



COMP30810

Intro to Text Analytics

Dr. Binh Thanh Le

thanhbinh.le@ucd.ie

Insight Centre for Data Analytics

School of Computer Science

University College Dublin

Today Goals:

- Understanding the main concept of Naive Bayes classifier
- Applying the classifiers to Text Mining

What is Naive Bayes theorem?

In probability theory and statistics, **Bayes' theorem** (alternatively **Bayes' law** or **Bayes' rule**) describes the probability of an event, based on prior knowledge of conditions that might be related to the event.

GAUSSIAN
NAIVE BAYES
CLASSIFIER

"Gaussian" because this is a normal distribution

This is our prior belief

Likelihood

Class Prior Prob.

Posterior Prob. $p(\text{class} | \text{data}) = \frac{p(\text{data} | \text{class}) \times p(\text{class})}{p(\text{data})}$

Predictor Prior Prob.

We don't calculate this in naive bayes classifiers

ChrisAlbon

Example

Weather	Play
Sunny	No
Overcast	Yes
Rainy	Yes
Sunny	Yes
Sunny	Yes
Overcast	Yes
Rainy	No
Rainy	No
Sunny	Yes
Rainy	Yes
Sunny	No
Overcast	Yes
Overcast	Yes
Rainy	No

We have a training data set of weather and corresponding target variable 'Play' (suggesting possibilities of playing). Now, we need to classify whether players will play or not based on weather condition.

Step 1: Convert the data set into a frequency table

Frequency Table		
Weather	No	Yes
Overcast		4
Rainy	3	2
Sunny	2	3
Grand Total	5	9

Example

Weather	Play
Sunny	No
Overcast	Yes
Rainy	Yes
Sunny	Yes
Sunny	Yes
Overcast	Yes
Rainy	No
Rainy	No
Sunny	Yes
Rainy	Yes
Sunny	No
Overcast	Yes
Overcast	Yes
Rainy	No

We have a training data set of weather and corresponding target variable 'Play' (suggesting possibilities of playing). Now, we need to classify whether players will play or not based on weather condition.

Step 2: Create Likelihood table by finding the probabilities like **Overcast probability = 0.29** and **probability of playing is 0.64**.

Likelihood table				
Weather	No	Yes		
Overcast		4	=4/14	0.29
Rainy	3	2	=5/14	0.36
Sunny	2	3	=5/14	0.36
All	5	9		
	=5/14	=9/14		
	0.36	0.64		

Example

Weather	Play
Sunny	No
Overcast	Yes
Rainy	Yes
Sunny	Yes
Sunny	Yes
Overcast	Yes
Rainy	No
Rainy	No
Sunny	Yes
Rainy	Yes
Sunny	No
Overcast	Yes
Overcast	Yes
Rainy	No

We have a training data set of weather and corresponding target variable 'Play' (suggesting possibilities of playing). Now, we need to classify whether players will play or not based on weather condition.

Step 3: Now, use Naive Bayesian equation to **calculate the posterior probability for each class**. The class with the highest posterior probability is the outcome of prediction.



Players will play if weather is sunny.
Is this statement is correct?

ANSWER

$$\begin{aligned} P(\text{Yes} \mid \text{Sunny}) &= P(\text{Sunny} \mid \text{Yes}) * P(\text{Yes}) / P(\text{Sunny}) \\ &= \left(\frac{3}{9} * \frac{9}{14}\right) / \left(\frac{5}{14}\right) = \mathbf{0.60} \end{aligned}$$

$$\begin{aligned} P(\text{No} \mid \text{Sunny}) &= P(\text{Sunny} \mid \text{No}) * P(\text{No}) / P(\text{Sunny}) \\ &= \left(\frac{2}{5} * \frac{5}{14}\right) / \left(\frac{5}{14}\right) = \mathbf{0.40} \end{aligned}$$

Classification Methods

Supervised Machine Learning

- Input:
 - a document d
 - a fixed set of classes $C = \{c_1, c_2, \dots, c_j\}$
 - A training set of m hand-labeled documents $(d_1, c_1), \dots, (d_m, c_m)$
- Output:
 - a learned classifier $\gamma: d \rightarrow c$

Positive or negative movie review?

- unbelievably disappointing
- Full of zany characters and richly applied satire, and some great plot twists
- this is the greatest screwball comedy ever filmed
- It was pathetic. The worst part about it was the boxing scenes.





Classification Methods

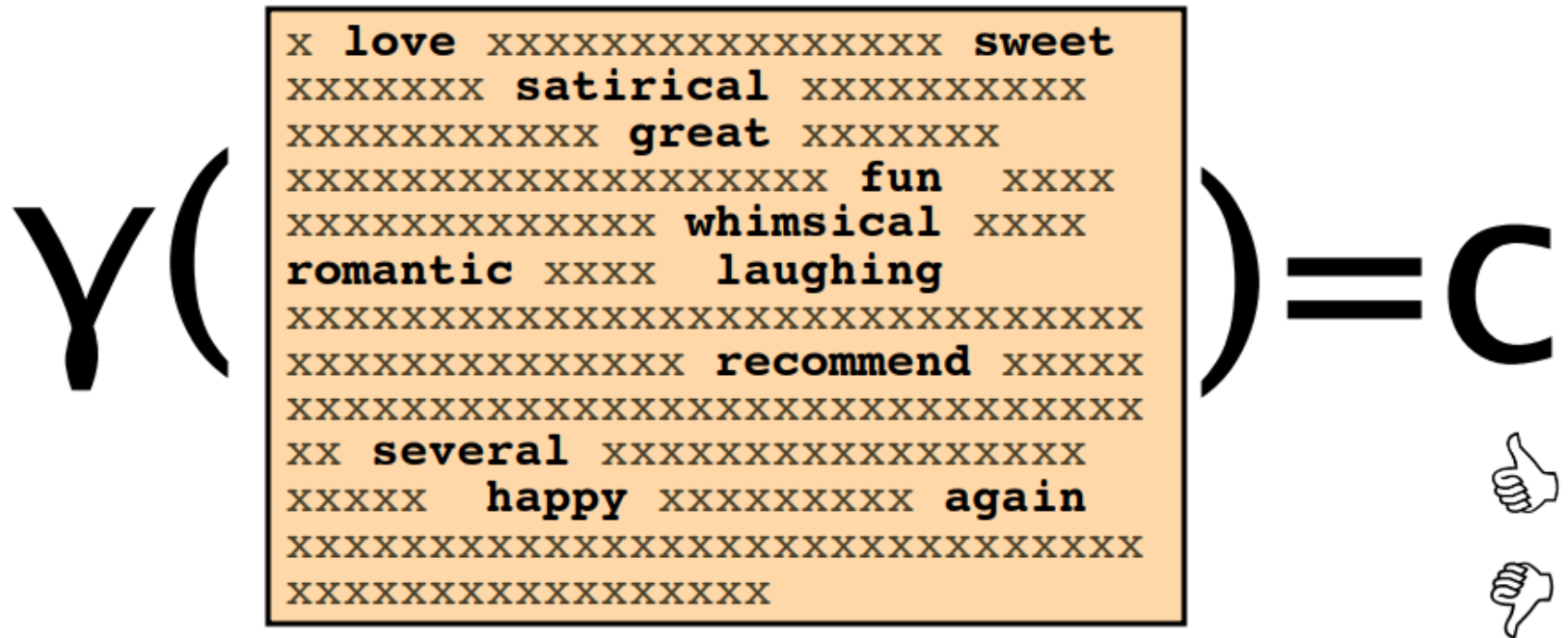
Supervised Machine Learning

- Any kind of classifier
 - Naïve Bayes
 - k-Nearest Neighbors
 - Logistic regression
 - Support-vector machines
 - ...

The bag of words representation

$$Y(\text{I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet.}) = C$$


The bag of words representation





The bag of words representation

$Y($

great	2
love	2
recommend	1
laugh	1
happy	1
...	...

$) = C$

Bayes' Rule Applied to Documents and Classes

- For a document ***d*** and a class ***c***

$$P(c | d) = \frac{P(d | c)P(c)}{P(d)}$$

Output

$$c_{MAP} = \operatorname{argmax}_{c \in C} P(c | d)$$

MAP is “maximum a posteriori” = most likely class

$$= \operatorname{argmax}_{c \in C} \frac{P(d | c)P(c)}{P(d)}$$

Bayes Rule

$$= \operatorname{argmax}_{c \in C} P(d | c)P(c)$$

Dropping the denominator

Bayes' Rule Applied to Documents and Classes

$$\begin{aligned}c_{MAP} &= \operatorname{argmax}_{c \in C} P(d | c)P(c) \\ &= \operatorname{argmax}_{c \in C} P(x_1, x_2, \dots, x_n | c)P(c)\end{aligned}$$

Document d
represented as
features
 $x_1 \dots x_n$

Independence Assumptions

- **Bag of Words assumption:** Assume position doesn't matter
- **Conditional Independence:** Assume the feature probabilities $P(x_i | c_j)$ are independent given the class c .

$$P(x_1, \dots, x_n | c) = P(x_1 | c) \cdot P(x_2 | c) \cdot P(x_3 | c) \cdot \dots \cdot P(x_n | c)$$

Bayes' Rule Applied to Documents and Classes

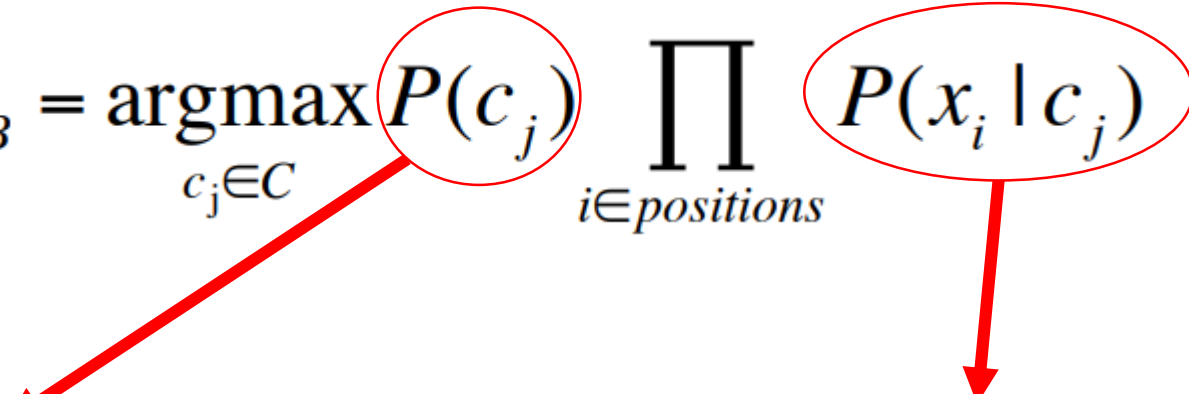
$$c_{MAP} = \operatorname{argmax}_{c \in C} P(x_1, x_2, \dots, x_n | c) P(c)$$

$$c_{NB} = \operatorname{argmax}_{c \in C} P(c_j) \prod_{x \in X} P(x | c)$$

$$c_{NB} = \operatorname{argmax}_{c_j \in C} P(c_j) \prod_{i \in \text{positions}} P(x_i | c_j)$$

positions \leftarrow all word positions in test document

Bayes' Rule Applied to Documents and Classes

$$c_{NB} = \operatorname{argmax}_{c_j \in C} P(c_j) \prod_{i \in \text{positions}} P(x_i | c_j)$$


$$\hat{P}(c_j) = \frac{\text{doccount}(C = c_j)}{N_{doc}}$$

$$\hat{P}(w_i | c_j) = \frac{\text{count}(w_i, c_j)}{\sum_{w \in V} \text{count}(w, c_j)}$$

**fraction of times word w_i appears
among all words in documents
of topic c_j**

Algorithm

- From training corpus, extract *Vocabulary (V)*

- Calculate $P(c_j)$ terms

- For each $c_j \in C$ do

$docs_j \leftarrow$ all docs with class $= c_j$

$$P(c_j) \leftarrow \frac{|docs_j|}{|\text{total \# documents}|}$$

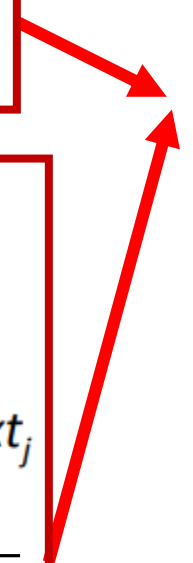
- Calculate $P(w_k | c_j)$ terms

- $Text_j \leftarrow$ single doc containing all $docs_j$

- For each word w_k in *V*

$n_k \leftarrow$ # of occurrences of w_k in $Text_j$

$$P(w_k | c_j) \leftarrow \frac{(n_k) + 1}{\sum_{w \in V} \text{count}(w | c_j) + |V|}$$





$P(c_1) * \prod P(w_k c_1)$	A
$P(c_2) * \prod P(w_k c_2)$	B


If $(A > B) \rightarrow$ Class 1
Else \rightarrow Class 2


Example

Train

 That is a **good** film, the content is really **great**, especially **great** ending.

 **Great** movie! **Excellent** ending, and **great** actions.

 **Great** film ever! But the ending is so **sad**.

 The ending is **boring**. The action is so **poor**. But the main actress's dress is so **great**!

Test

The film's ending is **great**. Main actor is super **great**, I love him.
Main actress is also **great**, but her action is a bit **boring** and **poor**.



	Doc	Words	Class
Train	1	Good Great Great	Positive
	2	Great Excellent Great	Positive
	3	Great Sad	Positive
	4	Boring Poor Great	Negative
Test	5	Great Great Great Boring Poor	?

Example

	Doc	Words	Class
Train	1	Good Great Great	Positive
	2	Great Excellent Great	Positive
	3	Great Sad	Positive
	4	Boring Poor Great	Negative
Test	5	Great Great Great Boring Poor	?

Prior Prob.

$$P\{+1\} = \frac{3}{4}$$
$$P\{-1\} = \frac{1}{4}$$



	Great	Good	Excellent	Sad	Boring	Poor	Class
Doc1	2	1	0	0	0	0	1
Doc2	2	0	1	0	0	0	1
Doc3	1	0	0	1	0	0	1
Doc4	1	0	0	0	1	1	-1

Example

	Great	Good	Excellent	Sad	Boring	Poor	Class
Doc1	2	1	0	0	0	0	1
Doc2	2	0	1	0	0	0	1
Doc3	1	0	0	1	0	0	1
Doc4	1	0	0	0	1	1	-1

Prior Probabilities:

$$P\{+1\} = \frac{3}{4}$$

$$P\{-1\} = \frac{1}{4}$$

Conditional Probabilities:

$$\begin{aligned}P(\text{Great} | 1) &= (5+1) / (8+6) = 6/14 \\P(\text{Good} | 1) &= (1+1) / (8+6) = 2/14 \\P(\text{Excellent} | 1) &= (1+1) / (8+6) = 2/14 \\P(\text{Sad} | 1) &= (1+1) / (8+6) = 2/14 \\P(\text{Boring} | 1) &= (0+1) / (8+6) = 1/14 \\P(\text{Poor} | 1) &= (0+1) / (8+6) = 1/14\end{aligned}$$

$$\begin{aligned}P(\text{Great} | -1) &= (1+1) / (3+6) = 2/9 \\P(\text{Good} | -1) &= (0+1) / (3+6) = 1/9 \\P(\text{Excellent} | -1) &= (0+1) / (3+6) = 1/9 \\P(\text{Sad} | -1) &= (0+1) / (3+6) = 1/9 \\P(\text{Boring} | -1) &= (1+1) / (3+6) = 2/9 \\P(\text{Poor} | -1) &= (1+1) / (3+6) = 2/9\end{aligned}$$

Example

$$P(w_k | c_j) \leftarrow \frac{(n_k) + 1}{\sum_{w \in V} \text{count}(w | c_j) + |V|}$$

Conditional Probabilities:

$$\begin{aligned} P(\text{Great} | 1) &= (5+1) / (8+6) = 6/14 \\ P(\text{Good} | 1) &= (1+1) / (8+6) = 2/14 \\ P(\text{Excellent} | 1) &= (1+1) / (8+6) = 2/14 \\ P(\text{Sad} | 1) &= (1+1) / (8+6) = 2/14 \\ P(\text{Boring} | 1) &= (0+1) / (8+6) = 1/14 \\ P(\text{Poor} | 1) &= (0+1) / (8+6) = 1/14 \end{aligned}$$

$$\begin{aligned} P(\text{Great} | -1) &= (1+1) / (3+6) = 2/9 \\ P(\text{Good} | -1) &= (0+1) / (3+6) = 1/9 \\ P(\text{Excellent} | -1) &= (0+1) / (3+6) = 1/9 \\ P(\text{Sad} | -1) &= (0+1) / (3+6) = 1/9 \\ P(\text{Boring} | -1) &= (1+1) / (3+6) = 2/9 \\ P(\text{Poor} | -1) &= (1+1) / (3+6) = 2/9 \end{aligned}$$

Prior Probabilities:

$$\begin{aligned} P\{+1\} &= \frac{3}{4} \\ P\{-1\} &= \frac{1}{4} \end{aligned}$$

Test	5	Great Great Great Boring Poor	?
------	---	-------------------------------	---

Choosing a class:

$$\begin{aligned} P(+1 | d5) &= (3/4) * (6/14) * (6/14) * (6/14) * (1/14) * (1/14) \\ &\approx 0.000301 \end{aligned}$$

$$\begin{aligned} P(-1 | d5) &= (1/4) * (2/9) * (2/9) * (2/9) * (2/9) * (2/9) \\ &\approx 0.000135 \end{aligned}$$

→ d5 is **Positive**

$$P(c_j) \prod_{i \in \text{positions}} P(x_i | c_j)$$

Python Implementation

```
t1 = 'good great great'
t2 = 'great excellent great'
t3 = 'great sad'
t4 = 'boring poor great'

t_test = ['great great great boring poor']
trainlabel = [1,1,1,-1]
trainset = pd.DataFrame({'data':[t1,t2,t3,t4], 'label':trainlabel})
```

◆ data ◆ label ◆

0	good great great	1
1	great excellent great	1
2	great sad	1
3	boring poor great	-1

```

from sklearn.feature_extraction.text import CountVectorizer
count_vect = CountVectorizer()
list_contents = trainset.data

count_matrix = count_vect.fit_transform(list_contents)
df_count = pd.DataFrame(count_matrix.toarray(), columns=[count_vect.get_feature_names()])
df_count['label'] = trainlabel
df_count

```

executed in 18ms, finished 15:00:54 2018-10-28

◆ boring ◆ excellent ◆ good ◆ great ◆ poor ◆ sad ◆ label ◆

0	0	0	1	2	0	0	1
1	0	1	0	2	0	0	1
2	0	0	0	1	0	1	1
3	1	0	0	1	1	0	-1

```

count_test_matrix = count_vect.transform(t_test)
df_count_test = pd.DataFrame(count_test_matrix.toarray(), columns=[count_vect.get_feature_names()])
df_count_test

```

executed in 14ms, finished 20:47:07 2018-10-28

◆ boring ◆ excellent ◆ good ◆ great ◆ poor ◆ sad ◆

0	1	0	0	3	1	0
---	---	---	---	---	---	---

Learning:

```
1 from sklearn.naive_bayes import MultinomialNB
2 model = MultinomialNB().fit(df_count.iloc[:, :6], df_count['label'])
```

Data

Label

Applying:

```
1 predicted = model.predict(df_count_test.values)
2 print('Class for test sample: ' + str(predicted))
```

executed in 4ms, finished 10:22:57 2018-10-31

Class for test sample: [1]

Probabilities:

```
1 probab_predict = model.predict_proba(df_count_test.values)
2 print('Probabilities for test sample: ' + str(probab_predict))
```

executed in 4ms, finished 10:50:06 2018-10-31

Probabilities for test sample: [[0.31024139 0.68975861]]

Class {-1} Class {+1}

Summary

- Naive Bayes is
 - Very Fast, low storage requirements
 - Robust to Irrelevant Features
 - Very good in domains with many equally important features
 - Optimal if the independence assumptions hold
 - A good dependable baseline for text classification