

Centro de Informática CIn - UFPE

Curso: PD em Larga Escala

Projeto Final

O Random Forest é um método de aprendizado ensemble que pode ser utilizado tanto para regressão como para classificação, ele constrói coleções de árvores de decisão no processo de aprendizado de forma a obter melhor desempenho que cada árvore de decisão poderia oferecer individualmente. No caso da classificação, como o problema abordado neste projeto, o resultado do random forest é a classe selecionada pela maioria das árvores de decisão. O método é bastante utilizado devido a sua flexibilidade, que permite trabalhar com problemas de regressão e classificação com desempenho satisfatório, outro ponto é a facilidade para determinar a importância das features e suas contribuições para o modelo. Apesar de muitas vantagens, o método geralmente é bem custoso, visto que está construindo muitas árvores de decisão por trás e isso pode ser problemático em conjuntos de dados maiores.

Tarefa 1 - Processamento ETL

Executando os passos descritos, você terá no HDFS dados no formato:

- ★ **label** - ao, br, pt, mz, mo, gw (no conjunto reduzido há apenas seis países)
- ★ **features** - vetor esparsa com a representação do texto de cada página

Passo-a-passo da Tarefa 1

- ★ Baixe o arquivo pt7-raw.zip
- ★ Copie a pasta descompactada para user_data/pt7-raw
- ★ Copie os arquivos do PT7 para o HDFS
 - `docker exec -it master /bin/bash`
 - `hadoop fs -mkdir -p /bigdata/`
 - `hadoop fs -put /user_data/pt7-raw hdfs://master:8020/bigdata/`
- ★ Processe o job labels-pt7-raw.scala
 - `spark-shell --master spark://master:7077 -i /user_data/labels-pt7-raw.scala`

Como resultado, será obtido um dataframe conforme imagens a seguir.

```
C:\> Administrador: Prompt de Comando - docker exec -it master /bin/bash

scala> tldDF.show

+-----+-----+-----+
|label|          url|      text64byte|
+-----+-----+-----+
|.ao|http://mercado.co...|DQpMaWtlcw0KU3Vic...|
|.ao|http://mercado.co...|DQpMaWtlcw0KU3Vic...|
|.ao|http://mercado.co...|DQpMaWtlcw0KU3Vic...|
|.ao|http://mercado.co...|DQpMaWtlcw0KU3Vic...|
|.ao|http://mercado.co...|DQpMaWtlcw0KU3Vic...|
|.ao|http://mercado.co...|DQpMaWtlcw0KU3Vic...|
|.ao|http://mercado.co...|DQpMaWtlcw0KU3Vic...|
|.ao|http://mercado.co...|DQpMaWtlcw0KU3Vic...|
|.ao|http://mercado.co...|DQpMaWtlcw0KU3Vic...|
|.ao|http://mercado.co...|DQpMaWtlcw0KU3Vic...|
|.ao|http://mercado.co...|DQpMaWtlcw0KU3Vic...|
|.ao|http://mercado.co...|DQpMaWtlcw0KU3Vic...|
|.ao|http://mercado.co...|DQpMaWtlcw0KU3Vic...|
|.ao|http://mercado.co...|DQpMaWtlcw0KU3Vic...|
|.ao|http://mercado.co...|DQpMaWtlcw0KU3Vic...|
|.ao|http://mercado.co...|DQpMaWtlcw0KU3Vic...|
|.ao|http://mercado.co...|DQpMaWtlcw0KU3Vic...|
|.ao|http://mercado.co...|DQpMaWtlcw0KU3Vic...|
|.ao|http://mercado.co...|DQpMaWtlcw0KU3Vic...|
|.ao|http://mercado.co...|DQpMaWtlcw0KU3Vic...|
+-----+-----+-----+
only showing top 20 rows
```

```
C:\> Administrador: Prompt de Comando - docker exec -it master /bin/bash

scala> tldDF.groupBy("label").count().show()

+-----+-----+
|label|count|
+-----+-----+
|.ao| 2122|
|.br| 7053|
|.mz| 2820|
|.pt| 3054|
|.gw| 1603|
|.mo|  362|
+-----+-----+
```

Processe o job etl-pt7.scala

★ *spark-shell --master spark://master:7077 -i /user_data/etl-pt7.scala*

Como resultado, será obtido um dataframe conforme imagens a seguir. Neste ponto, o dataframe multilabel com os vetores esparsos será gravado no seu HDFS no caminho `hdfs://master:8020/bigdata/pt7-hash.parquet`

```
Administrator: Prompt de Comando - docker exec -it master /bin/bash
scala> the_df.show()
22/07/08 13:20:13 WARN DAGScheduler: Broadcasting large task binary with size 4.0 MiB
+-----+-----+
|label|          features|
+-----+-----+
|.mz|(262144,[69,452,1...|
|.mz|(262144,[69,1004,...|
|.mz|(262144,[226,3170...|
|.mz|(262144,[1083,186...|
|.mz|(262144,[69,1004,...|
|.mz|(262144,[69,72,66...|
|.mz|(262144,[472,1004...|
|.mz|(262144,[188,452,...|
|.mz|(262144,[3704,376...|
|.mz|(262144,[69,1004,...|
|.mz|(262144,[69,452,1...|
|.mz|(262144,[3542,370...|
|.mz|(262144,[427,1252...|
|.mz|(262144,[452,1840...|
|.mz|(262144,[427,2209...|
|.mz|(262144,[2209,280...|
|.mz|(262144,[2778,370...|
|.mz|(262144,[202,827,...|
|.mz|(262144,[69,427,4...|
|.mz|(262144,[69,452,3...|
+-----+-----+
only showing top 20 rows
```

Tarefa 2 - Treinar e Testar Um Modelo Supervisionado

Passo 1:

- ★ colocar o arquivo “script.py” dentro da pasta “user_data”

Passo 2:

- ★ criar uma pasta “projeto” dentro da pasta “user_data”

Passo 3: Instalar a biblioteca numpy no master e nos 3 workers

- ★ docker exec -it master pip install numpy
- ★ docker exec -it worker-1 pip install numpy
- ★ docker exec -it worker-2 pip install numpy
- ★ docker exec -it worker-3 pip install numpy

Passo 4: Rodar o arquivo script.py (Salva métricas em /user_data/projeto/metricas.txt)

- ★ docker exec -it master /bin/bash
- ★ cd user_data
- ★ pyspark --master spark://master:7077 < script.py

Passo 5 (opcional): Exportar o modelo para o disco local (Salva o modelo em /user_data/projeto/modelo_rf)

- ★ hdfs dfs -copyToLocal hdfs://master:8020/bigdata/modelo_rf/user_data/projeto/

