

Trabalho de final de módulo – Mineração de Texto com ELK - INFNET

Professor: Felipe Fink Grael

Aluno: Marcelo Lopes da Silva

Dezembro 2020

Instalação dos pré requisitos

Para a presente solução serão utilizados o Elasticsearch e o Kibana executados via docker. O acesso ao elasticsearch será feito utilizando jupyter notebook. Outras dependências podem ser instaladas utilizando o gerenciador de pacotes pipenv. Os exercícios foram realizados em máquina com sistema operacional Ubuntu versão 20

1 – Instalar o docker/docker compose

<https://www.digitalocean.com/community/tutorials/how-to-install-and-use-docker-on-ubuntu-20-04-pt>

2 – Instalar o pipenv

Em um terminal digite o comando pip install pipenv

3 – Instalar as dependências utilizadas no exercício

No diretório onde serão executados os notebook digite o comando pipenv install. Note que os arquivo Pipfile e Pipfile.loc tem que estar no mesmo diretório.

A execução dos notebooks tem que ser feita no ambiente virtual criado pelo pipenv em que as dependências foram instaladas. Para isso no terminal digite o comando pipenv shell.

Execução do Elastic e Kibana

Como informado anteriormente a execução é feita utilizando docker. Um arquivo docker compose é fornecido junto a solução.

docker-compose.yml

execute

```
docker-compose.yml x
docker-compose.yml
1  version: '3'
2  services:
3    elasticsearch:
4      container_name: elasticsearch
5      image: docker.elastic.co/elasticsearch/elasticsearch:7.9.2
6      volumes:
7        - ./elastic:/usr/share/elasticsearch/data
8      environment:
9        - node.name=elastic_node01
10       - discovery.type=single-node
11      ports:
12        - 9200:9200
13      networks:
14        - elk
15    kibana:
16      container_name: kibana
17      image: docker.elastic.co/kibana/kibana:7.9.2
18      volumes:
19        - ./kibana/config:/usr/share/kibana/config:ro
20      ports:
21        - 5601:5601
22      networks:
23        - elk
24      depends_on:
25        - elasticsearch
26      links:
27        - elasticsearch
28    # logstash:
29    #   image: docker.elastic.co/logstash/logstash:7.9.2
30    #   volumes:
31    #     - ./logstash/config:/config-dir
32    #     - ./logstash/data:/logstash
33    #   links:
34    #     - elasticsearch
35    #   command: logstash -f /config-dir/logstash.conf
36    #   networks:
37    #     - elk
38    #   depends_on:
39    #     - elasticsearch
40    #   links:
41    #     - elasticsearch
42
43  networks:
44    elk:
45      driver: bridge
```

comando docker-compose up para iniciar os serviços

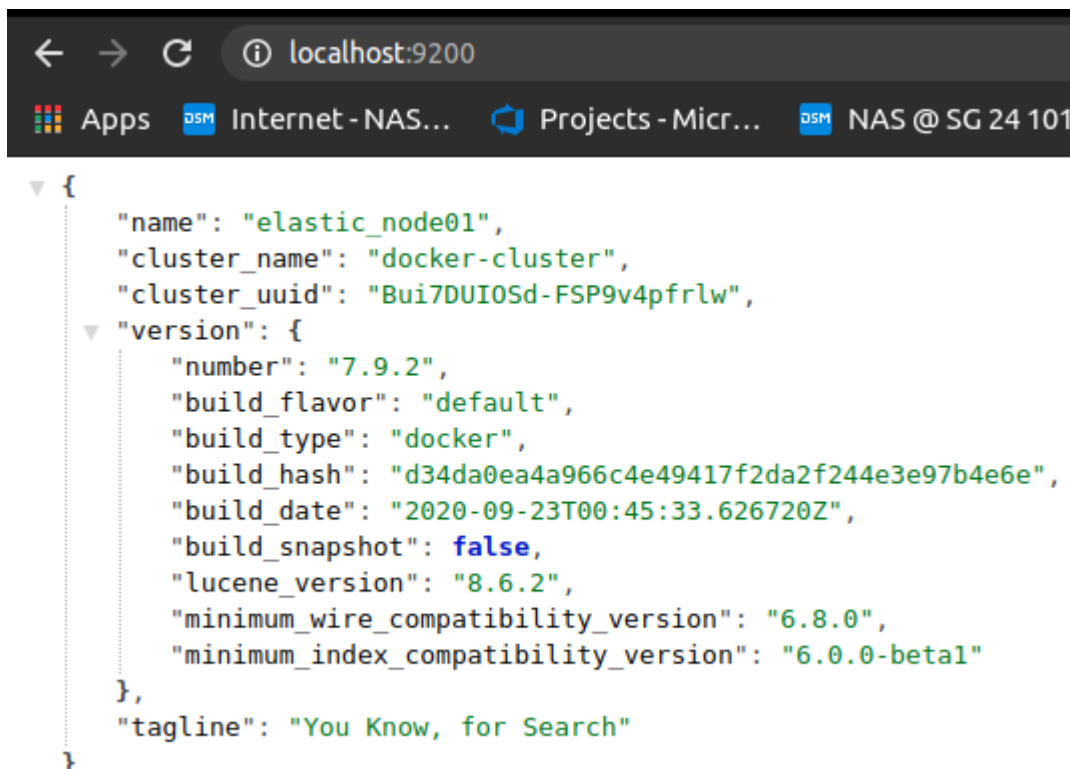
```

marcelo@BD-DESKTOP-UBUNTU:~$ cd Documentos/Infnet/mineracao_elk/
marcelo@BD-DESKTOP-UBUNTU:~/Documentos/Infnet/mineracao_elk$ docker-compose up
Starting elasticsearch ... done
Starting kibana ... done
Attaching to elasticsearch, kibana
elasticsearch | {"type": "server", "timestamp": "2020-12-08T10:07:19,454Z", "level": "INFO", "component": "o.e.n.Node",
", "cluster.name": "docker-cluster", "node.name": "elastic_node01", "message": "version[7.9.2], pid[6], build[default/do
cker/d34da0ea4a966c4e49417f2da2f244e3e97b4e6e/2020-09-23T00:45:33.626720Z], OS[Linux/5.4.0-56-generic/amd64], JVM[AdoptO
penJDK/OpenJDK 64-Bit Server VM/15/15+36]} }
elasticsearch | {"type": "server", "timestamp": "2020-12-08T10:07:19,484Z", "level": "INFO", "component": "o.e.n.Node",
", "cluster.name": "docker-cluster", "node.name": "elastic_node01", "message": "JVM home [/usr/share/elasticsearch/jdk]"
}
elasticsearch | {"type": "server", "timestamp": "2020-12-08T10:07:19,485Z", "level": "INFO", "component": "o.e.n.Node",
", "cluster.name": "docker-cluster", "node.name": "elastic_node01", "message": "JVM arguments [-Xshare:auto, -Des.network
kaddress.cache.ttl=60, -Des.networkaddress.cache.negative.ttl=10, -XX:+AlwaysPreTouch, -Xss1m, -Djava.awt.headless=true,
-Dfile.encoding=UTF-8, -Djna.nosys=true, -XX:-OmitStackTraceInFastThrow, -XX:+ShowCodeDetailsInExceptionMessages, -Dio
.netty.noUnsafe=true, -Dio.netty.noKeySetOptimization=true, -Dio.netty.recycler.maxCapacityPerThread=0, -Dio.netty.alloca
tor.numDirectArenas=0, -Dlog4j.shutdownHookEnabled=false, -Dlog4j2.disable.jmx=true, -Djava.locale.providers=SPI,COMPAT,
-Xms1g, -Xmx1g, -XX:+UseG1GC, -XX:G1ReservePercent=25, -XX:InitiatingHeapOccupancyPercent=30, -Djava.io.tmpdir=/tmp/ela
sticsearch-8339672847912995948, -XX:+HeapDumpOnOutOfMemoryError, -XX:HeapDumpPath=data, -XX:ErrorFile=logs/hs_err_pid%p.
log, -Xlog:gc*,gc+age=trace,safepoint:file=logs/gc.log:utctime,pid,tags:filecount=32,filesize=64m, -Des.cgroups.hierarch
y.override=/, -XX:MaxDirectMemorySize=536870912, -Des.path.home=/usr/share/elasticsearch, -Des.path.conf=/usr/share/elas
ticsearch/config, -Des.distribution.flavor=default, -Des.distribution.type=docker, -Des.bundled_jdk=true]} }
kibana | {"type": "log", "@timestamp": "2020-12-08T10:07:24Z", "tags": ["warning", "plugins-discovery"], "pid": 6, "mes

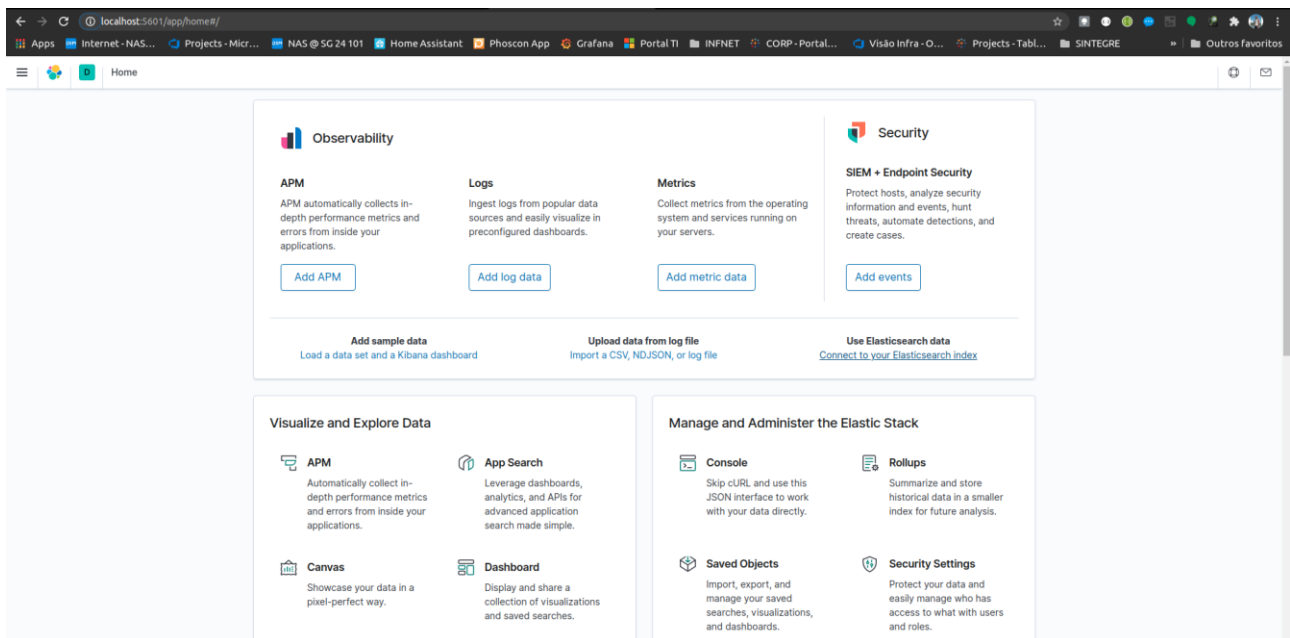
```

Após algum tempo verifique se o elasticsearch está respondendo, abrindo um browser e apontando para localhost:9200

Depois



verifique se o kibana está respondendo, apontando o browser para localhost:5601



Notebooks

Os notebooks de carga de informações e consultas podem ser observados nos arquivos: 1_carga_informacoes_dou_es.ipynb e 2_consultas_informacoes_dou_es.ipynb

Dashboard

Dashboard criado no kibana com visualizações dos top 10 registradores de informações, distribuição de documentos por tipo, frequência de termos e nuvem de palavras com os 70 termos mais relevantes.

