

ETL

- Steps) ① Drop duplicates & NA values -
- code: df.drop_duplicates(subset=[‘Text’], inplace=True)
df.dropna(axis=0, inplace=True)
- ② Preprocessing: drop all unwanted symbols, characters etc. from the text that do not effect the objective of problem.

① Contraction Mapping
Contraction_mapping = {“ain’t”: “is not”, “aren’t”: “are not”}

② Converting everything to lower case
Remove Stopwords

stop_words = set(stopwords.words(‘english’)) + + - - -

Cleaned_text = [“bought several vitality” -- product better, --]

Cleaned_summary = [“good quality food”, -- --]

Data[‘cleaned_text’] = Cleaned_text, Data[‘clean sum’] = Clean summary

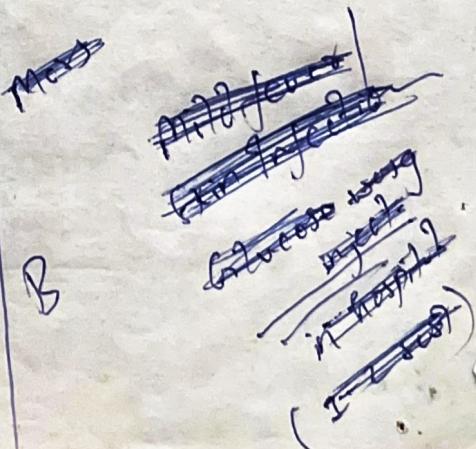
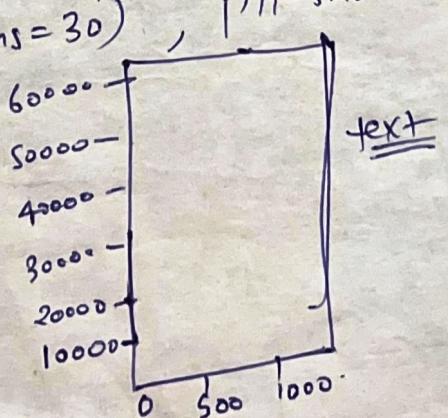
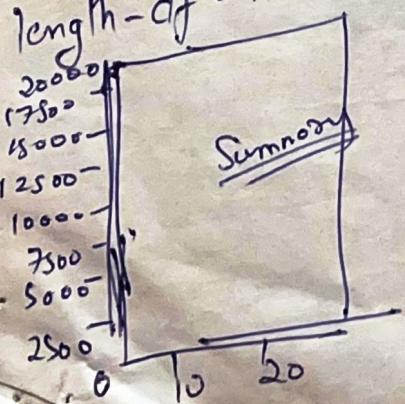
Visualizing Distribution
import matplotlib.pyplot as plt

text_word_count = [], Summary_word_count = []

for i in Data[‘cleaned_text’]:
text_word_count.append(len(i.split()))

length_df = pd.DataFrame({‘text’: text_word_count, ‘Summary’: Summary_word_count})

length_df.hist(bins=30), plt.show()



EDA

① Business Understanding.

② Data Understanding / Acquiring.

③ Identifying & handling the missing values.

④ Encoding the categorical data.

⑤ Checking the correlation b/w features.

⑥ feature importance (if features are many)

⑦ Target variable is balanced / imbalanced.

⑧ Splitting the dataset into train & test.

⑨ Feature Scaling.

If Model is not performing well on test data:

| Data Profiling (Test data).

Steps: ① Outliers check.

② Will check features / region where the values are abruptly changing (due to surround circumstances, pandemic..)

③ Will check/test feature statistics with train feature statistics

If All seems fine:
⇒ Regularization check.

⇒ Cross Validation if it is (powerful measure) due to overfitting.

⇒ Remove feature.
⇒ Ensembling