# Запуск приложения Apache Spark на облачном сервисе Amazon EC2

Мурашкин Вячеслав
2017

https://github.com/a4tunado/lectures-hse-spark/tree/master/aws

# 1. Создание ключа для авторизации

# 1. Создание ключа для авторизации

# 1. Создание ключа для авторизации

# 2. Создание хранилища для загрузки данных

# 2. Создание хранилища для загрузки данных

# 2. Создание хранилища для загрузки данных

# 2. Создание хранилища для загрузки данных

# 3. Создание EMR кластера

# 3. Создание EMR кластера

# 3. Создание EMR кластера

# 3. Создание EMR кластера

# 3. Создание EMR кластера

# 3. Создание EMR кластера

s3://aws-bigdata-blog/artifacts/aws-blog-emr-jupyter/install-jupyter-emr5.sh

# 3. Создание EMR кластера

- Укажите имя ключа, созданного на первом шаге

# 4. Настройка ssh прокси для подключения к EMR

# 4. Настройка ssh прокси для подключения к EMR

- Перед запуском необходимо обновить настройки *.pem файла, выполнив команду: chmod 400 <key-file>.pem

## Setup Web Connection

Hadoop, Ganglia, and other applications publish user interfaces as web sites hosted on the maste only available on the master node's local web server.

To reach the web interfaces, you must establish an SSH tunnel with the master node using either an SSH tunnel using dynamic port forwarding, you must also configure a proxy server to view the

**Step 1: Open an SSH Tunnel to the Amazon EMR Master Node -** [Learn more]

| Windows | Mac / Linux |

1. Open a terminal window. On Mac OS X, choose Applications > Utilities > Terminal. On oth at Applications > Accessories > Terminal.
2. To establish an SSH tunnel with the master node using dynamic port forwarding, type the the location and filename of the private key file (.pem) used to launch the cluster.

   ```
   ssh -i ~/hse.pem -ND 8157 hadoop@ec2-54-234-247-21.compute-1.amazonaws.com
   ```

   Note: Port 8157 used in the command is a randomly selected, unused local port.
3. Type yes to dismiss the security warning.

# 4. Настройка ssh прокси для подключения к EMR

Branch: master ▾    lectures-hse-spark / aws / chromeproxy.sh

a4tunado chromeproxy.sh

1 contributor

5 lines (3 sloc) | 163 Bytes

```bash
1    #!/bin/bash
2
3    "/Applications/Google Chrome.app/Contents/MacOS/Google Chrome" \
4      --user-data-dir="$HOME/chrome-with-proxy" --proxy-server="socks5://localhost:8157"
```

# 5. Jupyter notebook