

ĐẠI HỌC QUỐC GIA THÀNH PHỐ HỒ CHÍ MINH
TRƯỜNG ĐẠI HỌC KHOA HỌC TỰ NHIÊN
KHOA TOÁN - TIN HỌC



KTDL | BÀI 9: SỰ PHÂN LỚP DỮ LIỆU

Sinh viên thực hiện: Nguyễn Công Hoài Nam
Mã số sinh viên: 21280099

Ngày 9 tháng 6 năm 2024

MỤC LỤC

1	Thông tin về dữ liệu	1
2	Kiểm tra tính tuyến tính của dữ liệu	1
3	Xây dựng mô hình	3
3.1	Import các thư viện cần thiết	3
3.2	Chia dữ liệu thành tập train và test	3
3.3	Chuẩn hoá dữ liệu	4
3.4	Hyperparameter tuning	4
3.5	Evaluate	4
4	Trực quan hoá (vẽ hình)	5
5	Kết luận	7

1. Thông tin về dữ liệu

Tập dữ liệu **Banknote Authentication** gồm các đặc trưng được trích xuất từ các hình ảnh được chụp xác thực các tờ tiền giấy (hợp lệ hay giả mạo)

Chúng bao gồm các đặc trưng:

- **Variance**: độ biến thiên
- **Skewness**: độ lệch
- **Curtosis**: độ nhọn
- **Entropy**: mức độ hỗn loạn
- **Class**: nhãn đầu ra (1: hợp lệ, 0: giả mạo)

Ta sẽ import data từ **UCI Machine Learning Repository** bằng thư viện **ucimlrepo** của họ

```
1 # import libraries
2 import numpy as np
3 import pandas as pd
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 sns.set_style('whitegrid')
7
8 from ucimlrepo import fetch_ucirepo
9
10 # fetch dataset
11 banknote_authentication = fetch_ucirepo(id=267)
12
13 # data (as pandas dataframes)
14 X = banknote_authentication.data.features
15 y = banknote_authentication.data.targets
16 df = pd.concat([X, y], axis=1)
17 df.head()
```

	variance	skewness	curtosis	entropy	class
0	3.62160	8.6661	-2.8073	-0.44699	0
1	4.54590	8.1674	-2.4586	-1.46210	0
2	3.86600	-2.6383	1.9242	0.10645	0
3	3.45660	9.5228	-4.0112	-3.59440	0
4	0.32924	-4.4552	4.5718	-0.98880	0

2. Kiểm tra tính tuyến tính của dữ liệu

Để kiểm tra dữ liệu có tuyến tính hay không ta sẽ vẽ ma trận tương quan để kiểm tra sự tương của các biến

```
1 plt.figure(figsize=(6,6))
```

```

2 sns.heatmap(df.corr(), cmap="Blues_r", annot=True, fmt='.2',linewidths=.5, cbar=False)
3 # Show the plot
4 plt.show()

```

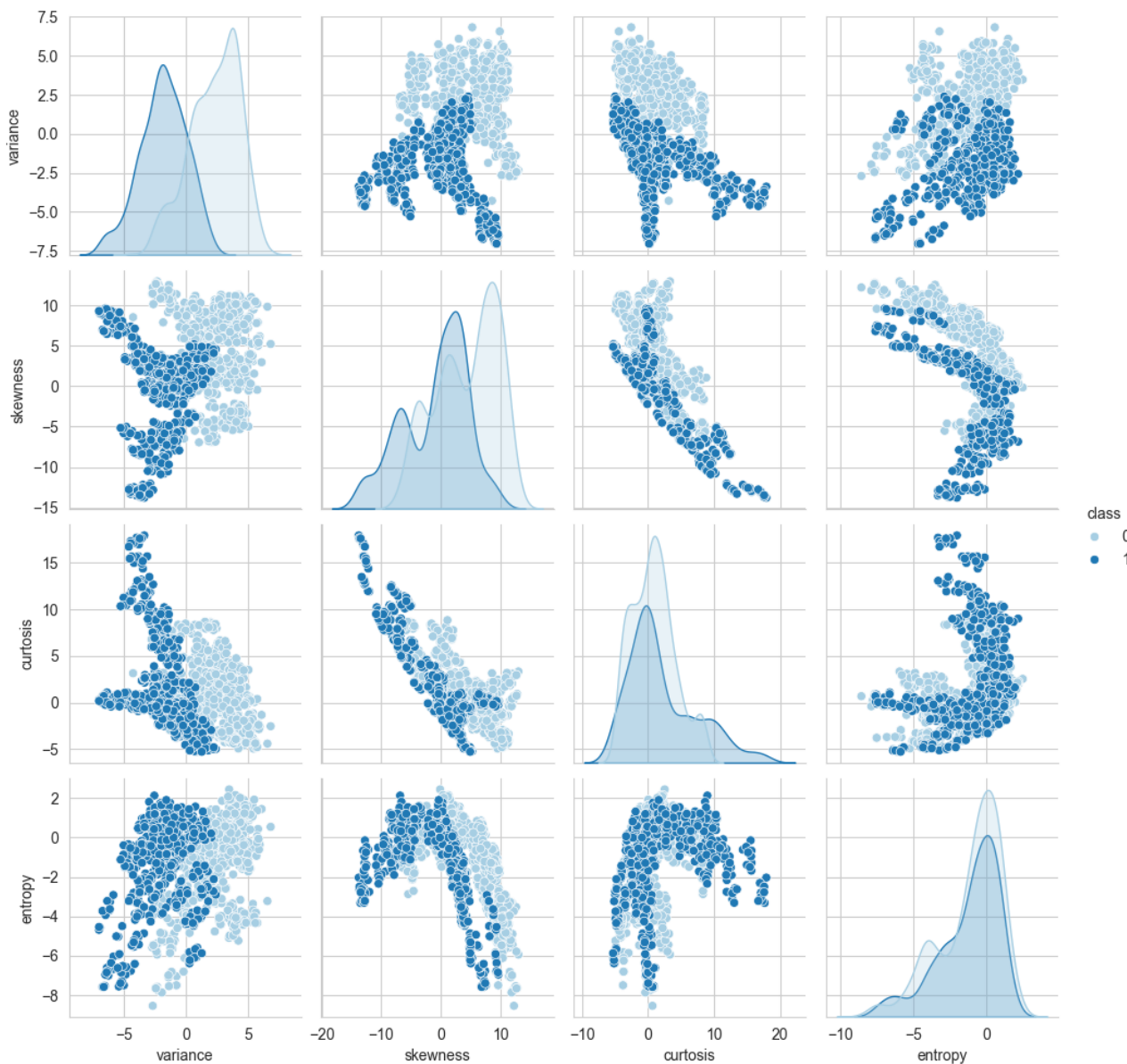


- Những cặp feature có tương quan gần cao (gần với -1 hoặc 1) thì chúng tuyến tính, ngược lại chúng phi tuyến.
 - Ta có thể thấy một số cặp feature có sự tuyến tính như (**skewness** với **kurtosis**) nhưng cũng có nhiều cặp features phi tuyến như (**kurtosis** và **class**).
 - Ta kết luận bộ dữ liệu không hoàn toàn tuyến tính
- Để kiểm chứng lại nhận định của mình ta sẽ vẽ scatter plot của từng cặp biến trong tập dữ liệu

```

1 sns.pairplot(df, hue = 'class', palette="Paired")

```



Dễ dàng thấy ở nhiều cặp biến các điểm dữ liệu không phân tán theo một đường thẳng nào => dữ liệu không tuyến tính

3. Xây dựng mô hình

3.1. Import các thư viện cần thiết

```

1 from sklearn.model_selection import train_test_split
2 from sklearn.preprocessing import StandardScaler
3 from sklearn.svm import SVC
4 from sklearn.model_selection import GridSearchCV
5 from sklearn.decomposition import PCA
6 from sklearn.metrics import classification_report, confusion_matrix
7 from sklearn.pipeline import make_pipeline

```

3.2. Chia dữ liệu thành tập train và test

```

1 # Split dataset
2 X = X.values
3 y = y.values.ravel()

```

```

4 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

```

3.3. Chuẩn hoá dữ liệu

```

1 # Standardization
2 scaler = StandardScaler()
3 X_train = scaler.fit_transform(X_train)
4 X_test = scaler.transform(X_test)

```

3.4. Hyperparameter tuning

Bởi vì bộ dữ liệu không tuyến tính hoàn toàn nên ta sẽ sử dụng kernel **rbf** dành cho data phi tuyến. Để tìm kiếm bộ siêu tham số tốt nhất cho model ta sẽ sử dụng thuật toán **GridSearchCV** để tìm ra tham số tốt nhất cho bài toán này.

```

1 param_grid = {'C': [0.1, 1, 10, 100, 1000],
2               'gamma': [1, 0.1, 0.01, 0.001, 0.0001],
3               'kernel': ['rbf']}
4
5 # Sử dụng GridSearchCV để tìm best hyperparameters cho model với kernel = 'rbf'
6 svc = SVC()
7 best_svc = GridSearchCV(svc, param_grid)
8 best_svc.fit(X_train, y_train)
9 print(f"Best param: {best_svc.best_params_}")

```

Kết quả bộ tham số tốt nhất là:

Best param: 'C': 10, 'gamma': 0.1, 'kernel': 'rbf'

3.5. Evaluate

Để đánh giá mô hình tốt hay tệ, ta sẽ đánh giá nó qua các metric, cụ thể ở đây là **confusion matrix** và **classification report** thể hiện thông số ở các class khác nhau.

```

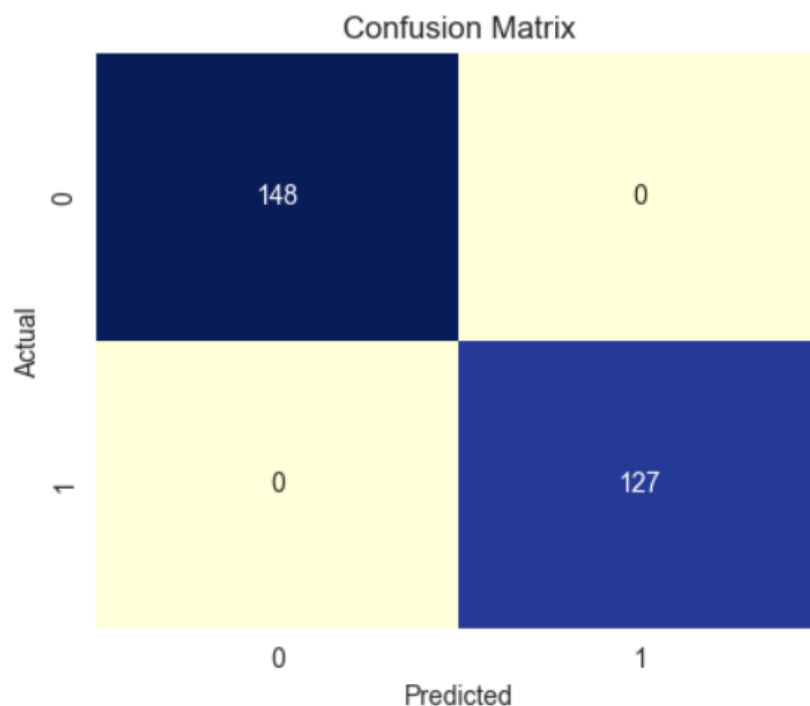
1 def evaluate(y_true, y_pred):
2
3     confusion = confusion_matrix(y_true, y_pred)
4     report = classification_report(y_true, y_pred)
5
6     print("Classification report:\n", report)
7     print("Confusion Matrix")
8
9     plt.figure(figsize=(5, 4))
10    sns.heatmap(confusion, annot=True, fmt='d', cmap='YlGnBu', cbar=False)
11    plt.xlabel('Predicted')
12    plt.ylabel('Actual')
13    plt.title('Confusion Matrix')
14    plt.show()
15
16    y_pred = best_svc.predict(X_test)
17    evaluate(y_test, y_pred)

```

Classification report:

	precision	recall	f1-score	support
0	1.00	1.00	1.00	148
1	1.00	1.00	1.00	127
accuracy			1.00	275
macro avg	1.00	1.00	1.00	275
weighted avg	1.00	1.00	1.00	275

Confusion Matrix



Kết quả cho thấy model mà ta xây dựng có điểm số tuyệt đối.

4. Trực quan hoá (vẽ hình)

Bởi vì bộ data này là dữ liệu nhiều chiều nên để có thể vẽ hình, thì ta chỉ có thể chọn 2 biến trong dataset để thể hiện.

Ở đây em chọn giảm chiều dữ liệu xuống còn 2 chiều để có cái nhìn tổng quát về dữ liệu

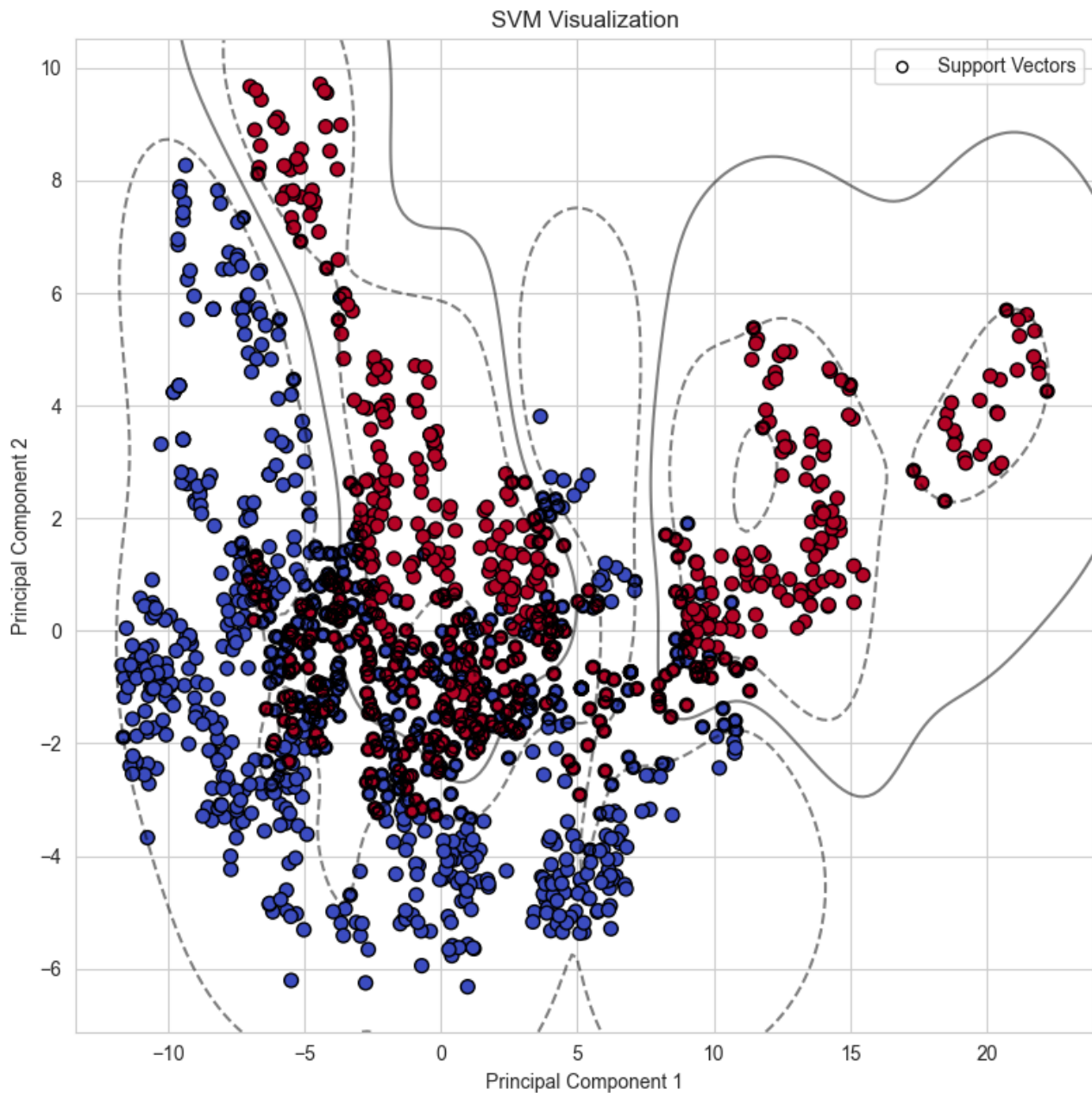
```

1 # Hàm tạo mô hình SVC với PCA
2 def make_2D_SVC(X, y):
3     pipeline = make_pipeline(
4         PCA(n_components=2), # PCA data xuống còn 2 chiều
5         SVC(C=10, gamma=0.1, kernel='rbf') # SVC với best hyperparameter tuning
6     )
7     pipeline.fit(X, y)
8     X_2D = pipeline.named_steps['pca'].transform(X)
9     model = pipeline.named_steps['svc']
10    return X_2D, model

```

Hàm vẽ hình gồm scatterplot, các đường biên,

```
1  def visualize_2d_svm(X, y):
2      X_2D, model = make_2D_SVC(X,y)
3      plt.figure(figsize=(8, 8))
4      plt.scatter(X_2D[:, 0], X_2D[:, 1], c=y, s=50, cmap='coolwarm', edgecolor='k')
5
6      # Lấy giới hạn của biểu đồ
7      ax = plt.gca()
8      xlim = ax.get_xlim()
9      ylim = ax.get_ylim()
10
11     # Tạo lưới để đánh giá mô hình
12     x = np.linspace(xlim[0], xlim[1], 200)
13     y = np.linspace(ylim[0], ylim[1], 200)
14     Y, X = np.meshgrid(y, x)
15     xy = np.vstack([X.ravel(), Y.ravel()]).T
16     P = model.decision_function(xy).reshape(X.shape)
17
18     # Vẽ đường biên quyết định và lề
19     ax.contour(X, Y, P, colors='k', levels=[-1, 0, 1], alpha=0.5, linestyles=['--', '-', '-.-',
20         ])
21
22     # Vẽ các support vectors
23     ax.scatter(model.support_vectors_[:, 0], model.support_vectors_[:, 1], s=30, linewidth=1,
24         edgecolors='black', facecolors='none', label='Support Vectors')
25
26     # Thêm nhãn trục và tiêu đề
27     plt.xlabel('Principal Component 1')
28     plt.ylabel('Principal Component 2')
29     plt.title('SVM Visualization')
30     plt.legend()
31     plt.tight_layout()
32     plt.show()
33
34 visualize_2d_svm(X,y)
```

Hình vẽ cho thấy sự phân hoá khá rõ rệt

5. Kết luận

Trong bài tập này, em đã áp dụng thuật toán **SVM** lên bộ dữ liệu **Banknote Authentication**. Các bước thực hiện bao gồm kiểm tra tính tuyến tính của dữ liệu, tuning siêu tham số, xây dựng mô hình, và trực quan hóa kết quả. Qua bài tập, em đã hiểu rõ hơn về thuật toán **SVM** và cách áp dụng nó vào thực tế.