

Penerapan Algoritma Logistic Regression pada Dataset weatherAUS untuk Prediksi Hujan Di Australia



Kelompok 6

- | | |
|--------------------------|-------------|
| 1. NURUL NAJWA SABILLA | (312110451) |
| 2. LAELA NUR ROHMAH | (312110425) |
| 3. ALVINA DAMAYANTI | (312110125) |
| 4. SARA KHUSNUL MUMTAZAH | (312110319) |

Kelas TI.21.A.3



Pendahuluan



01.

Gambaran singkat tentang Big Data Analytics dan PySpark

Big Data Analytics adalah proses yang digunakan untuk mengambil pola tersembunyi, korelasi yang tidak diketahui, tren pasar, dan preferensi pelanggan.

PySpark adalah library open-source Python yang menyediakan antarmuka untuk Apache Spark, sebuah kerangka kerja pengolahan data berdistribusi.

02.

Mengapa kita memilih dataset weatherAUS untuk analisis?

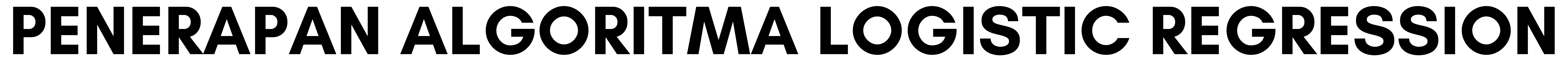
Karena kasusnya dataset weatherAUS yang digunakan untuk memprediksi apakah akan terjadi hujan besok berdasarkan data cuaca pada hari sebelumnya.

Metode machine learning yang digunakan adalah Logistic Regression karena metode ini sangat populer digunakan.

03.

Tujuan Presentasi

Dapat menghasilkan model yang memiliki akurasi yang baik dalam memprediksi hujan pada esok hari di wilayah Australia.



2 CARA LOGISTIC REGRESSION DIGUNAKAN DALAM PEMODELAN PREDIKSI APAKAH HUJAN BESOK (RAINTOMMOROW) PADA DATASET WEATHERAUS

1. Melakukan skema awal untuk dataset

3. Clening data sebelum dilakukan tindakan selanjutnya, yaitu mengisi missing value

Cleaning dilakukan dengan imputasi mean, median, modus, sehingga data bebas dari kekosongan data

4. Pada tahap cleaning ini variabel berganti dari “df” jadi “numerical”. Variabel numerical ini akan dijadikan features dalam model dan target itu yaitu “RainTommorow”

```
In [74]: df.printSchema()

root
|-- Date: string (nullable = true)
|-- Location: string (nullable = true)
|-- MinTemp: float (nullable = true)
|-- MaxTemp: float (nullable = true)
|-- Rainfall: float (nullable = true)
|-- Evaporation: float (nullable = true)
|-- Sunshine: float (nullable = true)
|-- WindGustDir: string (nullable = true)
|-- WindGustSpeed: float (nullable = true)
|-- WindDir9am: string (nullable = true)
|-- WindDir3pm: string (nullable = true)
|-- WindSpeed9am: float (nullable = true)
|-- WindSpeed3pm: float (nullable = true)
|-- Humidity9am: float (nullable = true)
|-- Humidity3pm: float (nullable = true)
|-- Pressure9am: float (nullable = true)
|-- Pressure3pm: float (nullable = true)
|-- Cloud9am: float (nullable = true)
|-- Cloud3pm: float (nullable = true)
|-- Temp9am: float (nullable = true)
|-- Temp3pm: float (nullable = true)
|-- RainToday: string (nullable = true)
|-- RainTomorrow: integer (nullable = false)
```

```
In [66]: assembler VectorAssembler(inputCols=numerical, outputCol='features')
output = assembler.transform(df)
output = output.select('features', 'RainTomorrow')
output = output.withColumnRenamed('RainTomorrow', 'label')

output.show(10)
# numerical
```

```
+-----+-----+
|               features | label |
+-----+-----+
|[13.39999996185302...]| 0 |
|[7.400000009536743...]| 0 |
|[12.89999996185302...]| 0 |
|[9.199999980926513...]| 0 |
|[17.5,32.29999923...]| 0 |
|[14.6000003814697...]| 0 |
|[14.3000001907348...]| 0 |
|[7.699999980926513...]| 0 |
|[9.699999980926513...]| 1 |
|[13.1000003814697...]| 0 |
+-----+-----+
only showing top 10 rows
```

1. Bagi DataFrame menjadi train_data (80%) dan test_data (20%)

2. Melakukan modelling Logistic Regression

Lakukan percobaan model dengan LogisticRegression

only showing top 10 rows

```
In [69]: my_eval = BinaryClassificationEvaluator()
my_final_roc = my_eval.evaluate(predictions_and_labels.predictions)
my_final_roc
```

Out[69]: 0.9999998149655362



HASIL ANALISIS

1. RINGKASAN HASIL ANALISIS LOGISTIC REGRESSION PADA DATASET WEATHERAUS

hasil dari evaluasi sebelumnya metrik ROC (Receiver Operating Characteristic) yang dihitung oleh BinaryClassificationEvaluator. Nilai ROC ini mendekati 1, yang menunjukkan bahwa model klasifikasi biner Anda memiliki kinerja yang sangat baik dalam membedakan antara kelas positif dan kelas negatif. Dalam konteks evaluasi metrik ROC, nilai AUC-ROC mendekati 1 (atau 100%) menunjukkan bahwa model klasifikasi memiliki kinerja yang sangat baik dalam membedakan antara dua kelas yang berbeda.

2. VISUALISASI HASIL, SEPERTI GRAFIK HUBUNGAN VARIABEL PREDIKTOR DAN VARIABEL TARGET

Visualization

```
In [62]: valid_predictions_count = predictions_and_labels.predictions[col("prediction") == col("label")].count()
invalid_predictions_count = predictions_and_labels.predictions[col("prediction") != col("label")].count()
accuracy_predictions_count = valid_predictions_count / predictions_and_labels.predictions.count() * 100

print("akurasi prediksi yang benar dari label dan prediction: " + str(valid_predictions_count))
print("akurasi prediksi yang salah dari label dan prediction: " + str(invalid_predictions_count))
print("akurasi untuk prediksi yang benar: " + str(accuracy_predictions_count))
```

```
akurasi prediksi yang benar dari label dan prediction: 28847
akurasi prediksi yang salah dari label dan prediction: 0
akurasi untuk prediksi yang benar: 100.0
```

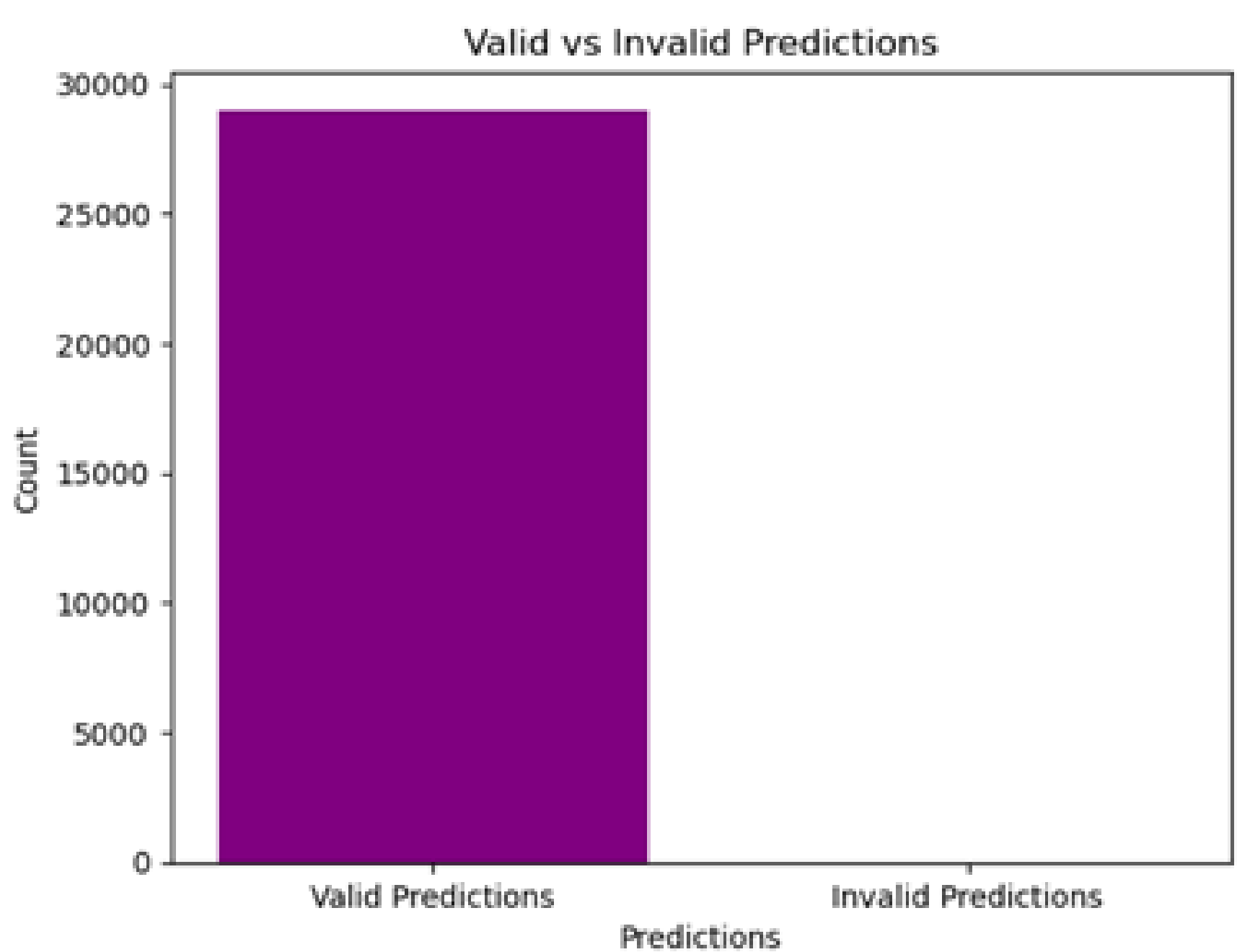
Hasil output menunjukkan bahwa dari seluruh prediksi yang telah dilakukan, tidak ada yang salah (0 prediksi yang salah). Seluruh prediksi (28847 prediksi) cocok atau sesuai dengan nilai label yang seharusnya. Dalam istilah akurasi, ini berarti bahwa model atau algoritma yang digunakan untuk membuat prediksi memiliki akurasi sebesar 100%.

Berikut ini visualisasinya

```
In [63]: valid_predictions_count = 29006
invalid_predictions_count = 0

# Menyiapkan data
categories = ['Valid Predictions', 'Invalid Predictions']
counts = [valid_predictions_count, invalid_predictions_count]

# Membuat grafik batang
plt.bar(categories, counts, color=['purple', 'grey'])
plt.xlabel('Predictions')
plt.ylabel('Count')
plt.title('Valid vs Invalid Predictions')
plt.show()
```



1. Label

```
In [64]: label_to_one = 6338
label_to_zero = 22668

# Menghitung total prediksi
total_predictions = label_to_zero + label_to_one

# Menghitung persentase
percentages = [label_to_zero / total_predictions * 100, label_to_one / total_predictions * 100]

categories = ['No', 'Yes']

# Membuat grafik pie
plt.pie(percentages, labels=categories, autopct='%1.1f%%', explode=[0, 0.1], colors=['grey', 'purple'], shadow=True)
plt.axis('equal') # Equal aspect ratio ensures that pie is drawn as a circle.
plt.title('Persentase Hari Besok Turun Hujan Dari Hasil Target')
plt.show()
```



2. Prediksi

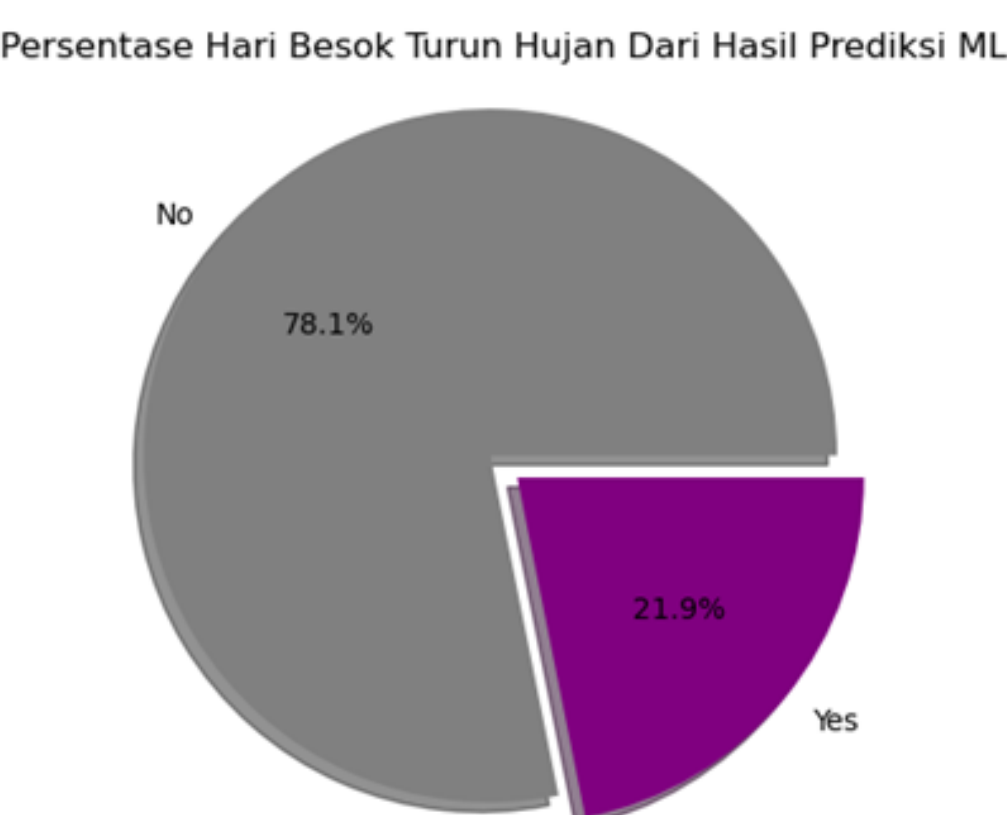
```
In [65]: prediction_to_one = 6338
prediction_to_zero = 22668

# Menghitung total prediksi
total_predictions = prediction_to_zero + prediction_to_one

# Menghitung persentase
percentages = [prediction_to_zero / total_predictions * 100, prediction_to_one / total_predictions * 100]

categories = ['No', 'Yes']

# Membuat grafik pie
plt.pie(percentages, labels=categories, autopct='%1.1f%%', explode=[0, 0.1], colors=['grey', 'purple'], shadow=True)
plt.axis('equal') # Equal aspect ratio ensures that pie is drawn as a circle.
plt.title('Persentase Hari Besok Turun Hujan Dari Hasil Prediksi ML')
plt.show()
```



3. INTERPRETASI HASIL DAN TEMUAN PENTING

Dengan kata lain, model ini berhasil memprediksi semua data dengan benar sesuai dengan label yang ada. Ini adalah hasil yang sangat baik dan menunjukkan bahwa model tersebut sangat andal dalam melakukan prediksi pada dataset yang digunakan. Akurasi sebesar 100% adalah tujuan yang sangat diinginkan dalam banyak kasus, tetapi juga perlu memeriksa apakah ini adalah hasil yang sesuai dengan ekspektasi atau apakah mungkin ada masalah seperti overfitting.



Hasil dan Kesimpulan

1

TEMUAN PENTING DAN INSIGHT (WAWASAN) YANG DIPEROLEH

- Mengetahui klasifikasi data pada logistic regression berguna memprediksi label data atau target
- Lalu, mengetahui probabilitas prediksi yang menjadi output. Hal ini memungkinkan kita untuk mengukur tingkat keyakinan atau kepastian dalam prediksi.
- Serta logistic regression itu membutuhkan metrik evaluasi untuk mengukur kinerja model salah satunya dengan metric ROC.

2

KESIMPULAN DAN SARAN

Berdasarkan hasil dataset wetherAUS ini, kemungkinan terjadi hujan pada hari esok hari di Australia diprediksi dengan logistic regression mendekati 1, yang menunjukkan bahwa model sangat baik. Lalu, melihat akurasi antara hubungan variabel prediktor/prediksi dan variabel target menjadi metrik evaluasi kinerja model untuk mengukur sejauh mana model dapat membedakan actual label dengan prediksi hasil ML adalah 100%. Namun, hal ini bisa menimbulkan masalah yaitu overfitting (algoritma ML terlalu spesifik terhadap data). Saran mengenai hal itu dapat melakukan penyederhanaan model dengan mengganti algoritma untuk mencari model yang terbaik.