

Business Objective

Input Data

- ```
display first five rows of the dataframe - df
inputDF.head(5)
```

[illegible]

 Generate
  Code
  Markdown

```
print the dimensions of dataframe - df
inputDF.shape
```

Python

```
print the information of the dataframe
inputDF.info()
```

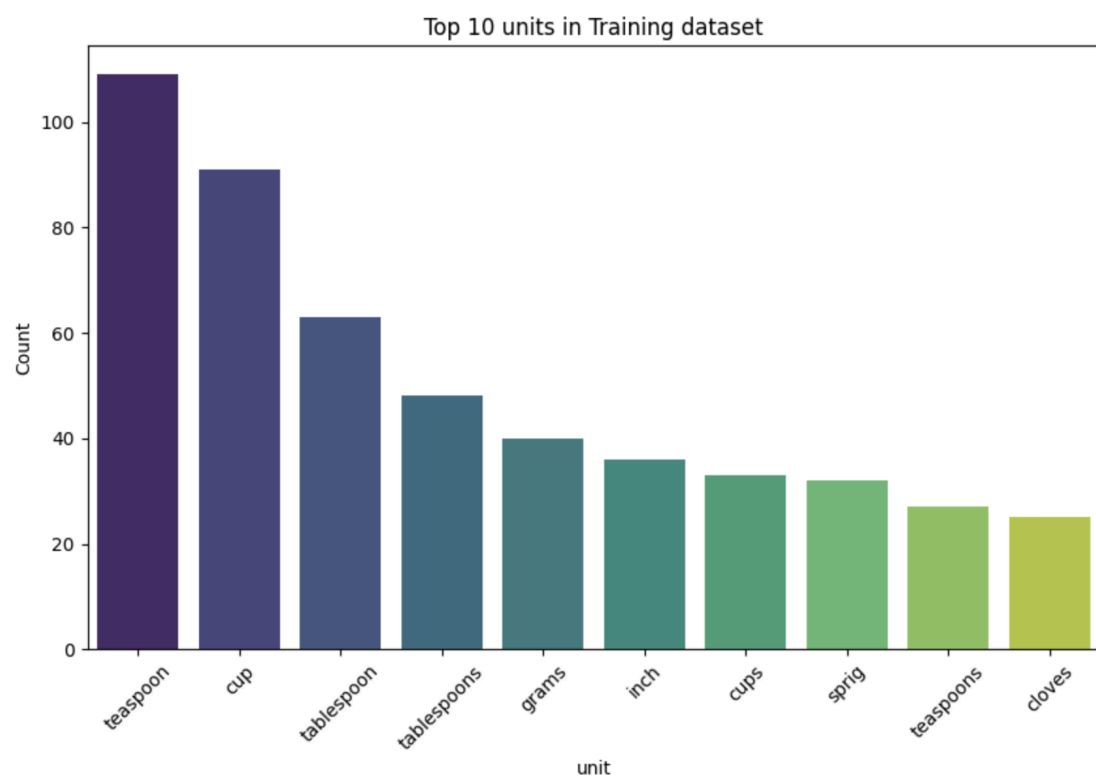
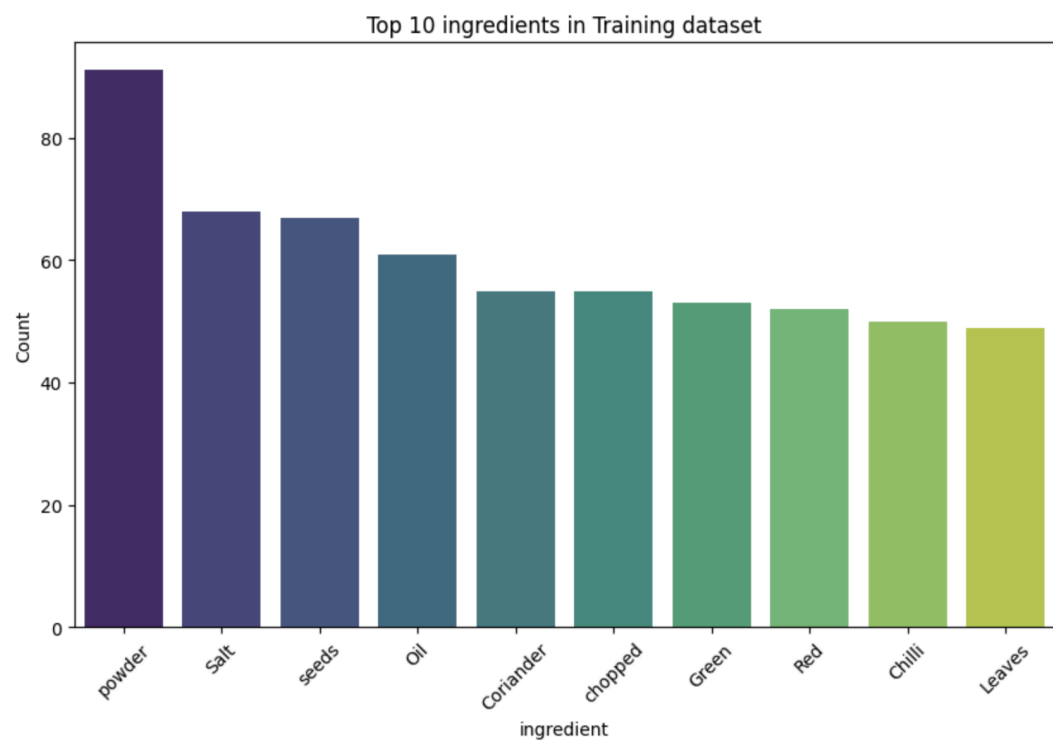
```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 285 entries, 0 to 284
Data columns (total 2 columns):
 # Column Non-Null Count Dtype
--- -
 0 input 285 non-null object
 1 pos 285 non-null object
dtypes: object(2)
memory usage: 4.6+ KB
```

- There were 81 records for which the pos\_tokens length was not equal to input\_text length. So, they were dropped.

## Data Analysis

```
First 10 records in Training dataset:
Input Token --> POS Token
1 --> quantity
cup --> unit
Gram --> ingredient
flour --> ingredient
besan --> ingredient
1/2 --> quantity
teaspoon --> unit
Ajwain --> ingredient
Carom --> ingredient
seeds --> ingredient
```

```
First 10 records in Validation dataset:
Input Token --> POS Token
1 --> quantity
cup --> unit
rice --> ingredient
12 --> quantity
small --> unit
onions --> ingredient
2 --> quantity
cloves --> ingredient
garlic --> ingredient
inch --> unit
```



## Observations

- Powder is the top most ingredient in the entire dataset
- Along with powder, Salt & Seeds are also mostly commonly used
- Teaspoon, cup & tablespoon are the most units

## Feature Engineering

- We identified common Unit and Quantity keywords by analyzing frequent terms.
- We defined regular expressions to match different forms of Units
- We also converted words to features to analyse the impact

## Core Features

- Token properties
- Character patterns

## Quantity & Unit Detection

- Regular Expressions for Quantities
- Numeric pattern recognition for Units

## Contextual Features

- Boundary Markers
- Previous & Next words/tokens

## Model Training

CRF model is used to train the recipe data. Below are the parameters initialized:

```
initialise CRF model with the specified hyperparameters and use weight_dict
crfModel = sklearn_crfsuite.CRF(
 algorithm='lbfgs',
 c1=0.5,
 c2=1.0,
 max_iterations=100,
 all_possible_transitions=True
)

train the CRF model with the weighted training data
crfModel.fit(XTrainWeightedFeatures, yTrainLabels)
```

CRF

```
CRF(algorithm='lbfgs', all_possible_transitions=True, c1=0.5, c2=1.0,
 max_iterations=100)
```

## Model Evaluation

### Training Data Results

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| ingredient   | 1.00      | 1.00   | 1.00     | 3411    |
| quantity     | 1.00      | 1.00   | 1.00     | 658     |
| unit         | 0.99      | 1.00   | 1.00     | 563     |
| accuracy     |           |        | 1.00     | 4632    |
| macro avg    | 1.00      | 1.00   | 1.00     | 4632    |
| weighted avg | 1.00      | 1.00   | 1.00     | 4632    |

Confusion Matrix:

Labels ['unit', 'ingredient', 'quantity']

```
[[562 0 1]
 [0 3411 0]
 [3 0 655]]
```

## Validation Data Results

```
specify flat classification report
yValTrueFlat = [label for sent in yValLabels for label in sent]
yValPredFlat = [label for sent in yValPred for label in sent]
print(classification_report(yValTrueFlat, yValPredFlat))
```

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| ingredient   | 1.00      | 1.00   | 1.00     | 1611    |
| quantity     | 1.00      | 0.99   | 1.00     | 294     |
| unit         | 0.99      | 1.00   | 1.00     | 244     |
| accuracy     |           |        | 1.00     | 2149    |
| macro avg    | 1.00      | 1.00   | 1.00     | 2149    |
| weighted avg | 1.00      | 1.00   | 1.00     | 2149    |

Confusion Matrix:

Labels ['unit', 'ingredient', 'quantity']

```
[[244 0 0]
 [0 1611 0]
 [2 0 292]]
```

Validation Accuracy is very good - **0.9906**

This indicates that the Model has been trained very well.

## Error Analysis

Label: quantity, ErrorCount: 2, Class Weight: 7.0395

Error DataFrame:

|   | token  | true_label | predicted_label | prev_token | next_token | class_weight |
|---|--------|------------|-----------------|------------|------------|--------------|
| 0 | little | quantity   | unit            | leaves     | Salt       | 7.039514     |
| 1 | taste  | quantity   | unit            | per        | 1/2        | 7.039514     |

## Conclusion

- Model shows strong performance in recognising recipe entities.
- Model Accuracy is very high on both Training & Test data sets
- Recommend this model to be tested on other real world practical data.