```
1   Scope of 'Contents' {
2
3
4        01   Introduction to Naive Bayes
5             < concept, characteristics,
6             flaws >
7
8                  02   Choice of data set
9                       < considerations,
10                      preprocessing >
11
12                            03   Code for Model
13                                 < Choice of Vectorizer/confusion
14   }                            matrix >
```

```
1
2    01 {
3
4
5
6        [Naive Bayes]
7
8            <introduction>
9
10
11
12    }
13
14
```

1  **What is Naive Bayes {**
2
3        < It is a family of probabilistic classifiers
         based on applying Bayes' theorem with a strong
4        (naive) assumption of independence between
         features. >
5
6  **}**
7
8
9  **What makes it 'naive'? {**
10
11       < It assumes that all features are independent
         given the class >
12
13
14 **}**

1  # Key concepts; {
2
3
4          'Bayes' Theorem:'
5
6                  <p It calculates the probability of a class
7                  $C$ given some evidence (features) $X$:
8                    >
9
10                                $$P(C|X) = \frac{P(X|C) \times P(C)}{P(X)}$$
11
12          </p>
13
14  }

```
1   Key concepts; {
2
3
4
5       Naive Assumptions:'
6           <p It assumes that all features are
7           independent given the class:
8            >
9
10
11
12          </p>
13
14  }
```

$$P(X|C) = \prod_i P(x_i|C)$$

1  How does it work? {
2
3          < The algorithm learns the prior probability of
4          each class P(C) and the likelihood of features
           given the class P(Xi|C) from the training data. >
5
6  }
7
8
9  Predictions {
10         < For a new input, it computes the probability
11         for each class using Bayes' theorem and
           predicts the class with the highest
12         probability. >
13
14 }

1  **Why is Naive Bayes popular for**
2  **'text classification?'** {
3      01    Works well with high-dimensional
4            data like text, where features are
5            word counts or frequencies.
6      02  Fast and scalable.
7
8
9      03  Often surprisingly accurate despite
              the strong independence assumption.
10
11          Commonly used for spam detection,
12     04   sentiment analysis, document
              classification, and more.
13
14  }

Created by Evan Tan

1  # Types of Naive Bayes classifiers;
2
3  {
4      [🔍]  Multinomial Naive Bayes
5            <Best for discrete features like word counts
6            (common in text).>
7      [🗒]  Bernoulli Naive Bayes
8            <For binary/boolean features (word presence
9            or absence).>
10
11     [☁]  Gaussian Naive Bayes
12           <For continuous data assuming normal
13           distribution.>
14  }

```
1    Examples About 'Naive Bayes for Spam Detection (Toy
2    Data)'{
3
4                "Free money now"        "Call me now"
5         [ ☁ ]                    [ 🛡 ]
6              < Spam >                 < Ham[1] >
7
8
9
10             "Win money win prize"   "Hello how are you"
11        [ ▦ ]                   [ ▣ ]
12             < spam >                < Ham >
13
14   }
```

```
1
2     02  {
3
4
5
6        [Choice of data set]
7
8            <considerations/
9            preprocessing>
10
11
12    }
13
14
```

```
1   Chosen data set {
2
3
        <SMS Spam Collection Dataset >
4
5
6
7                              *   The SMS Spam Collection is a set of SMS
        < /1 >                     tagged messages that have been
8                                  collected for SMS Spam research.
9
10                             *   The files contain one message per line.
11         < /2 >                  Each line is composed by two columns:
                                   v1 contains the label (ham or spam) and
12                                 v2 contains the raw text.
13
14  }
```

```
1
2  Content {
3
4      < This corpus has been collected from free or
5      free for research sources at the Internet >
6
7      425 SMS spam messages
8      < Manually extracted from the
9      Grumbletext Web site. >
10
11     3,375 SMS randomly chosen ham messages
12     <from the NUS SMS Corpus (NSC), which is a dataset of
13  }  about 10,000 legitimate messages collected for
14     research at the Department of Computer Science at the
       National University of Singapore >
```
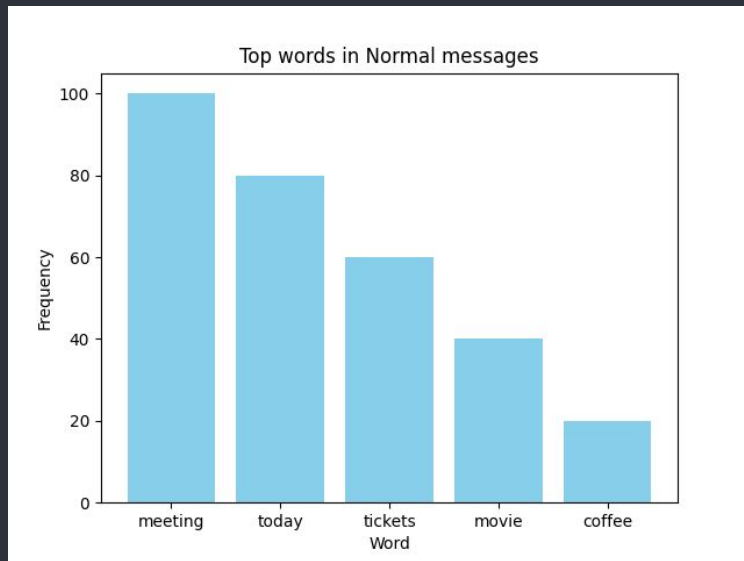
Created by Evan Tan

```
1
2
3   Bar chart of
4   'Normal Messages'
5
6   {
7
8       < Figure: Top words in ham
9       messages by frequency. The
10      training text has been cleaned
11      and aggregated to count words
12      (excluding generic terms).>
13
14  }
```
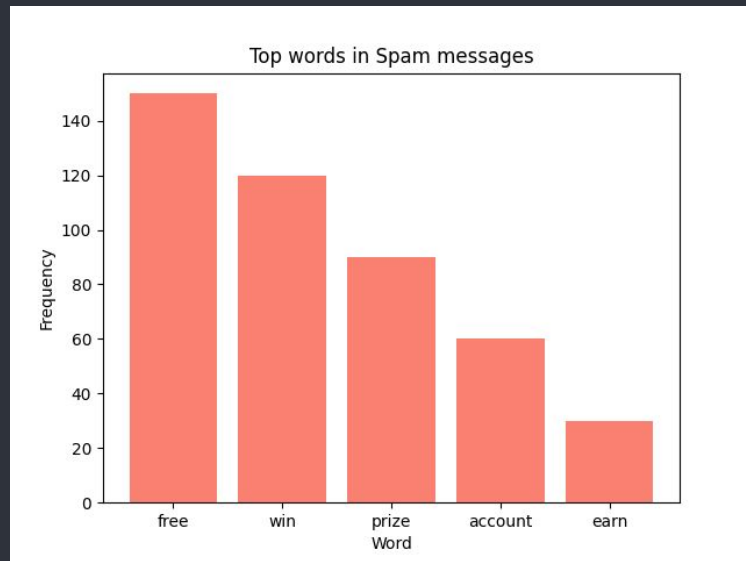


Top words in Normal messages

1
2
3
4 **Bar chart of**
5 **'Spam Messages'** {
6
7
8
9      < Top words in spam messages.
10     Here we see "free", "win",
11     "prize", etc., which are typical
12     in spam offers >
13
14 }



Top words in Spam messages

```
1
2
3    Distribution of
4    'ham vs spam'
5
6    messages {
7
8            <The dataset is heavily
9
10           imbalanced (many more ham than
11           spam) as shown by the bar
12           chart.>
13
14   }
```
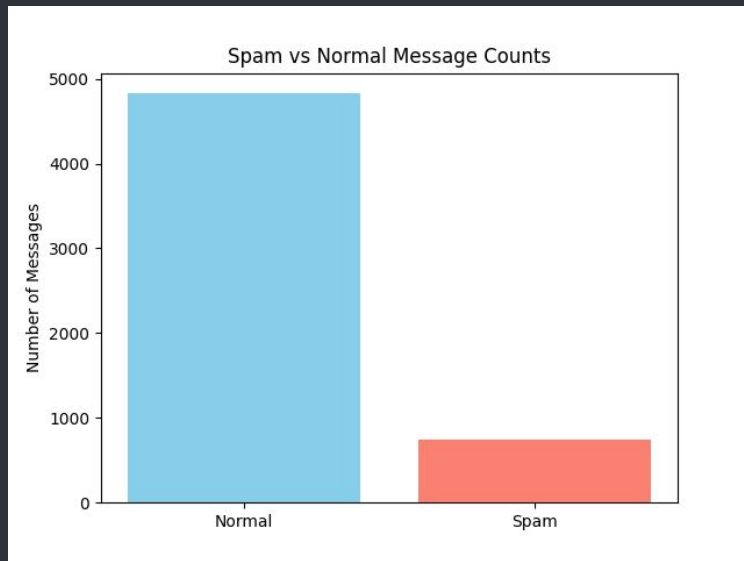


Spam vs Normal Message Counts

```
1
2
3    4,827 ham and 747 spam
4
5    messages {
6
7        < as documented in the Kaggle SMS Spam
8        Collection >
9
10   }
11
12
13
14
```

Created by Evan Tan

```
1
2    03 {
3
4
5
6        [Code for Model]
7
8            < Choice of Vectorizer/confusion
9            matrix >
10
11
12
13   }
14
```

```
 1    Libraries {
 2
 3        Numpy/Pandas                    sklearn
 4
 5        Reading of dataset             Usage of algorithm
 6
 7          seaborn/matplotlib
 8
 9        Data visualisation
10
11
12
13
14    } Each import line brings in functions/classes needed later.
```

Model (1); {

```python
#vectorize text into numerical values. Normal text => 0, spam => 1
#Binary classification problem. Preprocessing step in supervised learning
spam_df['spam'] = spam_df['Category'].map({'ham': 0, 'spam': 1})


#START OF MODELLING
#train/test split, splits the dataset into training and testing subsets
X_train, X_test, y_train, y_test = train_test_split(spam_df.Message,spam_df.spam, random_state = 23, test_size=0.3)
'''
```

Mapping 'ham'→0 and 'spam'→1 prepares the data for modeling, making it a binary classification problem. This label encoding is a common preprocessing step in supervised learning. Train-test split: The script splits the dataset into training and testing subsets:

}

Created by Evan Tan

## Model (2); {

```python
#vectorize text into numerical values. Normal text => 0, spam => 1
#Binary classification problem. Preprocessing step in supervised learning
spam_df['spam'] = spam_df['Category'].map({'ham': 0, 'spam': 1})


#START OF MODELLING
#train/test split, splits the dataset into training and testing subsets
X_train, X_test, y_train, y_test = train_test_split(spam_df.Message,spam_df.spam, random_state = 23, test_size=0.3)
'''
```

The train_test_split function randomly partitions the data into
70% training and 30% test sets. This separation allows the
model's performance to be evaluated on unseen data. The
random_state ensures reproducibility of the split.

}

# Model (3); {

```python
#Store word count as a matrix, text must be converted to numeric features
#CountVectorizer tokenizes the text and builds a vocabulary of all words
#Then transforms each message into a vector of word counts
cv = CountVectorizer()
X_train_count = cv.fit_transform(X_train.values)

#Training of model
#Multinomial Naive Bayes classifier is then trained on the vectorized features
model = MultinomialNB()
model.fit(X_train_count, y_train)
```

CountVectorizer tokenizes the text and builds a vocabulary of all
words then transforms each message into a vector of word counts
The fit_transform on training data learns the vocabulary and
produces a sparse count matrix.

}

Created by Evan Tan

# Model (4); {

```
1
2
3      #Store word count as a matrix, text must be converted to numeric features
       #CountVectorizer tokenizes the text and builds a vocabulary of all words
4      #Then transforms each message into a vector of word counts
       cv = CountVectorizer()
5      X_train_count = cv.fit_transform(X_train.values)
6
       #Training of model
7      #Multinomial Naive Bayes classifier is then trained on the vectorized features
8      model = MultinomialNB()
       model.fit(X_train_count, y_train)
9
```

10      Training the model: A Multinomial Naive Bayes classifier is then
11      trained on the vectorized features
        MultinomialNB is well-suited for discrete word count features. It
12      learns the probability of each word given the class (ham or
13      spam), which allows it to classify new messages. After this line,
14 }    nb is a trained model ready to predict labels.

# Prediction function{

```python
#turning test into a function for reusability
trained_model = model
trained_vectorizer = cv
def predict_spam(messages, model=trained_model, vectorizer=trained_vectorizer):
    message_count = vectorizer.transform(messages)
    prediction = model.predict(message_count)
    return prediction
```

This function takes a string msg, vectorizes it using the same CountVectorizer, and applies nb.predict to output 0/1. It returns "spam" or "ham" for convenience.
For example, calling predict_spam("Congratulations, you have won!") would print "spam" if the model predicts 1.

}

1
2
3
4
5
6
7
8
9
10
11
12
13
14

1  **Alternative 'vectorizer' {**
2
3          TF-IDF weighting
4
5
6          *  TF-IDF downweights very common words
7             across all messages and upweights
8             distinctive words. This often yields
             better performance by capturing term
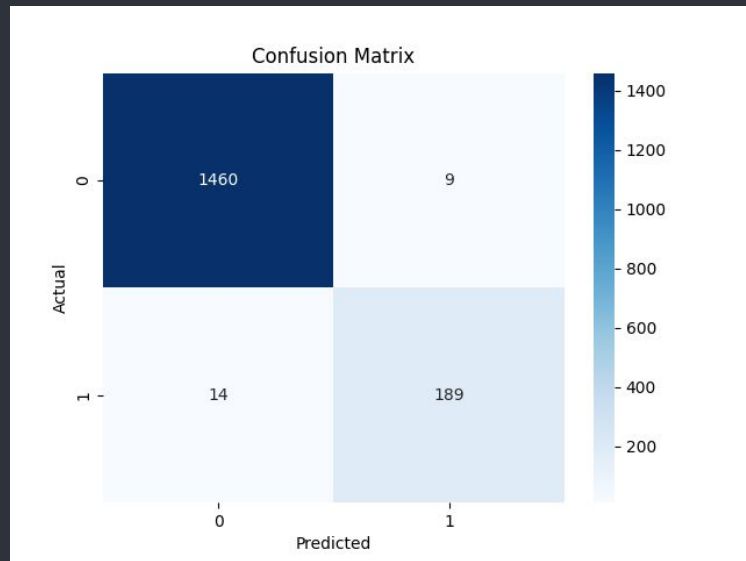9             significance rather than just frequency
10
11
12
13
14  **}**

```python
tfidf = TfidfVectorizer()
X_train_tfidf = tfidf.fit_transform(X_train.values)
```

Created by Evan Tan

1

2

3 # Comparison

4 'vectorizers'

5

6 (CountVectorizer){

7

8 Each row is the actual class,

9 each column the predicted class.

10 In this example, the model correctly identifies most ham

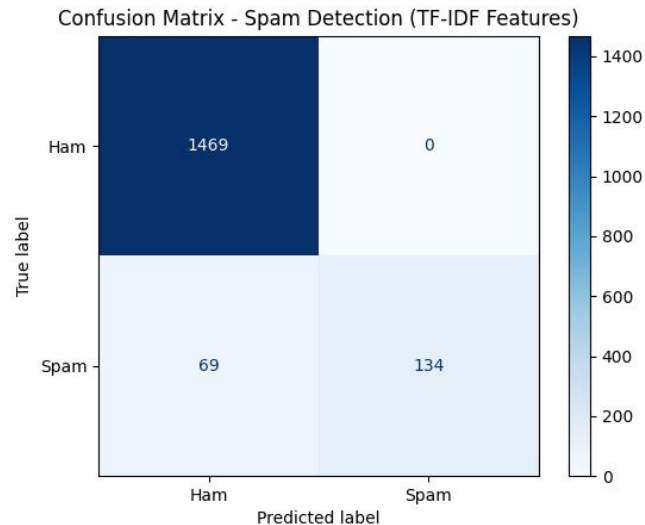11 and spam messages, but it misses
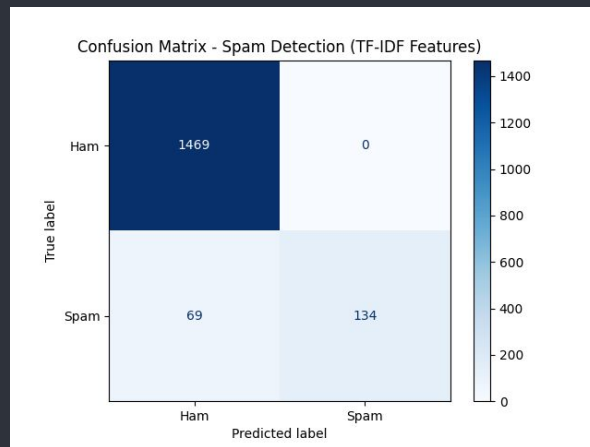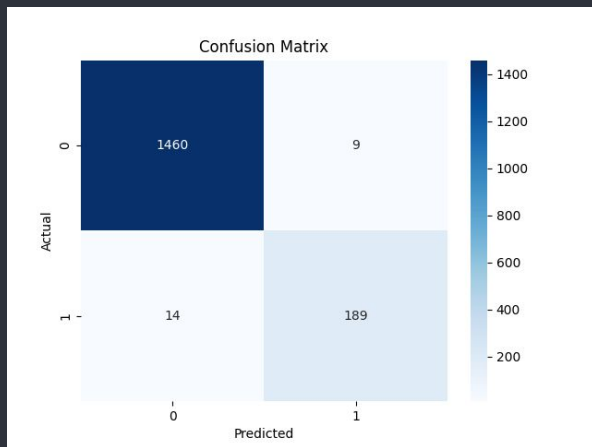
12 one spam (one false negative).

13

14 }



Confusion Matrix

1
2
3 # Comparison
4 'vectorizers'
5
6 (TF-IDF){
7
8     It perfectly classifies all ham
      messages (no false positives), but it
9     misses some spam messages —
      specifically, 69 spam messages are
10    incorrectly predicted as ham (false
      negatives). This indicates that while
11    the model is conservative in labeling
12    messages as spam (high precision), it
      could benefit from improvements in
13    recall to catch more actual spam.
14 }



Confusion Matrix - Spam Detection (TF-IDF Features)

Created by Evan Tan

1  ## conclusion{

2

3

4

5

6

7

8

9

10

11

12

13

14 }



Confusion Matrix



Confusion Matrix - Spam Detection (TF-IDF Features)

`<MultinomialNB assumes discrete counts of features. It's naturally suited to raw counts rather than continuous TF-IDF weights.`

`TF-IDF transforms the features into floats, which can violate the Naive Bayes assumptions, reducing accuracy.>`

Created by Evan Tan

```
 1    Resources {
 2
 3              Links:
 4
 5          *   scikit
 6          *   tokenizer
 7
 8              Data set:
 9          *   spam ham data set
10
11
12
13
14    }
```