



Scalable Knowledge Sharing Platform; An Artificial Intelligence Approach

by: Ibrahim O. Odufowora

Outlines



- Problem Description
- Understanding the data set
- Data Exploration
- Solution Methodology & Implementation
 - Data preprocessing
 - Model building & Parameter tuning
 - Model Validation
- Environment Setup
 - Development
 - Production
- References

Problem Description



- Identification of duplicate questions.
- A page for each logically distinct question.
- Classification Problem
- Supervised Learning
- $P \ll n$, then most models can be used.
- A 'scalable question platform' should provide answers to each distinct question in a single resource location. This will enhance proper resources utilization.

Data Exploration - Word Cloud

What is the step by step guide to invest in share market in india?



What is the step by step guide to invest in share market?

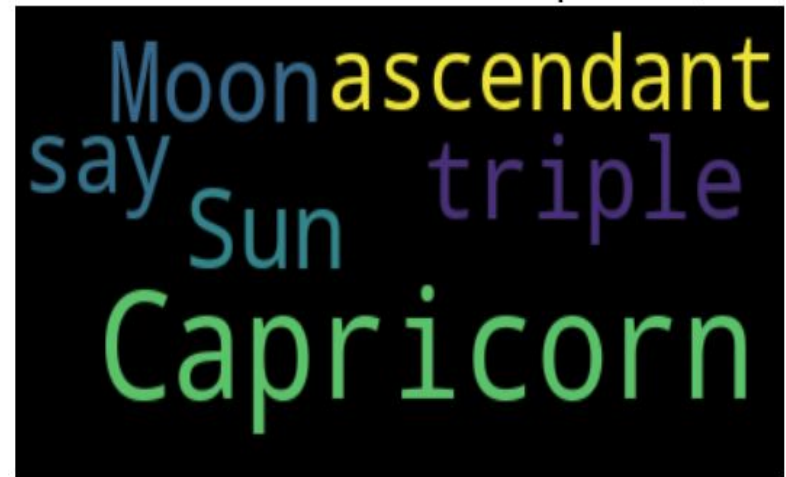


Data Exploration - Word Cloud

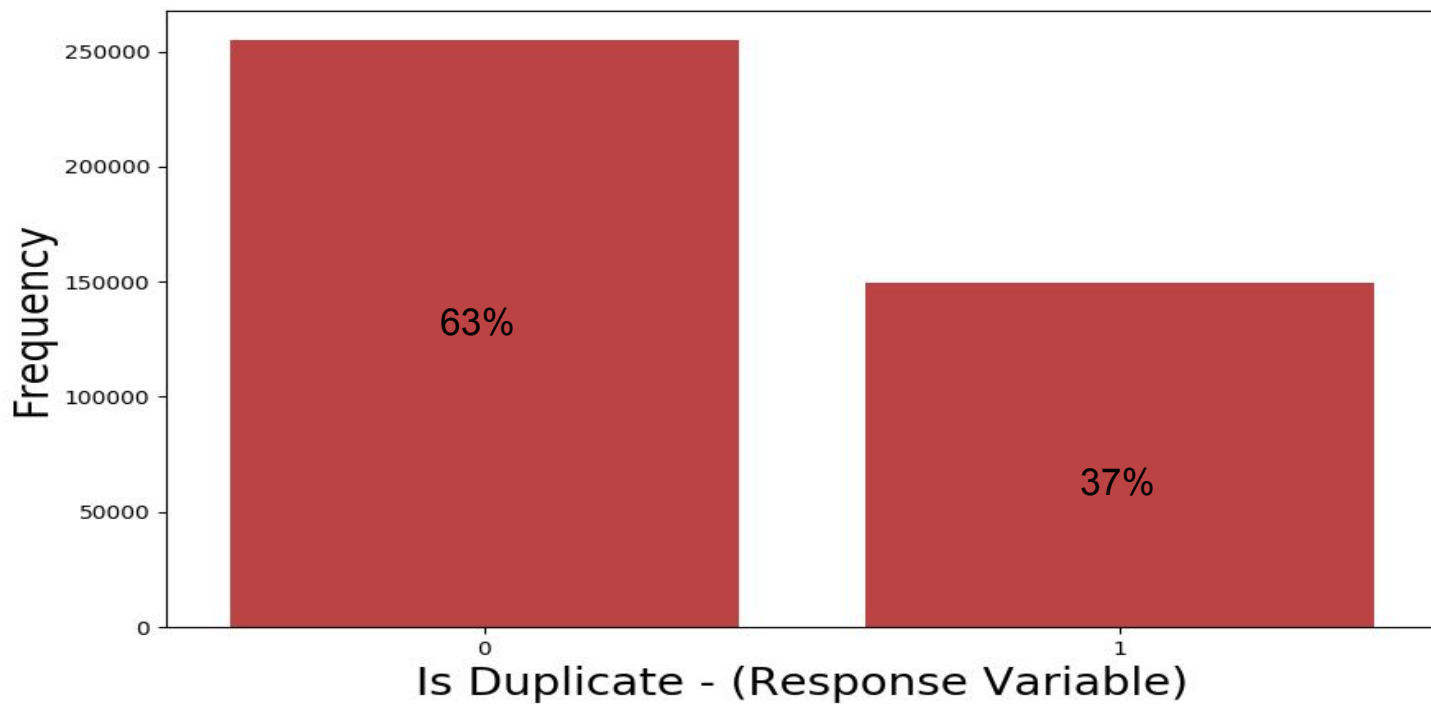
Astrology: I am a **Capricorn** **Sun** **Cap** **moon** and **cap** **rising**...what does that **say** about me?



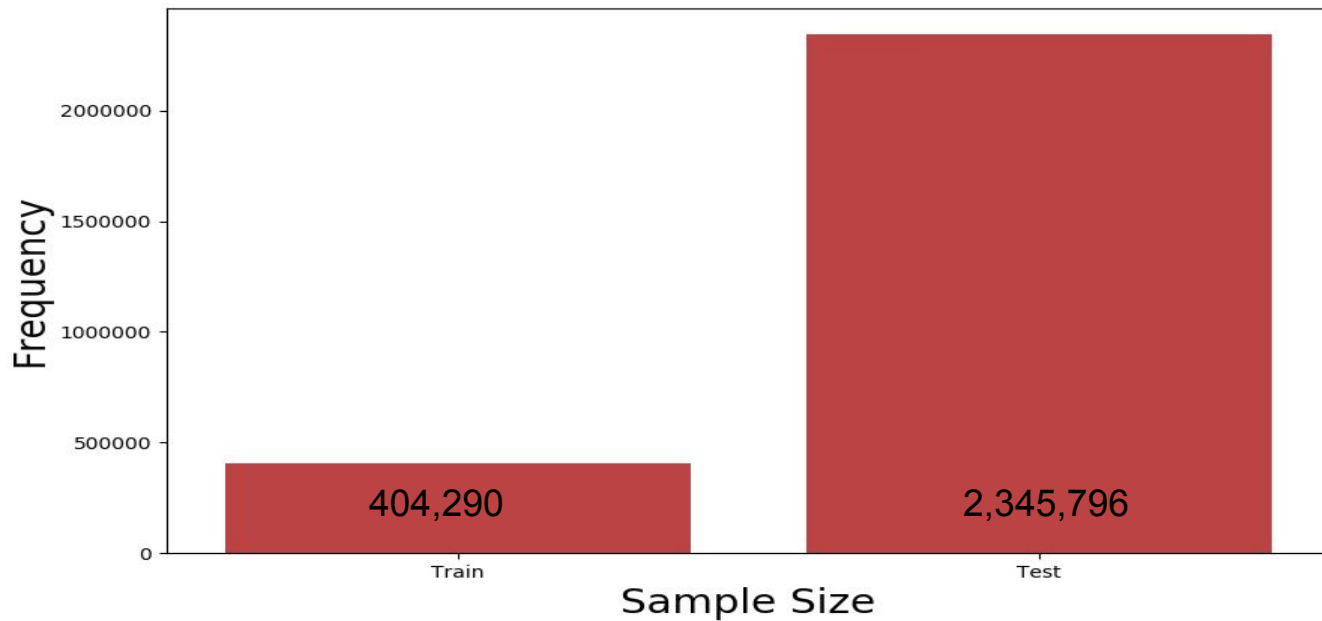
I'm a **triple** **Capricorn** (**Sun**, **Moon** and **ascendant** in Capricorn) What does this **say** about me?



Data Exploration



Data Exploration



Parameter Tuning



- Cross-validation or GridSearch would also be used to select best model parameters.
- Given the imbalance in the frequency distribution of response predictor, it might be suggested to use stratified random sampling method during CV.

Solution Methodology & Implementation



- Methodology:
 - Automate the process of detecting semantically equivalent questions.
 - Direct a user to an existing resource/page if the question already exist.
- Implementation:
 - Build effective and scalable machine learning and natural learning processing models.
 - Data preprocessing.
 - Models building & Parameters tuning.
 - Models Validation.
 - Interface a scalable, online real-time RESTful API with the best model.
 - Load balance the API in order to minimize response time and to avoid overload.

Environment Setup



- Development
 - Build the models on the train set.
 - Use test set for models validation.
 - Interface a Flask RESTful API with the best model.
- Production
 - Move the best model and its optimal parameters to production.
 - Interface a RESTful API with WSGI server with the model parameters.
 - Load balance the API.
 - As more data surfaces, rebuild the model occasionally to improve predictability.

References



- Quora. 2017. *First Quora Dataset Release: Question Pairs*. [ONLINE] Available at: <https://data.quora.com/First-Quora-Dataset-Release-Question-Pairs>. [Accessed 30 September 2018].