

## Problem Definition, Given Resources and Assumptions

The objective of the challenge is to *identify gender* based on a given dataset for training using an appropriate Machine Learning model. The model is required to run through on a fresh, unseen test dataset and predict whether the individual data lines belong to a default on a loan or not. The classes are provided in the training set based on previous observations while they need to be predicted in the test set.

Based on the above information, this is a supervised **binary classification** problem.

### General Assumptions

- Both Training & Testing data have been procured by similar processes/data sources and there is no bias in sampling
- A complete row of observation is uniquely composed of all its features
- A satisfactory **Training** accuracy (**above 90%**) will provide a well generalized model

### About the given dataset

The given dataset shows the following properties:

1. It is composed of 3 different type of files which are to be separately extracted and then joined together based on common keys
2. A few features like '**entities**' have multiple values for each observation. These values need to be unpacked and placed together to form a string which can be encoded later
3. As the number of features is small, we would like to retain as many good features as possible for prediction
4. The combined data set needs to be deduplicated to produce unique observations of interest
5. As the target variable – **gender**, is quite imbalanced we need to resample the records to have better data set which our algorithms can learn on

## Solution Methodology

### Technical Specifications & Environment:

**My Laptop specifications** → RAM, processor and Operating System

#### Windows edition

---

Windows 10 Pro

© 2018 Microsoft Corporation. All rights reserved.



#### System

---

Processor: Intel(R) Core(TM) i3-4005U CPU @ 1.70GHz 1.70 GHz

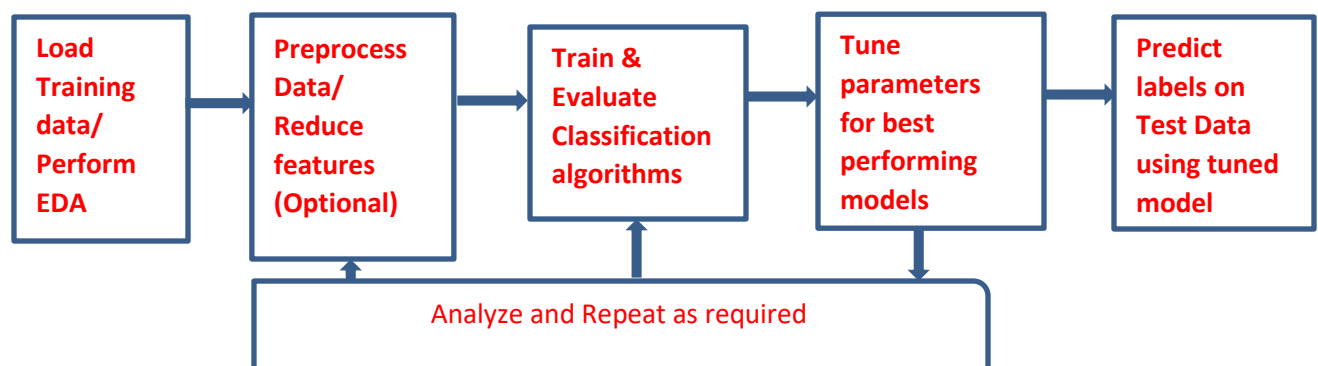
Installed memory (RAM): 4.00 GB

System type: 64-bit Operating System, x64-based processor

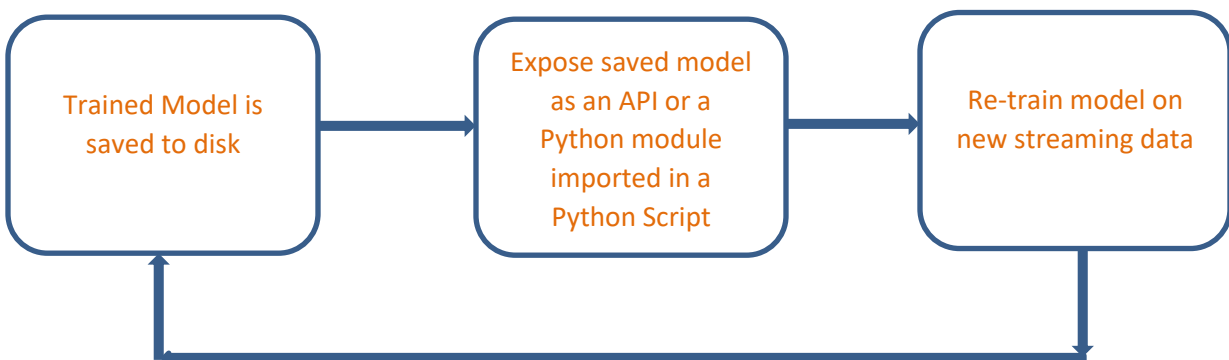
### Programming environment/libraries

Solution was created on Python3 kernel on *Anaconda Jupyter Notebook* environment. Scikit Learn libraries were used for Machine Learning and Pandas for data manipulation. Other supporting modules for time, plotting etc. were also imported as required.

## High Level Process Architecture Diagram



### *Training a first cut model on training data*



*Continually enrich the created model on new/streaming data*

Data Science is an experimentative/stochastic process and a classical binary classification model follows the steps below:

## 1. Pre-processing

Load Training Data → Deduplicate records → Check missing data and impute (most frequent value) → Encode categorical data (label encoder, mapping categories to numbers) → Scale continuous variables to be in range (0, 1) → Resample (up sample minority class and down sample majority class) based on class distribution of target variable → *Perform feature engineering to remove redundant predictors (optional)* → Segregate predictors versus target variables (X, y)

All steps above are standard and apply to any machine learning project. However, choosing the methods in each step is a matter of debate and in this case, they were chosen based on both experience and the data properties in question.

I have chosen not to use the **Pipeline** class from *scikit learn* as I wanted to customize the individual pre-processing steps and view the intermediate outputs if necessary. Production models should use Pipelines as it helps automate the entire process

## Feature engineering → Feature selection

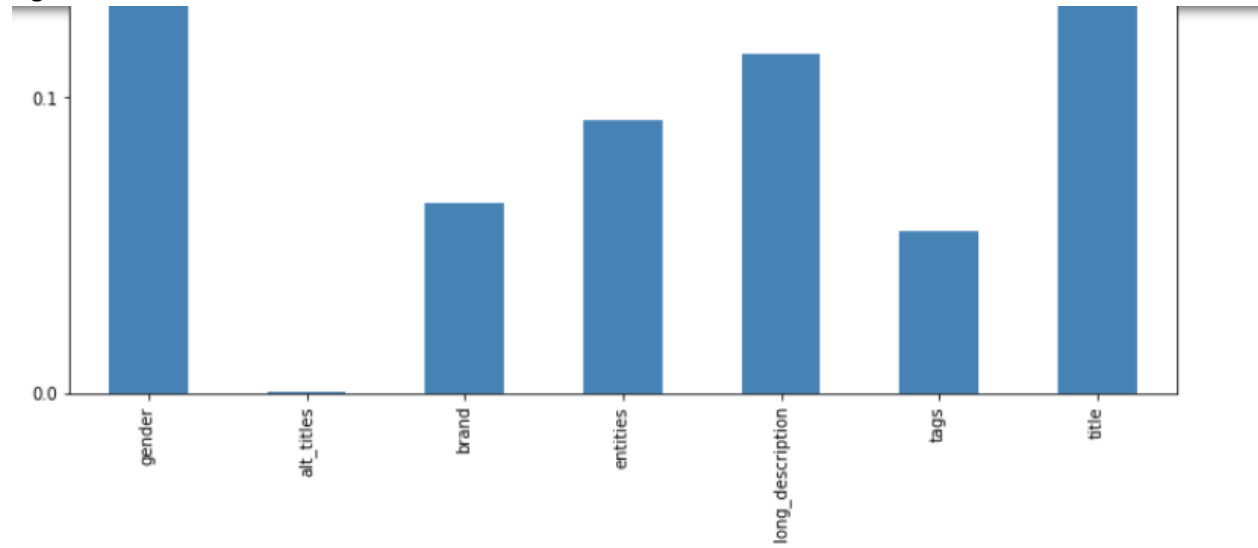
As per [Wikipedia](#), feature selection is a specific discipline of feature engineering which aims to find the most relevant features from a given set based on certain criteria. It is suggested that it be *applied when number of features/predictors is greater than 10*. In this case, we have 29 and hence it fits the case. Less features can also help reduce training time in most cases. It is still not an exact science and requires human interpretation. However, automated feature engineering using *evolutionary algorithms* is possible but doesn't guarantee the most optimal result always.

In my implementation, I have used **Tree Based Feature Importance using ExtraTreesClassifier** along with **correlogram of features**

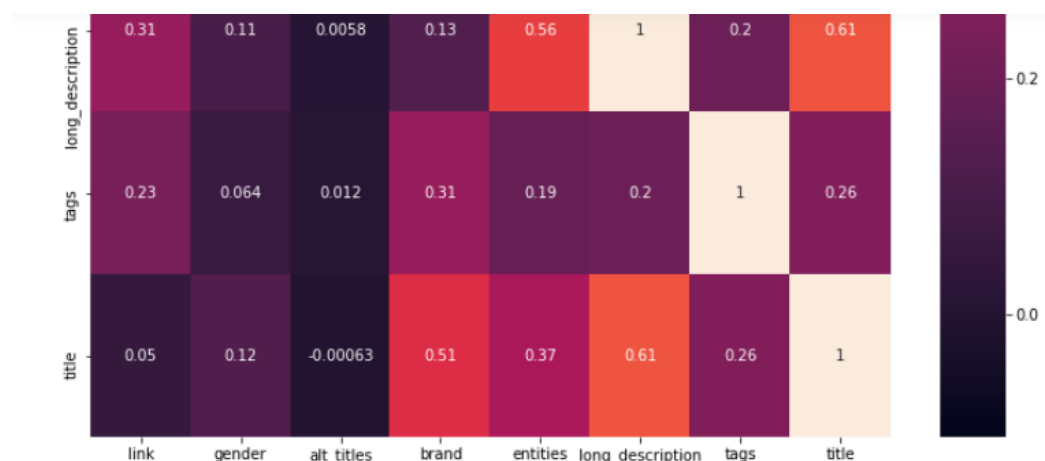
## Four step feature engineering/dimensionality reduction

**Step 1:** Remove any unique key or non-random variables *or noise variables*

**Step 2:** Look at a feature importance graph generated by a Tree based algorithm and remove the least significant



**Example:** From the graph below we can easily spot **alt\_titles** as a candidate for removal as it has almost 0 importance



**Step 4:** look at a correlation graph and remove features which are a) very lowly correlated with the TARGET or b) are mutually strongly correlated

**Example:** **brand** is strongly correlated to **title** but the latter has higher correlation with TARGET. Hence, we can remove **brand**

## 2. Training classifiers

Train → check score → keep model with best accuracy → perform a 10-fold F1 score Validated score for each model and pick the best

In order to evaluate training accuracy a **10-fold cross validated F1 score** is calculated for each of the 3 different types of classifiers and the best is retained. The best model is then tuned for best parameters.

### 3. Predict target class

Finally, the best model is used to predict the classes and a file a per specifications is produced for evaluation

## Obvious challenges and possible shortcomings of the solution

Training and tuning models is a **memory and time intensive** process. Some challenges and shortcomings as below can significantly reduce model performance and reliability of results

- Lack of domain expertise: Though data science is a general science, having knowledge about the domain of data helps to make knowledgeable decisions like the kind of imputation required, the expected number of features to be retained and the accuracy of prediction (precision versus recall) needed for the problem
- Making the right judgements: Most data science coding is ad-hoc at this stage and even though we have automatons like creating a pipeline process or using Grid Search or Random Search of the hyperparameter space, we still can't decide on the best approach without consulting the web. There is surely no 'perfect solution' till we try everything!!
- Overfitting: In our effort to get the best training accuracy, we might fit our model too close to the training dataset and thus increase variance. This is the reason testing accuracy post evaluation is always lower than training accuracy obtained.
- Avoiding memory Intensive algorithms: In the given problem, I couldn't run some more memory intensive algorithms like KNN Classifier for benchmarking. **Training Neural Networks was costly as well.**
- Execution Speed: **Speed** is also an issue as the total runtime of the notebook is around **40-50 minutes**
- Preprocessing mismatch: As data properties/distribution might vary between test and train the same set of preprocessing steps might not be ideal for both. For example, some new, unseen categories in the test data might not be recognized by the trained model and it might produce an incorrect classification

## Questions posed by challenges

- ➔ Is training accuracy a good indicator of comparable testing accuracy? This can **prevent overfitting** before it's too late
- ➔ Can we **modularize pre-processing** for a domain-based dataset (like a bank, insurance etc.)? This can help standardize process for a problem type