# Hands on Introduction to IBM's
## Data Science Experience

Power of data. Simplicity of design. Speed of innovation.

**Joel Patterson**

# Agenda

| Time | Description |
|------|-------------|
| 7:30 AM - 8:00 AM | **Registration and Coffee** |
| 8:00 AM - 8:30 AM | **Overview of the Watson Data Platform and IBM Data Science Experience (DSX)** |
| 8:30 AM - 10:00 AM | **Lab 1 -  Surviving the Titanic** |
| 10:00 AM - 11:00 AM | **Lab 2 - Machine Learning with Spark ML** |
| 11:00 AM – 11:15 AM | **Break** |
| 11:15 AM – 11:45 AM | **Lab 3 -  R, Shiny, and GUI Interfaces** |
| Extra time | **Lab 4 -  Choose From Several Options** |
| 11:45 AM – 12:00 PM | **Questions and Wrap Up** |

# Participant Background

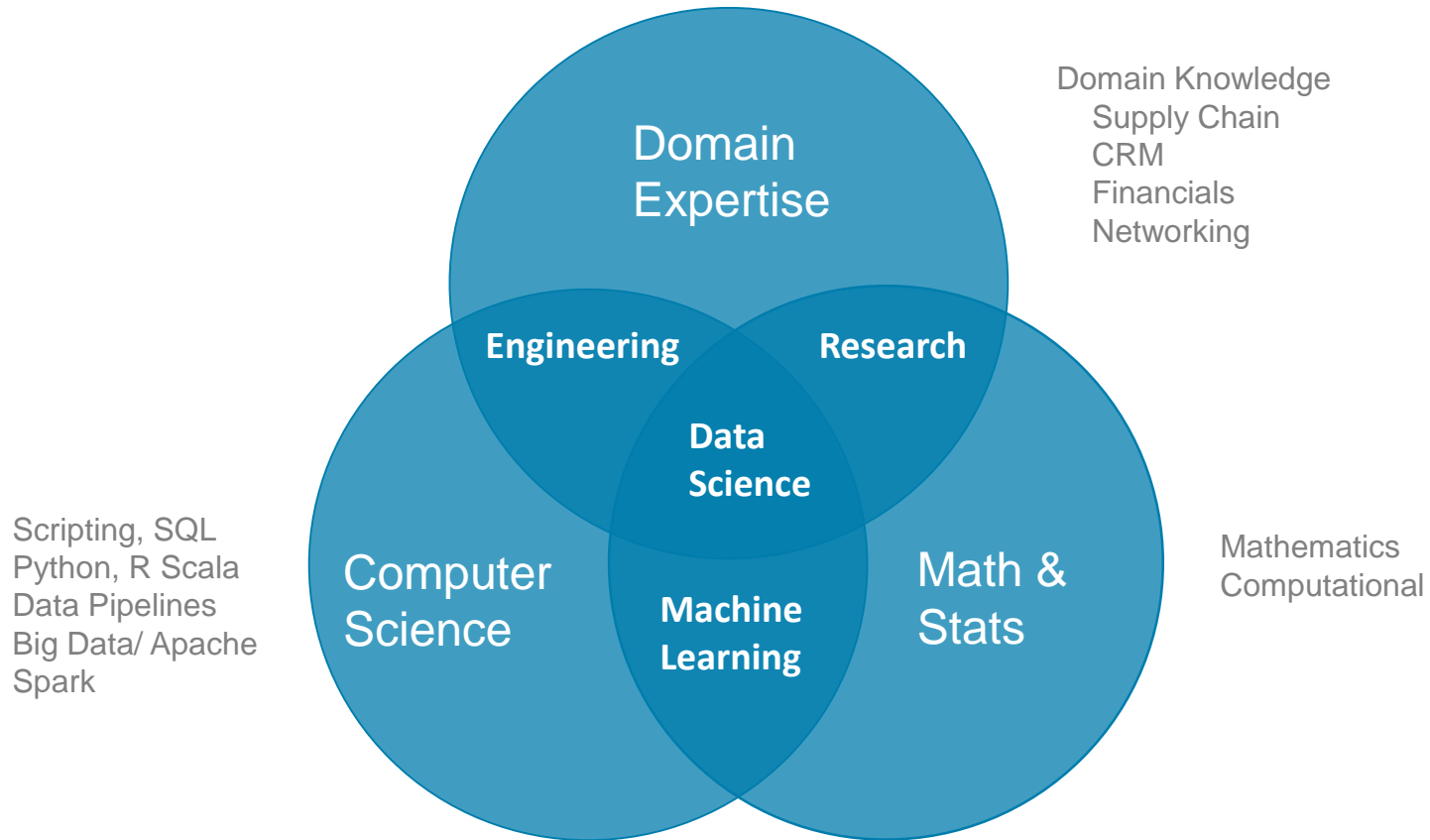## Open Source

- R/Python/Scala
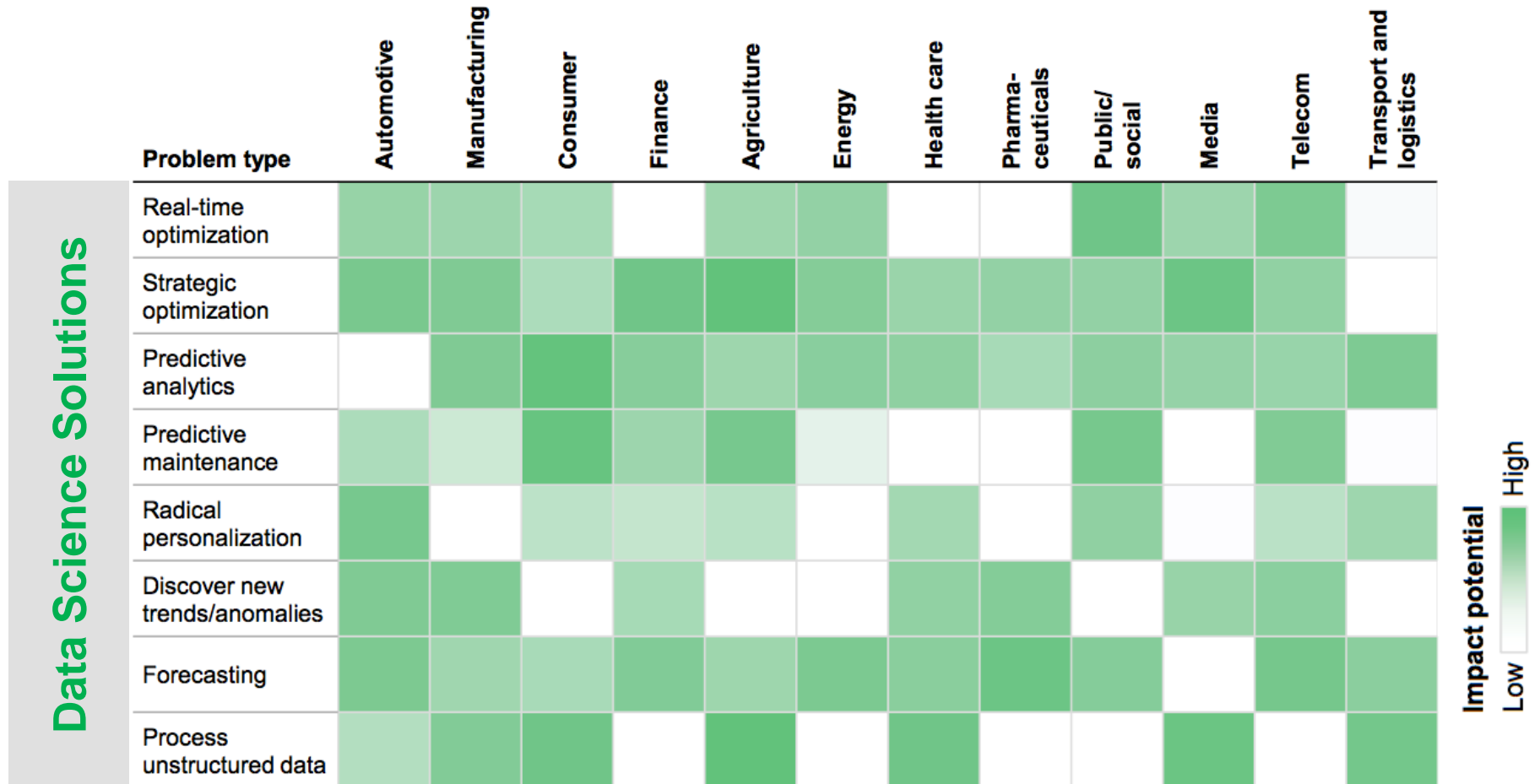- Jupyter Notebook
- Spark
- Hadoop

## IBM

- Bluemix

# What is a Data Scientist?



Domain
Expertise

Domain Knowledge
Supply Chain
CRM
Financials
Networking

**Engineering**

**Research**

**Data
Science**

Scripting, SQL
Python, R Scala
Data Pipelines
Big Data/ Apache
Spark

Computer
Science

**Machine
Learning**

Math &
Stats

Mathematics
Computational

*Data Science Projects Require Multiple Skills*

# Data Science Impact Across Industries and Use Cases

## $10s of Billions in each industry and use case



SOURCE: McKinsey Global Institute analysis

# Challenges in delivering value with Data Science

## Data

- Data resides in silos and difficult to access
- Detailed data was never stored
- Unstructured and external data wasn't considered

## Governance

- If the data isn't secure, self-service isn't a reality
- Understanding lineage and getting to a system of truth

## Skills

- Data Science skills are in low supply and high demand
- Nurturing new data professionals is challenging

## Infrastructure

- Need an environment that enables collaboration and deployment to production
- Discrete tools present barriers to progress

# Watson Data Platform

# IBM Watson Data Platform

## Mission: Make Data Simple and Accessible to All

Platform.          Method.          Ecosystem.

# Data Refinery – Open Beta

## *A Breakthrough Approach to Explore and Prepare Data*

**Data Engineer**



## Wrangle

Interactively explore, resolve quality issues, enrich, classify, standardize, and summarize data.

## Flow

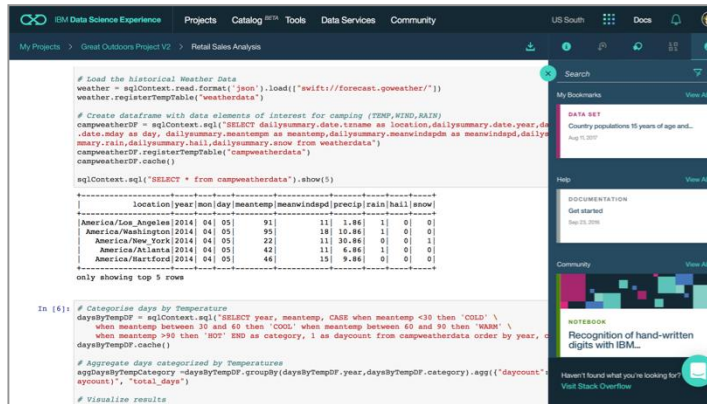Create data flows visually, schedule for repeatability, monitor and notify

## Adapt

Connect to 30+ cloud and on-premises stores and scale on demand with cataloging and governance

# Data Science Experience

*Brings together everything a Data Scientist needs to be successful*

**Data Scientist**



## Learn
Built-in learning to get started or go the distance with advanced tutorials

## Create
The best of open source and IBM value-add to create state-of-the-art data products

## Collaborate
Community and social features that provide meaningful collaboration

# Data Catalog – Open Beta

*Unlock tribal knowledge to unleash your data professionals*

## Discover

Intelligent discovery of data, advanced classification and profiling to provide context

## Catalog

A rich metadata index of all data, with social collaboration and enhanced findability

## Govern

Powerful governance policy tools to control and protect access to data with visibility to data use

# Intelligent data fabric provides consistent platform experience



This fabric remains consistent throughout the Watson Data Platform experience – regardless if you are ingesting data, shaping data, building algorithms, deploying models and more…

1

# How does WDP help fulfill the promise of your data?

## Data

Puts every important data source at the fingertips of the teams that need it wherever resides

## Governance

Enforces your policies without getting in the way of delivering insights

## Skills

Makes the most of the data professionals you have and helps them grow and learn from each other as a team

## Infrastructure

Delivers the foundation for your first data project through to the complete transformation of your business

# Data Science Experience

# Core Attributes of the Data Science Experience

**IBM Data Science Experience**

| **Community** | **Open Source** | **IBM Added Value** |
|---|---|---|
| • Find tutorials and datasets | • Code in Scala/Python/R/SQL | • IBM Machine Learning |
| • Read articles and papers | • Jupyter Notebooks | • SPSS Modeler Canvas |
| • Connect with Data Scientists | • RStudio IDE and Shiny | • Prescriptive Analytics - DOcplexcloud |
| • Share comments | • Apache Spark | • Projects and Version Control |
| • Copy and share notebooks | • Your favorite libraries | • Managed Spark Service |

Powered by IBM **Watson Data Platform**

# DSX Architecture

## DSX architecture

Last updated: **June 27, 2017**

Search this document

♡ 0  🔖

DSX provides you with the environment and tools to solve your business problems by collaboratively analyzing data. This

illustration shows how the architecture of DSX is centered around the project. A project is how you organize your resources for solving a business problem.

# Community Cards provide in-context learning

# Collaborate Using Projects

# Add Collaborators to a Project

# Add New Collaborator

Add users to your project for collaboration. Users with write access can add services to your project...

| Type name or email address |
|---|

**Select** ⌃

Viewer

Editor

Admin

*Cancel*  Add

# GitHub Integration

# Live chat on Intercom for support from the IBM team and to provide your feedback on how we can improve

# What is a "Notebook"?

## Pen and Paper

- **Pen and paper has long provided the rich experience that scientists need to document progress through notes and drawings:**
  - Expressive
  - Cumulative
  - Collaborative



## Notebooks

- **Notebooks are the digital equivalent of the "pen and paper" lab notebook, enabling data scientists to document reproducible analysis:**
  - Markdown and visualization
  - Iterative exploration
  - Easy to share

# Integrated Jupyter Notebooks for interactive and collaborative development - seamless execution on Spark

# From a Notebook in DSX you can use IBM's managed Spark Service to blend multiple data types, sources, and workloads

Execute SQL Statements

Streaming Analytics via Micro-batch

M.L. and Statistical Algorithms

Distributed Graph Processing Framework

| Spark SQL | Spark Streaming | MLlib Machine Learning | Graph |
|---|---|---|---|

- General compute engine
- Basic I/O functions
- Task dispatching
- Scheduling

**Spark Core**

**Data Sources**

| IBM Cloud | Public Cloud | Cloud Apps | On-Premises |
|---|---|---|---|

**IBM Cloud**

BigInsights (HDFS)  Cloudant (DBaaS)  dashDB (Analytics)

SQDB (Managed DB2)  Swift (Object Storage)

**Public Cloud**

amazon web services | S3   Cassandra
mongoDB   redis
rackspace   Microsoft Azure
MySQL   HIVE   PostgreSQL   HDFS
dBase   APACHE HBASE

**Cloud Apps**

NETSUITE   salesforce   CSV
JDBC   { JSON }   Parquet
elasticsearch.   AVRO

**On-Premises**

ORACLE
SAP
IBM DB2

# Benefits of Spark for Data Science



- General compute engine
- Basic I/O functions
- Task dispatching
- Scheduling

| Spark SQL | Spark Streaming | MLlib Machine Learning | GraphX Graphing |
|-----------|-----------------|------------------------|-----------------|

Spark Core

IBM

+

Spark

- Allows Data Scientists to code at scale
  - In-Memory processing that scales in a distributed architecture
- Supports multiple programing interfaces (Scala, Python, Java and R)
- Provides unified APIs (SQL, Streaming, Machine Learning, etc.)

# The Spark service uses Bluemix Object Storage as its preferred data store for building performant applications

- **Object storage provides inexpensive, scalable and self-healing retention of massive amounts of unstructured data**

- **Every object exists at the same level in a flat address space**

- **Bluemix Object Storage has a drag-and-drop upload and Swift API for programmatic access**

Object Storage

IBM

# Supported Data Sources/Targets for DSX via on- premises and cloud Connectors

| Cloud Sources | On-Premises Sources | Cloud Targets | On-Premises Targets |
|---|---|---|---|
| Amazon Redshift | Apache Hive | Amazon S3 | IBM DB2® LUW |
| Amazon S3 | Cloudera Impala | Bluemix Object Storage | IBM Pure Data for Analytics® |
| Apache Hive | IBM DB2® LUW | IBM Cloudant™ | Teradata |
| Bluemix Object Storage | IBM Informix® | IBM dashDB | |
| IBM BigInsights™ on Cloud * | IBM Pure Data for Analytics® | IBM BigInsights™ on Cloud * | |
| IBM Cloudant™ | Microsoft SQL Server | IBM DB2® on Cloud | |
| IBM dashDB | MySQL Enterprise Edition | IBM SQL Database | |
| IBM DB2® on Cloud | Oracle | IBM Watson™ Analytics | |
| IBM SQL Database | Pivotal Greenplum | PostgreSQL on Compose | |
| Microsoft Azure | PostgreSQL | SoftLayer Object Storage | |
| PostgreSQL on Compose | Sybase | | |
| Salesforce | Sybase IQ | | |
| SoftLayer Object Storage | Teradata | | |

All of the supported targets are compatible with each source

# DSX has RStudio built into the experience thanks to our strategic partnership

# With RStudio you can create Shiny web applications to make your analysis accessible to the business

# Operationalize insights with IBM Machine Learning

# IBM Machine Learning

**IBM Cloud** Object Storage

MySQL

hadoop HDFS

TERADATA  **IMS**

amazon web services | **S3**

Microsoft® SQL Server®

IBM DB2

Microsoft Azure

Data Science Experience

Validate model

Area Under ROC Curve

Web Service

**Data Access:**
- Easily connect to Behind-the-Firewall and Public Cloud Data

- Catalogued and Governed Controls through Watson Data Platform

**Creating Models:**
- Single UI and API for creating ML Models on various Runtimes

- Auto-Modeling and Hyperparameter Optimization

**Web Service:**
- Real-time, Streaming, and Batch Deployment

- Continuous Monitoring and Feedback Loop

**Intelligent Apps:**
- Integrate ML models with apps, websites, etc.

- Continuously Improve and Adapt with Self-Learning

# DSX Canvas

- DSX Canvas will have compatibility with legacy SPSS Modeler streams

- Multiple execution runtimes: SPSS Modeler, SparkML

- Planned support for R/Python/SQL code





- Pipeline deployment from DSX Canvas (left) via IBM Machine Learning

# Stream Designer – Open Beta

- Characteristics of Stream Processing
    - Continuous processing
    - Multiple varied data sources
    - High data rates/ data volumens
    - Near-real time action
- DSX
    - Design stream flow with new Stream Designer
    - Executes in Streaming Analytics Service (based on IBM Streams)
    - Can invoke stream within Jupyter notebooks using Stream API

# DSX Local

- **Very similar to the public cloud version of DSX**

- **Runs on hardware that is provided by the customer**
  - The DSX Local software and hardware are managed by the customer

- **DSX Local comes with all the software it needs to run, although it can integrate with existing customer systems such as**
  - Databases and HDFS storage
  - LDAP servers for authentication

# Get Started with Data Science Experience Today!

## Calling all Data Scientists!

▪Our mission is to win the **hearts and minds** of Data Scientists

▪IBM Data Science Experience is a **freemium model** with value-add features, pricing and up-sell in development

▪**Sign up** and encourage your colleagues to do so at **datascience.ibm.com**

**IBM Data Science Experience**
**https://www.youtube.com/watch?v=1HjzkLRdP5k&t=29s**