



**THE WORLD BANK**

**Industrial Project:  
An alternative  
approach to PPP**

Group 2: Mayank Loyalka, Preet Oza,  
Adhyatma Gautam and Tzuchi Chiu

## Table of Contents

<b>INTRODUCTION .....</b>	<b>2</b>
<b>APPROACH .....</b>	<b>2</b>
<i>Data Selection.....</i>	2
<i>Process Pipeline.....</i>	2
<b>METHODOLOGY &amp; ANALYSIS.....</b>	<b>3</b>
<i>Level 1 Introduction .....</i>	3
<i>Level 2 Introduction .....</i>	4
<i>Level 2 Algorithms.....</i>	4
<i>Level 3 Introduction .....</i>	5
<i>Level 3 Algorithms.....</i>	5
<b>USER INTERFACE OVERVIEW .....</b>	<b>6</b>
<b>FUTURE OVERVIEW.....</b>	<b>6</b>
<b>CONCLUSION .....</b>	<b>6</b>
<b>GLOSSARY.....</b>	<b>7</b>
<b>APPENDIX.....</b>	<b>8</b>

## INTRODUCTION

Our business objective was to create an alternative approach to traditional PPP (purchasing power parity) method under ICP program of the world bank, for the category of house food consumption goods and items, by eliminating the human factors utilizing different machine learning and text mining techniques available, performing real time analysis on the data, and presenting all our findings into an interactive user interface (UI).

## APPROACH

### Data Selection

To achieve the business objective, we wanted to test our approach on a smaller subset of food data and successfully scale it to all the food categories in the future. For this we decide to run the analysis on bakery products. This process was further categorized into three levels.

Level 1: Classification of bakery vs non bakery item.

Level 2: Classification of biscuit vs non-biscuit.

Level 3: Further categorization of biscuit vs non-biscuit based on their name.

### Process Pipeline

Our process pipeline can be divided into 3 phases as follows:

1. Data collection and data preparation
2. Model selection
3. Model prediction

For phase 1, we built on the web scraper of Cal Poly DxHub, which scrapes the product name from websites and then stores the data in DynamoDB. On the extracted data we performed the following steps Data Cleaning, Building N Grams and Computing TF-IDF.

For phase 2, to identify the best classification algorithm to differentiate between products we used the computed TF-IDF to train machine learning models like Decision trees (bagging trees), K-Nearest Neighbor (KNN), Linear Discriminant Analysis (LDA), Logistic regression, Neural Network and Random Forest. Furthermore, to understand which metric will determine the algorithm's success, we tested all the models against the accuracy and ROC (Receiver operating characteristics) metric.

For phase 3, After testing all the machine learning models successfully for each level, we used the model with the highest roc and accuracy values. We integrated the respective chosen models in our user interface to make future products predictions.

## METHODOLOGY & ANALYSIS

### Level 1 Introduction

Once we had a good understanding of what our data set was, we started to work on our Level 1 modeling. The goal here was to create a model that just by looking at a product name can classify if an item is Bakery or Non-Bakery. Since we already had data on Bakery products, we needed some data on Non-Bakery products as well. With the help of Non-Bakery item data, we can gain a better understanding of our Bakery data and create better models. So, we decided to include Rice data as our non-Bakery data just for our level 1 modeling purpose. When we started to perform data understanding. For Level 1 we decided to perform Text analytics. To perform text analytics, we had to make tokens of our product name and then performed the following 4 major steps

**Step 1 - Data Cleaning:** We had to remove the following words as they do not carry any prediction power.

- Stop words: Stop words are a set of commonly used words in any language. For example, in English, “the”, “is” and “and”, would easily qualify as stop words.
- Punctuations, Numbers and Symbols: Punctuation marks like (,?,” , Symbol and numbers like 1,2,3 needed to be removed.
- Rice and Cookie: Upon further analysis we noticed that words like Cookie and Rice were the most frequent words in each category. Being the most frequent words, they were the most dominant attribute in making predictions. We decided to remove them as we did not want to create a model that just relied on 1 word for making predictions.

**Step 2 - Identifying the most frequent words:** Our next goal was to identify the most within each category and create a document frequency matrix using them

- Bakery Product: In Bakery product we had 200 different words like Chocolate, Sugar, Sandwich, Chip which appeared a lot (Figure 1).
- Rice Product: In Rice product we had 300 different words like Brown, Grain, Chicken, Vegetables which appeared a lot (Figure 2)

**Step 3 - Creating n grams:** We created N grams with N=1 and N=2. This helped us to create a model that will make better predictions. After that we computed their prediction power with the help of TF-IDF.

**Step 4 - Singular Value Decomposition:** Once we calculated TF- IDF we performed Singular Value Decomposition to reduce our attributes from 1900 to 300. Once we had a reduced number of attributes, we ran 6 different machine learning algorithms (Figure 3) and tested all of them against ROC and accuracy metric. We decided to go with Random Forest as it had an accuracy and ROC of 89%.

## Level 2 Introduction

Our Level 1 model was able to successfully classify products like Oreo Biscuit, Maria Biscuit, Chocolate ice cream into Bakery products. However, we needed to come up with further classification of those products. Since almost 60% of our data contained details about Biscuit products, we decided to make further classification of those products as Biscuit or Non-Biscuit.

## Level 2 Algorithms

Our goal for Level 2 was to create a model that just by product name can classify if a product is Biscuit or Non-Biscuit. To perform text analytics, we made tokens of our product name and performed the following 4 major steps:

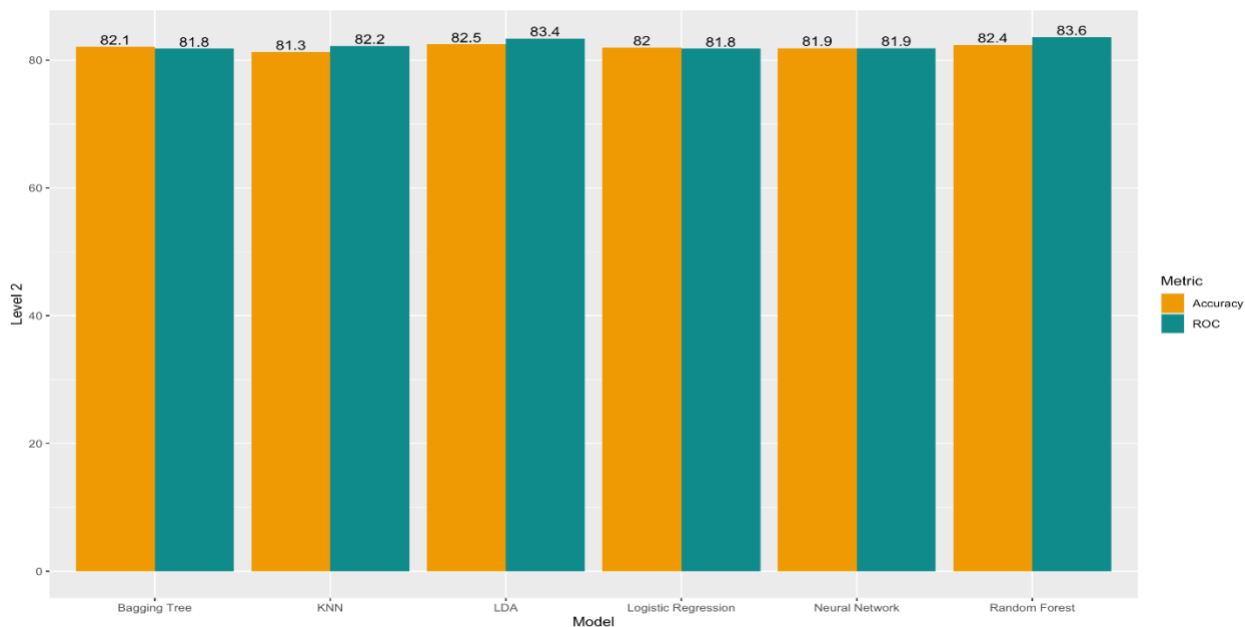
**Step 1 - Data Cleaning:** Since we were already using the same data set, we did not have to go over the same data cleaning process again.

**Step 2 - Identifying the most frequent words:** Our next goal was to identify the most within each category and create a document frequency matrix using them

- Biscuit Product: In Biscuit products we had 150 words like Chocolate, Sugar, Sandwich, Chip which appeared a lot (Figure 4).
- Non-Biscuit Product: In Non-Biscuit product we had 100 words like ice, chocolate, dough, cream which appeared a lot (Figure 5).

**Step 3 - Creating N grams:** We created N grams with N=1 and N=2. This helped us to create a model that will make better predictions. After that we computed their prediction power with the help of TF-IDF (Figure 6).

**Step 4 - Singular Value Decomposition:** Once we calculated TF-IDF we performed Singular Value Decomposition to reduce our attributes from 1500 to 300. We ran 6 different classification models (As depicted Below) and tested them against the ROC and accuracy metric. We choose Random Forest as the most suitable model because of its high accuracy value



### Level 3 Introduction

Our level 2 model was able to classify danish biscuits, Oreo biscuits, and maria biscuits as Biscuits but it was important for us to come up with a way to make further classification of those products. To do that we decide to take an unsupervised machine learning approach for Level 3. Also, Level 3 was the most important from the client perspective, as they will eventually use it to find comparable products. Those products within the same classification will be suitable for calculating the Purchasing Power Parities (PPP).

### Level 3 Algorithms

For level 3 the main goal here was to create a model where we'd simply be able to use product names to make our classifications. We first decided to implement clustering analysis to understand our data better. By implementing k means clustering, we would be able to classify products based on the names we gave our clusters. We explored outputs for many numbers of "k" clusters to see if the differences between each cluster would increase. In addition, we also needed to ensure that the cluster names made sense for our classification models. We discovered that the optimal number of clusters was 5. After running cluster analysis, based on our cluster word frequency we named cluster 1 as chocolate, cluster 2 as chocolate chip, cluster 3 as creamme, cluster 4 as peanut butter, and cluster 5 as sugar(Figure 7).

Once we were able to make 5 easily understandable clusters. We decided to now perform supervised machine learning. To perform text analytics, we had to make tokens of our product name and perform the following 4 major steps.

**Step 1 - Data Cleaning:** Since we were already using the same data set, we did not have to go over the same data cleaning process again

**Step 2 - Identifying the most frequent words:** Our next goal was to identify the most frequently used words within each cluster and create a document frequency matrix using them.

- Cluster 1 - Chocolate: Among the 100 most frequent words some of them were sandwich, dough, cream, shortbread (Figure 8)
- Cluster 2 - Chocolate Chip: Among the 110 most frequent words some of them were dough, cream, ice, milk (Figure 9)
- Cluster 3 - Creamme: Among the 90 most frequent words some of them were vanilla, fill, sandwich (Figure 10)
- Cluster 4 - Peanut Butter: Among the 80 most frequent words some of them were bar, Jelly (Figure 11)
- Cluster 5 - Sugar: Among the 120 most frequent words some of them were frost, fill, Bake (Figure 12)

**Step 3 - Creating n grams:** We created N grams with N=1 and N=2. This helped us to create a model that will make better predictions. After that we computed their prediction power with the help of TF-IDF.

**Step 4 - Singular Value Decomposition:** Once we calculated TF-IDF we performed Singular Value Decomposition to reduce our attributes from 1200 to 300 Once we had a reduced number of attributes, we ran 6 different machine learning algorithms and tested all of them against ROC and accuracy metric. We decided to go with Random Forest as it had an accuracy of 89% and roc of 88%. (Figure 13)

## USER INTERFACE OVERVIEW

To culminate what we have accomplished for this project, we created a user interface (UI) that allows for the World Bank to easily classify products for all 3 levels. Our user interface requires the user to upload a CSV file and it would display the dataset and the 3 levels of classification. In addition to this, we built our models to output probabilities that would also be displayed for the respective product in each level.

## FUTURE OVERVIEW

One of the biggest challenges that the world bank faced was taking language into account while making predictions. In future we believe 2 AWS services like AWS Rekognition and AWS Translate will help to solve this problem. After both the services will be incorporated this is what the new workflow will look like

1. By using the AWS Rekognition “Text in image” features, we can upload the product's packages images that we want to analyze on the Rekognition, and AWS Rekognition will identify the text on the image and display the detected texts on the image and transform them to the texts. Although they can detect the words in different languages, but it doesn't automatically translate the foreign languages into English, there are still some non-English words on the package.

2. With the AWS Translate, which is a text translation service that uses advanced machine learning technologies to provide high-quality translation on demands. So, after recognizing the texts on the product's package, we are going to use the AWS Real-Time translation service and copy paste the non-English words to the platform, and then with this smart service it will translate the detected language to English which we can eventually feed those English words into our model.

## CONCLUSION

Our goals for the project will be creating: Fast, Smart, and Automated process for PPP calculation. In this case, when the World Bank wants to find the products that are comparable across countries, they can simply drag and drop product labels into AWS Rekognition that will automatically recognize product descriptions and translate them into different languages if needed. Then using the User Interface, it will help to identify the following:

- Is the product a Bakery item or Non-Bakery item?
- If the product is a Bakery item, then is it a Biscuit or Non-Biscuit item?
- Making further classification of that product

Going forward, we hope that the User Interface designed in incorporation with AWS services will simplify the data collection process to meet the goals of the ICP program. This will be a big milestone for the World Bank and their ICP project to leap forward. We are very much looking forward to the World Bank to use our technology to better analyze the countries' GDPs and purchasing power in the future to reach their bigger goal.

## GLOSSARY

ROC	ROC curve, also known as Receiver Operating Characteristics Curve, is a metric used to measure the performance of a classifier model. IT depicts the rate of true positives with respect to the rate of false positives, therefore highlighting the sensitivity of the classifier model.
ACCURACY	Accuracy is one metric for evaluating classification models. Informally, accuracy is the fraction of predictions our model got right. Formally, accuracy has the following definition: $\text{Accuracy} = \frac{\text{Number of correct predictions}}{\text{Total number of predictions}}$ .
TF-IDF	In information retrieval Tf-idf, is short for term frequency-inverse document frequency, is a numerical statistic that is intended to reflect how important a word is to a document in a collection or corpus.
TOKENS	Tokens are the individual units of meaning you're operating on. This can be words, phonemes, or even full sentences. Tokenization is the process of breaking text documents apart into those pieces.
DOCUMENT FREQUENCY MATRIX	A document-term matrix is a mathematical matrix that describes the frequency of terms that occur in a collection of documents.
SINGULAR VALUE DECOMPOSITION	The goal of SVD is to find the optimal set of factors that best predict the outcome.
N GRAMS	They are basically a set of co-occurring words within a given window and when computing the n-grams you typically move one word forward (although you can move X words forward in more advanced scenarios)
BAGGING TREE	Bagging classification model uses a simple approach of improving the estimate of one by combining the estimates of many. Bagging constructs n classification trees using bootstrap sampling of the training data and then combines their prediction to produce a final metadata.
KNN	K nearest neighbors is a simple algorithm that stores all available cases and classifies new cases based on a similarity measure (e.g., distance functions).
LDA	Linear Discriminant Analysis or Normal Discriminant Analysis or Discriminant Function Analysis is a dimensionality reduction technique that is commonly used for supervised classification problems. It is used for modelling differences in groups i.e. separating two or more classes.
LRM	Logistic regression is a classification algorithm used to assign observations to a discrete set of classes. Logistic regression transforms its output using the logistic sigmoid function to return a probability value.
NEURAL NETWORK	A neural network is a series of algorithms that recognize underlying relationships in a set of data through a process that imitates the way the human brain operates. The artificial neural network (ANN) assimilates data in the same way the human brain processes information.
RANDOM FOREST	Random forest is a Supervised Machine Learning Algorithm that is used widely in Classification and Regression problems. It builds decision trees on different samples and takes their majority vote for classification and average in case of regression.





Figure 4.

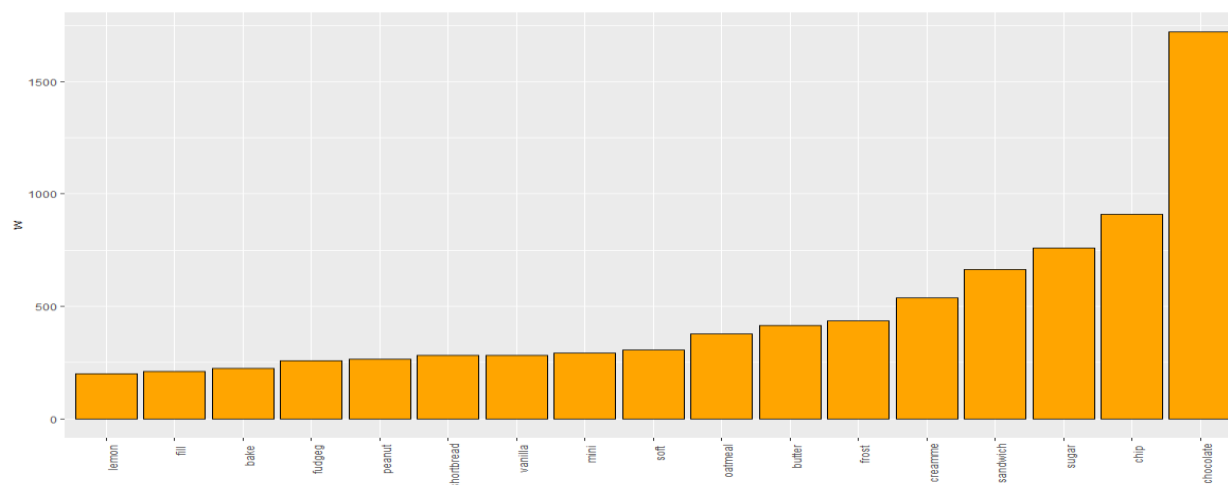


Figure 5.

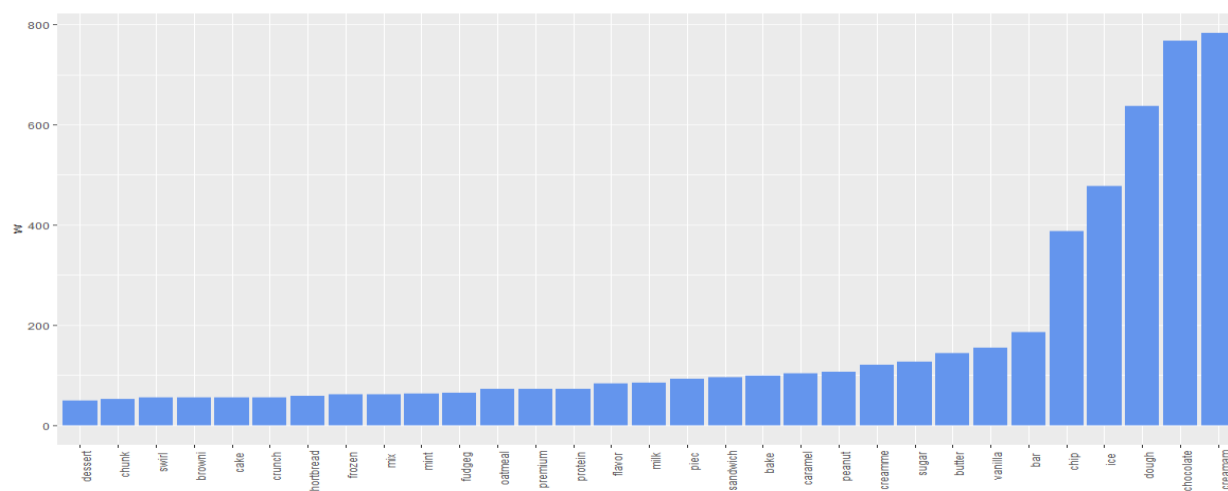


Figure 6.

```

[1] "cooki"
[4] "cream"
[7] "piec"
[10] "chocol"
[13] "crumbl"
[16] "candi"
[19] "ice"
[22] "cooki"
[25] "cooki_cream"
[28] "cream_cooki"
[31] "piec_wrap"
[34] "chocol_coat"
[37] "crumbl_top"
[40] "candi_ribbon"
[43] "ice_cream"
[46] "cooki_dough"

"cream"
"cooki"
"wrap"
"coat"
"top"
"ribbon"
"cream"
"dough"
"cream_ice"
"cooki_dough"
"wrap_dark"
"coat_cooki"
"top_white"
"ribbon_premium"
"cream_bar"
"dough_crunch"

"ice"
"dough"
"dark"
"cooki"
"white"
"premium"
"bar"
"crunch"
"ice_cream"
"dough_piec"
"dark_chocol"
"cooki_crumbl"
"white_candi"
"premium_ice"
"bar_cooki"

```



Figure 9.

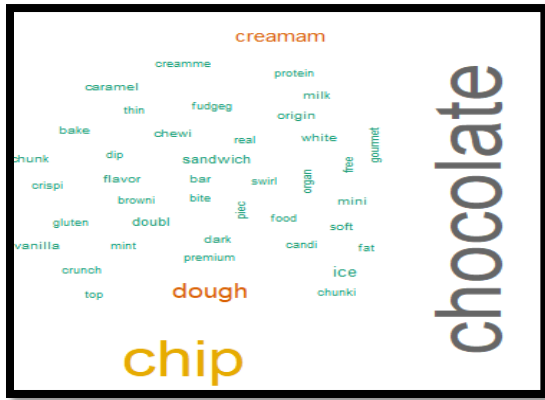


Figure 10.

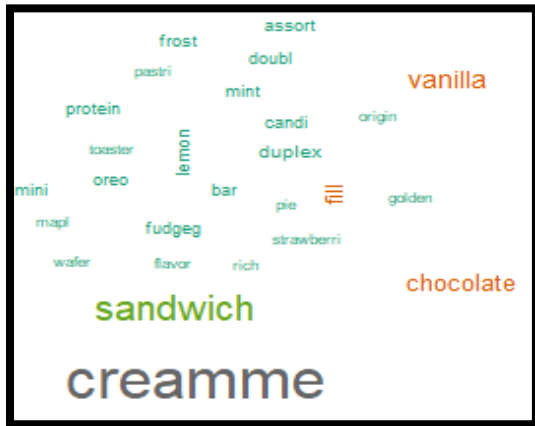


Figure 11.

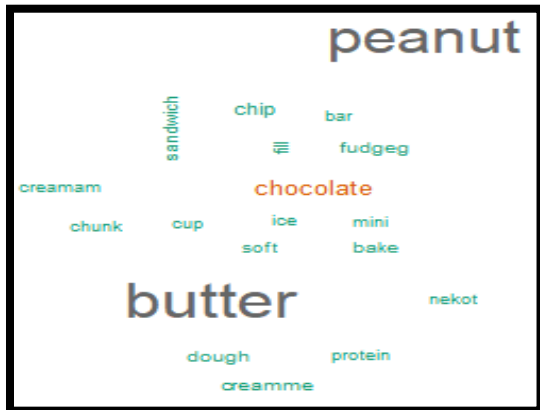


Figure 12.



Figure 13.

