

# Probabilistic Models: Spring 2014

## Document Classification Example, Solutions

We are given the following **corpus** and topics. Only the words in **bold** are in the vocabulary.

Topic	Text
Fantasy	The hobbit tricked the <b>troll</b> . He hid from the <b>dragon</b> . The <b>dragon</b> set the <b>town on fire</b> . The dwarf killed the <b>dragon</b> and became <b>king</b> .
Technology	Many people use a <b>fire wall</b> to increase their <b>security</b> . The <b>security forum</b> helps people configure their <b>fire wall</b> to prevent hackers from setting their computers on <b>fire</b> .
High Seas	The <b>pirate</b> sailed his <b>ship</b> into <b>town</b> . The <b>pirate</b> scaled the <b>wall</b> and took the <b>king</b> prisoner on the <b>ship</b> . He later set the <b>town on fire</b> .
Technology	A <b>troll</b> lives in this <b>forum</b> . Do not feed the <b>troll</b> ; he believes he is <b>king</b> of the <b>forum</b> and will set any thread on <b>fire</b> .
Fantasy	The <b>king</b> beyond the <b>wall</b> attacked a <b>town</b> . A <b>pirate</b> works for a different <b>king</b> . Yet another <b>king</b> has a <b>dragon</b> that set a <b>town on fire</b> .

1. Convert the documents into their bag of words representation. Use this order for the words: dragon, fire, forum, king, pirate, security, ship, town, troll, wall.

Table 1: The bag of words representation of each document,  $\mathbf{n}$ , and the number of times each word  $j$  occurs in documents of each topic  $k$ ,  $\sum_{i \in \mathbf{z}_k} \mathbf{n}_{i,j}$ . The “Total” gives  $\sum_{s=1}^d \sum_{i \in \mathbf{z}_k} \mathbf{n}_{i,s}$  for each topic.

Topic	Bag of words					Topic		
	Fantasy	Technology	High Seas	Technology	Fantasy	Fantasy	Technology	High Seas
dragon	3	0	0	0	1	4	0	0
fire	1	3	1	1	1	2	4	1
forum	0	1	0	2	0	0	3	0
king	1	0	1	1	3	4	1	1
pirate	0	0	2	0	1	1	0	2
security	0	2	0	0	0	0	2	0
ship	0	0	2	0	0	0	0	2
town	1	0	2	0	2	3	0	2
troll	1	0	0	2	0	1	2	0
wall	0	2	1	0	1	1	2	1
Total						16	14	9

2. Construct the naive Bayes classifier for the corpus.

- Prior probabilities for the topics

$$Pr(C = \text{Fantasy}) = \frac{N_{\text{Fantasy}} + 1}{N + T} = \frac{3}{8}$$

$$Pr(C = \text{Technology}) = \frac{N_{\text{Technology}} + 1}{N + T} = \frac{3}{8}$$

$$Pr(C = \text{High Seas}) = \frac{N_{\text{High Seas}} + 1}{N + T} = \frac{2}{8}$$

- Conditional probabilities for the words given the topics

We can calculate, for example,

$$Pr(w_t = \text{dragon} | C = \text{Fantasy}) = \frac{1 + \sum_{i \in \mathbf{Z}_{\text{Fantasy}}} \mathbf{n}_{i, \text{dragon}}}{d + \sum_{s=1}^d \sum_{i \in \mathbf{Z}_{\text{Fantasy}}} \mathbf{n}_{i, s}}$$

$$= \frac{1 + 4}{10 + 16}$$

$$= \frac{5}{26}$$

The rest of the conditional probabilities are calculated similarly.

Word	Topic		
	Fantasy	Technology	High Seas
dragon	$\frac{5}{26}$	$\frac{1}{24}$	$\frac{1}{19}$
fire	$\frac{3}{26}$	$\frac{5}{24}$	$\frac{2}{19}$
forum	$\frac{1}{26}$	$\frac{4}{24}$	$\frac{1}{19}$
king	$\frac{5}{26}$	$\frac{2}{24}$	$\frac{2}{19}$
pirate	$\frac{2}{26}$	$\frac{1}{24}$	$\frac{3}{19}$
security	$\frac{1}{26}$	$\frac{3}{24}$	$\frac{1}{19}$
ship	$\frac{1}{26}$	$\frac{1}{24}$	$\frac{3}{19}$
town	$\frac{4}{26}$	$\frac{1}{24}$	$\frac{3}{19}$
troll	$\frac{2}{26}$	$\frac{3}{24}$	$\frac{1}{19}$
wall	$\frac{2}{26}$	$\frac{3}{24}$	$\frac{2}{19}$

3. Calculate the likelihood, or **conditional distributions**, for each document in the corpus ( $Pr(\mathbf{n}_i | C = z_i)$ ).
4. Calculate the posterior probability, or **classification distribution**, for the following unlabeled documents ( $Pr(C = k | \mathbf{n}_i)$ ).

Topic	Text
?	The red <b>king</b> and his <b>troll</b> attacked the <b>town</b> by <b>ship</b> . Somehow, the red <b>king</b> still set the <b>town</b> on <b>fire</b> .
?	The <b>forum</b> is on <b>fire</b> with discussion of a <b>pirate</b> ship which bypassed the <b>security</b> of a cruise <b>ship</b> . The <b>pirate</b> uploaded a video to the <b>forum</b> ; naturally, the cruise <b>ship</b> was on <b>fire</b> .

## Some useful equations

$N$  := the number of documents

$T$  := the number of topics

$N_k$  := the number of documents from topic  $k$

$\mathbf{n}_{i,j}$  := the number of times word  $j$  appears in document  $i$

$z_i$  := the topic of document  $i$

$\mathbf{Z}_k$  := the indices of all documents from topic  $k$

$$Pr(C = k) = \frac{N_k + 1}{N + T}$$

$$Pr(w_t = j | C = k) = \frac{1 + \sum_{i \in \mathbf{Z}_k} \mathbf{n}_{i,j}}{d + \sum_{s=1}^d \sum_{i \in \mathbf{Z}_k} \mathbf{n}_{i,s}}$$

$$P(\mathbf{n}_i | C = k) = P(\text{drawing } \mathbf{n}_i \text{ one way} | C = k) \times \text{number of ways to draw } \mathbf{n}_i$$

$$Pr(C = k | \mathbf{n}_i) = \frac{Pr(\mathbf{n}_i | C = k) \times Pr(C = k)}{Pr(\mathbf{n}_i)}$$