

1-a) Bernoulli Naive Bayes model:

$$P(x_{\text{francisco}} = \text{true} \mid \text{class} = \text{SFO}) = \frac{2}{2} = 1.0$$

$$P(x_{\text{london}} = \text{true} \mid \text{class} = \text{SFO}) = \frac{1}{2} = 0.5$$

$$P(x_{\text{francisco}} = \text{true} \mid \text{class} = \text{JFK}) = \frac{1}{1} = 1.0$$

b) multinomial NB Model :-

$$P(x = \text{francisco} \mid \text{class} = \text{SFO}) = \frac{4}{14}$$

$$P(x = \text{london} \mid \text{class} = \text{SFO}) = \frac{1}{14}$$

$$P(x = \text{francisco} \mid \text{class} = \text{JFK}) = \frac{1}{8}$$

c) - i) Bernoulli model :- It is not very accurate, because it ignores frequency information

ii) more accurate, because it uses frequency information. However it ignores position information so doesn't distinguish between a city name occurring at the beginning/end of itinerary from one occurring in the middle.

d) we will use as a feature ~~the~~ the term that occurs in the last position of each document

2-a) It will never choose a category unless all words in a document were seen for that category for training set. It will rank between classes for which all words were seen similarly to two smoother classifiers.

b) It will be more likely to choose category for which some/many of the words in the document unseen.

3-a) the precision is  $\frac{TP}{TP+FP}$

$$= \frac{3}{2+3} = 3/5$$

The recall is  $= \frac{TP}{TP+FN}$

$$= 3/8$$

3-5)

① The IR system which always return no result will have high accuracy for most queries, since the corpus ~~query~~ contains only a few relevant documents. Documents that are truly relevant are the only ones that will be mistakenly classified as irrelevant, and the accuracy is close to 1. Recall and precision are two different measures that can jointly capture the trade-off between more relevant results and return fewer irrelevant results.

ii) there are of course many correct answers.

the small answer is:

Assume document 1 is only relevant document

$$A_2 = \{1, 2, 3\}$$

$$B_2 = \{3\}$$

Both  $A_2$  and  $B_2$  made 2 mistakes, so they have accuracy 80%.

The precision at  $A_2$  is  $\frac{1}{3}$ , the precision for  $B_2$  is 0. Since  $B_2$  didn't return any relevant documents, it is no utility.