# **Project Synopsis Information Retrieval 2022**

Name of Student: Sneegdh Krishnna

Registration No: 199302164

1. Project Title: Identifying Question Pairs Similarity

## 2. Overview of Project (8-10 lines):

Task is to predict whether a pair of questions that are already asked are duplicates or not. This could be useful to instantly provide answers to questions that have already been answered. It is a Real-life NLP problem, used by many companies like Quora, Stack overflow. Currently, Quora is facing the issue of Duplicate questions. For instance, consider a pair of questions as follows:

- Question1: What should I do to be a great Geologist?
- Question2: How can I be a good Geologist?
  These questions asked are very similar with the same intent, but with different wordings and different ways of writing. Hence, we can merge these questions. This can save a lot of time and improves the customer experience.

#### Approach:

- Will perform EDA, check class imbalance, null values, repetitive qts
- Performance Metric: log loss func
- Will do Text-preprocessing -> removing stop words, punctuation, conv into lower case, etc
- Convert into Vectors using BOW & Word2Vec/Glove(SpaCy) (IR tasks)

### Will Take a step further and make it and end-to-end project:

- Do Feature Extraction create new features in order to increase the model accuracy
- ML approach RF & XGBoost algorithms
- Streamlit simple interface
- Will deploy the models with Heroku

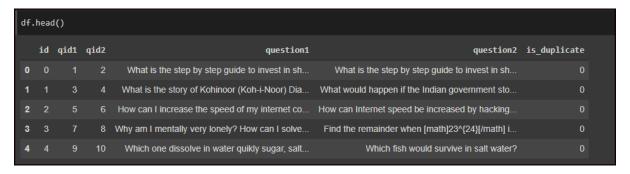
## 3. Any previous work done in same field (if any, mention the detail):

The dataset was taken from Kaggle. Quora had hosted a Kaggle Competition where one had to apply advanced techniques to classify whether question pairs are duplicates or not. Doing so will make it easier to find high quality answers to questions resulting in an improved experience for Quora writers, seekers, and readers.

#### 4. Dataset:

# i. Dataset Description:

- Train.csv file
- Size of Train.csv 60MB
- Number of rows in Train.csv = 404,290
- Contains 5 columns: qid1, qid2, question1, question2, is duplicate
- 'is\_duplicate' is the target label which is 0 for non-similar questions and 1 for similar questions



#### ii. Dataset Features:

It is a binary classification problem, for a given pair of questions we need to predict if they are duplicates or not. The classes are not perfectly balanced, but it's workable.

- Total number of question pairs for training: 404290
- Question pairs are not Similar (is\_duplicate = 0): 63.08%
- Question pairs are Similar (is\_duplicate = 1): 36.92%
- Total num of Unique Questions: 537933

iii. Dataset Link: Quora Question Pairs | Kaggle