

Лабораторная работа

Использование фреймворка Apache Spark для построения интеллектуальных систем, использующих большие данные

Выполнила: Короткова Инга Сергеевна

2020 год

**Цель работы** – получить навыки проектирования интеллектуальных моделей с помощью библиотеки MLlib фреймворка Apache для задач классификации пациентов с заболеваниями печени.

**Задачи:**

1. Создать объект SparkSession.
2. Загрузить набор данных, содержащий результаты лабораторных исследований пациентов Indian\_Liver\_Patient\_Dataset\_ILPDh.csv, у части из которых диагностировано заболевание печени.
3. Произвести исследовательский анализ данных.
4. Разделить набор данных на обучающую и тестовую выборку.
5. Обучить модели классификации на основе алгоритмов случайный лес и логистическая регрессия.
6. Оценить эффективность моделей на тестовой выборке с помощью матрицы неточностей, критериев полноты Recall и точности Precision

```
1 !pip install pyspark

Requirement already satisfied: pyspark in /usr/local/lib/python3.6/dist-packages (3.0.1)
Requirement already satisfied: py4j==0.10.9 in /usr/local/lib/python3.6/dist-packages (from pyspark) (0.10.9)
```

```
1 from pyspark.sql import SparkSession
2 from pyspark.ml.feature import VectorAssembler
3 from pyspark.ml.classification import RandomForestClassifier
4 from pyspark.ml.classification import LogisticRegression
5
```

▼ **Создание SparkSession**

Создадим объект SparkSession.

spark – объект SparkSession.

```
1 spark = SparkSession.builder.appName('Liver Patient Dataset AnalysisApp').getOrCreate()
```

▼ **Загрузка данных**

Загрузим датасет, содержащий результаты лабораторных исследований пациентов Indian\_Liver\_Patient\_Dataset\_ILPDh.csv, у части из которых диагностировано заболевание печени.

df – объект датафрейм, в который происходит загрузка данных из файла 'Indian\_Liver\_Patient\_Dataset\_ILPDh.csv'.

```
1 df = spark.read.csv('Indian_Liver_Patient_Dataset_ILPDh.csv', inferSchema=True, sep=";", header=True)
```

```
1 df.head()

Row(Age=65, Gender=1, Total_Bilirubin=0.7, Direct_Bilirubin=0.1, Alkaline_Phosphotase=187, Alamine_Aminotransferase=16, Aspartate_Ar
```

▼ **Произведем исследовательский анализ данных:**

- получим объём исследуемых данных;
- получим число атрибутов и их типы данных;
- посмотрим распределение числа примеров классов.

```
1 print(f"Total rows {df.count()} and columns {len(df.columns)}")
```

```
Total rows 579 and columns 11
```

```
1 df.printSchema()
```

```

root
|-- Age: integer (nullable = true)
|-- Gender: integer (nullable = true)
|-- Total_Bilirubin: double (nullable = true)
|-- Direct_Bilirubin: double (nullable = true)

```

```
1 df.groupby('Class_ID').count().show()
```

```

+-----+-----+
|Class_ID|count|
+-----+-----+
|      1|  414|
|      2|  165|
+-----+-----+

```

## ▼ Разделим набор данных на обучающую и тестовую выборку.

Для обучения модели необходимо выбрать признаки, применяемые в модели, сформировать вектор признаков features и добавить его в датафрэйм. В работе выборка делится на обучающую и тестовую в отношении 3:1 (75%: 25%).

```

1 df_assembler = VectorAssembler(inputCols=['Age', 'Gender', 'Total_Bilirubin', \
2 'Direct_Bilirubin', 'Alkaline_Phosphotase', \
3 'Alamine_Aminotransferase', 'Aspartate_Aminotransferase', \
4 'Total_Protiens', 'Albumin', 'Albumin_and_Globulin_Ratio'], outputCol="features")
5
6 df = df_assembler.transform(df)
7 df.printSchema()

```

```

root
|-- Age: integer (nullable = true)
|-- Gender: integer (nullable = true)
|-- Total_Bilirubin: double (nullable = true)
|-- Direct_Bilirubin: double (nullable = true)
|-- Alkaline_Phosphotase: integer (nullable = true)
|-- Alamine_Aminotransferase: integer (nullable = true)
|-- Aspartate_Aminotransferase: integer (nullable = true)
|-- Total_Protiens: double (nullable = true)
|-- Albumin: double (nullable = true)
|-- Albumin_and_Globulin_Ratio: double (nullable = true)
|-- Class_ID: integer (nullable = true)
|-- features: vector (nullable = true)

```

```

1 train_df,test_df=df.randomSplit([0.75,0.25])
2 print('Training set info:')
3 train_df.groupby('Class_ID').count().show()
4 print('Test set info:')
5 test_df.groupby('Class_ID').count().show()

```

```

Training set info:
+-----+-----+
|Class_ID|count|
+-----+-----+
|      1|  309|
|      2|  119|
+-----+-----+

```

```

Test set info:
+-----+-----+
|Class_ID|count|
+-----+-----+
|      1|   105|
|      2|    46|
+-----+-----+

```

## ▼ Обучение моделей

Обучим модели классификации на основе алгоритмов случайный лес и логистическая регрессия.

### Random Forest

```

1 rf_classifier=RandomForestClassifier(labelCol='Class_ID', numTrees=35).fit(train_df)
2 rf_predictions=rf_classifier.transform(test_df)
3 rf_predictions.select(['Class_ID', 'prediction']).show(30,False)
4 rf_predictions.printSchema()
5 print(rf_classifier.featureImportances)

```

Class_ID	prediction
2	1.0
1	1.0
2	1.0
2	2.0
1	1.0
2	1.0
2	1.0
1	1.0
1	2.0
1	1.0
1	1.0
1	1.0
2	1.0
2	1.0
2	1.0
2	2.0
2	2.0
1	1.0
1	1.0
1	2.0
2	1.0
2	1.0
2	2.0
1	1.0
1	2.0
2	1.0
1	1.0
1	1.0
1	1.0
2	2.0

only showing top 30 rows

```

root
|-- Age: integer (nullable = true)
|-- Gender: integer (nullable = true)
|-- Total_Bilirubin: double (nullable = true)
|-- Direct_Bilirubin: double (nullable = true)
|-- Alkaline_Phosphotase: integer (nullable = true)
|-- Alamine_Aminotransferase: integer (nullable = true)
|-- Aspartate_Aminotransferase: integer (nullable = true)
|-- Total_Protiens: double (nullable = true)
|-- Albumin: double (nullable = true)
|-- Albumin_and_Globulin_Ratio: double (nullable = true)
|-- Class_ID: integer (nullable = true)
|-- features: vector (nullable = true)
|-- rawPrediction: vector (nullable = true)
|-- probability: vector (nullable = true)
|-- prediction: double (nullable = false)

(10,[0,1,2,3,4,5,6,7,8,9],[0.1386747315496918,0.0055715843908345,0.09620664383827936,0.11052536619276124,0.1915840523012491,0.15981:

```

Logistic Regression

```

1 log_reg=LogisticRegression(labelCol='Class_ID').fit(train_df)
2 lr_predictions = log_reg.transform(test_df)
3 lr_predictions.select(['Class_ID', 'prediction']).show(30,False)
4 lr_predictions.printSchema()

```

Class_ID	prediction
2	2.0
1	2.0
2	2.0
2	2.0
1	1.0
2	2.0
2	2.0
1	1.0
1	2.0
1	1.0
1	2.0
1	1.0
2	2.0
2	2.0
2	1.0
2	2.0
2	2.0
1	2.0
1	1.0
1	2.0
2	1.0
2	1.0
2	2.0
1	1.0
1	2.0
2	1.0
1	1.0
1	2.0
1	1.0
2	1.0

only showing top 30 rows

-----

▼ Оценка моделей.

Оценим эффективность моделей на тестовой выборке с помощью матрицы неточностей, критериев полноты Recall и точности Precision.

```
|-- Aspartate_Aminotransferase: integer (nullable = true)
```

Для этого определим true\_positives, true\_negatives, false\_positives, false\_negatives – количество истинноположительных, истинноотрицательных, ложно-положительных и ложноотрицательных случаев соответственно.

При помощи переменных recall и precision, выведем значения полноты и точности классифицирующей модели.

```
|-- probability: vector (nullable = true)

1 true_positives = rf_predictions[(rf_predictions.Class_ID == 2) & (rf_predictions.prediction == 2)].count()
2 true_negatives = rf_predictions[(rf_predictions.Class_ID == 1) & (rf_predictions.prediction == 1)].count()
3 false_positives = rf_predictions[(rf_predictions.Class_ID == 1) & (rf_predictions.prediction == 2)].count()
4 false_negatives = rf_predictions[(rf_predictions.Class_ID == 2) & (rf_predictions.prediction == 1)].count()
5 print('TP = ',true_positives)
6 print('TN = ',true_negatives)
7 print('FP = ',false_positives)
8 print('FN = ',false_positives)
9 precision = float(true_positives) / (true_positives + false_positives)
10 recall = float(true_positives)/(true_positives + false_negatives)
11 print('recall = ',recall)
12 print('precision = ',precision)

TP = 7
TN = 95
FP = 10
FN = 10
recall = 0.15217391304347827
precision = 0.4117647058823529
```

Для логистической регрессии:

```
1 true_positives = lr_predictions[(lr_predictions.Class_ID == 2) & (lr_predictions.prediction == 2)].count()
2 true_negatives = lr_predictions[(lr_predictions.Class_ID == 1) & (lr_predictions.prediction == 1)].count()
3 false_positives = lr_predictions[(lr_predictions.Class_ID == 1) & (lr_predictions.prediction == 2)].count()
4 false_negatives = lr_predictions[(lr_predictions.Class_ID == 2) & (lr_predictions.prediction == 1)].count()
5 print('TP = ',true_positives)
6 print('TN = ',true_negatives)
7 print('FP = ',false_positives)
8 print('FN = ',false_positives)
9 precision = float(true_positives) / (true_positives + false_positives)
10 recall = float(true_positives)/(true_positives + false_negatives)
11 print('recall = ',recall)
12 print('precision = ',precision)
```

```
TP = 12
TN = 96
FP = 9
FN = 9
recall = 0.2608695652173913
precision = 0.5714285714285714
```

## Заключение

В данной работе мы познакомились с библиотекой MLlib фреймворка Apache для решения задачи классификации пациентов с заболеваниями печени.

Для этого:

- создали объект сессию SparkSession
- загрузили набор данных, содержащий результаты лабораторных исследований пациентов Indian\_Liver\_Patient\_Dataset\_ILPDh.csv, у части из которых диагностировано заболевание печени
- произвели обзорный исследовательский анализ данных
- разделили набор данных на обучающую и тестовую выборку
- обучили две модели машинного обучения для классификации пациентом: random forest и logistic regression
- оценили работу моделей на тестовой выборке при помощи матрицы неточностей, критериев полноты Recall и точности Precision.

Как итог, логистическая регрессия показала лучшие результаты, чем случайный лес.