

# 605.744: Information Retrieval

## Problem Set (Module 7)

Sabbir Ahmed

October 17, 2022

1. (40%) Briefly explain the following concepts in your own words and be sure to indicate how they are relevant to text classification:

- (a) bias-variance tradeoff

**Answer:** Bias is the difference between the average prediction of the model and the expected value, while variance is the measure of the spread of the data. High bias oversimplifies the model which leads to consistent errors, while high variance leads to overfitting models. The bias-variance tradeoff refers to the tradeoff between those 2 metrics such that they do not underfit or overfit the model.

- (b) k-fold cross-validation

**Answer:** K-fold cross-validation is typically required in models involved with limited data. It is a method where the data is shuffled and resampled into K-groups to evaluate the model.

- (c) macro-averaging vs. micro-averaging

**Answer:** Macro-averaging gives equal weight to each class by computing their average, while micro-averaging prioritizes each per-document classification decision by pooling them across classes to create contingency tables.

- (d) soft margin

**Answer:**

2. (10%) Why is 3-Nearest-Neighbor (3-NN) almost always a better choice than 1-NN for a binary (i.e, two-class) text classification problem.

**Answer:**

3. (50%) Naïve Bayes using the Binomial (also known as Bernoulli) model. First calculate estimates of  $P(c)$  and  $P(w|c)$  given the following training sentences. There are only three classes: Travel, Business, and Health. You should not use any smoothing. For features you should only use the following seven vocabulary terms: {denver, employers, florida, hospital, jobs, nurses, vacation} and you should ignore all other words and any punctuation. Next compute the probability of each class for the two test documents A & B below. Finally, indicate which is the predicted (i.e., the most likely) class for each test document. Please read these directions thoroughly, count carefully, and do show all of your work. Report probabilities using scientific notation (e.g.,  $1.563 \times 10^{-5}$ ) with three digits after the decimal point.

**Training data:**

- 1) Travel: denver hospital administrator takes vacation in florida
- 2) Travel: nurses plan a trip to florida
- 3) Travel: employers offering more jobs with vacation benefits
- 4) Business: employers see growth in information science
- 5) Business: hospital nurses in denver say high paying jobs are vanishing
- 6) Business: employers say florida is nice vacation spot and there are good jobs
- 7) Health: study: more hospital nurses need to take a vacation
- 8) Health: local doctors attend florida conference on diabetes
- 9) Health: hospital trains maternity ward nurses
- 10) Health: denver hospital says local employers have jobs for nurses

**Test documents:**

- 1) florida nurses take skiing vacation in denver
- 2) jobs available for experienced nurses at florida hospital

**Answer:** The probabilities,  $P(c)$  and  $P(w|c)$ , are calculated with the following equations:

$$P(c_i) = \frac{\text{count}(c_i)}{N}, \quad N = 10$$

$$P(w_j|c_i) = \frac{\text{count}(w_j = 1 \in c_i)}{\text{count}(c_i)}$$

Computing the probabilities for the classes, Travel, Business, and Health:

$$\begin{aligned} P(\text{Travel}) &= \frac{3}{10} \\ P(\text{Business}) &= \frac{3}{10} \\ P(\text{Health}) &= \frac{4}{10} \end{aligned}$$

Computing the probabilities for the terms, {denver, employers, florida, hospital, jobs, nurses, vacation}:

$$\begin{aligned} P(\text{denver}|\text{Travel}) &= \frac{1}{3} \\ P(\text{denver}|\text{Business}) &= \frac{1}{3} \\ P(\text{denver}|\text{Health}) &= \frac{1}{4} \end{aligned}$$

$$\begin{aligned} P(\text{employers}|\text{Travel}) &= \frac{1}{3} \\ P(\text{employers}|\text{Business}) &= \frac{2}{3} \\ P(\text{employers}|\text{Health}) &= \frac{1}{4} \end{aligned}$$

$$P(\text{florida}|\text{Travel}) = \frac{2}{3}$$

$$P(\text{florida}|\text{Business}) = \frac{1}{3}$$

$$P(\text{florida}|\text{Health}) = \frac{1}{4}$$

$$P(\text{hospital}|\text{Travel}) = \frac{1}{3}$$

$$P(\text{hospital}|\text{Business}) = \frac{1}{3}$$

$$P(\text{hospital}|\text{Health}) = \frac{3}{4}$$

$$P(\text{jobs}|\text{Travel}) = \frac{1}{3}$$

$$P(\text{jobs}|\text{Business}) = \frac{2}{3}$$

$$P(\text{jobs}|\text{Health}) = \frac{1}{4}$$

$$P(\text{nurses}|\text{Travel}) = \frac{1}{3}$$

$$P(\text{nurses}|\text{Business}) = \frac{1}{3}$$

$$P(\text{nurses}|\text{Health}) = \frac{3}{4}$$

$$P(\text{vacation}|\text{Travel}) = \frac{2}{3}$$

$$P(\text{vacation}|\text{Business}) = \frac{1}{3}$$

$$P(\text{vacation}|\text{Health}) = \frac{1}{4}$$

Summarizing the probabilities into a table:

Term	P(w Travel)	P(w Business)	P(w Health)
denver	0.33	0.33	0.25
employers	0.33	0.66	0.25
florida	0.66	0.33	0.25
hospital	0.33	0.33	0.75
jobs	0.33	0.66	0.25
nurses	0.33	0.33	0.75
vacation	0.66	0.33	0.25

To find the best class, the following equation can be used:

$$\text{BestClass} = P(c_i) \cdot \prod_{j=1}^{\# \text{ words}} P(w_j|c_i)$$

After normalization, the first test document becomes:

$$D_1 = \{\text{"florida", "nurses", "vacation", "denver"}\}$$

$$\begin{aligned}
P(\text{Travel}|D_1) &= P(\text{Travel}) \cdot P(\text{denver}|\text{Travel}) \cdot P(\text{employers}|\text{Travel}) \cdot P(\text{florida}|\text{Travel}) \\
&\quad \cdot P(\text{hospital}|\text{Travel}) \cdot P(\text{jobs}|\text{Travel}) \cdot P(\text{nurses}|\text{Travel}) \cdot P(\text{vacation}|\text{Travel}) \\
&= (0.30)(0.33)(1 - 0.33)(0.66)(1 - 0.33)(1 - 0.33)(0.33)(0.66) \\
&= 4.390 \times 10^{-3}
\end{aligned}$$

$$\begin{aligned}
P(\text{Business}|D_1) &= P(\text{Business}) \cdot P(\text{denver}|\text{Business}) \cdot P(\text{employers}|\text{Business}) \cdot P(\text{florida}|\text{Business}) \\
&\quad \cdot P(\text{hospital}|\text{Business}) \cdot P(\text{jobs}|\text{Business}) \cdot P(\text{nurses}|\text{Business}) \cdot P(\text{vacation}|\text{Business}) \\
&= (0.30)(0.33)(1 - 0.66)(0.33)(1 - 0.33)(1 - 0.66)(0.33)(0.33) \\
&= 2.743 \times 10^{-4}
\end{aligned}$$

$$\begin{aligned}
P(\text{Health}|D_1) &= P(\text{Health}) \cdot P(\text{denver}|\text{Health}) \cdot P(\text{employers}|\text{Health}) \cdot P(\text{florida}|\text{Health}) \\
&\quad \cdot P(\text{hospital}|\text{Health}) \cdot P(\text{jobs}|\text{Health}) \cdot P(\text{nurses}|\text{Health}) \cdot P(\text{vacation}|\text{Health}) \\
&= (0.40)(0.25)(1 - 0.25)(0.25)(1 - 0.75)(1 - 0.25)(0.75)(0.25) \\
&= 6.592 \times 10^{-4}
\end{aligned}$$

Therefore, the best class for  $D_1$  is Travel.

After normalization, the second test document becomes:

$$D_2 = \{\text{"jobs", "nurses", "florida", "hospital"}\}$$

$$\begin{aligned}
P(\text{Travel}|D_2) &= P(\text{Travel}) \cdot P(\text{denver}|\text{Travel}) \cdot P(\text{employers}|\text{Travel}) \cdot P(\text{florida}|\text{Travel}) \\
&\quad \cdot P(\text{hospital}|\text{Travel}) \cdot P(\text{jobs}|\text{Travel}) \cdot P(\text{nurses}|\text{Travel}) \cdot P(\text{vacation}|\text{Travel}) \\
&= (0.30)(1 - 0.33)(1 - 0.33)(0.66)(0.33)(0.33)(0.33)(1 - 0.66) \\
&= 1.097 \times 10^{-3}
\end{aligned}$$

$$\begin{aligned}
P(\text{Business}|D_2) &= P(\text{Business}) \cdot P(\text{denver}|\text{Business}) \cdot P(\text{employers}|\text{Business}) \cdot P(\text{florida}|\text{Business}) \\
&\quad \cdot P(\text{hospital}|\text{Business}) \cdot P(\text{jobs}|\text{Business}) \cdot P(\text{nurses}|\text{Business}) \cdot P(\text{vacation}|\text{Business}) \\
&= (0.30)(1 - 0.33)(1 - 0.66)(0.33)(0.33)(0.66)(0.33)(1 - 0.33) \\
&= 1.097 \times 10^{-3}
\end{aligned}$$

$$\begin{aligned}
P(\text{Health}|D_2) &= P(\text{Health}) \cdot P(\text{denver}|\text{Health}) \cdot P(\text{employers}|\text{Health}) \cdot P(\text{florida}|\text{Health}) \\
&\quad \cdot P(\text{hospital}|\text{Health}) \cdot P(\text{jobs}|\text{Health}) \cdot P(\text{nurses}|\text{Health}) \cdot P(\text{vacation}|\text{Health}) \\
&= (0.40)(1 - 0.25)(1 - 0.25)(0.25)(0.75)(0.25)(0.75)(1 - 0.25) \\
&= 5.933 \times 10^{-3}
\end{aligned}$$

Therefore, the best class for  $D_2$  is Health.