# Breaking Down Data Science Life-Cycle.

IBM **Developer**

# Breaking Down Data Science Life Cycle.

## Episode #1

23-Oct

- Discovery & Gathering of Unstructured Data.
- Data Understanding.

## Episode #2

21-Nov

- Data Preparation.
- ICP for Data.

## Episode #3

26-Dec

- Machine Learning on Data.
- Modeling and Deployment.

# Data Discovery & Gathering

*Mashael AlMuhanna*
*Data Governance Specialist*

IBM **Developer**

What is Data ?

How is it Generated ?

Data refers to information that is machine-readable as opposed to human-readable.

# The importance of Data

Data driven HR.

Companies target

Data driven Strategies
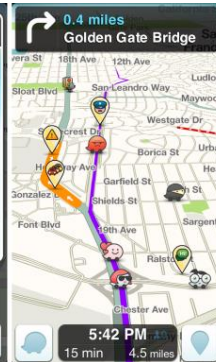
Relate to customers

# Data is Everywhere

Influencing What We Do



Netflix provides personalized recommendations

Waze provides a personalized driving experience

Uber delivers food that you like and is the right temperature

Self driving cars react to changing conditions

ALL based on DATA

# Data-driven cultures realize higher business returns

## Achieve Differentiation

**Manufacturing**
Predictive maintenance, production output & inventory

**Banking & Financial**
Reveal trading behavior, regulatory compliance

**Retail**
Dynamic pricing and predictive merchandising

**Healthcare**
Accuracy of diagnosis and regulatory compliance

**Telecom & Media**
Predictive customer Experience and loyalty

## Drive 6% Greater Productivity*

GC
Reduce Risk

CRO
Make
Money

CMO
Improve
Loyalty

CFO
Save Money

## Apply Technology 13:1* ROI

**Cloud**
data agility and efficiencies

**Mobile-IoT**
real-time and flexible data access

**Open Source**
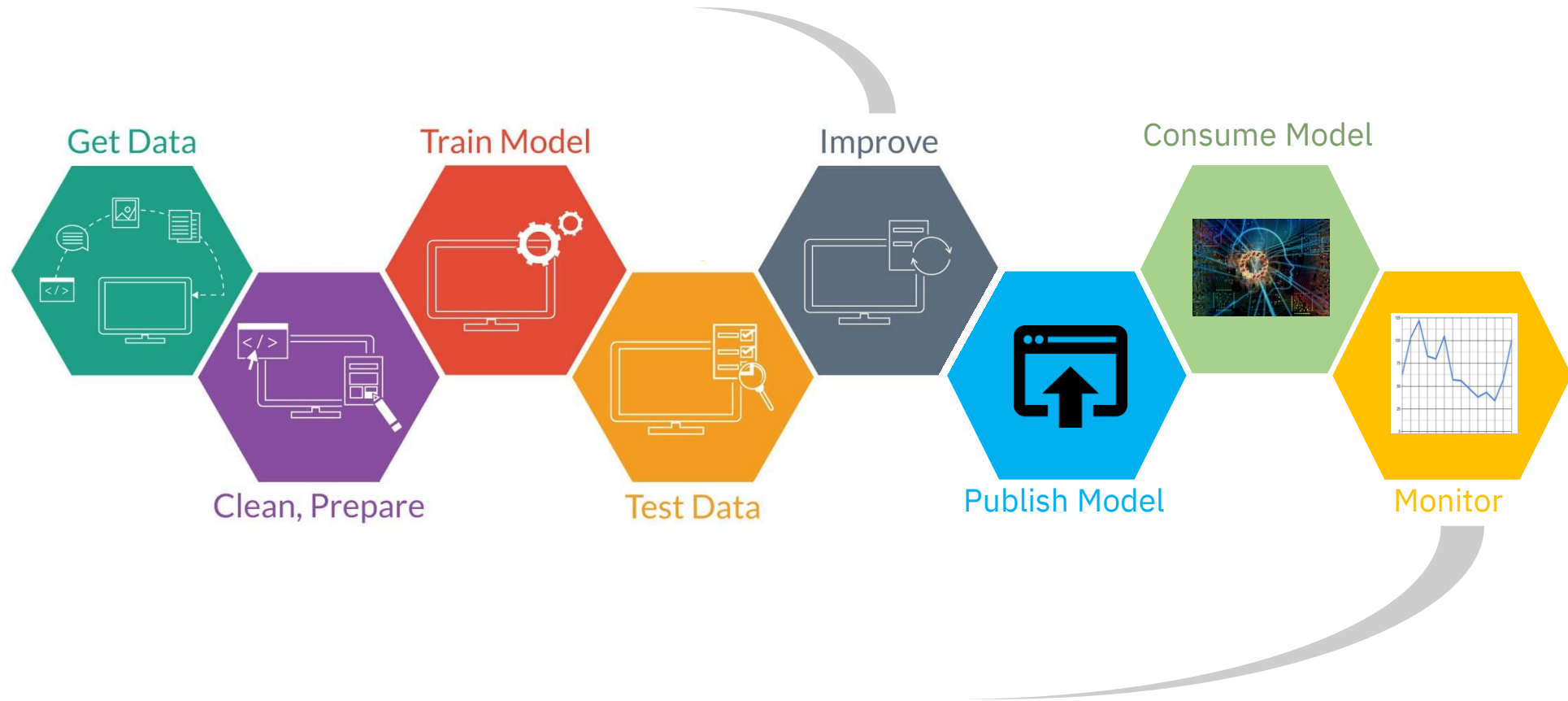speed innovation and data exploitation

**Artificial Intelligence**
scale discovery of hidden insights

* MIT Slone School of Management Study of 330 firms

Data-related challenges are hindering 96% of organizations from achieving AI

# What Real-World AI Actually Looks Like



Get Data

Clean, Prepare

Train Model

Test Data

Improve

Publish Model

Consume Model

Monitor

# We make data simple
# and accessible

# Discover & Search for Data, Refine & Act



**Select Data Sources**

Add source(s)

Index data in place

Infosets: group all volumes you want to search in

**Discover Your Data**

Visualizations show:
- Types of data stored
- Allocation by Date
- Allocation by Size

Discover where your oldest, biggest or least used data resides

Utilize overlays to highlight potential sensitive data

**Search & Refine**

Identify specific data based in your requirements and search criteries

Create own Filters based on Company Policies and Common sense to find datasets to act on

Create result data sets

**Report Findings**

Reports of result data sets

Notify data owner

**Act on Result Sets**

Manage in place
• Delete
• Move
• Export
• Retain

# Open Architecture



**Support for 100+ data sources and 450+ file types**

# Discover, Analyze & Act to Govern Information

- ## Discover
  - – Unstructured data across your enterprise

- ## Analyze
  - – Identify sensitive & critical information:
    - National ID
    - Passport numbers
    - Name
    - Addresses
    - Etc.
  - – Advanced search capabilities to help find country specific data

- ## Act
  - – Declare as records, delete, move…

# Governance Catalog

- Single interface for all governance activities
  - Establishes and promote common understanding for ALL enterprise users
  - Provides asset management and data lineage functionality

# Business Terms

## Assets

### Categories

- 📁 DiscoveryTerms
- 📁 Employee Benefits
- 📁 Global Life Glossary
  - 📁 Corporate Casualty
  - 📁 CrossDivision
  - 📁 Employee Be...
  - 📁 Individuals
- 📁 Management Pr...
  - 📁 Project Man...
    - 📁 Custome...
    - 📁 Employe...
    - 📁 Main dia...

**stewards**

**synonyms**

**deprecated**

**related terms**

### ▼ Terms (8)

1-8 of 8    ⏮ ◀ Page

- 📋 Base Annual Earnings
- 📋 Date of Birth
- 📋 Dependent Child
- 📋 Employee Date of Birth
- 📋 Employee Date of Employment
- 📋 Long Term Disability
- 📋 Occupational Accidental Death
- 📋 Short Term Disability

### MORGANHOST

- ⊟ 🖥 MORGANHOST
  - ⊟ 🟢 Database
    - ⊟ 🔲 Schema
      - ⊞ 🔲 IP
      - ⊟ 🔲 IP_ALT_NM
        - 📊 ALT_NM
        - 📊 EFF_DT
        - 📊 END_DT
        - 📊 IP_ID
        - 📊 IP_NM_TP_ID
        - 📊 LNG_ID
        - 📊 PPN_DT
        - 📊 PPN_TM
        - 📊 SRC_STM_ID

**tables**

**columns**

**reports**

33%  22%  34%  11%

**processes**

# Data Understanding

*Hissah AlMuneef*
*Developer Advocate*

**IBM Developer**

Data exploration, data cleaning, feature engineering, pre-processing, etc...

80

**20** Model building

# Data Understanding

**Task #1**

**Data Description**

**Task #2**

**Data Exploration**

# Data Description

Date Description contain:

- Source of Data

- The Number of Cases

- The Number of Fields

- Description of Fields

# Example for SaudiViz...

| | |
|---|---|
| airlines | *Airline names.* |

**Description**

Look up airline names from their carrier codes.

**Usage**

airlines

**Format**

Data frame with columns

**carrier** Two letter abbreviation

**name** Full name

**Source**

https://www.transtats.bts.gov/DL_SelectFields.asp?Table_ID=236

# Data Exploration

# Date Exploration Steps

1. Identify Data Types

2. Explore Each Variable

3. Find Correlation

# Step 1

Identify data types

IBM **Developer**

# Data Types



Other examples?

IBM **Developer**

# Identify data types

| Numerical | | Categorical |
|---|---|---|
| seats | speed | engine |
| 55 | NA | Turbo-fan |
| 182 | NA | Turbo-fan |
| 182 | NA | Turbo-fan |
| 182 | NA | Turbo-fan |
| 55 | NA | Turbo-fan |
| 182 | NA | Turbo-fan |
| 182 | NA | Turbo-fan |
| 182 | NA | Turbo-fan |

IBM **Developer**

# Step 2

Explore Variables One by One

IBM **Developer**

# Numerical Variable

| Central Tendency | Measure of Dispersion | Visualization Methods |
|---|---|---|
| Mean | Range | Histogram |
| Median | Quartile | Box Plot |
| Mode | IQR | |
| Min | Variance | |
| Max | Standard Deviation | |
| | Skewness and Kurtosis | |

# A. Central Tendency

- **Mode:** most frequent measure

- **Median:** mid-point of an array of measures

- **Mean:** arithmetic average (Sum/N)

IBM **Developer**

# B. Variance and Standard Deviation

Standard deviation is the square
root of the Variance

$$SD = \sqrt{\frac{\sum |x - \bar{x}|^2}{n}}$$

| x1 | 1 | 1 | 1 |
|----|---|------------|---|
| x2 | 1 | 2 | 2 |
| x3 | 1 | 2 | 3 |
| std | 0 | 0.57735027 | 1 |

**IBM Developer**

# C. Histogram & Boxplot

# C. Histogram & Boxplot

# Categorical Variable

| Students using each interface | Count | Percent |
|---|---|---|
| Interface 1 | 17 | 43.6% |
| Interface 2 | 4 | 10.3% |
| Interface 3 | 8 | 20.5% |
| Interface 4 | 10 | 25.6% |
| Total | 39 | 100% |



**IBM Developer**

# A. Numerical Summaries

- Count
- Count%
- Mode

| Students using each interface | Count | Percent |
|---|---|---|
| Interface 1 | 17 | 43.6% |
| Interface 2 | 4 | 10.3% |
| Interface 3 | 8 | 20.5% |
| Interface 4 | 10 | 25.6% |
| Total | 39 | 100% |

# B. Bar plot & Pie plot



What's your favorite ice cream flavor?

# Step 3

Find association between variables

# IBM **Developer**

# Associations between numerical variables

# Associations between categorical variables

# Associations between numerical-categorical variables

# Watson Studio

Exploration

# Watson Studio

# Watson Studio Dashboard

# Overview Page

My Projects / my-new-project

Add to project ▾

**Overview**    Assets    Environments    Bookmarks    Deployments    Collaborators    Settings

## my-new-project
Last Updated: May 03 2018

0
Assets

0
Bookmarks

1
Collaborators

**Date created**
May 03 2018

**Description**
No description available

**Storage**

0% of 25 GB used

**Collaborators**                    View all (1)

**SM**    Steve Martinelli
            Admin

**Bookmarks**                        View all (0)

## Recent activity

Alerts related to this project will show here
when the project is active.

**IBM Developer**

# Assets Page

# Community



IBM Developer

# Connections



IBM Developer

In Watson Studio, after you set up a project and add data to it, you can start analyzing and visualizing your data:

With Code
Or Without Code

# Jupyter Notebook

File   Edit   View   Insert   Cell   Kernel   Help

Format   Markdown

## Create a bar graph to visualize patterns

Create a bar graph to show the total number of collisions by borough:

```
In [13]:   borough = collisions_df.groupBy('BOROUGH').count().sort('count').toPandas()
           colors = ['g','0.75','y','k','b','r']
           plt.barh(range(6),borough.sort_values(by='count', ascending=True)['count'], color=colors)
           plt.xlabel('Collisions')
           plt.ylabel('Borough')
           plt.title('Total Number of Collisions by Borough', size=15)
           plt.yticks(range(6), borough['BOROUGH'])
           plt.show()
```



Total Number of Collisions by Borough

# RStudio



IBM **Developer**

# Dashboard



IBM **Developer**

# Titanic Dataset

The sinking of the Titanic is one of the most famous shipwrecks in history. The Titanic sank after colliding with an iceberg, killing 1502 out of 2224 passengers and crew.

This sensational tragedy shocked the international community and led to better safety regulations for ships.

# Data Description

**Description:**

Information about the passenger of titanic, 891 case.

**Format**

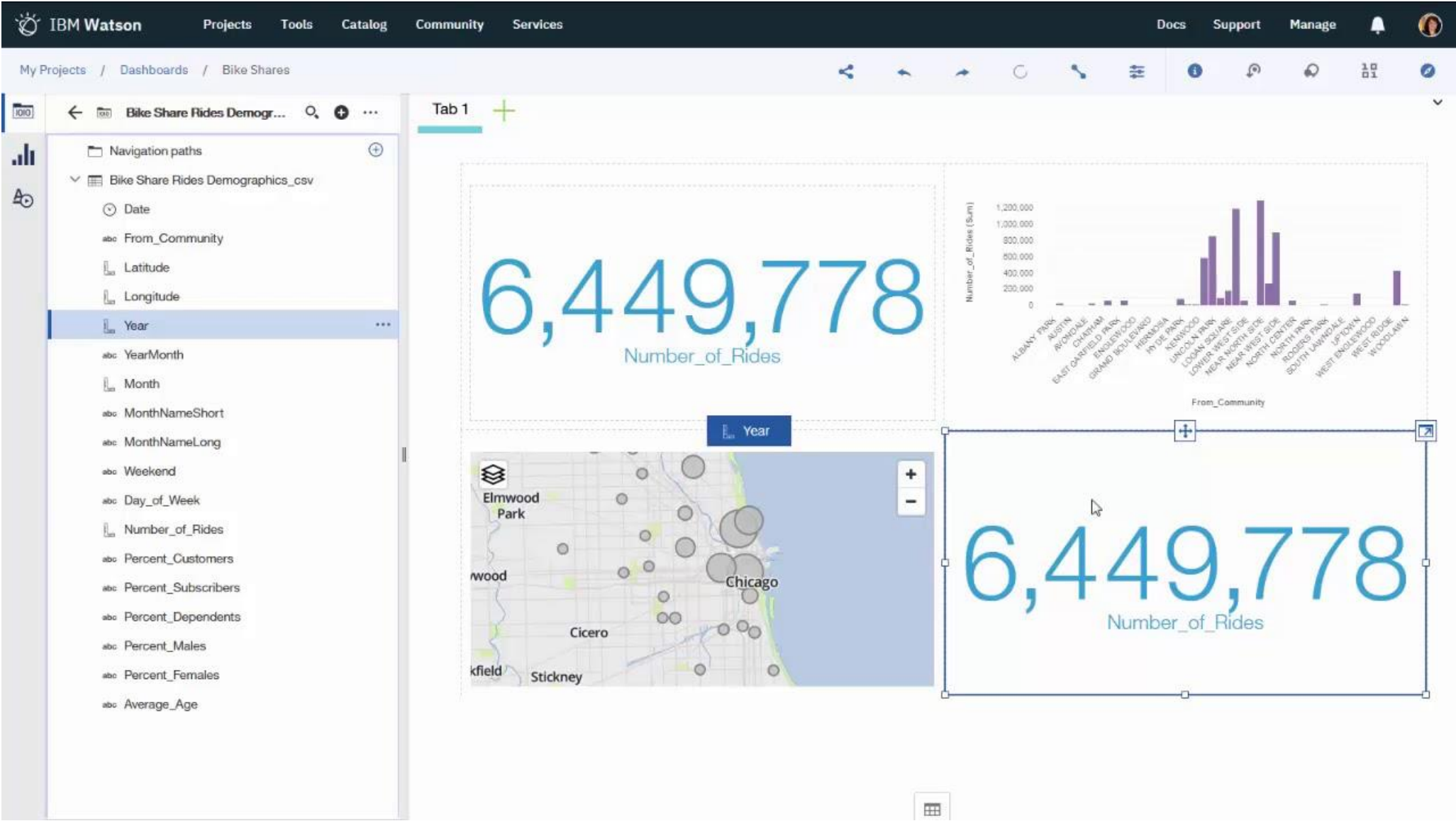| | |
|---|---|
| **survival** | Survival  (0 = No; 1 = Yes) |
| **pclass** | Passenger Class  (1 = 1st; 2 = 2nd; 3 = 3rd) |
| **name** | Name of passenger |
| **Gender** | Gender |
| **age** | Age |
| **sibsp** | Number of Siblings/Spouses Aboard |
| **parch** | Number of Parents/Children Aboard |
| **ticket** | Ticket Number |
| **fare** | Passenger Fare |
| **cabin** | Cabin |
| **embarked** | Port of Embarkation  (C = Cherbourg; Q = Queenstown; S = Southampton) |

**Source**　　　https://www.kaggle.com/c/titanic/data

IBM **Developer**

# Please, Sign Up for IBM Cloud (US)

https://ibm.biz/BdYpAP

# GitHub

https://github.com/DevExCodeHub/DataScienceSeries-Ep1

Wi-Fi Password: makeithappen

**IBM Developer**