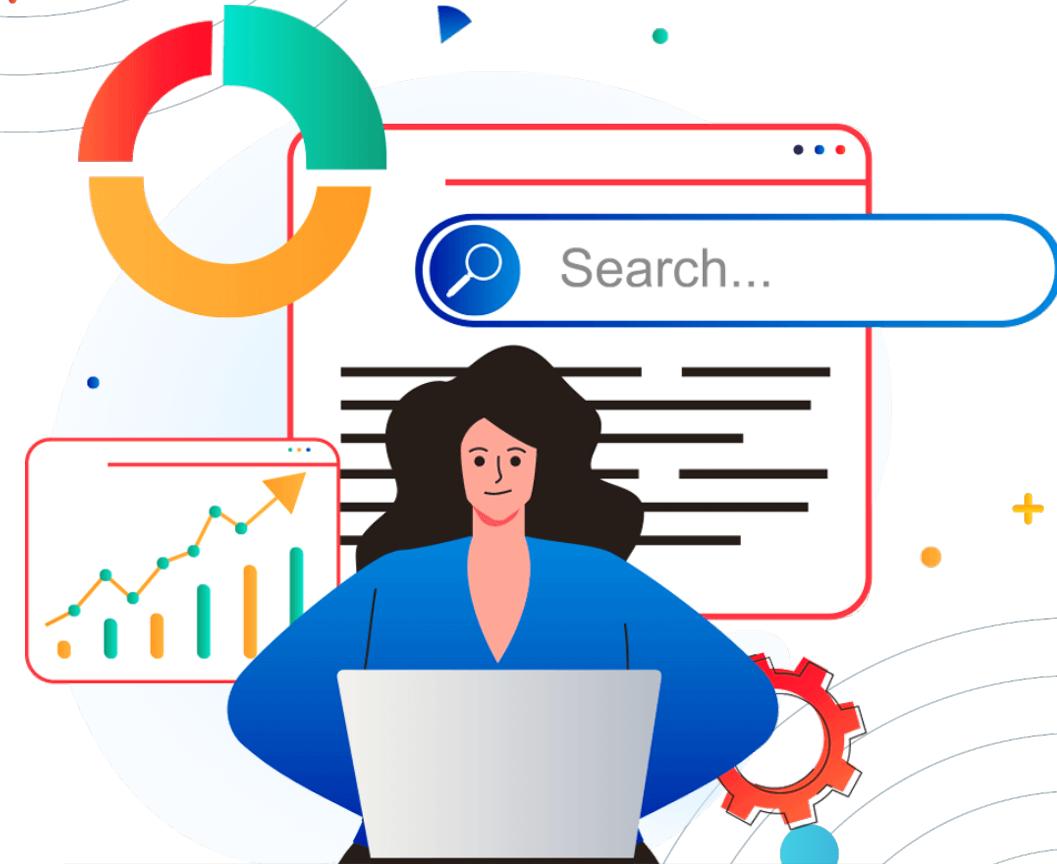


Data Analysis For Shark Tank India Pitches

**Using Python, R, Excel and
Tableau**

By: Shipra Mehandru

IITD BAO'21 Course – Feb 2022



What the Project Covers



Dataset understanding and Preprocessing on Python and Excel

- I have **eliminated the non correlated** variables using Excel regression runs
- Also **standardized the variables** before applying any ML model



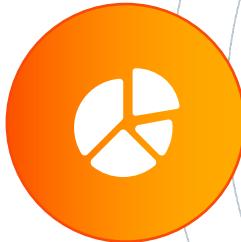
Predicting classes using ML models on Python

- All the **Supervised learning ML Classification** runs used have been showed in detail along with an explanation for why it was done
- And **links to output notebooks** are included on every page instead of copy pasting all screen shots



Bettering analysis through R Studio, Clustering

- In order to see **improved prediction accuracies** I have analyzed 2 more datasets, which includes:
 - A **bigger sample**;
 - A better **correlation** analysis through **R Studio**;
 - Finding better clusters through **unsupervised K Means**;
 - While avoiding the overfit model trap



Data Visualizations using Tableau & Python

- The deck also shows **charts for every run on Python** and R Studio using inbuilt visualization libraries
- A **Tableau Dashboard** has been added as a screen shot at the end, which can be accessed on the Public Tableau link. It analyzes data further

Understanding the Dataset

The dataset is from Kaggle:

Where a user compiled 121 pitches from Season 1 using Wikipedia

kaggle [Link](#)

This is the Description of Variables and the Meta Data of the File:

- 1. Episode_number - Number of the episode
 - 2. Pitch_number - Number of the Pitch
 - 3. Brand_name - Name of the brand Idea
 - 4. Idea - behind the brand building Deal
 - 5. Deal - done or not ; 1 - YES, 0 - NO
 - 6. Pitcheraskamount - Amount asked by the pitchers
 - 7. Ask_equity - Equity offered by the pitchers
 - 8. Ask_valuation - Valuation asked by pitchers
 - 9. Deal_amount - Final Deal Amount
 - 10. Deal_equity - Final Deal equity percentage
 - 11. Deal_valuation - Final Valuation of Company after Deal
- + 17 other Data points on the Investors or Sharks who were present and who invested in the pitches



Ghazal Alagh
Co-Founder & CIO of
Mamaearth



Anupam Mittal
Founder & CEO of
People Group



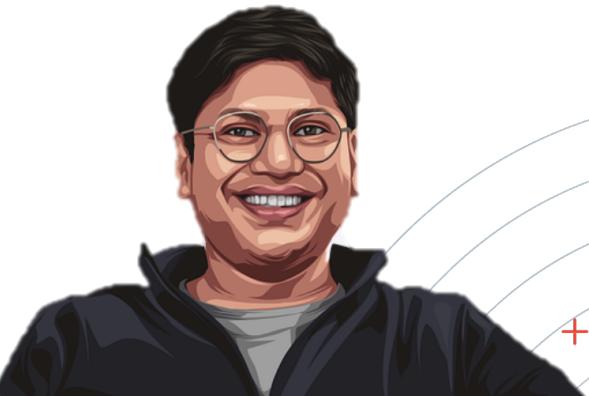
Namita Thapar
Executive Director at
Emcure Pharma



Vineeta Singh
CEO & Co-Founder of
SUGAR Cosmetics



Aman Gupta
Co-Founder & CMO
of boAt



Peyush Bansal
Founder & CEO of
Lenskart



Ashneer Grover
Co-Founder & MD of
BharatPe

Understanding Classifiers

Machine Learning has a lot of techniques that can be used to analyze a dataset

(For Supervised Learning where the output variable is known)

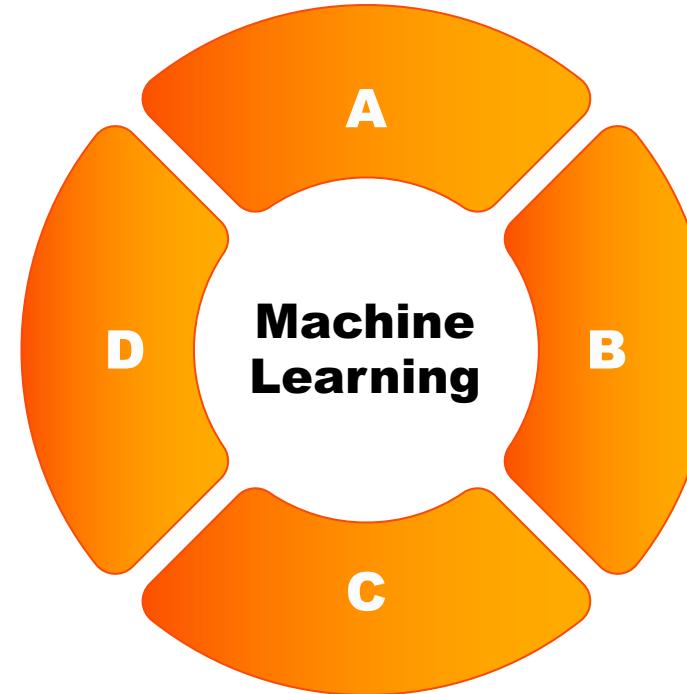
Classification

- Logistic Regression
- K-Nearest Neighbors
- Support Vector
- Naïve Bayes
- Decision Tree
- Random Forest

(Which learns from a reward mechanism and improves itself through hidden layers)

Deep Learning

- Artificial Neural Networks
- Convolutional/Recurrent Neural
- Reinforcement Neural



+ Dimension Reduction Analysis
Big Data techniques
NLP and other

Regression

- Linear and polynomial
- Support Vector (with/without Kernel)
- Decision Tree
- Random Forest

(For Continuous values, where there are no categorical output values)

Clustering

- K-Means
- Hierarchical

(For Unsupervised Learning where we don't know the output variable, it is not labeled as yet)

I will be using the following for the purposes of analyzing the Shark Tank data:



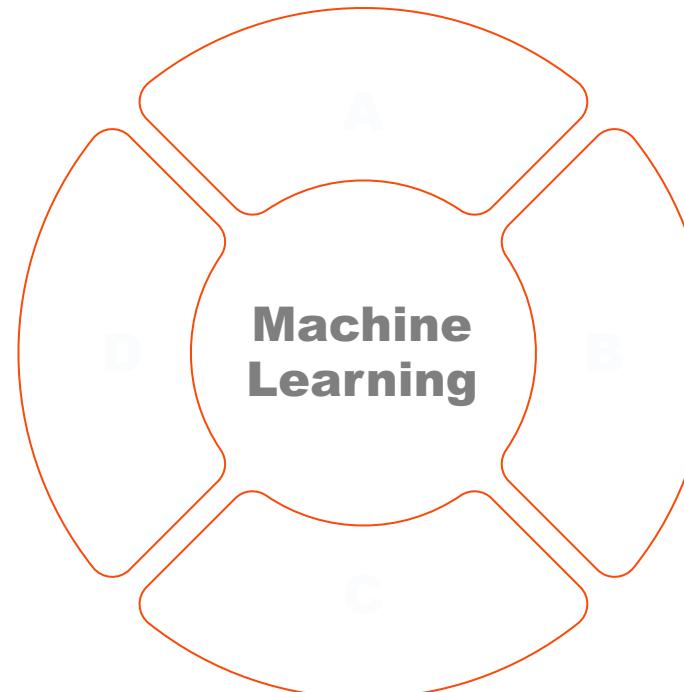
Classification

- Logistic Regression
- K-Nearest Neighbors
- Support Vector
- Naïve Bayes
- Decision Tree
- Random Forest

To predict Deal or No Deal outcomes i.e. Categorical Output

Deep Learning

- Artificial Neural Networks
- Convolutional/Recurrent Neural
- Reinforcement Neural



Regression

- Linear
- Polynomial
- Support Vector (with/without Kernel)
- Decision Tree
- Random Forest

To understand how much a variable is influencing the outcome

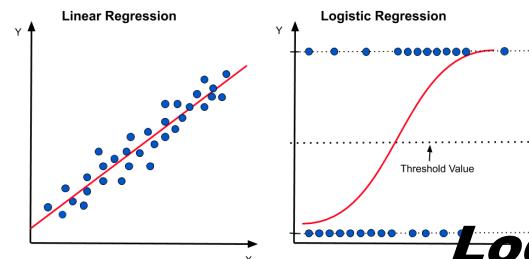
Clustering

- K-Means
- Hierarchical

To understand if there can be better categorical outcomes for the same data

+ Dimension Reduction Analysis
Big Data techniques
NLP and other

What each of these Supervised Classification models show:

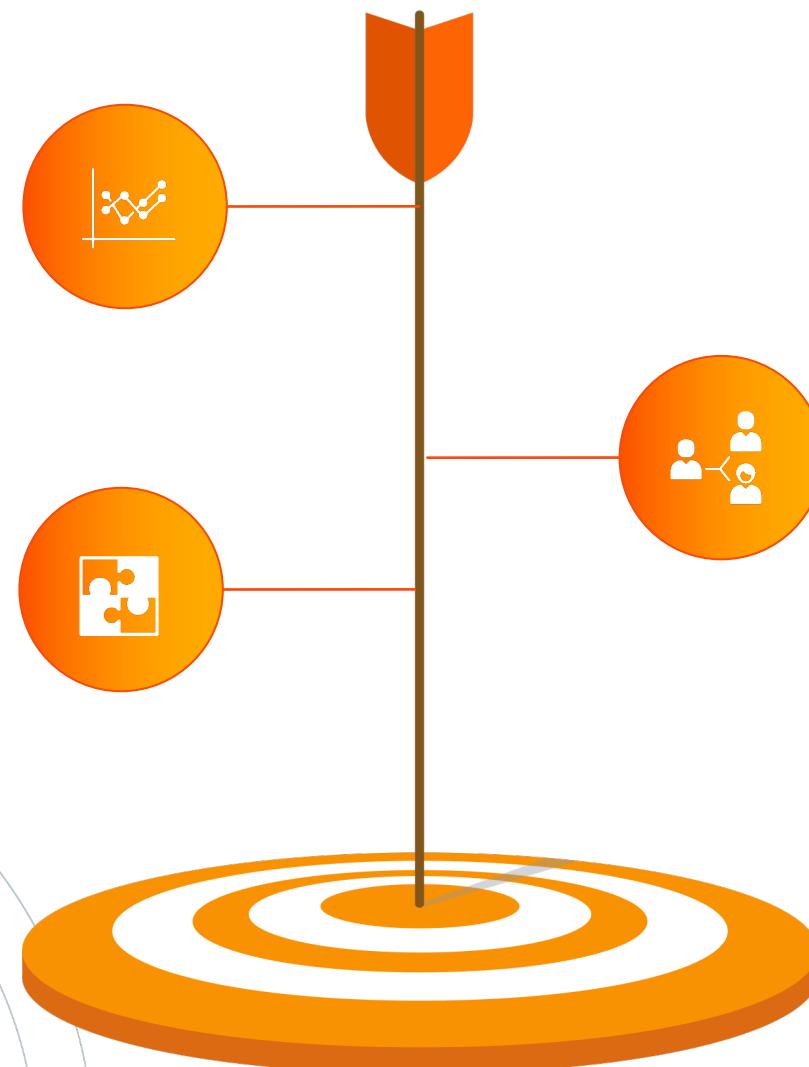
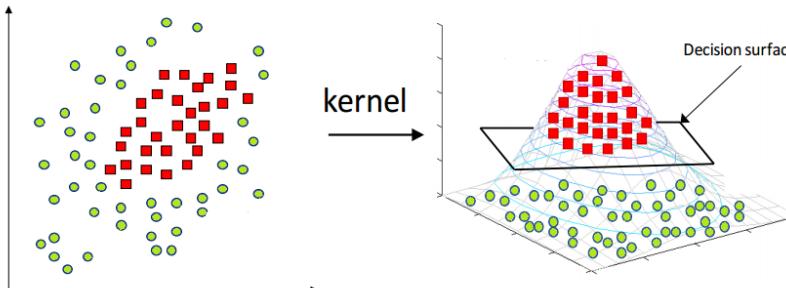


Logistic Reg.

Based on the concept of probability - from 0 to 1, of a datapoint falling into each class instead of Linear where the points can fall anywhere along a linear line

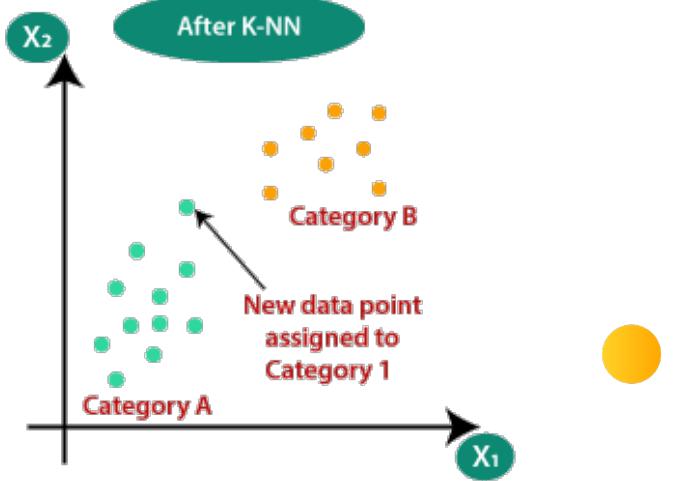
Kernel SVM

It either creates a hyperplane where the classes exist separately – in a Linear SVM or for Kernel method (RBF) / non linear data it separates via area on a higher plane

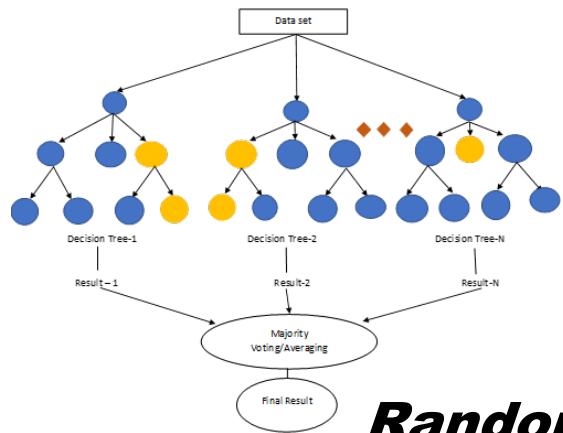


K-NN

Defines a dataset's class by finding their distances from the chosen variables. A point should be nearer to its neighbors and farther from other classes

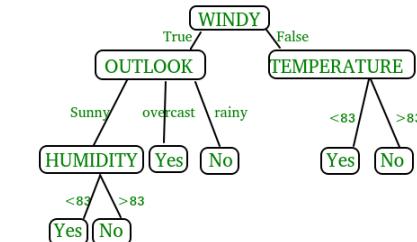
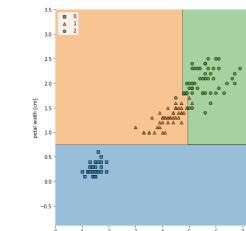
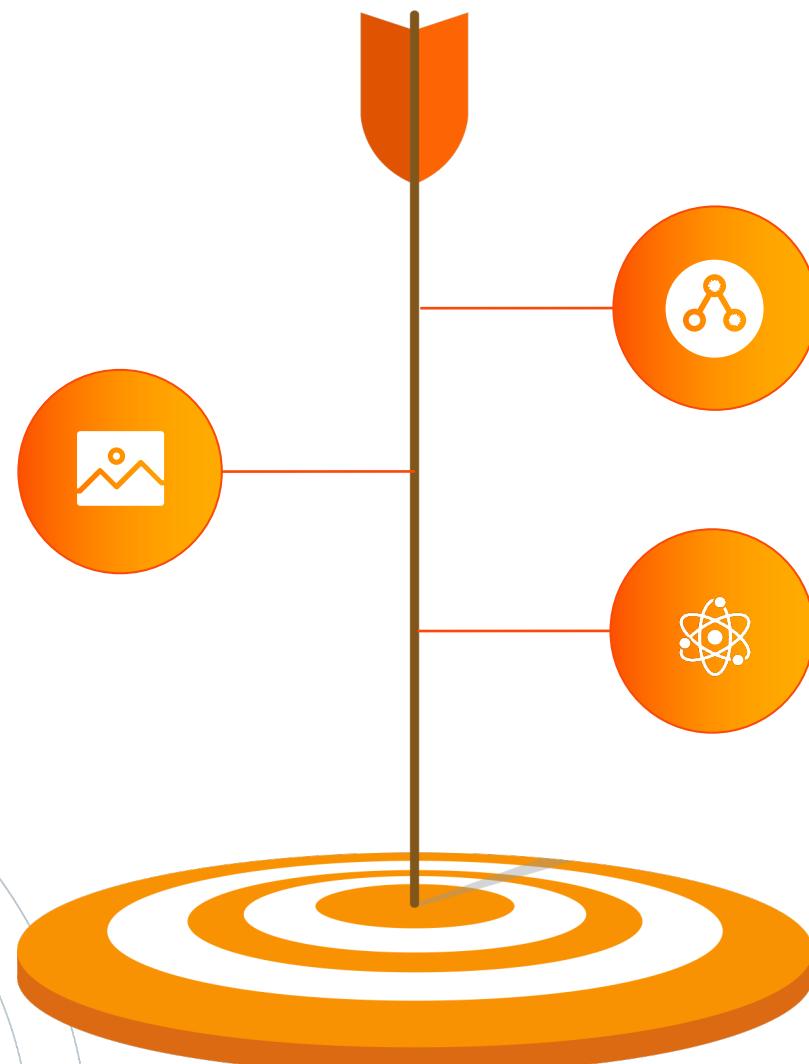
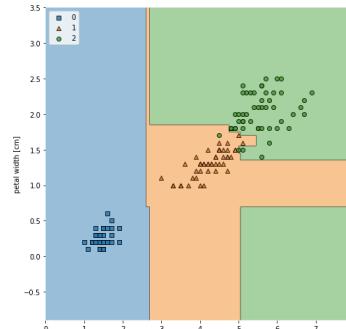


What each of these Supervised Classification models show:



Random Forest

As the name suggests it's a cluster of different decision trees and building a strong classifier using weaker classifiers



Decision Tree

Based on a series of events occurring, we go from the node that has the highest certainty and keep branching low certainty nodes, variables

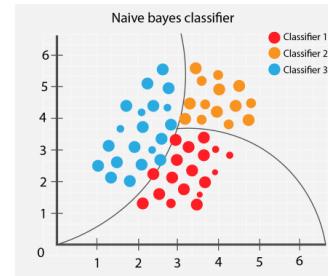
Naïve Bayes

It is also based on probability of a datapoint being in one class given the probability that other have been distributed to its neighboring classes. I.e. event A occurring while B has already occurred

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

using Bayesian probability terminology, the above equation can be written as

$$\text{Posterior} = \frac{\text{prior} \times \text{likelihood}}{\text{evidence}}$$



*Before we get into the
Classifiers, let us see
the pertinent variables
to run ML models by
using **Regression** &
Correlation in Excel*

I did a bit of data preprocessing in the excel file itself:

Excel screenshot showing a table of data. The 'Social Benefit' column is highlighted in yellow.

Excel screenshot showing a regression analysis output. A formula `=CORREL(A2:A122,O2:O122)` is highlighted in red.

From the pitches and industries, I categorized the pitches as offering some **Social Benefit** or not

Chose pitchers' requested **amount, equity, valuation, number and gender of pitchers & tech vs not tech** pitches as predicting Deal or no Deal in the **first Regression Run**

Excel ribbon showing the 'Data Analysis' tool selected. Below it, a regression analysis output is shown with several coefficients highlighted in yellow.

Did **3 rounds of regressions** and correlations, eliminating the variables with a high P value. Tool the ones with a high correlation **and low P value (<0.5)** to as the **most influential variables** for the ML predictor models

+ I used the **STANDARDISE** Formula with Mean and Standard Deviation to standardize all values before the regression runs

What is a Regression Equation and hence how are Coefficients of variables useful?

It helps examine each of the **independent** variables' **linear relationship** with the dependent variable

Multiple Regression Model with k Independent Variables:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k + \varepsilon$$

Y-intercept Population slopes Random Error

Whereas a Correlation model **directly gives us the relation** of one array of values in a variable to those of another variable. That is also a Linear relationship measure, **though not causal**

Variables Chosen for the ML Predictions of Deal or No Deal:

Interestingly most marketing data I've worked with usually has a low R Square and a high Error. We see the same thing here with low regression coefficients and correlation scores.

Tech or not Tech Pitch

Correlation = 0.01
Reg. Coeff. = -0.04



Pitchers willing to give % of Equity

Correlation = -0.06
Reg. Coeff. = -0.05



No of Female Pitchers in team

Correlation = -0.2
Reg. Coeff. = -0.11



Social Benefit of the Pitch product/ service (Eg: environment, health, city)

Correlation = 0.16
Reg. Coeff. = 0.2



Asked valuation by the pitchers for their product/ service

Correlation = -0.18
Reg. Coeff. = -0.21



I've chosen Social Benefit which has the strongest positive correlation and Asked for valuation which has the strongest negative correlation

Predicting Deal or No Deal for 25% of the pitches

*X_test and y_test
+ Data Preprocessing
Using Python*

Showing the **Complete** **ML Run for** **Logistic** **Regression**



[Link – to All classifiers](#)
python notebook



All Classifiers_Shark Tank India.ipynb

```
# Logistic Regression
# Importing the libraries
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd

# method for uploading data from local
from google.colab import files
uploaded = files.upload()

# Importing the dataset
import io
dataset = pd.read_csv(io.BytesIO(uploaded['SharkTank-Final_used_Ind.csv']))

dataset.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 121 entries, 0 to 120
Data columns (total 36 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   Episode          121 non-null    int64  
 1   Company          121 non-null    object  
 2   Idea              121 non-null    object  
 3   Incurred          121 non-null    float64
 4   SocialBenefit    121 non-null    float64
 5   Tech              121 non-null    object  
 .....
```

Used a Google Colab environment for this, as my organization does not allow other notebooks.
All Python code can be found on my [GitHub](#) page – [IITD-BAO-21-Notebooks](#)

```
# discontinuous columns https://stackoverflow.com/questions/50
X = dataset.iloc[:, [4,12]].values
y = dataset.iloc[:, 6].values

# Splitting the dataset into the Training set and Test set
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25, random_state=0)

# Feature Scaling
from sklearn.preprocessing import StandardScaler
sc = StandardScaler()
X_train = sc.fit_transform(X_train)
X_test = sc.transform(X_test)

# Fitting Logistic Regression to the Training set
from sklearn.linear_model import LogisticRegression
classifier = LogisticRegression(random_state = 0)
classifier.fit(X_train, y_train)

LogisticRegression(random_state=0)

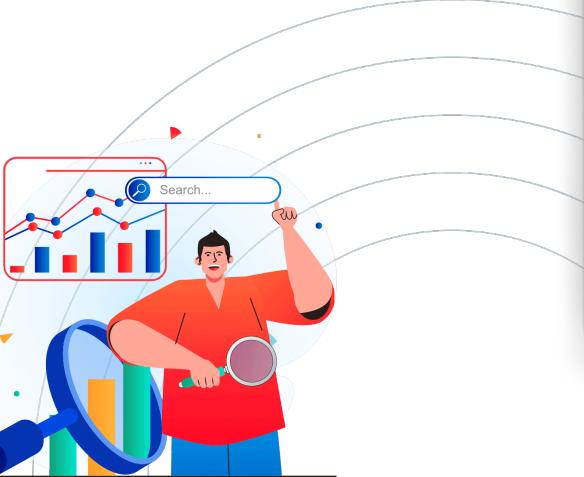
# Predicting the Test set results
y_pred = classifier.predict(X_test)

# Making the Confusion Matrix
from sklearn.metrics import confusion_matrix
cm = confusion_matrix(y_test, y_pred)

print(cm)
[[ 2 15]
 [ 1 13]]
```

Setting Social Benefit, Ask Valuation as X Array and Deal as Y Array through **iloc function**
Using the **StandardScaler** to transform all columns as Social benefit, deal and Asked Valuations are all in different ranges
And fitting the **Logistic regression** model for y-pred

Splitting the data into Train and Test and fitting different ML models



```
#Performance metrics
# Accuracy
from sklearn.metrics import accuracy_score
accuracy = accuracy_score(y_test, y_pred)
print('Accuracy is %f' % accuracy)

# Precision
from sklearn.metrics import precision_score
precision = precision_score(y_test, y_pred)
print('Precision is %f' % precision)

→ Accuracy is 0.483871
Precision is 0.464286

# Visualising the Training set results
from matplotlib.colors import ListedColormap
X_set, y_set = X_train, y_train
X1, X2 = np.meshgrid(np.arange(start = X_set[:, 0].min() - 1, stop = X_set[:, 0].max(), step = 0.01),
                     np.arange(start = X_set[:, 1].min() - 1, stop = X_set[:, 1].max(), step = 0.01))
plt.contourf(X1, X2, classifier.predict(np.array([X1.ravel(), X2.ravel()]).T), alpha = 0.75, cmap = ListedColormap(('red', 'green')))
plt.xlim(X1.min(), X1.max())
plt.ylim(X2.min(), X2.max())
for i, j in enumerate(np.unique(y_set)):
    plt.scatter(X_set[y_set == j, 0], X_set[y_set == j, 1], c = ListedColormap(('red', 'green'))(i), label = j)
plt.title('Logistic Regression (Training set)')
plt.xlabel('MultiVar')
plt.ylabel('Deal or No Deal')
plt.legend()
plt.show()
```

Logistic Regression (Training set)

Deal or No Deal

MultiVar

Red: Deal
Green: No Deal

Precision a measure of how the model is performing is .48
And the **Matplotlib** chart shows the errors the model has made in predicting False positives, i.e. in this case Deal = 1 instead of No Deal = 0

```
# Kernel SVM
# Fitting Kernel SVM to the Training set
from sklearn.svm import SVC
classifier = SVC(kernel = 'rbf', random_state = 0)
classifier.fit(X_train, y_train)

# Predicting the Test set results
y_pred = classifier.predict(X_test)

# Random Forest Classification
# Fitting Random Forest Classification to the Training set
from sklearn.ensemble import RandomForestClassifier
classifier = RandomForestClassifier(n_estimators = 10, criterion = 'entropy', random_state = 0)
classifier.fit(X_train, y_train)

# Predicting the Test set results
y_pred = classifier.predict(X_test)

# Decision Tree Classification
# Fitting Decision Tree Classification to the Training set
from sklearn.tree import DecisionTreeClassifier
classifier = DecisionTreeClassifier(criterion = 'entropy', random_state = 0)
classifier.fit(X_train, y_train)

# K-Nearest Neighbors (K-NN)
# Fitting K-NN to the Training set
from sklearn.neighbors import KNeighborsClassifier
classifier = KNeighborsClassifier(n_neighbors = 5, metric = 'minkowski', p = 2)
classifier.fit(X_train, y_train)

# Naive Bayes
# Fitting Naive Bayes to the Training set
from sklearn.naive_bayes import GaussianNB
classifier = GaussianNB()
classifier.fit(X_train, y_train)
```

SVM Kernel RBF (Training set)

Deal or No Deal

MultiVar

Random Forest (Training set)

Deal or No Deal

MultiVar

Decision Tree (Training set)

Deal or No Deal

MultiVar

KNN (Training set)

Deal or No Deal

MultiVar

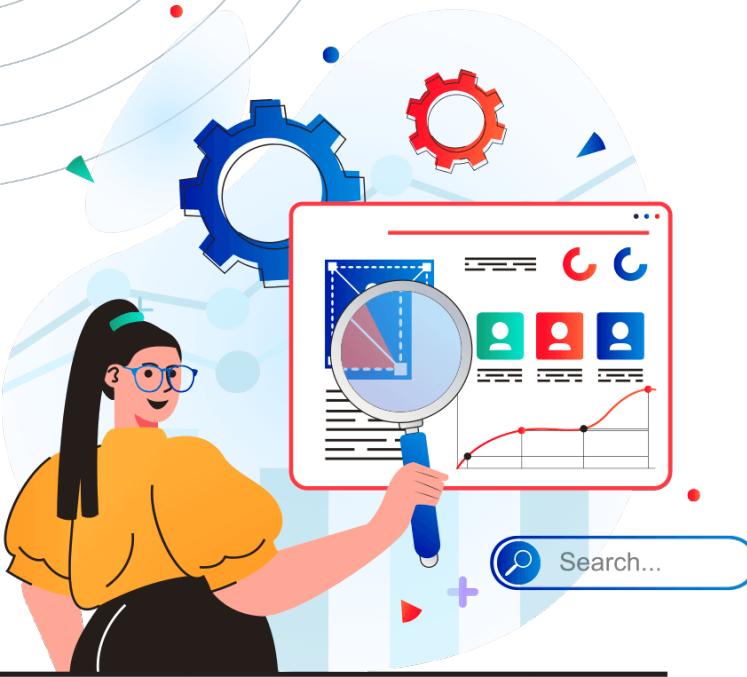
Naive Bayes (Training set)

Deal or No Deal

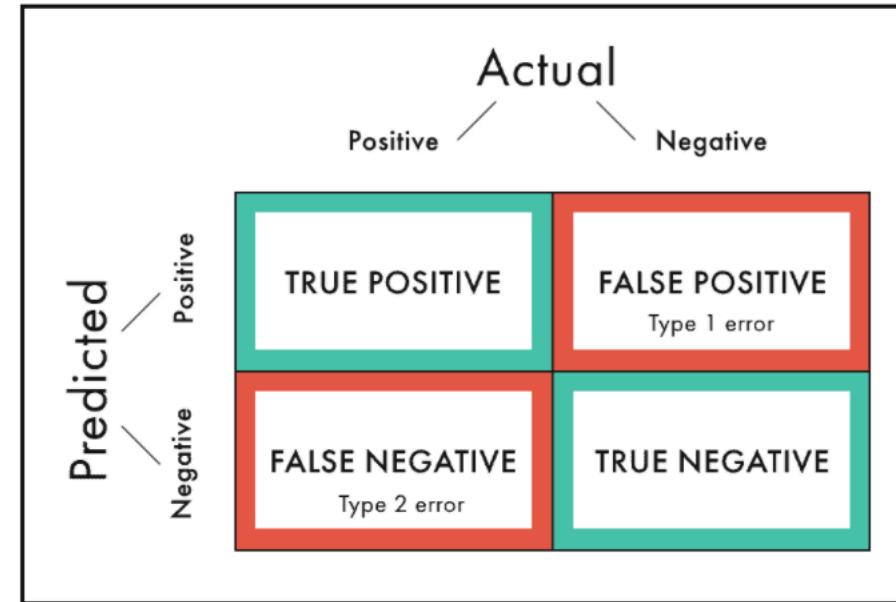
MultiVar

Repeating the fitting, predicting, accuracy and plotting commands for **SVM Kernel**, **Random Forest**, **Decision Tree**, **K-Nearest Neighbors** and **Naïve Bayes**

What does the Confusion Matrix show:



The Confusion Matrix



This is a way of checking how are our `y_predicts` compared to the `y_test` data outputs

These are then used to check the :

Accuracy= (True predictions)/(All Predictions)

Precision= True Positives/(False Positives + True Positives)

Accuracies that I got through Different Models Shark Tank India data

Confusion matrices

print(cm)

```
[[ 2 15]
 [ 1 13]]
```

Logistic
Accuracy: 0.48
Precision 0.46

print(cm)

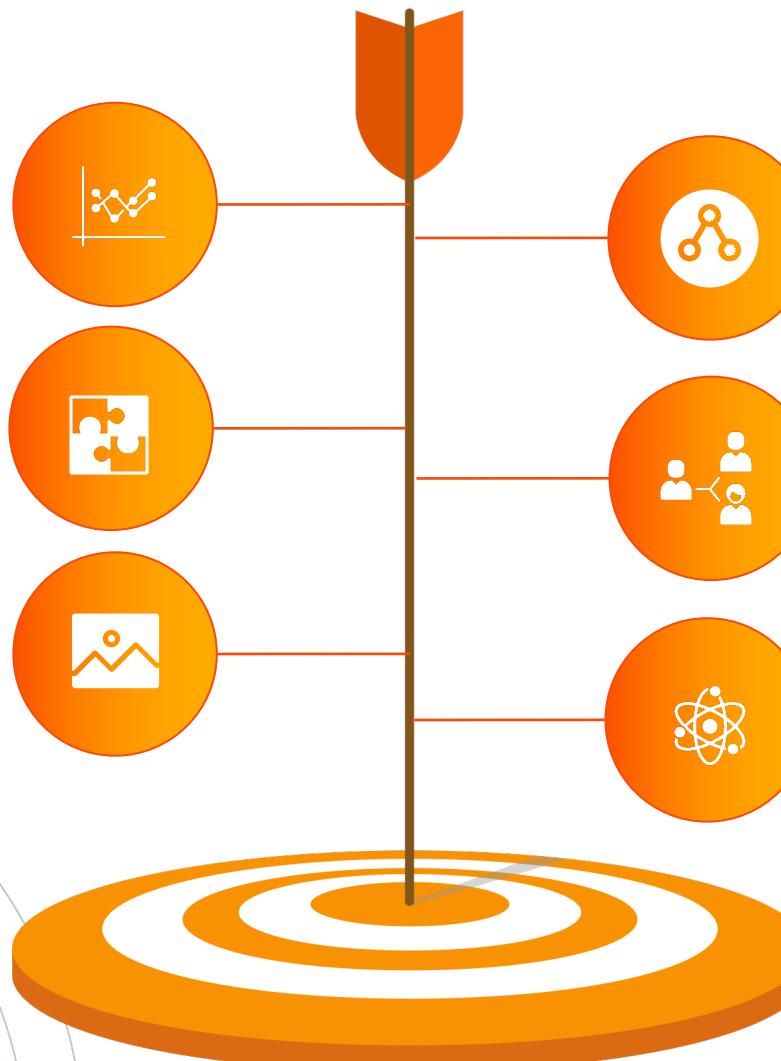
```
[[ 1 16]
 [ 0 14]]
```

Kernel SVM
Accuracy: 0.48
Precision 0.47

print(cm)

```
[[ 6 11]
 [ 4 10]]
```

Random Forest
Accuracy: 0.52
Precision 0.48



Confusion matrices

print(cm)

```
[[ 5 12]
 [ 5  9]]
```

Decision Tree
Accuracy: 0.45
Precision 0.43

K-NN
Accuracy: 0.54
Precision 0.50

print(cm)

```
[[ 3 14]
 [ 0 14]]
```

Naïve Bayes
Accuracy: 0.54
Precision 0.50



Given that the accuracies are at 50%, should we assume that the more data we compile from further seasons, the better prediction we can make?

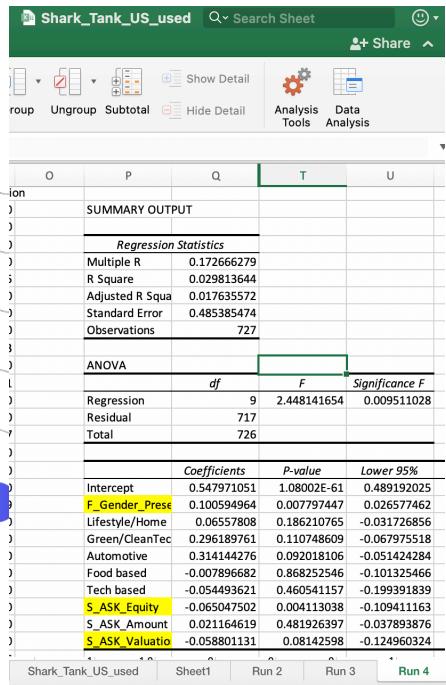
Testing our assumption by using Season 5-11, 2014-2021 data of the US Shark Tank



[Link to dataset](#)



[Link to Python Run on GitHub](#)



Regression & Correlation runs have been added to the project submission folder

Variables for the ML model chosen bases the regression and correlation coefficients again – by multiple rounds of P value based eliminations

```
# Importing the dataset
import io
dataset = pd.read_csv(io.BytesIO(uploaded['Shark_Tank_US_chang...'))
dataset.info()

# discontinuous columns https://stackoverflow.com/que...
X = dataset.iloc[:, [0-1]].values
y = dataset.iloc[:, 2].values

# Splitting the dataset into the Training set and Test set
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(...)

# Feature Scaling
from sklearn.preprocessing import StandardScaler
sc = StandardScaler()
X_train = sc.fit_transform(X_train)
X_test = sc.transform(X_test)

# Fitting Logistic Regression to the Training set
from sklearn.linear_model import LogisticRegression
classifier = LogisticRegression(random_state = 0)
classifier.fit(X_train, y_train)

LogisticRegression(random_state=0)
```

This has 727 cases up from the 121 that we had for India, and this time we used the Female entrepreneurs data and the Equity stake that the pitchers were willing to give up

Accuracies that I got through Different Models for the US Shark Tank data

Confusion matrices

```
[11] print(cm)
[[ 81  0]
 [ 0 138]]
```

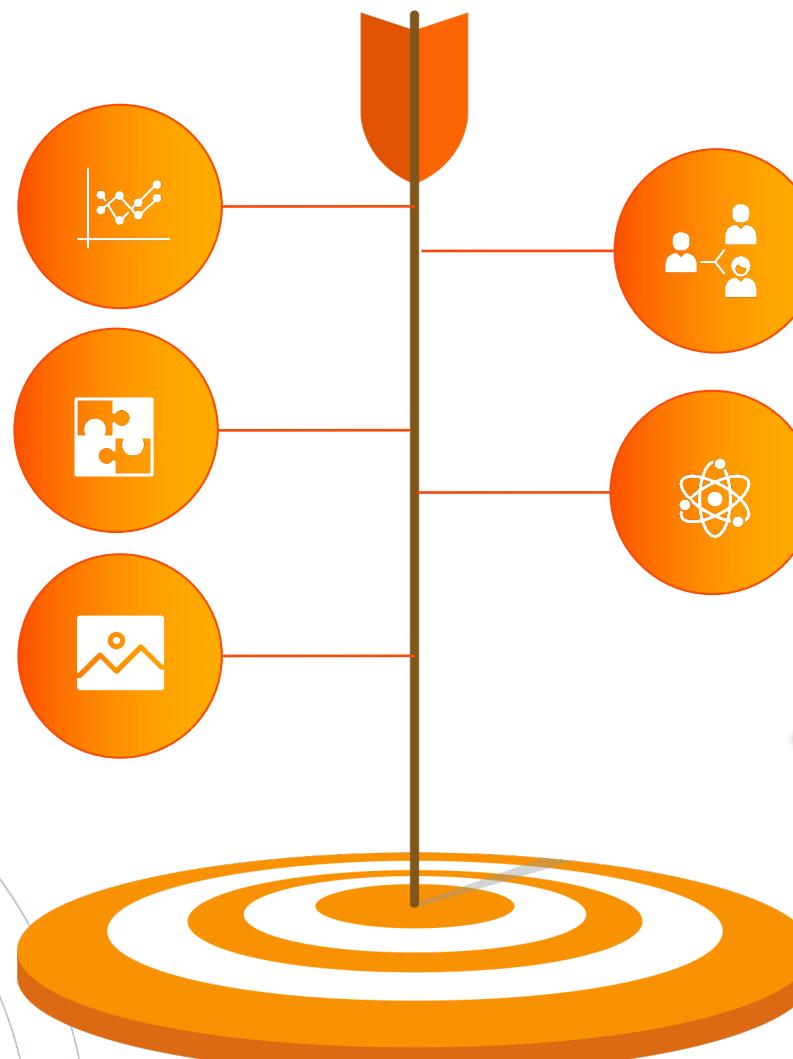
Logistic
Accuracy: 1.00
Precision 1.00

```
[[ 81  0]
 [ 0 138]]
```

SVM
Accuracy: 1.00
Precision 1.00

```
[[ 81  0]
 [ 0 138]]
```

Random Forest
Accuracy: 1.00
Precision 1.00



K-NN
Accuracy: 1.00
Precision 1.00

Naïve Bayes
Accuracy: 1.00
Precision 1.00

Confusion matrices

```
[[ 81  0]
 [ 0 138]]
```

```
[[ 81  0]
 [ 0 138]]
```

By this time and with the dataset that we already have for US Shark Tank – the prediction of all predictions is 100% accurate – specially when taken on the simplistic terms of Equity offered and number of Female Pitchers



However, I discovered that for some datasets if one of the prediction classes is in the ratio of 1:9 with the other class – the prediction model predicts all Negative predicts – resulting in high accuracies.

To solve for this we need to weight down the dominant class:

Also Using R Studio instead of Excel to check Correlations among all variables



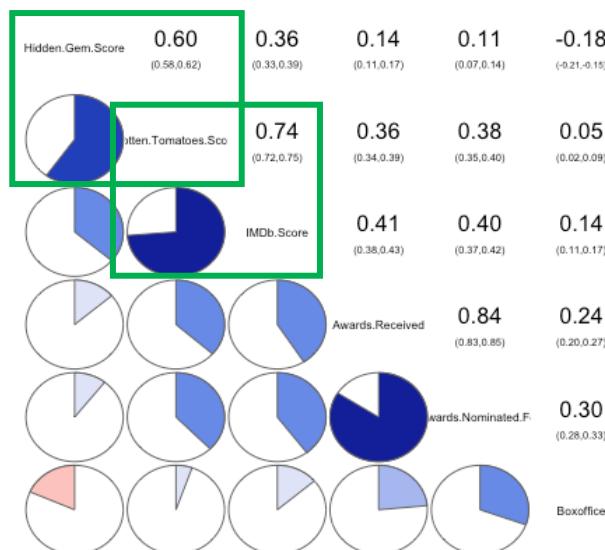
Testing this assumption by using a comprehensive Netflix database to predict:

- + **Hidden Gem scores categorized as “good” ($>=7$) and “not good”($<=6.9$)**

R Studio [Link](#) to R Studio Code on GitHub

kaggle [Link](#) to Dataset

Corrrgram depicting correlations



The **Corrrgram** analysis from **R** gave us a visual view of the Hidden gem scores being highly correlated to **Rotten Tomato** Scores and them in turn related to **IMDB Ratings**, so we discarded the Awards and Box Office data from this analysis

All Python and this R File and output has also been included as attachment in folder

```

Netflix Data.Rmd* 
16 Add a new chunk by clicking the *Insert Chunk* button on
23:4 (Top Level) R Markdown

Console Terminal Jobs
R 4.1.2 · ~/Downloads/Term Project/ 
y.csv")
> summary(netflix)
Series.or.Movie    IMDb.Score   Rotten.Tomatoes.Score
Length:3728        Min. :1.600   Min. : 0.00
Class :character   1st Qu.:6.000   1st Qu.: 38.00
Mode  :character   Median :6.600   Median : 65.00
                           Mean :6.585   Mean : 59.91
                           3rd Qu.:7.300   3rd Qu.: 83.00
                           Max. :9.300   Max. :100.00
> library(dplyr)
Attaching package: 'dplyr'
The following objects are masked from 'package:stats':
> library(corrgram)
> corrgram(netflix)
> corrgram(netflix, order=TRUE, lower.panel=panel.pie, upper.panel=panel.in="Corrrgram depicting correlations")

```

There were 3000+ datapoints but only 10% of these were categorized as “Good” Hidden Gem Score Movies

 [Link to Python Run on GitHub](#)



The screenshot shows a Jupyter Notebook interface with the following code and output:

```
CO Netflix_All Classifiers.ipynb
File Edit View Insert Runtime Tools Help Last saved at 6:36 PM
+ Code + Text

dataset.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3728 entries, 0 to 3727
Data columns (total 8 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   Series or Movie    3728 non-null   object  
 1   IMDb Score        3728 non-null   float64 
 2   Rotten Tomatoes Score 3728 non-null   int64 

# Fitting Simple Linear Regression to the Training set
from sklearn.linear_model import LinearRegression
regressor = LinearRegression()
regressor.fit(X_train, y_train)

LinearRegression()

# Predicting the Test set results
y_pred = regressor.predict(X_test)

# Performance metrics
from sklearn import metrics
print(metrics.mean_absolute_error(y_test, y_pred))
print(metrics.mean_squared_error(y_test, y_pred))
print(np.sqrt(metrics.mean_squared_error(y_test, y_pred)))

1.008717925860649
2.222980152969743
1.4909661810281758

print(cm)
[[1043    1]
 [ 74    1]]

# Performance metrics
# Accuracy
from sklearn.metrics import accuracy_score
accuracy = accuracy_score(y_test, y_pred)
print('Accuracy is %f' % accuracy)

# Precision
from sklearn.metrics import precision_score
precision = precision_score(y_test, y_pred)
print('Precision is %f' % precision)

Accuracy is 0.932976
Precision is 0.500000
```

A scatter plot titled "Hidden Gems Score Prediction" is displayed, showing a positive linear trend with red data points and a blue regression line.

First I ran the Linear Regression since the scores to be predicted were values ranging from 0.0 to 10.0
Then classifying $7 \geq$ as Good and others as not Good, I ran the classifiers

Accuracies that I got through Different Models for the Original Netflix data

Confusion matrices

Linear Regression

**Mean Absolute Error: 1.01
MSE 2.22**

$\begin{bmatrix} [1043 & 1] \\ [74 & 1] \end{bmatrix}$

Logistic Regression

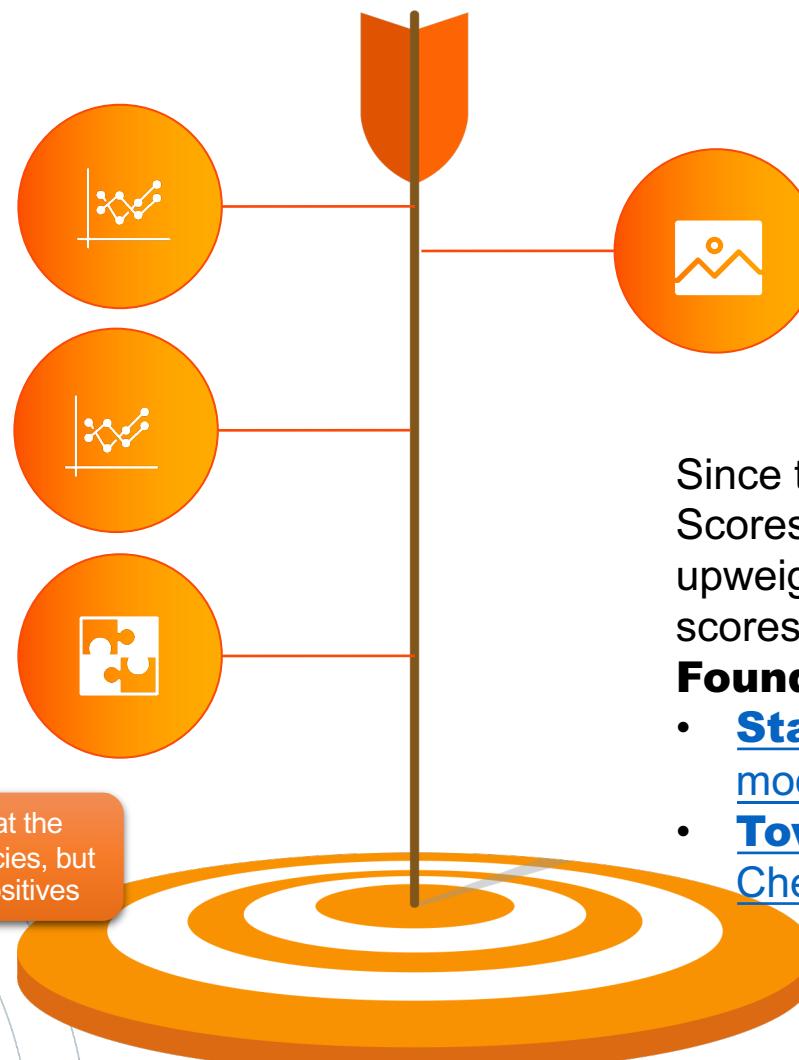
**Accuracy: 0.93
Precision 0.50**

$\begin{bmatrix} [1044 & 0] \\ [75 & 0] \end{bmatrix}$

SVM

**Accuracy: 1.00
Precision 1.00**

The issue here is not that the models have good accuracies, but that all predictions are Positives



Confusion matrices

Random Forest

**Accuracy: 1.00
Precision 1.00**

$\begin{bmatrix} [81 & 0] \\ [0 & 138] \end{bmatrix}$

Since the number of “Good” Hidden Gem Scores are 10% of the data, we need to upweight this and down weight the “Not Good” scores for a reliable predictive model

Found 2 relevant articles on this:

- [**StackExchange**](#): 100% accuracies in model
- [**Towards Data Science**](#) on Medium: Check if Classification model is Overfit



***Data corrected in the excel file
+ Checking if better clusters can
be made off the Hidden Gem
scores than just Good or Bad
classes using K-Means Clustering
Analysis***

I corrected the data in excel making it so that the “Good” scores were at least 25-30% of the total dataset and re-ran the analysis with 1000+ datapoints

 [Link](#) to corrected Python Run on GitHub

I also ran a **K Means** clustering to see if the data can be better clustered rather than just Good or Bad ratings, and sure enough the Elbow method threw up **3 possible clusters with unique centroids**:



The screenshot shows a Jupyter Notebook with the following content:

```
# Fitting Simple Linear Regression to the Training set
from sklearn.linear_model import LinearRegression
regressor = LinearRegression()
regressor.fit(X_train, y_train)

LinearRegression()

# Performance metrics
from sklearn import metrics
print(metrics.mean_absolute_error(y_test, y_pred))
print(metrics.mean_squared_error(y_test, y_pred))
print(np.sqrt(metrics.mean_squared_error(y_test, y_pred)))

1.7223621493440309
4.005375321421638
2.0013433791884983

# Support Vector Machine (SVM)
# Fitting SVM to the Training set
from sklearn.svm import SVC
classifier = SVC(kernel = 'linear', random_state = 0)
classifier.fit(X_train, y_train)

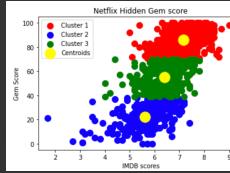
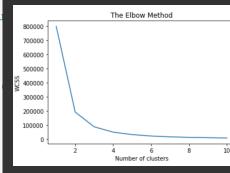
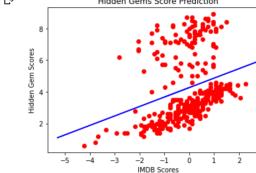
SVC(kernel='linear', random_state=0)

# Random Forest Classification
# Fitting Random Forest Classification to the Training set
from sklearn.ensemble import RandomForestClassifier
classifier = RandomForestClassifier(n_estimators = 10, criterion = 'entropy', random_state = 0)
classifier.fit(X_train, y_train)

# Predicting the Test set results
y_pred = classifier.predict(X_test)

# K-Means Clustering
# Using the elbow method to find the optimal number of clusters
from sklearn.cluster import KMeans
wcss = []
for i in range(1, 11):
    kmeans = KMeans(n_clusters = i, init = 'k-means++')
    kmeans.fit(X)
    wcss.append(kmeans.inertia_)
plt.plot(range(1, 11), wcss)
plt.title('The Elbow Method')
plt.xlabel('Number of clusters')
plt.ylabel('WCSS')
plt.show()

Hidden Gems Score Prediction
```



The revised runs from all model are now not skewed to a Positive predict and the Accuracies are:

Logistic Reg: 84%

SVM 82%

Random Forest 78%

Decision Tree: 78%
K-NN = 81%

Naïve Bayes: 82%

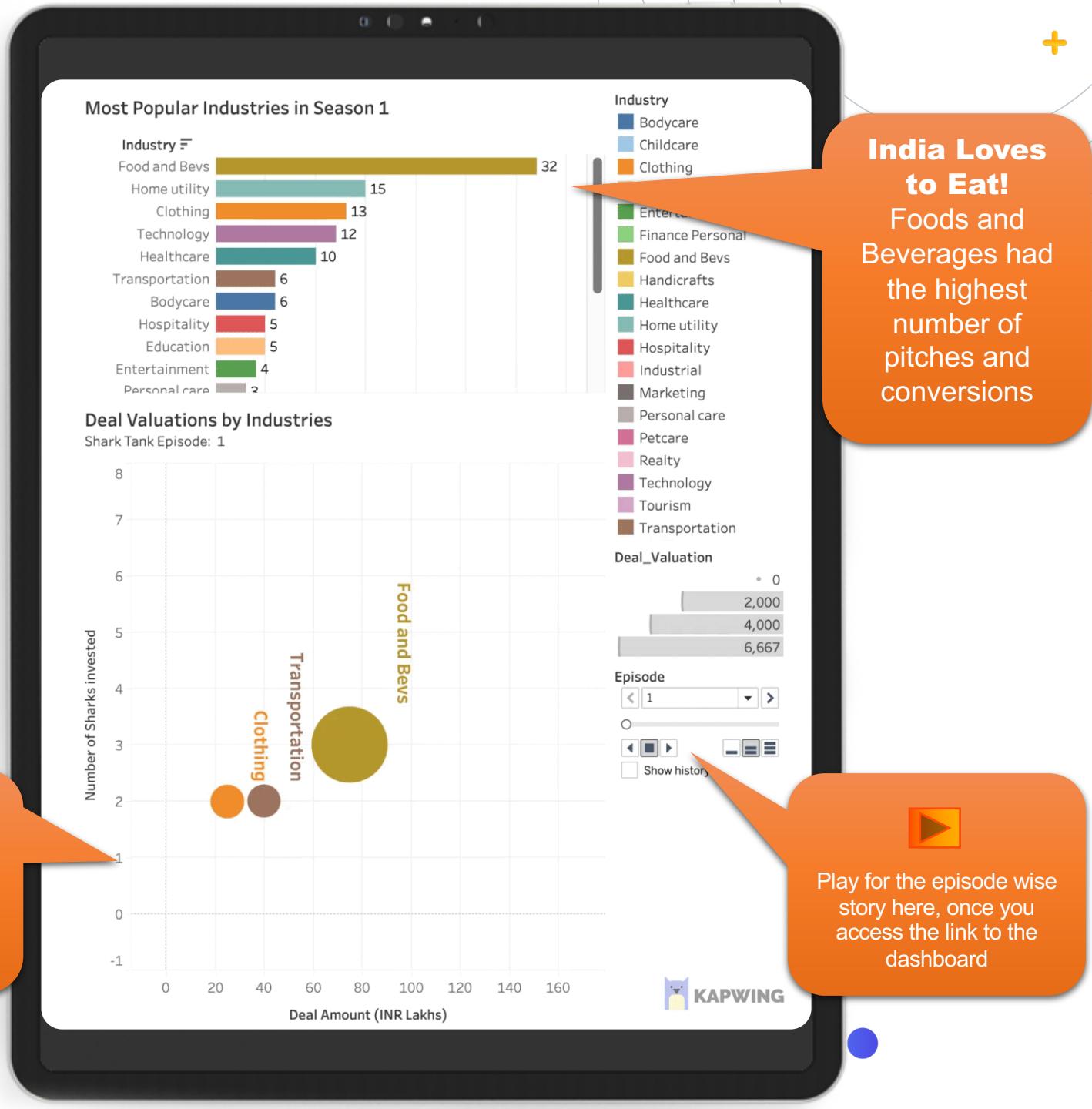
Understanding Shark Tank India data using Tableau Visualizations

A lot can be learned from the data through pure crosstabs and clear data visuals:



[Link](#) to Tableau Public Dashboard with Story Play across all Episodes of the show

The biggest Deal Valuations that got converted came through **Education and Healthcare – Social Benefit matters** when it comes to Start-up investors in India



Some useful Tableau tutorials:

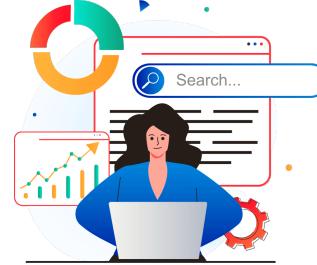
- String based calculations: [Link on Youtube](#)
- Creating a Storyboard: [Link](#)
- Motion charts in Tableau: [Link](#)

Key Nuggets from this Data Analysis



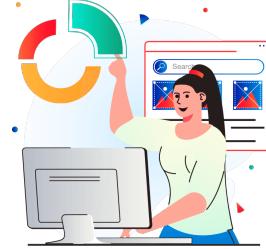
Data Preprocessing needs time

A lot of effort went into **sorting excel files**, **standardizing variables** in Python and understanding **overweighted data**, running **correlation** which help run a better ML model



The Bigger the Sample the Better the run

When we got a bigger dataset either through the US version or Netflix, the **accuracies** of the model **improved**, but do check for overweighted classes



Run all possible Classifiers to know the best

Naïve Bayes is not always the go to model as was my previous assumption. The Netflix data had better accuracy using the Logistic Regression



There is always a better way to cluster or view the data

Do not assume that the existing classes are the only classes, sometimes an **unsupervised cluster analysis** can give better **segments** of the data. Also **visualizations help learn more** about the data than even ML sometimes

***Thanks for taking the
time to read this!***

*Please reach out for further discussion – also I've
tried to keep all relevant material used in this GitHub
Notebook: [IITD-BAO-21-Notebooks](#)*