# Travel Review Ratings

**Subject: Machine Learning**

Learn

Apply

Understand

**Presented By,**

Rajaraman Ganesan – 251056279 (M.Eng)

Vatsal Shah – 251041322 (M.Eng)

# Objective & Goal

Understand the **process** & **difference** by applying algorithms with **custom models** and **pre-build libraries** in Python

Apply algorithm in most-demanded languages in the industry - **Python & R**

**Compare** process, results and time complexity

Use **Amazon SageMaker** - understand the flow

Objective & Goal

Dataset Overview

Approach Design

Visualization Overview

Machine Learning Model

Results

Amazon SageMaker

Future Work & *QR Code

# Dataset overview

- Travel planning has become one of important commercial use.

- Many tourists look for some places like malls, restaurants or vacation spots, etc. online in recent times.

- Online review plays a critical role in the tourism industry, which mainly offers services and focuses on customer satisfaction.

**No Kaggle Kernels**

**Oriented**: UCI Machine Learning Repository ([Link](Link))
**Attributes** : 25
**Tuples**: 5456

Attribute 1 : Unique user id

Attribute 2 : Average ratings on churches

Attribute 3 : Average ratings on resorts

Attribute 4 : Average ratings on beaches

Attribute 5 : Average ratings on parks

Attribute 6 : Average ratings on theatres

Attribute 7 : Average ratings on museums

Attribute 8 : Average ratings on malls

Attribute 9 : Average ratings on zoo

Attribute 10 : Average ratings on restaurants

Attribute 11 : Average ratings on pubs/bars

Attribute 12 : Average ratings on local services

Attribute 13 : Average ratings on burger/pizza shops

Attribute 14 : Average ratings on hotels/other lodgings

Attribute 15 : Average ratings on juice bars

Attribute 16 : Average ratings on art galleries

Attribute 17 : Average ratings on dance clubs

Attribute 18 : Average ratings on swimming pools

Attribute 19 : Average ratings on gyms

Attribute 20 : Average ratings on bakeries

Attribute 21 : Average ratings on beauty & spas

Attribute 22 : Average ratings on cafes

Attribute 23 : Average ratings on view points

Attribute 24 : Average ratings on monuments

Attribute 25 : Average ratings on gardens

```python
# Importing the dataset
data = pd.read_csv('C:/train.csv')
print(data.shape)
data.head()
```

(5456, 25)

Out[9]:

| | User | Category 1 | Category 2 | Category 3 | Category 4 | Category 5 | Category 6 | Category 7 | Category 8 | Catego |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | User 1 | 0.0 | 0.0 | 3.63 | 3.65 | 5.0 | 2.92 | 5.0 | 2.35 | 2.3 |
| 1 | User 2 | 0.0 | 0.0 | 3.63 | 3.65 | 5.0 | 2.92 | 5.0 | 2.64 | 2.3 |
| 2 | User 3 | 0.0 | 0.0 | 3.63 | 3.63 | 5.0 | 2.92 | 5.0 | 2.64 | 2.3 |
| 3 | User 4 | 0.0 | 0.5 | 3.63 | 3.63 | 5.0 | 2.92 | 5.0 | 2.35 | 2.3 |
| 4 | User 5 | 0.0 | 0.0 | 3.63 | 3.63 | 5.0 | 2.92 | 5.0 | 2.64 | 2.3 |

5 rows × 25 columns

| Objective & Goal | Dataset Overview |
| --- | --- |
| Approach Design | Visualization Overview |
| Machine Learning Model | Results |
| Amazon SageMaker | Future Work & *QR Code |

# Process Design

Unsupervised learning approach

Don't have a preliminary info on output values

Explore data to find some intrinsic structures in them

Use of clustering algorithm technique like k-means algorithm.

Objective & Goal

Dataset Overview

Approach Design
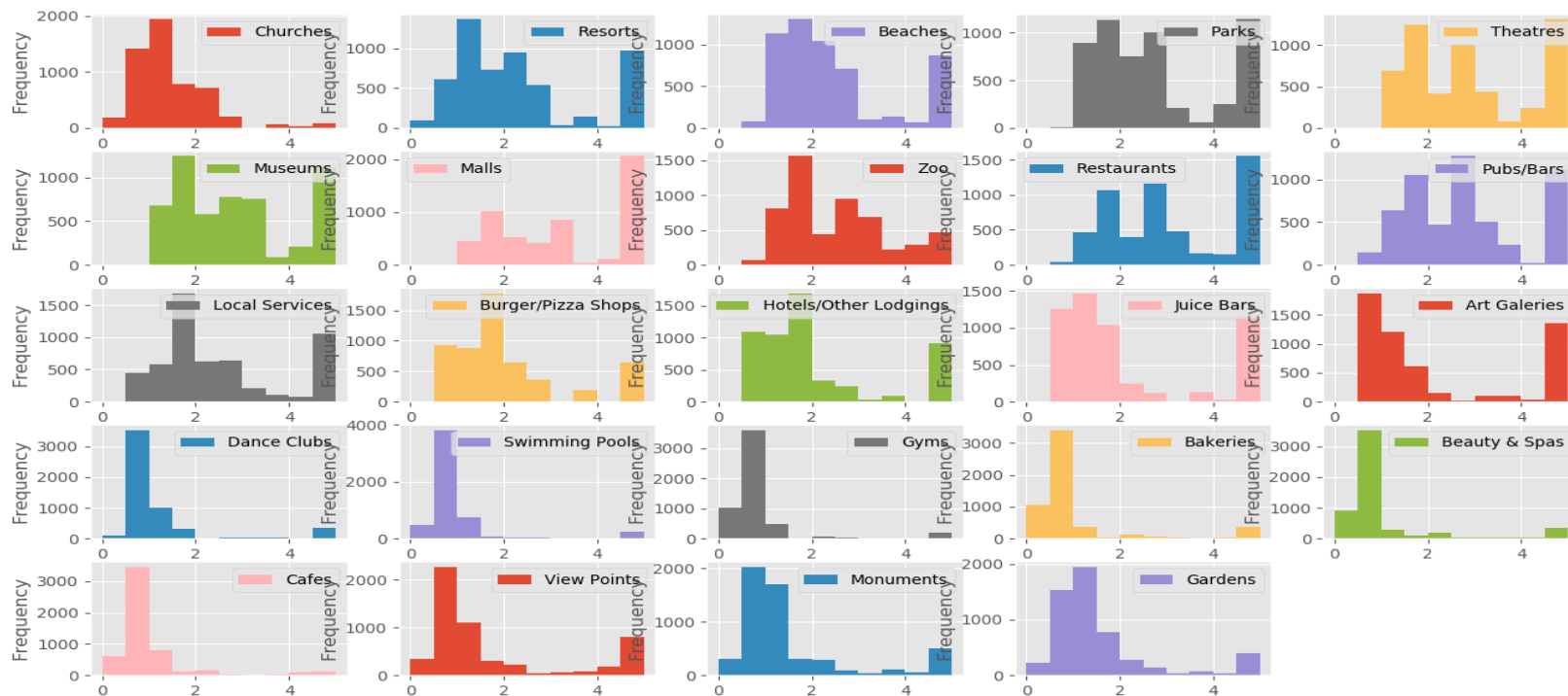
Visualization Overview

Machine Learning Model

Results

Amazon SageMaker

Future Work & *QR Code

# Visualization Overview

Shows 24 categories reviewed by users with hist diagram

# Continue…

## Shows 24 categories reviewed by users with density diagram

# Continue…

Shows 24 categories reviewed by users with box plot diagram

Objective & Goal

Dataset Overview

Approach Design

Visualization Overview

Machine Learning Model
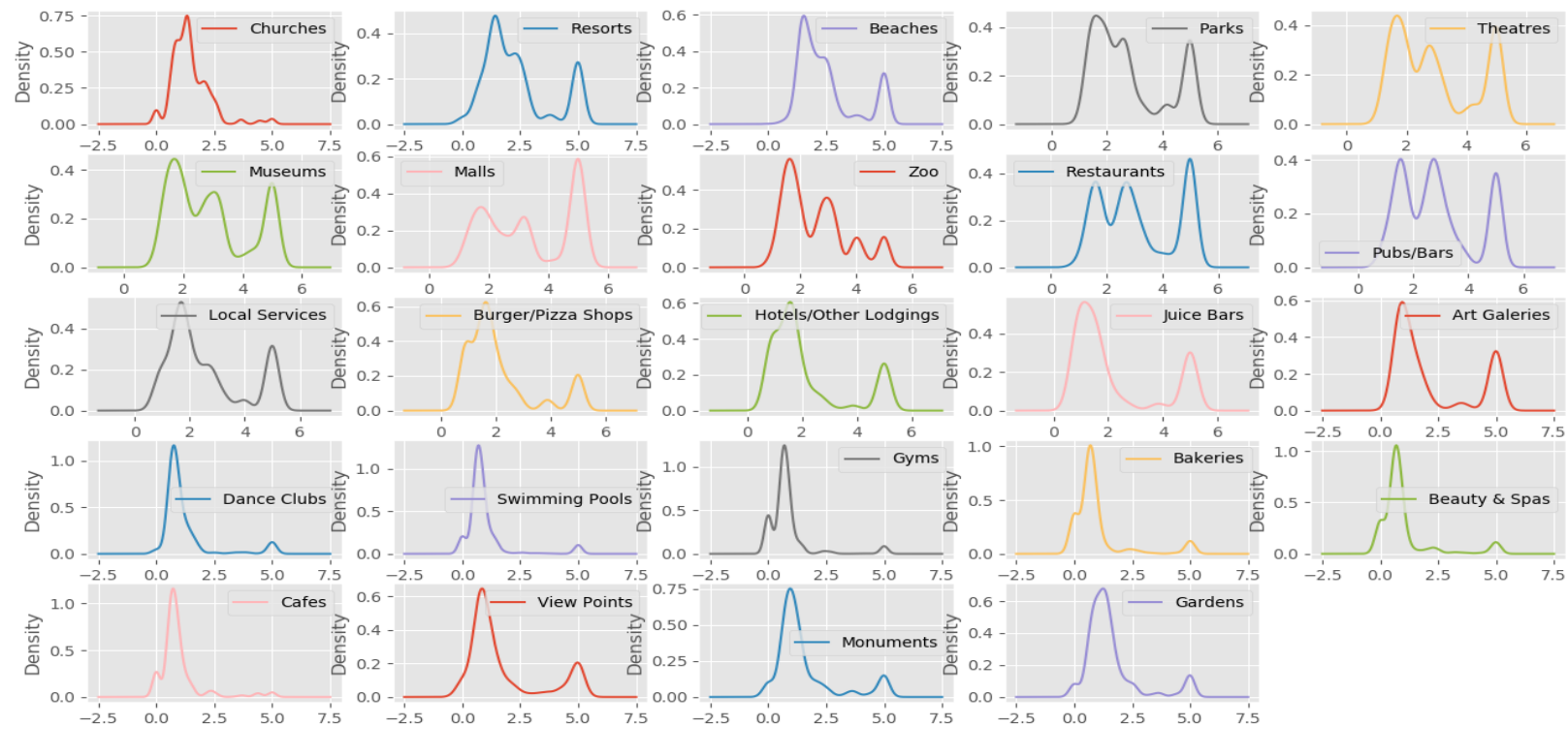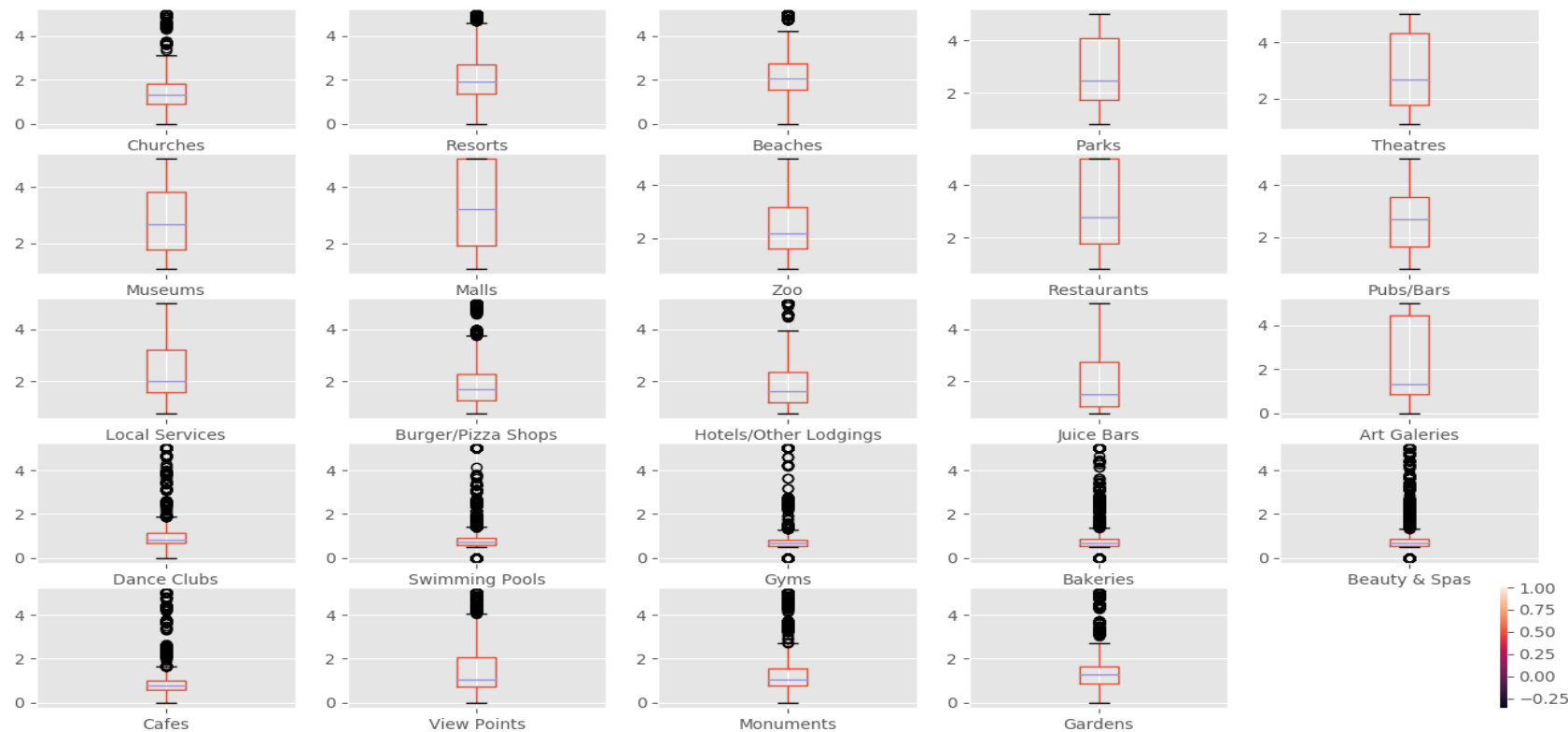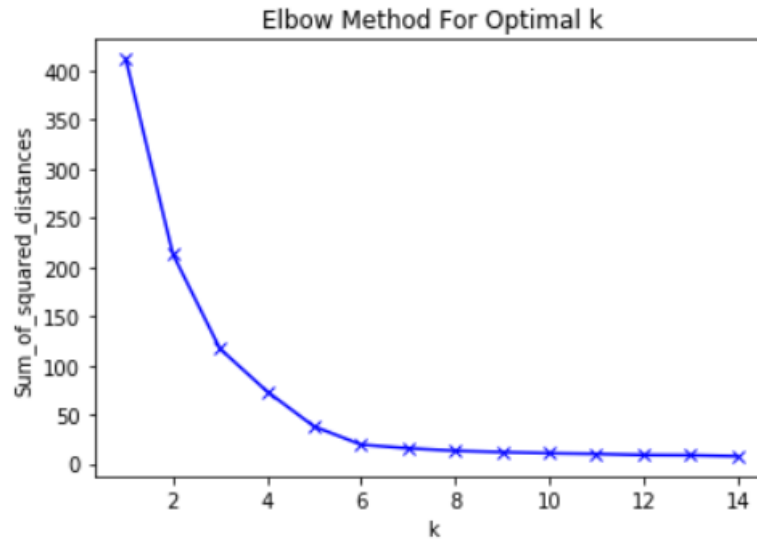
Results

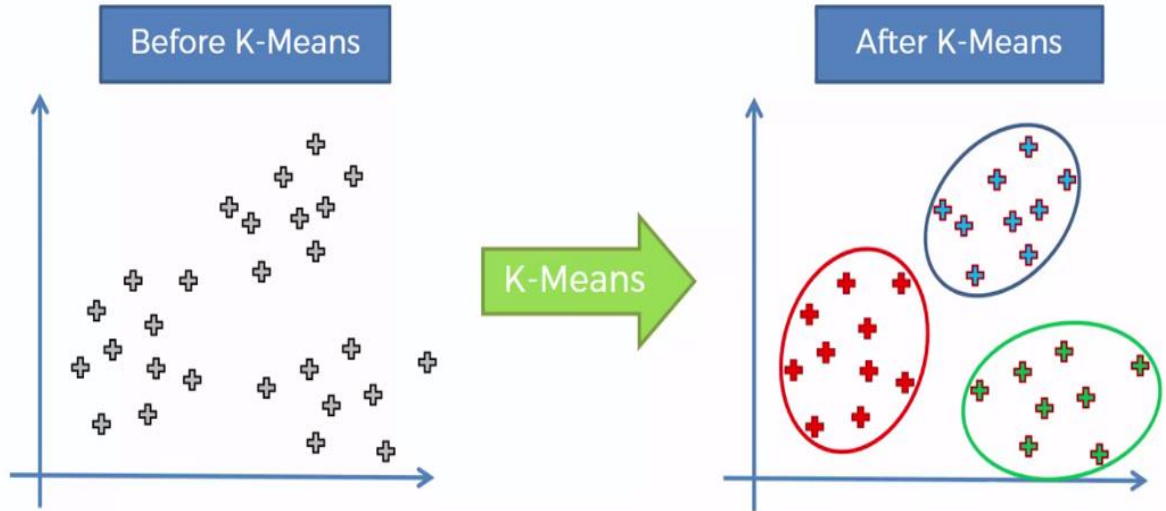Amazon SageMaker

Future Work & *QR Code

# K-Means Algorithm

- **Choose no: of clusters**

- Initialization

- Assign cluster

- Move centroid

- Optimization

- Convergence

- Within-cluster sum of squares is a measure of the variability of the observations within each cluster
- Idea behind using elbow method to choose after which **WSS** is almost constant.
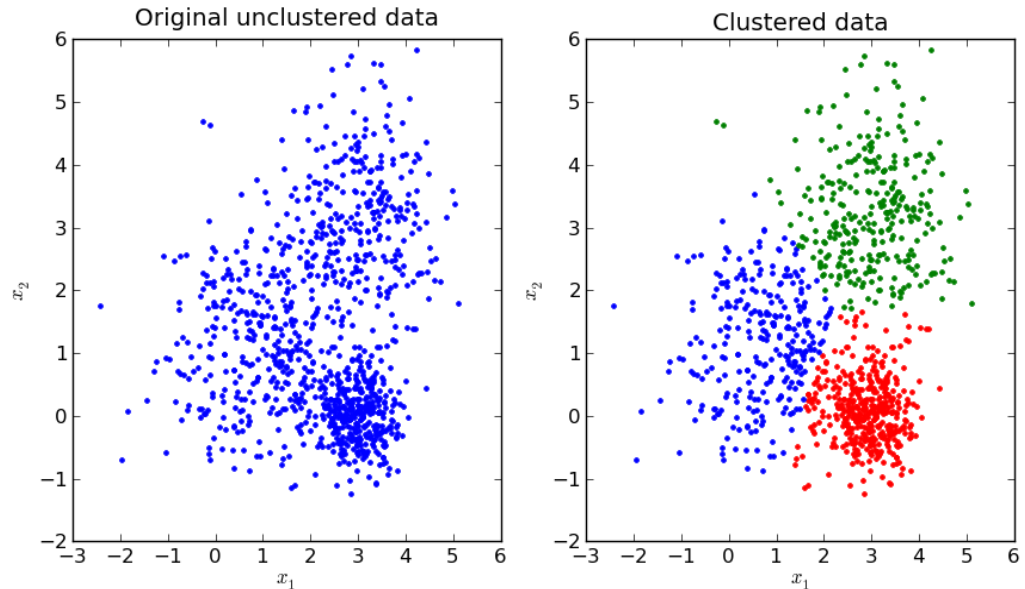


Elbow Method For Optimal k

# K-Means Algorithm

- Choose no: of clusters
- **Initialization**
- Assign cluster
- Move centroid
- Optimization
- Convergence

- Initialize k points, randomly
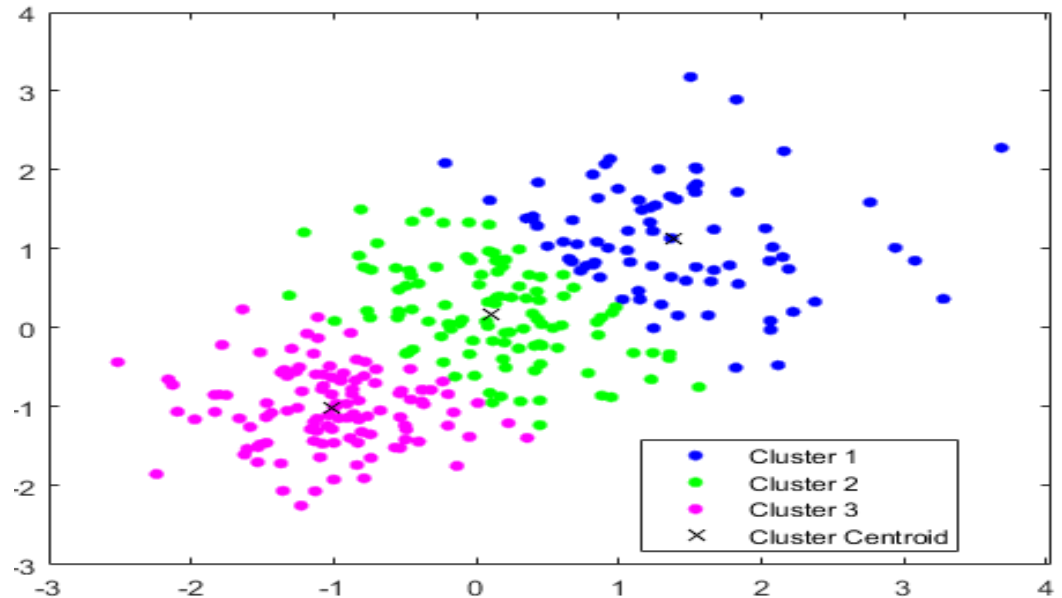- Value of clusters are determined by elbow curve.

# K-Means Algorithm

- Choose no: of clusters
- Initialization
- **Assign cluster**
- Move centroid
- Optimization
- Convergence

- Dist. between data points and centroid are computed.
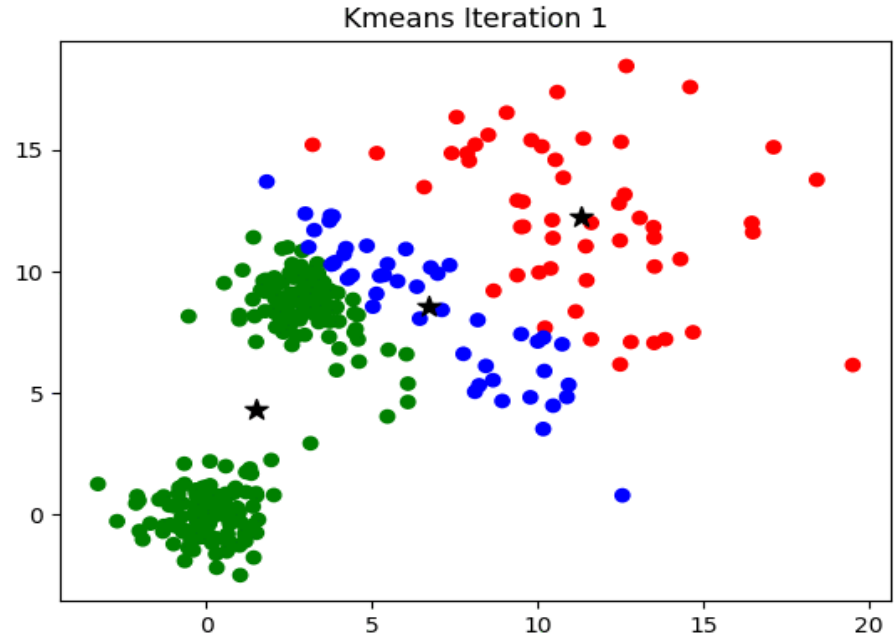- Based on min. Distance, data are divided into groups.

# K-Means Algorithm

- Choose no: of clusters

- Initialization

- Assign cluster

- **Move centroid**

- Optimization

- Convergence

- Compute the mean of all three dots.
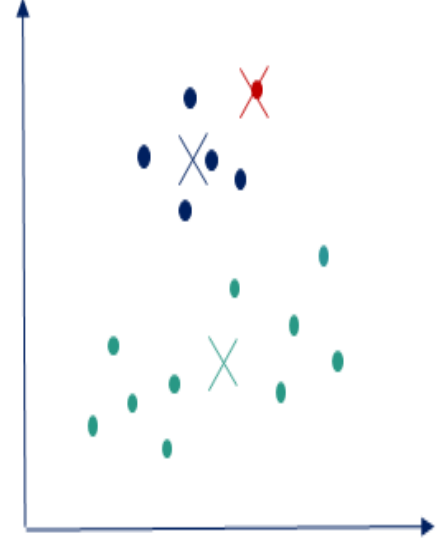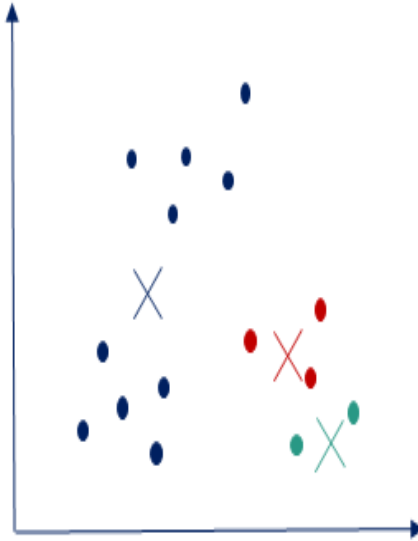- Reposition blue, green, pink cluster centroid to this mean.

# K-Means Algorithm

- Choose no: of clusters

- Initialization

- Assign cluster

- Move centroid

- **Optimization**

- Convergence

- Redo the above steps till cluster stop changing their positions.



Kmeans Iteration 1

# K-Means Algorithm

- Choose no: of clusters
- Initialization
- Assign cluster
- Move centroid
- Optimization
- **Convergence**

- Divide the data points to clusters, once the algorithm converges.

Objective & Goal

Dataset Overview

Approach Design

Visualization Overview

Machine Learning Model
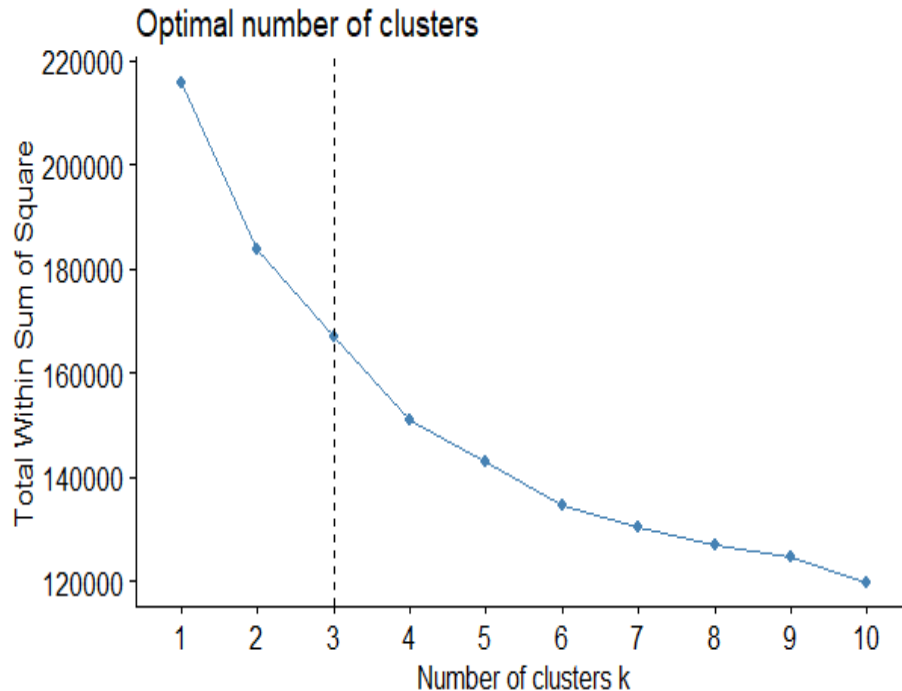
Results

Amazon SageMaker

Future Work & *QR Code

# R - Library

- library (cluster) –
  - Finding groups in data

- library (ggplot2) –
  - system for declaratively creating graphs
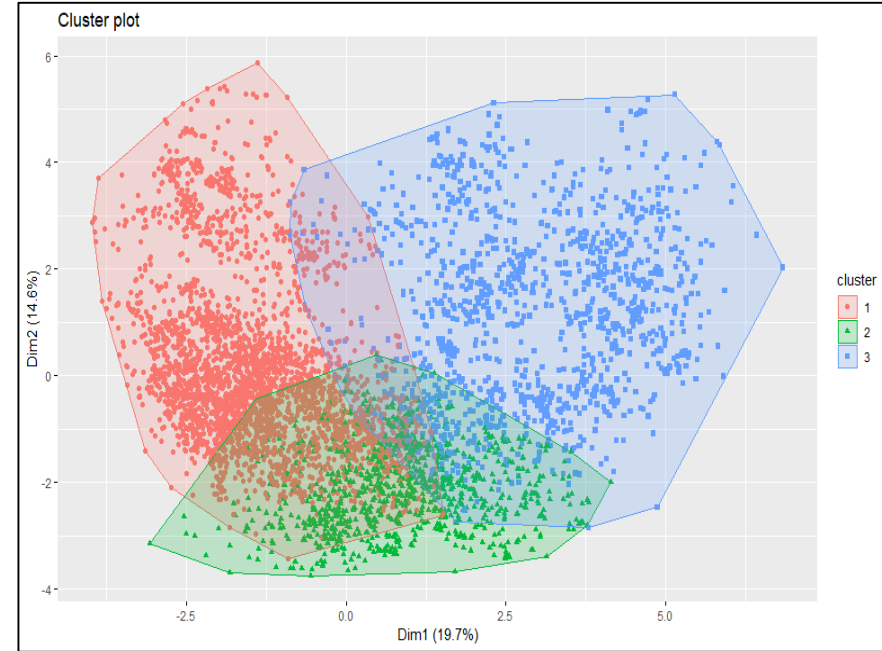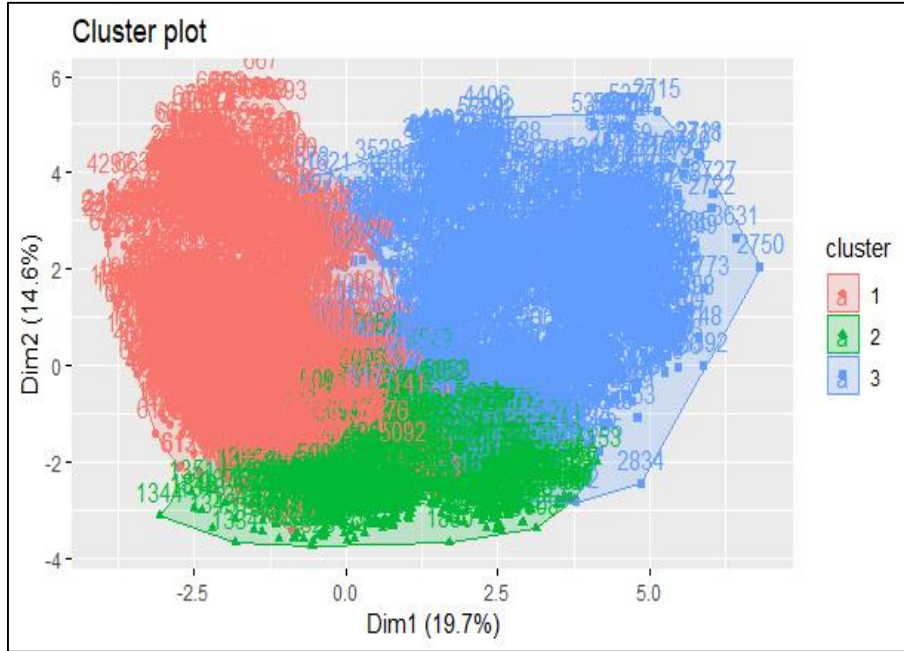
- library (factoextra) –
  - Used to visualize clusters

```
>
>
> km.res$size
[1] 1093 2460 1903
>
>
>
>
```
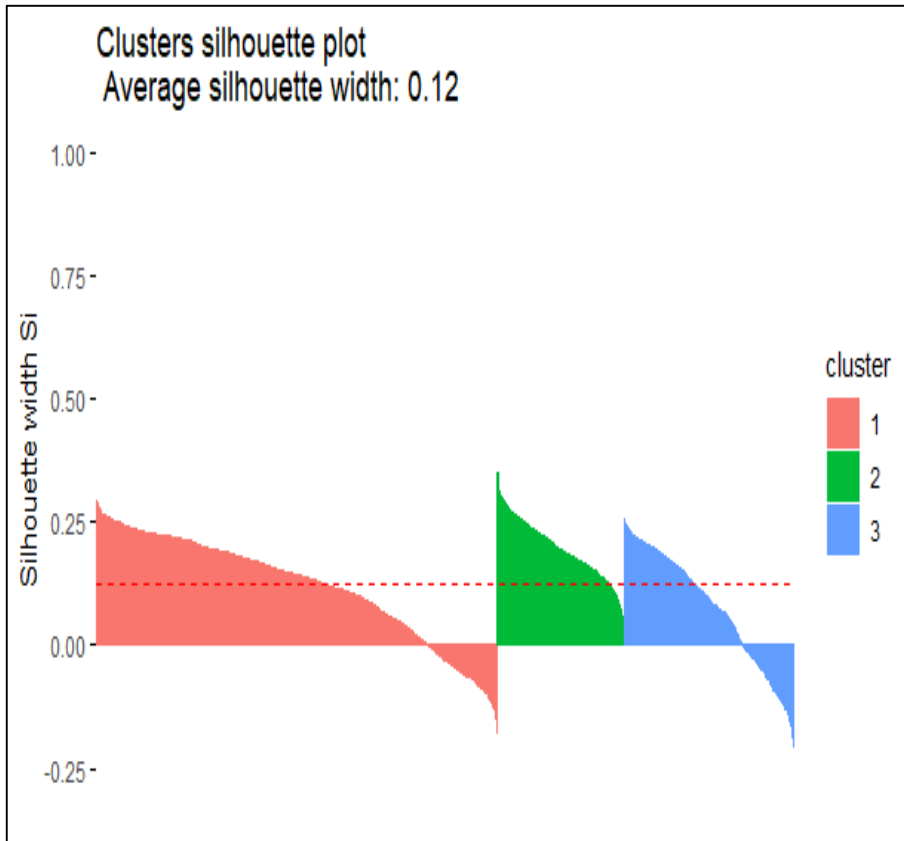
# Results - R



Optimal number of clusters

- Method to validate the number of clusters is the elbow method.
- Help finding the appropriate number of clusters in a dataset.
- Here, 3 clusters are suggested.

# Continue…



- As it a multidimensional dataset, to visualize cluster plot PCA is used.
- Reduce the dimensionality of a data set consisting of many variables correlated with each other
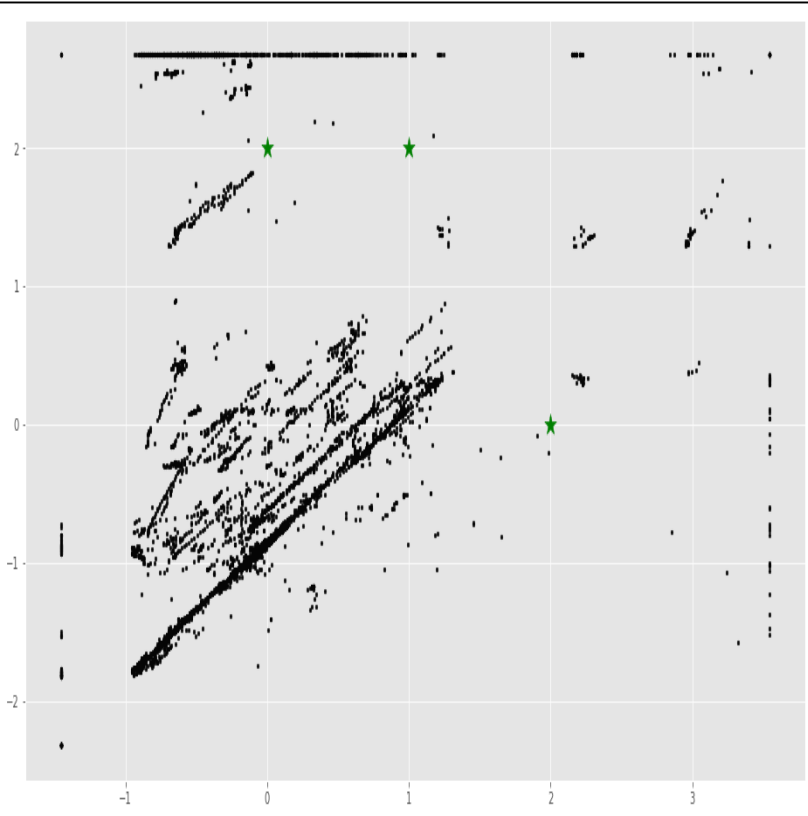
# Continue…



**Silhouette Plot** shows for each cluster:
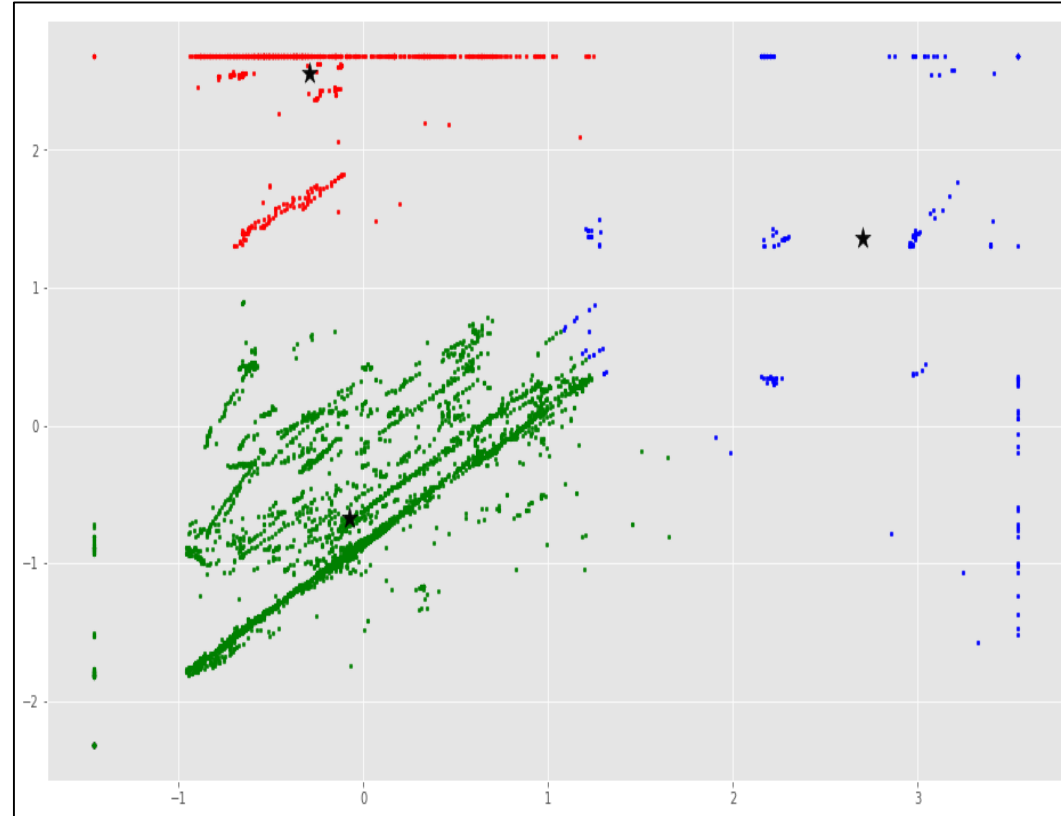
- Each vertical line corresponds to an element.
- The average silhouette width

```
  cluster size ave.sil.width
1       1 3140           0.12
2       2  992           0.19
3       3 1324           0.08
>
```

# Results - Python



Plot along with random centroid

Assigning new centroid

(Screenshot taken during evaluation process)

Objective & Goal

Dataset Overview

Approach Design

Visualization Overview
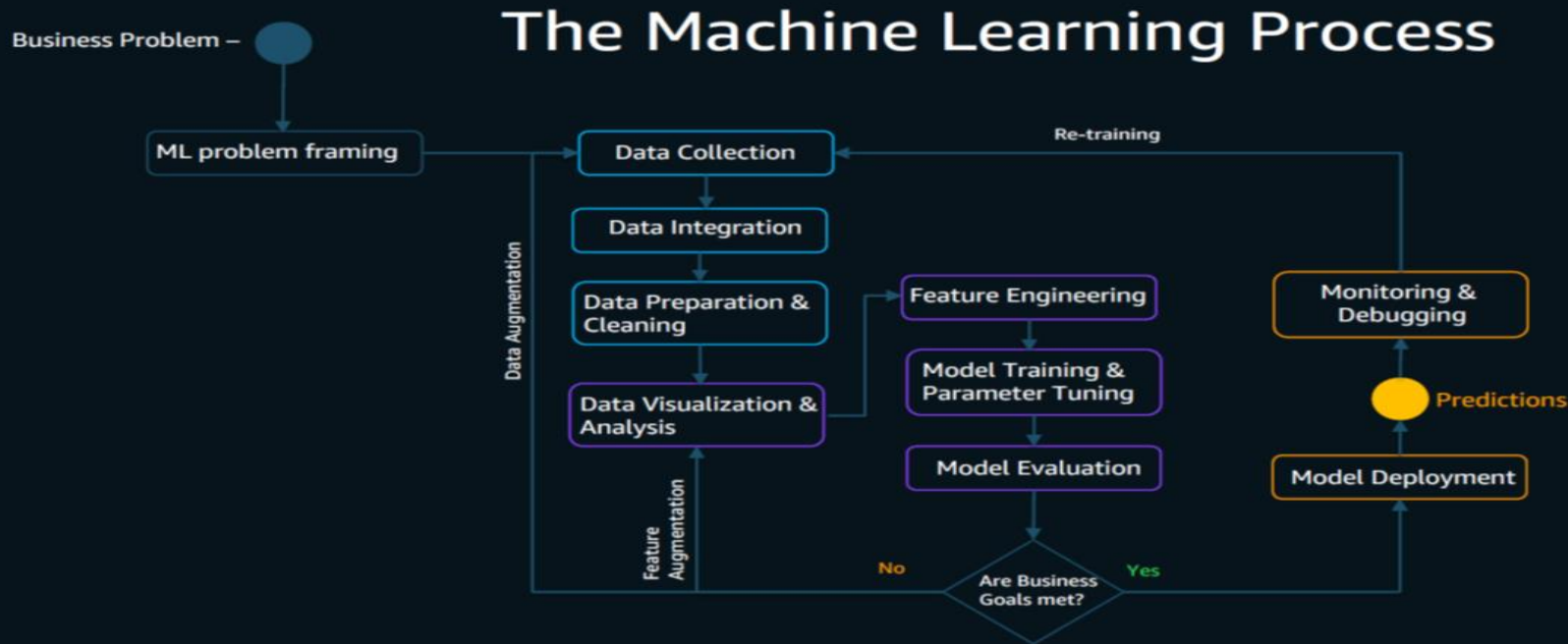
Machine Learning Model

Results
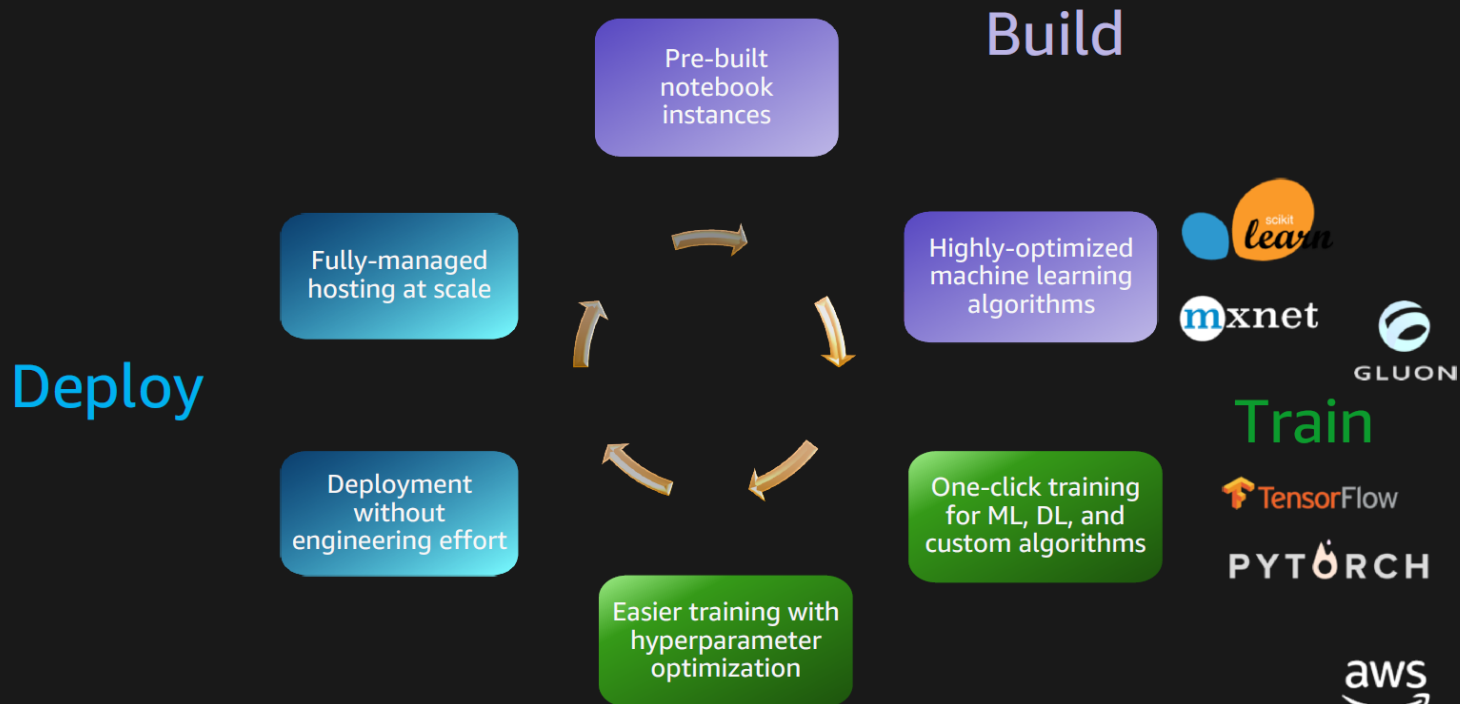
Amazon SageMaker

Future Work & *QR Code

# ML Process on AWS



The Machine Learning Process

Business Problem –

ML problem framing → Data Collection → Data Integration → Data Preparation & Cleaning → Data Visualization & Analysis → Feature Engineering → Model Training & Parameter Tuning → Model Evaluation → Are Business Goals met? — No / Yes → Model Deployment → Monitoring & Debugging → Predictions

Data Augmentation

Feature Augmentation

Re-training

Amazon SageMaker

# Dashboard



aws | Services ▾ | Resource Groups ▾ | ⭐ | 🔔 vatsal ▾ | N. Virginia ▾ | Support ▾

**Amazon SageMaker** ✕

Dashboard

Search^Beta

**Ground Truth**
Labeling jobs
Labeling datasets
Labeling workforces

**Notebook**
Notebook instances
Lifecycle configurations
Git repositories

**Training**
Algorithms
Training jobs
Hyperparameter tuning jobs

**Inference**
Compilation jobs
Model packages
Models
Endpoint configurations

Set up and manage labeling jobs for highly accurate training datasets using active learning and human labeling.

[ Labeling jobs ]

Availability of AWS and SageMaker SDKs and sample notebooks to create training Jobs and deploy models.

[ Notebook instances ]

Train and tune models at any scale. Leverage high performance AWS algorithms or bring your own.

[ Training jobs ]

[ Hyperparameter tuning jobs ]

Create models from training jobs or import external models for hosting to run inferences on new data.

[ Models ]

[ Endpoints ]

[ Batch transform jobs ]

## Recent activity

Recent activity within the [ Last 7 days ▾ ]

**Ground Truth**
Labeling jobs
No recent activity.

**Notebook**
Notebook instances
⊘ 1 Created

**Training**
Training jobs
No recent activity.

Hyperparameter tuning jobs
No recent activity.

**Inference**
Models
No recent activity.

Endpoints
No recent activity.

Batch transform jobs
No recent activity.

💬 Feedback | 🌐 English (US)

© 2008 - 2019, Amazon Internet Services Private Ltd. or its affiliates. All rights reserved.    Privacy Policy    Terms of Use

# Notebook instances

# Instance in pending status

# In Service

# Jupyter Lab

# Jupyter

Open JupyterLab     Quit

Files     Running     Clusters     SageMaker Examples

Select items to perform actions on them.

| ☐ 0 ▾ | 📁 / machine-learning-kmeans | Name ↓ | Last Modified | File size |
|---|---|---|---|---|
| | 📁 .. | | seconds ago | |
| ☐ | 📓 project.ipynb | | 6 days ago | 566 B |
| ☐ | 📄 main.py | | 6 days ago | 149 B |
| ☐ | 📄 README.md | | 6 days ago | 76 B |
| ☐ | 📄 train.csv | | 6 days ago | 637 kB |

Upload     New ▾

# Examples

# On AWS

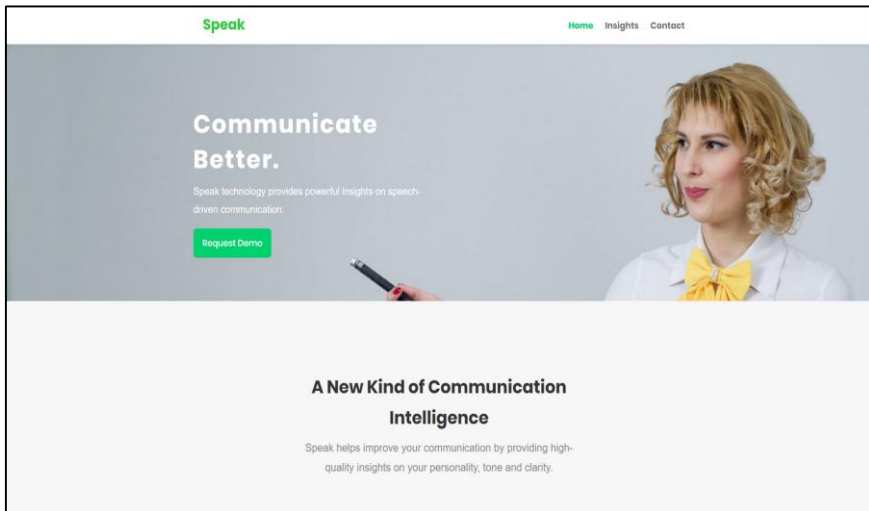| Objective & Goal | Dataset Overview |
| --- | --- |
| Approach Design | Visualization Overview |
| Machine Learning Model | Results |
| Amazon SageMaker | Future Work & *QR Code |

# Technologies Used

# Future work

Find | Find more data related to text, audio & video

Explore | Explore more with AWS SageMaker

Apply | Implement for real world problems

# QR Code

**Access resources and source code**

SCAN ME

Issues  Marketplace  Explore

vatsal2210 / Machine-Learning-Project-Western

Unwatch ▾ 2    ★ Star 0    Fork 0

<> Code    Issues 0    Pull requests 0    Projects 0    Wiki    Insights    Settings

Machine learning subject project

Manage topics

Edit

22 commits    1 branch    0 releases    1 contributor

Branch: master ▾    New pull request    Create new file    Upload files    Find File    Clone or download ▾

vatsal2210 Delete vatsal2210 Machine-Learning-Project-Western Machine learning s... ...    Latest commit 0658071 16 seconds ago

| | | |
|---|---|---|
| Amazon SageMaker | Updates | 18 hours ago |
| Python | Updates | 21 minutes ago |
| R | Updates | 21 minutes ago |
| Resources | Merge branch 'master' of https://github.com/vatsal2210/Machine-Learni... | 30 seconds ago |
| Group 16 - Presentation.pptx | Updates | a minute ago |
| Project_Information.pdf | Updates | 18 hours ago |
| README.md | Updates | 18 hours ago |
| train.csv | changes | 20 days ago |

README.md

## Machine-Learning-Project-Western

A project is on user's rating of travelling. There are various 24 attractions included in this database.

### Dataset description

- Total categories: 24
- Total Number of users: 5456
- ID Column: Users ID

### Columns:

Attribute 1 : Unique user id Attribute 2 : churches Attribute 3 : resorts Attribute 4 : beaches Attribute 5 : parks Attribute 6 : theatres Attribute 7 : museums Attribute 8 : malls Attribute 9 : zoo Attribute 10 : restaurants Attribute 11 : pubs/bars

# Resources:

- Getting started with Amazon SageMaker: Link
- Use the Amazon SageMaker SDK: Python: Link
- 'Evaluation of Partitioning Clustering Algorithms for Processing Social Media Data in Tourism domain', Shini Renjith, A. Sreekumar, M. Jathavedan, 2018 IEEE Recent Advances in Intelligent Computational Systems (RAICS) | December 06 - 08, 2018 , Trivandrum Link
- 'Use and Impact of Online Travel Reviews', Markus Schuckert , Liu XianweiRob Law Link
- 'Hospitality and Tourism Online Reviews: Recent Trends and Future Directions', Ulrike Gretzel , Kyung-Hyan Yoo Link
- SageMaker Examples: Link
- Python vs R Comparison: Link
- Python Libraries: Link
- R Libraries : Link
- K-means Algorithm: Link

time for questions

THANK YOU