**Sanitas**

**FIU**

**Robert Stempel College
of Public Health
& Social Work**

SeasonCaster
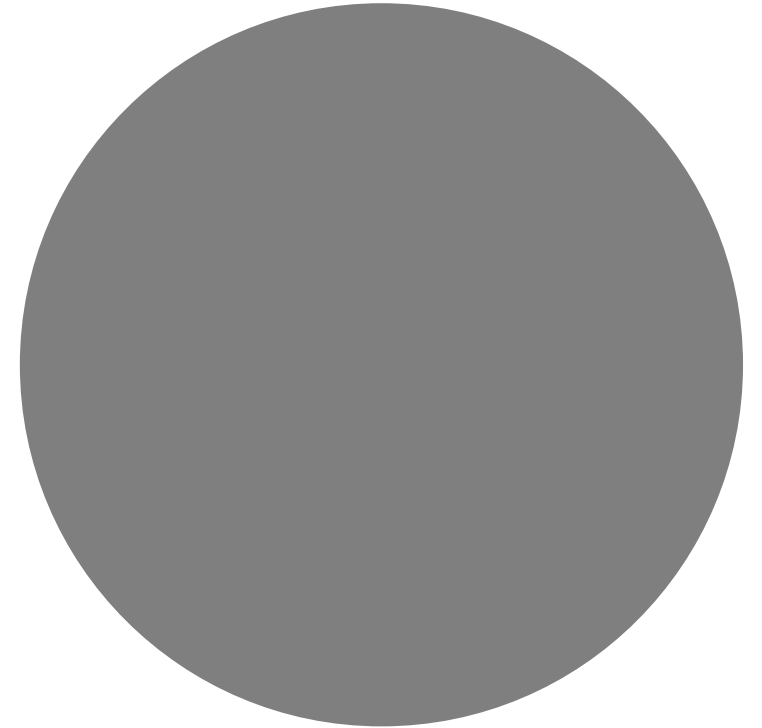*Business Forecasting App*

# Data Science, Big Data and Analytics: Present and Future

Perspectives from Academia, Industry and Consulting

Zoran Bursac, PhD, MPH

Josh Callaway, MS, MPH

Fernando Lopez, MS, PhD Candidate

# Data

- Health care, business, technology -> data

- Big data -> voluminous data sets (structured or unstructured)

- Produced every day all around us

- Analytics -> examining data to detect patterns

Different sources, different sizes

High variety, volume, velocity

Online networks, web pages, audio/video, social media, logs

Techniques-> machine learning, data mining, natural language processing, statistics

Extraction, preparation, storage/warehousing, blending, analytics

# Big Data Analytics and Data Science

Current Trends and Common Data Science Tools
Process, Perform and Visualize

# Free Open Source

| | | | |
|---|---|---|---|
| Hadoop->distributed processing of large data across clusters | Hive->warehouse to manage large data in distributed SQL storage | Kafka->real time pipeline of streaming data | Pig->large data analytics |
| R, Rstudio, ggplot-> analytics and data visualization | Python, Julia -> high level programming with efficient algorithms and speed for large data processing | Jupyter notebook -> manage documents such as code, explanatory and shared | RapidMiner -> data preparation, machine learning and model deployment |

# Do you need to know all? NO

Hadoop     R     SQL     Python     Hive     Pig etc…

# Big data analytics for personalized medicine and pharmacogenomics



Current Opinion in Biotechnology

- Cirillo and Valencia (2019). Machine learning algorithms for multi-view data analysis. Data from multiple sources (genomic, proteomic, metabolomic) used to identify associations within and between multiple sets of patients, and generate models for patient clustering.

# The Latest Buzzwords

| | | | | |
|---|---|---|---|---|
| Data science | Artificial intelligence | Machine learning | Data mining | Big data |
| Data warehouse | Data lake | Cloud computing | Hadoop | Internet of things |

www.Kaggle.com

Diabetes

# Step 1.
# Local Environment

# Step 2.
# Insert into
# Remote Database

```r
db_diabetic <- dbPool(
  RMySQL::MySQL(),
  dbname = "db_diabetic",
  host = "globalhealth.cunhdm3u6041.us-east-2.rds.amazonaws.com",
  username = "db_diabetic",
  password = "uh29vawh4tnWHFHcJ"
)
```

```r
demographic <- read_csv("demographic.csv")
ref_encounter <- read_csv("ref.encounter.csv")
admission <- read_csv("admission.csv")
diagnosis <- read_csv("diagnosis.csv")
lab <- read_csv("lab.csv")
procedure <- read_csv("procedure.csv")
medication <- read_csv("medication.csv")
payment <- read_csv("payment.csv")
```
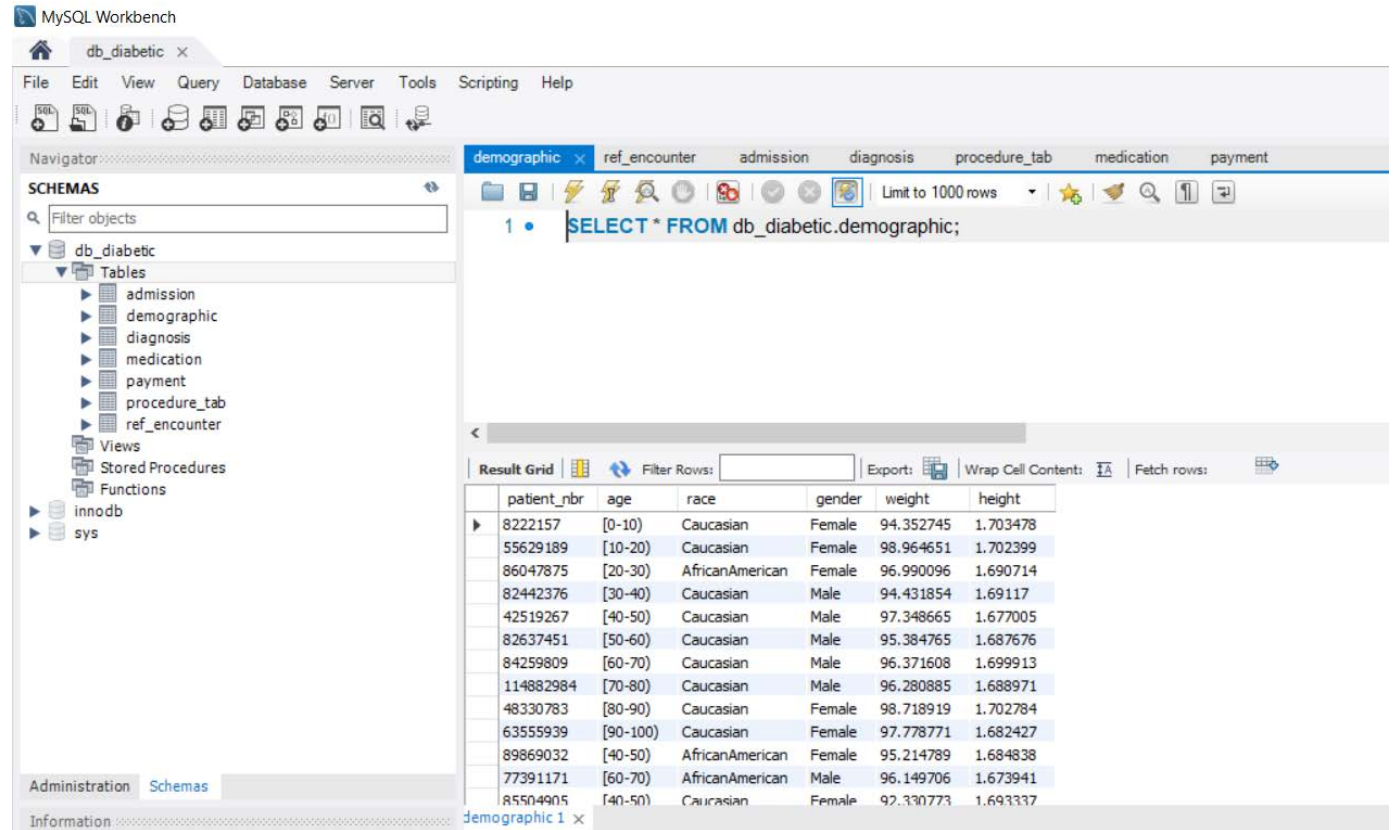
```r
query_ref_encounter1 <- paste(
  "CREATE TABLE ref_encounter (
  patient_nbr DOUBLE,
  encounter_id DOUBLE
);"
)
# query2 <- "INSERT INTO Product_Names
#   VALUES (?, ?);"

dbExecute(db_diabetic, query_ref_encounter1)

# Begin the query
query_ref_encounter2 <- "INSERT into ref_encounter (patient_nbr, encounter_id) VALUES"

# Finish it with
query_ref_encounter3 <- paste0(
  query_ref_encounter2,
  paste(sprintf("('%s', '%s')",
                ref_encounter$patient_nbr,
                ref_encounter$encounter_id
  ),
  collapse = ","))
)

dbExecute(db_diabetic, query_ref_encounter3)
```
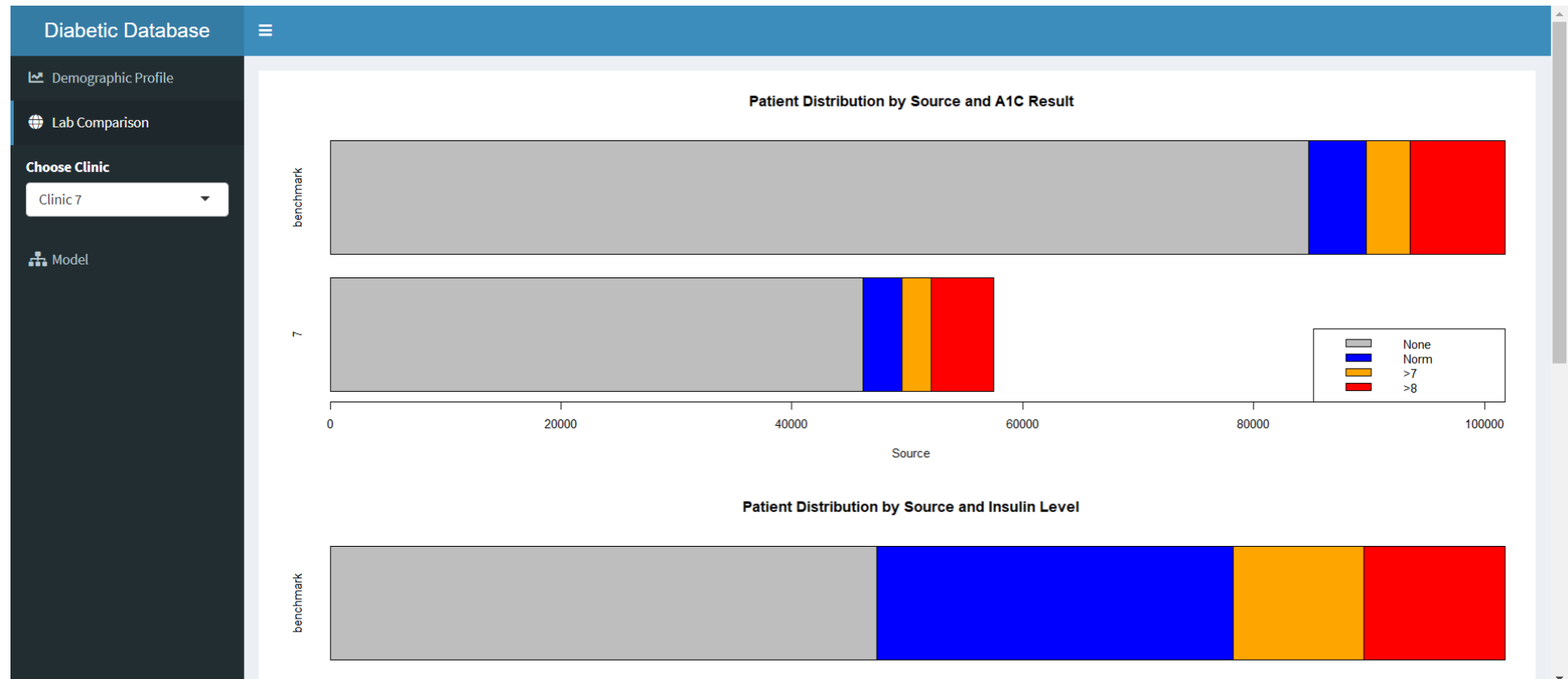
# Step 3.
# SQL Database

# Step 4.
# R Shiny Web App

# Step 4.
# R Shiny Web App

# Local Storage

Text Files (.txt)

Comma Separated Value Files (.csv)

Excel Database (.xlsx)

Microsoft Access Database (.accdb)

# Remote Storage

**Dropbox**

**Google Sheets**

**AWS S3 Bucket**

**Cloud Database**
- MySQL
- MongoDB
- PostgreSQL
- NoSQL
- Oracle

# Cloud Platforms

| Devices | Device Connectivity | Storage | Analytics | Presentation & Action |
|---|---|---|---|---|
| | Event Hubs | SQL Database | Machine Learning | App Service |
| | Service Bus | Table/Blob Storage | Stream Analytics | Power BI |
| | External Data Sources | DocumentDB | HDInsight | Notification Hubs |
| | | External Data Sources | Data Factory | Mobile Services |
| | | | | BizTalk Services |

Microsoft Azure

Current Content Ignite Infrastructure

# Distributed Computing Technology

All processing jobs (scripts; i.e. R, Python, Scala, etc.) are divvied up among all available processing units (computers, cores, threads, etc.)

Hadoop

Spark

Microsoft R Open

Microsoft R Open Multithreaded Performance

# Machine Learning

# Microsoft Azure Power BI

# Microsoft Azure Power BI

# Microsoft Azure Power BI

# Solution Architectures at Sanitas

- Analytics on big data

- Data warehouse

- Real-time analytics

- Population Health Management for Healthcare

# Sanitas

**BUSINESS INTELLIGENCE STRATEGY**
**Sanitas USA**

**Data Sources**

**Data Governance**

- ✓ - Standards and Policies
- ✓ - Data Quality (DQ)
- ✓ - Data Security and Privacy
- ✓ - Architecture/Integration
- ✓ - DW and Business Intelligence (BI)
- ✓ - Self-service Architectures

**Information Analysis for Business units**
**Clinical – Financial - Operational**

**Visualize and Report**

**Dashboards, KPIs and Metrics**

**Health Care Team**

**Data Types**

Structured Data

Unstructured Data

Print, Scan, Fax

Real Time Data

Streaming Data

**Training and Predictive Experimentation**

Machine Learning Studio

**Visualized Analytics**

**Ensure Data Quality**

Data Factory

**ETL/ELT workflows**

SQL Data Warehouse

**Cloud data warehouse**

Power BI

Azure Databricks

**Apache Spark-based analytics**

Azure Analysis Services

Azure HDInsight

**Hadoop, Spark, Hive, LLAP, Kafka, Storm, R**
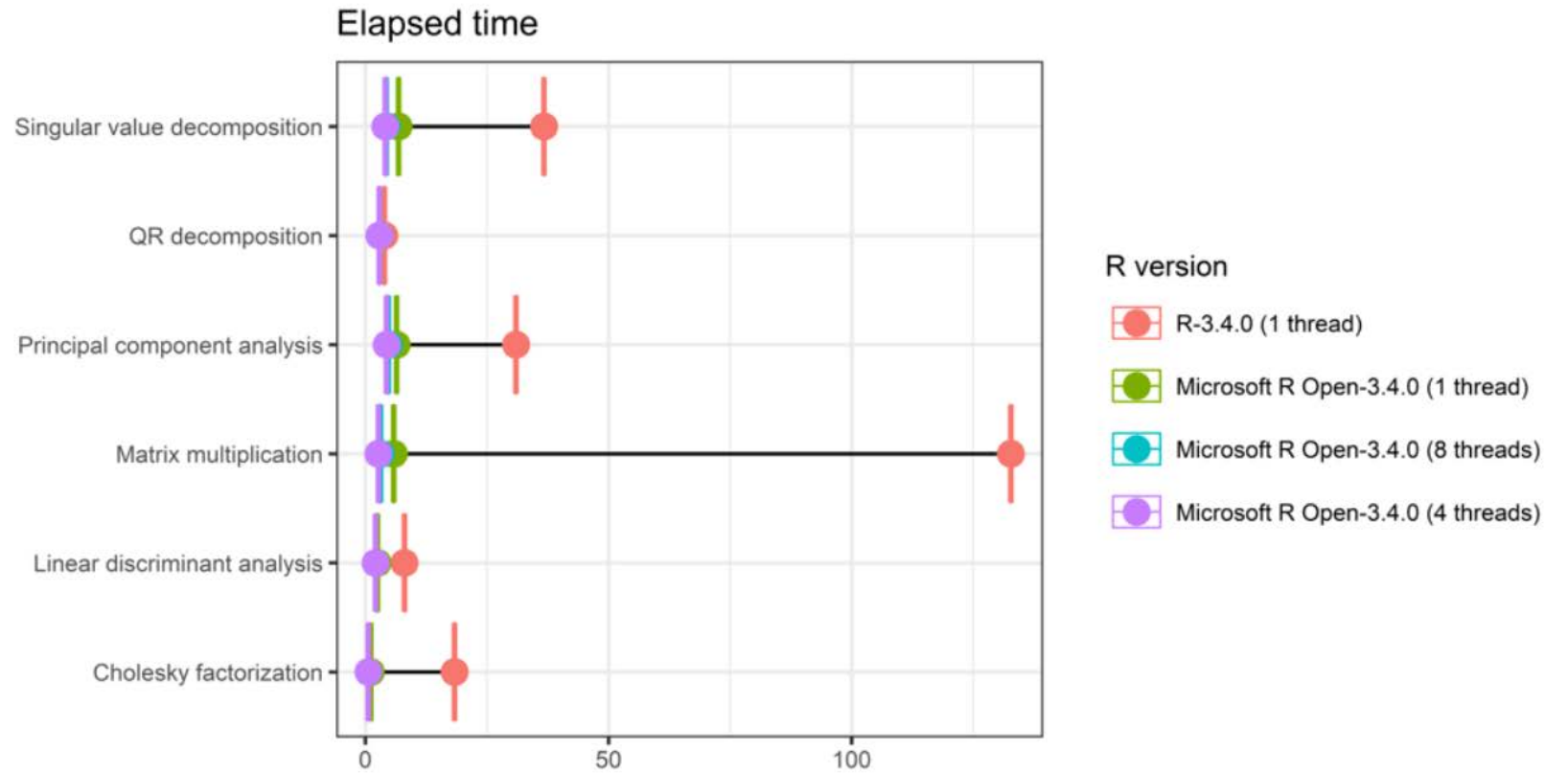
**Reports**

eCW Report Delivery

**Operations**

POPULATION HEALTH MANAGEMENT

**Security, Roles and Audit**

**IT Alignment**

**SANITAS USA – Strategic plan**

ML Studio Experiment

ML Studio Model Train

ML Studio Model Score

Diabetes Analysis ‣ Score Model ‣ Scored dataset

rows 4500  columns 14

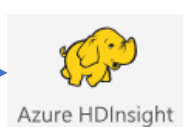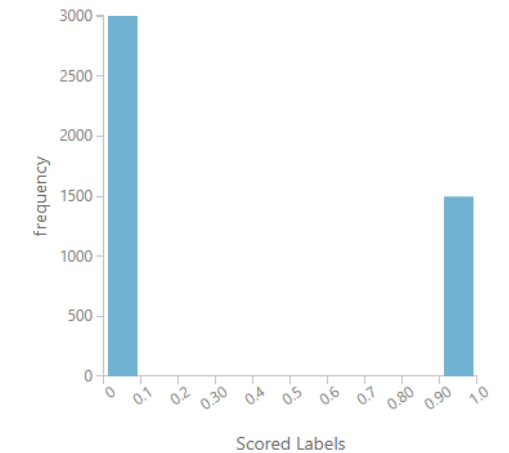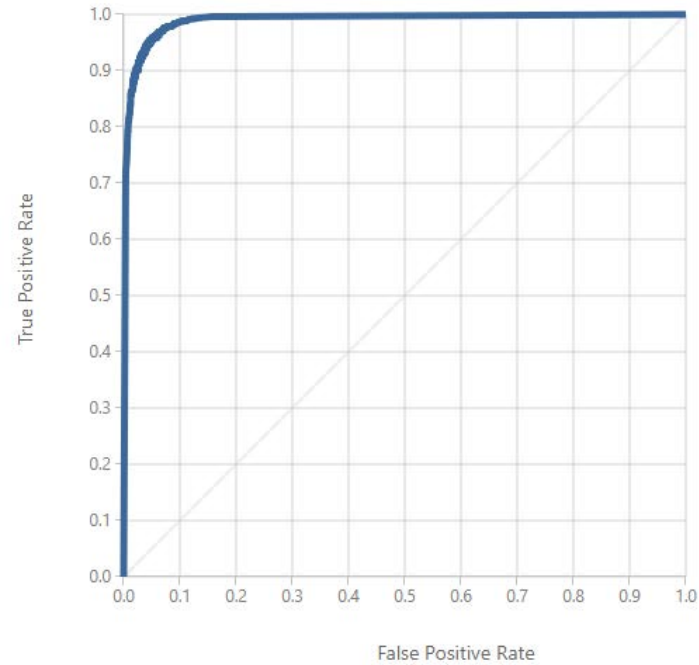| ...se | DiastolicBloodPressure | TricepsThickness | SerumInsulin | BMI | DiabetesPedigree | Age | Diabetic | Physician | Ln(Age) | Scored Labels | Scored Probabilities |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | 0.643229 | -0.399444 | 1.586817 | 1.439866 | 0.476071 | 0.017857 | 1 | Wade Munger | 0.035804 | 1 | 0.998918 |
| | -0.55022 | 1.11204 | 1.812272 | -0.156933 | 0.006973 | 0.017857 | 1 | Niew Leekpai | 0.035804 | 1 | 1 |
| | 0.583557 | -1.292594 | -0.660226 | 0.775076 | 0.225251 | 0.017857 | 0 | Roman Pilcher | 0.035804 | 0 | 0.123397 |
| | -0.967927 | -0.742963 | -0.772953 | 0.664622 | 0.388068 | 0.035714 | 0 | Mara Rasmussen | 0.070017 | 0 | 0.000004 |
| | -0.192185 | -0.949075 | 0.271658 | -0.662061 | 0.021152 | 0.017857 | 1 | Ethan Rincon | 0.035804 | 1 | 0.999995 |
| | 1.060936 | -0.33074 | -0.742893 | -1.338164 | 0.033424 | 0.035714 | 0 | Vaughn Oquendo | 0.070017 | 0 | 0.000011 |
| | -1.683996 | 0.56241 | -0.667741 | 1.764616 | 0.019228 | 0.232143 | 0 | Neandro Baeza | 0.370849 | 0 | 0.000025 |
| | -0.908254 | 0.287594 | 1.368876 | -1.073757 | 0.039045 | 0.178571 | 0 | Delmar Pelchat | 0.299754 | 0 | 0.000011 |
| | 1.299626 | -1.22389 | -0.412224 | 0.696152 | 0.040121 | 0.321429 | 1 | Nazzareno Piccio | 0.476447 | 1 | 0.998053 |
| | -0.729237 | 0.974632 | -0.562528 | 0.016523 | 0.202185 | 0.464286 | 1 | Billie Stonge | 0.620054 | 1 | 0.999978 |
| | 0.941591 | -1.430001 | -0.915742 | -1.275964 | 0.293586 | 0.285714 | 0 | Jimmie Turman | 0.435929 | 0 | 0.000005 |
| | 1.478643 | 2.692228 | 3.525736 | 0.4355 | 0.304365 | 0.160714 | 1 | Deanna Ball | 0.274517 | 1 | 0.999623 |
| | 0.822246 | -1.22389 | -0.833075 | -1.212273 | 0.056504 | 0 | 0 | Jenny Norgaard | 0 | 0 | 0.000028 |
| | 1.180281 | 0.837225 | -0.65271 | 0.701326 | 0.260088 | 0.232143 | 0 | Daitaro Ishida | 0.370849 | 0 | 0.000919 |

Statistics

Visualizations

Scored Labels
Histogram

# Rejoinder

Data -> our biggest asset        Emphasis on "good" data

Processing streams -> wrangling, carpenting

Storage –> lakes, warehousing, data bases

Analytics -> mining, machine learning

Output -> new knowledge, information, inferences

Feed back to the users -> gather more data

# Question to the global audience

- What are your needs and where are you currently with respect to?
  - Data collection, quality, storage, analytical and computing power
    - Where is data coming from, single or multiple sources
    - Who is maintaining data quality and fidelity
    - Do you have adequate storage with proper security; planning for the future
    - Are you investing in resources and trained personnel with data science skills