

```
## Lower Casing
zomato["reviews_list"] = zomato["reviews_list"].str.lower()
```

```
## Removal of Punctuations
import string
PUNCT_TO_REMOVE = string.punctuation
def remove_punctuation(text):
    """custom function to remove the punctuation"""
    return text.translate(str.maketrans('', '', PUNCT_TO_REMOVE))
```

```
zomato["reviews_list"] = zomato["reviews_list"].apply(lambda text: remove_punctuation(text))
```

```
import nltk
nltk.download('stopwords')
```

```
## Removal of Stopwords
from nltk.corpus import stopwords
STOPWORDS = set(stopwords.words('english'))
def remove_stopwords(text):
    """custom function to remove the stopwords"""
    return " ".join([word for word in str(text).split() if word not in STOPWORDS])
```

```
zomato["reviews_list"] = zomato["reviews_list"].apply(lambda text: remove_stopwords(text))
```

```
## Removal of URLs
def remove_urls(text):
    url_pattern = re.compile(r'https?://\S+|www\.\S+')
    return url_pattern.sub(r'', text)
```

```
zomato["reviews_list"] = zomato["reviews_list"].apply(lambda text: remove_urls(text))
```

```
zomato[['reviews_list', 'cuisines']].sample(5)
```

TF-IDF Vectorization

TF-IDF (Term Frequency-Inverse Document Frequency) vectors for each document. This will give you a matrix where each column represents a word in the general vocabulary (all words that appear in at least one document) and each column represents a restaurant, as before.

TF-IDF is the statistical method of assessing the meaning of a word in a given document.

```
df_percent.set_index('name',  
inplace=True)
```

```
indices =  
pd.Series(df_percent.index)
```

```
# Creating tf-idf matrix  
tfidf =  
TfidfVectorizer(analyzer='word',  
ngram_range=(1, 2), min_df=0,  
stop_words='english')
```

```
tfidf_matrix =  
tfidf.fit_transform(df_percent['reviews_list'])
```

```
cosine_similarities =  
linear_kernel(tfidf_matrix,  
tfidf_matrix)
```