# Workshop Outline

- Introduction Data Science
- Predicting Employee Churn
- Preparing SQL Server Environment
- What is SQL?
- SQL Server Management Studio
- Preparing to write SQL
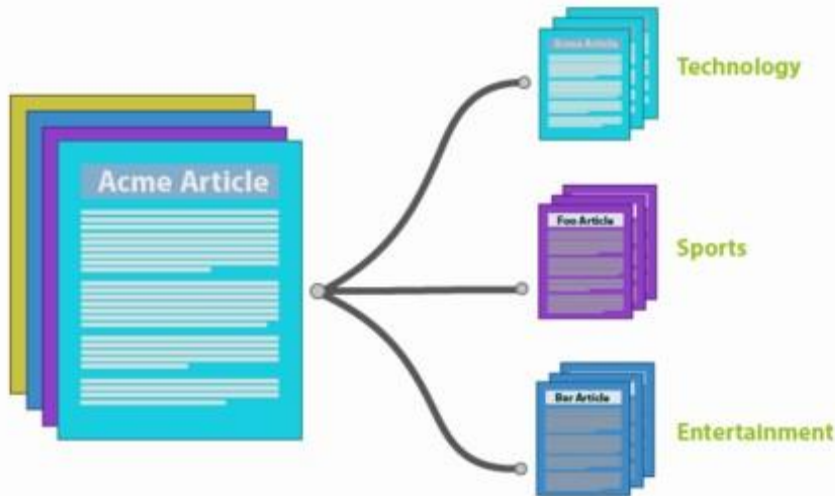- Machine Learning Services in SQL
- Building AI App

| | A | B | Members | Annual Fee |
|---|---|---|---|---|
| 1 | Affiliate | State | 205 | 50 |
| 2 | Norfolk | VA | 65 | 35 |
| 3 | Houston | TX | 657 | 75 |
| 4 | Manhattan | NY | 336 | 60 |
| 5 | Albany | NY | 453 | 50 |
| 6 | Washington | DC | 432 | 50 |
| 7 | Richmond | VA | 77 | 25 |
| 8 | Memphis | TN | 578 | 70 |
| 9 | Brooklyn | NY | 153 | 65 |
| 10 | Boston | MA | 32 | 65 |
| 11 | Waltham | MA | 43 | 35 |
| 12 | Schenectady | NY | 235 | 85 |
| 13 | Newark | NJ | 68 | 75 |
| 14 | Morristown | NJ | | |
| 15 | | | | |

# Types of Machine Learning

- Supervised learning
  - Output labels are known
  - Learn the ML model that produces the prediction closest to the output

- Unsupervised Learning
  - Output labels are not known
  - Divides data into clusters of similar data points

- Reinforcement Learning
  - An agent interacts with an environment and performs action
  - Learns through experience (reward mechanism)

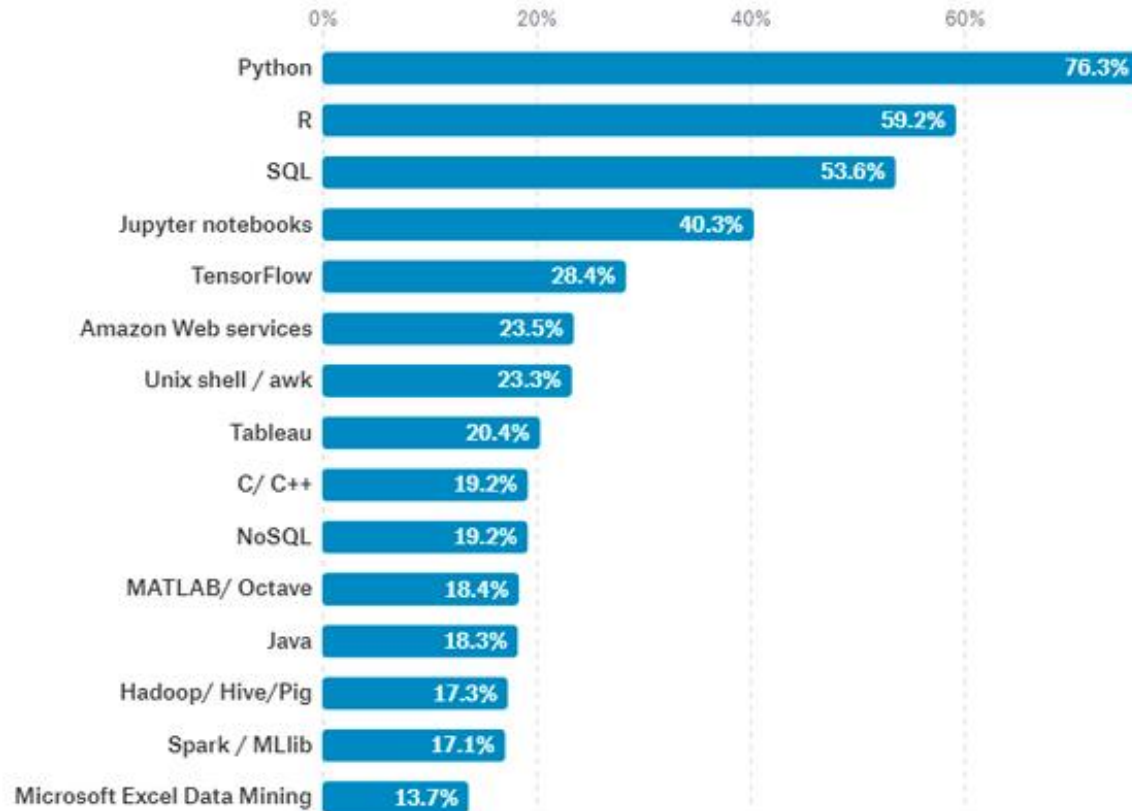DEEP REINFORCEMENT LEARNING
IN PACMAN

TYCHO VAN DER OUDERAA

# Jupyter Notebook

- Open-source web application

- To create and share documents that contain live code, equations, visualizations and narrative text.

- Uses include
  - data cleaning and transformation
  - numerical simulation
  - statistical modeling
  - data visualization
  - machine learning
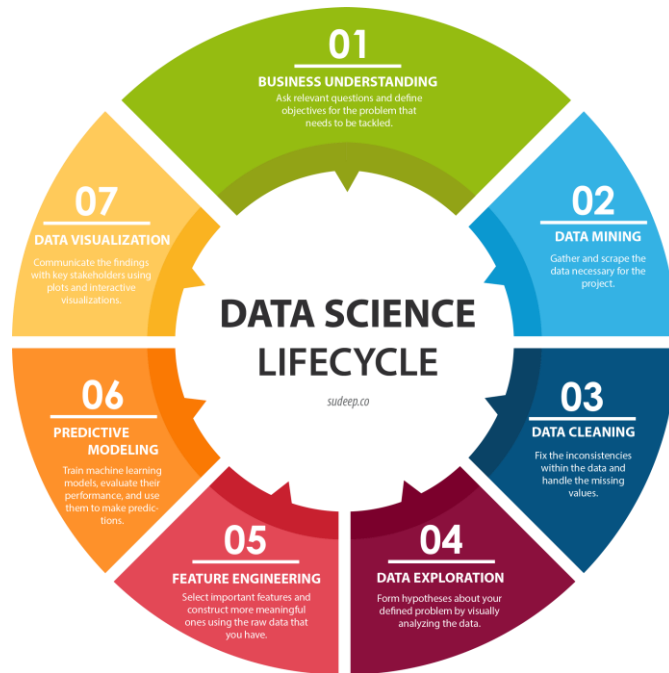  - and much more

# Top Data Science Technologies

| Technology | Percentage |
|---|---|
| Python | 76.3% |
| R | 59.2% |
| SQL | 53.6% |
| Jupyter notebooks | 40.3% |
| TensorFlow | 28.4% |
| Amazon Web services | 23.5% |
| Unix shell / awk | 23.3% |
| Tableau | 20.4% |
| C/ C++ | 19.2% |
| NoSQL | 19.2% |
| MATLAB/ Octave | 18.4% |
| Java | 18.3% |
| Hadoop/ Hive/Pig | 17.3% |
| Spark / MLlib | 17.1% |
| Microsoft Excel Data Mining | 13.7% |

## Data Format

### Structured

01234
56789

### Unstructured

**Data Source**

**Internal**

**Human-Generated**
- Survey ratings
- Aptitude testing

**Machine-Generated**
- Web metrics from Web logs
- Product purchase from sales Records
- Process control measures

**Human-Generated**
- Emails, letters, text messages
- Audio transcripts
- Customer comments
- Voicemails
- Corporate video/communications
- Pictures, illustrations
- Employee reviews

**External**

**Human-Generated**
- Number of Retweets, Facebook likes, Google Plus +1s
- Ratings on Yelp
- Patient ratings ratings

**Machine-Generated**
- GPS for tweets
- Time of tweet/updates/postings

**Human-Generated**
- Content of social media updates
- Comments in online forums
- Comments on Yelp
- Video reviews
- Pinterest images
- Surveillance video

# Data Science Life Cycle

**PRECISION, RECALL AND F1**
Uses **POSITIVES** & **NEGATIVE** to measure a model's
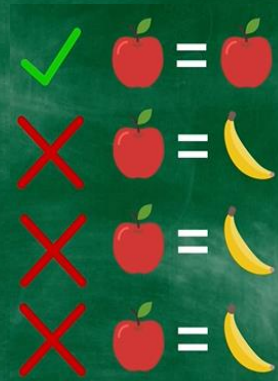**ACCURACY** when making predictions

**PRECISION**
**RECALL**
**F1**



**ACCURACY**

\+

\+

\-

\-

# Download Scripts

**http://github.com/BeirutAI/ML_in_SQL**
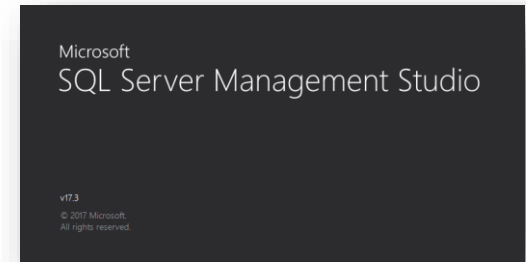
# What are databases?

- Holds data
- Organizes data
- Retrieve/Search data through DBMS

*A usually large collection of data organized especially for rapid search and retrieval.*

# SQL Server Management Studio (SSMS)

- SQL Server Management Studio (SSMS) is a powerful graphical DB management tool
  - Administrate databases (create, modify, backup / restore DB)
  - Create and modify Entity Relationship (E/R) diagrams
  - View / modify table data and other DB objects
  - Execute SQL queries
  - Free and easy to use tool
  - Works with all SQL Server versions
  - And much more

Microsoft
SQL Server Management Studio

v17.3
© 2017 Microsoft.
All rights reserved.

**Structured**: database schema

- Relational database

**Semi-structured**

```
{ "key": "value"}
```

- JSON

**Unstructured**: schemaless, more like files

- Videos, photos

# SQL and NoSQL

## SQL

- Tables
- Database schema
- Relational databases

## NoSQL

- Non-relational databases
- Structured or unstructured
- Key-value stores (e.g. caching)
- Document DB (e.g. JSON objects)

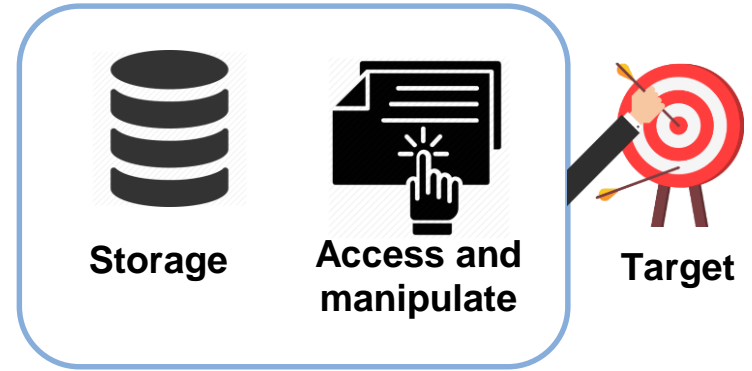**Structured Query Language**

**Storage**

**Access and manipulate**

**Target**

**Developed by**

**Like All relational database management system**

**Easy to learn**

**Result**

| ID | EmployeeName |
|----|--------------|
| 1 | Guy Gilbert |
| 2 | Kevin Brown |
| 3 | Roberto Tamburello |

```
Select
    *
    From DimEmployee
```

**Several SQL dialects exist**

```
Select
    ID,
    EmployeeName
    From DimEmployee
```

**Basic command Vocabulary less than 100 words**

# SQL

CREATE    READ    UPDATE    DELETE

# C R U D

- Open database
- Think about use case
- Define business problems
- Extract , Transform, Load (ETL)
- Data wrangling
- Production scripts

# Descriptive Data

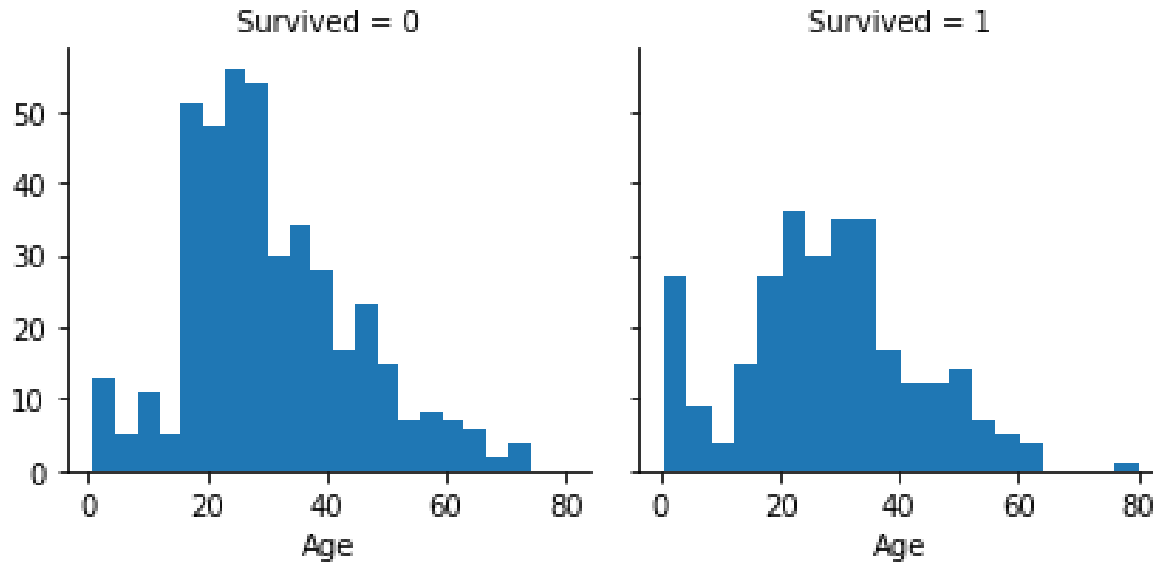| | PassengerId | Survived | Pclass | Age | SibSp | Parch | Fare |
|---|---|---|---|---|---|---|---|
| count | 891.000000 | 891.000000 | 891.000000 | 714.000000 | 891.000000 | 891.000000 | 891.000000 |
| mean | 446.000000 | 0.383838 | 2.308642 | 29.699118 | 0.523008 | 0.381594 | 32.204208 |
| std | 257.353842 | 0.486592 | 0.836071 | 14.526497 | 1.102743 | 0.806057 | 49.693429 |
| min | 1.000000 | 0.000000 | 1.000000 | 0.420000 | 0.000000 | 0.000000 | 0.000000 |
| 25% | 223.500000 | 0.000000 | 2.000000 | 20.125000 | 0.000000 | 0.000000 | 7.910400 |
| 50% | 446.000000 | 0.000000 | 3.000000 | 28.000000 | 0.000000 | 0.000000 | 14.454200 |
| 75% | 668.500000 | 1.000000 | 3.000000 | 38.000000 | 1.000000 | 0.000000 | 31.000000 |
| max | 891.000000 | 1.000000 | 3.000000 | 80.000000 | 8.000000 | 6.000000 | 512.329200 |

- Passengers
- Fares variance (8\$ → 512\$)
- Males 65%
- Cabin

| | Name | Sex | Ticket | Cabin | Embarked |
|---|---|---|---|---|---|
| count | 891 | 891 | 891 | 204 | 889 |
| unique | 891 | 2 | 681 | 147 | 3 |
| top | Sharp, Mr. Percival James R | male | CA. 2343 | B96 B98 | S |
| freq | 1 | 577 | 7 | 4 | 644 |

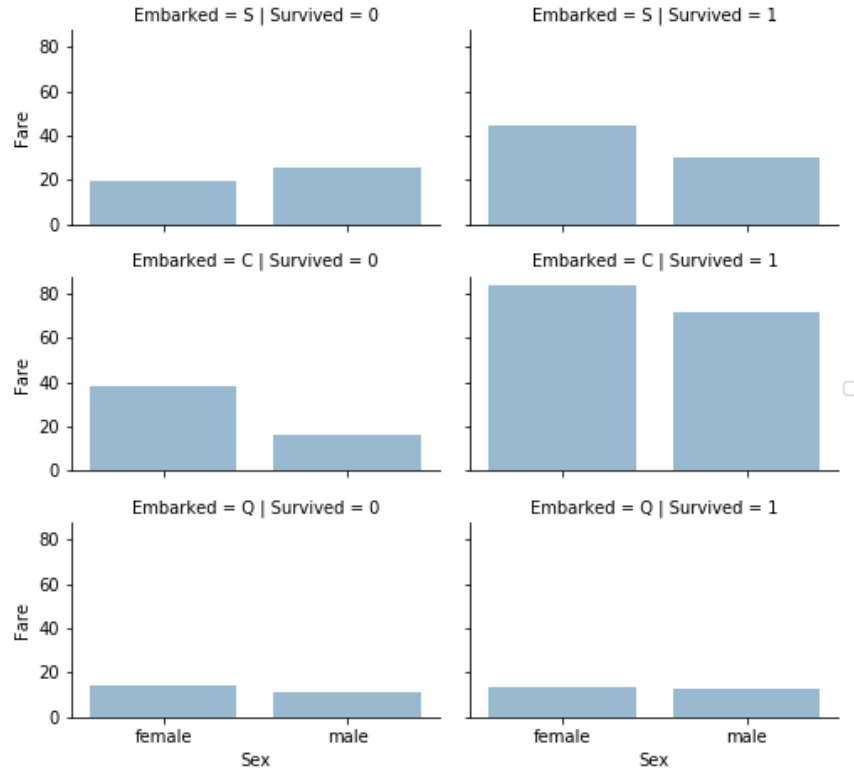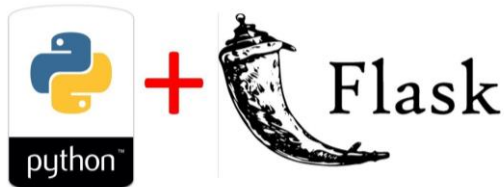# Analyze by visualizing data

# Machine Learning Services

**Rest API using Flask**



**Why Python in SQL**



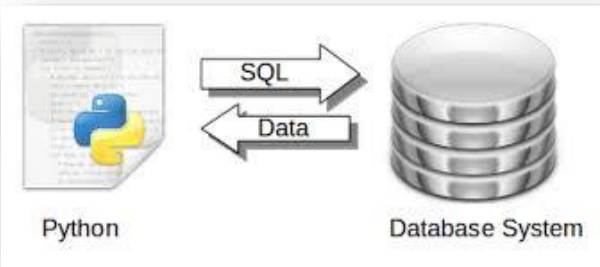**Is it SQL in Python ?**
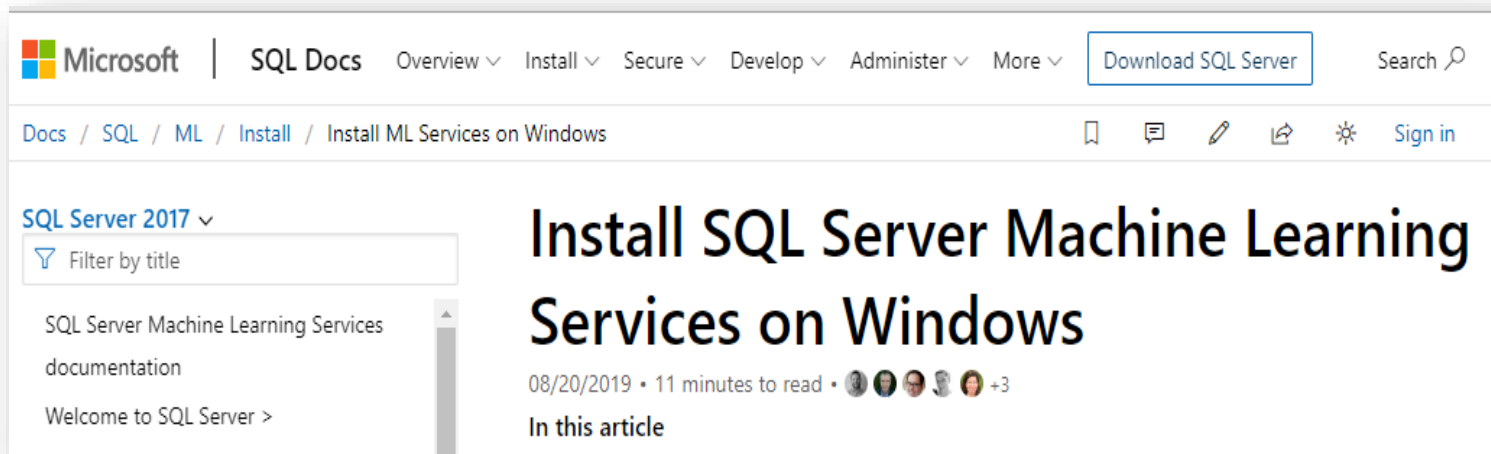


**I can use SQLAlchemy ?**



**Run external script**

# Using Python Inside SQL

- Eliminate data movement
- Easily operate Python code inside SQL
- Achieve enterprise grade performance and scale (much faster mechanism than ODBC)

# Configuration

- To enable SQL Instance to run Python scripts:

```
sp_configure

EXEC sp_configure 'external scripts enabled', 1
RECONFIGURE WITH OVERRIDE
```



|   | name | minimum | maximum | config_value | run_value |
|---|------|---------|---------|--------------|-----------|
| 25 | default language | 0 | 9999 | 0 | 0 |
| 26 | default trace enabled | 0 | 1 | 1 | 1 |
| 27 | disallow results from triggers | 0 | 1 | 0 | 0 |
| 28 | EKM provider enabled | 0 | 1 | 0 | 0 |
| 29 | external scripts enabled | 0 | 1 | 1 | 1 |
| 30 | filestream access level | 0 | 2 | 0 | 0 |
| 31 | fill factor (%) | 0 | 100 | 0 | 0 |

# Execute External Script

```
EXECUTE sp_execute_external_script
      @language = N'Python'
    , @script = @PythonScript
    , @input_data_1 = N'SELECT CONVERT(VARCHAR, Year) AS Year, Quarter, Client,
Revenue FROM Sales;'
    , @input_data_1_name = N'data'
    , @output_data_1_name = N'output'
WITH RESULT SETS ((
        Year NVARCHAR(10),
        Client NVARCHAR(10),
        Q1 INT,
        Q2 INT,
        Q3 INT,
        Q4 INT,
        Total INT
));
```

```
Declare @PythonScript nvarchar(max)
Set @PythonScript =N'
import pandas as pd

table = pd.crosstab(
    [data.Year, data.Client], # group by in rows
    data.Quarter, # group by in columns
    values = data.Revenue, # values to aggregate
    aggfunc= sum,
    margins= True
)

table.reset_index(inplace=True)

print(table)

output = table
'
```

| | Year | Quarter | Client | Revenue |
|---|---|---|---|---|
| 1 | 2014 | Q1 | Wallmart | 1000 |

| Year | Client | Q1 | Q2 | Q3 | Q4 | Total |
|---|---|---|---|---|---|---|
| 2014 | Fox | 6593 | 4332 | 123 | 6504 | 17552 |
| 2014 | Wallmart | 1000 | 560 | 2341 | 4000 | 7901 |
| 2015 | Fox | 34333 | 431 | 6665 | 4443 | 45872 |
| 2015 | Wallmart | 654 | 4555 | 8760 | 1233 | 15202 |
| All | NULL | 42580 | 9878 | 17889 | 16180 | 86527 |

# Install Libraries

# Install Libraries

Build AI App in 10min