# Introduction to Data Science and IBM's Data Science Experience

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

**Joel Patterson**
**Bernie Beekman**
**Davin Shearer**

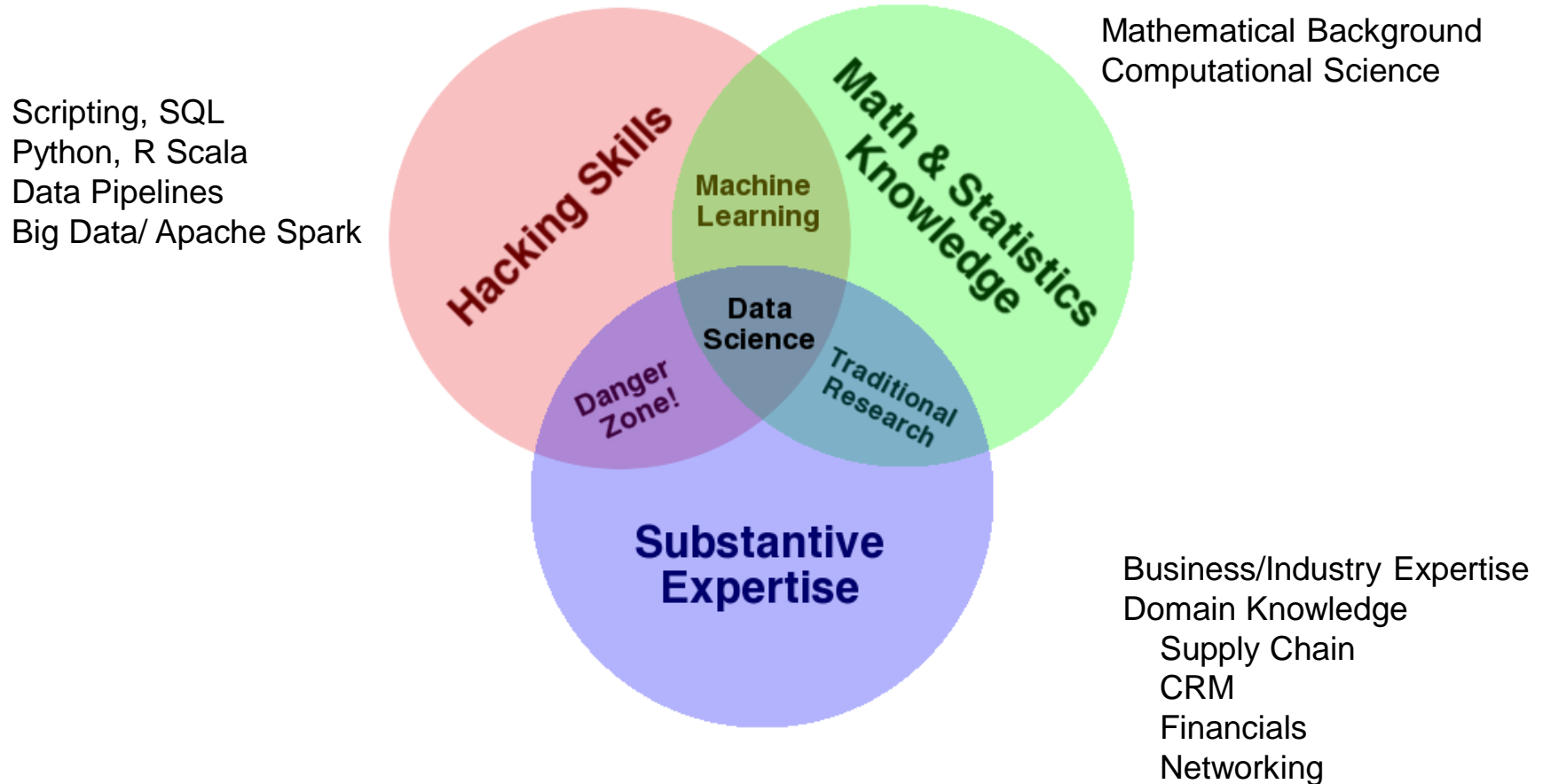# The digital age is changing the way we live, play, learn and work…

# What is the Data Scientist?

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

Mathematical Background
Computational Science



**Hacking Skills**

**Math & Statistics Knowledge**

Machine Learning

Data Science

Danger Zone!

Traditional Research

**Substantive Expertise**

Business/Industry Expertise
Domain Knowledge
   Supply Chain
   CRM
   Financials
   Networking

*Drew Conway's Data Science Venn Diagram*

# Google Trends – Data Science Languages



SPSS  SAS

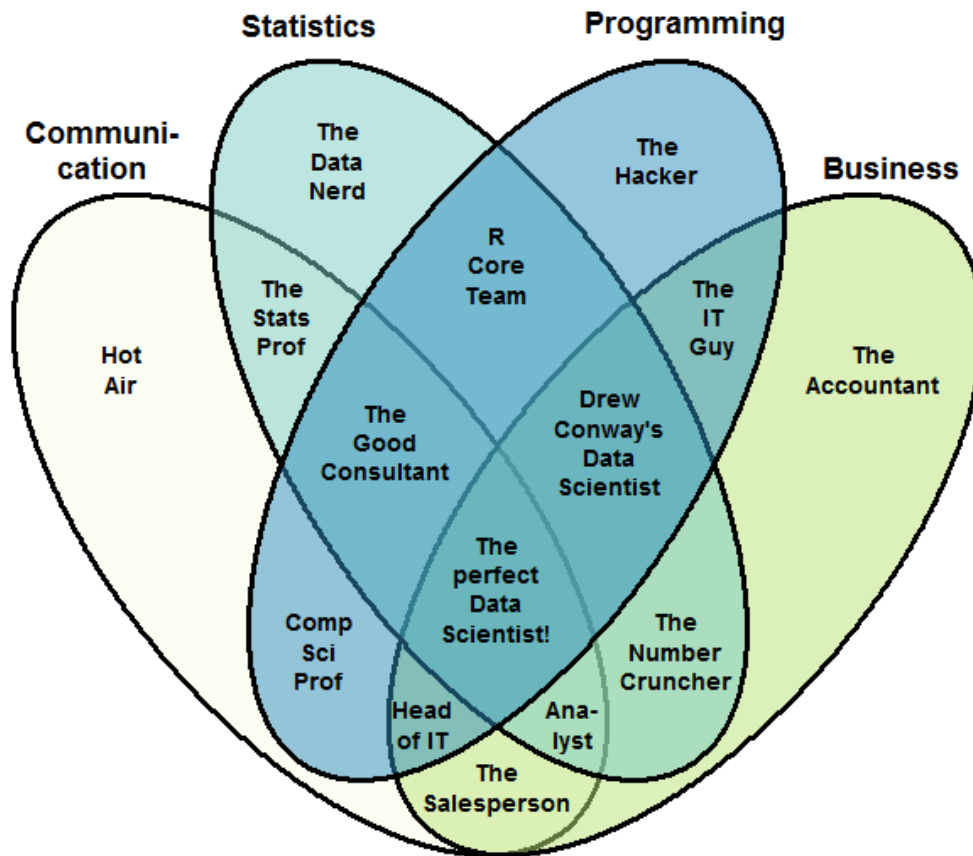Python  R  Scala

Trends in Google Searches (September 2nd 2016)

# The perfect Data Science Team
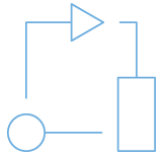
**The Data Scientist Venn Diagram**



Normally not all the skills are in one single person but rather in a data science team
In IBM Data Science Experience we include tools   to make the perfect Data Science Team
All in a collaborative, cloud environment that scales in demand

# IBM Watson Data Platform
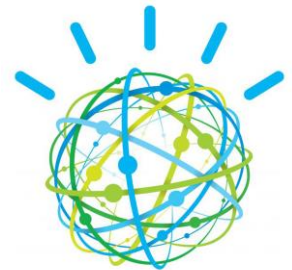
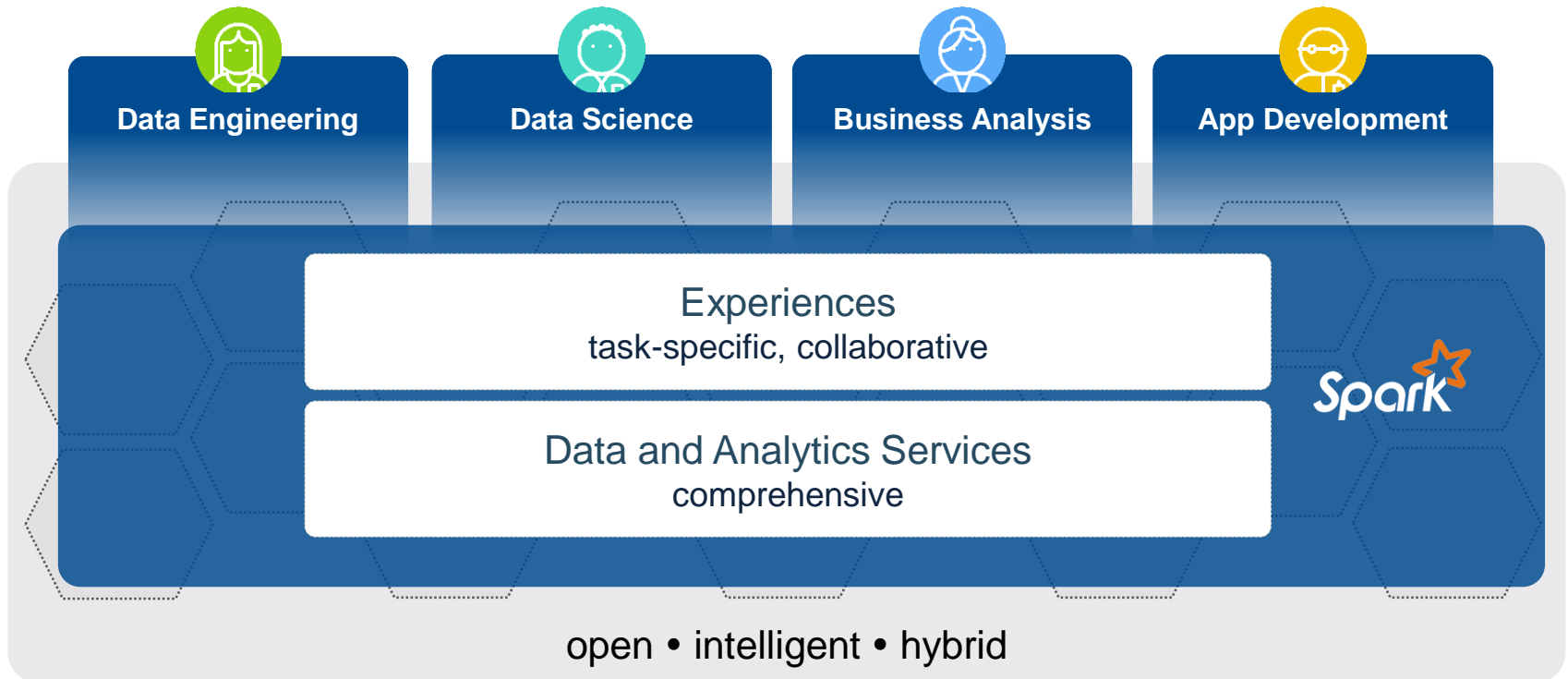## Mission: Make Data Simple and Accessible to All
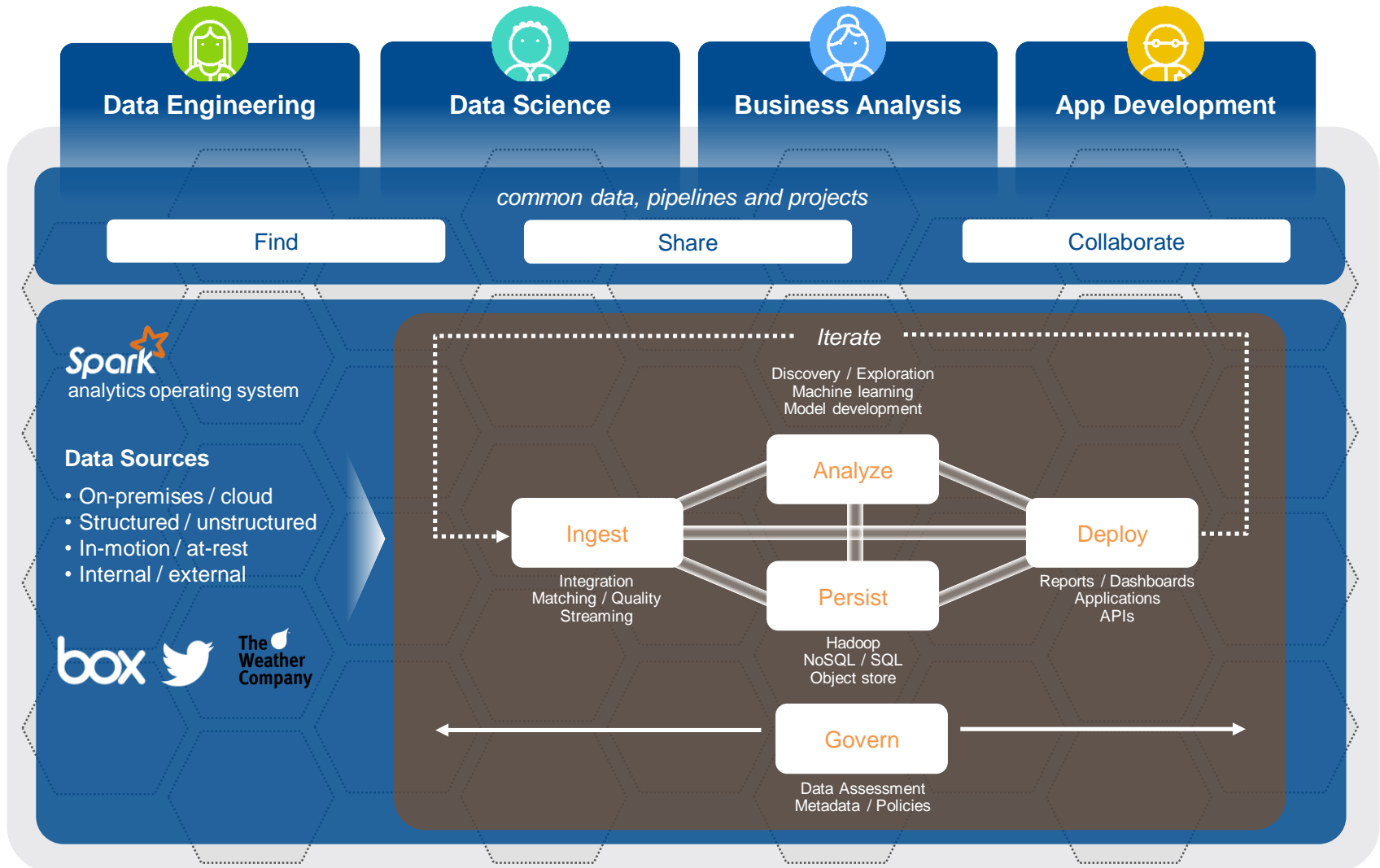
Platform.    Method.    Ecosystem.

*http://ibm.co/makedatasimple*

# IBM Watson Data Platform

## Experience New Ways To Put Data To Work

**Data Engineering**

**Data Science**

**Business Analysis**

**App Development**

### Experiences
task-specific, collaborative

### Data and Analytics Services
comprehensive

*Spark*

open • intelligent • hybrid

# Data Scientist Challenges

- **Rigid toolset**
  - Have to choose one and only one approach
  - Cannot easily connect all of the capabilities needed
  - Difficult to navigate between the various tools used

- **Fragmented and time consuming**
  - Using multiple disjointed environments
  - Separate on-ramp/community for each tool/environment
  - Does not have meta data or data lineage

- **Analytical Silo**
  - Difficult to maintain and version control project assets
  - Limited means of collaborating with team
  - Results are difficult to share

# Data Science Experience

Brings together popular Data Science **Open Source tools** with
IBM value-add functionalities coupled with **community and social** features

### 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

External URL: http://datascience.ibm.com

# 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**

* Closed beta

# DSX Architecture

# Community Cards provide in-context learning for users

# 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 DSX

# 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

# 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**

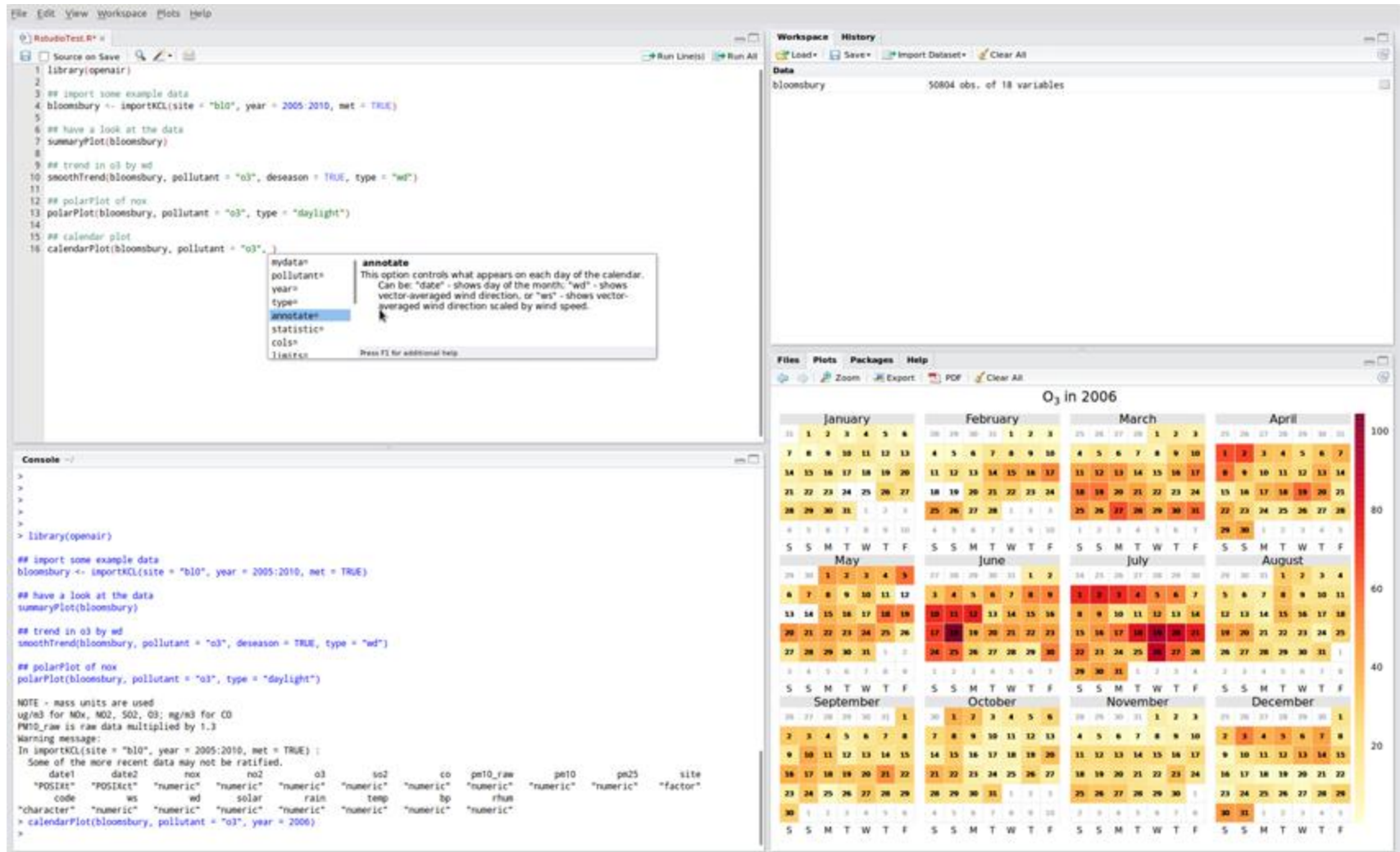| Cloud Sources | On-Premises Sources | Cloud Targets | On-Premises Targets |
|---|---|---|---|
| Amazon Redshift | Apache Hive | Amazon S3 | IBM DB2® LUW |
| Amazon S3 | Cloudera Impala | Bluemix Object Storage | IBM Pure Data for Analytics® |
| Apache Hive | IBM DB2® LUW | IBM Cloudant™ | Teradata |
| 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

# 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