(name_file, metricas, metodo)
```

Acúria: 0.9473684210526315
Recall: 0.9473684210526315
Precisão: 0.9500891265597149
F-Measure: 0.9473684210526315
Índice Kappa: 0.9186295503211992
Acúria Balanceada: 0.9500891265597149

```
[122]: # metodo = 'TreeSELECT'
# metricas = {'acc': 0, 'recall': 0, 'precision': 0, 'f1': 0, 'roc': 0, 'kappa': 0, 'balanced_acc': 0}
#
# for train_index, test_index in zip(train, test):
#     x_train, x_test = data_scaler2[train_index], data_scaler2[test_index]
#     y_train, y_test = data_results[train_index], data_results[test_index]
#
#     tree = DecisionTreeClassifier()
#     tree.fit(x_train, y_train)
#
#     y_predict = tree.predict(x_test)
#
#     calcula_metricas(metricas, y_test, y_predict)
#
# for metrica, value in metricas.items():
#     metricas[metrica] = value/10
#
# show_metricas(metricas)
# write_metricas(name_file, metricas, metodo)
```

2.15 Aplicando SVM com K-fold

```
[123]: metodo = 'SVMSELECT'
metricas = {'acc': 0, 'recall': 0, 'precision': 0, 'f1': 0, 'roc': 0, 'kappa': 0, 'balanced_acc': 0}

x_train, x_test, y_train, y_test = train_test_split(data_scaler, data_results)

svm = SVC()
svm.fit(x_train, y_train)

y_predict = svm.predict(x_test)

calcula_metricas(metricas, y_test, y_predict)

show_metricas(metricas)
write_metricas(name_file, metricas, metodo)
```

Acuária: 0.9473684210526315
Recall: 0.9473684210526315
Precisão: 0.9555555555555556
F-Measure: 0.947223828802776
Indice Kappa: 0.9208333333333334
Acuária Balanceada: 0.9523809523809524

```
[124]: # metodo = 'SVMSELECT'
# metricas = {'acc': 0, 'recall': 0, 'precision': 0, 'f1': 0, 'roc': 0, 'kappa': 0, 'balanced_acc': 0}
#
# for train_index, test_index in zip(train, test):
#     x_train, x_test = data_scaler2[train_index], data_scaler2[test_index]
#     y_train, y_test = data_results[train_index], data_results[test_index]
#
#     svm = SVC()
#     svm.fit(x_train, y_train)
#
#     y_predict = svm.predict(x_test)
#
#     calcula_metricas(metricas, y_test, y_predict)
#
# for metrica, value in metricas.items():
#     metricas[metrica] = value/10
#
# show_metricas(metricas)
# write_metricas(name_file, metricas, metodo)
```

```
[125]: analyse = './metricas.csv'
metricas = pd.read_csv(
    analyse,
)
metricas
```

```
[125]:
```

	Acurácia	Recall	Precisão	F1	Kappa	\
KNN	1.000000	1.000000	1.000000	1.000000	1.000000	
Gauss	0.947368	0.947368	0.945887	0.947368	0.920668	
Tree	1.000000	1.000000	1.000000	1.000000	1.000000	
SVM	0.947368	0.947368	0.955556	0.946856	0.920833	
KNNSELECT	1.000000	1.000000	1.000000	1.000000	1.000000	
GaussSELECT	0.947368	0.947368	0.944056	0.947368	0.920668	
TreeSELECT	0.947368	0.947368	0.950089	0.947368	0.918630	
SVMSELECT	0.947368	0.947368	0.955556	0.947224	0.920833	

	Acurácia Balanceada
KNN	1.000000
Gauss	0.945887
Tree	1.000000
SVM	0.944444
KNNSELECT	1.000000
GaussSELECT	0.944056
TreeSELECT	0.950089
SVMSELECT	0.952381