

BÀI TUẦN 04: EVALUATING REGRESSION MODELS PERFORMANCE

1. Thông tin sinh viên

DƯƠNG MINH LƯỢNG-18521071

2. Source

1. Linear Regression

```
1. import pandas as pd
2. import numpy as np
3. from sklearn.linear_model import LinearRegression
4. from sklearn.preprocessing import OneHotEncoder
5. from sklearn.compose import ColumnTransformer
6.
7. # dataset = pd.read_csv('Position_Salaries.csv')
8. dataset_train = pd.read_csv('city_day_AQI_train.csv')
9. dataset_test = pd.read_csv('city_day_AQI_test.csv')
10. del dataset_train['Date']
11. del dataset_test['Date']
12. dataset_train.fillna(dataset_train.groupby('City').transform('mean'),inplace=True)
13. dataset_test.fillna(dataset_test.groupby('City').transform('mean'),inplace=True)
14. dataset_train = dataset_train.fillna(dataset_train.mean())
15. dataset_test = dataset_test.fillna(dataset_test.mean())
16. X = dataset_train.iloc[:, :-2].values
17. Y = dataset_train.iloc[:, -2].values
18. X_test = dataset_test.iloc[:, :-2].values
19. Y_test = dataset_test.iloc[:, -2].values
20. a=X_test
21. ct = ColumnTransformer(transformers=[('encoder', OneHotEncoder(), [0])],
    remainder='passthrough')
22. X = np.array(ct.fit_transform(X))
23. X_test = np.array(ct.transform(X_test))
24.
25. lin_reg = LinearRegression()
```

```
26. lin_reg.fit(X, Y)
27. Y_test_pred = lin_reg.predict(X_test)
28.
29. from sklearn.metrics import mean_squared_error
30. from sklearn.metrics import r2_score
31. from math import sqrt
32. print("SSE", len(X_test)*mean_squared_error(Y_test, Y_test_pred))
33. print("RMSE", sqrt(mean_squared_error(Y_test, Y_test_pred)))
34. r2 = r2_score(Y_test, Y_test_pred)
35. print("r2=", r2)
36. adjusted_r_squared = 1 - (1-r2)*((len(Y_test)-1)/(len(Y_test)-a.shape[1]-1))
37. print("adjusted_r_squared= ", adjusted_r_squared)
```

2. **Polynomial Regression**

```
1. import pandas as pd
2. import numpy as np
3. from sklearn.linear_model import LinearRegression
4. from sklearn.preprocessing import OneHotEncoder
5. from sklearn.compose import ColumnTransformer
6.
7. # dataset = pd.read_csv('Position_Salaries.csv')
8. dataset_train = pd.read_csv('city_day_AQI_train.csv')
9. dataset_test = pd.read_csv('city_day_AQI_test.csv')
10. del dataset_train['Date']
11. del dataset_test['Date']
12. dataset_train.fillna(dataset_train.groupby('City').transform('mean'),inplace=True)
13. dataset_test.fillna(dataset_test.groupby('City').transform('mean'),inplace=True)
14. dataset_train = dataset_train.fillna(dataset_train.mean())
15. dataset_test = dataset_test.fillna(dataset_test.mean())
16. X = dataset_train.iloc[:, :-2].values
17. Y = dataset_train.iloc[:, -2].values
18. X_test = dataset_test.iloc[:, :-2].values
19. Y_test = dataset_test.iloc[:, -2].values
20. a=X_test
21. ct = ColumnTransformer(transformers=[('encoder', OneHotEncoder(), [0])],
    remainder='passthrough')
22. X = np.array(ct.fit_transform(X))
23. X_test = np.array(ct.transform(X_test))
24.
25. from sklearn.preprocessing import PolynomialFeatures
26. poly_transform = PolynomialFeatures(degree=2)
27. X_poly = poly_transform.fit_transform(X)
28. X_poly_test = poly_transform.transform(X_test)
29. poly_lin_reg = LinearRegression()
30. poly_lin_reg.fit(X_poly, Y)
31. from sklearn.metrics import mean_squared_error
32. from sklearn.metrics import r2_score
33. from math import sqrt
34. Y_poly_pred_test = poly_lin_reg.predict(X_poly_test)
35. print("SSE",len(X_test)*mean_squared_error(Y_test, Y_poly_pred_test))
36. print("RMSE", sqrt(mean_squared_error(Y_test, Y_poly_pred_test)))
37. r2=r2_score(Y_test, Y_poly_pred_test)
38. print("r2=",r2)
39. adjusted_r_squared = 1 - (1-r2)*((len(Y_test)-1)/(len(Y_test)-a.shape[1]-1))
40. print("adjusted_r_squared= ",adjusted_r_squared)
```

3. Support Vector Regression (SVR)

```
1. import pandas as pd
2. import numpy as np
3. from sklearn.preprocessing import OneHotEncoder
4. from sklearn.compose import ColumnTransformer
5.
6. # dataset = pd.read_csv('Position_Salaries.csv')
7. dataset_train = pd.read_csv('city_day_AQI_train.csv')
8. dataset_test = pd.read_csv('city_day_AQI_test.csv')
9. del dataset_train['Date']
10. del dataset_test['Date']
11. dataset_train.fillna(dataset_train.groupby('City').transform('mean'),inplace=True)
12. dataset_test.fillna(dataset_test.groupby('City').transform('mean'),inplace=True)
13. dataset_train = dataset_train.fillna(dataset_train.mean())
14. dataset_test = dataset_test.fillna(dataset_test.mean())
15. X = dataset_train.iloc[:, :-2].values
16. Y = dataset_train.iloc[:, -2].values
17. X_test = dataset_test.iloc[:, :-2].values
18. Y_test = dataset_test.iloc[:, -2].values
19. a=X_test
20. ct = ColumnTransformer(transformers=[('encoder', OneHotEncoder(), [0])],
    remainder='passthrough')
21. X = np.array(ct.fit_transform(X))
22. X_test = np.array(ct.transform(X_test))
23. Y = Y.reshape(-1,1)
24. Y_test = Y_test.reshape(-1,1)
25. from sklearn.preprocessing import StandardScaler
26. sc_X = StandardScaler()
27. sc_y = StandardScaler()
28. X_trans = sc_X.fit_transform(X)
29. Y_trans = sc_y.fit_transform(Y)
30. X_trans_test = sc_X.transform(X_test)
31. Y_trans_test = sc_y.transform(Y_test)
32. from sklearn.svm import SVR
33. regressor = SVR(kernel = 'rbf')
34. regressor.fit(X_trans, Y_trans)
35. def predict(model, X, SC_X, SC_Y):
36.     X_trans = SC_X.transform(X)
37.     Y_trans_pred = model.predict(X_trans)
```

```

38. Y_pred = SC_Y.inverse_transform(Y_trans_pred)
39. return Y_pred
40. Y_pred_train = predict(regressor, X, sc_X, sc_y)
41. Y_pred_test = predict(regressor, X_test, sc_X, sc_y)
42. from sklearn.metrics import mean_squared_error
43. from sklearn.metrics import r2_score
44. from math import sqrt
45. print("SSE",len(X_test)*mean_squared_error(Y_test,Y_pred_test))
46. print("RMSE", sqrt(mean_squared_error(Y_test, Y_pred_test)))
47. r2 = r2_score(Y_test, Y_pred_test)
48. print("r2=",r2)
49. adjusted_r_squared = 1 - (1-r2)*((len(Y_test)-1)/(len(Y_test)-a.shape[1]-1))
50. print("adjusted_r_squared=",adjusted_r_squared)

```

4. Decision Tree Regression

```

1. import pandas as pd
2. import numpy as np
3. from sklearn.preprocessing import OneHotEncoder
4. from sklearn.compose import ColumnTransformer
5.
6. # dataset = pd.read_csv('Position_Salaries.csv')
7. dataset_train = pd.read_csv('city_day_AQI_train.csv')
8. dataset_test = pd.read_csv('city_day_AQI_test.csv')
9. del dataset_train['Date']
10. del dataset_test['Date']
11. dataset_train.fillna(dataset_train.groupby('City').transform('mean'),inplace=True)
12. dataset_test.fillna(dataset_test.groupby('City').transform('mean'),inplace=True)
13. dataset_train = dataset_train.fillna(dataset_train.mean())
14. dataset_test = dataset_test.fillna(dataset_test.mean())
15. X = dataset_train.iloc[:, :-2].values
16. Y = dataset_train.iloc[:, -2].values
17. X_test = dataset_test.iloc[:, :-2].values
18. Y_test = dataset_test.iloc[:, -2].values
19. a=X_test
20. ct = ColumnTransformer(transformers=[('encoder', OneHotEncoder(), [0])],
    remainder='passthrough')
21. X = np.array(ct.fit_transform(X))
22. X_test = np.array(ct.transform(X_test))
23. Y = Y.reshape(-1,1)
24. Y_test = Y_test.reshape(-1,1)
25. from sklearn.tree import DecisionTreeRegressor

```

```

26. regressor = DecisionTreeRegressor()
27. regressor.fit(X, Y)
28. Y_pred_test=regressor.predict(X_test)
29.
30. from sklearn.metrics import mean_squared_error
31. from math import sqrt
32. print("SSE",len(X_test)*mean_squared_error(Y_test, Y_pred_test))
33. print("RMSE", sqrt(mean_squared_error(Y_test, Y_pred_test)))
34. from sklearn.metrics import r2_score
35. r2=r2_score(Y_test, Y_pred_test)
36. print("r2=",r2)
37. adjusted_r_squared = 1 - (1-r2)*((len(Y_test)-1)/(len(Y_test)-a.shape[1]-1))
38. print("adjusted_r_squared= ",adjusted_r_squared)

```

5. **Random Forest Regression**

```

1. import pandas as pd
2. import numpy as np
3. from sklearn.preprocessing import OneHotEncoder
4. from sklearn.compose import ColumnTransformer
5.
6. # dataset = pd.read_csv('Position_Salaries.csv')
7. dataset_train = pd.read_csv('city_day_AQI_train.csv')
8. dataset_test = pd.read_csv('city_day_AQI_test.csv')
9. del dataset_train['Date']
10. del dataset_test['Date']
11. dataset_train.fillna(dataset_train.groupby('City').transform('mean'),inplace=True)
12. dataset_test.fillna(dataset_test.groupby('City').transform('mean'),inplace=True)
13. dataset_train = dataset_train.fillna(dataset_train.mean())
14. dataset_test = dataset_test.fillna(dataset_test.mean())
15. X = dataset_train.iloc[:, :-2].values
16. Y = dataset_train.iloc[:, -2].values
17. X_test = dataset_test.iloc[:, :-2].values
18. Y_test = dataset_test.iloc[:, -2].values
19. a=X_test
20. ct = ColumnTransformer(transformers=[('encoder', OneHotEncoder(), [0])],
    remainder='passthrough')
21. X = np.array(ct.fit_transform(X))
22. X_test = np.array(ct.transform(X_test))
23. Y = Y.reshape(-1,1)
24. Y_test = Y_test.reshape(-1,1)
25. from sklearn.ensemble import RandomForestRegressor
26. regressor = RandomForestRegressor(n_estimators = 7,random_state=0 )

```

```

27. regressor.fit(X, Y)
28. Y_pred_test=regressor.predict(X_test)
29. from sklearn.metrics import mean_squared_error
30. from sklearn.metrics import r2_score
31. from math import sqrt
32. print("SSE",len(X_test)*mean_squared_error(Y_test,Y_pred_test))
33. print("RMSE", sqrt(mean_squared_error(Y_test, Y_pred_test)))
34. r2 = r2_score(Y_test, Y_pred_test)
35. print("r2=",r2)
36. adjusted_r_squared = 1 - (1-r2)*((len(Y_test)-1)/(len(Y_test)-a.shape[1]-1))
37. print("adjusted_r_squared= ",adjusted_r_squared)

```

3. Kết quả

	Squared Sum (SSE)	RMSE	R^2	R_adjusted^2
Linear Regression	26800886.85	55	0.838938357454	0.838701662751
Polynomial Regression	18739039.877	46	0.887386542	0.88722104
Support Vector Regression (SVR)	19250528.9379	46.613	0.8843127159	0.88414270
Decision Tree Regression	35066732.743865	63	0.789264	0.788954532
Random Forest Regression	22161420.8821	50	0.8668195	0.8666238

Nhận xét:

- Với Hàm lỗi Squared Sum (SSE) thì phương pháp Polynomial Regression là tốt nhất.
- Với Hàm lỗi Root Mean Squared (RMSE) thì phương pháp Polynomial Regression là tốt nhất.
- Với Hàm đánh giá R^2 thì phương pháp Polynomial Regression là tốt nhất.
- Với Hàm đánh giá R_{adjusted}^2 thì phương pháp Polynomial Regression là tốt nhất.
- Với tập dữ liệu và test trên thì với phương pháp Polynomial Regression thì cho tất cả các đánh giá độ đo là tốt nhất.
- Với tập dữ liệu và test trên thì với phương pháp Decision Tree Regression thì cho tất cả các đánh giá độ đo là tệ nhất.
- Nếu Random Forest Regression không có `random_state = 0` thì sẽ cho những kết quả khác nhau sau mỗi lần chạy.