## Preparation

- Install "Virtual Box"
- Go to <a href="http://hortonworks.com/downloads/#sandbox">http://hortonworks.com/downloads/#sandbox</a>
  - Search for "Archive" and download 2.4 version if you have less than 10GB memory
- OR: search for "Tutorial: Sandbox on Azure" and follow instructions
- Start virtual machine and open
  - http://127.0.0.1:8080
  - http://127.0.0.1:9995





#### Omics-based Cancer Diagnostics

Tim Conrad, Freie Universität Berlin

Forschungs Campus MODAL, MedLAB



SPONSORED BY THE









Introduction

SPARK & HDFS

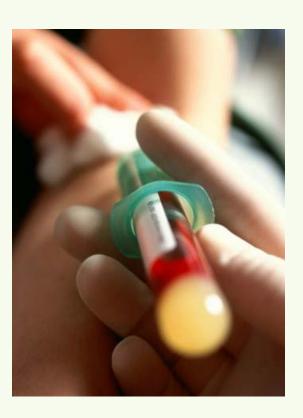
**Getting Started** 

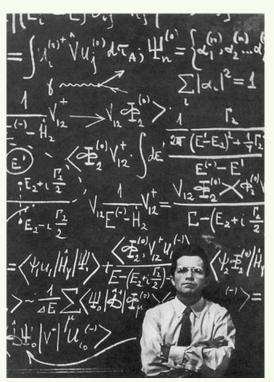
Simple Example

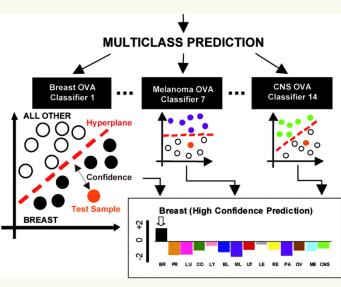
**Application** 

Classification of biomedical data

ACQUIRE DATA ANALYSE DATA CLASSIFY DATA

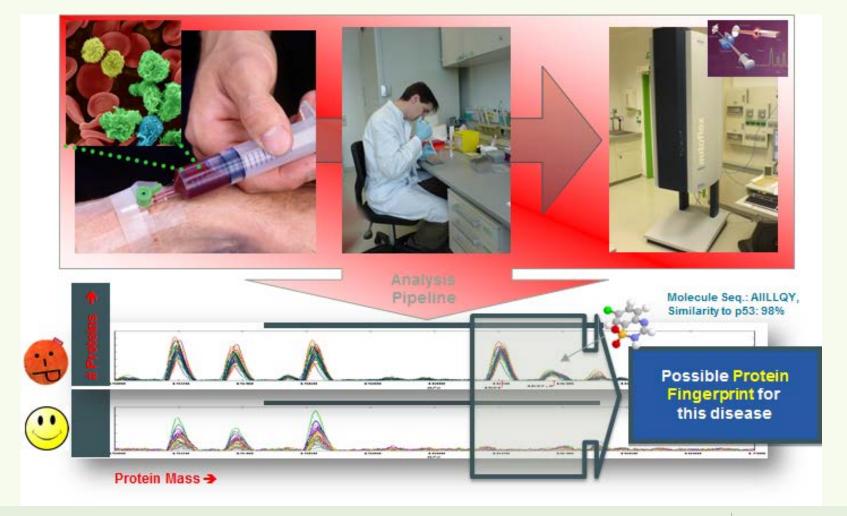




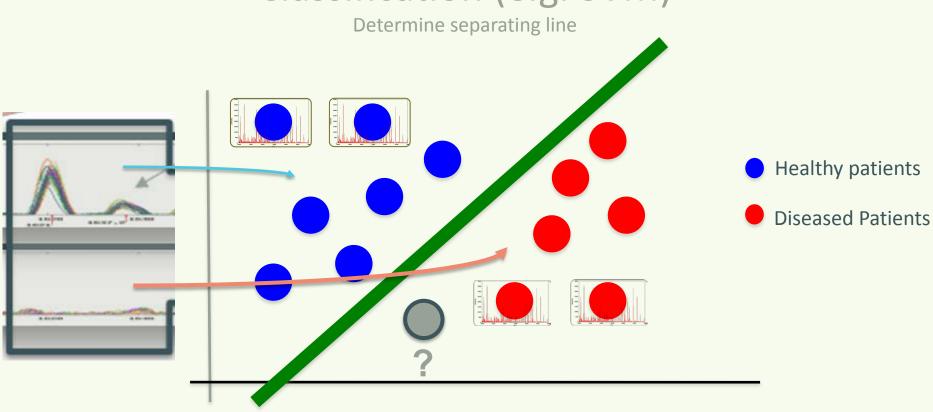








#### Classification (e.g. SVM)



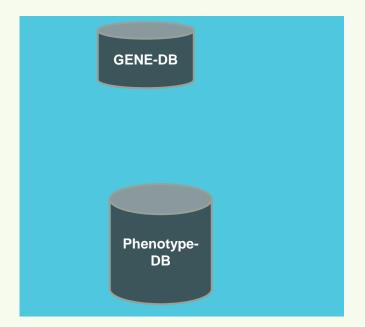


#### **Biomarkers Discovery Workflow**

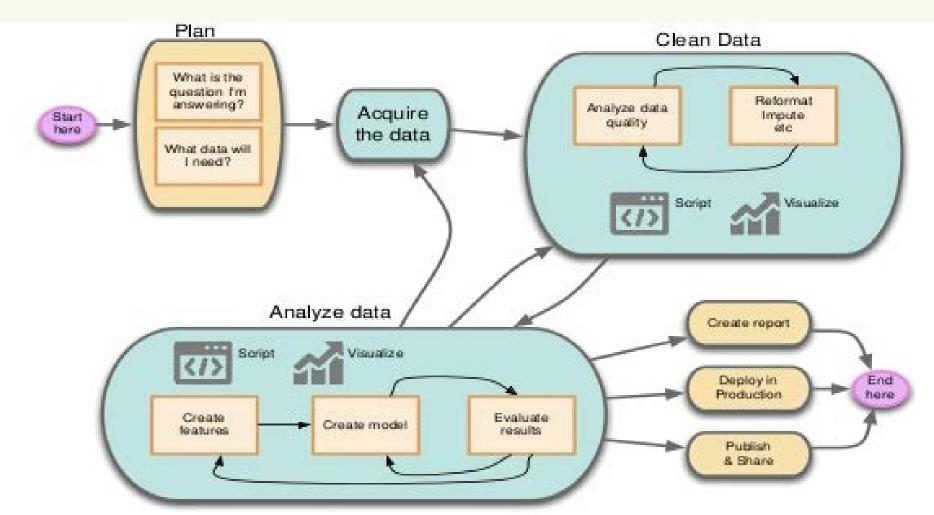






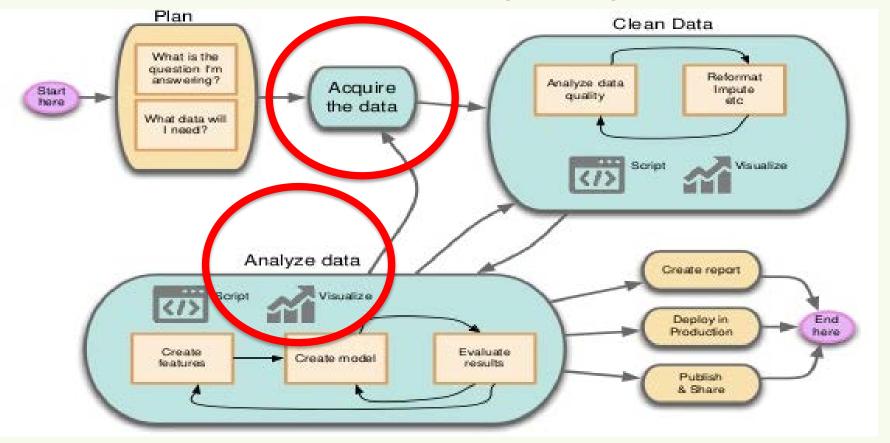




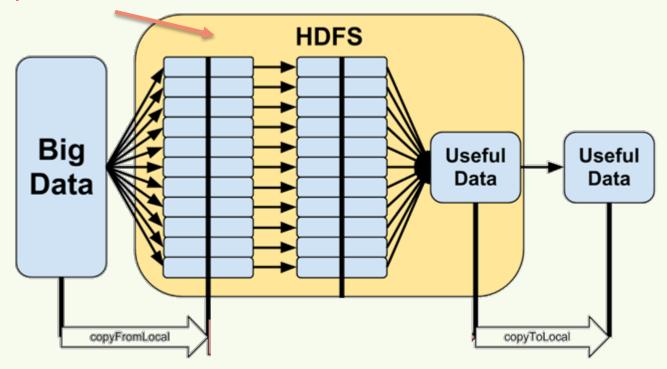




#### What if this is too big for a single machine?



# Split data AND analysis to many machines!



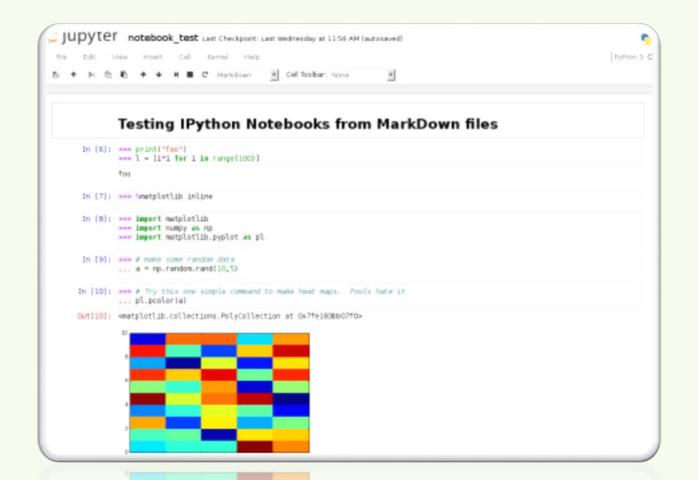


### Machine learning with Spark



### What do we learn?

- How to use a LOCAL SPARK system with HDFS
- How to change to a DISTRIBUTED SPARK system
- Examples:
  - k-Means
  - Predicting breast cancer from proteomics data





- Cassandra
- Shell
- Elasticsearch
- Spark

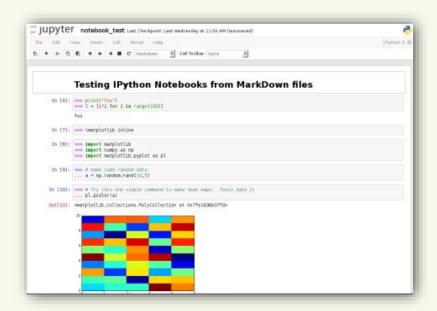
• Flink

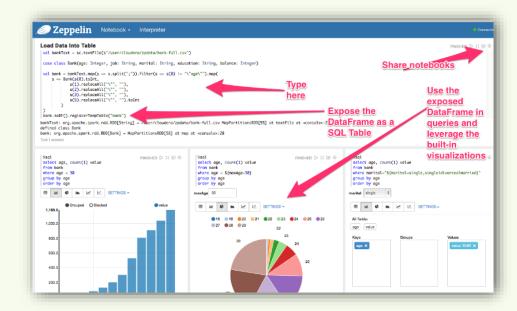
Tajo

• Geode

• Build your own!

- Hive
- Ignite
- Lens
- Markdown
- Postgresql, Hawq
- Scalding







### What can be done next?

Flink \$Spark Algorithm • Use a DISTRIBUTED system Resource pool Algorithm View Zeppelin Server Web browser Resource pool "Algorithm runs where resource exists" Interpreter Process



### Hands-on

• Task 1 (201):

bit.ly/dsss201

Diagnose breast cancer patients

• Task 2 (101):

bit.ly/dsss101

Get more familiar with Spark & Zeppelin





Notebook +

Interpreter

Configuration

#### Welcome to Zeppelin!

Zeppelin is web-based notebook that enables interactive data analytics.

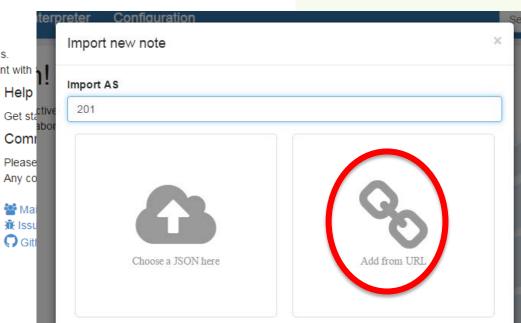
You can make beautiful data-driven, interactive, collaborative document with

#### TOLEDOOK !

🎎 Import note

to order new note

- © 201
- AON Demo
- Australian Dataset (Hive example)
- Australian Dataset (SparkSQL example)
- Hello World Tutorial
- 🖰 IoT Data Analysis (Keynote Demo)
- nagellan-blog
- Phoenix demo
- Predicting airline delays
- Sensors & Machines Predictive Analysis
- Single view demo
- twitter the
- Zeppelin Tutorial



- Open IN YOUR BROWSER: bit.ly/dsss201
- 2. Copy new link
- 3. Use new link for importing

