A

Synopsis/Project Report

On

# Creating a WordCloud using NLP and TF-IDF in Python

Submitted in partial fulfilment of the requirement for the IV semester

**Bachelor of Computer Science**

By

**Lalit Dumka**

Under the Guidance of

**Mr. Ravindra Koranga**

**Assistant Professor**

**Department of CSE**



# DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING

# GRAPHIC ERA HILL UNIVERSITY, BHIMTAL CAMPUS

## SATTAL ROAD, P.O. BHOWALI,

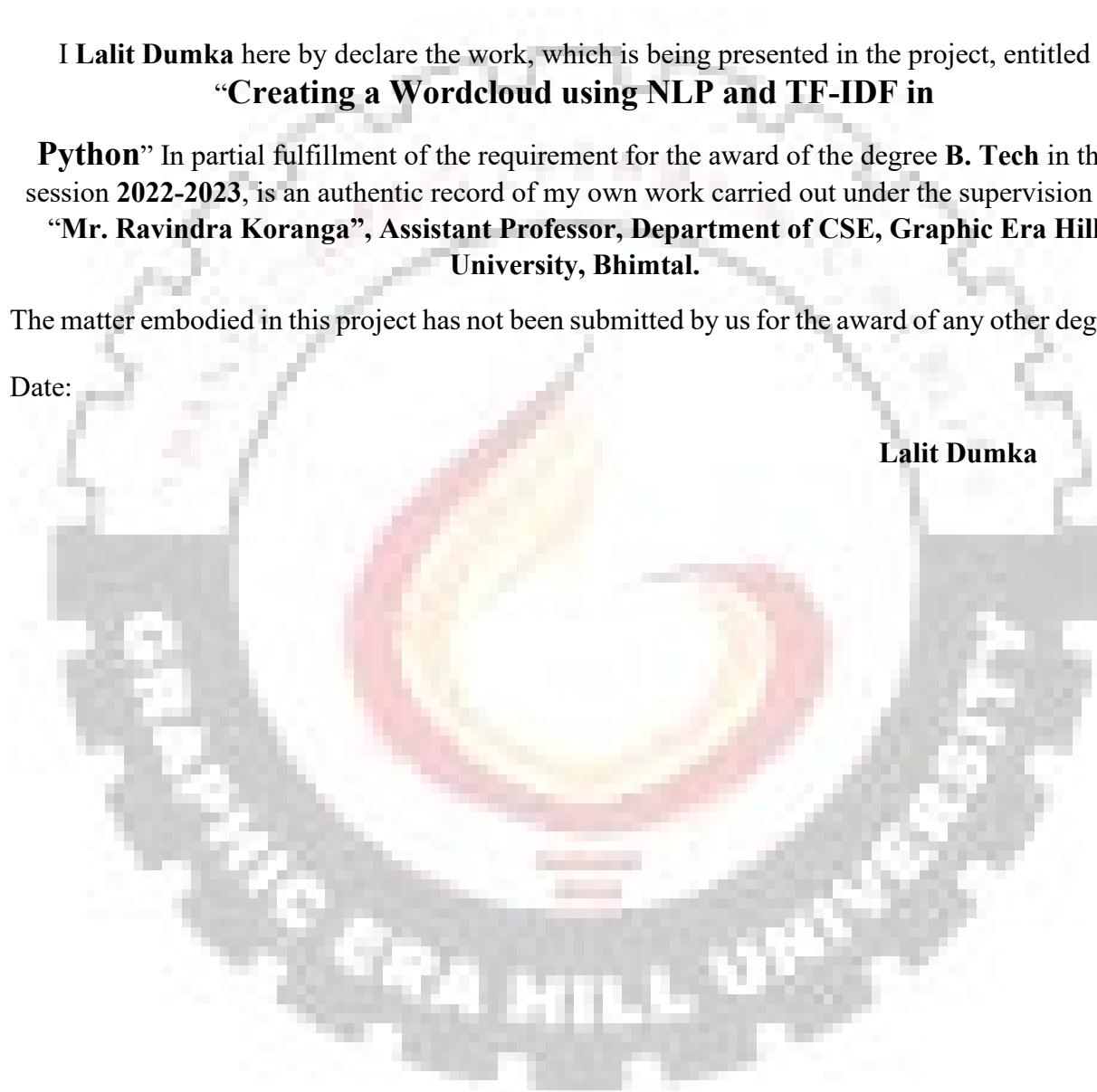## DISTRICT- NAINITAL-263132

## 2022- 2023

# STUDENT'S DECLARATION

I **Lalit Dumka** here by declare the work, which is being presented in the project, entitled
"**Creating a Wordcloud using NLP and TF-IDF in**

**Python**" In partial fulfillment of the requirement for the award of the degree **B. Tech** in the session **2022-2023**, is an authentic record of my own work carried out under the supervision of "**Mr. Ravindra Koranga**", **Assistant Professor, Department of CSE, Graphic Era Hill University, Bhimtal.**

The matter embodied in this project has not been submitted by us for the award of any other degree.

Date:

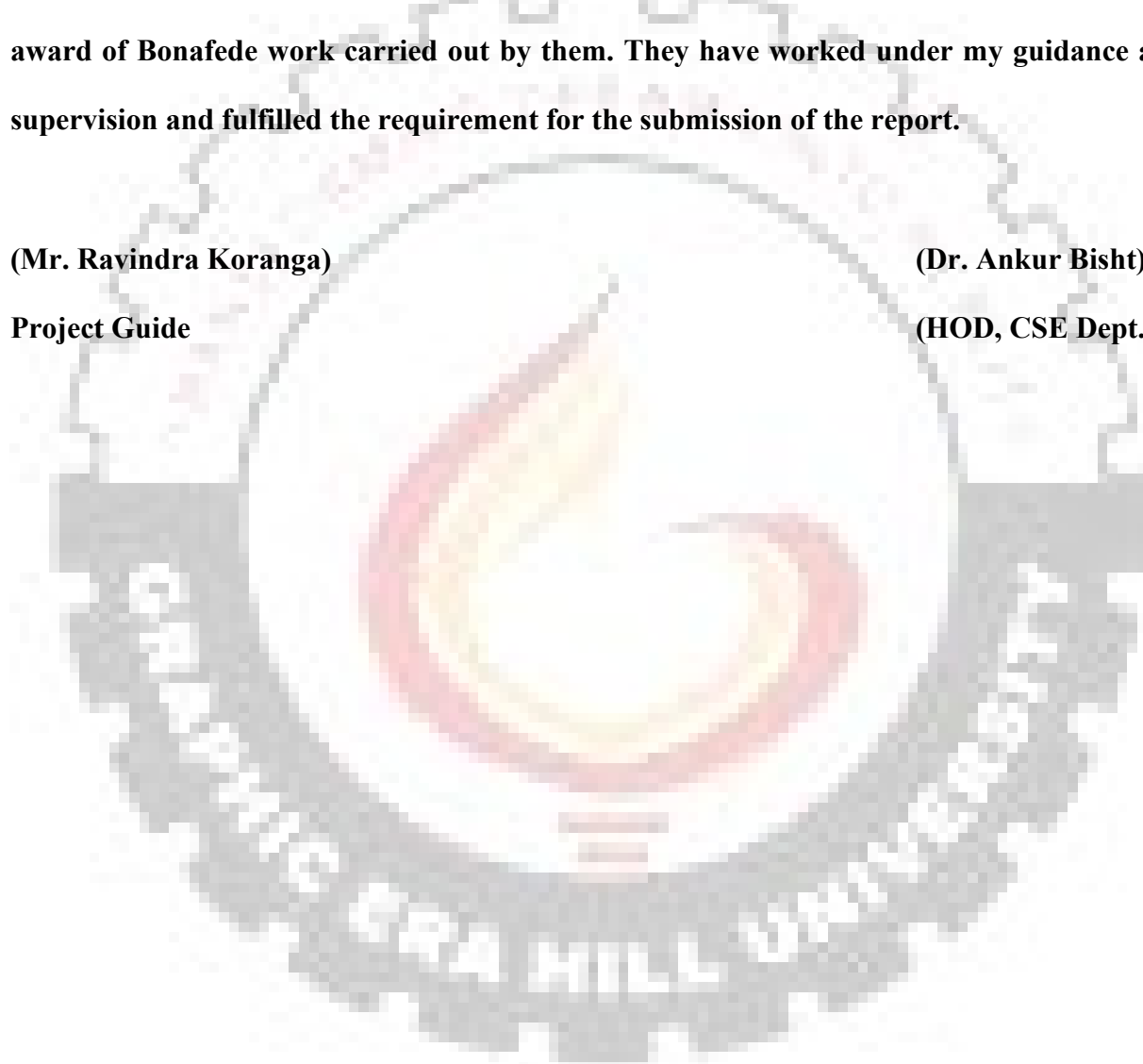**Lalit Dumka**

# CERTIFICATE

The project report entitled "Creating a wordcloud using NLP and TF-IDF in Python " Being submitted by <u>Lalit Dumka</u> to Graphic Era Hill University Bhimtal Campus for the award of Bonafede work carried out by them. They have worked under my guidance and supervision and fulfilled the requirement for the submission of the report.


**(Mr. Ravindra Koranga)**                                       **(Dr. Ankur Bisht)**

**Project Guide**                                               **(HOD, CSE Dept.)**

# ACKNOWLEDGEMENT

# TABLE OF CONTENTS

# PROJECT ABSTRACT

In this small project, I aim to leverage Natural Language Processing (NLP) techniques and the TF-IDF algorithm to generate an insightful word cloud from a given text corpus. WordClouds are powerful visual representations that highlight the most significant words in a document, providing a quick overview of its main themes.

To accomplish this, I will employ Python and its popular NLP libraries, such as NLTK (Natural Language Toolkit) and scikit-learn. The project will involve several key steps. Firstly, I will preprocess the text data by removing stop words, punctuation, and performing tokenization. Then, utilizing the TF-IDF algorithm, I will compute the importance score for each word in the corpus. The TF-IDF score takes into account both the term frequency and inverse document frequency, helping to identify words that are prevalent in a specific document while being relatively rare across the entire corpus.

# INTRODUCTION

The objective of this small project is to create a wordcloud using Natural Language Processing (NLP) techniques and the TF-IDF algorithm in Python. NLP is a field of artificial intelligence that focuses on the interaction between computers and human language. TF-IDF, short for Term Frequency-Inverse Document Frequency, is a widely used numerical statistic in NLP that evaluates the importance of a word within a document relative to a larger corpus.

By combining the power of NLP and TF-IDF, we can create a wordcloud that not only showcases the most prominent words in a text corpus but also emphasizes words that are unique to a particular document while being infrequent in the overall corpus. This approach allows us to capture the essence of the text and visualize its main themes at a glance.

The resulting wordcloud can be a valuable tool in various applications, including sentiment analysis, content summarization, and data exploration. It enables researchers, analysts, and data scientists to quickly identify the significant keywords in a document and gain valuable insights from large volumes of text.

# OBJECTIVE

The objective of this project is to leverage Natural Language Processing (NLP) techniques and the TF-IDF algorithm in Python to create a wordcloud that effectively summarizes the key themes and important words within a given text corpus.

This project aims to provide a practical demonstration of how NLP techniques and the TF-IDF algorithm can be applied to create a wordcloud that effectively summarizes the main themes and important keywords within a given text corpus. The resulting wordcloud can serve as a valuable tool for data exploration, content summarization, and information visualization in various domains.

# PROBLEM STATEMENT

In a large text corpus, there is a need to identify and extract the most significant keywords that are representative of the main themes within the corpus. Traditional keyword extraction methods based solely on word frequency may not accurately capture the importance and relevance of specific terms within the corpus.

The problem is to develop a solution that effectively analyzes the text corpus and utilizes NLP and TF-IDF to compute the importance score for each word. By considering both the term frequency and inverse document frequency, the solution should accurately identify keywords that are highly relevant to individual documents and relatively rare across the entire corpus.

The solution should provide a means to preprocess the text data by removing stop words, punctuation, and performing tokenization. It should then calculate the TF-IDF scores for each word in the corpus, using appropriate weighting schemes to ensure the significance of each term is accurately represented.

# Block Diagram of TF-IDF weight Calculation

```
                    ┌─────────┐
                    │  Start  │
                    └────┬────┘
                         │
                         ▼
                  ╱─────────────╲
                 ╱   Stemming    ╲
                 ╲    Result      ╱
                  ╲─────┬────────╱
          ┌────────────┼────────────┐
          ▼            ▼            ▼
 ┌──────────────┐ ┌──────────────┐ ┌──────────────┐
 │ calculate    │ │ calculate    │ │ count the    │
 │ term         │ │ term         │ │ number of    │
 │ frequency(tf)│ │ document     │ │ frequencies  │
 │ document     │ │ frequency(df)│ │              │
 └──────┬───────┘ └──────────────┘ └──────────────┘
        │
        │                    ┌──────────────┐
        │                    │ count Invers │
        │◄───────────────────│ Document     │
        │                    │ frequency(tf)│
        ▼                    └──────────────┘
 ┌──────────────┐
 │ calculate    │
 │ tf-idf       │
 │ Wdt=TFdt x   │
 │ IDFt         │
 └──────┬───────┘
        ▼
 ╱──────────────╲      ┌─────────┐
 │ save          │────►│ Finish  │
 │ calculation   │     └─────────┘
 │ result        │
 ╲──────────────╱
```

# Steps in pre-processing of data

# Resources and Technology used

## Language and Libraries Used

Our project will be a wordcloud and it should be responsive for each accessible device. Hence, we are using these Technologies for making the project easier.

- ✓ **Python:** Python is a widely used programming language for NLP tasks due to its simplicity and availability of numerous libraries.
- ✓ **NLTK (Natural Language Toolkit):** NLTK is a popular Python library for natural language processing. It provides various tools and resources for tokenization, stemming, lemmatization, and other NLP tasks.
- ✓ **scikit-learn:** scikit-learn is a powerful machine learning library in Python. It includes efficient implementations of the TF-IDF vectorization algorithm.
- ✓ **WordCloud:** WordCloud is a Python library that allows you to create visually appealing word clouds from text data. It provides various customization options such as color palettes, mask shapes, and font styles.
- ✓ **Matplotlib:** Matplotlib is a widely used plotting library in Python. It can be used to visualize the generated word cloud.

## Editor/Software Used

- ✓ VS code
- ✓ Jupyter Web

# Basic modules of the project

Our project is divided into four modules:

- ✓ NLTK
- ✓ scikit-learn
- ✓ WordCloud
- ✓ Matplotlib

- **nltk:** The Natural Language Toolkit (NLTK) is a widely used library for natural language processing tasks. It provides various modules for text preprocessing, tokenization, stemming, lemmatization, and more.
- **sklearn:** The scikit-learn library is a powerful machine learning library that includes a Tf idf Vectorizer class. This class is used to convert text data into TF-IDF feature vectors, which are essential for computing TF-IDF scores.
- **wordcloud:** The wordcloud module is used to create visually appealing word clouds from text data. It provides the WordCloud class, which allows you to customize the appearance of the word cloud.
- **matplotlib:** Matplotlib is a popular plotting library in Python. It is used to visualize the generated word cloud.

# CODE AND OUTPUT



```
import matplotlib.pyplot as plt
import pandas as pd
import re
from sklearn.feature_extraction.text import TfidfVectorizer
from nltk import download
# download('stopwords')
# download('wordnet')
from nltk.corpus import stopwords
from wordcloud import WordCloud
import json
from nltk.stem import WordNetLemmatizer
```

```
# task 1 - loading
data = pd.read_json('recipes.json', lines=True)
data.head(5)
data.info()
data['Ingredients'].head(5)
```



```
# task 2 - cleaning
data['Ingredients']=data['Ingredients'].map(lambda x: re.sub(r'[^a-zA-Z]',' ',
str(x)))
data = data.dropna(subset=['Ingredients'])
stop = stopwords.words('english') + ['tsp','tbsp','finely','extra','chopped']
stop
```

```python
[9]:  # task 3 - remove encoding
      def remove_encoding_word(word):
          word=str(word)
          word=word.encode('ASCII','ignore').decode('ASCII')
          return word

[10]: def remove_encoding_text(text):
          text = str(text)
          text = ' '.join(remove_encoding_word(word) for word in text.split() if word not in stop)
          return text

[11]: # task 4 - define lemmatizing
      data['Ingredients']=data['Ingredients'].apply(remove_encoding_text)

[12]: text=" ".join(words for words in data['Ingredients'])
      len(text)

[12]: 334272

[13]: lemma= WordNetLemmatizer().lemmatize

[14]: lemma('leaves')

[14]: 'leaf'

[15]: # task 5 - fit and transform text, with or without lemmatizing
      def tokenize(document):
          tokens = [lemma(w) for w in document.split() if len(w)>3 and w.isalpha()]
          return tokens

[16]: vectorizer = TfidfVectorizer(tokenizer = tokenize, ngram_range = (2,2), stop_words = stop, strip_accents='unicode')

[17]: tdm=vectorizer.fit_transform(data['Ingredients'])
```

```
C:\Users\LalitDumka\AppData\Local\Programs\Python\Python311\Lib\site-packages\sklearn\feature_extraction\text.py:528: UserWarning: The parameter 'token_pattern' will not be used since 'tokenizer' is not None'
  warnings.warn(
C:\Users\LalitDumka\AppData\Local\Programs\Python\Python311\Lib\site-packages\sklearn\feature_extraction\text.py:409: UserWarning: Your stop_words may be inconsistent with your preprocessing. Tokenizing the stop words generated tokens ['doe'] not in stop_words.
  warnings.warn(
```

```python
[18]: vectorizer.vocabulary_.items()

[18]: dict_items([('olive knob', 9699), ('knob butter', 7502), ('butter onion', 1627), ('onion sausagemeat', 9857), ('sausagemeat skinned', 13214), ('skinned sausage', 13992), ('sausage grated', 13191), ('grated zest', 6345), ('zest lemon', 17881), ('lemon fresh', 8023), ('fresh white', 5613), ('white breadcrumb', 17437), ('breadcrumb ready', 1295), ('ready dried', 12199), ('dried apricot', 4588), ('apricot chestnut', 290), ('chestnut canned', 2277), ('canned vacuum', 1859), ('vacuum packed', 16845), ('packed fresh', 10288), ('fresh dried', 5558), ('dried thyme', 4655), ('thyme cranberry', 16185), ('cranberry fresh', 3531), ('fresh frozen', 5563), ('frozen boneless', 5651), ('boneless skinless', 1056), ('skinless chicken', 13970), ('chicken breast', 2317), ('breast pack', 1336), ('pack ready', 10250), ('ready made', 12209), ('made shortcrust', 8524), ('shortcrust pastry', 13744), ('pastry beaten', 10590), ('beaten glaze', 689), ('butter dark', 1572), ('dark muscovado', 4150), ('muscovado sugar', 9299), ('sugar ...
```

```python
# task 3 - remove encoding
def remove_encoding_word(word):
    word=str(word)
    word=word.encode('ASCII','ignore').decode('ASCII')
    return word


def remove_encoding_text(text):
    text = str(text)
    text = ' '.join(remove_encoding_word(word) for word in text.split() if word not in stop)
    return text
```

```python
# task 4 - define lemmatizing
data['Ingredients']=data['Ingredients'].apply(remove_encoding_text)
text=" ".join(words for words in data['Ingredients'])
len(text)
lemma= WordNetLemmatizer().lemmatize
lemma('leaves')
```

```python
# task 5 - fit and transform text, with or without lemmatizing
def tokenize(document):
    tokens = [lemma(w) for w in document.split() if len(w)>3 and w.isalpha()]
```

```python
    return tokens
vectorizer = TfidfVectorizer(tokenizer = tokenize, ngram_range = (2,2),
stop_words = stop, strip_accents='unicode')
tdm=vectorizer.fit_transform(data['Ingredients'])
vectorizer.vocabulary_.items()
```



```python
# task 6 -  get word frequencies and create wordcloud
tfidf_weights=[(word,tdm.getcol(idx).sum()) for word, idx in
vectorizer.vocabulary_.items()]
tfidf_weights=dict(tfidf_weights)
w=WordCloud(width=1500,height=1200, mode='RGBA', background_color='white',
max_words=2000).fit_words(dict(tfidf_weights))
plt.figure(figsize=(20,15))
plt.imshow(w)
plt.axis('off')
plt.savefig('recipes_wordcloud.png')
txtFile = open("output.txt","w")
sortedWeights = sorted(tfidf_weights.items(), key=lambda x:x[1],reverse=True)
txtFile.write("Ingredient        :            Weight \n")
for key,values in sortedWeights:
    txtFile.write(str(key)+" : " +str(values) + "\n")
txtFile.close()
```

# Word Cloud and Text File Output



## Ingredient : Weight

caster sugar : 33.92679359870039
plain flour : 26.1646998160027802
double cream : 24.47528017761138
golden caster : 23.665605324844833
icing sugar : 23.56736189792369
garlic clove : 20.5616648091688
muscovado sugar : 15.157442191497505
butter softened : 14.33284930512973
thinly sliced : 13.434276096688432
unsalted butter : 13.02225291911623
dark chocolate : 12.759546623175648
wine vinegar : 12.328195823929985
ground cinnamon : 11.836742570209484
large egg : 11.655002686716756
zest juice : 11.596074513899923
white wine : 11.50957441262325

plus dusting : 11.200198008258184
vanilla extract : 11.146409447587125
brown sugar : 11.135413088674557
light muscovado : 10.955134306343428
self raising : 10.824187358979268
raising flour : 10.782791060584968
plus greasing : 10.619582545840558
chicken stock : 10.54995369400423
clove crushed : 10.449428841281208
grated zest : 10.238750697266983
mixed spice : 9.882048828461663
cinnamon stick : 9.826696050610387
flour plus : 9.761752979175906
baking powder : 9.575576364877861
juice lemon : 9.218973352201182
brussels sprout : 8.760995074060489

egg beaten : 8.733133620281524
streaky bacon : 8.561729322633997
ground almond : 8.527966748756377
clear honey : 8.485596404076897
good quality : 8.45323965947161
sugar plus : 8.389880240431221
butter plus : 8.36841543181027
onion sliced : 8.300531309453822
thyme leaf : 8.263535749988357
zest orange : 8.260474585826856
zest lemon : 8.076860974104306
ground ginger : 7.923132778011525
dried cranberry : 7.616847676526653
maple syrup : 7.6098421312007485
lemon juice : 7.543426313041776
large onion : 7.530992783461033
golden syrup : 7.528255046680915
puff pastry : 7.3514946689895995
wholegrain mustard : 7.346843934074862
small bunch : 7.334888420314011
star anise : 7.186835264657585
fresh frozen : 7.16753473144454
vegetable stock : 7.055420180052844
plus serve : 6.950395768563531
butter melted : 6.935057863546623
sugar dusting : 6.877052007208772
freshly grated : 6.803062087205156
sugar large : 6.6841164198277
flaked almond : 6.615035320411994
soft brown : 6.589275376256462
apple peeled : 6.546330918694762
white chocolate : 6.537824204026434
knob butter : 6.475600035271045
thyme sprig : 6.473124535072162
small handful : 6.448344956228847
orange juice : 6.370504009714113
frozen cranberry : 6.356445501846716
peeled cored : 6.231054584465458
small pack : 6.134358043381125
grated nutmeg : 6.125389597761781
light brown : 6.070506759090989
cocoa powder : 6.067324945597459
bicarbonate soda : 6.054598831265656
mari piper : 5.970117673900321
dijon mustard : 5.943162595544985
juice orange : 5.929455258515979

flour dusting : 5.909574844108866
dried apricot : 5.908431948446099
dark muscovado : 5.889335563862861
plus little : 5.853880818141382
flour baking : 5.8113100082383
serve optional : 5.773989544631382
leaf picked : 5.653649277862263
small onion : 5.599796689010847
white bread : 5.4801556178746225
cider vinegar : 5.450944692684173
spring onion : 5.395562647257052
juice lime : 5.316331217423093
black peppercorn : 5.28922224312433
flour ground : 5.266958383465993
glac cherry : 5.264095966063564
sage leaf : 5.212046569688317
ground clove : 5.2104242762605635
pinch ground : 5.183622411466276
go well : 5.163376835303158
bramley apple : 5.134819390644449
potato peeled : 5.085394888209205
virgin olive : 5.08369970102348
smoked salmon : 5.03575046896877
whole milk : 4.993033020220411
small piece : 4.9839941319157495
salted butter : 4.943177210530218
ready roll : 4.934405070468728
sugar vanilla : 4.907390618514263
handful parsley : 4.9070226440403495
onion halved : 4.902345455249703
stem ginger : 4.893432850989686
white breadcrumb : 4.856425990554216
milk chocolate : 4.855886269249453
dried fruit : 4.820237069127879
chocolate broken : 4.776806296688027
food colouring : 4.760692389939962
ready rolled : 4.667419817096415
granulated sugar : 4.639177024147694
room temperature : 4.6368899141339295
pack ready : 4.543476184561189
softened butter : 4.539718513473184
sprout trimmed : 4.533009874572549
balsamic vinegar : 4.510122051580977
mixed peel : 4.497783151079771
smoked streaky : 4.490831354827359
Check the output.txt for more…

# <u>BENEFITS</u>

- **Data Visualization:** Word clouds provide an intuitive and visually appealing representation of textual data. They help in quickly identifying the most frequent and important terms in a document or corpus.

- **Text Analysis:** By utilizing NLP techniques, such as tokenization, TF-IDF computation, and word cloud generation, the project allows you to gain insights into the underlying text data. It helps in identifying key themes, significant terms, or prominent features within the text.

- **Communication and Presentation:** Word clouds are often used for effective communication and presentation purposes. They condense large amounts of text into a visually engaging format, making it easier to convey information to an audience.

- **Feature Extraction:** The TF-IDF values computed during the project can serve as feature vectors for further analysis or machine learning tasks. These vectors can be used to represent documents or text snippets, enabling various text classification or clustering applications.

## **CONCLUSION**

The project of creating a word cloud using NLP and TF-IDF in Python allows you to analyze and visualize text data in a meaningful way. By leveraging NLP techniques, such as tokenization and TF-IDF computation, you can identify the most important terms within the text and generate a visually appealing word cloud. The project offers benefits in terms of data visualization, text analysis, communication, and feature extraction. It can be a valuable tool for gaining insights from textual data and presenting information in a concise and visually appealing manner.