

# **Data Chefs**

## **Cuisine Classification and Recipe Summarization Final Report**

*Group 4*

*Bhavya Melana, Diksha Gulati, Harshdeep Mac, Utkarsh Kekre*

*Georgia State University, J Mack Robinson College of Business*

## INDEX

● Introduction .....	3
● The Solution .....	3
● Technology .....	3
● Proposed Architecture .....	4
● Business Use Case.....	5
I.    Cuisine Classification .....	5
II.    Recipe Summarization .....	5
● Data Sources .....	6
I.    AllRecipe .....	6
II.    Epicurious .....	7
● Data Scraping .....	7
● Data Storage .....	8
● Azure Cosmos DB .....	10
● Data Cleaning .....	11
● Data Science Solution .....	12
I.    Recipe Summarization .....	13
II.    Cuisine Classification.....	20
● Visualization .....	30
● Nutritional Value .....	31
● Business Value .....	35
● References .....	36

## **Introduction**

Cooking is one of the most central exercises in one's life. It isn't just associated with the delight of eating yet in addition profoundly influences different parts of human life, for example, wellbeing, dietary, culinary craftsmanship, diversion, human correspondence, etc.

Henceforth, the number of websites uploading different kinds of recipes and their users are increasing rapidly, day by day. According to DMR, in the year 2020, Allrecipes have 60.3 million unique visitors and 370 million recipes saved on the website. Allrecipes is available in 24 countries in 12 different languages. Also, one of other famous website Epicurious have 9.50 million users with around 73.5% traffic from the United States. These statistics shows the interest of the users.

All our team members are food enthusiasts which inspired us to come up with an idea for such food websites, to make life of a user easy. User experience has become a preordained part of any company. Hence, saving user's time and providing them a smooth experience yields great value to retain the customers and building trust.

As a user food website user, we thought of few challenges we face, they are as follows:

- Not all the recipes uploaded by users have their corresponding cuisines labeled.
- Not all the recipes have a meaningful title.

When solution is built, users will have the following advantages:

- Cuisine classification based on the ingredients of the recipe.
- Assigning a meaningful title to the recipes with irrelevant titles, using the directions mentioned in the recipes on the websites.

## **The Solution**

We aim to enhance user experience and save their time by improvising their journey on the website in use. Hence, ensuring customer retention by providing a positive experience, reducing the burden, gaining credibility, and giving them productive screen time.

## **Technology**

To achieve our solution, we have jotted down few important concepts which we are going to use to make it possible:

- Stemming
- Term Document Matrix
- Classification Algorithm
- Sequence2Sequence Encoder-Decoder
- Natural Language Processing (NLP)

To implement the above-mentioned concepts, a huge amount of data will be trained on machine learning models and furthermore, tested to check the accuracy to provide a valuable result. Models to be implemented:

- Recurrent Neural Networks (RNN)
- Long Short-Term Memory (LSTMs)

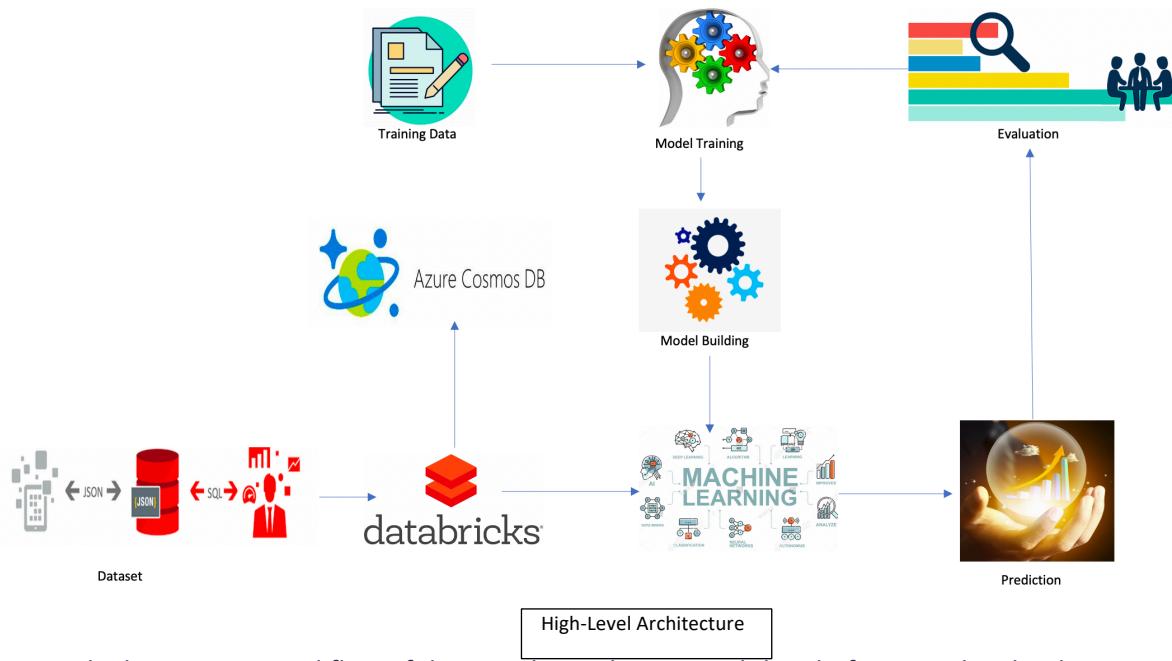
We are going to use Microsoft Azure platform, to carry out the above enumerated tasks, as we need an effective portal. Some of the requirements for the proposed solution are:

- Data storage
- Executing the machine learning models
- Azure Databricks

We shall discuss all the above-mentioned tasks further in detail as part of this report, along with results and the challenges that were faced.

## Proposed Architecture

The graphical representation below illustrates the flow of data in a process based on the inputs and outputs. All through our task there are different stages and procedures through which the data will stream and accordingly it is basic that there ought to be a defined flow of data to get the desired result.



Along with above-mentioned flow of data, we have also covered the platforms and technologies used at distinctive stages to process the data. This architecture defines the execution of tasks that is required for the application to run smoothly and provide desired results.

# Business Use Case

## Cuisine Classification:

With Cuisine Classification, we intend to provide a solution where-in all the recipes have the appropriate cuisine marked in the recipe, this will help in the smooth filtering process, based on cuisine type. This solution will target in improvising customer's journey on the application, and hence result in better user experience. As you can see below, there are 2 screenshots, the first screenshot (Fig. 1) displays the recipe without the tagged cuisine. Based on the ingredients stated in the recipe, we are going to recommend the cuisine type and designate it with the recipe, as shown in Fig. 2. The mentioned solution will also be a recommendation system used in web searches, and when user jot down the ingredients, it will recommend the cuisine type.

**Dal Makhani (Indian Lentils)**  
★★★★★ 130 made it | 55 reviews | 10 photos  
Recipe by: SOGOLONDJATA  
"Ever go to an Indian restaurant and wonder how they make their lentil dishes? I did, and I wanted to duplicate Indian food. Then I scoured the internet to figure out how they achieved them, and through mixing and matching different versions, I think I've arrived at this recipe, which I think is pretty close. This version is lighter than the ones I've had in restaurants, but still has the rich flavor. Kasuri methi (fenugreek leaves) is almost impossible to find in the U.S., even in NYC, but it gives this dish something very special!"  
Save | I Made It | Rate It | Print | Pin | Share

**Ingredients**  
4 h 15 m | 6 servings | 390 cal  
1 cup lentils  
1/4 cup dry kidney beans (optional)  
water to cover  
5 cups water  
salt to taste  
2 tablespoons vegetable oil  
1 tablespoon cumin seeds  
1 1/2 tablespoons ginger paste  
1 1/2 tablespoons garlic paste  
1/2 teaspoon ground turmeric  
1 pinch cayenne pepper, or more to taste  
1 cup canned tomato puree, or more to taste  
1 tablespoon chili powder  
2 tablespoons ground coriander

**Dal Makhani (Indian Lentils)**  
★★★★★ 130 made it | 55 reviews | 10 photos  
Recipe by: SOGOLONDJATA  
"Ever go to an Indian restaurant and wonder how they make their lentil dishes? I did, and I wanted to duplicate Indian food. Then I scoured the internet to figure out how they achieved them, and through mixing and matching different versions, I think I've arrived at this recipe, which I think is pretty close. This version is lighter than the ones I've had in restaurants, but still has the rich flavor. Kasuri methi (fenugreek leaves) is almost impossible to find in the U.S., even in NYC, but it gives this dish something very special!"  
Save | I Made It | Rate It | Print | Pin | Share

**Cuisine Type: Indian**  
4 h 15 m | 6 servings | 390 cal  
1 cup lentils  
1/4 cup dry kidney beans (optional)  
water to cover  
5 cups water  
salt to taste  
2 tablespoons vegetable oil  
1 tablespoon cumin seeds  
1 1/2 tablespoons ginger paste  
1 1/2 tablespoons garlic paste  
1/2 teaspoon ground turmeric  
1 pinch cayenne pepper, or more to taste  
1 cup canned tomato puree, or more to taste  
1 tablespoon chili powder  
2 tablespoons ground coriander

## Recipe Summarization:

Here, we aim to achieve a solution where-in all the recipes have meaningful titles, using the directions mentioned in the recipe. As you can see below, there are 2 screenshots, the first screenshot (Fig. 1) displays the recipe with an irrelevant title causing filtering and sorting of recipes difficult for the users, our objective is to provide a meaningful title to before-mentioned recipes to enhance the user experience, as presented in the second screenshot (Fig. 2).

**Mother's Day Pie**  
★★★★★ - 144 Ratings | 115 Reviews | 33 Photos  
This is a crustless coconut custard pie. The recipe was given to me by an elderly lady I knew. It's never failed me!  
By Jo Pelletier  
Save | Pin | Print | Share | Watch

**Prep:** 15 mins | **Cook:** 35 mins | **Total:** 2 hrs 50 mins | **Additional:** 2 hrs | **Servings:** 8 | **Yield:** 1 9-inch pie | **Nutrition Info**

**Coconut Custard Pie**  
★★★★★ - 144 Ratings | 115 Reviews | 33 Photos  
This is a crustless coconut custard pie. The recipe was given to me by an elderly lady I knew. It's never failed me!  
By Jo Pelletier  
Save | Pin | Print | Share | Watch

**Prep:** 15 mins | **Cook:** 35 mins | **Total:** 2 hrs 50 mins | **Additional:** 2 hrs | **Servings:** 8 | **Yield:** 1 9-inch pie | **Nutrition Info**

Fig. 1

Fig. 2

# Data Sources

A data source may be the initial location where data is born or where physical information is first digitized, however even the most refined data may serve as a source, as long as another process access utilizes it. We scraped the data from external data sources in html format which is the text file format. The extracted data contains Ingredients, directions and Nutritional Fact.

There are 2 data sources, basically food websites:

## 1. All Recipes

It is a food focused online social networking service. Recipes are categorized by season, type (such as appetizer or dessert), and ingredients. The website allows users to search recipes and include specification in their search (such as type of meal, nutrition, key ingredients, and time needed to prepare the dish).

The screenshot shows a recipe page for "Addictive Sweet Potato Burritos" on allrecipes.com. The page includes a title, a star rating of 4.5 stars, and a brief description. It features a large image of the dish and several smaller images of related dishes. Below the image, there are sections for "Directions", "Ingredients", and "Nutrition Facts". The "Directions" section provides step-by-step instructions with timing for prep, cook, and total time. The "Ingredients" section lists the required items with their quantities. The "Nutrition Facts" section provides nutritional information per serving. There are also links for "You might also like" and "Get the magazine".

## 2. Epicurious

It is a digital brand for consumers interested in food and cooking.

The screenshot shows a recipe page for "Swiss Chard Ribs and Stems" on epicurious.com. The page includes a title, a star rating of 4.5 stars, and a brief description. It features a large image of the dish and several smaller images of related dishes. Below the image, there are sections for "PREPARATION", "INGREDIENTS", and "ASSEMBLY". The "PREPARATION" section provides step-by-step instructions for making Béchamel. The "INGREDIENTS" section lists the required items with their quantities. The "ASSEMBLY" section provides instructions for combining the ingredients and cooking the dish. There are also links for "Do Ahead" and "Assembly".

## Data Scraping

Instead of taking a built-in data set to solve our business problems, we have come up with a solution to build our own dataset by scrapping the food websites. Two scrapper scripts have been written in python using “BeautifulSoup” package to scrape the websites “Epicurious” and “AllRecipes” to fetch ingredients, instructions, and nutrition from the webpage and automatically stores them in azure cosmos database using Azure SQL API. The scrapper codes are shown below.

### Epicurious Scrapper Code

```
17  for cuisine in all_cuisines:
18      total_pages = int(BeautifulSoup(requests.get(all_cuisines[cuisine]).text, 'lxml').find('nav', class_="common-pagination")['data-total-pages'])
19      for page in range(total_pages):
20          current_page = BeautifulSoup(requests.get(f'https://www.epicurious.com/search/?cuisine={cuisine}&content=recipe&page={page+1}').text, 'lxml')
21          recipe_index = 0
22          while recipe_index <= 17:
23              link = current_page.find('article', {"data-index":str(recipe_index)})
24              if link != None:
25                  url = "https://www.epicurious.com" + current_page.find('article', {"data-index":str(recipe_index)}).find('a')['href']
26                  try:
27                      current_recipe = BeautifulSoup(requests.get(url).text, 'lxml')
28                      recipe_title = current_recipe.find('h1', {"itemprop":"name"}).text
29
30                      # Fetch the ingredients from the recipe
31                      ingredients = []
32                      for ingredient in current_recipe.find_all('li', {"class":"ingredient", "itemprop":"ingredients"}):
33                          ingredients.append(ingredient.text)
34
35                      # Clean and Parse the ingredients
36                      #ingredients = Ing_Cleaner(ingredients)
37
38                      # Putting the ingredients in ingredient corpus
39                      for i in ingredients:
40                          ingredient_corpus.add(i)
41
42                      ingredients_data_set = {"title":recipe_title, "cuisine": cuisine, "ingredients": ingredients}
43                      ingredients_data_set = json.dumps(ingredients_data_set)
44
45                      # Fetch the instructions from the recipe
46                      instructions = []
47                      for inst in current_recipe.find_all('li', class_="preparation-step"):
48                          instructions.append(inst.text)
49                      instructions_data_set = {"title":recipe_title, "cuisine": cuisine, "instructions": instructions}
50                      instructions_data_set = json.dumps(instructions_data_set)
51
52                      #Fetch the nutritions from the recipe
53                      nutritions = {"title":recipe_title, "cuisine": cuisine}
54                      nutri_label = current_recipe.find_all('span', class_="nutri-label")
55                      nutri_data = current_recipe.find_all('span', class_="nutri-data")
56                      length = len(nutri_label)
57                      for l in range(length):
58                          nutritions[nutri_label[l].text] = nutri_data[l].text
59                      nutritions_data_set = json.dumps(nutritions)
60
61                      upsert_document(client, database_name, "Ingredients", ingredients_data_set)
62                      upsert_document(client, database_name, "Instructions", instructions_data_set)
63                      upsert_document(client, database_name, "Nutritions", nutritions_data_set)
64                  except:
65                      pass
66                  recipe_index += 1
67                  dummy += 1
68              else:
69                  ....
```

## AllRecipes Scraper Code:

```
6
7     for cuisine in world_cuisines.find_all('a', class_="grid-col--subnav"): # This loop is to iterate through cuisines
8         cuisine_url = cuisine['href'] # url of the current cuisine
9         cuisine_page = BeautifulSoup(requests.get(cuisine_url).text, "lxml")
10        cuisine_name = str.lower(cuisine.find('span').text.replace(" recipes", ""))
11        all_cuisines[cuisine_name] = cuisine_url # Use this dictionary to create the labels dictionary
12        for sub_cuisine in cuisine_page.find_all('a', class_="grid-col--subnav"): # This loop is to iterate thorough sub cuisines
13            sub_cuisine_url = sub_cuisine['href'] # url of the sub-cuisine
14            sub_cuisine_page = BeautifulSoup(requests.get(sub_cuisine_url).text, "lxml")
15            page_number = 0
16            while page_number <= 10000: # This loop is to iterate through pages
17                current_page_url = sub_cuisine_url + f'?internalSource=hub%20nav&referringId=728&referringContentType=Recipe%20Hub&linkName=hub%20nav%20daughter&clickId=hub%20nav%202&page={(page_number+1)}'
18                page_number += 1
19                current_page = BeautifulSoup(requests.get(current_page_url).text, "lxml")
20                if current_page.find("title").text == "Allrecipes - File Not Found":
21                    break
22                for recipe in current_page.find_all("article", class_="fixed-recipe-card"): # This loop is to iterate through all recipies in a page
23                    try:
24                        recipe_url = recipe.find("a", class_="fixed-recipe-card__title-link")['href']
25                        recipe_page = BeautifulSoup(requests.get(recipe_url).text, "lxml")
26                        recipe_name = recipe_page.find("h1").text
27
28                        ingredients = []
29                        for ing in recipe_page.find_all("span", {"class": ["recipe-ingred_txt added", "ingredients-item-name"]}):
30                            ingredients.append(ing.text)
31                        ingredient_data_set = {"source": "AllRecipes", "cuisine": cuisine_name, "title": recipe_name, "ingredients": ingredients}
32                        ingredient_data_set = json.dumps(ingredient_data_set)
33
34                        instructions = []
35                        for inst in recipe_page.find_all('li', {"class": ["step", "subcontainer instructions-section-item"]}):
36                            instructions.append(inst.text)
37
38                        instruction_data_set = {"source": "AllRecipes", "cuisine": cuisine_name, "title": recipe_name, "instructions": instructions}
39                        instruction_data_set = json.dumps(instruction_data_set)
40
41                        upsert_document(client, database_name, "Ingredients", ingredient_data_set)
42                        upsert_document(client, database_name, "Instructions", instruction_data_set)
43                    except:
44                        pass
45                    print(cuisine, " is complete")
```

## Data Storage

The meaningful information from the data needs to be derived for the data to be useful. Unstructured data such as text is major input in this project and since the volume is high, so we stored the data in Azure Cosmos DB.

We have primarily 3 types of data i.e. ingredients, directions and nutritional values of all recipes which will be stored in their respective containers. As we have scrapped the data from 2 sources, the scrapping script for the 2 sources is attached below. The raw form of data will be stored in containers named as “Ingredients”, “Instructions” and “Nutritons”. The data in these containers are stored in form of documents and have a JSON schema shown below.

## Ingredients JSON Schema:

```
1  {
2      "title": "Vietnamese Rice Noodle Salad",
3      "cuisine": "vietnamese",
4      "ingredients": [
5          "4 oz thin rice noodles",
6          "1/4 cup rice vinegar (not seasoned)",
7          "1 tablespoon sugar",
8          "1 tablespoon Asian fish sauce",
9          "1/4 teaspoon salt",
10         "1 carrot, coarsely shredded",
11         "2 scallions, thinly sliced crosswise",
12         "1 cup loosely packed mixed fresh cilantro, mint, and/or basil leaves, torn if large",
13         "1/4 cup chopped unsalted dry-roasted peanuts"
14     ],
15     "id": "40084444-09e7-415a-b88d-af968139dad2",
16     "_rid": "12w6AMXcXfndSwAAAAAAA==",
17     "_self": "dbs/12w6AA==/colls/12w6AMXcXfk=/docs/12w6AMXcXfndSwAAAAAAA==/",
18     "_etag": "\ed004475-0000-0200-0000-5e9e98880000"",
19     "_attachments": "attachments/",
20     "_ts": 1587452040
21 }
```

## Instructions JSON Schema:

```
1  {
2      "source": "AllRecipes",
3      "cuisine": "mexican",
4      "title": "Tofu Tacos I",
5      "instructions": [
6          "\nOver a medium heat fry the tofu, oil, garlic and onion in a large skillet for 5 minutes. Add the chili powder, paprika, cayenne, i
7          "\nSpoon the tofu mixture into taco shells. Top the mixture with lettuce, tomatoes, avocado, Cheddar cheese and salsa.\n
8          "\n\n"
9      ],
10     "id": "3ded8057-9ae0-4471-b848-50350b8d5cc2",
11     "_rid": "12w6AMkGx8xKAAAAAAA==",
12     "_self": "dbs/12w6AA==/colls/12w6AMkGx8w=/docs/12w6AMkGx8xKAAAAAAA==/",
13     "_etag": "\e600ce60-0000-0200-0000-5e9e61080000",
14     "_attachments": "attachments/",
15     "_ts": 1587437832
16 }
```

## Nutrition JSON Schema:

```
1  {
2      "title": "Efo Riro With Kale and Whitefish",
3      "cuisine": "african",
4      "Calories": "447",
5      "Carbohydrates": "11 g(4%)",
6      "Fat": "37 g(57%)",
7      "Protein": "20 g(40%)",
8      "Saturated Fat": "3 g(15%)",
9      "Sodium": "129 mg(5%)",
10     "Polyunsaturated Fat": "6 g",
11     "Fiber": "3 g(14%)",
12     "Monounsaturated Fat": "25 g",
13     "Cholesterol": "49 mg(16%)",
14     "id": "dd4d39a6-2db0-4d0c-9188-59a546a8674b",
15     "_rid": "12w6ANdpSMIBAAAAAAA==",
16     "_self": "dbs/12w6AA==/colls/12w6ANdpSMI=/docs/12w6ANdpSMIBAAAAAAA==/",
17     "_etag": "\e6008c68-0000-0200-0000-5e9e613e0000",
18     "_attachments": "attachments/",
19     "_ts": 1587437886
20 }
```

## Azure Cosmos DB

Cosmos Db is a database service that is globally distributed. It permits you to deal with your information regardless of whether you keep them in data centers that are dissipated all through the world. It gives the tools you have to scale both worldwide distribution design and computational assets, and these devices are given by Microsoft Azure. The database is used for document, key value, relational, and graph models. We have stored our data in the document format. A document database is a kind of nonrelational database that is intended to store and query data as JSON-like documents. Document databases make it simpler for engineers to store and query data in a database by utilizing a similar document model configuration they use in their application code.

Using Cosmos DB has following advantages:

### 1. Global Distribution

It is available in every location where Azure is available by default. So, you can setup instances of your Cosmos DB at any location that you want simply by activating the desired location from the Azure portal. This will ensure that the data is replicated and available for the users in the region with guaranteed low latency. Additionally, it also provides automatic and manual failover that enables high availability and disaster recovery.

### 2. Performance

It ensures accessing data with faster response time. It also provides guaranteed throughput based on the provisioned output capacity.

### 3. Pricing

It is dependent on the required throughput and the storage necessary for your data.

### 4. Multi-Model and Multi APIs

It enables to create containers that can store data in Key-Value, Columnar, Document or Graph data stores. Along with this, it also provides users the flexibility to choose different APIs to access the data such as:

- SQL API or MongoDB API (for Document databases)
- Table API (Key-Value databases)
- Cassandra API (Columnar databases)
- Gremlin API (Graph databases)

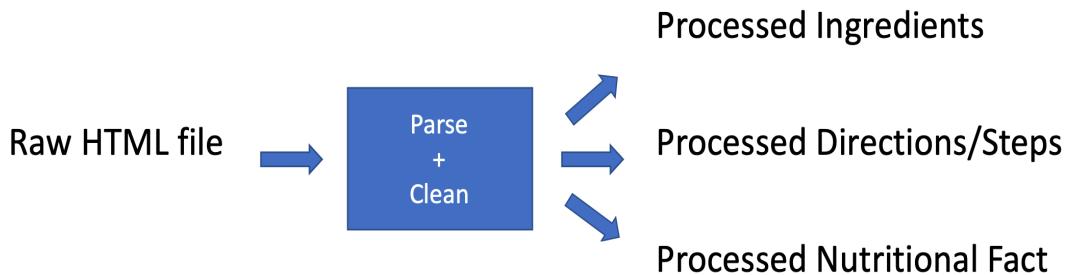
### 5. Defined Consistency Levels

There are 5 consistency levels that strikes a balance between latency, throughput, and availability.

- **Strong Consistency:** It ensures consistency across all nodes, in all regions, but this comes at the cost of overall performance.
- **Bounded Staleness Consistency:** It provides a mean to set the level of freshness of data.
- **Session Consistency:** It ensures that there is no dirty read for the writer but it's possible to have dirty reads for other users. This is the default consistency level.
- **Consistent Prefix:** It ensures that the read data has been updated to all replicas.
- **Eventual Consistency:** It provides no guarantees on the freshness of the data or on the order. However, this provides the fastest performance.

## Data Cleaning

Process of preparing data for analysis by removing or modifying data that is incorrect, incomplete, irrelevant, duplicated, or improperly formatted.



After we store our data in the respective containers, we need to clean the data so that it can be used in modeling out the proposed solution. For this we need to clean our ingredients to get only the root ingredient out of the whole line i.e.

**"2 cups drained canned diced tomatoes (from 28-ounce can)" → "tomato".**

To obtain this, we have written a script which cleans the ingredients as shown above, cleans the directions to remove any special characters and same for nutritional values. The script includes functions like Stemmer, stop word removal, Part of speech tagger, quantity removal, measurements removal, color removal, additional words removal and regex for special chars and number removal. Below is an example of our cleaning.



Now as the schema for both the sources will be uniform that is the output of the cleaning script will have the same schema for both the sources, the data sources will be automatically stitched together into the respective containers.

## Data Science Solution

Different data science solution has been chosen to solve business problems.

### Problem 1: Cuisine Classification

**Solution:** The complete step by step solution to solve the problem of cuisine classification is shown below:

1. **Ingredient Corpus:** A dictionary to save all the ingredients to be later used in creating a final matrix for the classification model.
2. **Label Dictionary:** A dictionary to save all the cuisines to be later used in creating labels for the data matrix.
3. **Term Document Matrix:** Written a python code to convert ingredients into a term document matrix (a binary matrix) to be later used in training and testing of classification model.
4. **Classification Algorithm:** Plan is to implement an auto-ML pipeline in azure to automatically train all the classification models mainly neural-network, support vector machine, logistic regression, random forest, decision tree etc. Auto-ML can pick the best model based on evaluation parameter.
5. **Evaluation:** Confusion matrix works well for most of the classification problems

#### Sample Outputs:

##### Example 1:

- Predicted Cuisine: African or Not African (It could be any other cuisine)
- Original Cuisine: African
- Ingredients: onion, plum tomato, bell pepper, clove, butter, tilapia, bay season, oil, tomato, thyme, garam masala, turmeric, Maggi shrimp, shrimp, kale, yam

## Problem 2: Recipe Summarization

**Solution:** The problem that we have in our hand is an abstractive summarization problem. We are using sequence2sequence encoder-decoder which is an advanced deep learning model. Model is a combination of multiple LSTMs. The complete model is divided into encoder, decoder and attention mechanism.

- **Encoder:** Bi-directional LSTM layer that extracts information from the original text. The bidirectional LSTM reads one word at a time and since it is a LSTM, it updates its hidden state based on the current word and the words it has read before.
- **Decoder:** Uni-directional LSTM layer that generates summaries one word at a time. The decoder LSTM starts working once it gets the signal than the full source text has been read. It uses information from the encoder as well as what is has written before to create the probability distribution over the next word.
- **Attention Mechanism:** Through attention mechanism, the decoder can access the intermediate hidden states in the encoder and use all that information to decide which word is next.

## Sample Outputs:

### Example 1:

- **Generated:** Chicken cake
- **Original:** Chicken French – Rochester, NY Style
- **Recipe:** all-purpose flour; salt; eggs; white sugar; grated parmesan cheese; olive oil; skinless; butter ; minced garlic ; dry sherry ; lemon juice ; low sodium chicken base ; ;Mix together the flour , salt , and pepper in a shallow bowl . In another bowl, whisk beaten eggs, sugar, and Parmesan cheese until the mixture is thoroughly blended and the sugar has dissolved. Heat olive oil in a large skillet over medium heat until the oil shimmers. Dip the chicken breasts into the flour mixture, then into the egg mixture, and gently lay them into the skillet. Pan-fry the chicken breasts until golden brown and no longer pink in the middle, about 6 minutes on each side. Remove from the skillet and set aside. In the same skillet over medium-low heat, melt the butter, and stir in garlic, sherry, lemon juice, and chicken base.

## Recipe Summarization

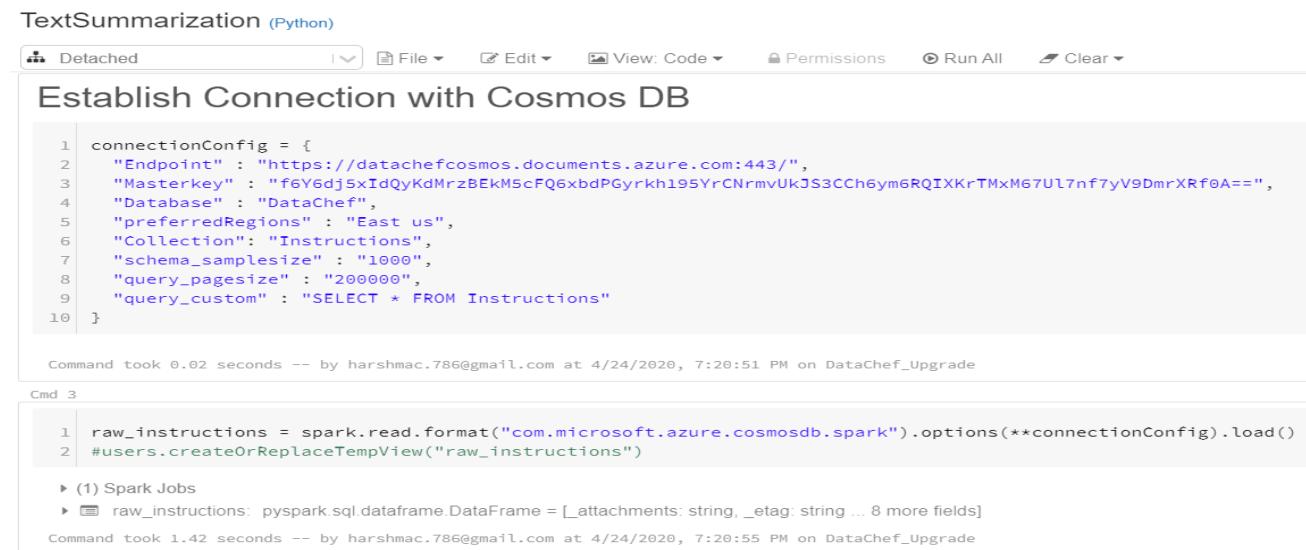
Recipe summarization refers to generating recipe titles from summarizing the recipe instructions and ingredients which relates more to the recipe. This concept is known as text summarization which uses deep learning capabilities to train a recurrent neural network model using multiple LSTM layers with combination of dense and dropout layers, to learn the language context and tries to summarize the recipe.

### 1. Data Pre-Processing

As we have the cleaned ingredients and instructions in their respective containers stored in Cosmos DB, we need to establish a DB connection with Azure Databricks where all the preprocessing will take place and the model will be trained.

To do this we use the Azure-CosmosDB-Spark connector provided by Azure to do the same. Installing this on the cluster enables us to create a connection link between Azure Databricks and the container where data is stored in CosmosDB.

In the connection configuration, we need to specify the DB endpoint, masterkey, Database name, DB region where it is hosted, container name and the query to fetch the documents. Then this connection config is passed as an argument to the **spark.read.format** function to create a Spark dataframe of the fetched documents from cosmos DB. Repeating this same process for both the containers that is cleaned ingredients and instructions. To play with the data in python, we created a Pandas dataframe from spark dataframe.



```
TextSummarization (Python)
Detached File Edit View: Code Permissions Run All Clear

Establish Connection with Cosmos DB

1 connectionConfig = {
2     "Endpoint" : "https://datachefcosmos.documents.azure.com:443/",
3     "Masterkey" : "f6Y6dj5xIdQyKdMrzBEK5cFQ6xbdPGyrkh195YrCNrmvUkJS3CCh6ym6RQIXKrTMxM67Ul7nf7yV9DmrXRf0A==",
4     "Database" : "DataChef",
5     "preferredRegions" : "East us",
6     "Collection": "Instructions",
7     "schema_samplesize" : "1000",
8     "query_pagesize" : "200000",
9     "query_custom" : "SELECT * FROM Instructions"
10 }

Command took 0.02 seconds -- by harshmac.786@gmail.com at 4/24/2020, 7:20:51 PM on DataChef_Upgrade

Cmd 3

1 raw_instructions = spark.read.format("com.microsoft.azure.cosmosdb.spark").options(**connectionConfig).load()
2 #users.createOrReplaceTempView("raw_instructions")

▶ (1) Spark Jobs
▶ raw_instructions: pyspark.sql.dataframe.DataFrame = [attachments: string, _etag: string ... 8 more fields]

Command took 1.42 seconds -- by harshmac.786@gmail.com at 4/24/2020, 7:20:55 PM on DataChef_Upgrade
```

Now we clean the instruction which are currently a list, making it a single string by concatenating the list and removing punctuation and other advertisement and steps which are irrelevant to our processing. Same cleaning is done for recipe titles.

We need to join the two Pandas dataframes based on title as it is the common key between 2 DF's. So, left joining the 2 DF's, we obtain the final DF. We create a custom column which has all ingredients first

and then the instructions following that. This is needed to make sure our Vocab contains all possible words which should be present in our Superset of words which will be explained in coming steps of this paper. Dropping nulls and showing the results of our final DF.

```

TextSummarization (Python)
Detached File View: Code Permissions Run All Clear
Schedule Comments R
Cmd 16
Final DF to model
1 | final_df
Out[344]:
cuisine_x Clean_inst Clean_titles joined_ing concated_data
0 mexican combine the lime juice and onion slice in a bo... fresh sangrita lime paper onion orang style valentina salt lime paper onion orang style valentina salt co...
1 mexican combine the tomatillos serrano peppers garlic ... came en su jugo meat in its juices tomatillo chile pepper garlic water bacon flan... tomatillo chile pepper garlic water bacon flan...
2 mexican combine the tomatillos serrano peppers garlic ... came en su jugo meat in its juices bacon tomatillo chile pepper garlic water flan... bacon tomatillo chile pepper garlic water flan...
3 mexican in a saucepan combine chicken breast with chic... authentic enchiladas verdes bone chicken breast chicken broth onion garlic... bone chicken breast chicken broth onion garlic...

```

Now selecting only headlines that is “Clean\_titles” and mixture of ingredients and instruction which is “concated\_data” which will be part of our vocab for our model.

Next, we create the vocab and vocab\_count for our model which is all the words in our corpus and its frequency which is vocab\_count as a dictionary. Now we create a custom method which takes the vocab as an argument and returns 2 dictionary one being indexes with respect to words and other being vice-versa i.e. words with respect to indexes.

Next, we initialize the **GloVe** dictionary.

GloVe is an unsupervised learning algorithm for obtaining vector representations for words. Training is performed on aggregated global word-word co-occurrence statistics from a corpus, and the resulting representations showcase interesting linear substructures of the word vector space.

There are various pre-trained word vectors in GloVe mainly:

- Wikipedia 2014 + Gigaword 5 (6B tokens, 400K vocab, uncased, 50d, 100d, 200d, & 300d vectors).
- Common Crawl (42B tokens, 1.9M vocab, uncased, 300d vectors)
- Common Crawl (840B tokens, 2.2M vocab, cased, 300d vectors).
- Twitter (2B tweets, 27B tokens, 1.2M vocab, uncased, 25d, 50d, 100d, & 200d vectors).

In this project, we have used the Wikipedia + Gigaword format of GloVe dictionary with 100 dimensional vectors. This dictionary is a TXT file of 330MB which needs to be installed on our cluster for processing.

To install this dictionary, we use the data upload feature of Databricks which uploads the file in FileStore location under tables container. Our code can directly reference this location to read this text file and can be processed.

Next, we read this dictionary in our cluster and extract all words present in it with indexes and the same indexes with the weights of the respective words in another variable in the “get\_glove” method.

We need to also create embeddings for our model which needs to be added, so we create in the next step.

First initialize our vocab\_size to our total words in vocab i.e. 21374 and our GloVe dimension as embedding dimension to 100. We create a matrix of the same shape as our vocab size and embedding dimension i.e. 21374 X 100 and assign random values in the matrix values.

Next, we copy the GloVe weights of each word found in the corpus to the embedding. With this we see that around 14K has been found in GloVe out of our corpus of 21K which is nearly 60%.

Now, for the remaining words, we try to map the GloVe vectors by finding the similar words using Cosine-Similarity to map these into our embedding. For ex:

#### TextSummarization (Python)

```

57
58 # map vocab to GloVe using cosine similarity
59 glove_idx2idx = build_word_to_glove(embedding, word2idx, idx2word, glove_index_dict, glove_embedding_weights)

# of GloVe substitutes found: 43
0.5911 dressings => sauces
0.5679 lek => soi
0.5530 montrachet => nobu
0.5488 paula => lisa
0.5464 concentration => resistance
0.5418 allergic => bacterial
0.5285 pinks => oranges
0.5230 supremes => blondie
0.5129 communal => communities
0.5116 mise => tipo

Command took 0.17 seconds -- by harshmac.786@gmail.com at 4/25/2020, 1:17:43 AM on DataChef_Upgrade

```

Next, we create dense vectors for our X and Y that is X for descriptions and Y for titles. To do this, we just take the indexes of words in titles and description and map it with its index for its place.

To split our data into train and test, we import the **Sklearn package's test\_train\_split** method to create 2 DF's of 80/20%.

Viewing our data randomly:

```

TextSummarization (Python)
MachineLearning1

1 def prt(label, word_idx, idx2word):
2     """Map 'word_idx' list to words and print it with its associated 'label'. """
3     words = [idx2word[word] for word in word_idx]
4     print('{0}: {1}\n{2}'.format(label, ' '.join(words)))
5
6 print('Random head, description:')
7 i = 0
8 prt('H', Y_train[i], idx2word)
9 prt('D', X_train[i], idx2word)
10

Random head, description:
H: tom kha gai chicken coconut soup

D: stalk lemongrass ginger kaffir lime sodium chicken broth shiitak coconut milk fish sugar chili oil cilantro lime wedg using the back of a knife lightly smash lemongrass and gin ger cut lemongrass into 4" pieces bring lemongrass ginger lime leaves and broth to a boil in a large saucenpan reduce heat and simmer until flavors are melded 8 10 minutes strain b roth into clean saucenpan discard solids add chicken and return to a boil reduce heat add mushrooms and simmer skimming occasionally until chicken is cooked through and mushrooms a re soft 20 25 minutes mix in coconut milk fish sauce and sugar divide soup among bowls serve with chili oil cilantro and lime wedges

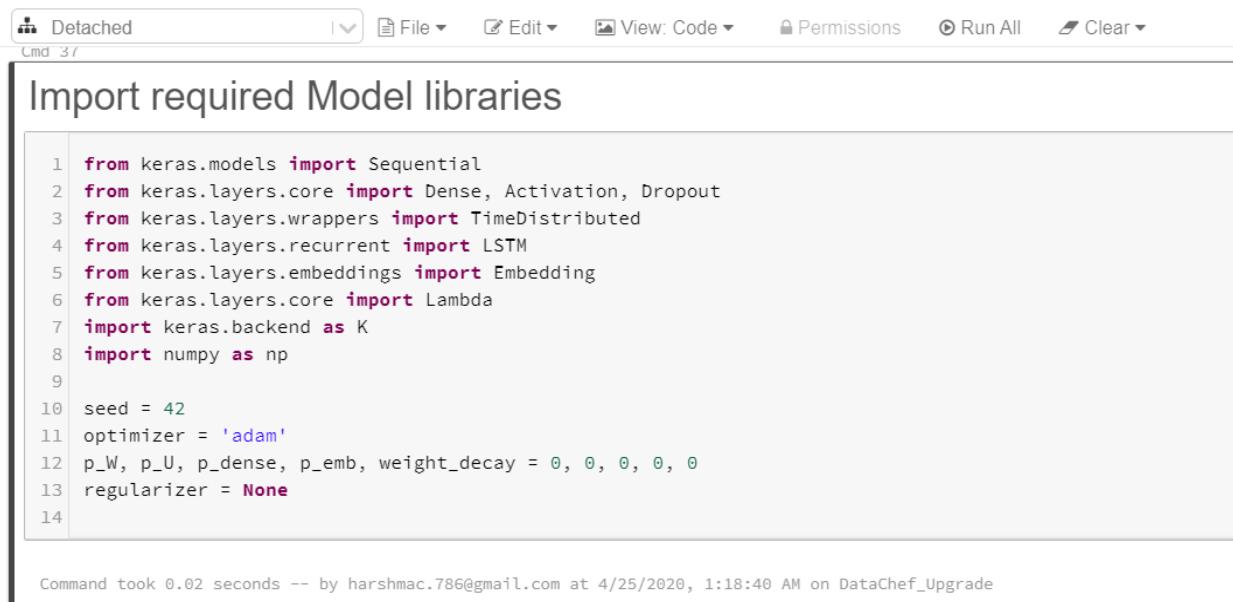
Command took 0.02 seconds -- by harshmac.786@gmail.com at 4/25/2020, 1:18:21 AM on DataChef_Upgrade

```

## 2. Modelling our RNN with LSTM's

Now, importing all the required libraries to create our RNN model with multiple LSTM layers. We will be using activation functions and optimizers in our model to train a good performing model.

TextSummarization ([Python](#))



The screenshot shows a Jupyter Notebook interface with the following details:

- Toolbar: Detached, File ▾, Edit ▾, View: Code ▾, Permissions, Run All, Clear ▾.
- Cell Type: Cmd 3/7
- Section Title: Import required Model libraries
- Code Content:

```
1 from keras.models import Sequential
2 from keras.layers.core import Dense, Activation, Dropout
3 from keras.layers.wrappers import TimeDistributed
4 from keras.layers.recurrent import LSTM
5 from keras.layers.embeddings import Embedding
6 from keras.layers.core import Lambda
7 import keras.backend as K
8 import numpy as np
9
10 seed = 42
11 optimizer = 'adam'
12 p_W, p_U, p_dense, p_emb, weight_decay = 0, 0, 0, 0, 0
13 regularizer = None
14
```
- Message Bar: Command took 0.02 seconds -- by harshmac.786@gmail.com at 4/25/2020, 1:18:40 AM on DataChef\_Upgrade

Now, we create our model which will consider the embedding, vocab, corpus, X and Y and try to predict a summary for our recipe aka instructions. To do this, we initialize a Sequential model which will enable us to add multiple layers into our model same way as we append variables into a list.

Firstly, we add our embedding layer into our model. The embedding layer is used to create word vectors for incoming words. It sits between the input and the LSTM layer, i.e. the output of the Embedding layer is the input to the LSTM layer. Our embedding should be of same dimension as the vocab size and embedding size.

Secondly, we add our LSTM layer with dropout layer into our model thrice. For our 1 layer of LSTM, we have 512 neurons mapped into it. Following it we add a dropout layer to prevent the LSTM to overfit during training.

Next, we add a Time distributed dense layer, which is actually an important component in our LSTM because our model tries to predict multiple words for our summarization, so as LSTM generates more than 1 output for our final output, our RNN model needs to wait until the LSTM has processed all the possible words and predicted its output for the next layer. This is where Time distributed layer comes into the picture as it waits for the LSTM to complete its processing to feed the consolidated output of LSTM into the next layer that is the dense layer. The dense layer is made up of neurons equal to the number of vocab words that are present in our corpus which is around 21K so that it can activate any word which has a higher probability for output.

Finally, we add the activation function in our model which is "SoftMax" to predict the output words.

And to compare the correctness of the predicted words, we add a compiler layer with loss function as “categorical\_crossentropy” and an optimizer as “Adam”.

```

1 def create_model(vocab_size, embedding_size, LR, rnn_layers, rnn_size, embedding=None):
2     """Construct and compile LSTM model."""
3     # Create a standard stacked LSTM
4     if embedding is not None:
5         embedding = [embedding]
6     model = Sequential()
7     model.add(Embedding(vocab_size, embedding_size,
8                         input_length=maxlen,
9                         W_regularizer=regularizer, dropout=p_emb, weights=embedding, mask_zero=True,
10                        name='embedding_1'))
11    for i in range(rnn_layers):
12        lstm = LSTM(rnn_size, return_sequences=True,
13                    W_regularizer=regularizer, U_regularizer=regularizer,
14                    b_regularizer=regularizer, dropout_W=p_W, dropout_U=p_U,
15                    name='lstm_{0}'.format(i + 1))
16        model.add(lstm)
17        model.add(Dropout(p_dense, name='dropout_{0}'.format(i + 1)))
18
19    if activation_rnn_size:
20        model.add(SimpleContext(simple_context, rnn_size, name='simplecontext_1'))
21
22    model.add(TimeDistributed(Dense(
23        vocab_size,
24        W_regularizer=regularizer,
25        b_regularizer=regularizer,
26        name='timedistributed_1')))
27    model.add(Activation('softmax', name='activation_1'))
28
29    # opt = Adam(lr=LR) # keep calm and reduce learning rate
30    model.compile(loss='categorical_crossentropy', optimizer=optimizer)
31

```

Our model has been initialized; we need to create our model parameters which include the learning rate which will enable our training speed for our model. Here we also initialize the RNN layers to be 3 and RNN size to be 512 for each layer.

Inspect our model to be accurate before training it for our dataset.

```

1 def inspect_model(model):
2     """Print the structure of Keras `model`."""
3     for i, l in enumerate(model.layers):
4         print(i, 'cls={} name={}'.format(type(l).__name__, l.name))
5         weights = l.get_weights()
6         print_str = ''
7         for weight in weights:
8             print_str += str(weight) + ' '
9         print(weights)
10        print()
11
12 inspect_model(model)
13
14 cls=Embedding name=embedding_1
15 [array([[-0.02869378, -0.03123973, 0.04808905, ..., -0.03416672,
16       0.04811154, -0.05177343], ...),
17  [-0.04584654, 0.04310323, 0.03905785, ..., 0.01259857,
18   0.03101949, 0.0602401], ...),
19  [-0.0071953, 0.023127, 0.0023731, ..., -0.071895,
20   0.066894, 0.019539], ...),
21  ...,
22  [-0.12421, 0.081948, 0.041412, ..., -0.021196,
23   0.022603, 0.013364], ...),
24  [0.017221, -0.032282, 0.08381, ..., -0.12928,
25   -0.024153, 0.018702], ...),
26  [-0.056554, 0.0020514, 0.017659, ..., -0.0088773,
27   0.021854, 0.002858]], dtype=float32]
28
29 cls=LSTM name=lstm_1
30 [array([[ -0.04708928, 0.04788486, -0.01991399, ..., 0.00173479,
```

Finally, we create a generator matrix for our training dataset and our validation dataset with all the various variables we have created so far like, vocab\_size, indexes to words, model, GloVe indexes, our batch size of 32.

Lastly, we train our model on training and validation dataset with 4 epochs with 640 steps in each epoch with 32 samples per step.

```

TextSummarization (Python)
Detached File View Code Permissions Run All Clear
Command took 0.08 seconds -- by harshmac.786@gmail.com at 4/25/2020, 1:20:19 AM on DataChef_Upgrade
Cmd 47
Training Model and Fitting
1 h = model.fit_generator(
2     traingen, samples_per_epoch=nb_train_samples,
3     nb_epoch=1, validation_data=valgen, nb_val_samples=nb_val_samples
4 )

/local_disk0/tmp/1587737721093-0/PythonShell.py:3: UserWarning: The semantics of the Keras 2 argument `steps_per_epoch` is not the same as the Keras 1 argument `samples_per_epoch`.  
`steps_per_epoch` is the number of batches to draw from the generator at each epoch. Basically `steps_per_epoch = samples_per_epoch/batch_size`. Similarly `nb_val_samples`->'val_steps' and `val_samples`->'steps' arguments have changed. Update your method calls accordingly.
/local_disk0/tmp/1587737721093-0/PythonShell.py:3: UserWarning: Update your 'fit_generator' call to the Keras 2 API: 'fit_generator(<generator..., validation_data=<generator..., steps_per_epoch=640, epochs=1, validation_steps=640)'
/databricks/python/lib/python3.7/site-packages/tensorflow_core/python/framework/indexed_slices.py:433: UserWarning: Converting sparse IndexedSlices to a dense Tensor of unknown shape. This may consume a large amount of memory.
  "Converting sparse IndexedSlices to a dense Tensor of unknown shape."
Epoch 1/1
1/640 [=====] - ETA: 5:13:17 - loss: 3.5994
2/640 [=====] - ETA: 3:35:58 - loss: 3.4952
3/640 [=====] - ETA: 2:36:07 - loss: 3.5296
4/640 [=====] - ETA: 2:06:08 - loss: 3.4843
5/640 [=====] - ETA: 1:48:18 - loss: 3.5153
6/640 [=====] - ETA: 1:36:13 - loss: 3.4977
7/640 [=====] - ETA: 1:27:28 - loss: 3.5088
8/640 [=====] - ETA: 1:20:53 - loss: 3.5039
9/640 [=====] - ETA: 1:15:45 - loss: 3.5092
622/640 [=====] - ETA: 1:02 - loss: 2.1929
623/640 [=====] - ETA: 59s - loss: 2.1929
624/640 [=====] - ETA: 55s - loss: 2.1928
625/640 [=====] - ETA: 52s - loss: 2.1928
626/640 [=====] - ETA: 48s - loss: 2.1925
627/640 [=====] - ETA: 45s - loss: 2.1923
628/640 [=====] - ETA: 41s - loss: 2.1919
629/640 [=====] - ETA: 38s - loss: 2.1918
630/640 [=====] - ETA: 34s - loss: 2.1914
631/640 [=====] - ETA: 31s - loss: 2.1913
632/640 [=====] - ETA: 27s - loss: 2.1911
633/640 [=====] - ETA: 24s - loss: 2.1909
634/640 [=====] - ETA: 20s - loss: 2.1909
635/640 [=====] - ETA: 17s - loss: 2.1910
636/640 [=====] - ETA: 13s - loss: 2.1909
637/640 [=====] - ETA: 10s - loss: 2.1909
638/640 [=====] - ETA: 6s - loss: 2.1911
639/640 [=====] - ETA: 3s - loss: 2.1908
640/640 [=====] - 2689s 4s/step - loss: 2.1905 - val_loss: 2.4281

```

Command took 44.87 minutes -- by harshmac.786@gmail.com at 4/25/2020, 1:20:33 AM on DataChef\_Upgrade

This screenshot above is model trained for 1 epoch with the same configuration of 640 steps and 32 samples in each step. This model took 45 minutes to train and had a training loss of 2.1905 and a validation loss of 2.4281.

We had trained 2 more models with 4 epochs which took 5 hours to train which had a training loss of 1.3609 and a validation loss of 1.5202. Though the loss for this model was low as compared to the above model, but the predictions were grammatically correct with our 1 epoch model.

Another model was trained on 2 epochs which took 2.25 hours which has a training loss of 1.78 and a validation loss of 1.93.

### 3. Prediction on sample

Now, we test our model to see our summarization of our model.

We select a random recipe from our data and try to predict the summary of recipe as a title.

TextSummarization (Python)

Detached | File | Edit | View: Code | Permissions | Run All | Clear | Schedule | Comments | Runs | Revision history

Cod 49

### Testing on Random Recipes

Show code

```
Randomly sampled recipe:
chipotle chili sauce
mayonnais chipotl chile place all the ingredients in a food processor or blender and purée until smooth transfer to a serving bowl if using right away or transfer to an airtight container and refrigerate if storing
```

## Prediction:

TextSummarization (Python)

Detached | File | Edit | View: Code | Permissions | Run All | Clear | Schedule | Comments | Runs | Revision history

Cod 51

### Prediction

```
samples = gensamples(
    skips=2,
    k=1,
    batch_size=2,
    short=False,
    temperature=1.,
    use_unk=True,
    model=model,
    data=(x, y),
    idx2word=idx2word,
    oov=oov0,
    glove_idx2idx=glove_idx2idx,
    vocab_size=vocab_size,
    nb_unknown_words=nb_unknown_words,
)

HEAD: chipotle chili sauce
DESC: mayonnais chipotl chile place all the ingredients in a food processor or blender and purée until smooth transfer to a serving bowl if using right away or transfer to an airtight container and refrigerate if storing
HEADS:
7.5237438678741455 bean
93.93346548880444 quick stew curry tarch enchiladas curried salad in he seca 22x3
Command took 1.82 seconds -- by harshmac.786@gmail.com at 4/25/2020, 2:20:23 AM on DataChef_Upgrade
```

We can see here that the second nearest prediction is grammatically correct but is no where near to our recipe title. Another example below.

TextSummarization (Python)

Detached | File | Edit | View: Code | Permissions | Run All | Clear | Schedule | Comments | Runs | Revision history

Cod 52

```
samples2 = gensamples(
    skips=2,
    k=1,
    batch_size=2,
    short=False,
    temperature=1.,
    use_unk=True,
    model=model,
    data=(x2, y2),
    idx2word=idx2word,
    oov=oov0,
    glove_idx2idx=glove_idx2idx,
    vocab_size=vocab_size,
    nb_unknown_words=nb_unknown_words,
)

HEAD: spaghetti pie ii
DESC: spaghetti butter chees egg beef onion bell pepper garlic tomato sugar oregano cottage cheese mozzarella cheese cook and drain spaghetti stir in margarine parmesan cheese and eggs while spaghetti is hot form spaghetti mixture into a crust in a buttered 10 inch pie plate preheat oven to 350 degrees f 175 degrees c in a skillet cook the beef onion green pepper and garlic drain off the fat and stir in the unstrained tomatoes tomato paste sugar and oregano heat through spread cottage cheese over the spaghetti crust then pour in the beef and tomato mixture bake at 350 degrees
HEADS:
20.687145447731018 chestnuts nutmeg
32.38890564441681 pancetta macaroni with beans salsa
Command took 1.12 seconds -- by harshmac.786@gmail.com at 4/25/2020, 2:26:46 AM on DataChef_Upgrade
```

So overall, we see that this model needs hyperparameter tuning to a great extent so that it can perform better with a loss value of lower than 1.

## Cuisine Classification

Cuisine classification is a technique of classifying cuisine based on the ingredients present in it. The concept or the idea can also be considered as an extension of classification problem. The extra layer of preprocessing converts the ingredients into numbers and generates a final data matrix. Later, different classification models are applied on the data to classify cuisines based on the ingredients. Finally, the best model based on evaluation matrix will be deployed.

### 1. Data Pre-Processing

As we have the cleaned ingredients and instructions in their respective containers stored in Cosmos DB, we need to establish a DB connection with Azure Databricks where all the preprocessing will take place and the model will be trained.

To do this we use the Azure-CosmosDB-Spark connector provided by Azure to do the same. Installing this on the cluster enables us to create a connection link between Azure Databricks and the container where data is stored in CosmosDB.

In the connection configuration, we need to specify the DB endpoint, master key, Database name, DB region where it is hosted, container name and the query to fetch the documents. Then this connection config is passed as an argument to the **spark.read.format** function to create a Spark data frame of the fetched documents from cosmos DB. Repeating this same process for both the containers that is cleaned ingredients and instructions. To play with the data in python, we created a Pandas data frame from spark data frame.

```
TextSummarization (Python)
Detached File Edit View: Code Permissions Run All Clear
Establish Connection with Cosmos DB
1 connectionConfig = {
2     "Endpoint" : "https://datachefcosmos.documents.azure.com:443/",
3     "Masterkey" : "f6Y6dj5xIdQyKdMrzBEkM5cFQ6xbdPGyrkh195YrCNrmvUkJ53CCh6ym6RQIXKrTMxM67UL7nf7yV9DmrXRf0A==",
4     "Database" : "DataChef",
5     "preferredRegions" : "East us",
6     "Collection" : "Instructions",
7     "schema_samplesize" : "1000",
8     "query_pagesize" : "200000",
9     "query_custom" : "SELECT * FROM Instructions"
10 }

Command took 0.02 seconds -- by harshmac.786@gmail.com at 4/24/2020, 7:20:51 PM on DataChef_Upgrade
Cmd 3
1 raw_instructions = spark.read.format("com.microsoft.azure.cosmosdb.spark").options(**connectionConfig).load()
2 #users.createOrReplaceTempView("raw_instructions")
▶ (1) Spark Jobs
▶ raw_instructions: pyspark.sql.dataframe.DataFrame = [attachments: string, _etag: string ... 8 more fields]
Command took 1.42 seconds -- by harshmac.786@gmail.com at 4/24/2020, 7:20:55 PM on DataChef_Upgrade
```

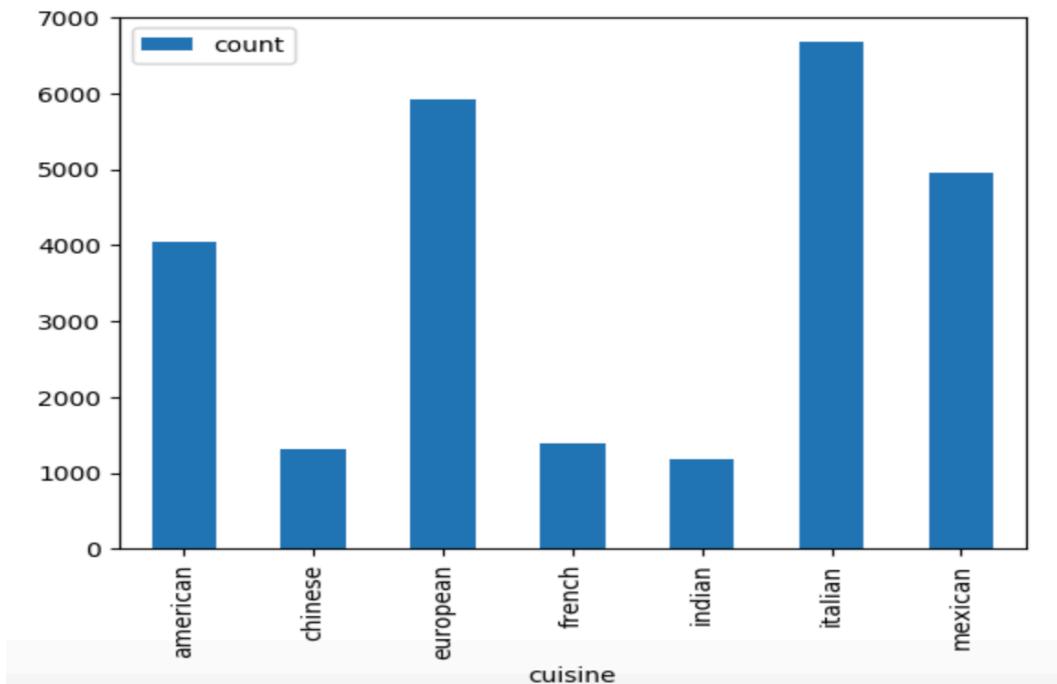
After converting the spark object to python's pandas dataframe, we are originally left with 40K entries. To keep our label count low to make our model makes better predictions, we have slices some of the cuisines from the main data frame. We have only kept seven major cuisines for example **Mexican, Italian, Chinese, Indian, European, American and French**. Alongside, we have also removed the duplicated entries present in our data to kill the bias in our model. Finally, we have left with 18,869 entries for our model to be trained on.

Our final panda's data frame looks as below:

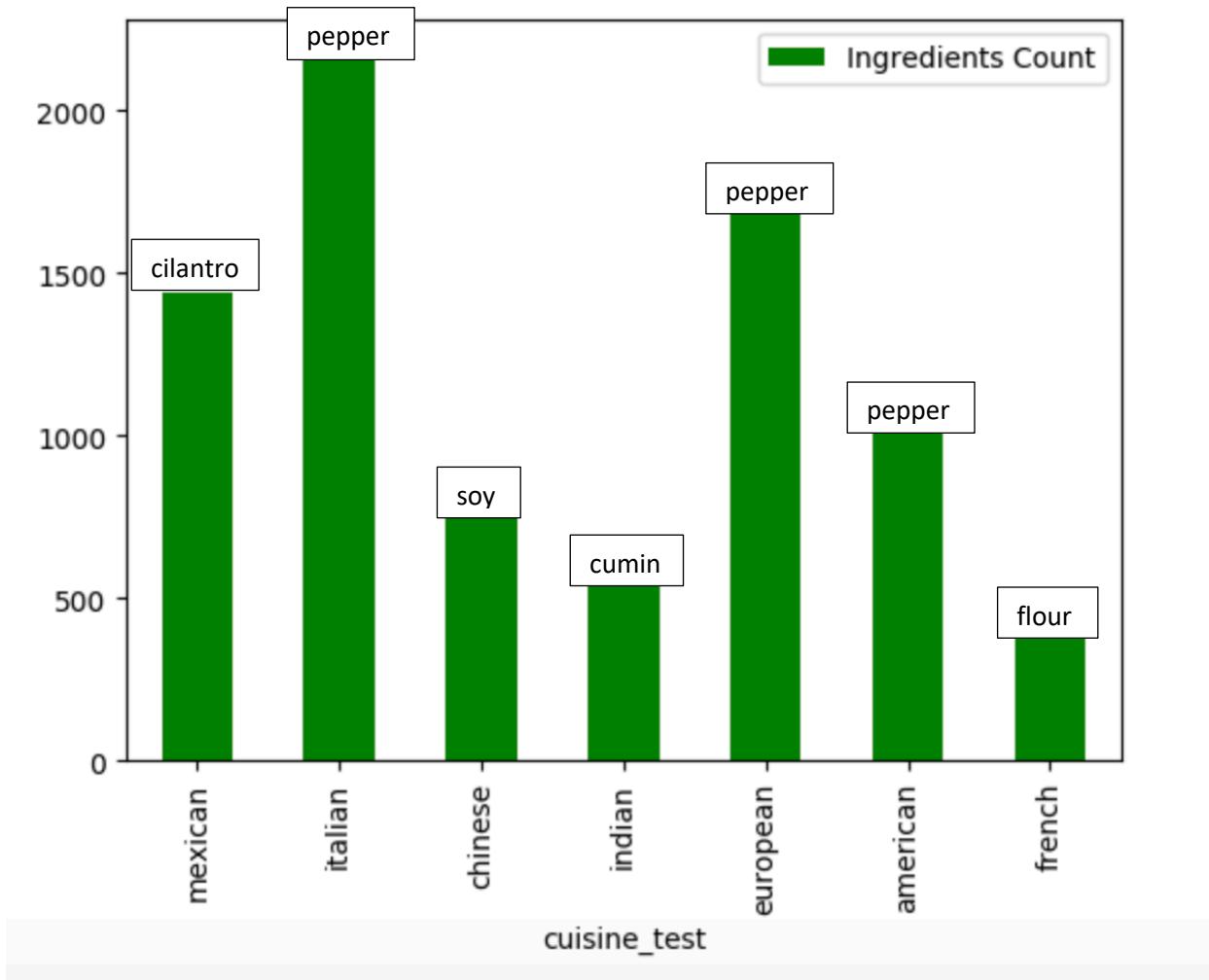
cuisine	id	ingredients	source	title
mexican	a73872f8-eb3e-4a50-9c01-12a7f3d7b9e6	[water, sugar, mint, lime, zest, salt]	AllRecipes	Watermelon-Mint Paletas
mexican	1fd48ed8-9af7-47ed-8e13-65b343a59c1f	[tomatillo, chile pepper, garlic, water, bacon...]	AllRecipes	Carne en su Jugo (Meat in its Juices)
mexican	b689840e-cd48-403f-822b-7a1b559d0040	[bone chicken breast, chicken broth, onion, ga...]	AllRecipes	Authentic Enchiladas Verdes
mexican	b72e97bb-20c7-49ca-a73e-0affa0a61d20	[water, cinnamon, clove, star anis, milk, tabl...]	AllRecipes	Champurrado

### Visualization:

- As We have kept seven major cuisines which are Mexican, Italian, Chinese, Indian, European, American and French. The below graph shows the count of each cuisine type in the data frame.



- The below graph shows the ingredients count in each cuisine type which is majorly present. For e.g.- Mexican use the cilantro most, Chinese use the soy most and Indian use the cumin.



**Label Dictionary:** Now we have created a dictionary that will save all the unique cuisines present in our data. We have taken only seven cuisines so the length of the dictionary would be seven. In the dictionary, name of cuisines will be the keys and their numerical encoding like 0,1,2 etc would be the value as per the cuisine. This dictionary will be later used to create labels for the classification model

The code below explains the part of creating and populating dictionary:

```

1 # a label dictionary to be used later for creating labels
2 label_dictionary = {}
3
4 unique_cuisine = raw_df['cuisine'].unique() # picking only unique cuisines
5 for index, uc in enumerate(unique_cuisine): # populating label dictionary
6     label_dictionary[uc] = index
7

```

**Ingredients Dictionary:** Here we have created a dictionary to save all the unique ingredients present in our data. As per the processing, there are approximately 8K unique ingredients in our data set. In the dictionary, name of ingredients is represented as keys and their numerical encoding such as 0,1,2 etc would be the values. This ingredients dictionary will be later used to create the features columns for the final matrix of classification model.

The code below explains the part of creating and populating dictionary:

```

1 def createIngredientDict(ingredient_corpus): # function that creates ingredient dictionary
2     ingredient_dictionary = {} # ingredients along with their indexes
3     for index, ing in enumerate(ingredient_corpus):
4         ingredient_dictionary[ing] = index
5
6     return ingredient_dictionary
7

```

**Ingredient-To-Vector Function:** Now as we know that a computer doesn't understand text and hence there was need to convert the entire list of ingredients to number. So, we have created a function called Ingredient\_to\_vector function that takes list of ingredients as an input and convert it into a binary vector by matching the value with the ingredient dictionary. This binary vector will be of shape (number of ingredients, 1). These different vectors for all the entries will act as a unique row in final data matrix.

The functional code to implement ingredientToVector is as follows:

```

def ingredientToVector(ingredient_list, ingredient_corpus): # function to convert a list of ingredient to a binary vector
    ingredient_dictionary = createIngredientDict(ingredient_corpus)
    ingredient_vector = np.zeros((len(ingredient_dictionary), 1))
    indices_list = []
    for i in ingredient_list:
        if i in ingredient_dictionary:
            indices_list.append(ingredient_dictionary[i])
    for i in indices_list:
        ingredient_vector[i] = 1

    return ingredient_vector.T

```

**Final Matrix:** At last we have created final matrix using all the helper code written above. In this final matrix, ingredients will as feature in columns, cuisines will go a label for every entry, and ingredients list will be converted to binary vector which will than be later inserted into the matrix as an entry. The approximate dimension of our final matrix is (18K, 8K).

The code to create final matrix is shown below:

```

1 # creating a final matrix our from the data frame
2 final_matrix = np.zeros((len(raw_df), len(ingredient_corpus)+1))
3 for i in range(len(raw_df)):
4     final_matrix[i, 0:len(ingredient_corpus)] = ingredientToVector(raw_df.iloc[i, :]['ingredients'], ingredient_corpus)
5     final_matrix[i, len(ingredient_corpus)] = label_dictionary[raw_df.iloc[i, :]['cuisine']]

```

**Data Splitting:** As a pre-defined practice in machine learning, we always divide the data into train and test split to evaluate the model. In this part we have split our entire data matrix into train and test part. 80% of our data goes in training set and 20% in testing set.

Code to split the data is shown below:

```
1 features = final_matrix[:, 0:final_matrix.shape[1]-1] # feature matrix
2 labels = final_matrix[:, -1] # label vector
3
4 X_train, X_test, Y_train, Y_test = train_test_split(features, labels, test_size = 0.20) # splitting the data into training and testing set
```

**Dimensionality Reduction:** Model with large number of features inversely effects the computation and speed of algorithm. So, we have decided to reduce the number of features to optimize our algorithms. We have used principal component analysis by retaining 95% of the variance in the original data. Once done, total number of features have been reduced to 1.8K from 8K where at the same time keeping 95% of the variance retained.

```
1 pca = PCA(0.95) # instantiating pca while retaining 95% of variance in the data
2 pca.fit(X_train)
3 X_train = pca.transform(X_train) # transforming training set
4 X_test = pca.transform(X_test) # transforming testing set
```

## 2. Model Selection:

Once our data is cleaned and pre-processed. We have decided to apply the different classification models on it. We have chosen few classification models such as Support vector machine with both linear and gaussian kernel, Neural network, logistic regression, random forest, decision tree etc.

Let's go through each model one by one and see how they performed on our data:

1. **Support Vector Machine – Linear Kernel:** First we have applied the support vector machine with linear kernel and evaluated the model based on training and testing set accuracy

**Training accuracy: 83%**

**Testing accuracy: 72%**

Support vector machine without any regularization and optimization gave us decent result with training and testing accuracy of 83% and 72% respectively. This explains that our data is might be linearly separable and nonlinear functions and algorithms would not do well on it.

The code is shown below:

## Support Vector Machine (Linear Kernel)

```
1 svm = SVC(kernel="linear")
2 svm.fit(X_train, Y_train)
```

[Show result](#)

Cmd 13

```
1 train_score = accuracy_score(Y_train, svm.predict(X_train))
2 print("Accuracy on training set: ", train_score)
3 test_score = accuracy_score(Y_test, svm.predict(X_test))
4 print("Accuracy on testing set: ", test_score)
```

Accuracy on training set: 0.8306724080821464

Accuracy on testing set: 0.7231054583995761

Command took 6.57 minutes -- by utkarshkekre@gmail.com at 4/25/2020, 10:22:08 AM on DataChef\_Upgrade

2. **Support Vector Machine – Gaussian Kernel:** As we have sensed in our first model that we might have a linear separation in our data and so the nonlinear functions might not perform well on our data. Let's see, what happened when we applied gaussian kernel SVM on our data.

**Training accuracy: 70%**

**Testing accuracy: 47%**

As expected, support vector machine with non-linear kernel didn't perform well on both training and testing accuracy.

## Support Vector Machine (Gaussian Kernel)

```
1 svm_rbf = SVC(kernel="rbf")
2 svm_rbf.fit(X_train, Y_train)
```

[Show result](#)

Cmd 15

```
1 train_score = accuracy_score(Y_train, svm_rbf.predict(X_train))
2 print("Accuracy of training set: ", train_score)
3 test_score = accuracy_score(Y_test, svm_rbf.predict(X_test))
4 print("Accuracy of testing set: ", test_score)
```

Accuracy of training set: 0.4741967538920172

Accuracy of testing set: 0.4785373608903021

Command took 10.50 minutes -- by utkarshkekre@gmail.com at 4/25/2020, 10:42:25 AM on DataChef\_Upgrade

3. **Neural Network:** Fit the neural network with 3 hidden layers of 100 units each. Activation function we have used is the relu activation. We have used random initialization, but we haven't used any regularization in this. Let's see the evaluation matrix.

**Training accuracy: 99%**

**Testing accuracy: 75%**

Neural network without any hyperparameter tuning gave us the satisfying result. We have expected that along with "He" initialization and L2 regularization can dramatically improve our testing accuracy as well.

## Neural Network

```
1 nn = MLPClassifier(hidden_layer_sizes = (100, 100, 100))
2 nn.fit(X_train, Y_train)
```

[Show result](#)

Cmd 17

```
1 train_score = accuracy_score(Y_train, nn.predict(X_train))
2 print("Accuracy of training set: ", train_score)
3 test_score = accuracy_score(Y_test, nn.predict(X_test))
4 print("Accuracy of testing set: ", test_score)
```

Accuracy of training set: 0.9942365021530308

Accuracy of testing set: 0.7011128775834659

Command took 0.25 seconds -- by utkarshkekre@gmail.com at 4/25/2020, 11:13:05 AM on DataChef\_Upgrade

4. **Decision Tree:** Fit the decision tree on our data and got the following results:

**Training accuracy: 99%**

**Testing accuracy: 45%**

Decision tree without any regularization and optimization is showing a problem of overfitting. So, even with regularization it is not possible to have a significant rise in testing accuracy. So, we have not decided to go ahead with this model.

## Decision Tree

```
1 decision_tree_clf = DecisionTreeClassifier()  
2 decision_tree_clf.fit(X_train, Y_train)
```

[Show result](#)

---

Cmd 19

```
1 train_score = accuracy_score(Y_train, decision_tree_clf.predict(X_train))  
2 print("Accuracy of training set: ", train_score)  
3 test_score = accuracy_score(Y_test, decision_tree_clf.predict(X_test))  
4 print("Accuracy of testing set: ", test_score)
```

Accuracy of training set: 0.9986088108645247

Accuracy of testing set: 0.43614202437731847

Command took 0.16 seconds -- by utkarshkekre@gmail.com at 4/25/2020, 11:15:07 AM on DataChef\_Upgrade

5. **Random Forrest:** Fitting random forest with an n\_estimators value of 100 gave us the following accuracy on training and testing data.

**Training accuracy: 99%**

**Testing accuracy: 65%**

Tough this model has performed better than some of the models like decision tree and SVM with gaussian kernel still it is not better than logistic regression or neural network. So, we thought of dropping this model

## Random Forest

```
1 random_forest = RandomForestClassifier(n_estimators = 100)
2 random_forest.fit(X_train, Y_train)
```

Show result

Cmd 21

```
1 train_score = accuracy_score(Y_train, random_forest.predict(X_train))
2 print("Accuracy of training set: ", train_score)
3 test_score = accuracy_score(Y_test, random_forest.predict(X_test))
4 print("Accuracy of testing set: ", test_score)
```

Accuracy of training set: 0.9986088108645247  
Accuracy of testing set: 0.624006359300477

Command took 1.42 seconds -- by utkarshkekre@gmail.com at 4/25/2020, 11:17:00 AM on DataChef\_Upgrade

6. **Logistic Regression:** At last we have applied logistic regression hoping that it would give us some satisfying result as our data is linearly separable.

**Training accuracy: 80%**

**Testing accuracy: 72%**

Most of the model showed overfitting for the data. Whereas logistic regression is the only model that showed nearly good fit without regularization. And, this is a no surprise as our data is linearly separable and logistic regression is best tool for fitting linearly separable data. So, we have decided to include this model.

## Logistic Regression

```
1 log_reg = LogisticRegression(penalty="l2")
2 log_reg.fit(X_train, Y_train)
```

Show result

Cmd 23

```
1 train_score = accuracy_score(Y_train, log_reg.predict(X_train))
2 print("Accuracy of training set: ", train_score)
3 test_score = accuracy_score(Y_test, log_reg.predict(X_test))
4 print("Accuracy of testing set: ", test_score)
```

Accuracy of training set: 0.7953627028817489  
Accuracy of testing set: 0.7286698463169051

Command took 0.12 seconds -- by utkarshkekre@gmail.com at 4/25/2020, 11:18:52 AM on DataChef\_Upgrade

### **3. Python Modules:**

We have used built-in python modules in this project for cleaning, pre-processing and modelling of data. This list of the modules is here as follows:

1. **BeautifulSoup**: To scrape the data from the website
2. **Pandas**: Python's built in library to handle large data sets
3. **Numpy**: Python's library for scientific computation
4. **Sci-kit Learn**: Python's module for model selection, data modelling and evaluation

## **Visualizations**

### **Building Connection:**

As we have the cleaned nutrition's and the content in the nutrients in the containers stored in Cosmos DB, we need to establish a DB connection with Azure Databricks where all the visualization will take place.

To do this we use the Azure-CosmosDB-Spark connector provided by Azure to do the same. Installing this on the cluster enables us to create a connection link between Azure Databricks and the container where data is stored in CosmosDB.

In the connection configuration, we need to specify the DB endpoint, master key, Database name, DB region where it is hosted, container name and the query to fetch the documents. Then this connection config is passed as an argument to the **spark.read.format** function to create a Spark data frame of the fetched documents from cosmos DB. Repeating this same process for both the containers that is cleaned nutrition's. To play with the data in python, we created a Pandas data frame from spark data frame.

### **Nutritional Evaluation:**

Nutritional Assessment is important to have a complete picture of the food-related health and total body wellness. It is recommended to have information of "why and when to eat" and not just for "what to eat".

According to recent studies, most Americans don't get enough nutrients but are over-the-top on the amount of sodium and fat in their diets. Researchers found that many American adults fell short of consuming enough vitamin A, vitamin C, vitamin D, calcium, and iron. In addition, Americans have a tendency to eat far more saturated fat, cholesterol, and sodium than is normal or recommended.

There are few symptoms which depicts that you have a dietary issue:

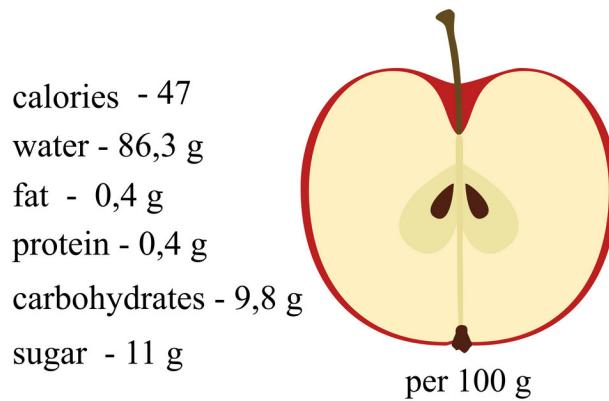
- Ongoing fatigue that keeps you from feeling rested.
- Constipation or recurrent diarrhea.
- Blood sugar levels that fluctuate greatly throughout the day.
- Low grade aches and pains with no other cause.
- Lowered immune system leaving you susceptible to illnesses.
- Irritability, depression, and other mood issues.
- Obesity and excess weight Dizziness.

- Headaches, migraines, and brain fog.
- Daily heartburn with no other cause

To keep a track on which recipe has what nutrition content, we suggest using the below suggested evaluations to display on web pages, as it will result in a better user experience and users can track the nutritional value of the dish they are going to consume. It will be beneficial for:

- Health Conscious Customers.
- Calorie Counting Customers.
- Keeping it Legal, as many Canadian food companies have recently been subjected to a healthy ban on trans-fats.
- Professionalism.
- Doing the Right Thing as you are being honest and accountable.

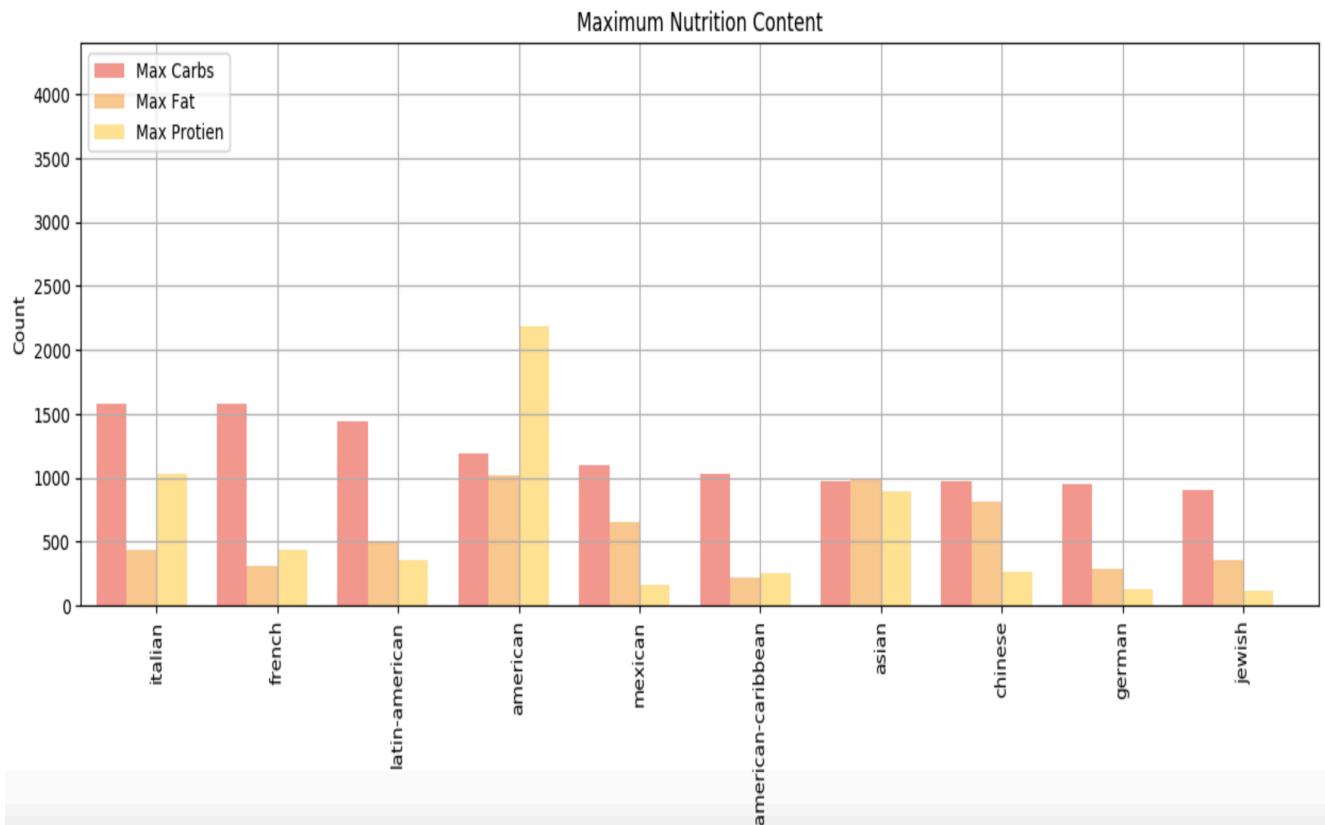
## Apple. Nutritional value.



## Recommended Evaluations to Display on Website:

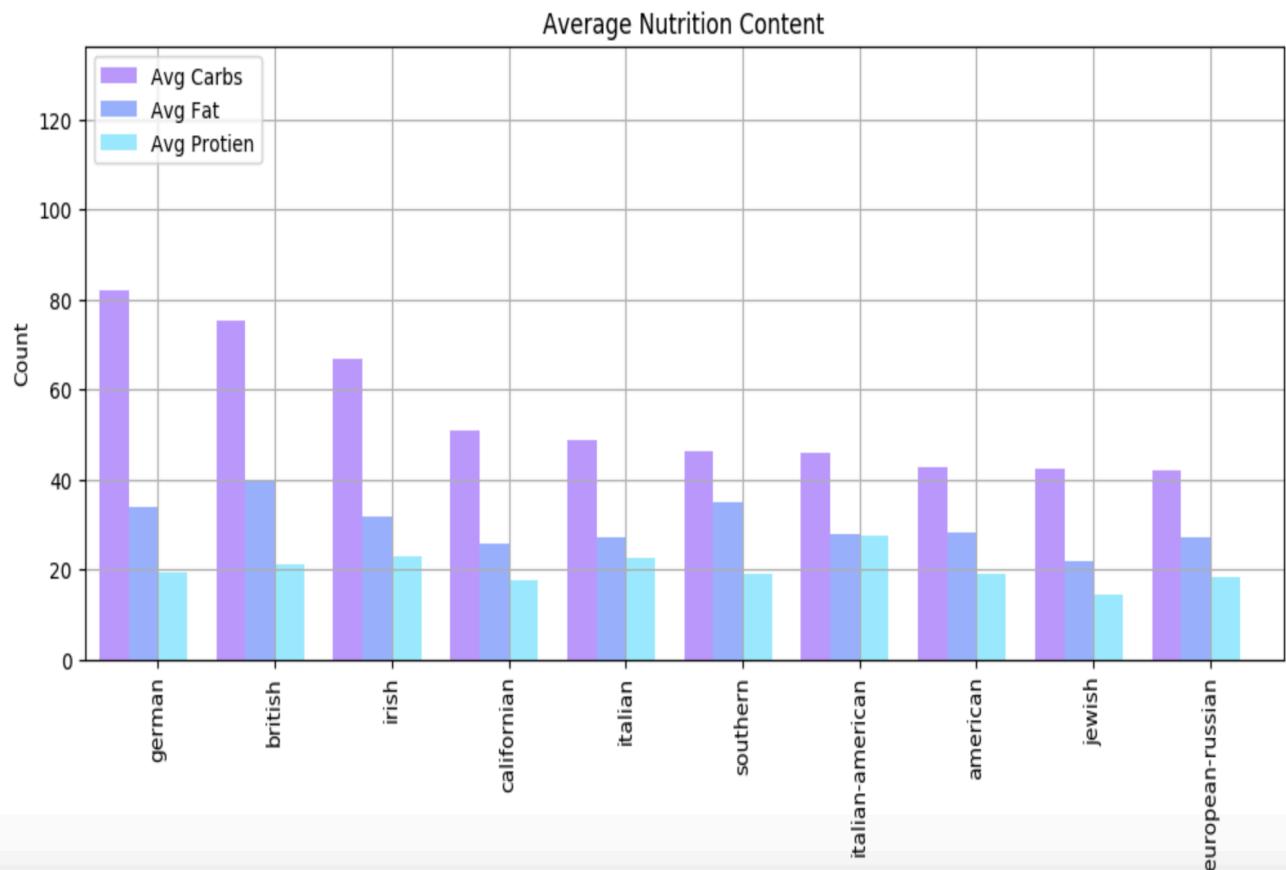
### 1. Maximum Nutrition Content

The graph shows the top 10 cuisine type with maximum carbs, maximum fat and maximum Protein. This will help customer to calculate their nutrition content and decide what to eat, whether they want to eat something which has maximum carbs (cuisine type like- Italian, French) or which has maximum fat (cuisine type like- American, Asian) or which has maximum protein (cuisine type like- American, Italian).



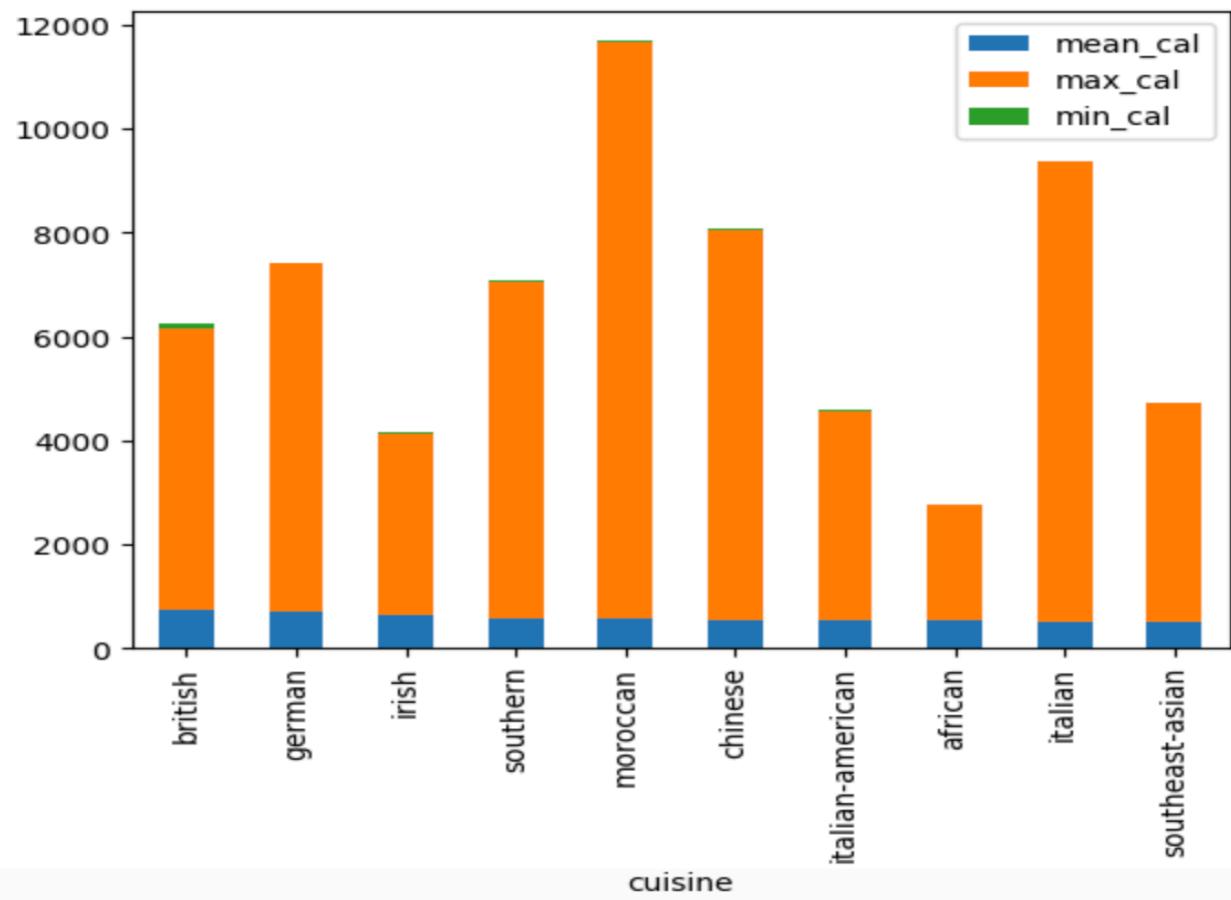
## 2. Average Nutrition Content

The graph shows the top 10 cuisine type with average carbs, average fat, and average Protein. It is strange to know that cuisine type like German and British has the most stable average content of all of the three nutrient content which are carbs, fat and protein.



### 3. Calories Content

This graph shows the maximum, minimum and average calories content in the top 10 cuisine type. Having a calorie count is extremely important as this will help one to keep track of calorie in and calorie out concept. Cuisine type like Moroccan, Italian and Chinese has the maximum calorie count.



## **Business Values:**

We are targeting to deliver the best to our customers and this solution will generate business values in 5 ways:

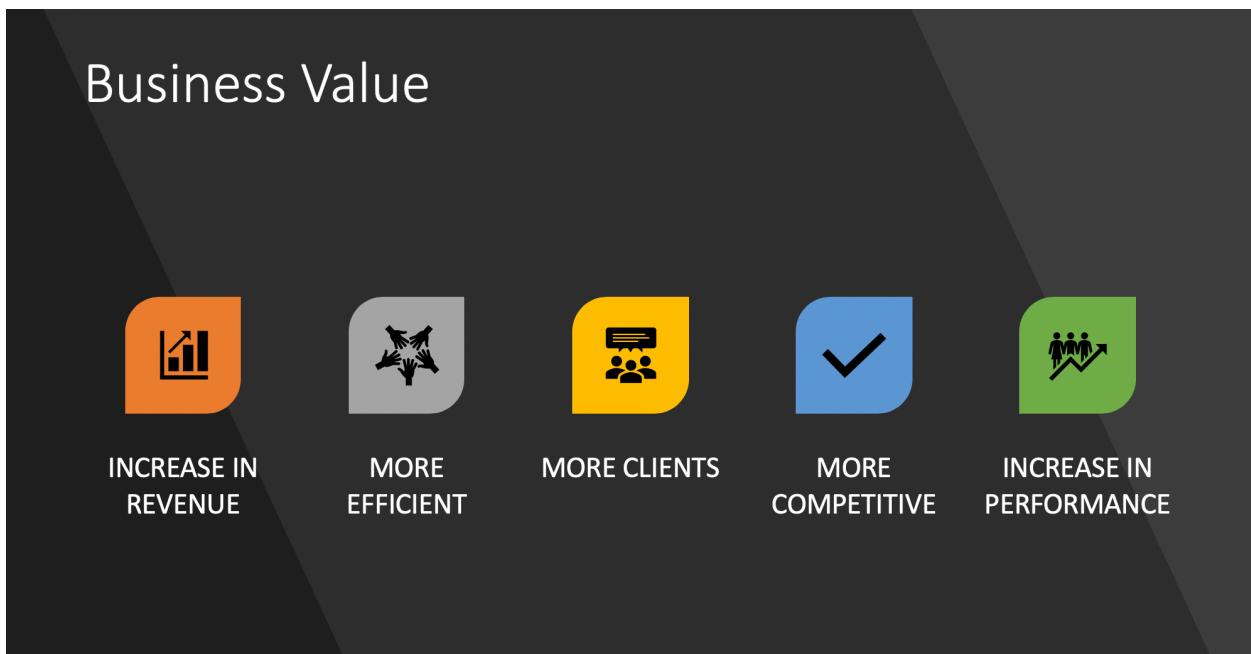
**Increase in Revenue:** As we provide more facilities to the users, it will attract more users due to mouth to mouth publicity and also, marketing from online platforms. The increase in customer base will result in increase of the revenue.

**More Efficient:** We are targeting to give a better user experience when anyone uses the website, we are giving a meaningful title to the recipe, which will reduce the screen time for the user and user can be more productive in the time they spent on the website. This will earn trust for the company.

**More Clients:** Increasing the user satisfaction will result in attracting more clients and increase the customer base of the company.

**More Competitive:** As we are also going to provide the nutritional values evaluation for cuisines and recipes, users can use it as their calorie counter and can make use of it for planning their diet. This feature will give competition to the companies like HealthifyMe, Fitocracy etc., henceforth, it will be a strong competitor in the market.

**Increase in Performance:** As the above-mentioned values will contribute in the success of the food website company, shortly the reduced screen time, less clicks, and more information will result in alleviating the performance of the website, hence increasing the customer satisfaction.



## **References:**

- <https://expandedramblings.com/index.php/allrecipes-facts-statistics/>
- <https://www.similarweb.com/website/epicurious.com>
- <https://en.wikipedia.org/wiki/Allrecipes.com>
- <https://en.wikipedia.org/wiki/Epicurious>
- <https://www.netwoven.com/2018/12/04/a-complete-walkthrough-of-azure-cosmos-db-and-why-should-you-use-it/>
- <https://milehighspine.com/the-importance-of-a-nutritional-assessment/>
- <https://www.menusano.com/nutrition-blog/6-ways-providing-nutritional-value-information-increases-sales-and-grows-your-food-service-business/>
- <https://nlp.stanford.edu/projects/glove/>