PHIẾU GHI NHẬN ĐIỀU CHỈNH ĐỒ ÁN MÔN HỌC

Môn học: Hệ Thống Thương Mại Thông Minh (504049)

Học kỳ: 1

Năm học: 2021 - 2022

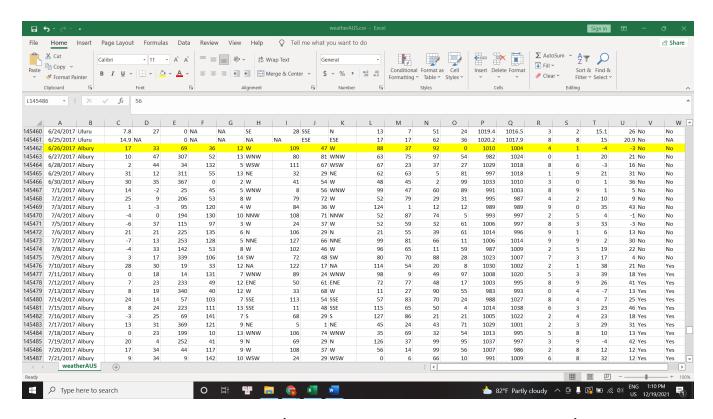
Mã nhóm: 12

Điều chỉnh 1: Thêm dữ liệu mới cho dataset weather AUS.csv

```
    [ ] data.shape##Tổng giá trị dữ liệu và thuộc tính dataset cũ (145460, 23)
    [ ] data.shape##Tổng giá trị dữ liệu và thuộc tính thêm mới (161154, 23)
```

Dữ liệu weather AUS.csv ban đầu được nghiên cứu bao gồm: Nhiệt độ, lượng mưa, độ bốc hơi, độ ẩm, áp xuất, hướng gió, vận tốc gió và mật độ mây được thu thập trong 10 năm (từ 1/12/2008 đến 25/6/2017) ở 49 trạm khí tượng của nước Úc bao gồm 145460 dòng với 23 thuộc tính.

Tổng giá trị được thêm mới gồm 15694 dòng với 23 thuộc tính vẫn được chỉ định ở 49 trạm ở Úc. Được thu thập thêm từ 26/6/2017 đến 25/6/2018.



Vậy ta sẽ có dataset mới bao gồm được thu thập từ 11 năm (1/12/2008 đến 25/6/2018) với 161154 dòng và 23 thuộc tính.

Thêm dòng giá trị mới và so sánh độ tương thích dữ liệu so với dataset cũ:

- Date: ngày tháng năm thu thập dòng dữ liệu (mm/dd/yyyy) sẽ được thêm 6/26/2017 đến 6/25/2018 (321 dòng).
- Location: địa điểm (name) được thêm 49 địa điểm ở nước Úc tương ứng với 6/26/2017 đến
 6/25/2018 (mỗi địa điểm sẽ được thêm 321 dòng dữ liệu mới).
- MinTemp: nhiệt độ tối thiểu (oC) được dựa trên giá trị của dataset có sẵn giao động từ -8.5 đến 33.9 oC và được lấy tương thích nếu nhiệt độ cao gây ra sự bốc hơi thì khả năng mưa ngày hôm nay và ngày mai sẽ cao.
- MaxTemp: nhiệt độ tối đa (oC) được dựa trên giá trị của dataset có sẵn giao động từ -4.8 đến 44.1 oC và được lấy tương thích nếu nhiệt độ cao gây ra sự bốc hơi thì khả năng mưa ngày hôm nay và ngày mai sẽ cao.

- Rainfall: lượng mưa được ghi nhận trong ngày (mm) được thêm vào dựa trên giá trị cũ 0 đến 371 mm, nếu ngày hôm đó mưa sẽ được lấy lượng mưa rơi trong ngày và dự đoán cho ngày tiếp theo.
- Evaporation: sự bốc hơi được ghi nhận từ 24 giờ đêm đến 9h sáng (mm) được lấy từ 0 đến
 145 mm dự trên data cũ nếu ngày hôm đó có nhiệt độ cao khả năng bốc hơi cũng sẽ cao.
- Sunshine: số giờ nắng trong ngày (hour) được lấy từ 0 đến 14.5hour dựa trên data cũ sẽ phụ thuộc vào nhiệt độ trong ngày và sự bốc hơi.
- WindGustDir: hướng gió (N-E-W-S).
- WindGustSpeed: vận tốc gió (km/h) được lấy từ 0 đến 135 km/h của data cũ.
- WindDir9am: hướng gió lúc 9 giờ sáng (N-E-W-S).
- WindDir3pm: hướng gió lúc 3 giờ chiều (N-E-W-S).
- WindSpeed9am: vận tốc gió lúc 9 giờ sáng (km/h) được lấy 0 đến 130km/h của data cũ dựa trên hướng gió 9h sáng và 3h chiều.
- WindSpeed3pm: vận tốc gió lúc 3 giờ chiều (km/h) được lấy 0 đến 87km/h của data cũ dựa trên hướng gió 9h sáng và 3h chiều.
- Humidity9am: độ ẩm không khí lúc 9 giờ sáng (g/m3) được lấy từ 0 đến 100 g/m3 của data cũ dựa trên nhiệt độ tối đa và tối thiểu trong ngày, sự bốc hơi của ngày hôm đó.
- Humidity3pm: độ ẩm không khí lúc 3 giờ chiều (g/m3) được lấy từ 0 đến 100 g/m3 của data cũ dựa trên nhiệt độ tối đa và tối thiểu trong ngày, sự bốc hơi của ngày hôm đó.
- Pressure9am: áp suất không khí lúc 9 giờ sáng (N/m2) được lấy từ 980.5 1041 dựa trên data cũ áp suất không khí cao thì ghi nhận mưa trong ngày là giá trị No, ấp suất không khí thấp thì hôm đó ghi nhận được trời sẽ đổ mưa trong ngày dựa trên các trị nhiệt độ tối đa tối thiểu và độ ẩm hướng gió vận tốc gió.
- Pressure3pm: áp suất không khí lúc 3 giờ chiều (N/m2) được lấy từ 977.1 1039.6 dựa trên data cũ áp suất không khí cao thì ghi nhận mưa trong ngày là giá trị No, ấp suất không khí thấp thì hôm đó ghi nhận được trời sẽ đổ mưa trong ngày dựa trên các trị nhiệt độ tối đa tối thiểu và độ ẩm hướng gió vận tốc gió
- Cloud9am: mật độ mây theo các mức độ lúc 9 giờ sáng được lấy từ (0-9).

- Cloud3pm: mật độ mây theo các mức độ lúc 3 giờ chiều được lấy từ (0-9).
- Temp9am: nhiệt độ ghi nhận lúc 9 giờ sáng (oC) được lấy từ -7.2 dến 40.2 oC ảnh hưởng đến sự bốc hơi và độ ẩm trong không khí dẫn đến tích tụ nước và gây mưa
- Temp3pm: nhiệt độ ghi nhận lúc 3 giờ chiều (oC) được lấy từ -5.4 đến 46.7 oC ảnh hưởng đến sự bốc hơi và độ ẩm trong không khí dẫn đến tích tụ nước và gây mưa
- RainToday: thuộc tính ghi nhận mưa trong ngày (yes/no).
- RainTomorrow: thuộc tính dự đoán mưa trong ngày tiếp theo (yes/no).

Điều chỉnh 2: Sử dụng dataset mới (đã thêm dữ liệu mới vào dataset cũ)

Với các dữ liệu bị thiếu trong dataset weatherAUS.csv sẽ được xử lý missing value như sau:

- Bước 1: Loading dataset
- Bước 2: Phân tích các thuộc tính thành 2 loại (kiểu số và kiểu kí tự).
- Bước 3: Có nhiều cách xử lý dữ liệu bị thiếu: Sử dụng trung vị hoặc giá trị trung bình, chọn giá trị ngẫu nhiên để đưa vào hoặc giá trị có tần suất xuất hiện nhiều nhất của mỗi thuộc tính. Ở trường hợp dataset này sử dụng điền các giá trị còn thiếu bằng giá trị có tần suất xuất hiện nhiều nhất của mỗi thuộc tính. Sử dụng hàm fillna() để tìm dữ liệu xuất hiện nhiều nhất lấp vào chỗ trống.

Về độ tương thích của dữ liệu sao khi missing value: sẽ được so sánh bằng cách sử dụng điền các giá trị còn thiếu bằng giá trị có tần suất xuất hiện nhiều nhất của mỗi thuộc tính. Bằng cách dự đoán và kiểm tra kết quả đánh giá theo thuật toán Sử dụng thuật toán Logistic Regression.

• Trước khi missing value

```
y_pred_test = logreg.predict(X_test)
y_pred_test
array(['No', 'No', 'No', 'No', 'No', 'Yes'], dtype=object)
```

Kiểm tra đánh giá độ chính xác

```
[35]: from sklearn.metrics import accuracy_score

print('Model accuracy score: {0:0.4f}'. format(accuracy_score(y_test, y_pred_test)))

Model accuracy score: 0.8498
```

• Sau khi missing value:

```
model = LogisticRegression(max_iter=500)
model.fit(x_train, y_train)
predicted=model.predict(x_test)
conf = confusion_matrix(y_test, predicted)
print ("The accuracy of Logistic Regression is : ", accuracy_score(y_test, predicted)*100, "%")
The accuracy of Logistic Regression is : 84.00494267865723 %
```

The accuracy of Logistic Regression is: 84.00494267865723 %

Kết luận: vậy giá trị điền vào NaN trống của data về độ tương thích là không chênh lệch kết quả quá nhiều so với trước. Khi dự đoán xác suất của trước khi missing là 84.98% còn sau khi missing dự đoán xác suất là 84%.

Điều chỉnh 3: Đưa ra đánh giá với từng thuật toán sử dụng trong bài.

Với đề tài dự báo thời tiết ở nước Úc, quá trình phân tích train và test tập dữ liệu. Với thuật toán Logistic Regression cho ra độ chính xác 84%, thuật toán XGBoost 84.75%, thuật toán Gaussian Naive Bayes 80.1%, thuật toán Bernoulli Naive Bayes 76.67% và thuật toán Random Forest Model cho đô chính xác 85.48%.

Vậy với chủ đề dự báo thời tiết ở nước Úc thuật toán Random Forest Model cho ra dự đoán chính xác cao nhất. Ngược lại thuật toán Bernoulli Naive Bayes cho ra dự đoán thấp nhất.

Thuật toán Random Forest Model cho ra kết quả chính xác nhất tuy nhiên tốn nhiều thời gian hơn so với các thuật toán còn lại.

Thuật toán XGBoost cho ra kết quả dự đoán thời gian ngắn nhất mà độ chính xác so với Thuật toán Random Forest Model cũng không chênh lệch nhiều.

Thuật toán Logistic Regression cho ra độ chính xác cao đứng thứ 3 tuy nhiên thời gian dự đoán nhanh mà độ chính xác so với Thuật toán Random Forest Model cũng không chênh lệch nhiều.

Thuật toán Gaussian Naive Bayes cùng với thuật toán Bernoulli Naive Bayes cho ra độ chính xác đứng thứ 4 và thứ 5 thời gian dự đoán ngắn tuy nhiên kết quả dự đoán chênh lệch cao.

Kết luận: đối với dataset dự báo thời tiết ở nước Úc, chúng em sẽ cho thuật toán XGBoost là phù hợp nhất vì thời gian dự đoán nhanh và cho ra kết quả chính xác không quá chênh lệch với Random Forest Model. Random Forest Model cho ra kết quả chính xác cao nhất nhưng thời gian dự đoán khá lâu nên chúng em sẽ không chọn thuật toán này là tốt nhất.