

Text classification using Naive Bayes

Naive Bayes

- Naive bayes is a probabilistic algorithm

Naive Bayes

- Naive bayes is a probabilistic algorithm
- Based on bayes theorem

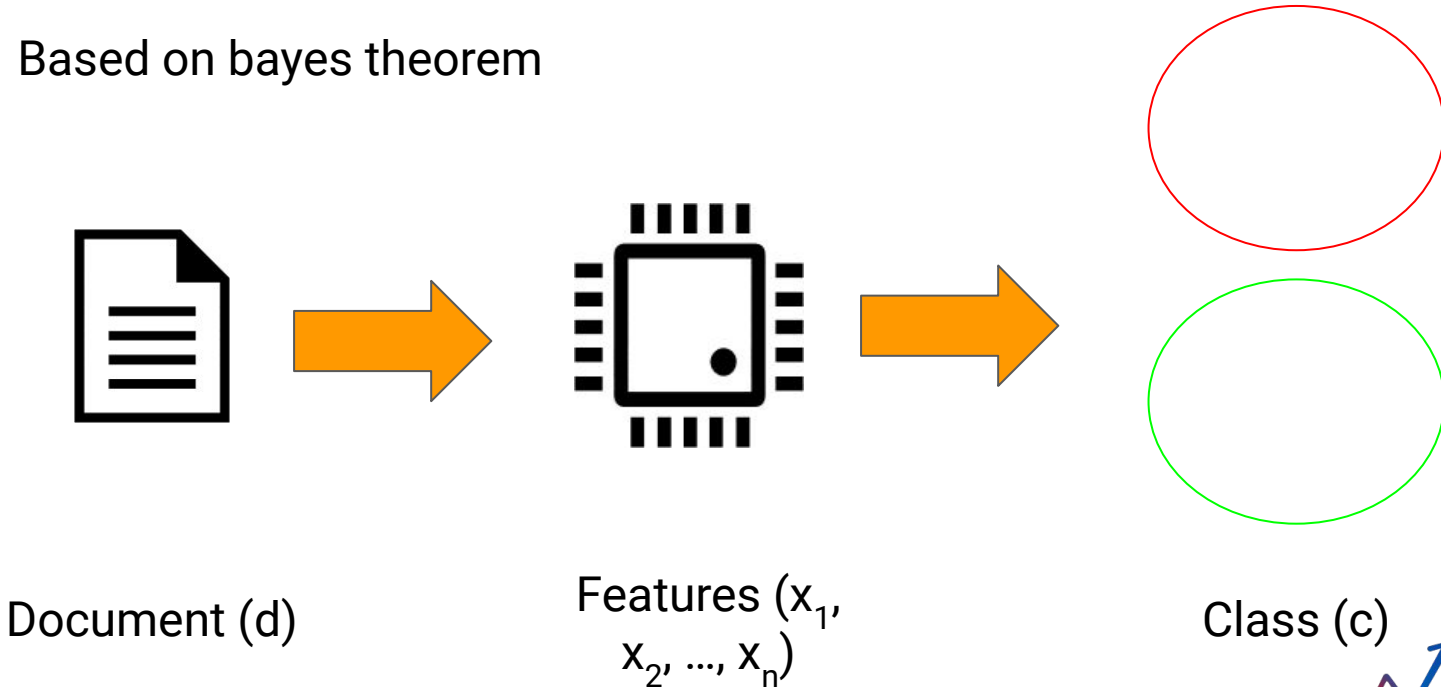
Naive Bayes

- Naive bayes is a probabilistic algorithm
- Based on bayes theorem

$$P(E_1 | E_2) = \frac{P(E_2 | E_1) * P(E_1)}{P(E_2)}$$

Naive Bayes for text classification

- Naive bayes is a probabilistic algorithm
- Based on bayes theorem



Naive Bayes for text classification

- Naive bayes is a probabilistic algorithm
- Based on bayes theorem

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

Naive Bayes for text classification

- Naive bayes is a probabilistic algorithm
- Based on bayes theorem

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

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

Naive Bayes for text classification

- Naive bayes is a probabilistic algorithm
- Based on bayes theorem

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

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

How likely the document will occur

Naive Bayes for text classification

- Naive bayes is a probabilistic algorithm
- Based on bayes theorem

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

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

Naive Bayes for text classification

- Naive bayes is a probabilistic algorithm
- Based on bayes theorem

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

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

Naive Bayes for text classification

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

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

Naive Bayes for text classification

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

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

Naive Bayes for text classification

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

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

$$P(c) = \frac{\text{Number of documents belonging to class } c}{\text{Total number of documents}}$$

Naive Bayes for text classification

- Calculate the probability of each class $P(c)$

Naive Bayes for text classification

- Calculate the probability of each class $P(c)$
 - For each class c in C :

Naive Bayes for text classification

- Calculate the probability of each class $P(c)$
 - For each class c in C :
 - $P(c) = | \text{docs} | / | \text{total number of documents} |$
 - docs is the number of documents belonging to class c

Naive Bayes for text classification

- Calculate the probability of each class $P(c)$
 - For each class c in C :
 - $P(c) = | \text{docs} | / | \text{total number of documents} |$
 - docs is the number of documents belonging to class c
- Calculate the conditional probabilities $P(w|c)$

Naive Bayes for text classification

- Calculate the probability of each class $P(c)$
 - For each class c in C :
 - $P(c) = | \text{docs} | / | \text{total number of documents} |$
 - docs is the number of documents belonging to class c
- Calculate the conditional probabilities $P(w|c)$
 - $P(w|c) = \text{Count}(w,c) / \text{Count}(c)$

Naive Bayes for text classification

- Calculate the probability of each class $P(c)$
 - For each class c in C :
 - $P(c) = | \text{docs} | / | \text{total number of documents} |$
 - docs is the number of documents belonging to class c
- Calculate the conditional probabilities $P(w|c)$
 - $P(w|c) = \text{Count}(w,c) / \text{Count}(c)$
 - Laplace add 1 smoothing to deal with new words

Naive Bayes for text classification

- Calculate the probability of each class $P(c)$
 - For each class c in C :
 - $P(c) = | \text{docs} | / | \text{total number of documents} |$
 - docs is the number of documents belonging to class c
- Calculate the conditional probabilities $P(w|c)$
 - $P(w|c) = \text{Count}(w,c) / \text{Count}(c)$
 - Laplace add 1 smoothing to deal with new words
 - $P(w|c) = [\text{Count}(w,c) + 1] / [\text{Count}(c) + |V|]$

Naive Bayes for text classification

- Calculate the probability of each class $P(c)$
 - For each class c in C :
 - $P(c) = | \text{docs} | / | \text{total number of documents} |$
 - docs is the number of documents belonging to class c
- Calculate the conditional probabilities $P(w|c)$
 - $P(w|c) = \text{Count}(w,c) / \text{Count}(c)$
 - Laplace add 1 smoothing to deal with new words
 - $P(w|c) = [\text{Count}(w,c) + 1] / [\text{Count}(c) + |V|]$
 - $\text{Count}(w,c)$ = Number of times w occurs in documents of class c
 - $\text{Count}(c)$ = Number of words in documents of class c

Naive Bayes for text classification

- Calculate the probability of each class $P(c)$
 - For each class c in C :
 - $P(c) = | \text{docs} | / | \text{total number of documents} |$
 - docs is the number of documents belonging to class c
- Calculate the conditional probabilities $P(w|c)$
 - $P(w|c) = \text{Count}(w,c) / \text{Count}(c)$
 - Laplace add 1 smoothing to deal with new words
 - $P(w|c) = [\text{Count}(w,c) + 1] / [\text{Count}(c) + |V|]$
 - $\text{Count}(w,c)$ = Number of times w occurs in documents of class c
 - $\text{Count}(c)$ = Number of words in documents of class c
 - $|V|$ = Vocabulary size

Naive Bayes for text classification

	Doc_ID	Document	Class
Train	1	Chennai Delhi Mumbai	i
	2	Delhi Delhi Kolkata	i
	3	Delhi Gurgaon	i
	4	Beijing Shanghai Delhi	c

Naive Bayes for text classification

	Doc_ID	Document	Class
Train	1	Chennai Delhi Mumbai	i
	2	Delhi Delhi Kolkata	i
	3	Delhi Gurgaon	i
	4	Beijing Shanghai Delhi	c
Test	5	Delhi Delhi Kolkata Kolkata Beijing	

Naive Bayes for text classification

	Doc_ID	Document	Class
Train	1	Chennai Delhi Mumbai	i
	2	Delhi Delhi Kolkata	i
	3	Delhi Gurgaon	i
	4	Beijing Shanghai Delhi	c
Test	5	Delhi Delhi Kolkata Kolkata Beijing	

$$P(i) = 3/4$$

Naive Bayes for text classification

	Doc_ID	Document	Class
Train	1	Chennai Delhi Mumbai	i
	2	Delhi Delhi Kolkata	i
	3	Delhi Gurgaon	i
	4	Beijing Shanghai Delhi	c
Test	5	Delhi Delhi Kolkata Kolkata Beijing	

$$P(i) = 3/4$$

$$P(c) = 1/4$$

Naive Bayes for text classification

	Doc_ID	Document	Class
Train	1	Chennai Delhi Mumbai	i
	2	Delhi Delhi Kolkata	i
	3	Delhi Gurgaon	i
	4	Beijing Shanghai Delhi	c
Test	5	Delhi Delhi Kolkata Kolkata Beijing	

$$P(i) = 3/4$$

$$P(c) = 1/4$$

$$P(w|c) = [\text{Count}(w,c) + 1] / [\text{Count}(c) + |V|]$$

Naive Bayes for text classification

	Doc_ID	Document	Class
Train	1	Chennai Delhi Mumbai	i
	2	Delhi Delhi Kolkata	i
	3	Delhi Gurgaon	i
	4	Beijing Shanghai Delhi	c
Test	5	Delhi Delhi Kolkata Kolkata Beijing	

$$P(i) = 3/4$$

$$P(c) = 1/4$$

$$P(w|c) = [\text{Count}(w,c) + 1] / [\text{Count}(c) + |V|]$$

$$P(\text{Delhi} | i) =$$

Naive Bayes for text classification

	Doc_ID	Document	Class
Train	1	Chennai Delhi Mumbai	i
	2	Delhi Delhi Kolkata	i
	3	Delhi Gurgaon	i
	4	Beijing Shanghai Delhi	c
Test	5	Delhi Delhi Kolkata Kolkata Beijing	

$$P(i) = 3/4$$

$$P(c) = 1/4$$

$$P(w|c) = [\text{Count}(w,c) + 1] / [\text{Count}(c) + |V|]$$

$$P(\text{Delhi} | i) = [4 + 1] /$$

Naive Bayes for text classification

	Doc_ID	Document	Class
Train	1	Chennai Delhi Mumbai	i
	2	Delhi Delhi Kolkata	i
	3	Delhi Gurgaon	i
	4	Beijing Shanghai Delhi	c
Test	5	Delhi Delhi Kolkata Kolkata Beijing	

$$P(i) = 3/4$$

$$P(c) = 1/4$$

$$P(w|c) = [\text{Count}(w,c) + 1] / [\text{Count}(c) + |V|]$$

$$P(\text{Delhi} | i) = [4 + 1] / [8 +$$

Naive Bayes for text classification

	Doc_ID	Document	Class
Train	1	Chennai Delhi Mumbai	i
	2	Delhi Delhi Kolkata	i
	3	Delhi Gurgaon	i
	4	Beijing Shanghai Delhi	c
Test	5	Delhi Delhi Kolkata Kolkata Beijing	

$$P(i) = 3/4$$

$$P(c) = 1/4$$

$$P(w|c) = [\text{Count}(w,c) + 1] / [\text{Count}(c) + |V|]$$

$$P(\text{Delhi} | i) = [4 + 1] / [8 + 7]$$

Naive Bayes for text classification

	Doc_ID	Document	Class
Train	1	Chennai Delhi Mumbai	i
	2	Delhi Delhi Kolkata	i
	3	Delhi Gurgaon	i
	4	Beijing Shanghai Delhi	c
Test	5	Delhi Delhi Kolkata Kolkata Beijing	

$$P(i) = 3/4$$

$$P(\text{Delhi} \mid i) = 5/15$$

$$P(c) = 1/4$$

Naive Bayes for text classification

	Doc_ID	Document	Class
Train	1	Chennai Delhi Mumbai	i
	2	Delhi Delhi Kolkata	i
	3	Delhi Gurgaon	i
	4	Beijing Shanghai Delhi	c
Test	5	Delhi Delhi Kolkata Kolkata Beijing	

$$P(i) = 3/4$$

$$P(c) = 1/4$$

$$P(\text{Delhi} \mid i) = 5/15$$

$$P(\text{Kolkata} \mid i) = 2/15$$

Naive Bayes for text classification

	Doc_ID	Document	Class
Train	1	Chennai Delhi Mumbai	i
	2	Delhi Delhi Kolkata	i
	3	Delhi Gurgaon	i
	4	Beijing Shanghai Delhi	c
Test	5	Delhi Delhi Kolkata Kolkata Beijing	

$$P(i) = 3/4$$

$$P(c) = 1/4$$

$$P(\text{Delhi} \mid i) = 5/15$$

$$P(\text{Kolkata} \mid i) = 2/15$$

$$P(\text{Beijing} \mid i) = 1/15$$

Naive Bayes for text classification

	Doc_ID	Document	Class
Train	1	Chennai Delhi Mumbai	i
	2	Delhi Delhi Kolkata	i
	3	Delhi Gurgaon	i
	4	Beijing Shanghai Delhi	c
Test	5	Delhi Delhi Kolkata Kolkata Beijing	

$$P(i) = 3/4$$

$$P(c) = 1/4$$

$$P(\text{Delhi} \mid i) = 5/15$$

$$P(\text{Kolkata} \mid i) = 2/15$$

$$P(\text{Beijing} \mid i) = 1/15$$

$$P(\text{Delhi} \mid c) = 2/10$$

Naive Bayes for text classification

	Doc_ID	Document	Class
Train	1	Chennai Delhi Mumbai	i
	2	Delhi Delhi Kolkata	i
	3	Delhi Gurgaon	i
	4	Beijing Shanghai Delhi	c
Test	5	Delhi Delhi Kolkata Kolkata Beijing	

$$P(i) = 3/4$$

$$P(c) = 1/4$$

$$P(\text{Delhi} \mid i) = 5/15$$

$$P(\text{Kolkata} \mid i) = 2/15$$

$$P(\text{Beijing} \mid i) = 1/15$$

$$P(\text{Delhi} \mid c) = 2/10$$

$$P(\text{Kolkata} \mid c) = 1/10$$

Naive Bayes for text classification

	Doc_ID	Document	Class
Train	1	Chennai Delhi Mumbai	i
	2	Delhi Delhi Kolkata	i
	3	Delhi Gurgaon	i
	4	Beijing Shanghai Delhi	c
Test	5	Delhi Delhi Kolkata Kolkata Beijing	

$$P(i) = 3/4$$

$$P(c) = 1/4$$

$$P(\text{Delhi} \mid i) = 5/15$$

$$P(\text{Kolkata} \mid i) = 2/15$$

$$P(\text{Beijing} \mid i) = 1/15$$

$$P(\text{Delhi} \mid c) = 2/10$$

$$P(\text{Kolkata} \mid c) = 1/10$$

$$P(\text{Beijing} \mid c) = 2/10$$

Naive Bayes for text classification

$$\begin{array}{lll} P(i) = 3/4 & P(\text{Delhi} | i) = 5/15 & P(\text{Delhi} | c) = 2/10 \\ P(c) = 1/4 & P(\text{Kolkata} | i) = 2/15 & P(\text{Kolkata} | c) = 1/10 \\ & P(\text{Beijing} | i) = 1/15 & P(\text{Beijing} | c) = 2/10 \end{array}$$

	Doc_ID	Document	Class
Test	5	Delhi Delhi Kolkata Kolkata Beijing	

Naive Bayes for text classification

$$\begin{array}{lll} P(i) = 3/4 & P(\text{Delhi} | i) = 5/15 & P(\text{Delhi} | c) = 2/10 \\ P(c) = 1/4 & P(\text{Kolkata} | i) = 2/15 & P(\text{Kolkata} | c) = 1/10 \\ & P(\text{Beijing} | i) = 1/15 & P(\text{Beijing} | c) = 2/10 \end{array}$$

	Doc_ID	Document	Class
Test	5	Delhi Delhi Kolkata Kolkata Beijing	

$$P(i | D5) =$$

Naive Bayes for text classification

$$\begin{array}{lll} P(i) = 3/4 & P(\text{Delhi} | i) = 5/15 & P(\text{Delhi} | c) = 2/10 \\ P(c) = 1/4 & P(\text{Kolkata} | i) = 2/15 & P(\text{Kolkata} | c) = 1/10 \\ & P(\text{Beijing} | i) = 1/15 & P(\text{Beijing} | c) = 2/10 \end{array}$$

	Doc_ID	Document	Class
Test	5	Delhi Delhi Kolkata Kolkata Beijing	

$$P(i | D5) = P(i)$$

Naive Bayes for text classification

$$\begin{array}{lll} P(i) = 3/4 & P(\text{Delhi} | i) = 5/15 & P(\text{Delhi} | c) = 2/10 \\ P(c) = 1/4 & P(\text{Kolkata} | i) = 2/15 & P(\text{Kolkata} | c) = 1/10 \\ & P(\text{Beijing} | i) = 1/15 & P(\text{Beijing} | c) = 2/10 \end{array}$$

	Doc_ID	Document	Class
Test	5	Delhi Delhi Kolkata Kolkata Beijing	

$$P(i | D5) = P(i) * P(\text{Delhi} | i) * P(\text{Delhi} | i) * P(\text{Kolkata} | i) * P(\text{Kolkata} | i) * P(\text{Beijing} | i)$$

Naive Bayes for text classification

$$\begin{array}{lll} P(i) = 3/4 & P(\text{Delhi} | i) = 5/15 & P(\text{Delhi} | c) = 2/10 \\ P(c) = 1/4 & P(\text{Kolkata} | i) = 2/15 & P(\text{Kolkata} | c) = 1/10 \\ & P(\text{Beijing} | i) = 1/15 & P(\text{Beijing} | c) = 2/10 \end{array}$$

	Doc_ID	Document	Class
Test	5	Delhi Delhi Kolkata Kolkata Beijing	

$$P(i | D5) = P(i) * P(\text{Delhi} | i) * P(\text{Delhi} | i) * P(\text{Kolkata} | i) * P(\text{Kolkata} | i) * P(\text{Beijing} | i)$$

$$P(i | D5) = 3/4 * 5/15 * 5/15 * 2/15 * 2/15 * 1/15$$

Naive Bayes for text classification

$$\begin{aligned}P(i) &= 3/4 & P(\text{Delhi} \mid i) &= 5/15 & P(\text{Delhi} \mid c) &= 2/10 \\P(c) &= 1/4 & P(\text{Kolkata} \mid i) &= 2/15 & P(\text{Kolkata} \mid c) &= 1/10 \\& & P(\text{Beijing} \mid i) &= 1/15 & P(\text{Beijing} \mid c) &= 2/10\end{aligned}$$

	Doc_ID	Document	Class
Test	5	Delhi Delhi Kolkata Kolkata Beijing	

$$P(i \mid D5) = P(i) * P(\text{Delhi} \mid i) * P(\text{Delhi} \mid i) * P(\text{Kolkata} \mid i) * P(\text{Kolkata} \mid i) * P(\text{Beijing} \mid i)$$

$$P(i \mid D5) = 3/4 * 5/15 * 5/15 * 2/15 * 2/15 * 1/15 = 0.000098$$

Naive Bayes for text classification

$$\begin{aligned}P(i) &= 3/4 & P(\text{Delhi} \mid i) &= 5/15 & P(\text{Delhi} \mid c) &= 2/10 \\P(c) &= 1/4 & P(\text{Kolkata} \mid i) &= 2/15 & P(\text{Kolkata} \mid c) &= 1/10 \\& & P(\text{Beijing} \mid i) &= 1/15 & P(\text{Beijing} \mid c) &= 2/10\end{aligned}$$

	Doc_ID	Document	Class
Test	5	Delhi Delhi Kolkata Kolkata Beijing	

$$P(i \mid D5) = 0.000098$$

$$P(c \mid D5) = P(c) * P(\text{Delhi} \mid c) * P(\text{Delhi} \mid c) * P(\text{Kolkata} \mid c) * P(\text{Kolkata} \mid c) * P(\text{Beijing} \mid c)$$

Naive Bayes for text classification

$$\begin{array}{lll} P(i) = 3/4 & P(\text{Delhi} | i) = 5/15 & P(\text{Delhi} | c) = 2/10 \\ P(c) = 1/4 & P(\text{Kolkata} | i) = 2/15 & P(\text{Kolkata} | c) = 1/10 \\ & P(\text{Beijing} | i) = 1/15 & P(\text{Beijing} | c) = 2/10 \end{array}$$

	Doc_ID	Document	Class
Test	5	Delhi Delhi Kolkata Kolkata Beijing	

$$P(i | D5) = 0.000098$$

$$P(c | D5) = P(c) * P(\text{Delhi} | c) * P(\text{Delhi} | c) * P(\text{Kolkata} | c) * P(\text{Kolkata} | c) * P(\text{Beijing} | c)$$

$$P(c | D5) = 1/4 * 2/10 * 2/10 * 1/10 * 1/10 * 2/10$$

Naive Bayes for text classification

$$\begin{aligned}P(i) &= 3/4 & P(\text{Delhi} \mid i) &= 5/15 & P(\text{Delhi} \mid c) &= 2/10 \\P(c) &= 1/4 & P(\text{Kolkata} \mid i) &= 2/15 & P(\text{Kolkata} \mid c) &= 1/10 \\& & P(\text{Beijing} \mid i) &= 1/15 & P(\text{Beijing} \mid c) &= 2/10\end{aligned}$$

	Doc_ID	Document	Class
Test	5	Delhi Delhi Kolkata Kolkata Beijing	

$$P(i \mid D5) = 0.000098$$

$$P(c \mid D5) = P(c) * P(\text{Delhi} \mid c) * P(\text{Delhi} \mid c) * P(\text{Kolkata} \mid c) * P(\text{Kolkata} \mid c) * P(\text{Beijing} \mid c)$$

$$P(c \mid D5) = 1/4 * 2/10 * 2/10 * 1/10 * 1/10 * 2/10 = 0.00002$$

Naive Bayes for text classification

$$\begin{array}{lll} P(i) = 3/4 & P(\text{Delhi} | i) = 5/15 & P(\text{Delhi} | c) = 2/10 \\ P(c) = 1/4 & P(\text{Kolkata} | i) = 2/15 & P(\text{Kolkata} | c) = 1/10 \\ & P(\text{Beijing} | i) = 1/15 & P(\text{Beijing} | c) = 2/10 \end{array}$$

	Doc_ID	Document	Class
Test	5	Delhi Delhi Kolkata Kolkata Beijing	

$$P(i | D5) = 0.000098$$

$$P(c | D5) = 0.00002$$

Naive Bayes for text classification

$$P(i) = 3/4$$

$$P(c) = 1/4$$

$$P(\text{Delhi} | i) = 5/15$$

$$P(\text{Kolkata} | i) = 2/15$$

$$P(\text{Beijing} | i) = 1/15$$

$$P(\text{Delhi} | c) = 2/10$$

$$P(\text{Kolkata} | c) = 1/10$$

$$P(\text{Beijing} | c) = 2/10$$

	Doc_ID	Document	Class
Test	5	Delhi Delhi Kolkata Kolkata Beijing	

$$P(i | D5) = 0.000098$$

$$P(i | D5) > P(c | D5)$$

$$P(c | D5) = 0.00002$$

Naive Bayes for text classification

$$P(i) = 3/4$$

$$P(c) = 1/4$$

$$P(\text{Delhi} | i) = 5/15$$

$$P(\text{Kolkata} | i) = 2/15$$

$$P(\text{Beijing} | i) = 1/15$$

$$P(\text{Delhi} | c) = 2/10$$

$$P(\text{Kolkata} | c) = 1/10$$

$$P(\text{Beijing} | c) = 2/10$$

	Doc_ID	Document	Class
Test	5	Delhi Delhi Kolkata Kolkata Beijing	i

$$P(i | D5) = 0.000098$$

$$P(i | D5) > P(c | D5)$$

$$P(c | D5) = 0.00002$$

Thank You