# Probabilistic Models: Spring 2014
# Document Classification Example

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

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

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.

2. Construct the naive Bayes classifier for the corpus.

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

# Some useful equations

$$N := \text{the number of documents}$$
$$T := \text{the number of topics}$$
$$N_k := \text{the number of documents from topic } k$$
$$\mathbf{n}_{i,j} := \text{the number of times word } j \text{ appears in document } i$$
$$z_i := \text{the topic of document } i$$
$$\mathbf{Z}_k := \text{the indices of all documents from topic } k$$

$$Pr(C = k) = \frac{N_k + 1}{N + T}$$
$$Pr(w_t = j | C = k) = \frac{1 + \sum_{i \in \mathbf{Z}_k} \mathbf{n}_{i,j}}{d + \sum_{s=1}^{d} \sum_{i \in \mathbf{Z}_k} \mathbf{n}_{i,s}}$$
$$P(\mathbf{n}_i | C = k) = P(\text{drawing } \mathbf{n}_i \text{ one way} | C = k) \times \text{number of ways to draw } \mathbf{n}_i$$
$$Pr(C = k | \mathbf{n}_i) = \frac{Pr(\mathbf{n}_i | C = k) \times Pr(C = k)}{Pr(\mathbf{n}_i)}$$