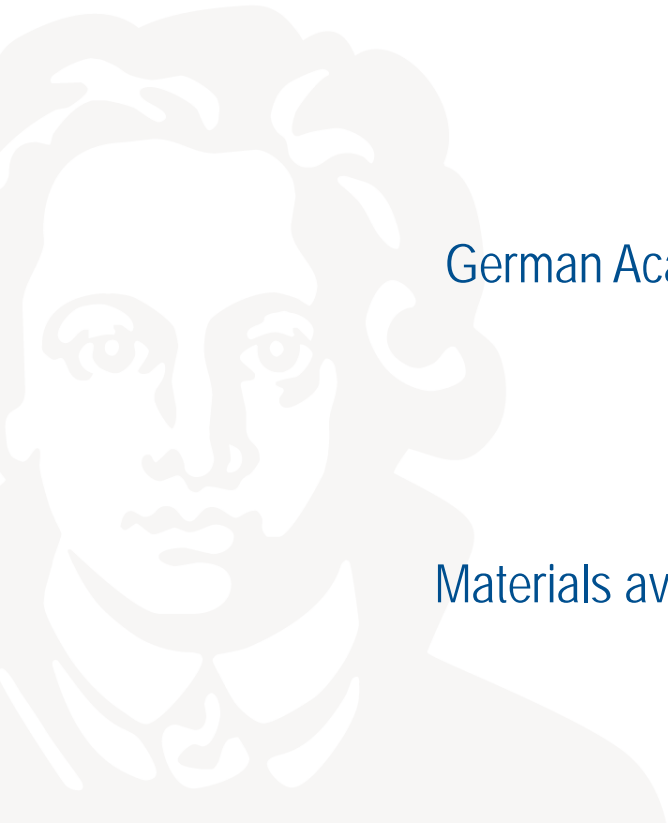# Managing and Analyzing Big Data in the Cloud

**Klaus M. Miller**
Goethe University Frankfurt

German Academic Association for Business Research
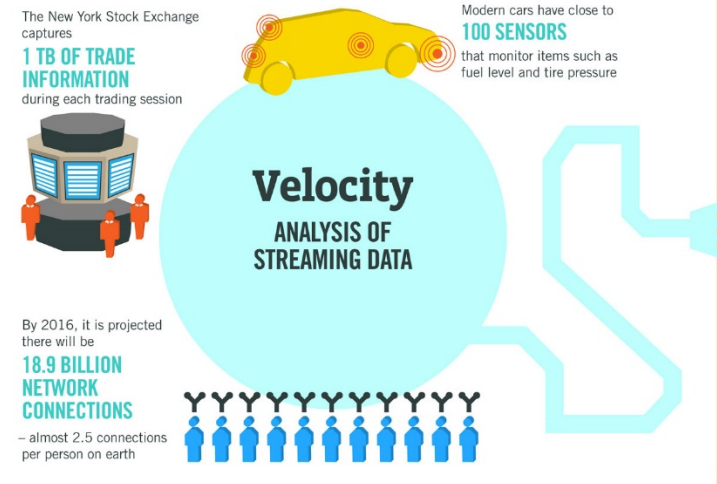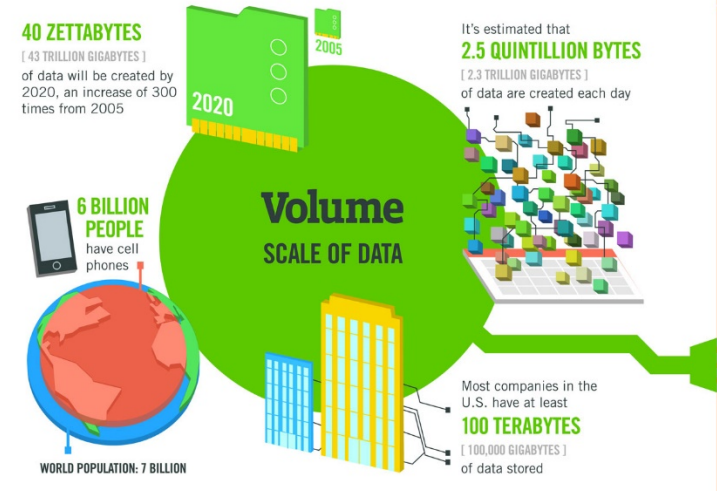Annual Meeting

Frankfurt, March 18th, 2020

Materials available at: https://github.com/stm/vhb_2020

# Description of Problem

# Definition of Big Data

**40 ZETTABYTES**
[ 43 TRILLION GIGABYTES ]
of data will be created by 2020, an increase of 300 times from 2005

**6 BILLION PEOPLE**
have cell phones

WORLD POPULATION: 7 BILLION

It's estimated that
**2.5 QUINTILLION BYTES**
[ 2.3 TRILLION GIGABYTES ]
of data are created each day

Most companies in the U.S. have at least
**100 TERABYTES**
[ 100,000 GIGABYTES ]
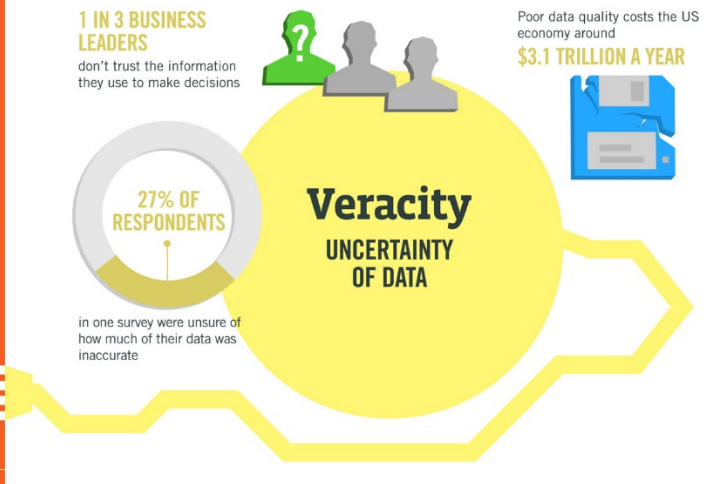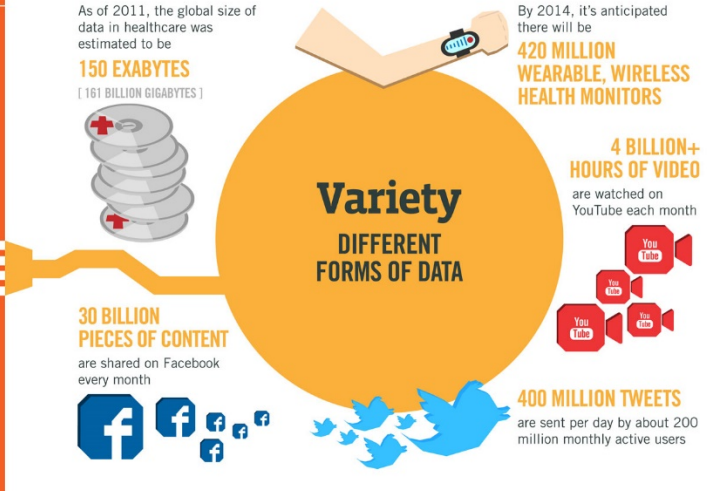of data stored

**Volume**
SCALE OF DATA

## The FOUR V's of Big Data

From traffic patterns and music downloads to web history and medical records, data is recorded, stored, and analyzed to enable the technology and services that the world relies on every day. But what exactly is big data, and how can these massive amounts of data be used?

As a leader in the sector, IBM data scientists break big data into four dimensions: **Volume, Velocity, Variety and Veracity**

Depending on the industry and organization, big data encompasses information from multiple internal and external sources such as transactions, social media, enterprise content, sensors and mobile devices. Companies can leverage data to adapt their products and services to better meet customer needs, optimize operations and infrastructure, and find new sources of revenue.

By 2015
**4.4 MILLION IT JOBS**
will be created globally to support big data, with 1.9 million in the United States

As of 2011, the global size of data in healthcare was estimated to be
**150 EXABYTES**
[ 161 BILLION GIGABYTES ]

**30 BILLION PIECES OF CONTENT**
are shared on Facebook every month

**Variety**
DIFFERENT FORMS OF DATA

By 2014, it's anticipated there will be
**420 MILLION WEARABLE, WIRELESS HEALTH MONITORS**

**4 BILLION+ HOURS OF VIDEO**
are watched on YouTube each month

**400 MILLION TWEETS**
are sent per day by about 200 million monthly active users

The New York Stock Exchange captures
**1 TB OF TRADE INFORMATION**
during each trading session

**Velocity**
ANALYSIS OF STREAMING DATA

By 2016, it is projected there will be
**18.9 BILLION NETWORK CONNECTIONS**
– almost 2.5 connections per person on earth

Modern cars have close to
**100 SENSORS**
that monitor items such as fuel level and tire pressure

**1 IN 3 BUSINESS LEADERS**
don't trust the information they use to make decisions

**27% OF RESPONDENTS**
in one survey were unsure of how much of their data was inaccurate

**Veracity**
UNCERTAINTY OF DATA

Poor data quality costs the US economy around
**$3.1 TRILLION A YEAR**

**Sources:** McKinsey Global Institute, Twitter, Cisco, Gartner, EMC, SAS, IBM, MEPTEC, QAS

IBM

# Definition of Big Data

Data is big anytime it makes you feel it is.

# Example of Big Data: Cookie Data Set (I)

**Data Provider: Large European Ad Exchange**
- 84% reach of internet users in relevant market
- Desktop and mobile browsing traffic
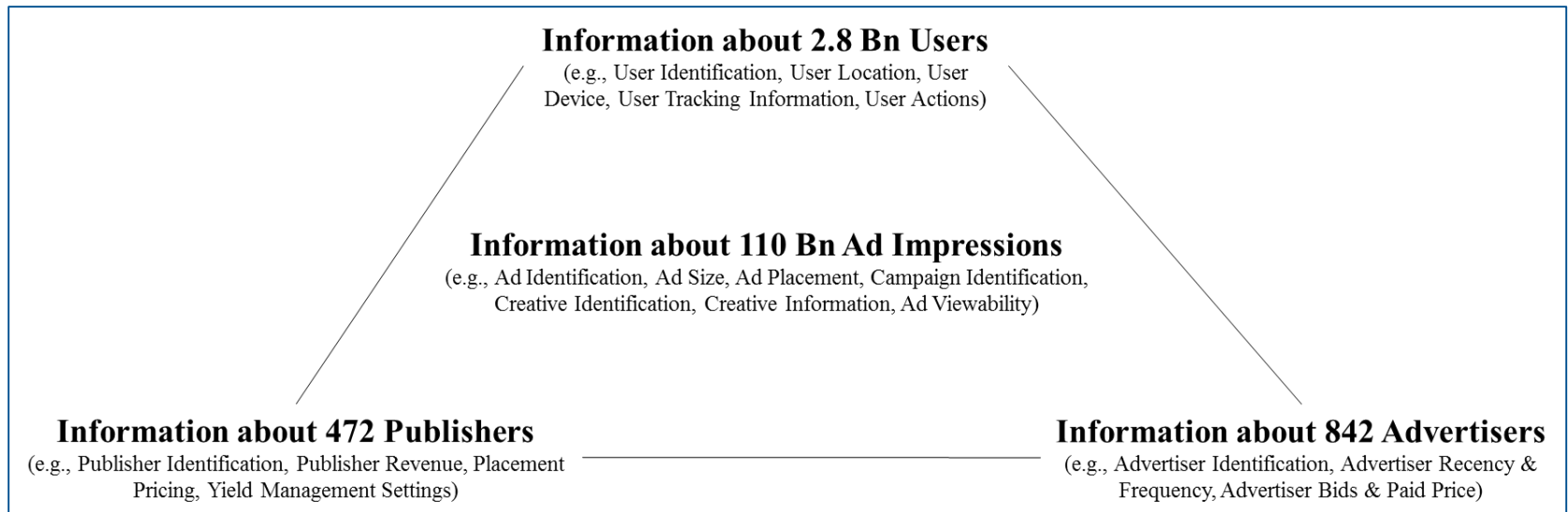
**Observation Period:**
- ~2.5 years

**Dimensions**
- Log-level data set
- ~ 130 columns
- 550 million auctions (= rows) per day
- Total size: 60-65 TB

# Example of Big Data: Cookie Data Set (II)

## Overview of Type of Information Embedded in the Data

**Information about 2.8 Bn Users**
(e.g., User Identification, User Location, User
Device, User Tracking Information, User Actions)

**Information about 110 Bn Ad Impressions**
(e.g., Ad Identification, Ad Size, Ad Placement, Campaign Identification,
Creative Identification, Creative Information, Ad Viewability)

**Information about 472 Publishers**
(e.g., Publisher Identification, Publisher Revenue, Placement
Pricing, Yield Management Settings)

**Information about 842 Advertisers**
(e.g., Advertiser Identification, Advertiser Recency &
Frequency, Advertiser Bids & Paid Price)

# Possibilities to Deal with Large Data Sets in R

1. Allocate more memory („memory.limit")

2. Vectorize (use „apply" family instead of „for loops")

3. Collect garbage („gc")

4. Parallelize („parApply", „doParallel")

5. Use Command Line Interface (e.g., GIT Bash)

6. **Scale to the Cloud**

# Scaling to the Cloud

# Cloud Services

# Creating an AWS Account

- Free Basic Account and Credits
- Choose Region (e.g., EU (Frankfurt))
- Access AWS Services:
  - Data Storage (S3)
  - Elastic Map Reduce (EMR)

- AWS Educate for Usage in Class



https://aws.amazon.com/

# Data Storage: AWS S3

# Data Access

- AWS Console Access (Browser)

- API Access
  - Data sharing with collaborators worldwide
  - Obtaining data from data providers

- Cloud Storage Browsers
  (e.g., S3 Browser, Cyberduck)

# Create EMR Cluster

# Step 1: Software and Steps

Create Cluster - Advanced Options   **Go to quick options**

**Step 1: Software and Steps**

Step 2: Hardware

Step 3: General Cluster Settings

Step 4: Security

## Software Configuration

**Release** | emr-5.23.0 | ▼ | ⓘ

| | | |
|---|---|---|
| ☑ Hadoop 2.8.5 | ☐ Zeppelin 0.8.1 | ☐ Livy 0.5.0 |
| ☐ JupyterHub 0.9.4 | ☐ Tez 0.9.1 | ☐ Flink 1.7.1 |
| ☐ Ganglia 3.7.2 | ☐ HBase 1.4.9 | ☑ Pig 0.17.0 |
| ☑ Hive 2.3.4 | ☐ Presto 0.215 | ☐ ZooKeeper 3.4.13 |
| ☐ MXNet 1.3.1 | ☐ Sqoop 1.4.7 | ☐ Mahout 0.13.0 |
| ☑ Hue 4.3.0 | ☐ Phoenix 4.14.1 | ☐ Oozie 5.1.0 |
| ☑ Spark 2.4.0 | ☐ HCatalog 2.3.4 | ☐ TensorFlow 1.12.0 |

Multi-master support

☐ Enable multi-master support  ⓘ

AWS Glue Data Catalog settings (optional)

☐ Use for Hive table metadata  ⓘ

☐ Use for Spark table metadata  ⓘ

Edit software settings  ⓘ

🔵 Enter configuration  ⚪ Load JSON from S3

```
classification=config-file-name,properties=[myKey1=myValue1,myKey2=myValue2]
```

## Add steps (optional)  ⓘ

**Step type** | Select a step | ▼ | Configure

☐ Auto-terminate cluster after the last step is completed

Cancel    **Next**

# Step 2: Hardware



Create Cluster - Advanced Options   **Go to quick options**

Step 1: Software and Steps

**Step 2: Hardware**

Step 3: General Cluster Settings

Step 4: Security

## Hardware Configuration ⓘ

If you need more than 20 EC2 instances, see this topic 🗗.

Instance group configuration  ● **Uniform instance groups**
*Specify a single instance type and purchasing option for each node type.*

○ **Instance fleets**
*Specify target capacity and how Amazon EMR fulfills it for each node type. Mix instance types and purchasing options. Learn more 🗗*

Network  `vpc-8f2a6fe6 (172.31.0.0/16) (default) ▼`   Create a VPC 🗗 ⓘ

EC2 Subnet  `subnet-8ca58fc6 | Default in eu-central-1c ▼`

Root device EBS volume size  `10` GiB ⓘ

Choose the instance type, number of instances, and a purchasing option. You can choose to use On-Demand Instances, Spot Instances, or both. The instance type and purchasing option apply to all EC2 instances in each instance group, and you can only specify these options for an instance group when you create it. Learn more about instance purchasing options 🗗

| Node type | Instance type | Instance count | Purchasing option | Auto Scaling | |
|---|---|---|---|---|---|
| **Master**<br>Master - 1 ✏ | **m4.large** ✏<br>4 vCore, 8 GiB memory, EBS only storage<br>EBS Storage: 32 GiB ⓘ ✏<br>Add configuration settings ✏ | 1  Instances | ● On-demand ⓘ<br>○ Spot ⓘ<br>`Use on-demand as max price ▼` | Not available for Master | ❓ |
| **Core**<br>Core - 2 ✏ | **m4.large** ✏<br>4 vCore, 8 GiB memory, EBS only storage<br>EBS Storage: 32 GiB ⓘ ✏<br>Add configuration settings ✏ | `2`  Instances | ● On-demand ⓘ<br>○ Spot ⓘ<br>`Use on-demand as max price ▼` | Not enabled ✏ | ❓ |
| **Task** ✖<br>Task - 3 ✏ | **m4.large** ✏<br>4 vCore, 8 GiB memory, EBS only storage<br>EBS Storage: 32 GiB ⓘ ✏<br>Add configuration settings ✏ | `0`  Instances | ● On-demand ⓘ<br>○ Spot ⓘ<br>`Use on-demand as max price ▼` | Not enabled ✏ | ❓ |

# Step 2: Hardware

**Instance types**

| | Instance type | vCores | Memory (GB) | Storage (GiB) |
|---|---|---|---|---|
| ○ | c3.xlarge | 4 | 7.5 | 80 SSD |
| ○ | c3.2xlarge | 8 | 15 | 160 SSD |
| ○ | c3.4xlarge | 16 | 30 | 320 SSD |
| ○ | c3.8xlarge | 32 | 60 | 640 SSD |
| ○ | c4.large | 2 | 3.8 | EBS only |
| ○ | c4.xlarge | 4 | 7.5 | EBS only |
| ○ | c4.2xlarge | 8 | 15 | EBS only |
| ○ | c4.4xlarge | 16 | 30 | EBS only |
| ○ | c4.8xlarge | 36 | 60 | EBS only |
| ○ | c5.xlarge | 4 | 8 | EBS only |
| ○ | c5.2xlarge | 8 | 16 | EBS only |

Cancel    Save

# Step 2: Hardware

# Step 2: Hardware

## Pricing for Amazon EMR and Amazon EC2 (On-Demand)

Region: EU (Frankfurt) ⇅

| | Amazon EC2 Price | Amazon EMR Price |
|---|---|---|
| **General Purpose - Current Generation** | | |
| m5.xlarge | $0.23 per Hour | $0.048 per Hour |
| m5.2xlarge | $0.46 per Hour | $0.096 per Hour |
| m5.4xlarge | $0.92 per Hour | $0.192 per Hour |
| m5.12xlarge | $2.76 per Hour | $0.27 per Hour |
| m5.24xlarge | $5.52 per Hour | $0.27 per Hour |
| m5a.xlarge | $0.208 per Hour | $0.043 per Hour |
| m5a.2xlarge | $0.416 per Hour | $0.086 per Hour |
| m5a.4xlarge | $0.832 per Hour | $0.172 per Hour |
| m5a.12xlarge | $2.496 per Hour | $0.27 per Hour |
| m5a.24xlarge | $4.992 per Hour | $0.27 per Hour |
| m5d.xlarge | $0.272 per Hour | $0.057 per Hour |
| m5d.2xlarge | $0.544 per Hour | $0.113 per Hour |
| m5d.4xlarge | $1.088 per Hour | $0.226 per Hour |
| m5d.12xlarge | $3.264 per Hour | $0.27 per Hour |
| m5d.24xlarge | $6.528 per Hour | $0.27 per Hour |
| m4.large | $0.12 per Hour | $0.03 per Hour |
| m4.xlarge | $0.24 per Hour | $0.06 per Hour |

## Prices range between $.03 - $.27 for EMR

# Step 3: Cluster Settings

Create Cluster - Advanced Options    **Go to quick options**

Step 1: Software and Steps

Step 2: Hardware

**Step 3: General Cluster Settings**

Step 4: Security

**General Options**

Cluster name    VHB_2020_Cluster

☑ Logging ⓘ

    S3 folder    s3://aws-logs-480030033378-eu-central-1/elasticmapred 📂

☑ Debugging ⓘ

☑ Termination protection ⓘ

**Tags** ⓘ

| Key | Value (optional) | |
|-----|------------------|--|
| Add a key to create a tag | | |

**Additional Options**

☐ EMRFS consistent view ⓘ

**Custom AMI ID** | None ▼ | ⓘ

▼ Bootstrap Actions

Bootstrap actions are scripts that are executed during setup before Hadoop starts on every cluster node. You can use them to install additional software and customize your applications. Learn more 🔗

**Add bootstrap action** | Custom action ▼ | **Configure and add**

Cancel    Previous    Next

# Step 3: Cluster Settings



Install RStudioServer on Cluster

# Step 3: Cluster Settings



```
install-rstudio-server - short.sh - Editor                                    _ □ ×
Datei  Bearbeiten  Format  Ansicht  ?

#Installation of the specified RStudio Server Version with the defined User and Password.
#These Variables can be changed above, if needed.

grep -Fq "\"isMaster\": true" /mnt/var/lib/info/instance.json
if [ $? -eq 0 ];
then
    while [[ $# > 1 ]]; do
        key="$1"

        case $key in
            # The above specified RStudio Server version
                    --sd-version)
                VERSION="$2"
                shift
                ;;
            # The above specified user
            --sd-user)
                USER="$2"
                shift
                ;;
            # The password for the above specified user
            --sd-user-password)
                PASS="$2"
                shift
                ;;

            *)
                echo "Unknown option: ${key}"
                exit 1;
        esac
        shift
    done
    echo "********************************************"
    echo "  1. Download RStudio Server ${VERSION}   "
    echo "********************************************"
    wget https://s3.amazonaws.com/rstudio-dailybuilds/rstudio-server-rhel-${VERSION}-x86_64.rpm
    echo "         2. Install dependencies          "
    echo "********************************************"
    # This is needed for installing devtools
    sudo yum -y install libcurl libcurl-devel 1>&2
    echo "         3. Install RStudio Server         "
    echo "********************************************"
    sudo yum -y install --nogpgcheck rstudio-server-rhel-${VERSION}-x86_64.rpm 1>&2
    echo "     4. Create R Studio Server user        "
    echo "********************************************"
    epass=$(perl -e 'print crypt($ARGV[0], "password")' ${PASS})
    sudo useradd -m -p ${epass} ${USER}
    # This is to allow access to HDFS
    sudo usermod -a -G hadoop ${USER}
    echo "  5. Create environment variables file     "
    echo "********************************************"
```

## Bootstrap Script to Install RStudioServer

# Step 4: Security

# Cluster Ready to Use

Clone    Terminate    AWS CLI export

Cluster: EMAC_Hamburg_Cluster    **Waiting**  Cluster ready after last step completed.

| Summary | Application history | Monitoring | Hardware | Configurations | Events | Steps | Bootstrap actions |

**Connections:**    Enable Web Connection – Hue, Spark History Server, Resource Manager ... (View All)

**Master public DNS:**    ec2-18-184-169-221.eu-central-1.compute.amazonaws.com  SSH

**Tags:**    -- View All / Edit

**Summary**
- **ID:** j-2MB53YH15M535
- **Creation date:** 2019-05-17 10:42 (UTC+2)
- **Elapsed time:** 15 minutes
- **Auto-terminate:** No
- **Termination protection:** On  Change

**Configuration details**
- **Release label:** emr-5.23.0
- **Hadoop distribution:** Amazon 2.8.5
- **Applications:** Hive 2.3.4, Pig 0.17.0, Hue 4.3.0, Spark 2.4.0
- **Log URI:** s3://aws-logs-480030033378-eu-central-1/elasticmapreduce/ 📁
- **EMRFS consistent view:** Disabled
- **Custom AMI ID:** --

**Network and hardware**
- **Availability zone:** eu-central-1c
- **Subnet ID:** subnet-8ca58fc6 ↗
- **Master:** Running  1  m4.large
- **Core:** Running  2  m4.large
- **Task:** --

**Security and access**
- **Key name:** ffm_master
- **EC2 instance profile:** EMR_EC2_DefaultRole
- **EMR role:** EMR_DefaultRole
- **Auto Scaling role:** EMR_AutoScaling_DefaultRole
- **Visible to all users:** All  Change
- **Security groups for Master:** sg-3d3ce557 ↗ (ElasticMapReduce-master: master)
- **Security groups for Core & Task:** sg-3d3ee757 ↗ (ElasticMapReduce-Core & Task: slave)

# Hadoop Cluster Manager Overview



http://ec2-18-184-169-221.eu-central-1.compute.amazonaws.com:8088

# Accessing RStudioServer



http://ec2-18-184-169-221.eu-central-1.compute.amazonaws.com:8787

# Run Analysis

# Sparklyr: R Interface for Apache Spark

Source: https://spark.rstudio.com/

# Sparklyr Demo

## Sparklyr Demo: R Interface with Apache Spark

*Klaus Miller, Goethe University Frankfurt*

*May 2019*

Example: Cluster Analysis Using Spark.ML to predict cluster membership with the iris dataset

Slightly adapted from source: https://spark.rstudio.com/

Load Packages

```r
library(tidyverse)
```

Installation

```r
#install.packages("sparklyr")

# Upgrade to latest version
#devtools::install_github("rstudio/sparklyr")
```

Connecting to Spark

```r
library(sparklyr)
```

```
##
## Attaching package: 'sparklyr'
```

```
## The following object is masked from 'package:purrr':
##
##     invoke
```

```r
sc <- spark_connect(master = "local")
```

Source: https://github.com/stm/vhb_2020

# Sparklyr Machine Learning Library

## Algorithms

Spark's machine learning library can be accessed from sparklyr through the `ml_*` set of functions:

| Function | Description |
| --- | --- |
| ml_kmeans | K-Means Clustering |
| ml_linear_regression | Linear Regression |
| ml_logistic_regression | Logistic Regression |
| ml_survival_regression | Survival Regression |
| ml_generalized_linear_regression | Generalized Linear Regression |
| ml_decision_tree | Decision Trees |
| ml_random_forest | Random Forests |
| ml_gradient_boosted_trees | Gradient-Boosted Trees |
| ml_pca | Principal Components Analysis |
| ml_naive_bayes | Naive-Bayes |
| ml_multilayer_perceptron | Multilayer Perceptron |
| ml_lda | Latent Dirichlet Allocation |
| ml_one_vs_rest | One vs Rest |

# Introduction to Spark in R

# Terminating Cluster

# Terminating Cluster

Clone | Terminate | AWS CLI export

Cluster: EMAC_Hamburg_Cluster  **Terminating**  Terminated by user request

| Summary | Application history | Monitoring | Hardware | Configurations | Events | Steps | Bootstrap actions |

**Connections:** --

**Master public DNS:** ec2-18-184-169-221.eu-central-1.compute.amazonaws.com  SSH

**Tags:** --

### Summary

**ID:** j-2MB53YH15M535
**Creation date:** 2019-05-17 10:42 (UTC+2)
**Elapsed time:** 20 minutes
**Auto-terminate:** No
**Termination protection:** Off

### Configuration details

**Release label:** emr-5.23.0
**Hadoop distribution:** Amazon 2.8.5
**Applications:** Hive 2.3.4, Pig 0.17.0, Hue 4.3.0, Spark 2.4.0
**Log URI:** s3://aws-logs-480030033378-eu-central-1/elasticmapreduce/
**EMRFS consistent view:** Disabled
**Custom AMI ID:** --

### Network and hardware

**Availability zone:** eu-central-1c
**Subnet ID:** subnet-8ca58fc6
**Master:** Terminating  1  m4.large
**Core:** Terminating  2  m4.large
**Task:** --

### Security and access

**Key name:** ffm_master
**EC2 instance profile:** EMR_EC2_DefaultRole
**EMR role:** EMR_DefaultRole
**Auto Scaling role:** EMR_AutoScaling_DefaultRole
**Visible to all users:** All  Change
**Security groups for** sg-3d3ce557 (ElasticMapReduce-**Master:** master)
**Security groups for** sg-3d3ee757 (ElasticMapReduce-**Core & Task:** slave)

# Big Data Not So „Big" After All

# Thank You for Your Attention!

**Klaus Miller**

Theodor-W.-Adorno-Platz 4

60323 Frankfurt am Main

Tel: +49-(0)69-798-33-865

Email: klaus.miller@wiwi.uni-frankfurt.de

Web: fromdatatodecisions.com

 @klausmiller