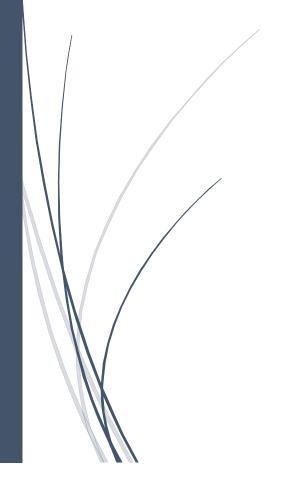
10/27/2020

# STATEMENT OF WORK-1002

MARCOS BITTENCOURT



Roshandeep Singh Saini 100766638

# Contents

1.1 Executive Summary	2
1.2 Problem Statement	2
1.3 Analytics Rationale Statement	2
1.4 Data	2
1.5 Output Variable Structure	3
1.6 Data Analysis Approach	3
1.7 Assumptions And Constraints	3
1 & Project Plan	1

# 1.1 Executive Summary

An analysis is conducted to predict whether a given dish is vegetarian or not, depending on the ingredients used in the dish. In simpler terms it is expected that that a dish with same ingredients, would result in a similar dish. Predictive modelling techniques would be used to find a solution which best classifies the dish as vegetarian or non-vegetarian in terms of diet.

### 1.2 Problem Statement

The problem statement is "To classify a dish as vegetarian and non-vegetarian based on the ingredients".

Based on the ingredients used to prepare an Indian dish, it can be classified as a part of a vegetarian diet or a non-vegetarian diet. To gain insights on how dishes with similar ingredients can be classified similarly or differently based on region, flavor, and course.

# 1.3 Analytics Rationale Statement

The rationale of the analysis is to accurately predict whether a dish can be called a vegetarian dish or a non-vegetarian dish depending on the contents, flavor profile, or the course of meal. Dishes with similar ingredients can be profiled as a different course of meal as per Indian cuisine. Hence, the aim is to quantify the similarity and disparity.

### 1.4 Data

The data has been acquired from Kaggle open datasets. It is a raw dataset named "indian\_food", which represents entirety of the testing data for August 5<sup>th</sup> to October 5<sup>th</sup>, 2020. It contains details about Indian food, the features used to classify them.

It has 255 records. These variables are: -

- Independent Variables- Actual measurement parameters of an Indian Dish which are name, ingredients, prep\_time, cook\_time, flavor\_profile, course, state, region.
- Dependent Variable- Classification of dish as part of the diet(vegetarian and non-vegetarian).

# 1.5 Output Variable Structure

The output variable will have data divided in 2 classes "Vegetarian" and "Non-Vegetarian" as this is a binomial classification problem. The hypothesis is that the ingredients contained in Vegetarian and Non-Vegetarian dishes would have close cosine similarity.

# 1.6 Data Analysis Approach

The methodology used will be to develop a classification model using machine learning approach to segregate Indian dishes into "vegetarian" and "non-vegetarian". It is a scenario of Supervised Learning. This would require classification algorithms SVM, Decision Tree and an advanced modelling through TensorFlow(Keras) etc. depending on the complexity and the type of data received. This would be determined in the EDA and modelling phase of the project. The software tools used will be: -

- Python for EDA (Exploratory Data Analysis), Data Cleaning, Model Building and Testing
- Jupyter Notebook IDE for development
- Tableau/Power BI for detailed visualizations

# 1.7 Assumptions And Constraints

The assumptions made for the given problem statement are: -

- The dataset acquired from Kaggle is a valid dataset, it is assumed that the records correspond to actual features and accurate measurements/description of the Indian dishes that were considered.
- All the features in the dataset corresponding to "indian\_food", acquired from Kaggle are required for the analysis and the features are independent of each other and are collectively required for tackling the problem statement.
- It is an assumption that the model is restricted to a binomial problem(Vegetarian and Non-Vegetarian) rather than multi class. Thus, this model can only be used to this specific problem statement.

The constraint in the entire scenario are: -

• No more data can be obtained. For the analysis and model building is limited to the given data only. This is a hard constraint that we have no control over.

- The above analysis is limited to the number of features available in the dataset, no more features can be obtained or used for the modelling purpose. This is a hard constraint that we have no control over.
- The classification model being built is a binomial classification problem, thus it cannot be used for any other study. The entire data set is classified in terms of "Vegetarian" and "Non-Vegetarian" Indian dish, so the model answers specifically to the above problem statement. This is a soft constraint, but it also depends on the problem statement.

# 1.8 Project Plan

The table below contains the project tasks and their estimated completion dates.

Phase	Task	Deliverables	Delivery Date
Project Organization	Setup	Setting up a GitHub repository for the project deliverables.  Setting up the project structure.	October 28, 2020
Business Understanding and Problem Discovery	Statement of Work(V1)	Developing a problem statement and a rational statement for the business problem.  Refining the problem, identifying the accurate requirements, assumptions, constraints for the business problem.  Acquiring the data sources relevant for solving the business problem, analysing the data source for restrictions and limitations.  Defining parameters to test the results and setting up a benchmark for the analysis.	November 1, 2020
Data Acquisition and Understanding	Data Acquisition and Understanding, Statement of Work(V2)	Performing preliminary EDA on the data set for basic statistics.  Identification of the steps required for feature engineering, normalization, scaling, cleaning, outlier removal, conversion to numerical features from strings and categories, correlation.	November 23, 2020

		Data preparation for a base model prediction. Selection of models required for advance classification.  Detailed documentation of the complete EDA process.	
ML Modelling and Evaluation	Advanced Predictive Modelling	Developing the base model and analyzing the results. Evaluation of base results.  Development of architecture and pipelines for advance classification algorithms. Using optimization techniques for further improvement and evaluation of results.  Prototyping the model, tweaking the parameter to attaining the proposed results for the evaluation metric selected.	November 23, 2020
Delivery and Acceptance	Deployment of Software Pipeline	Development of a software pipeline to consume the results from the model solution.  Providing and end-to-end solution for the final user to be consumed.  To submit the results, solution, and a detailed presentation of the entire analysis to the instructor.	December 18, 2020