



# Text Classification

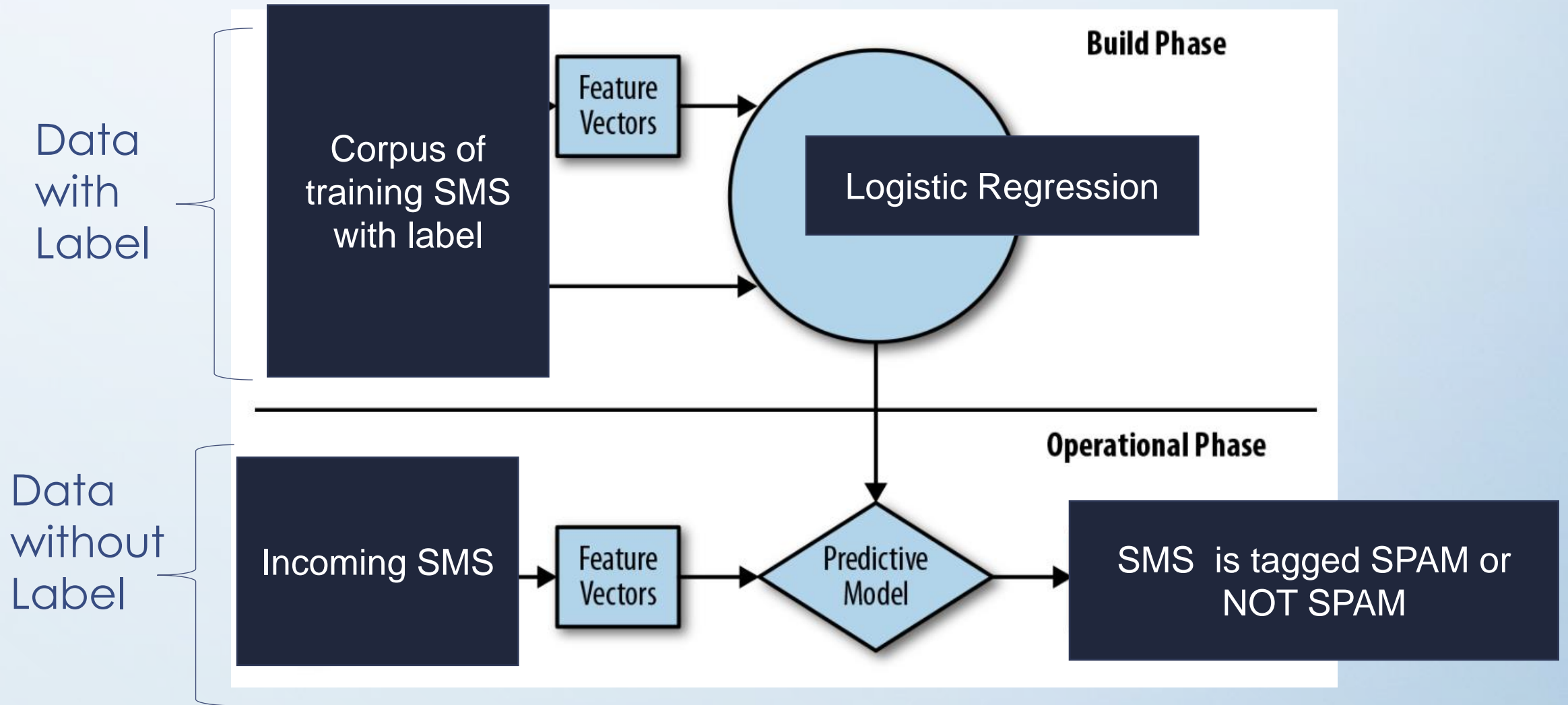
# Learning Objectives

- Describe each component in Text Classification System
- Calculate classification evaluation metric
- Build and use Logistic Regression model in Text Classification System
- Build and use Naive Bayes model in Text Classification System

# Logistic Regression



# Example



# Exercise 1: Build Logistic Regression model to predict ham or spam for SMS message

## Step 1:

Read the SMSSpamCollection.txt and pre-process the data

- Remove all the words with numbers and pure numbers. e.g. 21<sup>st</sup>, 2005
- Remove all the punctuation. e.g. !, ?
- Convert capital letter to small letter

	label	new_text
0	ham	ok lar joking wif u oni
1	spam	free entry in a wkly comp to win fa cup final tkts may text fa to to receive entry questionstd txt ratetcs apply s
2	ham	u dun say so early hor u c already then say
3	ham	nah i dont think he goes to usf he lives around here though
4	spam	freemsg hey there darling its been weeks now and no word back id like some fun you up for it still tb ok xxx std chgs to send £ to rcv

# Exercise 1: Build Logistic Regression model to predict ham or spam for SMS message

## Step 2:

Split the dataset into training set and test set

- Test dataset = 30% of observation and Training dataset = 70% of observation
- Random state =42, so we all get the same random train and test split

```
The size of original dataset: (5571,)
The size of training dataset: (3899,)
The size of test dataset: (1672,)
```



# Exercise 1: Build Logistic Regression model to predict ham or spam for SMS message

## Step 3:

Convert the text to vectors using count vectorizer

```
The dimensions of the training set: (3899, 6663)
```

```
The dimensions of the test set: (1672, 6663)
```

```
The features:
```

```
['aa' 'aah' 'aaoooooright' ... 'zoom' 'zouk' 'ülll']
```

# Exercise 1: Build Logistic Regression model to predict ham or spam for SMS message

## Step 4:

Fit Logistic Regression model on training data and apply the fitted model to test data. Predict the test data.

```
1 print(list(y_test[:10]))  
2 print(list(y_pred_cv[:10]))
```

```
['ham', 'spam', 'ham', 'spam', 'ham', 'ham', 'ham', 'ham', 'ham', 'ham']  
['ham', 'spam', 'ham', 'spam', 'ham', 'ham', 'ham', 'ham', 'ham', 'ham']
```



# Exercise 1: Build Logistic Regression model to predict ham or spam for SMS message

## Step 5:

Evaluate the mode. Decide how good the model is by calculating various metrics

- Confusion Metrix
- Precision
- Recall
- F1-score
- Accuracy

# Exercise 1: Build Logistic Regression model to predict ham or spam for SMS message

## Step 5 (cont.):

Evaluate the model. Decide how good the model is by calculating various metrics

```
array([[1451,    2],  
       [  32,  187]])
```

	precision	recall	f1-score	support
ham	0.98	1.00	0.99	1453
spam	0.99	0.85	0.92	219
accuracy			0.98	1672
macro avg	0.98	0.93	0.95	1672
weighted avg	0.98	0.98	0.98	1672

# Exercise 1: Build Logistic Regression model to predict ham or spam for SMS message

## Step 6:

Save the model and counter vectorizer

- Save count vectorizer so that we can retain the vocabulary list and other setting to get the features. New text will have to be transformed through the count vectorizer
- Save the model for predict the new text
- The model name is: `lr-2022-<MM>-<DD>.pkl`
- The vectorizer name: `countvectoriser-2022-<MM>-<DD>.pkl`

# Exercise 1 Answers:

**Step 1:** Read the SMSSpamCollection.txt and pre-process the data

```
1 import pandas as pd
2 import os
3 from google.colab import drive
4 drive.mount('/content/drive')
5
6 data_path = ['drive', 'MyDrive', 'ITB241', 'data']
```

Mounted at /content/drive

```
1 pd.set_option('display.max_colwidth', None) # Set display options so that the text aren't cut off
2
3 filepath = os.sep.join(data_path + ['SMSSpamCollection.txt'])
4 df=pd.read_csv(filepath,sep='\t')
5 df.columns=['label','text']
6 df.head()
```

```
1 import re
2 import string
3
4 pattern_alphanumeric="\w*\d\w*"
5 pattern_punctuation="["+re.escape(string.punctuation)+"]"
6
7 df['new_text']="empty"
8
9 for ind in df.index:
10     temp=re.sub(pattern_alphanumeric, '',df['text'][ind])
11     df['new_text'][ind]=re.sub(pattern_punctuation, '',temp).lower()
12 df_cleaned=df[['label', 'new_text']].copy()
13 df_cleaned
```

# Exercise 1 Answers

Step 2: Split the dataset into training set and test set

```
1 # split the data into inputs and outputs
2 X = df_cleaned['new_text']
3 Y = df_cleaned['label']
```

```
1 # split the data into a training and test set
2 from sklearn.model_selection import train_test_split
3
4 # test size = 30% of observations, which means training size = 70% of observations
5 # random state = 42, so we all get the same random train / test split
6 X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.3, random_state=42)
```

```
1 print("The size of original dataset:", X.shape)
2 print("The size of training dataset:", X_train.shape)
3 print("The size of test dataset:", X_test.shape)
```

# Exercise 1 Answers

## Step 3: Convert the text to vectors using count vectorizer

```
1  from sklearn.feature_extraction.text import CountVectorizer
2
3  cv = CountVectorizer(stop_words='english', ngram_range=(1,1))
4
5  X_train_cv = cv.fit_transform(X_train)
6  X_test_cv  = cv.transform(X_test)
7
8  # print the dimensions of the training set (text messages, terms)
9  print("The dimensions of the training set:", X_train_cv.toarray().shape)
10 print("The dimensions of the test set:", X_test_cv.toarray().shape)
11
12 # print the features that has been create as a result of fit_transform()
13 print("The features:\n", cv.get_feature_names_out())
```



# Exercise 1 Answers

## Step 4:

Fit Logistic Regression model on training data and apply the fitted model to test data. Predict the test data.

```
1  # Use a logistic regression model
2  from sklearn.linear_model import LogisticRegression
3  lr = LogisticRegression(solver='lbfgs')
4
5  # Train the model
6  lr.fit(X_train_cv, y_train)
7
8  # Take the model that was trained on the X_train_cv data and apply it to the X_test_cv data
9  y_pred_cv = lr.predict(X_test_cv)
10 # The output is all of the predictions
11 y_pred_cv
```

```
array(['ham', 'spam', 'ham', ..., 'ham', 'ham', 'spam'], dtype=object)
```

```
1  print(list(y_test[:10]))
2  print(list(y_pred_cv[:10]))
```

```
['ham', 'spam', 'ham', 'spam', 'ham', 'ham', 'ham', 'ham', 'ham', 'ham']
['ham', 'spam', 'ham', 'spam', 'ham', 'ham', 'ham', 'ham', 'ham', 'ham']
```

# Exercise 1 Answers

**Step 5:** Evaluate the model. Decide how good the model is by calculating various metrics

```
1  from sklearn.metrics import confusion_matrix
2
3  cm = confusion_matrix(y_test, y_pred_cv)
4  cm
```

```
from sklearn.metrics import classification_report
print(classification_report(y_test,y_pred_cv, target_names=['ham','spam']))
```

# Exercise 1 Answers

## Step 6: Save the model and counter vectorizer

```
1  # There are two objects that we need to save.
2
3  # First, the Count Vectorizer so that we can retained the vocabulary list and
4  # other settings to get the features. New text will have to be transformed through the
5  # Count Vectorizer.
6
7  # Second, the LR model. This will be used for prediction.
8
9  import pickle
10 from datetime import datetime
11
12 model_path = ['drive', 'MyDrive', 'ITB241', 'models']
13
14 time = datetime.now().strftime("%Y-%m-%d")
15 filename = 'lr-{}.pkl'.format(time)
16 tempList=[]
17 tempList.append(filename)
18 path1 = os.sep.join(model_path + tempList)
19
20 filename = 'countvectoriser-{}.pkl'.format(time)
21 tempList=[]
22 tempList.append(filename)
23 path2 = os.sep.join(model_path + tempList)
24
25 with open(path1, 'wb') as f1:
26     pickle.dump(lr, f1)
27
28 with open(path2, 'wb') as f2:
29     pickle.dump(cv, f2)
30
```

# Exercise 2: Use Logistic Regression model to predict ham or spam for SMS message

## Step 1:

Load in the regression model that was saved during the modelling stage

## Step 2:

Load in the vectorizer that was used to encode the training set

# Exercise 2: Use Logistic Regression model to predict ham or spam for SMS message

## Step 3:

Pre-process the new text: Clean the input data in the SAME way as it was done during modelling

- Create function preprocess(text) to clean the data

```
1 new_text="SIX chances to win CASH! From 100 to 20,000 pounds txt> CSH11 and send to 87575. Cost 150p/day, 6days, 16+ \
2 TsandCs apply Reply HL 4 info"
3 new_text=preprocess(new_text)
4 new_text
```

```
'six chances to win cash from to pounds txt and send to cost day tsandcs apply reply hl info'
```

# Exercise 2: Use Logistic Regression model to predict ham or spam for SMS message

## Step 4:

Numerically encode the input : Convert the text to vectors using the previous Vectorizer

- Create function `encode_text_to_vector(cv, text)` to encode the new text

```
1 new_text_vector=encode_text_to_vector(trained_cv,new_text)
2 print(new_text_vector)
```

```
(0, 248)      1
(0, 859)      1
(0, 900)      1
(0, 1187)     1
(0, 1346)     1
(0, 2521)     1
(0, 2746)     1
(0, 4335)     1
(0, 4685)     1
(0, 4972)     1
(0, 5971)     1
(0, 5997)     1
(0, 6400)     1
```



# Exercise 2: Use Logistic Regression model to predict ham or spam for SMS message

## Step 5:

Predict the label

```
1 predicted_label = (model.predict(new_text_vector))[0]
2 print ("The new text:\n",new_text)
3 print("Predicted label is:\n", predicted_label)
```

The new text:

```
SIX chances to win CASH! From 100 to 20,000 pounds txt> CSH11 and send to 87575. Cost 150p/day, 6days, 16+ TsandCs apply Reply HL 4 info
Predicted label is:
spam
```

# Exercise 2 Answers

## Step 1:

Load in the regression model that was saved during the modelling stage

```
1  import os
2  import pickle
3  from google.colab import drive
4
5  drive.mount('/content/drive')
6
7  model_path = ['drive', 'MyDrive', 'ITB241', 'models']
8  filename=['lr-2022-10-10.pkl']
9
10 path1 = os.sep.join(model_path + filename)
11 with open(path1, 'rb') as f:
12     model = pickle.load(f)
```

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force\_remount=True).

```
1  model
```

```
LogisticRegression()
```

# Exercise 2 Answers

## Step 2:

Load in the vectorizer that was used to encode the training set

```
1 model_path = ['drive', 'MyDrive', 'ITB241', 'models']
2 filename=['countvectoriser-2022-10-10.pkl']
3
4 path2 = os.sep.join(model_path + filename)
5 with open(path2, 'rb') as f:
6     trained_cv= pickle.load(f)
```

```
1 trained_cv
```

```
CountVectorizer(stop_words='english')
```

# Exercise 2 Answers

## Step 3: Pre-process the new text

```
1 import re
2 import string
3
4 def preprocess(text):
5     pattern_alphanumeric="\w*\d\w*"
6     pattern_punctuation="["+re.escape(string.punctuation)+"]"
7
8     text=re.sub(pattern_alphanumeric, '',text)
9     text=re.sub(pattern_punctuation, '',text).lower()
10
11     return text
```

```
1 new_text="SIX chances to win CASH! From 100 to 20,000 pounds txt> CSH11 and send to 87575. Cost 150p/day, 6days, 16+ \
2 TsandCs apply Reply HL 4 info"
3 new_text_processed=preprocess(new_text)
4 new_text_processed
```

```
'six chances to win cash from 100 to 20000 pounds txt and send to 87575 cost day tsandcs apply reply hl info'
```

# Exercise 2 Answer

## Step 4:

Numerically encode the input : Convert the text to vectors using the previous Vectorizer

```
1  def encode_text_to_vector(cv, text):  
2      text_vector = cv.transform( [text] )  
3      return text_vector
```

```
1  new_text_vector=encode_text_to_vector(trained_cv,new_text_processed)  
2  print(new_text_vector)
```

# Exercise 2 Answers

## Step 5: Predict the label

```
1 predicted_label = (model.predict(new_text_vector))[0]
2 print ("The new text:\n",new_text)
3 print("Predicted label is:\n", predicted_label)
```

The new text:

SIX chances to win CASH! From 100 to 20,000 pounds txt> CSH11 and send to 87575. Cost 150p/day, 6days, 16+ TsandCs apply Reply HL 4 info

Predicted label is:

spam

1



# Test other SMS messages

Message	Predicted Label
I'm gonna be home soon and i don't want to talk about this stuff anymore tonight, k? I've cried enough today.	ham
I've been searching for the right words to thank you for this breather. I promise i wont take your help for granted and will fulfil my promise. You have been wonderful and a blessing at all times.	ham
Oh k...i'm watching here:)	ham
Eh u remember how 2 spell his name... Yes i did. He v naughty	ham
Fine if that's the way u feel. That's the way its gota b	ham
Is that seriously how you spell his name?	ham

# Test other SMS messages

Message	Predicted Label
SIX chances to win CASH! From 100 to 20,000 pounds txt> CSH11 a nd send to 87575. Cost 150p/day, 6days, 16+ TsandCs apply Reply HL 4 info	spam
URGENT! You have won a 1 week FREE membership in our £100,000 Prize Jackpot! Txt the word: CLAIM to No: 81010 T&C www.dbuk.net LCCLTD POBOX 4403LDNW1A7RW18	spam
XXXMobileMovieClub: To use your credit, click the WAP link in the next txt message or click here>> <a href="http://wap.xxxmobilemovieclub.com?n=QJKGIGHJJGCBL">http://wap.xxxmobilemovieclub.com?n=QJKGIGHJJGCBL</a>	spam
England v Macedonia - dont miss the goals/team news. Txt ur national team to 87077 eg ENGLAND to 87077 Try:WALES, SCOTLAND 4txt/ú1.20 POBOXox365 04W45WQ 16+	spam

# Naïve Bayes



# Exercise 3: Training dataset

Text	Label
A great game	sports
The election was over	not sports
Very clean match	sports
It was a close election	not sports
A clean but forgettable game	Sports
A very close game	??

What is the label for a new text “A very close game”?  
Would it be “Sports” or “Not Sports”?

# Naïve Bayes Approach

Text	Label
a great game	sports
the election was over	not sports
very clean match	sports
it was a close election	not sports
a clean but forgettable game	sports

## Find:

$P(\text{sports} \mid \text{"a very close game"})$

$P(\text{not sports} \mid \text{"a very close game"})$

- If  $P(\text{sports} \mid \text{"a very close game"})$  is the larger value, then the label is sports
- If  $P(\text{not sports} \mid \text{"a very close game"})$  is the larger value, then the label is not sports

# Laplace Smoothing

Bayes Theorem

$$\frac{P(\text{sports} \mid \text{"A very close game"})}{P(\text{"a very close game"})} = \frac{P(\text{"a very close game"} \mid \text{sports}) \times P(\text{sports})}{P(\text{"a very close game"})}$$

Naïve Bayes

$$\frac{P(\text{"a"} \mid \text{sports}) \times P(\text{"very"} \mid \text{sports}) \times P(\text{"close"} \mid \text{sports}) \times P(\text{"game"} \mid \text{sports}) \times P(\text{sports})}{c}$$

0 + 1  
11 + 14

easier to calculate by breaking up

Text	Label
a great game	sports
The election was over	not sports
very clean match	sports
it was a close election	not sports
a clean but forgettable game	sports

The word "close" isn't present in any *sports* text. This means  $P(\text{"close"} \mid \text{sports})$  is 0 and the end result will be 0 !

- The above issue can be eliminated with Laplace smoothing, where every count is added by 1
- The possible number of words is also added to the divisor, and the division will not be more than 1

In this case, the set of possible words are:

['a', 'great', 'very', 'over', 'it', 'but', 'game', 'match', 'clean', 'election', 'close', 'the', 'was', 'forgettable'].

always add 1

The possible number of words is 14. Therefore:

$$P(\text{"close"} \mid \text{sports}) = (0+1) / (11+14) = 0.04 \quad P(\text{"game"} \mid \text{sports}) = (2+1) / (11+14) = 0.12, \text{ and etc}$$

total features for entire dataset (NON-REPEATED) --> 14



# Naïve Bayes

$P(\text{sports} \mid \text{"A very close game"})$

Naïve Bayes

$$P(\text{"a"} \mid \text{sports}) \times P(\text{"very"} \mid \text{sports}) \times P(\text{"close"} \mid \text{sports}) \times P(\text{"game"} \mid \text{sports}) \times P(\text{sports})$$

$$P(\text{"a"}) \times P(\text{"very"}) \times P(\text{"close"}) \times P(\text{"game"}) \quad \leftarrow \text{can ignore denominator}$$

$P(\text{not sports} \mid \text{"A very close game"})$

Naïve Bayes

$$P(\text{"A"} \mid \text{not sports}) \times P(\text{"very"} \mid \text{not sports}) \times P(\text{"close"} \mid \text{not sports}) \times P(\text{"game"} \mid \text{not sports}) \times P(\text{not sports})$$

$$P(\text{"A"}) \times P(\text{"very"}) \times P(\text{"close"}) \times P(\text{"game"}) \quad \leftarrow \text{can ignore denominator}$$

As the dominator are the same, we can **compare only the numerators** to decide which probability is higher

$$\begin{aligned} P(\text{"A"} \mid \text{sports}) &= (2+1) / (11+14) = 0.12 \\ P(\text{"very"} \mid \text{sports}) &= (1+1) / (11+14) = 0.08 \\ P(\text{"close"} \mid \text{sports}) &= (0+1) / (11+14) = 0.04 \\ P(\text{"game"} \mid \text{sports}) &= (2+1) / (11+14) = 0.12 \\ P(\text{sports}) &= 3/5 \end{aligned}$$

$$\begin{aligned} P(\text{"A"} \mid \text{not sports}) &= (1+1) / (9+14) = 0.087 \\ P(\text{"very"} \mid \text{not sports}) &= (0+1) / (9+14) = 0.043 \\ P(\text{"close"} \mid \text{not sports}) &= (1+1) / (9+14) = 0.087 \\ P(\text{"game"} \mid \text{not sports}) &= (0+1) / (9+14) = 0.043 \\ P(\text{not sports}) &= 2/5 \end{aligned}$$

calculate probability P based on the table

Text	Label
a great game	sports
The election was over	not sports
very clean match	sports
it was a close election	not sports
a clean but forgettable game	sports

# Exercise 3 answer

$$P(\text{"a"} \mid \text{sports}) \times P(\text{"very"} \mid \text{sports}) \times P(\text{"close"} \mid \text{sports}) \times P(\text{"game"} \mid \text{sports}) \times P(\text{sports})$$

$$= 0.12 \times 0.08 \times 0.04 \times 0.12 \times 0.6 = \underline{2.76 \times 10^{-5}}$$

$$P(\text{"a"} \mid \text{not sports}) \times P(\text{"very"} \mid \text{not sports}) \times P(\text{"close"} \mid \text{not sports}) \times P(\text{"game"} \mid \text{not sports}) \times P(\text{not sports})$$

$$= 0.087 \times 0.043 \times 0.087 \times 0.043 \times 0.4 = 0.56 \times 10^{-5}$$

Hence,  $P(\text{sports} \mid \text{"a very close game"})$  gives a higher probability, suggesting that the sentence belongs to the sports category.

## Exercise 4:

Can you implement the **Build Phase and Operational Phase** using Naïve Bayes to predict ham or spam for SMS messages

Hint: Modify the code from “Logistic Regression”. A very small changes only!

# Exercise 4 Answers:

Step 4: Fit Naive Bayes model on training data and apply the fitted model to test data

```
[ ] 1  # Use a Naive Bayes model
    2  from sklearn.naive_bayes import MultinomialNB
    3  nb = MultinomialNB()
    4
    5  # Train the model
    6  nb.fit(X_train_cv, y_train)
    7
    8  # Take the model that was trained on the X_train_cv data and apply it to the X_test_cv data
    9  y_pred_cv = nb.predict(X_test_cv)
   10  # The output is all of the predictions
   11  y_pred_cv
```

- Only shows the difference from the code Exercise 1 and Exercise 2
- Change all the variable **lr** to **nb** in the subsequent code