

(a) ① case-folding: change upper case to lower case
 ② removing punctuations and non-alphanumeric chars.
 more: canonization, remove of stop word, remove of page metadata, splitting into tokens.

(b) ① $/[A-Z]/ \rightarrow$ corresponding lower case

② $/\backslash W/ \rightarrow \emptyset$

(c) australia $\rightarrow 1(4) \rightarrow 5(1) \rightarrow 6(1)$
 election $\rightarrow 1(4) \rightarrow 2(1) \rightarrow 3(4) \rightarrow 4(1) \rightarrow 5(1) \rightarrow 6(2)$
 federal $\rightarrow 1(2) \rightarrow 3(3) \rightarrow 4(2) \rightarrow 5(1) \rightarrow 6(2)$
 vote $\rightarrow 4(2) \rightarrow 6(5)$
 ausrlian $\rightarrow 5(1)$

(d) f_t , and pointers (given a token, point to the inverted index)

(e) $W_{1,t} = \langle 4, 4, 2, 0, 0 \rangle$
 $W_{2,t} = \langle 0, 1, 0, 0, 0 \rangle$
 $W_{3,t} = \langle 0, 4, 3, 0, 0 \rangle$
 $W_{4,t} = \langle 0, 1, 1, 2, 0 \rangle$
 $W_{5,t} = \langle 1, 1, 1, 0, 1 \rangle$
 $W_{6,t} = \langle 1, 2, 2, 5, 0 \rangle$
 $W_{7,t} = \langle 2, 1, 0, 0, 0 \rangle$

$$\begin{aligned}\cos(1, q) &= \frac{8+4}{\sqrt{34} \cdot \sqrt{5}} \approx 0.920 \\ \cos(2, q) &= \frac{1}{\sqrt{5}} \approx 0.447 \\ \cos(3, q) &= \frac{4}{5 \cdot \sqrt{5}} \approx 0.358 \\ \cos(4, q) &= \frac{1}{3\sqrt{5}} \approx 0.149 \\ \cos(5, q) &= \frac{3}{2\sqrt{5}} \approx 0.671 \\ \cos(6, q) &= \frac{4}{\sqrt{34} \cdot \sqrt{5}} \approx 0.307\end{aligned}$$

r. the rank is $(d_1, d_5, d_2, d_3, d_6, d_4)$

(f) i: means the document is useful to the user, which can resolve the user's information need; we discover this by checking whether the user click into it or not.

ii $P@2 = 0$

Cg)? $1 + \log_2 f_{dt} = \log_2 2 f_{dt} > \log_2 f_{dt}$

the relative weight of f_{dt} increase means the terms appear in a document frequently gets a higher weight so the documents which contain more query term may get a higher rank

(h) Spelling Correction :

(i) parsing

(ii) a standard dictionary should be included and we need to use approximate string matching to compare the misspell word to the words in the dictionary to find the correct word.

(iii) Some document with misspell word would also be considered, like document 5.

$$\begin{aligned}
 (a) \quad & P(\text{FACEBOOK} \mid \text{ele} = N, \text{fed} = N, \text{aust} = N, \text{rot} = N, \text{ansr} = Y) \\
 & = P(\text{ele} = N \mid \text{FB}) \cdot P(\text{fed} = N \mid \text{FB}) \cdot P(\text{aust} = N \mid \text{FB}) \cdot P(\text{rot} = N \mid \text{FB}) \cdot P(\text{ansr} = Y \mid \text{FB}) \cdot P(\text{FB}) \\
 & = \frac{1}{2} \cdot \frac{1}{2} \cdot 1 \cdot \frac{1}{2} \cdot \frac{1}{3} = \frac{\frac{1}{2}^2}{12}
 \end{aligned}$$

$$\begin{aligned}
 & P(\text{TWITTER} \mid \text{ele} = N, \text{fed} = N, \text{aust} = N, \text{rot} = N, \text{ansr} = Y) \\
 & = P(\text{ele} = N \mid \text{TW}) \cdot P(\text{fed} = N \mid \text{TW}) \cdot P(\text{aust} = N \mid \text{TW}) \cdot P(\text{rot} = N \mid \text{TW}) \cdot P(\text{ansr} = Y \mid \text{TW}) \cdot P(\text{TW}) \\
 & = \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{4} \cdot \frac{2}{3} = \frac{\frac{1}{2}^2}{24} \\
 & \because \frac{\frac{1}{2}^2}{12} > \frac{\frac{1}{2}^2}{24} \\
 & \therefore \text{classify it to FACEBOOK}
 \end{aligned}$$

(b) The solution is based on GR (GINI not in lecture)

$$H(R) = -\left(\frac{1}{3}\log_2 \frac{1}{3} + \frac{2}{3}\log_2 \frac{2}{3}\right) = 0.918$$

$$IG(\text{ele} \mid R) = 0$$

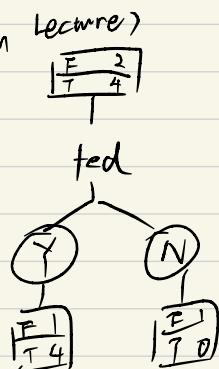
$$SI(\text{ele} \mid R) = 0$$

$$GR(\text{ele} \mid R) = 0$$

$$\begin{aligned}
 IG(\text{fed} \mid R) &= H(R) - \left(-\left(\frac{5}{6}\left(\frac{1}{5}\log_2 \frac{1}{5} + \frac{4}{5}\log_2 \frac{4}{5}\right)\right)\right) \\
 &= 0.918 - 0.602 \approx 0.316
 \end{aligned}$$

$$\begin{aligned}
 SI(\text{fed} \mid R) &= -\left(\frac{5}{6}\log_2 \frac{5}{6} + \frac{1}{6}\log_2 \frac{1}{6}\right) \\
 &= 0.650
 \end{aligned}$$

$$GR(\text{fed} \mid R) = 0.486$$



最后用 GR, 选 GR 最大的做 root.

(c) i. training data not enough, model is under fit

ii. NB not change $\because P(\text{aust} = Y \mid \text{TWITTER}) = P(\text{aust} = N \mid \text{TWITTER}) = \frac{1}{2}$

$$P(\text{aust} = Y \mid \text{FACEBOOK}) = P(\text{aust} = N \mid \text{FACEBOOK}) = \frac{1}{2}$$

(d) ele, fed, aust, vot, ausr is conditional independent given the label.
The assumption is invalid, for example aust and ausr is relevant since aust may be the misspelling form of ausr

(e) No, the dataset is small, and none of the training data is similar to the test data. Bagging is mainly to reduce sampling bias.

(f) $\{\text{ele} = Y\} \rightarrow \{\text{fed} = Y\}$ support: $\frac{5}{6}$
confident: $\frac{5}{6}$

- (a) ① Training data not enough
② Not effective to calculate distance in Boolean attribute.

Change: crawl more data and record f_t of the the data instead of whether the term appear.

- (b) the data can be separated into two side by a hyperplane, data on one side belongs to FB, other side belongs to TW

Since the data is already linear-seperable, we don't need to use kernel method to change the data, but soft-margin is still need to avoid being too susceptible to noise.

- (c) f_t , get more information, now we can get the frequency of terms, with frequency, we can calculate distance more accurately.

- (f) use cross-validation, (i) the dataset is the data with label
(ii) use part of the data set for testing, the rest of it for training, each time, and iterate the process, and calculate the calculate the average f_1 score of them at the end, f_1 score formula:
$$f_1 \text{ score} = \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$
 (iii) cross-validation