# Capstone Project-2
# Ted Talk Views Prediction

By-Prabir Debnath

# Points of Discussion

**AI**

# The Problem Statement

We are provided with a labeled dataset on the details of Ted Talk Videos. The task is to explore and analyze the data and to build a regression model for Ted Talk Videos Views Prediction.

# Concept of ML Model

The performance of a machine learning model depends on three factors:

**i. Quality of Data**

(cleaner experiences for better learning)

**ii. Quantity of Data**

(more experiences for better learning)
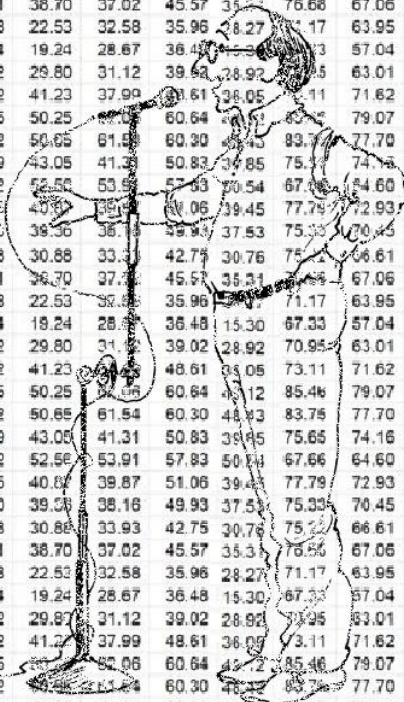
**iii. Quality of Model**

(right model and right hyperparameters for better learning)

# Summary of the Experience Set

Here, the dataset has 4,005 rows, which means 4,005 experiences about Ted Talk Videos and

It has 19 columns, which means each experience is observed along 19 features or dimensions.

# Let's Decode the Experiences !

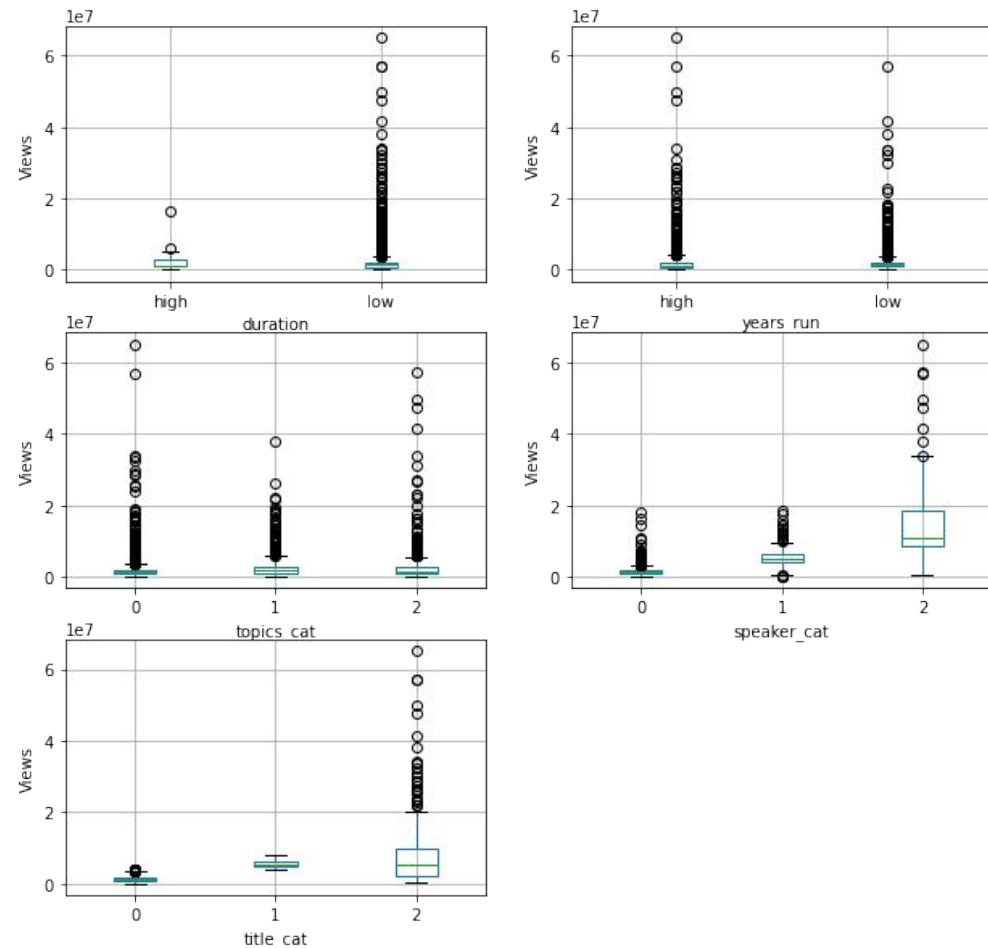# Exploration and Pre-processing of Data

We have done the exploration and pre-processing in seven steps to transform raw data into quality data for our ml model.

1. Connection with the Data
2. First Feelings of the Data
3. Deeper Understanding of the Data
4. Cleaning the Data
5. Treating Anomalies in the Data
6. Final Feature Selection from the Data
7. Preparation of Input and Output Data

AI

# Cleaning the Data

> We have handled 522 null values in 'occupations'  column with imputation, thus there is no loss of data

> We have done feature engineering on 'published_date' to extract 'years_run' feature

> There are very few experiences (only 1-27) for most of the 'native_lang' as compared to 'en'. We have removed those experiences and lost 48 experiences. Thus our model will predict the ted talk views only for english language

> We have encoded 'topics' column with three categories: Highly Favourite:2, Medium Favourite:1, Least Favourite:0

> We have encoded 'speaker_1' column with three categories: Highly Famous:2, Medium Famous:1, Least Famous:0

> We have encoded 'title' column with three categories: Highly Attractive:2, Medium Attractive:1, Least Attractive:0

Boxplot grouped by title_cat

# Treating Anomalies in the Data

We have removed 89 experiences with more than 10 million views as they are exceptional experiences for all of the features.

# Overall Feature Understanding

Here, the distribution of 'views', and 'duration' are positively skewed. Thus we have done log transformation on these features to normalize their distribution.

Correlation between Variables

# Looking for Truly Independent Features

Here, all our input variables are truly independent as all the VIF values are below 10

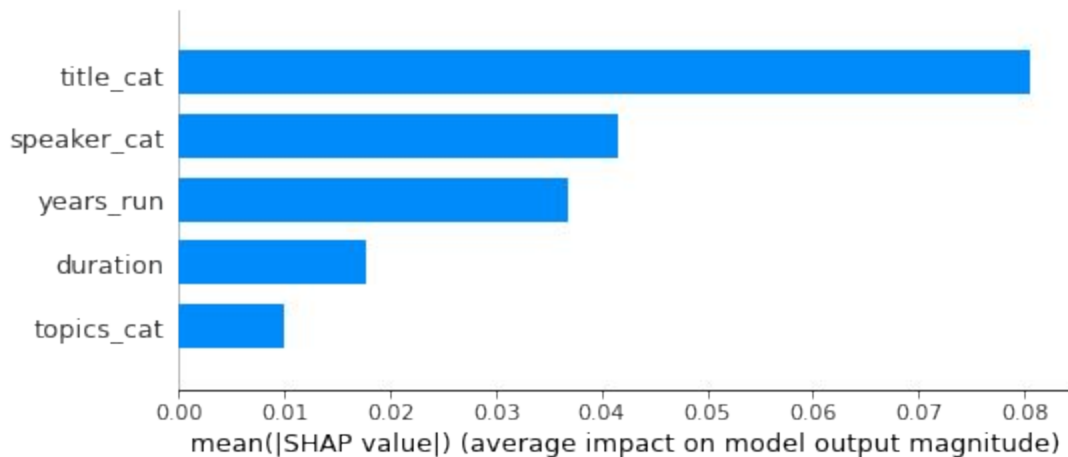| | variables | VIF |
|---|---|---|
| 0 | duration | 3.558988 |
| 1 | years_run | 3.552476 |
| 2 | topics_cat | 1.335001 |
| 3 | speaker_cat | 1.566756 |
| 4 | title_cat | 1.610806 |

# Let's Train the Models !

# Building and Evaluation of Model-1
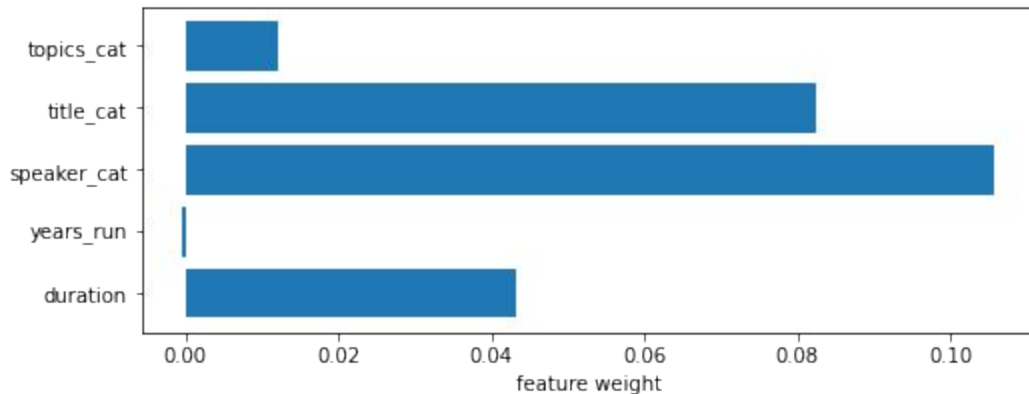
**Final Random Forest Model:**

After cross validation and hyperparameter tuning, the best parameters are {'max_depth': 10, 'max_features': 'sqrt', 'min_samples_leaf': 10, 'min_samples_split': 5, 'n_estimators': 200} (test r2_score is 63% and variance in prediction is 3%)

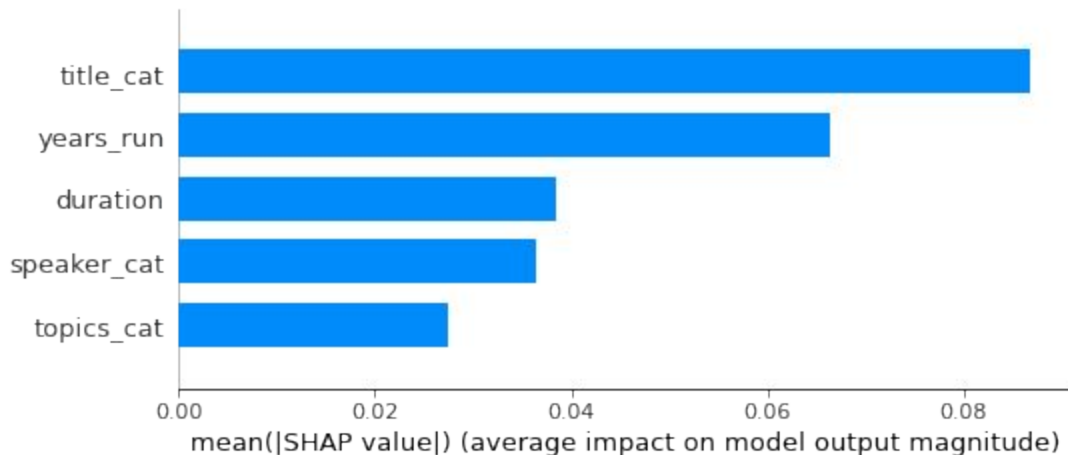# Building and Evaluation of Model-2

**Final Ridge Model:**

After cross validation and hyperparameter tuning, the best parameters are {'alpha': 100, 'fit_intercept': True, 'max_iter': 150, 'tol': 0.2} (test r2_score is 50% and variance in prediction is 4%)

# Building and Evaluation of Model-3
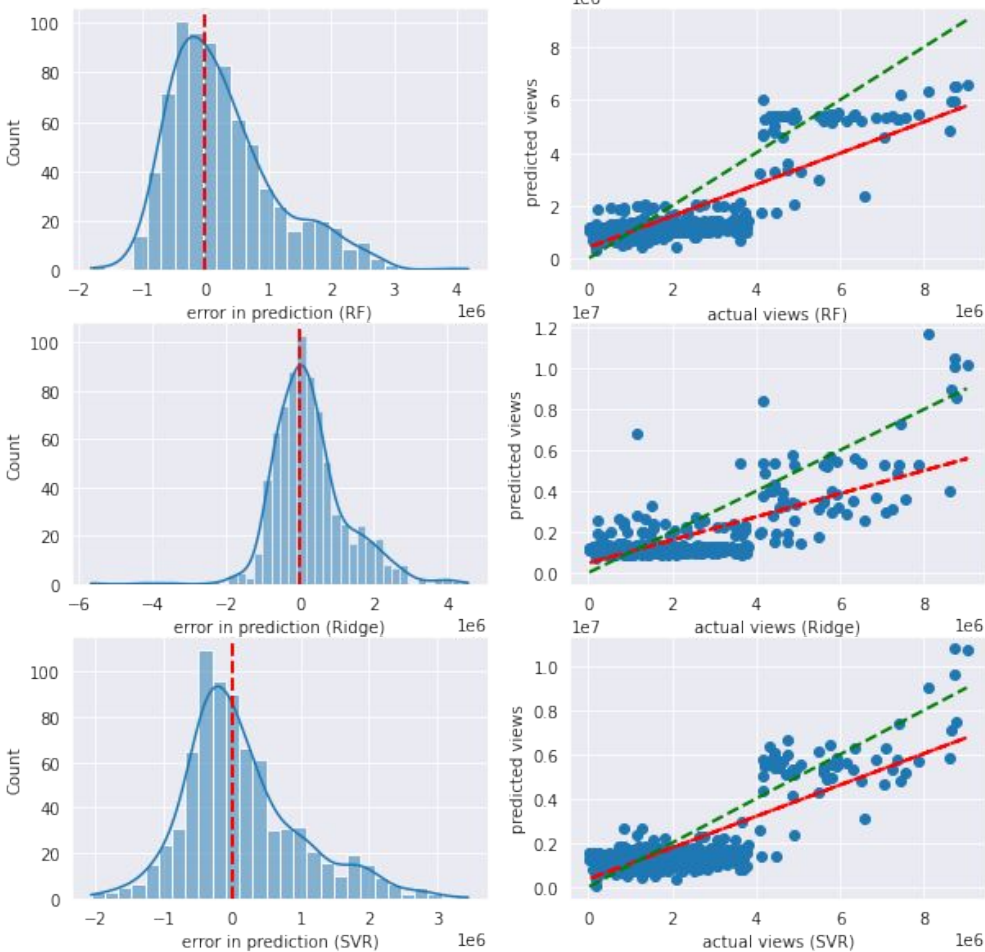
**Final SVR Model:**

After cross validation and hyperparameter tuning, the best parameters are {'C': 6, 'gamma': 0.1} (test r2_score is 66% and variance in prediction is 0%)

RF- Train_R2: 0.6, Test_R2: 0.63, Test_RMSE: 906040

Ridge- Train_R2: 0.45, Test_R2: 0.5, Test_RMSE: 1051682

SVR- Train_R2: 0.66, Test_R2: 0.66, Test_RMSE: 865653

# Final Conclusion

On the basis of the performance study of our three models, we are selecting **Support Vector Regressor** (*the best warrior*) for predicting ted_talk views, as it has the lowest variance in prediction and highest r2_score among all three models

**Thank you !**