# SmartMailGuard: AI-Powered Email Classification

Rupak R. Gupta (RRG)

B. Tech Information Technology (IT)

June 2024

**Abstract**

In today's age, our inboxes are constantly being flooded with emails from all kinds of sources, which makes finding a legitimate email in a sea of spam challenging. The aim of this project is to develop an Artificial Intelligence (AI) powered tool that can parse a given email and classify whether it is legitimate or spam.
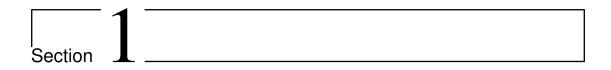
**Acknowledgements**

# Contents

**Acronyms**            **22**

**References**          **23**

# Preface

## About me

I am Rupak R. Gupta, though most people know me as RRG. I am a (soon-to-be) second year student at VJTI in IT. I have a deep interest in statistics, probability, and mathematics in general, and I am very curious about their applications in the field of computer science for problem-solving.

I have had a strong grasp on the programming languages Java , C++ and Python ; I have practised some standard Data Structures and Algorithms using each of these languages.

I have also acquired skills in writing LaTeX and creating digital diagrams, some of which are included in this proposal. I wish to utilize these skills as an educator/explainer myself to contribute to the education of others.

**Tasks repository** : `https://github.com/aitwehrrg/Project-X/`

## Motivation for this topic

As mentioned earlier, I have a lot of interest in topics such as probability and statistics, which encouraged me to select a project based on Machine Learning (ML). Another reason is that I am very much a patron of the best practices in cybersecurity, such as password management and email aliasing. Avoiding spam and phishing emails is yet another important cybersecurity practice, so this topic felt coherent to me.

# 2 Section

## Theory and Approach

### 2.1 Preprocessing

We will convert the text data from the emails into a format that can be fed into the neural network. This includes:

1. Cleaning up unwanted words and characters (like punctuation).

2. Converting all text to lowercase.

3. Extracting individual *tokens* from the data.

### 2.2 Neural Networks



Figure 2.1: Neural Network

A prototype Recurrent Neural Network (RNN) is depicted in fig. 2.1.

- **Input or embedding layer:** This layer has as many nodes as there are tokens. Each token input is vectorized using techniques such as `Word2Vec` and fed into the network.

- **LSTM layer:** Long Short-Term Memory (LSTM) [1] is a type of RNN architecture which is designed to remember information for long periods of time, allowing it to learn semantics of a some text using the context of words (tokens) present long back in the email.

- **Hidden layers:** These layers perform the computation to generate the required probabilities. The algorithm used is discussed in § 2.4.

- **Output layer:** This layer has two neurons that hold their respective final probability of the given email being not spam and the given email being spam. An alternative idea would be to use a *sigmoid function* [2] to output just one binary output; the email being spam or not.

## 2.3  Transformers

The architecture of a transformer is depicted in fig. 2.2.



Figure 2.2: Transformer Architecture [1]

- We shall select a pretrained transformer model to use such as Bidirectional Encoder Representations from Transformers (BERT), or Generative Pretrained Transformer (GPT) like GPT-2 [1].

- Using attention mechanism [3], understand semantics of text in the email to properly encode the tokens into vectors.

- The transformer is integrated with the LSTM described in § 2.2 using methods such as Concatenation, Stacking *etc.*

## 2.4 Naïve Bayes Algorithm

### 2.4.1 Bayes' Theorem

We will obtain a sample of emails that we can use for training. Our training data would be used to determine the *prior probability* [4] of an email being spam $P(\text{Spam})$.

After preprocessing, we can determine the *posterior probability* [5] for each token by using Bayes' theorem.

$$P(\text{Token} \,|\, \text{Spam}) := \frac{P(\text{Spam} \,|\, \text{Token}) P(\text{Token})}{P(\text{Spam})} \tag{2.1}$$

$P(\text{Spam} \,|\, \text{Token})$ and $P(\text{Token})$ are probabilities that can be estimated from the training data.

### 2.4.2 Naïve Bayes

Using Naïve Bayes, we can determine the total probability that an email is spam by multiplying the posterior probabilities of each token [4].

$$
\begin{aligned}
P(\text{Spam} \,|\, \text{Tokens}) &= \frac{P(\text{Spam}) \times \prod P(\text{Token} \,|\, \text{Spam})}{\prod P(\text{Token})} \\
&\propto \boxed{P(\text{Spam}) \times \prod_{\text{all Tokens}} P(\text{Token} \,|\, \text{Spam})}
\end{aligned}
\tag{2.2}
$$

We can ignore the prior probability $\prod P(\text{Token})$ since it is independent of the legitimacy of the email. Therefore, we only care about the relative value of the expression in eq. (2.2).

### 2.4.3 Smoothing

To account for exactly zero posterior probabilities of certain tokens in one of the categories, we will use *Laplace smoothing* [6]: $\alpha = 1$ to add one point of frequency to every token (figs. 2.3 and 2.4). Another option is using *additive smoothing* for other values of $\alpha$.



Figure 2.3: Before smoothing

Figure 2.4: After smoothing ($\alpha = 1$)

# Section 3

# Workflow

## 3.1   Week 1

Gather a large enough sample for training and testing the model. Study the sample and determine the techniques to preprocess the data in order for it to be fed to the neural network.

## 3.2   Week 2

Preprocess the data: clean the text to remove unwanted information, split into individual tokens. Learn the implementations of RNN (LSTM) and the Naïve Bayes Algorithm.

## 3.3   Weeks 3 and 4

Implement the neural network and and train it on the training data. Study the transformer model to use such as BERT, GPT or GPT-2, as well as the technique to integrate it to the network such as Concatenation or Stacking.

## 3.4   Week 5 onwards

Integrate the transformer model with the neural network and evaluate its performance on the testing sample. Grade its accuracy, fix bugs and tweak parameters to maximise accuracy and efficiency.

# Appendix A

# Task 1: Difficult dating

**Link:** https://github.com/aitwehrrg/Project-X

## A.1  Code

This task is completed as a Jupyter notebook.

**Importing packages**

```
[1]: import matplotlib.pyplot as plt
     import numpy as np
     import pandas as pd
     from sklearn.svm import SVC
     from sklearn.inspection import DecisionBoundaryDisplay
     from sklearn.metrics import accuracy_score
     from sklearn.model_selection import train_test_split
     from sklearn.preprocessing import StandardScaler, OneHotEncoder
```

**Defining constants**

```
[2]: TEST_SIZE: float = 0.2  # Train size would be 1 - TEST_SIZE
     TEST_CASE: str = 'dating.csv'
```

**Importing the `csv` as a `pandas DataFrame`**

```
[3]: df: pd.DataFrame = pd.read_csv(TEST_CASE)
     df.drop(df.columns[0], axis=1, inplace=True)  # Dropping the first
       ↪column (serial number)
```

## Preprocessing

### Preprocessing `career`

```
[4]: df['career'] = df['career'].fillna('undecided')  # Filling empty␣
       ↪cells with 'undecided'
     df['career'] = df['career'].str.lower()  # Turining everything␣
       ↪lowercase to prepare for encoding
```

### Preprocessing using One Hot Encoder

```
[5]: encoder = OneHotEncoder()

     # Encoding the string column
     career_encoded = encoder.fit_transform(df[['career']])

     # Convert back to DataFrame
     career_encoded_df = pd.DataFrame(career_encoded.toarray(),␣
       ↪columns=encoder.get_feature_names_out(['career']))

     # Combining with original DataFrame
     df = df.reset_index(drop=True)
     df = df.join(career_encoded_df)
     df.drop('career', axis=1, inplace=True)  # No more need for the␣
       ↪original string column
```

### Preprocessing the numeric columns (`age` and `income`)

```
[6]: df = df.fillna(df.mean())  # Replacing empty cells with mean instead
```

### Standard Scaling all numeric columns (Mean = 0, Standard Deviation = 1)

```
[7]: scaler = StandardScaler()
     # Attractiveness and fun are not scaled because we will need them␣
       ↪for plotting
     df[['age', 'income', 'sinc', 'intel', 'amb', 'like']] = scaler.
       ↪fit_transform(df[['age', 'income', 'sinc', 'intel', 'amb',␣
       ↪'like']])
```

## Splitting into training and testing data

```
[8]: X = df.drop('dec', axis=1)  # The SVM must not train on the output
     y = df['dec']  # We will test the output
     X_train, X_test, y_train, y_test = train_test_split(X, y,␣
       ↪test_size=TEST_SIZE, random_state=69)
```

**Training an SVM classifier**

```
[9]: svm = SVC()
     svm.fit(X_train, y_train)
```

```
[9]: SVC()
```

**Predicting the output of the test data and calculate the accuracy**

```
[10]: y_pred = svm.predict(X_test)
      accuracy: float = accuracy_score(y_test, y_pred)
      accuracy
```

```
[10]: 0.7975852272727273
```

## A.2   Output

```
[11]: # Setting up the plot
      svm.fit(X_train[['attr', 'fun']], y_train)  # Fit the plot based on␣
        ↪attractiveness and fun
      plt.scatter(X_train[y_train == 0]['attr'],
              X_train[y_train == 0]['fun'],
              c='r',
              label='Not Matched')
      plt.scatter(X_train[y_train == 1]['attr'],
              X_train[y_train == 1]['fun'],
              c='g',
              label='Matched')
      # Adding labels and title
      plt.xlabel('Attractiveness')
      plt.ylabel('Fun')
      plt.title('Accuracy Score = ' + str(round(accuracy, 4)))
      plt.legend()
      # Plotting the decision boundary
      DecisionBoundaryDisplay.from_estimator(
      svm,
      X_train[['attr', 'fun']],
      response_method='predict',
      alpha=0.5,
      ax=plt.gca()
      )
      # Final output
      plt.show()
```

Accuracy Score = 0.7976

# References

- *SVM* by StatQuest with Josh Starmer [7]

# Appendix B

# Task 2: Where my diamonds at?

**Link:** https://github.com/aitwehrrg/Project-X

## B.1 Code

This task is completed had a really long codebase, so it has been split into four Python files.

### Main

```python
import cv2
import numpy as np
import process_image as pi
import start_finder as sf
import vertices_finder as vf

BLACK: int = 0
WHITE: int = 255

TEST_CASES_1: tuple[str, str, str] = ('pics/tc1-1.png', 'pics/tc1-2.png',
    'pics/tc1-3.png')
TEST_CASES_2: tuple[str, str, str] = ('pics/tc2-1.png', 'pics/tc2-2.png',
    'pics/tc2-3.png')

SUIT_DICT: dict[int, str] = {
    0: 'Hearts',
    1: 'Clubs',
    2: 'Spades',
    3: 'Diamonds'
}

VALUE_DICT: dict[int, str] = {
    0: 'Ace',
    1: '2',
    2: '3',
    3: '4',
    4: '5',
    5: '6',
}
```

```python
30    # Traverse the top row
31    def horizontal_ratios(image: np.ndarray, target: int) -> list[float]:
32        start_x, start_y = sf.StartFinder(image, target).left()
33        ratios: list[float] = []
34        _, width = image.shape
35        for x in range(start_x, width):
36            if image[start_y][x] == target and image[start_y][x - 1] != target:
37                ratio: float = vf.VerticesFinder(image, target, (x, start_y),
                        'left').aspect_ratio()
38                if ratio != 0:
39                    ratios.append(ratio)
40        return ratios


43    # Traverse the left column
44    def vertical_ratios(image: np.ndarray, target: int, start_x: int = 0, start_y: int =
        0) -> list[float]:
45        ratios: list[float] = []
46        height, _ = image.shape
47        for y in range(start_y, height):
48            if image[y][start_x] == target and image[y - 1][start_x] != target:
49                ratio: float = vf.VerticesFinder(image, target, (start_x, y),
                        'top').aspect_ratio()
50                if ratio != 0:
51                    ratios.append(ratio)
52        return ratios


55    # Required function
56    def output(image: np.ndarray, target: int, suit: str | None = None, value: str | None
        = None) -> str:
57        if suit is None:
58            suit_ratios: list[float] = horizontal_ratios(image, BLACK)
59            suit_min: float = min(suit_ratios)
60            suit_max: float = max(suit_ratios)
61
62            # Check the minimum and maximum ratios of the rop row to determine the nature
                of the card
63            if suit_min == 1 and suit_max == 1:  # the card is upside down and not a 6
64                image = cv2.flip(image, -1)  # hold the card upright
65                return output(image, target, suit, value)  # try again
66
67            if suit_min < 1 and suit_max == 1:  # the card is upside down and a 6
68                value = VALUE_DICT[5]  # the card is a 6 (since it is upside down)
69                image = cv2.flip(image, -1)  # hold the card upright
70                return output(image, target, suit, value)  # try again
71
72            if suit_min < 1 < suit_max:  # the card is upright and an ace
73                suit = SUIT_DICT[suit_ratios.index(suit_max) - 1]  # minus 1 because
                        there are 5 diamonds but only 4 suits
74                value = VALUE_DICT[0]  # the card is an ace
75                return output(image, target, suit, value)
76
77            # the card is upright and not an ace
78            suit = SUIT_DICT[suit_ratios.index(suit_max) - 1]  # minus 1 because there
                are 5 diamonds but only 4 suits
79
80        # if the upright card is not an ace and the upside down card is not a 6
81        if value is None:
```

```
82      top_x, top_y = sf.StartFinder(image, target).top()
83      value_ratios = vertical_ratios(image, target, top_x, top_y)
84      value = VALUE_DICT[value_ratios.index(min(value_ratios))]
85
86   output_: str = f'{value} of {suit}'
87   print(output_)
88   cv2.imshow(output_, image)
89   cv2.waitKey(0)
90   return output_
91
92
93 def main() -> None:
94   # Select the test case from here
95   # test_cases: tuple[str, str, str] = TEST_CASES_1
96   # test_cases: tuple[str, str, str] = TEST_CASES_2
97
98   test_cases: tuple = TEST_CASES_1 + TEST_CASES_2
99   result: str = ''
100  for test_case in test_cases:
101      image = cv2.imread(test_case, cv2.IMREAD_GRAYSCALE)
102      image = pi.Process(image, BLACK, WHITE).process()
103      result += f'{output(image, BLACK)}, '
104  result += '\b\b'   # remove trailing comma and space; on some terminals it may not
         work
105  print(f'Final result: {result}')
106
107
108 if __name__ == '__main__':
109     main()
```

Listing B.1: main.py

## Process Image

```
1  import cv2
2  import numpy as np
3
4  BLACK: int = 0
5  GRAY: int = 28
6  WHITE: int = 255
7
8
9  class Process:
10     def __init__(self, image: np.ndarray, border_value: int, bg_value: int,
          tolerance: int = 2) -> None:
11         self._image = image
12         self._border_value = border_value
13         self._bg_value = bg_value
14         self._TOLERANCE = tolerance   # 2 is the tightest tolerance
15
16     def _remove_border(self) -> None:
17         height, width = self._image.shape
18         for y in range(height):
19             for x in range(width):
20                 if self._image[y][x] == self._border_value:
21                     self._image[y][x] = self._bg_value
22         # cv2.imshow('Border removed', self._image)
23         # cv2.waitKey(0)
```

```
24
25      # Crops the image to remove the remaining pixels that are not the background
26      def _crop(self) -> None:
27          height, width = self._image.shape
28          # Left crop
29          for x in range(width):
30              if self._image[0][x] != self._bg_value:
31                  self._image = self._image[:, x + 1:]
32                  _, width = self._image.shape
33                  break
34          # Right crop
35          for x in range(width - 1, -1, -1):
36              if self._image[height - 1][x] != self._bg_value:
37                  self._image = self._image[:, :x - 1]
38                  _, width = self._image.shape
39                  break
40          # Top crop
41          for y in range(height):
42              if self._image[y][0] != self._bg_value:
43                  self._image = self._image[y + 1:, :]
44                  height, _ = self._image.shape
45                  break
46          # Bottom crop
47          for y in range(height - 1, -1, -1):
48              if self._image[y][width - 1] == self._bg_value:
49                  self._image = self._image[:y - 1, :]
50                  break
51          # cv2.imshow('Cropped', self._image)
52          # cv2.waitKey(0)
53
54      def process(self) -> np.ndarray:
55          self._remove_border()
56
57          # Remove remnants of the border (salt and pepper noise)
58          self._image = cv2.medianBlur(self._image, 3)
59          # cv2.imshow('Median blur', self._image)
60          # cv2.waitKey(0)
61
62          self._crop()
63
64          # Remove any pixels that are not the required color
65          self._image[self._image > GRAY + self._TOLERANCE] = WHITE
66          self._image[self._image < GRAY - self._TOLERANCE] = WHITE
67          # cv2.imshow('Pixel color filter', self._image)
68          # cv2.waitKey(0)
69
70          # Convert the image to binary
71          _, self._image = cv2.threshold(self._image, 254, 255, cv2.THRESH_BINARY)
72          return self._image
```

Listing B.2: process_image.py

## Start Finder

```python
import numpy as np


class StartFinder:
```

```
 5    def __init__(self, image: np.ndarray, target: int) -> None:
 6        self._image: np.ndarray = image
 7        self._height, self._width = image.shape
 8        self._target: int = target
 9
10    def top(self) -> tuple[int, int]:
11        for y in range(self._height):
12            for x in range(self._width):
13                if self._image[y][x] == self._target:
14                    return x, y
15        raise ValueError('Target not found')
16
17    def left(self) -> tuple[int, int]:
18        for x in range(self._width):
19            for y in range(self._height):
20                if self._image[y][x] == self._target:
21                    return x, y
22        raise ValueError('Target not found')
```

Listing B.3: start_finder.py

## Vertices Finder

```
 1  import numpy as np
 2
 3  # To account for pixel inaccuracies
 4  MIN_RATIO: float = 0.96
 5  MAX_RATIO: float = 1.04
 6
 7
 8  class VerticesFinder:
 9      def __init__(self, image: np.ndarray, target: int, initial_vertex: tuple[int,
           int], direction: str) -> None:
10          self._image: np.ndarray = image
11          self._height, self._width = image.shape
12          self._target: int = target
13          self._initial_vertex = initial_vertex
14          self._start_x, self._start_y = initial_vertex
15          self._direction: str = direction
16
17      def _top(self) -> tuple[tuple[int, int], tuple[int, int], tuple[int, int]]:
18          # Find bottom vertex
19          y: int = self._start_y
20          for y in range(self._start_y, self._height):
21              if self._image[y, self._start_x] != self._target:
22                  break
23          bottom_vertex: tuple[int, int] = self._start_x, y - 1
24
25          mid_y: int = (self._start_y + bottom_vertex[1]) // 2
26
27          # Find right vertex
28          x = self._start_x
29          for x in range(self._start_x, self._width):
30              if self._image[mid_y, x] != self._target:
31                  break
32          right_vertex: tuple[int, int] = x - 1, mid_y
33
34          # Find left vertex
```

```python
        left_vertex: tuple[int, int] = 2 * self._start_x - right_vertex[0], mid_y

        return bottom_vertex, left_vertex, right_vertex

    def _left(self) -> tuple[tuple[int, int], tuple[int, int], tuple[int, int]]:
        # Find right vertex
        x: int = self._start_x
        for x in range(self._start_x, self._width):
            if self._image[self._start_y, x] != self._target:
                break
        right_vertex: tuple[int, int] = x - 1, self._start_y

        mid_x: int = (self._start_x + right_vertex[0]) // 2

        # Find bottom vertex
        y = self._start_y
        for y in range(self._start_y, self._height):
            if self._image[y, mid_x] != self._target:
                break
        bottom_vertex: tuple[int, int] = mid_x, y - 1

        # Find top vertex
        top_vertex: tuple[int, int] = mid_x, 2 * self._start_y - bottom_vertex[1]

        return top_vertex, bottom_vertex, right_vertex

    def vertices(self) -> tuple[tuple[int, int], tuple[int, int], tuple[int, int],
        tuple[int, int]]:
        # top, bottom, left, right
        match self._direction:
            case 'top':
                bottom_vertex, left_vertex, right_vertex = self._top()
                return self._initial_vertex, bottom_vertex, left_vertex, right_vertex

            case 'left':
                top_vertex, bottom_vertex, right_vertex = self._left()
                return top_vertex, bottom_vertex, self._initial_vertex, right_vertex

            case _:
                raise ValueError(f'Invalid direction: {self._direction}')

    def aspect_ratio(self) -> float:
        top_vertex, bottom_vertex, left_vertex, right_vertex = self.vertices()
        width = right_vertex[0] - left_vertex[0]
        height = bottom_vertex[1] - top_vertex[1]
        try:
            ratio: float = round(width / height, 2)
        except ZeroDivisionError:
            return 0

        # To account for pixel inaccuracies
        if MIN_RATIO < ratio < MAX_RATIO:
            return 1

        return ratio
```
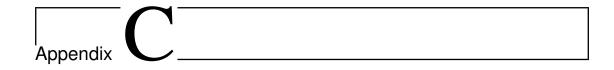
Listing B.4: vertices_finder.py

## B.2 Output

```
2 of Spades
3 of Clubs
4 of Hearts
3 of Spades
6 of Hearts
6 of Clubs
Final result: 2 of Spades, 3 of Clubs, 4 of Hearts, 3 of Spades, 6 of Hearts, 6 of Clubs

Process finished with exit code 0
```

## References

- Computer Vision (CV) Xplore Workshop materials [8]

# Task 3: This Bleu me away

**Link:** https://github.com/aitwehrrg/Project-X

## C.1 Code

This task is completed using one Python file.

```python
from collections import Counter
import math

TEST_CASE_1: tuple[list[str], str] = (['It is a guide to action that ensures that the
    military will forever heed Party '
                                      'commands',
                                      'It is the guiding principle which guarantees
                                          the military forces always being '
                                      'under the command of the Party',
                                      'It is the practical guide for the army always
                                          to heed the directions of the '
                                      'party'],
                                      'It is a guide to action which ensures that the
                                          military always obeys the '
                                      'commands of the party')

TEST_CASE_2: tuple[list[str], str] = (['It is a guide to action that ensures that the
    military will forever heed Party '
                                      'commands',
                                      'It is the guiding principle which guarantees
                                          the military forces always being '
                                      'under the command of the Party',
                                      'It is the practical guide for the army always
                                          to heed the directions of the '
                                      'party'],
                                      'It is the to action the troops forever hearing
                                          the activity guidebook that '
                                      'party direct')

# Maximum number of n-grams
N: int = 4


# Compute n-grams
```

```python
27  def n_grams(sentence: str, n: int = 1) -> list[str]:
28      words: list[str] = sentence.split()
29      if n == 1:
30          return words
31      n_grams_: list[str] = []
32      for i in range(len(words) - n + 1):
33          n_grams_.append(' '.join(words[i:i + n]))
34      return n_grams_
35
36
37  # Count frequency of n-grams
38  def count(n_gram: list[str]) -> int:
39      return sum(Counter(n_gram).values())
40
41
42  # Count frequency of clipped n-grams
43  def count_clip(candidate: str, references: list[str], n: int) -> int:
44      n_gram_candidate: list[str] = n_grams(candidate, n)
45      n_gram_references: list[list[str]] = [n_grams(reference, n) for reference in
            references]
46      counts: Counter = Counter(n_gram_candidate)
47      for key in counts:
48          max_reference_count: int = 0
49          for n_gram_reference in n_gram_references:
50              # n_gram_reference_count: Counter = Counter(n_gram_reference)
51              max_reference_count = max(max_reference_count,
                    n_gram_reference.count(key))
52          counts[key] = min(counts[key], max_reference_count)
53      return sum(counts.values())
54
55
56  # Calculate p_n
57  def precision(candidate: str, references: list[str], n: int) -> float:
58      count_: int = count(n_grams(candidate, n))
59      count_clip_: int = count_clip(candidate, references, n)
60      # print(n, f'{count_clip_}/{count_}')
61      return count_clip_ / count_ if count_ != 0 else 0
62
63
64  # Calculate BP
65  def brevity_penalty(candidate: str, references: list[str]) -> float:
66      c: int = len(n_grams(candidate))  # total length of candidate
67      word_count_references: list[int] = [len(n_grams(reference)) for reference in
            references]
68      deviations: list[int] = [abs(c - reference) for reference in
            word_count_references]
69      r: int = word_count_references[deviations.index(min(deviations))]  # best match
            reference length
70      if c > r:
71          return 1
72      return math.exp(1 - r / c)
73
74
75  # Required function
76  def bleu_score(references: list[str], candidate: str) -> float:
77      bleu: float = 1
78      bp: float = brevity_penalty(candidate, references)
79      print(f'BP: {bp}')
80      for n in range(1, N + 1):
```

```
81        p_n: float = precision(candidate, references, n)
82        bleu *= p_n
83    bleu **= 1 / N
84    bleu *= bp
85    return bleu
86
87
88 def main():
89    test_cases: tuple = (TEST_CASE_1, TEST_CASE_2)
90    for test_case in test_cases:
91        bleu: float = bleu_score(test_case[0], test_case[1])
92        print(f'BLEU score: {bleu}\n')
93
94
95 if __name__ == '__main__':
96    main()
```

Listing C.1: main.py

## C.2 Output

```
BP: 1.0
BLEU score: 0.5045666840058485

BP: 0.9355069850316178
BLEU score: 0.0



Process finished with exit code 0
```

## References

- "*Bleu: a Method for Automatic Evaluation of Machine Translation*" [9]
- *C5W3L06 Bleu Score* by DeepLearningAI [10]

# Acronyms

# References

[1] Michael Phi. Illustrated guide to transformers- step by step explanation. `https://towardsdatascience.com/illustrated-guide-to-transformers-step-by-step-explanation-f74876522bc0`.

[2] Neural networks. `http://www.youtube.com/playlist?list=PLZHQObOWTQDNU6R1_67000Dx_ZCJB-3pi`.

[3] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. `http://arxiv.org/abs/1706.03762`.

[4] StatQuest with Josh Starmer. Naive bayes, clearly explained!!! `https://www.youtube.com/watch?v=O2L2Uv9pdDA`.

[5] Petuum Inc. Intro to modern bayesian learning and probabilistic programming. `https://petuum.medium.com/intro-to-modern-bayesian-learning-and-probabilistic-programming-c61830df5c50`.

[6] Scikit-learn. 1.9. Naive Bayes. `https://scikit-learn/stable/modules/naive_bayes.html`.

[7] StatQuest with Josh Starmer. Support vector machines part 1 (of 3): Main ideas!!! `https://www.youtube.com/watch?v=efR1C6CvhmE`.

[8] Computer vision Xplore workshop - google drive. `https://drive.google.com/drive/folders/1R2WOCz8mDvmbE0qSEnq51uNh8EZigTU9`.

[9] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. `https://aclanthology.org/P02-1040`.

[10] DeepLearningAI. C5w3l06 bleu score (optional). `https://www.youtube.com/watch?v=DejHQYAGb7Q`.

$$\begin{bmatrix} \mathbb{R} & \mathfrak{R} \\ g & \infty \end{bmatrix}$$

The End.