



## Residência em Engenharia e Ciência de Dados

# Processamento de Dados em Larga Escala

### **PROJETO**

[ETL / treinamento / teste sobre Corpus PT7 multiclasse]

O PT7 Web (https://ieee-dataport.org/open-access/pt7-web-annotatedportuguese-language-corpus) é um Corpus anotado em língua portuguesa construído a partir de amostras coletadas de setembro de 2018 a março de 2020 de sete países de língua portuguesa: Angola, Brasil, Portugal, Cabo Verde, Guiné-Bissau, Macau e Moçambique. Os registros foram filtrados do Common Crawl — um conjunto de dados em escala de petabytes de domínio público de páginas da Web em vários idiomas, misturados em instantâneos temporais da Web, disponíveis mensalmente [1]. As páginas brasileiras foram rotuladas como classe positiva e as demais como classe negativa (português não brasileiro). O conjunto de dados totalizou 249,74 GB de texto HTML bruto relacionado a 16.346.693 páginas da web exclusivas. Os dados foram pré-processados para produzir vetores de distribuição de palavras de alta dimensionalidade (2 elevado a 18 = 262.144 características) como entrada para as fases de treinamento e teste. Uma demonstração do uso desses dados pode ser verificada em um projeto fracionário de dois níveis para investigar o desempenho do cluster no Spark [2].

[1] G.Wenzek, M.A.Lachaux, A.Conneau, V.Chaudhary, F.Guzman, A.Joulin, E.Grave, arXiv preprint arXiv:1911.00359 (2019).

[2] Rodrigues, J.; Vasconcelos, G.; Maciel, P. Time and cost prediction models for language classification over a large corpus on spark.

In: 2020 IEEE Symposium Series on Computational Intelligence (SSCI). [S.l.: s.n.], 2020. p. to appear.

Será utilizado um extrato reduzido do PT7 Web, equivalente 17014 páginas ~ 0.1% do Corpus original. Foram disponibilizados cinco arquivos (pt7-raw.zip), separados pelo domínio de nível superior de cada país (.br, .pt, .mo, .gw, .mz, .ao e .pt).

```
CC-MAIN-2018-39_ao
CC-MAIN-2018-39_br
CC-MAIN-2018-39_gw
CC-MAIN-2018-39_mo
CC-MAIN-2018-39_mz
CC-MAIN-2018-39_pt
```

Os dados se encontram rotulados como 1:pt\_BR e 0:pt\_OTHERS e estão disponível no formato a seguir, onde:

- label rótulo
- · url endereço original completo da página
- · digest uma função de bash do conteúdo da página
- raw os dados brutos do texto da página após limpeza de tags HTML

raw	digest	url	label
\r\nEmpresas\r\nP	JJYINTQR7DRFBMWAI	http://1-1.pt/300	0
\r\nEmpresas\r\nP	AUP6PWDHXTGEHMW4M	http://1-1.pt/300	0
\r\nEmpresas\r\nP	TMHH62CWHZRPRRQK3	http://1-1.pt/300	0
\r\nEmpresas\r\nP	BZXW5XKSAQUQUF3FM	http://1-1.pt/300	0
\r\nEmpresas\r\nP	TR4JYZTG563VKFQBY	http://1-1.pt/300	0
<pre>\r\nEmpresas\r\nP</pre>	6XKFIZTSVDN4BYZTE	http://1-1.pt/300	0
\r\nEmpresas\r\nP	3Y3MQ7H2ZF4B4Z7ST	http://1-1.pt/300	0
\r\n100-DJ - ao r	BPW5IV333KOAYNQOP	http://100-dj.pt/	0
\r\nLoading\r\	32I2HPTFDWTJ5P677	http://1000olhos	0
\r\n \n100PAVOR\r	7LEAGIOKBGPYZRGLT	http://100pavor.pt/	0
\r\nInício\n(curr	WAJQ27Y2I4SJ7SLN3	http://100solucoe	0
<pre>»\n \nGaleria\nCa</pre>	MIMURTS2WQLS3NLSO	http://100trilhos	0
»\n \nEscalada\r\	WSB2HR6434YB3YW3H	http://100trilhos	0
\r\nContacto: +35	BIRLLHEMAPI7C60KY	http://13luas.pt/	0
Porto de Recreio	OE5UZKCTVW3SGHIUM	http://13yachtbro	0
Porto de Recreio	R2CWSY7HDB5TEH4FV	http://13yachtbro	0
Porto de Recreio	NQ7YRBNLYIAFASEBN	http://13yachtbro	0
	RDZM2NLWH3P4Q7EZB		
	HHVCWGRAFQYXDLS4X	·	
\nCranchi Atlanti	5FWCBW3H6DRTK5EIV	http://13yachtbro	0

A tarefa preliminar do projeto consiste em realizar o processo de ETL sobre os dados brutos, transformando o conteúdo de cada página web em um vetor esparso de características no formato exigido pelo Spark. A base deve separar os dados em novos rótulos, de acordo com cada país, formando uma base rotulada multiclasse.

#### Tarefa 1 << disponibilizado >> : 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 esparso 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
  - \$ hadoop fs -put /user\_data/pt7-raw hdfs://master:8020/bigdata/
- Processe o jobs labels-pt7-raw.scala
  - \$ spark-shell --master spark://master:7077 -i /user\_data/labels-pt7-raw.scala

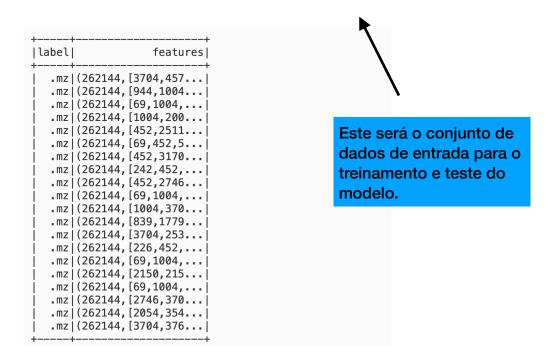
Como resultado, será obtido um dataframe conforme imagens a seguir.

label	l  url	text64byte		'
. ao   . ao	http://mercado.co	DQpMaWtlcw0KU3Vic	.pt  30   .ao  21   .gw  10   .mz  28	unt   + 053   054   122   503   320   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 <a href="https://master:8020/bigdata/pt7-hash.parquet">https://master:8020/bigdata/pt7-hash.parquet</a>



#### Tarefa 2 << a implementar >>: treinar e testar um modelo supervisionado.

Investigar as possibilidades de modelos supervisionados em <a href="https://spark.apache.org/docs/latest/ml-classification-regression.html#classification">https://spark.apache.org/docs/latest/ml-classification-regression.html#classification</a>. Implementar em Scala, Python, R ou Java o processo de treinamento e teste do modelo escolhido.

Dica: Para construção do modelo, baseie-se no slide 12-spark-mllib-kddcup.pdf. E, para ler o dataframe com os vetores esparsos do HDFS, utilize o comando (em Scala) a seguir.

```
val df = {
    spark.read
    .format("parquet")
    .load("hdfs://master:8020/bigdata/pt7-hash.parquet")
}
```

Não se esqueça de salvar o seu modelo treinado no HDFS para posterior uso.

#### Artefatos da entrega:

- Descrição breve do modelo supervisionado escolhido
- Código-fonte
- Métricas Acurácia, PR e F-Measure do modelo sobre os dados de teste
- Arquivo README.txt com instruções para execução no cluster (para correção)

Dúvidas serão tiradas por e-mail e nos momentos de acompanhamento de projeto definidos no cronograma.

Bons estudos.