**Especialização Lato Sensu em Ciência de Dados e Analytics**

# Soluções em Processamento para Big Data

## { Prática Google Data Proc }

Prof. Jairson Rodrigues
jairson.rodrigues@univasf.edu.br

# { google data proc }

## AGENDA

Cluster Google

# { roteiro }

- Cadastro no Google DataProc
- Construir um cluster de quatro máquinas
- Acessar o Name Node
- Executar algoritmos de ML
    - ETL
    - Regressão Logística
    - Árvores de Decisão
    - Kmeans

# { google data proc }

## CLOUD DATAPROC

A faster, easier, more cost-effective way to run Spark and Hadoop

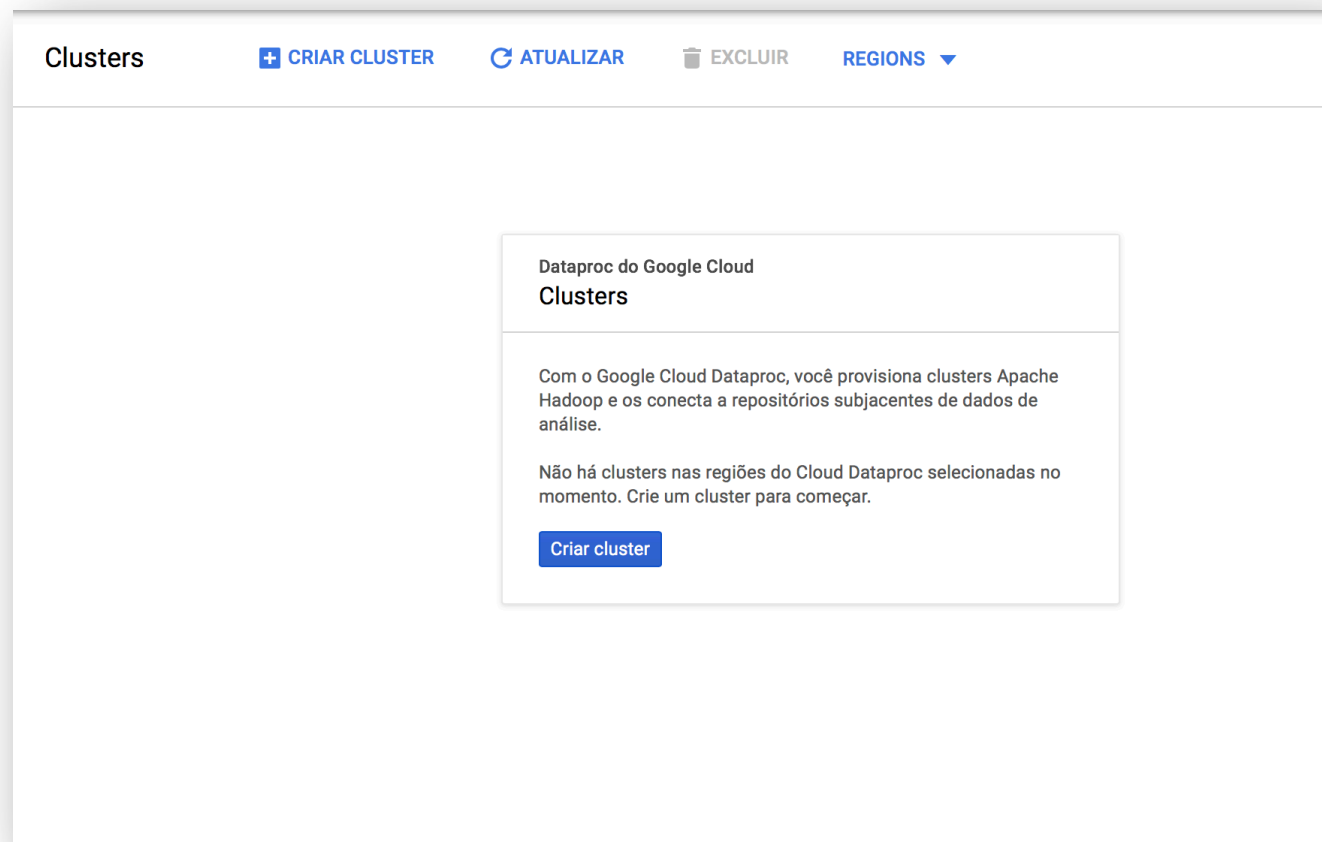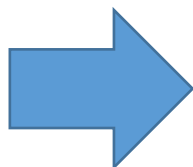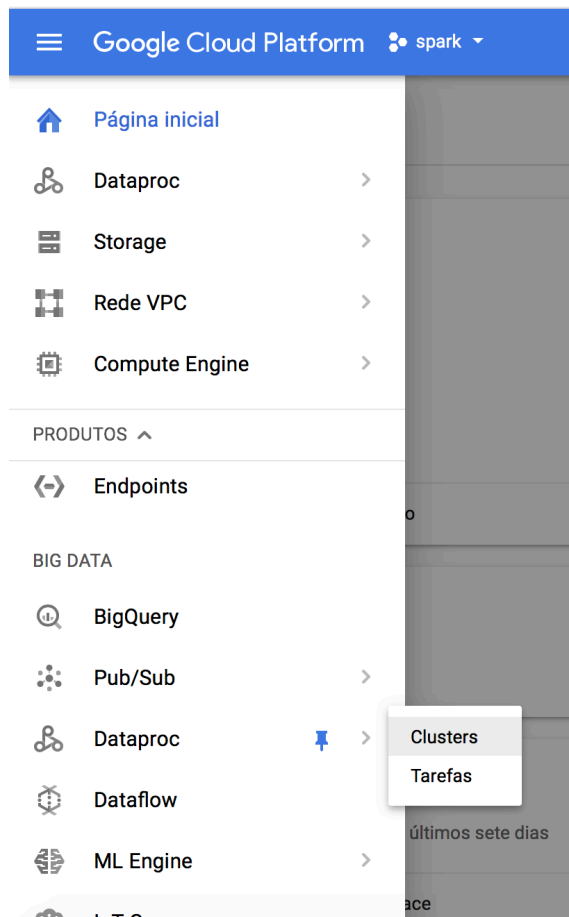[ VIEW CLOUD DATAPROC DOCS ]  [ VIEW MY CONSOLE ]

## Cloud-native Hadoop & Spark

Cloud Dataproc is a fast, easy-to-use, fully-managed cloud service for running **Apache Spark** and **Apache Hadoop** clusters in a simpler, more cost-efficient way. Operations that used to take hours or days take seconds or minutes instead, and you pay only for the resources you use (with per-second billing). Cloud Dataproc also easily integrates with other Google Cloud Platform (GCP) services, giving you a powerful and complete platform for data processing, analytics and machine learning.
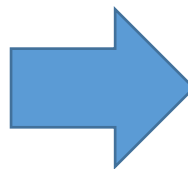
# { google data proc }

# { google data proc }

Tipos de Máquina Virtual



n1-standard-1 (1 vCPU, 3.75 GB de memória)

n1-standard-2 (2 vCPU, 7.50 GB de memória)

n1-standard-4 (4 vCPU, 15.0 GB de memória)

n1-standard-8 (8 vCPU, 30.0 GB de memória)

n1-standard-16 (16 vCPU, 60.0 GB de memória)

n1-standard-32 (32 vCPU, 120 GB de memória)

n1-standard-64 (64 vCPU, 240 GB de memória)

n1-highcpu-4 (4 vCPU, 3.60 GB de memória)

n1-highcpu-8 (8 vCPU, 7.20 GB de memória)

n1-highcpu-16 (16 vCPU, 14.4 GB de memória)

n1-highcpu-32 (32 vCPU, 28.8 GB de memória)

n1-highcpu-64 (64 vCPU, 57.6 GB de memória)

n1-highmem-2 (2 vCPU, 13.0 GB de memória)

n1-highmem-4 (4 vCPU, 26.0 GB de memória)

n1-highmem-8 (8 vCPU, 52.0 GB de memória)

n1-highmem-16 (16 vCPU, 104 GB de memória)

n1-highmem-32 (32 vCPU, 208 GB de memória)

n1-highmem-64 (64 vCPU, 416 GB de memória)

# { configuração de nós }

- Master: Google VM Machine
  - n1-highmem-2
  - CPU: 2
  - RAM: 13 Gb RAM
  - Disco local: 500 Gb
  - # instâncias: 1
- Slave: Google VM Machine
  - n1-highmem-2
  - CPU: 2
  - RAM: 13 Gb RAM
  - Disco local: 500 Gb
  - Disco SSD 375 Gb
  - # instâncias: 3

| Nome | Papel |
|------|-------|
| ✅ cluster-upe-m | Principal |
| ✅ cluster-upe-w-0 | Trabalho |
| ✅ cluster-upe-w-1 | Trabalho |
| ✅ cluster-upe-w-2 | Trabalho |

# { google data proc – master }

# { google data proc – slaves }

# { detalhes do cluster }

# { conexão SSH }

# { configurações iniciais }

- wget https://www.dropbox.com/s/r5xg2hi28g4s51f/kddcup_2.11-1.0.jar?dl=0
- wget https://github.com/SparkTC/spark-bench/releases/download/v55/spark-bench_2.1.1_0.2.2-RELEASE_55.tgz
- wget https://www.dropbox.com/s/0hi816m22uia2cw/spark-bench-env.sh?dl=0
- wget https://www.dropbox.com/s/98udl9vmi1fayzq/genkmeans.conf?dl=0
- wget https://www.dropbox.com/s/mwnad29s3zdtutx/kmeans.conf?dl=0
- wget http://kdd.ics.uci.edu/databases/kddcup99/kddcup.data.gz
- mv ~/kddcup_2.11-1.0.jar\?dl\=0 kddcup_2.11-1.0.jar
- mv ~/kmeans.conf?dl=0 ~/kmeans.conf
- mv ~/genkmeans.conf?dl=0 ~/genkmeans.conf
- gunzip kddcup.data.gz
- tar -xvf ~/spark-bench_2.1.1_0.2.2-RELEASE_55.tgz
- mv ~/spark-bench_2.1.1_0.2.2-RELEASE ~/spark-bench
- mv ~/spark-bench-env.sh\?dl\=0 ~/spark-bench/bin/spark-bench-env.sh
- mv kddcup.data kddcup.data.10
- hadoop fs -mkdir -p /kddcup/input/
- hadoop fs -copyFromLocal kddcup.data.10 /kddcup/input/

wget https://www.dropbox.com/s/hnfuinu3oose116/init.sh?dl=0
mv init.sh\?dl\=0 init.sh
chmod a+x init.sh; ./init.sh

# { geração de dados - kmeans -> 25 gb }

- ./spark-bench/bin/spark-bench.sh ~/genkmeans.conf



Proxy MAC OS: https://justpaste.it/1e30d

# { execução kmeans -> 8 gb }

- ./spark-bench/bin/spark-bench.sh ~/kmeans.conf



- 40 milhões de amostras
- 24 dimensões R$^{24}$
- # clusters: 5

# { envio de jobs / ETL / LogReg / Naïve Bayes }

- spark-submit --class "KDDCupETL" --master yarn --conf spark.serializer=org.apache.spark.serializer.KryoSerializer kddcup_2.11-1.0.jar

- spark-submit --class "KDDCupRL" --master yarn --conf spark.serializer=org.apache.spark.serializer.KryoSerializer kddcup_2.11-1.0.jar

- spark-submit --class "KDDCupNaiveBayes" --master yarn --conf spark.serializer=org.apache.spark.serializer.KryoSerializer kdkddcup_2.11-1.0.jar

# { envio de jobs
# - interface web }

- KDDCupETL
- KDDCupRL
- KDDCupNaiveBayes

← **Enviar um job**

**Região** ⓘ

us-central1 ▾

**Cluster**

cluster-upe ▾

**Tipo de tarefa**

Spark ▾

**Classe principal ou jar** ⓘ

KDDCupETL ▾

**Argumentos** (Opcional) ⓘ

Pressione <Retornar> para adicionar mais argumentos

**Arquivos jar** (Opcional) ⓘ

hdfs://user/kddcup_2.11-1.0.jar ✕

Insira o caminho do arquivo, por exemplo, hdfs://exemplo/exemplo.jar

**Propriedades** (Opcional) ⓘ

executor-memory | 4G ✕

spark.serializer | che.spark.serializer.KryoSerializer ✕

**+ Adicionar item**

# { envio de jobs - interface web }