

Basic Probability

- Conditional – or Posterior - Probability
 - $P(A|B)$ == the probability that A will occur given that B has occurred
 - $P(A|B) = P(A \text{ and } B \text{ both occur}) / P(B)$
 - For example: in a sentence from the “Wizard of Oz” what is the probability of the next word being “witch” given that the previous word was “wicked”
 - $P(\text{witch} | \text{wicked})$
- Prior Probability
 - The probability of something given no information
 - $P(\text{“the”}) = 0.07$ in English text

Bayes Rule

- Bayes Theorem

- $P(A|B) = (P(B|A)/P(B)) P(A)$

- $P(B|A)/P(B)$ is called the Likelihood of A given B

- Example

- 2 bowls of cookies

- Bowl A: 30 Chocolate chip, 10 plain
 - Bowl B: 20 Chocolate chip, 20 plain

- Suppose randomly pick a bowl and then randomly select a cookie from that bowl. When you do so, you get a chocolate chip

- What is the probability that you picked from bowl A?

- $P(A | cc)$

- Know: $P(cc|A) = 0.75$, $P(cc|B) = 0.5$, $P(cc)=0.675$, $P(A) = 0.5$

- $P(A|cc) = P(cc|A)*P(A)/P(cc) = 0.75*0.5/0.675 = 0.6$

Bayes and IR

- Want to compute
 - $P(Q=\text{true} \mid \text{document } D)$
 $= (P(D \mid Q=\text{true}) * P(Q=\text{true})) / P(D)$
- So, by Bayes would need:
 - $P(D \mid Q=\text{true}), P(Q=\text{true}), P(D)$
 - $P(Q=\text{true})$ = probability of picking a relevant document from among all documents. (This is the same for all documents)
 - $P(D) = P(t_1) * \dots$ for each word in D
 - *This we know, but it drops out as a constant term*
 - *Fortunate as VERY near 0.*

- $P(D | Q=\text{true})$
 - Assumes independence
 - So $P(D|Q=\text{true}) = P(t_1|Q=\text{true}) * P(t_2|Q=\text{true}) * \dots$
 - For each word t_i in D
 - $p_i = P(Q=\text{true} | t_i)$, $q_i = P(Q \neq \text{true} | t_i)$
 - Then let $c_i = \log \left(\frac{p_i(1-q_i)}{q_i(1-p_i)} \right)$
 - Roughly, c_i is how predictive of relevancy a given word is.
 - The **relevance weight**
 - Simple model says $p_i=1$ iff word t_i in Q ,
 - problem – this leads to division by 0 and is just wrong
 - When would this be correct?
 - Problem – words are all equally weighted.
 - Is this really a problem?
 - Summing the c_i gives an indicator of relevance of the document. Roughly, if >0 then relevant
 - Problem, where do you get all of the p_i and q_i
 - The sum gives a total ordering of documents.

Naïve Bayes Model

- A very popular approach to spam filtering
- Partially because it has a simple mechanism for incorporating user feedback
- Partially because it works pretty well
- Partially because this model allows a tunable threshold for splitting spam and non-spam

Naïve Bayes

- $$\frac{(P(isSpam) * P(D \vee isSpam))}{\sum_{x \in yes, no} (P(isSpam=x) * P(D \vee isSpam=x))}$$

– Note that this is not quite the Bayes formulation I gave earlier. Why?

- Assuming independence this is

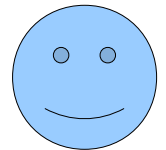
- $$\frac{(P(isSpam) * \prod_{w \in D} P(w \text{ given } isSpam))}{(\sum_{x \in yes, no} (P(isSpam=x) * \prod_{w \in D} P(w \text{ given } isSpam=x)))}$$

Naïve Bayes for Spam

- By Heaps law the number of terms in the product will be very large, and by Zipfs law most of these terms will be rare.
- Statistics about rare terms will be noisy (implied by the law of large numbers)
- Therefore, if use every term the estimate will be very noisy
- Therefore, should use only those terms that are reliable
 - Which ones are reliable?

Feature Selection

- Terms that are sufficiently frequent
 - This gives reliable terms but many of these terms will be uninformative (Luhn's idea)
- Use only terms that are highly predictive
 - Mutual Information
 - $\sum_{y \in \text{present, absent}} \sum_{x \in (\text{yes, no})} P(\text{word} = y, \text{isSpam} = x) * \log\left(\frac{P(\text{word} = y, \text{isSpam} = x)}{P(\text{word} = y) * P(\text{isSpam} = x)}\right)$
 - $P(\text{word} = y, \text{isSpam} = x)$
 - For discrete variables
 - = $P(\text{word} = y \mid \text{isSpam} = x) * P(\text{isSpam} = x)$
 - = $P(\text{isSpam} = x \mid \text{word} = y) * P(\text{word} = y)$
 - Note that MI gives a probabilistic foundation for Luhn's idea.



TF-IDF vs Probabilistic Models

- IDF term looks a lot like divisor in Bayes rule
- TF term looks a lot like – what?
 - It is related to $P(t_i | D)$ but this is not used.
- What in Probabilistic model corrects for the benefit of big documents?
 - Unlike cosine distance which is strictly non-negative the c_i terms may be less than 0. How?
 - So, VS only accumulates evidence in favor according to the intersection of document and query. Probabilistic accumulates evidence for and against according to all words in document.