# Hands on Introduction to
## IBM Data Science Experience

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

**Asad Mahmood**
**Jean Bright**
**Charles Morrison**

# Agenda

**Overview of DSX, Watson Data Platform and Use Case (Wednesday Feb. 28, 2018)**

| | |
|---|---|
| 3:00 PM - 4:00 PM | **Introduction to IBM Data Science Experience and Watson Data Platform** |
| 4:00 PM - 5:00 PM | **Overview of Use Case and Solution Approach** |

**Hands-On DSX Labs (Thursday March 1, 2018)**

| | |
|---|---|
| 9:00 AM - 10:30 AM | **Lab 1 -  Setting Up Your DSX Environment and Exploratory Data Analysis** |
| 10:30 AM - 12:00 PM | **Lab 2 - Data Visualization with R Studio and Shiny** |
| 12:00 PM - 12:30 PM | **Lunch** |
| 12:30 PM - 1:30 PM | **Lab 3 -  Building a Predictive Model with Watson Machine Learning** |
| 1:30 PM - 2:00 PM | **Next steps regarding POC and project timeline** |

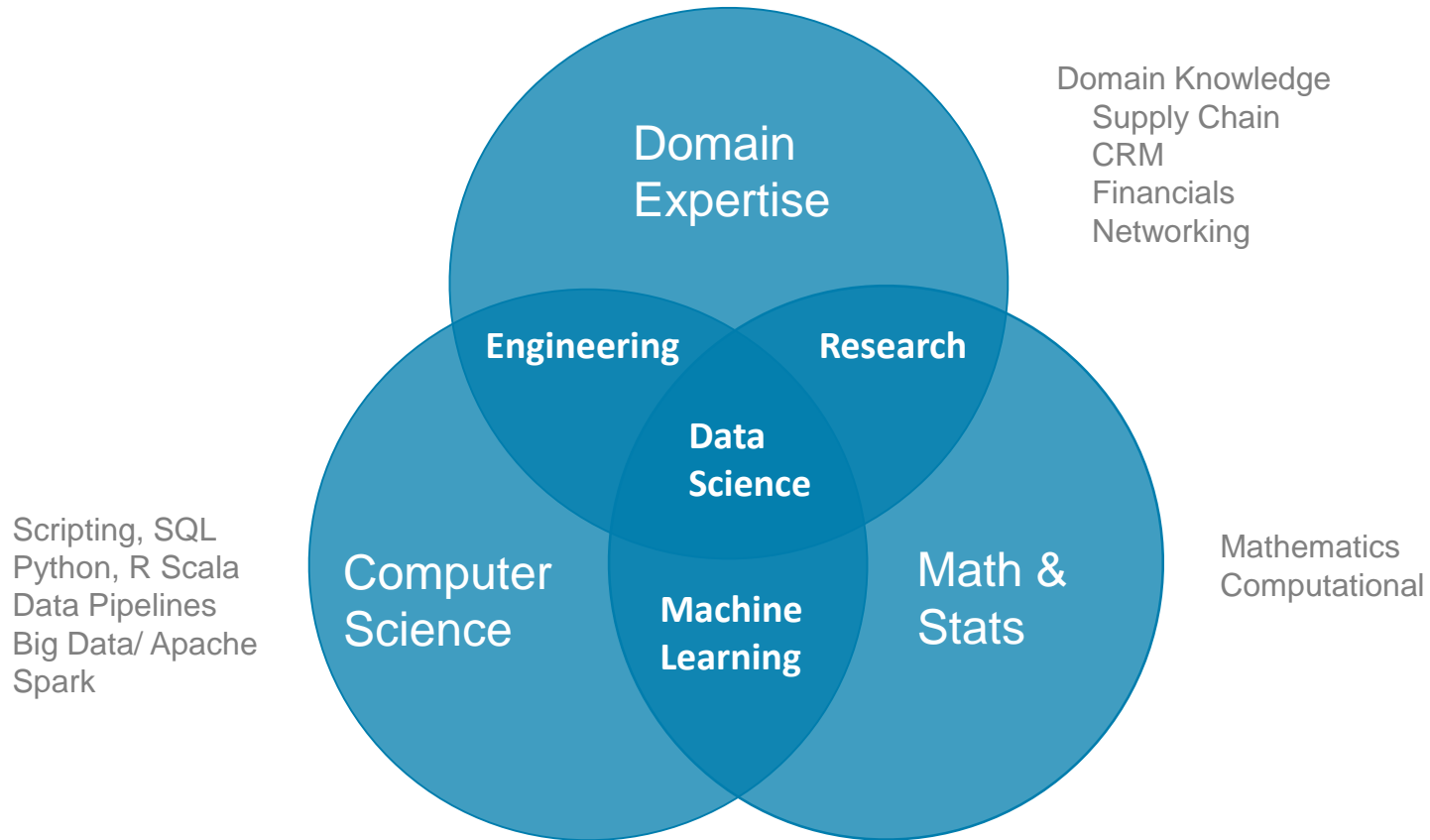# Participant Background

## Open Source

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

## IBM

- Bluemix
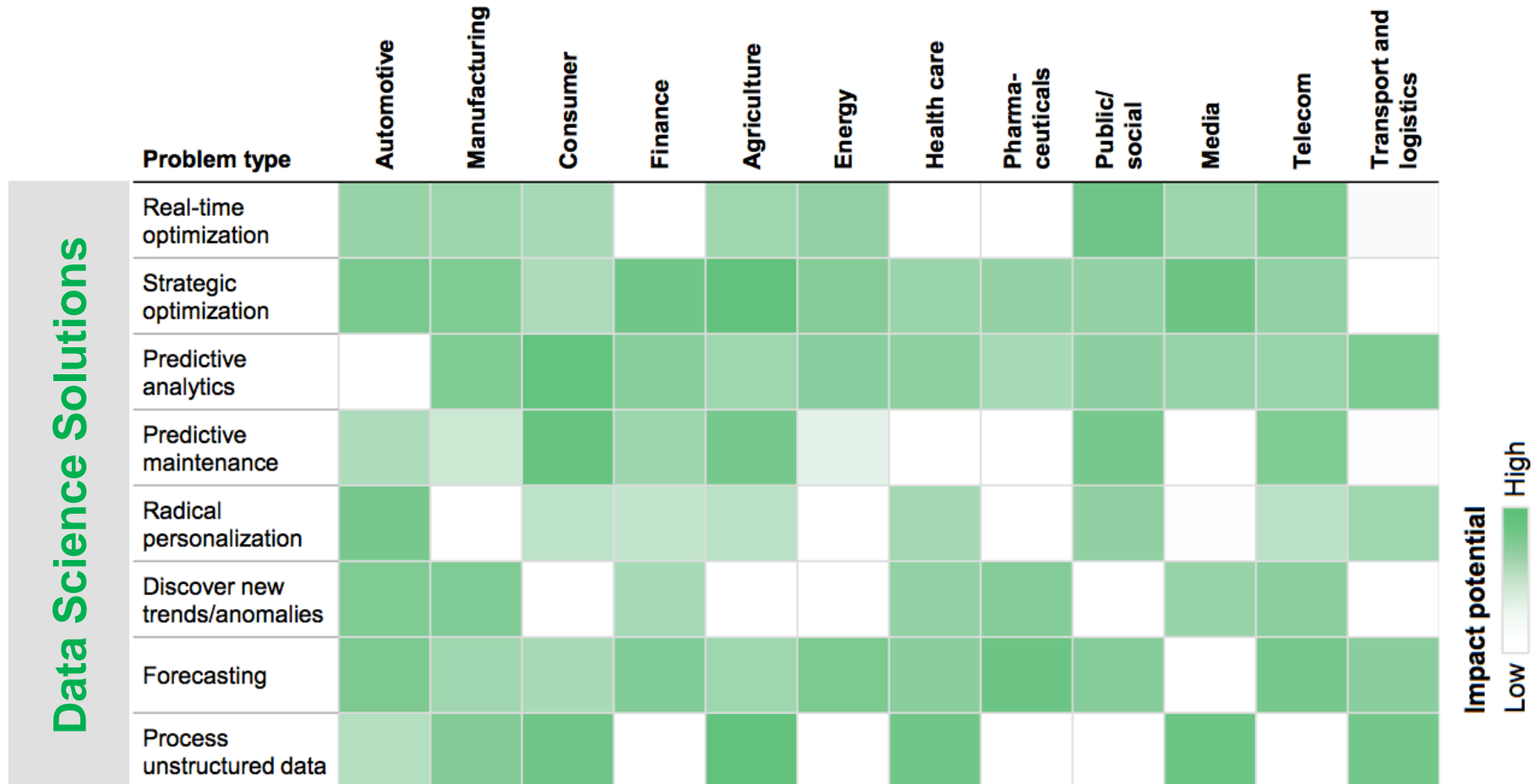- Data Science Experience
- Watson Machine Learning

# What is Data Science?



Domain Knowledge
Supply Chain
CRM
Financials
Networking

Domain Expertise

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

**Data Science Solutions**

| Problem type | Automotive | Manufacturing | Consumer | Finance | Agriculture | Energy | Health care | Pharma-ceuticals | Public/social | Media | Telecom | Transport and logistics |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Real-time optimization | | | | | | | | | | | | |
| Strategic optimization | | | | | | | | | | | | |
| Predictive analytics | | | | | | | | | | | | |
| Predictive maintenance | | | | | | | | | | | | |
| Radical personalization | | | | | | | | | | | | |
| Discover new trends/anomalies | | | | | | | | | | | | |
| Forecasting | | | | | | | | | | | | |
| Process unstructured data | | | | | | | | | | | | |

Impact potential — High / Low

SOURCE: McKinsey Global Institute analysis

5

# 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

- Self-service isn't a reality, if the data isn't secure
- 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.

*http://ibm.co/makedatasimple*

# IBM Watson Data Platform
## Experience New Ways To Put Data To Work

| Data Engineering | Data Science | Business Analysis | Application Development | Data Stewardship |
|---|---|---|---|---|

### Experiences
task-specific, collaborative

### Data and Analytics Services
comprehensive

Spark

open • intelligent • hybrid

# IBM Watson Data Platform
## Connects Users to Data and Analytics



**Data Engineering**

**Data Science**

**Business Analysis**

**Application Development**

**Data Stewardship**

*common data, pipelines and projects*

Find

Share

Collaborate

**Spark**
analytics operating system

**Data Sources**
- On-premises / cloud
- Structured / unstructured
- In-motion / at-rest
- Internal / external

box  Twitter  The Weather Company

*Iterate*

Discovery / Exploration
Machine learning
Model development

Analyze

Ingest

Deploy

Integration
Matching / Quality
Streaming

Persist

Reports / Dashboards
Applications
APIs

Hadoop
NoSQL / SQL
Object store

Govern

Data Assessment
Metadata / Policies

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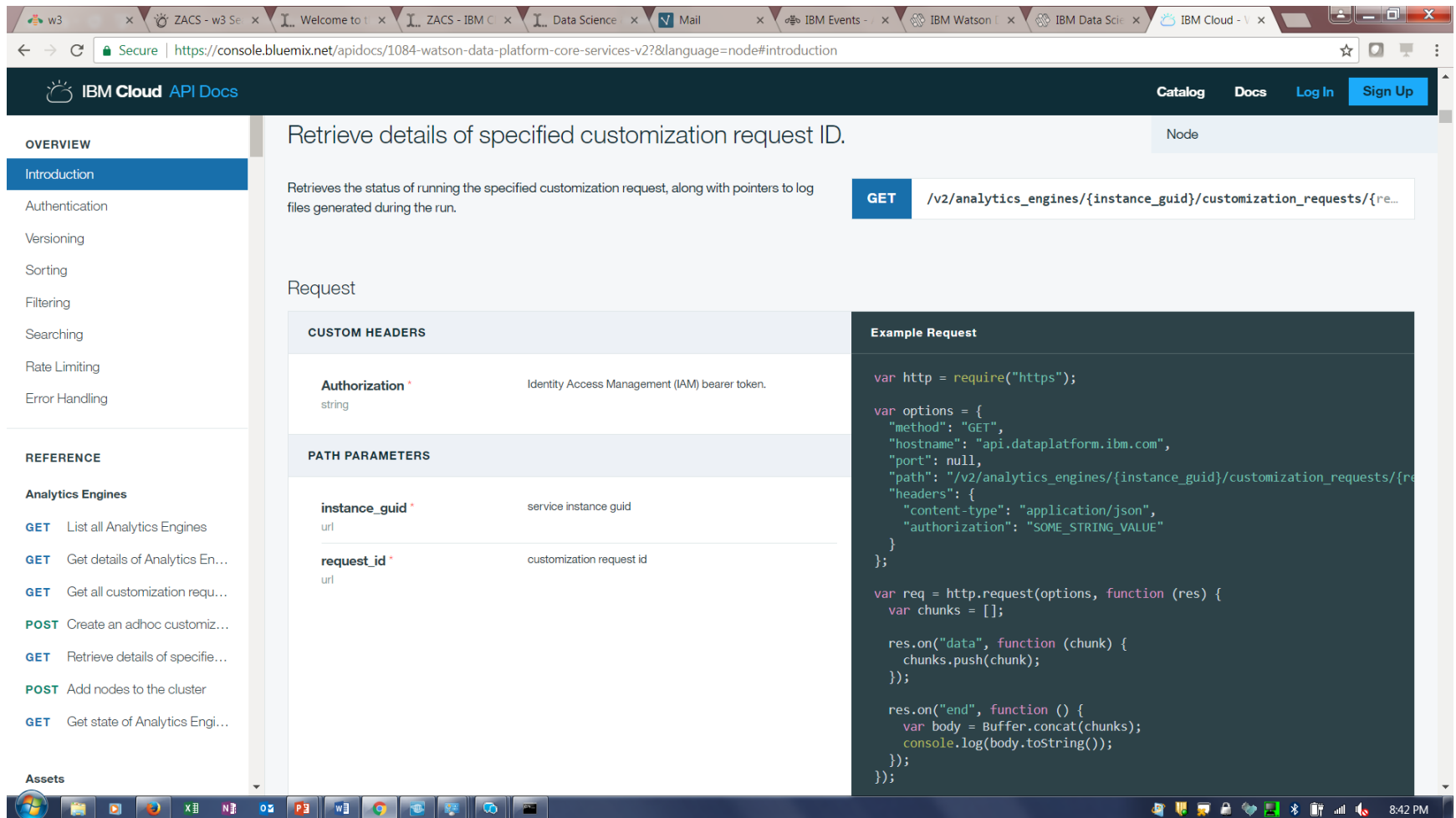The best of open source and IBM value-add to create state-of-the-art data products

## Collaborate
Data and Analytic assets are contained within projects which can be shared with other users.

# IBM Cloud PaaS

## *Rich Platform and Service APIs for your developers*

**Application Developer**

# 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…

# 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**

- Find tutorials and datasets
- Read articles and papers
- Connect with Data Scientists
- Share comments
- Copy and share notebooks

**Open Source**

- Code in Scala/Python/R/SQL
- Jupyter Notebooks
- RStudio IDE and Shiny
- Apache Spark
- Your favorite libraries

**IBM Added Value**

- IBM Machine Learning
- SPSS Modeler Canvas
- Prescriptive Analytics - DOcplexcloud
- Projects and Version Control
- Managed Spark Service

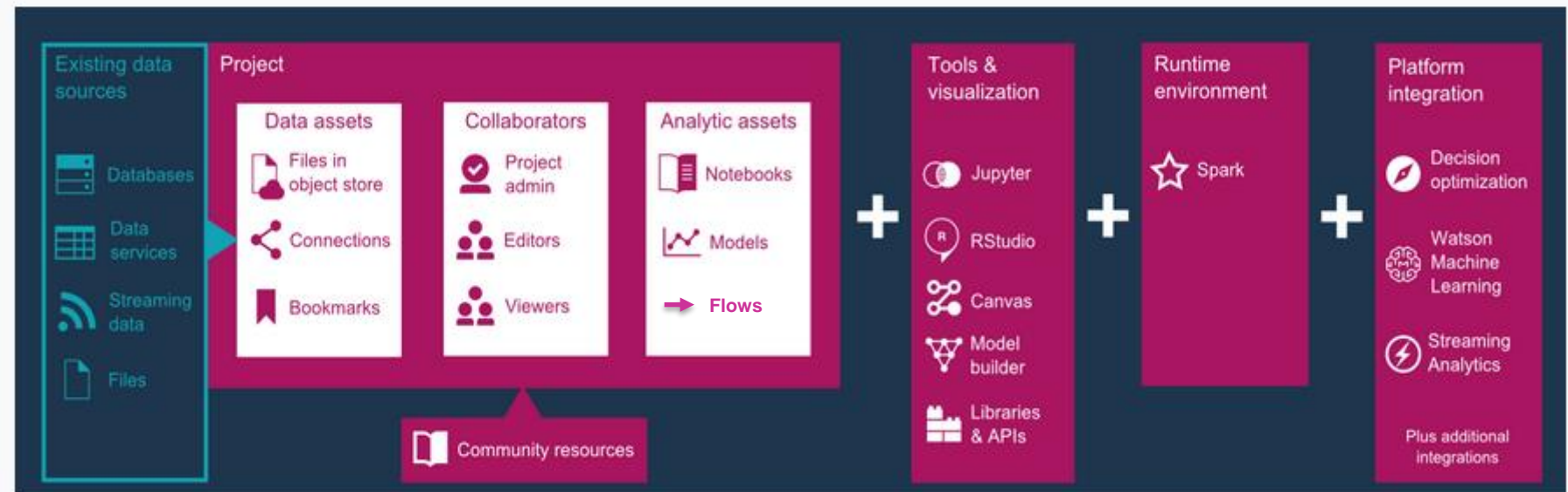Powered by IBM **Watson Data Platform**

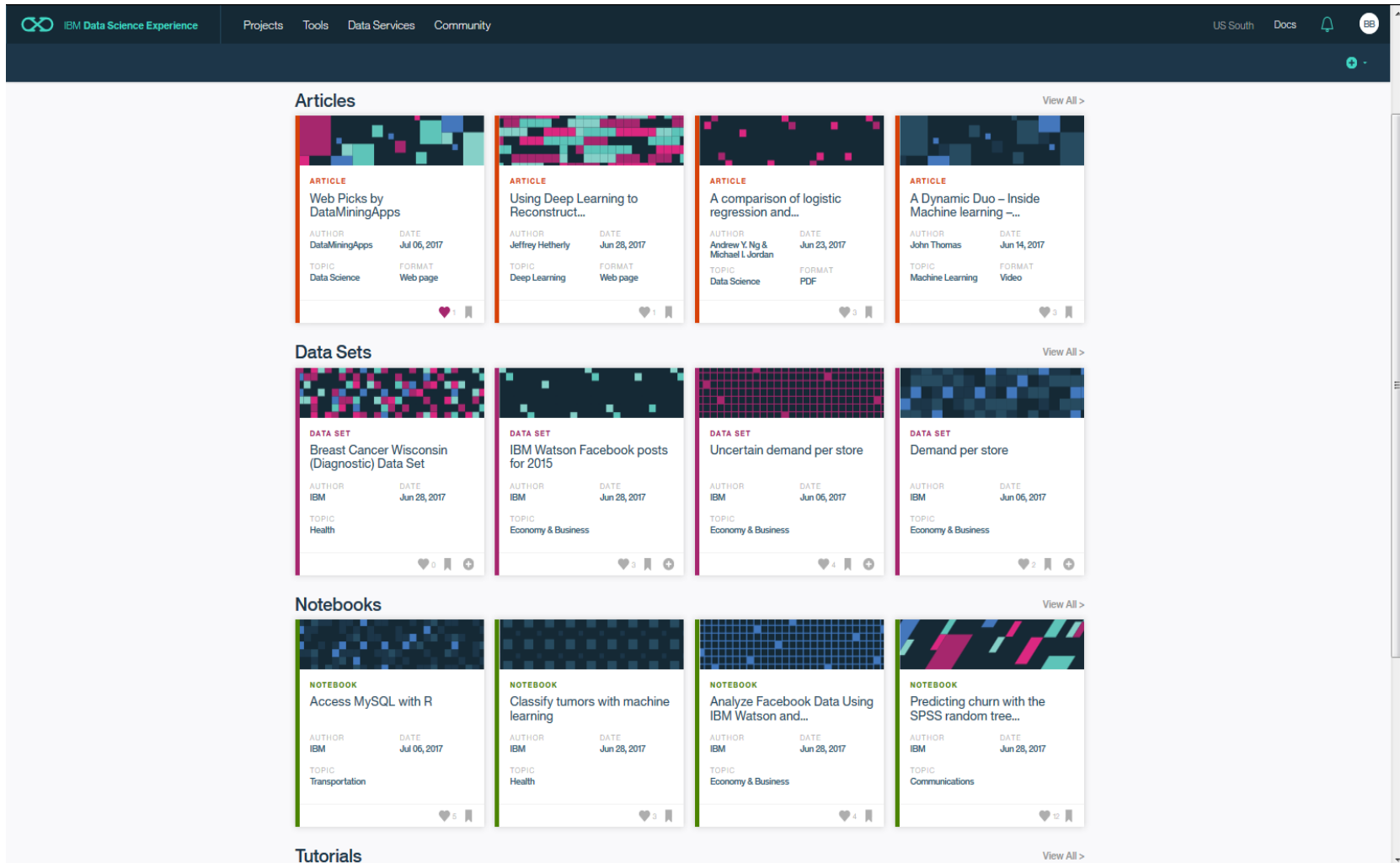# DSX Architecture



## DSX architecture

Last updated: **June 27, 2017**

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

≡  ∞  **Data Science** Experience  ⌄  Settings

Integrations

Profile     Services     Integrations

## GitHub Integration

Want to publish your notebooks on GitHub?
Before you can publish to GitHub, you need to create an access token. Visit GitHub personal access tokens, select repo scope and generate a token.

*Paste generated personal access token here*

40

Clear          Save

After the access token is saved, a GitHub repository can be connected to a project on the project's Settings page.

# 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

25

# 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

- **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 via
# on- premises and cloud Connectors

## IBM services in IBM Cloud

| | | | |
|---|---|---|---|
| IBM Informix | PostgreSQL on Compose | MySQL on Compose | Cloud Object Storage |
| IBM Db2 for i | IBM Cloudant | Cloud Object Storage (IaaS) | IBM Db2 on Cloud |
| Object Storage OpenStack Swift for IBM Cloud | IBM Db2 | IBM BigInsights HDFS | IBM Db2 Hosted |
| Object Storage OpenStack Swift (IaaS) | IBM PureData for Analytics | IBM Db2 for z/OS | IBM Db2 Warehouse on Cloud |

## Third-party services

| | | | |
|---|---|---|---|
| Cloudera Impala | Salesforce.com | Apache Hive | Amazon Redshift |
| Microsoft SQL Server | Sybase IQ | Sybase | Oracle |
| Amazon S3 | MySQL | Hortonworks HDFS | PostgreSQL |
| Pivotal Greenplum | Microsoft Azure SQL Database | | |

# 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 Machine Learning

# Watson 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

# Use SPSS Modeler to Visually Create ML Flows

- This 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 SPSS Modeler

# 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

# Labs

# Lab Overview

Use IBM's Data Science Experience (DSX) and IBM cloud services to create a working cloud-based application from start to finish. Participants will be led through a series of three labs. The three labs build upon one another so it is important that they are completed in order.

- Lab 1 - The first lab will begin with loading raw delimited data into DB2 Warehouse for Cloud and interacting with that data from a Jupyter notebook in DSX with python.

- Lab 2 - The second lab will guide participants in creating an R notebook and Shiny UI in DSX using RStudio.

- Lab 3 - The third lab will show how to use the Watson Machine Learning capability to create a machine learning model based on the supply chain data set. The machine learning model, deployed in the IBM Cloud, will be used to predict the severity level of each discrepancy based on action request characteristics.

# Lab 1

This lab will begin with loading raw delimited data into DB2 Warehouse for Cloud and interacting with that data from a Jupyter notebook in DSX with Python.

**Objectives:**

- Upon completing the lab, you will know how to:
  - Create a Jupyter IPython notebook from a URL
  - Establish a connection to DB2 Warehouse on Cloud
  - Use a dataframe to read and manipulate tables
  - Use Spark to explore and analyze the dataset
  - Write the modified dataset back to DB2 Warehouse on

# Lab 2

In this lab, you will learn some of the fundamentals of using RStudio and Shiny in DSX to work and interact with data in DB2 Warehouse and then create a fully operational "reactive" web application that you can enhance further.

## Objectives:

- Upon completing the lab, you will:
  - Create an RStudio project from a Git repository
  - Establish a connection to DB2 Warehouse
  - Query, explore and visualize data in an R notebook
  - Use ggplot2 to create bar plots of several columns in an R dataframe
  - Close the database connection
  - Leverage shiny to create and run a web application
  - Interact with the shiny web application by running it externally

# What is Machine Learning?

*"Computers that learn without being explicitly programmed"*
*"Using algorithms to understand patterns in data"*

Data

$$f(x) = \sum_{i=0}^{n} \alpha_i \, y_i x_i^T x + b$$

Algorithms

Predictions
& Insight

# Categories of Machine Learning

- **Supervised learning**
  - The program is "trained" on a pre-defined set of "training examples", which then facilitate its ability to reach an accurate conclusion when given new data
  - The algorithm is presented with example inputs and their outcomes (labels)
  - The goal is to learn a general rule that maps inputs to outputs

- **Unsupervised learning**
  - No labels are given to the learning algorithm, leaving it on its own to find structure (patterns and relationships) in its input

# Categories of Machine Learning

| Technique | Usage | Algorithms |
|---|---|---|
| Classification (or prediction) | • Used to predict group membership (e.g., will this employee leave?) or a number (e.g., how many widgets will I sell?) | • Decision Trees<br>• Logistic Regression<br>• Random Forests<br>• **Naïve Bayes**<br>• Linear Regression<br>• Lasso Regression<br>etc |
| Segmentation | • Used to classify data points into groups that are internally homogenous and externally heterogeneous.<br>• Identify cases that are unusual | • K-means<br>• Gaussian Mixture<br>• Latent Dirichlet allocation<br>etc |
| Association | • Used to find events that occur together or in a sequence (e.g., market basket) | • FP Growth |

# Training, testing, & validation sets

- **During the model development process, supervised learning techniques employ training and testing sets and sometimes a validation set.**
  - Historical data with known outcome
  - Data is randomly split into training, testing, and/or validation sets (mutually exclusive records)
- **Why?**
  - Training set
    - Build the model
    - Tune the parameters
  - Testing set
    - Assess model quality during training/tuning process
    - Avoid overfitting the model to the training set
  - Validation set
    - Estimate accuracy or error rate of model after tuning
    - Used to compare multiple models

# Spark ML

- **Spark ML is Spark's machine learning (ML) library**

- **Goal is to make machine learning scalable and easy**
  - No need to understand the detailed math!

- **Divides into two packages:**
  - spark.mllib contains the original API built on top of RDDs
  - spark.ml provides higher-level API built on top of DataFrames for constructing ML pipelines
  - A pipeline is a series of stages where each stage either transforms, or runs through a machine learning algorithm.

- **Using spark.ml is recommended because with DataFrames the API is more versatile and flexible**
  - spark.mllib will continue to be supported

# Spark ML Pipeline Terminology

Spark ML standardizes APIs for machine learning algorithms to make it easier to combine multiple algorithms into a single pipeline, or workflow

- **DataFrame**: Spark ML uses DataFrame from Spark SQL as an ML dataset, which can hold a variety of data types

- **Transformer**: A Transformer is an algorithm which can transform one DataFrame into another DataFrame

- **Estimator**: An Estimator is an algorithm which can be fit on a DataFrame to produce a Transformer

- **Pipeline**: A Pipeline chains multiple Transformers and Estimators together in a sequence to specify an ML workflow

- **Parameter**: All Transformers and Estimators share a common API for specifying parameters

# Lab 3 - Watson Machine Learning

In this lab, you will use IBM's Watson Machine Learning GUI to train, evaluate, and deploy a Watson Machine Learning model based on the modified supply chain dataset.

## Objectives:

- Upon completing the lab, you will:
  - Become familiar with the Watson Machine Learning GUI.
  - Train/Evaluate a machine learning model
  - Deploy a machine learning model.
  - Use the deployed machine learning model to make predictions.