Assignment-1

Q1: Given p=[('noun', 0.441),('verb', 0.255),('adj', 0.132),('adv', 0.172)], please use python code to calculate information, entropy, cross-entropy, and KL divergence.

A1. Python code to calculate information, entropy, cross-entropy, and KL divergence.

Information Theory Calculations

```
[] import math
[] p = [('noun', 0.441),('verb', 0.255),('adj', 0.132),('adv', 0.172)]
    # Calculate the information content of each event
    information = {event: round(-math.log2(prob), 3) for event, prob in p}
[ ] # Calculate the entropy of the distribution
    entropy = sum(-prob * math.log2(prob) for _, prob in p)
▶ # Define the target probability distribution q as a uniform distribution
    total_categories = len(p)
    uniform_prob = 1 / total_categories
    q_uniform = [(label, uniform_prob) for label, _ in p]
    print("Uniform Distribution:", q_uniform)
    # Calculate cross-entropy
    cross\_entropy = -sum(prob\_p * math.log2(prob\_q) \ for \ (label\_p, prob\_p), \ (\_, prob\_q) \ in \ zip(p, \ q\_uniform))
Uniform Distribution: [('noun', 0.25), ('verb', 0.25), ('adj', 0.25), ('adv', 0.25)]
[ ] # Calculate KL divergence
    kl_divergence = sum(prob_p * math.log2(prob_p / prob_q) for (label_p, prob_p), (_, prob_q) in zip(p, q_uniform))
[ ] print(f'Information content: {information}')
     print(f'Entropy: {entropy:.3f}')
     print(f'Cross-entropy with a uniform distribution: {cross_entropy:.3f}')
     print(f'KL divergence from a uniform distribution: {kl_divergence:.3f}')
     Information content: {'noun': 1.181, 'verb': 1.971, 'adj': 2.921, 'adv': 2.54}
     Entropy: 1.846
     Cross-entropy with a uniform distribution: 2.000
     KL divergence from a uniform distribution: 0.154
```

Problem Statement:

The task is to calculate information-theoretic metrics (entropy, cross-entropy, KL divergence) for probability distributions represented by lists of tuples, p and q, where each tuple is a speech part and its probability. These metrics provide insights into the properties of the data by quantifying the uncertainty, similarity, and difference between the distributions.

Information Content:

- The quantity of information that each category in the probability distribution carries is represented by its information content. It is calculated as the negative logarithm (base 2) of the probability.
- According to the values and the provided probabilities, adjectives carry the most information and nouns
 the least.

Entropy:

- Entropy quantifies the probability distribution's average level of surprise or uncertainty.
- The average level of uncertainty in the distribution is shown by the entropy value, which is 1.846 bits.

Cross-Entropy with a Uniform Distribution:

- If we utilize a different, usually simpler, distribution for encoding, then cross-entropy estimates the average number of bits required to encode the true distribution.
- If we assume a uniform distribution, the average number of bits required is indicated by the cross-entropy value with a uniform distribution of 2.000 bits.

KL Divergence from a Uniform Distribution:

- This technique calculates the difference in probability distributions between two sets of data. It quantifies the degree to which the true distribution deviates from a uniform distribution in this instance.
- The difference between the true distribution and a uniform distribution is indicated by the KL divergence from a uniform distribution, which is 0.154 bits.

In conclusion, these metrics shed light on the information content, degree of uncertainty, and distinctions between a uniform distribution and the supplied probability distribution. It appears that there is no substantial difference between the original distribution and a uniform distribution, based on the reduced KL divergence.

Q2. Build Bayes' classifiers using CountVectorizer:

- Model 0 Create a baseline model using unigram with other optional parameters
- Model 1 uses only nouns with other optional parameters
- Model 2 uses only verbs with other optional parameters
- Model 3 use your choice of configuration.

1. Build Bayes' classifiers using CountVectorizer.

The code that is provided shows numerous important steps in the model construction process. To ensure consistency and relevancy, data is first cleaned by lowercasing, tokenization, and stopword elimination. After cleaning, text data is converted into numerical representations using CountVectorizer, which is then used for feature extraction and model training. After that, the divided dataset is used to initialize and train Bayesian classifiers. The trained model then makes predictions for the test dataset's labels, allowing performance to be assessed using measures such as accuracy, precision, recall, and F1-score. By using a systematic approach, it is ensured that robust text classification models will be developed that can correctly classify text data into predefined classes.

Problem Statement:

The task involves building Bayes' classifiers using CountVectorizer with various configurations.

Model 0 serves as the baseline, employing unigram features with default parameters.

Model 1 focuses solely on nouns with additional optional parameters

Model 2 is designed for verbs with similar optional configurations.

Model 3 represents a customizable configuration, allowing for unique parameter choices tailored to specific requirements.

Data Processing Pipeline:

The given code goes through a number of steps in the data cleaning process that are meant to get the text data ready for feature extraction and classification.

Dataset Overview: The dataset we are using is the "Spam Dataset" consisting of electronic mail messages categorized into two distinct classes: "spam" and "ham" (non-spam). The dataset contains 2 columns with 5572 observations.

Data Cleaning:

• Lowercasing Text: The lower() technique is used by the clean_text function to convert all text to lowercase. This guarantees consistency in text representation, avoiding problems with case sensitivity in tokenization and analysis.

- **Tokenization:** The word_tokenize function from NLTK splits the text into individual words or tokens once it has been converted to lowercase. In order to facilitate additional processing and analysis, this stage deconstructs the text into its individual components.
- Elimination of Stopwords: The list of frequent stopwords in the English language is provided by the NLTK's stopwords corpus. Using list comprehension and the stopword set, the code eliminates these stopwords from the tokenized text. Stopwords, such "and", "the", and "is" are usually excluded from analyses since they have minimal semantic value and may affect the findings.
- One-hot coding: To further clean the data, we performed one-hot coding on the dataset becasue it had distinct classes, and this will make it easier for classification. We performed this using get_dummies function.

Feature Extraction and Model Training:

Model 0 (Baseline Model using Unigram with Default Parameters):

- Makes use of MultinomialNB from scikit-learn, train_test_split, and the CountVectorizer function with default parameters to extract unigram features from text data.
- Demonstrates strong classification performance on the test set, achieving a high accuracy rate of 97.97%.
- Creates a basic framework that can be used in later iterations for benchmarking and comparison with more complicated models.

Provides insights into the inherent characteristics and challenges of the dataset.

Model0 - Baseline Model CountVectorizer with default parameters

```
import pandas as pd
import nltk
from nltk.import word_tokenize, pos_tag
from nltk.tokenize.treebank import TreebankWordDetokenizer
from nltk.corpus import stopwords
nltk.download('stopwords')
nltk.download('wordnet')
nltk.download('aweraged_perceptron_tagger')
nltk.download('aweraged_perceptron_tagger')
nltk.download('punkt')
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.model_selection import train_test_split
from sklearn.model_selection.test_split
from sklea
```

```
[ ]
           v1
                                                      v2
      0 ham
                  Go until jurong point, crazy.. Available only ...
      1
         ham
                                  Ok lar... Joking wif u oni...
      2 spam Free entry in 2 a wkly comp to win FA Cup fina...
         ham
                U dun say so early hor... U c already then say...
         ham
                  Nah I don't think he goes to usf, he lives aro...
    # Function to clean text
     def clean_text(text):
          # Convert text to lowercase
         text = text.lower()
         # Tokenize text
         tokens = word_tokenize(text)
         # Remove stop words
         stop_words = set(stopwords.words('english'))
          filtered_tokens = [word for word in tokens if word not in stop_words]
         # Rejoin tokens into a string
          cleaned_text = ' '.join(filtered_tokens)
         return cleaned_text
     # Clean the 'v2' column text data
     df['cleaned_text'] = df['v2'].apply(clean_text)
[ ] #further cleaning, renaming and dropping columns
     df['target'] = df['v1']
     df['text'] = df['cleaned_text']
     df_target = pd.get_dummies(df['target']) #get_dummies is used for one-hot-encoding
     df = df.join(df_target)
     df_final = df.drop(columns=['v1', 'v2', 'target', 'cleaned_text'])
    ar_rinal = ar.arop(columns=['vi', 'vz', 'target', 'cleaned_text'])
[ ] df_final.head()
                                     text ham spam
    0 go jurong point , crazy .. available bugis n g... 1
                     ok lar ... joking wif u oni ...
    2 free entry 2 wkly comp win fa cup final tkts 2... 0 1
           u dun say early hor ... u c already say ... 1
         nah n't think goes usf , lives around though 1 0
[ ] #Vectorize the model
    vectorizer = CountVectorizer(ngram_range=(1,1),stop_words= 'english')
    text_count = vectorizer.fit_transform(df_final['text']).toarray()
    print(text_count)
    [[0 0 0 ... 0 0 0]
[0 0 0 ... 0 0 0]
     [0 0 0 ... 0 0 0]
     [0 0 0 ... 0 0 0]
     [0 0 0 ... 0 0 0]
     [0 0 0 ... 0 0 0]]
#training the data on 30% test data and 70% train data
    train_x, test_x, train_y, test_y = train_test_split(text_count, df_final['ham'], test_size= 0.3, stratify= df_final['ham'])
                                                                       + Code + Text
[] #choosing a model
    model0 = MultinomialNB()
    model0.fit(train_x, train_y)
```

```
[ ] #choosing a model
    model0 = MultinomialNB()
    model0.fit(train_x, train_y)
    ▼ MultinomialNB
    MultinomialNB()
[ ] predicted = model0.predict(test_x)
[ ] print(classification_report(test_y, predicted))
                 precision recall f1-score
                      0.91
                                0.95
                                          0.93
                                                    224
                      0.99
                                0.98
                                                    1448
                                          0.99
                                          0.98
                                                    1672
                      0.95
                                0.97
                                                    1672
    weighted avg
                      0.98
                                0.98
                                          0.98
                                                    1672
print(metrics.accuracy_score(predicted, test_y))
    0.9796650717703349
```

Model 1 (Utilizes only nouns with additional optional parameters):

- Maintains a competitive accuracy rate of 97.61% on the test set, despite narrowing down feature selection to nouns.
- Demonstrates the impact of feature selection strategies on model performance and efficiency.
- Highlights the importance of selecting relevant linguistic features for classification tasks.
- Illustrates the versatility of the CountVectorizer method in tailoring feature extraction to specific linguistic elements.

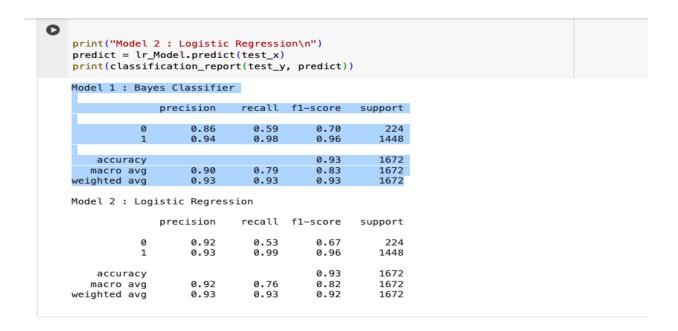
```
import pandas as pd
import nltk
from nltk import word_tokenize, pos_tag
from nltk.tokenize.treebank import TreebankWordDetokenizer
from nltk.corpus import stopwords
nltk.download('stopwords')
nltk.download('wordnet')
nltk.download('averaged_perceptron_tagger')
nltk.download('averaged_perceptron_tagger')
nltk.download('averaged_perceptron_tagger')
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.model_selection import train_test_split
from sklearn.model_selection.import train_test_
```

```
[ ] df = pd.read_csv("/content/drive/MyDrive/02 NLP/Assignments/spam.csv")
      df.head()
             v1
       0 ham
                    Go until jurong point, crazy.. Available only ...
                                     Ok lar... Joking wif u oni...
       1 ham
       2 spam Free entry in 2 a wkly comp to win FA Cup fina...
       3 ham U dun say so early hor... U c already then say...
       4 ham Nah I don't think he goes to usf, he lives aro...
 def clean_text(text):
           # Convert text to lowercase
           text = text.lower()
           # Tokenize text
           tokens = word_tokenize(text)
           #Tagging tokens
           tagged_tokens = pos_tag(tokens)
           #Filtering only nouns
           nouns = [word for word, pos in tagged_tokens if pos.startswith('NN') or pos == 'NNP' or pos == 'NNS']
           # Remove stop words
           stop_words = set(stopwords.words('english'))
           filtered_tokens = [word for word in nouns if word not in stop_words]
           # Rejoin tokens into a string
cleaned_text = ' '.join(filtered_tokens)
           return cleaned_text
      # Clean the 'v2' column text data
      df['cleaned_text'] = df['v2'].apply(clean_text)
[ ] #further cleaning, renaming and dropping columns
df('target'] = df('v1')
df('text') = df('cleaned_text')
df_target = pd.get_dummles(df('target')) #get_dummles is used for one-hot-encoding
     df = df.join(df_target)
df = df.join(df_target)
df_final = df.drop(columns=['v1','v2','target','cleaned_text'])
df_final.head()
     0 point .. bugis world la buffet cine wat 1 0
                                        lar oni
     2 entry comp tkts text fa entry question txt rate c 0 1
                                     dun hor c 1 0
                                    nah 1 0
#Vectorize the model
vectorizer = CountVectorizer(ngram_range=(1,1),stop_words= 'english')
     text_count = vectorizer.fit_transform(df_final['text']).toarray()
[0 0 0 ... 0 0 0]]
[0 0 0 ... 0 0 0]
[0 0 0 ... 0 0 0]
[ ] train_x, test_x, train_y, test_y = train_test_split(text_count, df_final['ham'], test_size= 0.3, stratify= df_final['ham'])
 [ ] #choosing a model
    model1 = MultinomialNB()
       model1.fit(train_x, train_y)
       w MultinomialNB
MultinomialNB()
[ ] predicted = model1.predict(test_x)
[ ] print(classification_report(test_y, predicted))
                                            recall f1-score
                          precision
                                                                       support
                                               0.92
0.99
                                                              0.91
0.99
                                                              0.98
0.95
0.98
                                  0.95
0.98
 print(metrics.accuracy_score(predicted, test_y))
       0.9760765550239234
```

Model 2 (Utilizes solely verbs with other optional parameters):

- This model achieves a 93% accuracy rate on the test set by only focusing on verbs as major features.
- Highlights how important linguistic elements are for accurately classifying texts, especially verbs.
- Provides information about the importance of various speech segments in classification tasks.
- Provides possibilities for more investigation and enhancement of feature selection tactics for particular text data attributes.

```
text ham spam
                 Go until jurong point, crazy.. Available only ... 1 0
                                                  Ok lar... Joking wif u oni...
         2 Free entry in 2 a wkly comp to win FA Cup fina... 0 1
                    U dun say so early hor... U c already then say... 1
         4 Nah I don't think he goes to usf, he lives aro... 1 0
          5567 This is the 2nd time we have tried 2 contact u... 0 1
          5568
                                   Will I b going to esplanade fr home? 1
                                                                                                        0
          5569
                      Pity, * was in mood for that. So...any other s... 1
                                                                                                       0
                      The guy did some bitching but I acted like i'd...
                                  Rofl. Its true to its name 1 0
         5572 rows x 3 columns
         #Keep only word where POS starts with 'V' POS = 'V'
         new_text =list()
                 i in data.text:
             # Convert text to lowercase
             i = i.lower()
            1 = 1.tower()
temp = word_tokenize(i)
tags = nltk.pos_tag(temp)
verb = [tag[0] for tag in tags if tag[1].startswith(POS)]
stop_words = set(stopwords.words('english'))
filtered_tokens = [word for word in verb if word not in stop_words]
new_text.append(" ".join(filtered_tokens))
      #Keep only word where POS starts with 'V' POS = 'V'
      POS = 'V'
new_text =list()
for i in data.text:
    # Convert text to lowercase
    i = i.lower()
         i = i.lower()
temp = word_tokenize(i)
tags = nltk.pos_tag(temp)
verb = [tag[0] for tag in tags if tag[1].startswith(POS)]
stop_words = set(stopwords.words('english'))
filtered_tokens = [word for word in verb if word not in stop_words]
new_text.append(" ".join(filtered_tokens))
[ ] # Replace new text in dataset data.text = new_text
      #use CountVectorizer to coner sentence into list of words
       cv = CountVectorizer(
      text_count = cv.fit_transform(data['text']).toarray()
#split dataset into train test set
train_x, test_x, train_y, test_y = train_test_split(text_count, data['ham'], test_size=0.3, stratify=data['ham'])
      #calling Multinomial model and training the train dataset model = MultinomialNB() model.fit(train_x, train_y)  
      #calling Logistic regression model and training the train dataset
lr_Model = LogisticRegression()
lr_Model.fit(train_x, train_y)
      #make predictio on test data on both models
print("Model 1 : Bayes Classifier\n")
predict = model.predict(test_x)
print(classification report(test v. predict))
```



Model 3 (Customizable configuration with user-defined parameters):

- Using TfidfVectorizer, text data was transformed into TF-IDF features up to 5000 features, using bigrams and unigrams.
- Using the SVC function, an SVM model with a linear kernel and C=1.0 was trained.
- The predict method of the SVM model was used to create model predictions for the test set.
- The accuracy score, which was calculated using accuracy_score, produced a high result of
 0.9809, which represents the percentage of test set occurrences that were properly classified.

```
# Function to clean text
[ ]
     def clean_text(text):
    # Convert text to lowercase
           text = text.lower()
           # Tokenize text
           tokens = word_tokenize(text)
# Remove stop words
           stop_words = set(stopwords.words('english'))
filtered_tokens = [word for word in tokens if word not in stop_words]
          # Rejoin tokens into a string
cleaned_text = ' '.join(filtered_tokens)
           return cleaned_text
     # Clean the 'v2' column text data
df['cleaned_text'] = df['v2'].apply(clean_text)
[ ] #further cleaning, renaming and dropping columns
   df['target'] = df['v1']
   df['text'] = df['cleaned_text']
      df_target = pd.get_dummies(df['target']) #get_dummies is used for one-hot-encoding
      df = df.join(df_target)
df_final = df.drop(columns=['v1','v2','target','cleaned_text'])
      df final.head()
                                                 text ham spam
      0 go jurong point , crazy .. available bugis n g... 1 0
                             ok lar ... joking wif u oni ...
      2 free entry 2 wkly comp win fa cup final tkts 2... 0 1
      3
               u dun say early hor ... u c already say ...
      4 nah n't think goes usf , lives around though 1 0
[ ] X = vectorizer.fit_transform(df['cleaned_text'])
y = df['v1']
[ ] # Split the data into train and test sets
     X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=8.3, random_state=42)
[ ] # Initialize and train the SVM model
    svm_model = SVC(kernel='linear', C=1.0, random_state=42)
    svm_model.fit(X_train, y_train)
                        SVC
     SVC(kernel='linear', random_state=42)
[ ] # Predictions using SVM model
     svm_y_pred = svm_model.predict(X_test)
[ ] print(classification_report(y_test, svm_y_pred))
                     precision recall f1-score support
             5 pam
                          0.97
                                     0.88
                                                  0.92
                                                              219
                                                 0.98
                                     0.94
     macro avg
weighted avg
                                                             1672
                                                  0.98
                                                             1672
O print("SVM Accuracy:", accuracy_score(y_test, svm_y_pred))
SVM Accuracy: 0.9000612440191388
```

Performance Comparison:

Model Accuracy Precision Recall F1-Score
--

Model 0	0.9797	0.91	0.95	0.93
Model 1	0.9761	0.91	0.92	0.91
Model 2	0.9300	0.86	0.59	0.70
Model 3	0.9809	0.98	0.94	0.96

In summary up, the four models show different levels of performance on several parameters. With its excellent accuracy, precision, recall, and F1-Score for both classes, Model 0 is clearly the most reliable. Model 1, which just considers nouns, does well overall despite a little decline in recall indicates a slight decrease in its ability to identify all relevant instances of the positive class. On the other hand, Model 2, which only emphasizes verbs, shows an apparent decrease in performance, especially when it comes to identifying class 0 cases. Model 3, which uses SVM, shows great precision, recall, and F1-Score for both classes, making it a strong competitor.

Comparatively, Model 0 and Model 3 stand out for their high accuracy and balanced metrics, making them reliable choices for spam detection. While Model 0 offers a solid baseline with its use of a Multinomial Naive Bayes classifier and unigram features, Model 3's SVM approach and possibly richer feature set allow it to slightly edge out in performance, especially in precision and F1-score.

Recommendations for Optimizing Model Parameters:

In order to improve the performance of each model overall, a thorough grid search can be used to optimize the model parameters. This involves looking into different hyperparameter setups. Their prediction power might be increased by fine-tuning variables including feature weights, kernel type, and regularization strength. In the end, selecting one of these models should be based on the particular needs and goals of the assigned task, taking precision, recall, and total accuracy into account. Furthermore, using cross-validation methods like k-fold cross-validation can yield estimates of each model's performance over various data subsets that are more trustworthy. Moreover, choosing the best model to use in practical situations requires careful consideration of the possible trade-offs between interpretability and model complexity. Thus, in order to guarantee the best possible performance and generalizability, a comprehensive assessment of several model architectures and parameter settings is necessary.

Q3. Build four logistic regression models that have the comparable parameters as the four models above. Compare the performance of all four models. Write a report about your data preprocessing pipeline, modeling, performance and discussion about the model performance. Make recommendation on optimizing model parameters.

A. Logistic Regression

Logistic regression stands as a fundamental and powerful statistical method for its application in classification tasks. Logistic regression is designed to predict the probability that a given input belongs to a certain category. This makes it especially suited for binary classification problems, such as spam detection, where each piece of text must be categorized into one of two groups: spam or not spam. However, its extends beyond binary classification through the use of multinomial logistic regression, which can handle multiple categories, making it versatile across a wide array of classification challenges.

Problem Statement

The task involves building logistic regression models with various configurations.

Model 0 serves as the baseline, employing unigram features with default parameters.

Model 1 focuses solely on nouns with additional optional parameters

Model 2 is designed for verbs with similar optional configurations.

Model 3 represents a customizable configuration, allowing for unique parameter choices tailored to specific requirements.

Data Processing Pipeline (Same as what we did for Bayes' Classifier)

The given code goes through a number of steps in the data cleaning process that are meant to get the text data ready for feature extraction and classification.

Dataset Overview: The dataset we are using is the "Spam Dataset" consisting of electronic mail messages categorized into two distinct classes: "spam" and "ham" (non-spam). The dataset contains 2 columns with 5572 observations.

Data Cleaning:

• Lowercasing Text: The lower() technique is used by the clean_text function to convert all text to lowercase. This guarantees consistency in text representation, avoiding problems with case sensitivity in tokenization and analysis.

- **Tokenization:** The word_tokenize function from NLTK splits the text into individual words or tokens once it has been converted to lowercase. In order to facilitate additional processing and analysis, this stage deconstructs the text into its individual components.
- Elimination of Stopwords: The list of frequent stopwords in the English language is provided by the NLTK's stopwords corpus. Using list comprehension and the stopword set, the code eliminates these stopwords from the tokenized text. Stopwords, such "and", "the", and "is" are usually excluded from analyses since they have minimal semantic value and may affect the findings.
- One-hot coding: To further clean the data, we performed one-hot coding on the dataset becasue it had distinct classes, and this will make it easier for classification. We performed this using get dummies function.

Feature Extraction & Model Training

Model 0 (Baseline Model using Unigram with Default Parameters):

- Model 0 achieved an overall accuracy of 97.66%, with a precision of 98% for both classes, which means model is highly reliable in its predictions across the dataset.
- The model has a perfect recall of 1.00 for class 1, indicating it correctly identified all instances of this class, the recall for class 0 is lower at 0.84. This discrepancy suggests that while the model is excellent at identifying the majority class, it struggles with the minority class to some extent.
- The F1-scores of 0.91 for class 0 and 0.99 for class 1 indicate a strong balance between precision and recall for the model.

Model 0 - Baseline Model CountVectorizer with default parameters Performing logistic regression

```
[1]
      1 import pandas as pd
       2 import nltk
       3 from nltk import word tokenize, pos tag
       4 from nltk.tokenize.treebank import TreebankWordDetokenizer
       5 from nltk.corpus import stopwords
       6   nltk.download('stopwords')
         nltk.download('wordnet')
      7
      8 nltk.download('omw-1.4')
      9 nltk.download('averaged_perceptron_tagger')
      10     nltk.download('punkt')
      11 from sklearn.feature_extraction.text import CountVectorizer
          from sklearn.model_selection import train_test_split
          from sklearn.naive_bayes import GaussianNB, CategoricalNB, MultinomialNB
     13
     14 from sklearn import metrics
     15 from sklearn.metrics import classification report, accuracy score
      16 from sklearn.linear model import LogisticRegression
     [nltk data] Downloading package stopwords to /root/nltk data...
      [nltk data] Unzipping corpora/stopwords.zip.
      [nltk data] Downloading package wordnet to /root/nltk data...
      [nltk data] Downloading package omw-1.4 to /root/nltk data...
      [nltk data] Downloading package averaged perceptron tagger to
      [nltk data]
                     /root/nltk data...
      [nltk_data] Unzipping taggers/averaged_perceptron_tagger.zip.
      [nltk data] Downloading package punkt to /root/nltk data...
      [nltk data] Unzipping tokenizers/punkt.zip.
       1 #Loading the datset
[3]
          df = pd.read csv("spam.csv")
[4]
          #Displaying the first 5 rows of the dataset
      1
          df.head(n=5)
            v1
                   Go until jurong point, crazy.. Available only ...
          ham
                                  Ok lar... Joking wif u oni...
          ham
      2 spam Free entry in 2 a wkly comp to win FA Cup fina...
                U dun say so early hor... U c already then say...
          ham
                  Nah I don't think he goes to usf, he lives aro...
          ham
```

LogisticRegression()

```
\frac{\checkmark}{5s} [5] 1 # Function to clean text
             def clean_text(text):
                # Convert text to lowercase
         3
                 text = text.lower()
         4
                # Tokenize text
         5
                tokens = word_tokenize(text)
         6
         7
                # Remove stop words
                stop_words = set(stopwords.words('english'))
                filtered_tokens = [word for word in tokens if word not in stop_words]
         9
                # Rejoin tokens into a string
        10
                cleaned_text = ' '.join(filtered_tokens)
        11
               return cleaned_text
        12
        14 # Clean the 'v2' column text data
        df['cleaned_text'] = df['v2'].apply(clean_text)
        16
   0
        1 #further cleaning, renaming and dropping columns
         2 df['target'] = df['v1']
         3 df['text'] = df['cleaned_text']
         4 df_target = pd.get_dummies(df['target']) #get_dummies is used for one-hot-encoding
         5 df = df.join(df_target)
         6 df_final = df.drop(columns=['v1','v2','target','cleaned_text'])
         7 df_final.head()
_{0s} [7] 1 #Vectorize the model
       vectorizer = CountVectorizer(ngram_range=(1,1),stop_words= 'english')
      4 text count = vectorizer.fit transform(df final['text']).toarray()
      5 print(text_count)
      [[000...000]
      [000...000]
      [000...000]
\frac{\checkmark}{0s} [8] 1 #training the data on 30% test data and 70% train data
      2 train_x, test_x, train_y, test_y = train_test_split(text_count, df_final['ham'], test_size= 0.3, stratify= df_final['ham'])
y choosing a model
         model0 = LogisticRegression()
      4 model0.fit(train_x, train_y)
```

```
predicted = model0.predict(test x)
          print(classification_report(test_y, predicted))
[11]
                   precision
                                recall f1-score
                                                    support
                0
                        0.98
                                  0.84
                                            0.91
                                                       224
                1
                        0.98
                                  1.00
                                            0.99
                                                      1448
         accuracy
                                            0.98
                                                      1672
                        0.98
                                  0.92
                                            0.95
                                                      1672
        macro avg
                        0.98
                                  0.98
                                            0.98
     weighted avg
                                                      1672
          print(metrics.accuracy_score(predicted, test_y))
     0.9766746411483254
```

Model 1 (Utilizes only nouns with additional optional parameters):

- The model achieved a very high overall **accuracy of 96.59%**, which means that it is capable of correctly classifying the majority of instances in the test set.
- With a perfect recall of 1.00 and an F1-score of 0.98 for class 1, the model demonstrates exceptional performance in identifying the positive class. This means that nouns are highly indicative of the positive class in the dataset, and the model rarely misses any positive instances.
- The recall for class 0 is significantly lower at 0.75 compared to class 1. This suggests that while the model is very good at identifying the positive class, it may be overlooking a quarter of the negative class instances.

```
7s D
        1 import pandas as pd
         2 import nltk
         3 from nltk import word_tokenize, pos_tag
4 from nltk.tokenize.treebank import TreebankWordDetokenizer
         5 from nltk.corpus import stopwords
         6 nltk.download('stopwords')
         7 nltk.download('wordnet')
         8 nltk.download('omw-1.4')
         9 nltk.download('averaged_perceptron_tagger')
        10 nltk.download('punkt')
        11   nltk.download('tagsets')
12   nltk.download('averaged_perceptron_tagger')
        13 from sklearn.feature_extraction.text import CountVectorizer
        14 from sklearn.model_selection import train_test_split
        from sklearn.linear_model import LogisticRegression
from sklearn import metrics
        17 from sklearn.metrics import classification_report, accuracy_score
   [nltk_data] Downloading package stopwords to /root/nltk_data...
                     Unzipping corpora/stopwords.zip.
        [nltk data]
        [nltk_data] Downloading package wordnet to /root/nltk_data...
        [nltk_data] Downloading package omw-1.4 to /root/nltk_data...
        [nltk_data] Downloading package averaged_perceptron_tagger to
        [nltk data]
                        /root/nltk data...
        [nltk data]
                     Unzipping taggers/averaged_perceptron_tagger.zip.
        [nltk_data] Downloading package punkt to /root/nltk_data...
        [nltk data] Unzipping tokenizers/punkt.zip.
        [nltk_data] Downloading package tagsets to /root/nltk_data...
        [nltk data] Unzipping help/tagsets.zip.
 [ ] 1 def clean_text(text):
        2
             # Convert text to lowercase
               text = text.lower()
        3
        4
               # Tokenize text
              tokens = word_tokenize(text)
#Tagging tokens
        5
        6
        7
               tagged_tokens = pos_tag(tokens)
        8
        9
               #Filtering only nouns
       10
                nouns = [word for word, pos in tagged_tokens if pos.startswith('NN') or pos == 'NNP' or pos == 'NNS']
                # Remove stop words
       11
       12
                stop_words = set(stopwords.words('english'))
                filtered_tokens = [word for word in nouns if word not in stop_words]
       13
       14
                # Rejoin tokens into a string
                cleaned text = ' '.join(filtered_tokens)
       15
       16
                return cleaned_text
       1 # Clean the 'v2' column text data
            df['cleaned_text'] = df['v2'].apply(clean_text)
            df.head()
        3
             v1
                                                        v2
                                                                                      cleaned_text
       0 ham
                    Go until jurong point, crazy.. Available only ...
                                                                    point .. bugis world la buffet cine wat
                                    Ok lar... Joking wif u oni...
```

2 spam Free entry in 2 a wkly comp to win FA Cup fina... entry comp tkts text fa entry question txt rate c

```
df_target = pd.get_dummies(df['v1']) #get_dummies is used for one hot encoding
    2 df = df.join(df_target)
    4 df_final = df.drop(columns=['v1', 'v2'])
    5 df final.head()
(2)
                         cleaned_text ham spam
            point .. bugis world la buffet cine wat 1
    1
    2 entry comp tkts text fa entry question txt rate c 0
                            dun hor c 1
[ ] 1 vectorizer = CountVectorizer(ngram_range=(1,1),stop_words='english')
     3 text_count = vectorizer.fit_transform(df_final['cleaned_text']).toarray()
[ ] 1 train_x, test_x, train_y, test_y = train_test_split(text_count, df_final['ham'], test_size= 0.3, stratify= df_final['ham'])
               #choosing a model
  [ ]
                model = LogisticRegression()
           2
           3
                model.fit(train x, train y)
           4
          ▼ LogisticRegression
          LogisticRegression()
           1
                predicted = model.predict(test_x)
           2
           3
                print(classification_report(test_y, predicted))
           4
                print(metrics.accuracy score(predicted, test y))
  (2)
                            precision
                                              recall f1-score
                                                                        support
                        0
                                   0.99
                                                0.75
                                                              0.86
                                                                             224
                        1
                                   0.96
                                                1.00
                                                              0.98
                                                                            1448
              accuracy
                                                              0.97
                                                                            1672
                                                              0.92
                                                                            1672
             macro avg
                                   0.98
                                                0.88
         weighted avg
                                  0.97
                                                0.97
                                                              0.96
                                                                            1672
         0.9659090909090909
```

Model 2 (Utilizes solely verbs with other optional parameters):

- The overall accuracy of the model is 0.93, which means that **93%** of the predictions made by the model are correct.
- The model has high precision for both classes (0.91 for class 0 and 0.93 for class 1), which means that when it predicts an instance to be in a particular class, it is likely to be correct. This indicates that the verbs present in the text are reliable indicators for the model's predictions.
- The F1-score for class 0 is 0.67, which is moderate and suggests a balance between precision and recall that is leaning towards precision. The F1-score for class 1 is very high at 0.96, reflecting the high precision and recall for this class.



```
dataset["target"] = dataset["v1"]
dataset["text"] = dataset["v2"]
# one hot encoding
data_target = pd.get_dummies(dataset["target"])
dataset = dataset.join(data_target)
dataset1 = dataset.drop(columns=["v1","v2","target"])
data = dataset1
data
```

	text	ham	spam
0	Go until jurong point, crazy Available only	1	0
1	Ok lar Joking wif u oni	1	0
2	Free entry in 2 a wkly comp to win FA Cup fina	0	1
3	U dun say so early hor U c already then say	1	0
4	Nah I don't think he goes to usf, he lives aro	1	0
5567	This is the 2nd time we have tried 2 contact u	0	1
5568	Will i_b going to esplanade fr home?	1	0
5569	Pity, * was in mood for that. Soany other s	1	0
5570	The guy did some bitching but I acted like i'd	1	0
5571	Rofl. Its true to its name	1	0

```
#calling Logistic regression model and training the train dataset
lr_Model = LogisticRegression()
lr_Model.fit(train_x, train_y)

print("Model 2 : Logistic Regression\n")
predict = lr_Model.predict(test_x)
print(classification_report(test_y, predict))
```

Model 2 : Logistic Regression

	precision	recall	f1-score	support
0	0.91	0.53	0.67	224
1	0.93	0.99	0.96	1448
accuracy			0.93	1672
macro avg	0.92	0.76	0.82	1672
weighted avg	0.93	0.93	0.92	1672

Model 3 (Customizable configuration with user-defined parameters):

- After Data cleaning, The cleaned text data is transformed into a numerical format using TF-IDF (Term Frequency-Inverse Document Frequency) vectorization. This step converts the text into a sparse matrix of token counts, weighted by the importance of each word within the dataset and across documents, with considerations for both unigrams and bigrams.
- We have used Logistic Regression model with L2 regularization which is trained on the vectorized text data.
- The model has a high **accuracy rate of 95.04%**, meaning that it correctly identifies spam and ham messages the majority of the time.
- The model has a very high precision rate of 97% for spam detection, indicating that when it predicts a message is spam. However, the recall is 64%, meaning it only correctly identifies 64% of all actual spam messages.

4th Logistic Model

```
# Importing necessary libraries
         import pandas as pd
     3 import nltk
     4 from nltk.corpus import stopwords
      5 from sklearn.feature_extraction.text import TfidfVectorizer
      6 from sklearn.model_selection import train_test_split
         from sklearn.linear model import LogisticRegression
     8 from sklearn.metrics import classification_report, accuracy_score, confusion_matrix
     9 from nltk.tokenize.treebank import TreebankWordDetokenizer
    10 from nltk import word_tokenize, pos_tag
    11
    12 # Download NLTK resources
    13    nltk.download('stopwords')
    14   nltk.download('wordnet')
    15    nltk.download('punkt')
[nltk_data] Downloading package stopwords to /root/nltk_data...
    [nltk data] Unzipping corpora/stopwords.zip.
    [nltk_data] Downloading package wordnet to /root/nltk_data...
    [nltk_data] Downloading package punkt to /root/nltk_data...
    [nltk data] Unzipping tokenizers/punkt.zip.
    True
[ ] 1 # Loading the dataset
         df = pd.read_csv("spam.csv")
         #Displaying the first 5 rows of the dataset
         df.head(n=5)
          v1
                                                    v2
     0 ham
                 Go until jurong point, crazy.. Available only ...
         ham
                                 Ok lar... Joking wif u oni...
     2 spam Free entry in 2 a wkly comp to win FA Cup fina...
         ham
               U dun say so early hor... U c already then say...
         ham
                 Nah I don't think he goes to usf, he lives aro...
[] 1
         # Function to clean text
      2
         def clean_text(text):
             # Convert text to lowercase
      3
             text = text.lower()
      4
      5
             # Tokenize text
      6
             tokens = word tokenize(text)
             # Remove stop words
              stop_words = set(stopwords.words('english'))
      8
      9
             filtered_tokens = [word for word in tokens if word not in stop_words]
     10
              # Rejoin tokens into a string
```

```
[ ] 1 # Vectorize the text using TF-IDF
      vectorizer = TfidfVectorizer(ngram_range=(1, 2), max_features=5000)
     3 X = vectorizer.fit_transform(df['cleaned_text'])
     4 \quad y = df['v1']
     1 # Split the data into train and test sets
      2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
1 l2_model = LogisticRegression(penalty='l2', solver='liblinear', random_state=42)
     3 l2_model.fit(X_train, y_train)
     5 y_pred = l2_model.predict(X_test)
     7 # Evaluate the model
     8 accuracy = accuracy_score(y_test, y_pred)
     9 conf_matrix = confusion_matrix(y_test, y_pred)
    10 classification_rep = classification_report(y_test, y_pred)
    11
    12 # Print the results
    13 print(f"Accuracy: {accuracy}")
    14 print(f"Confusion Matrix:\n{conf_matrix}")
     15 print(f"Classification Report:\n{classification_rep}")
```

```
Accuracy: 0.9503588516746412
   Confusion Matrix:
   [[1449
            41
    [ 79 140]]
   Classification Report:
                precision recall f1-score support
            ham
                     0.95
                              1.00
                                       0.97
                                                1453
           spam
                     0.97
                              0.64
                                       0.77
                                                 219
       accuracy
                                       0.95
                                                1672
      macro avg
                     0.96
                              0.82
                                       0.87
                                                1672
   weighted avg
                     0.95
                              0.95
                                       0.95
                                                1672
```

Performance Comparision

Model	Accuracy	Precision	Recall	F1-Score
Model 0	0.976	0.98	0.84	0.91
Model 1	0.965	0.99	0.75	0.86
Model 2	0.92	0.91	0.53	0.67
Model 3	0.950	0.95	1.00	0.97

- The baseline unigram model (Model 0) seems to be quite strong, indicating that even a simple approach can be very effective for spam detection tasks.
- The nouns-only model (Model 1), while being the most precise, sacrifices recall. This could be problematic in spam detection where missing spam messages (false negatives) can be more than mistakenly filtering out non-spam messages (false positives).
- The verbs-only model (Model 2) performs the poorest, which suggests that verbs alone do not carry enough discriminative information for spam detection in this context.
- Model 3 demonstrates the importance of comprehensive text vectorization and regularization. It
 does not have the highest accuracy, but it achieves a perfect recall score, which is often more
 important in spam detection scenarios to ensure that all spam messages are caught.

In conclusion, Model 3, despite slightly lower accuracy, is arguably the best model out of the four due to its perfect recall and highest F1-score, ensuring that no spam messages are overlooked.

Recommendations for Optimizing Model Parameters:

To optimize a logistic regression model for text classification, the recommendations are as follows:

- Extend grid search with cross-validation to fine-tune hyperparameters such as the regularization strength and the solver algorithm, which can impact model performance.
- Feature engineering techniques, especially experimenting with different n-gram ranges in TF-IDF vectorization, to better capture contextual information from the text data.
- Consider class weights adjustment to address class imbalances, which can improve the recall for underrepresented classes and lead to a more balanced model performance.

Q4. Write a conclusion about the performance of the eight models.

	lassifiers		

Model	Accuracy	Precision	Recall	F1-Score
Model 0	0.9797	0.91	0.95	0.93
Model 1	0.9761	0.91	0.92	0.91
Model 2	0.9300	0.86	0.59	0.70
Model 3	0.9809	0.98	0.94	0.96

- Model 0 shows strong overall performance with high accuracy (0.9797), precision (0.91), recall (0.95), and F1-score (0.93).
- Model 1 has a slight drop in performance compared to Model 0, with a small decrease across all
 metrics.
- Model 2 has a significant drop in recall (0.59), which drags down its F1-score (0.70), meaning it struggles to identify all positive instances.
- Model 3 shows the best performance within the Bayes' classifiers, with the highest accuracy (0.9809), precision (0.98), and F1-score (0.96), and a high recall (0.94).

The Bayes' classifiers with the combination of TF-IDF vectorization and SVM (Model 3) is the most effective approach, especially in balancing precision and recall.

Logistic Regression Models:

Model	Accuracy	Precision	Recall	F1-Score
Model 0	0.976	0.98	0.84	0.91
Model 1	0.965	0.99	0.75	0.86
Model 2	0.92	0.91	0.53	0.67
Model 3	0.950	0.95	1.00	0.97

- Model 0 starts strong with high accuracy (0.976) and precision (0.98), but the recall is lower (0.84) compared to its Bayes' counterpart.
- Model 1 maintains high precision (0.99) but shows a significant decrease in recall (0.75), leading to a lower F1-score (0.86).
- Model 2 shows the weakest performance in this set with the lowest recall (0.53) and F1-score (0.67), indicating a difficulty in identifying true positives.
- Model 3 has high accuracy (0.950), good precision (0.95), and an excellent recall of 1.00, resulting in the highest F1-score (0.97) among the logistic regression models.

For logistic regression, Model 3 stands out with perfect recall, indicating that every positive instance has been identified. However, this perfect recall should be scrutinized for potential overfitting in evaluation.

Comparing both the model types performance:

- Models using only nouns (Model 1) tend to perform better than those using only verbs (Model 2),
 which suggests that nouns carry more discriminative power for the classification tasks at hand.
- TF-IDF Vectorization: Model 3's use of TF-IDF vectorization is clearly beneficial, as it consistently is in top performance in both sets of models, highlighting the value of weighting terms based on their frequency across documents.
- The Bayes' classifiers seem more consistent across different feature extraction methods compared to logistic regression models. However, the logistic regression model with TF-IDF vectorization and L2 regularization (Model 3) achieves the highest F1-score overall, which is a balanced measure of precision and recall.
- The use of L2 regularization in logistic regression (for Model 3) seems to have a positive impact on the model's ability to generalize, as evidenced by the high F1-score.

Links to source code

Question 1: Python code for information, cross-entropy, KL-Divergenge

Bayes' classifiers Model 0

Bayes' classifiers Model 1

Bayes' classifiers Model 2

Bayes' classifiers Model 3

Logistic Regression Model 0

Logistic Regression Model 1

Logistic Regression Model 2

Logistic Regression Model 3