# Speech and Text Processing

- Dan Jurafsky
- James Martin

- Text Normalization
  - → tokenizing
  - → lemmatization
  - → stemming
  - → Sentence segmentation

# Exercises

(Q2.1)

(a) `\b [a-z A-Z]+ \b`

(b) `\b [a-z]* b \b`

(c) `\b b+ (ab+)+ \b`

(Q2.2)

(a) `(.+) ⎵ \1`

(b) `^[0-9]+.*[A-Za-z]*\.`

(c) `\b grotto \b .* \b raven \b |`
    `\b raven \b .* \b grotto \b`

(Q2.4)

| # | d | e | a | l |
|---|---|---|---|---|
| # | 0 | 1 | 2 | 3 | 4 |
| l | 1 | 1 | 2 | 3 | 4 |
| e | 2 | 2 | 1 | 2 | 3 |
| d | 3 | 2 | 2 | 2 | 3 |
| a | 4 | 3 | 3 | 3 | 3 |

(Q2.5)

| # | b | r | i | e | f |
|---|---|---|---|---|---|
| # | 0 | 1 | 2 | 3 | 4 | 5 |
| d | 1 | 2 | 3 | 4 | 5 | 6 |
| r | 2 | 3 | 2 | 3 | 4 | 5 |
| i | 3 | 4 | 3 | 2 | 3 | 4 |
| v | 4 | 5 | 4 | 3 | 4 | 5 |
| e | 5 | 6 | 5 | 4 | 5 | 6 |

| # | d | i | v | e | r | s |
|---|---|---|---|---|---|---|
| # | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| d | 1 | 0 | 1 | 2 | 3 | 4 | 5 |
| r | 2 | 1 | 2 | 3 | 4 | 5 | 6 |
| i | 3 | 2 | 1 | 2 | 3 | 4 | 5 |
| v | 4 | 3 | 2 | 1 | 2 | 3 | 4 |
| e | 5 | 4 | 3 | 2 | 1 | 2 | 3 |

# Chapter 3 : N-gram language Models

## Exercises

(Q 3.1) $P(w_n | w_{n-2} w_{n-1}) = \dfrac{C(w_{n-2} w_{n-1} w_n)}{C(w_{n-2} w_{n-1})}$

&lt;s&gt;&lt;s&gt; I am Sam &lt;/s&gt;
&lt;s&gt;&lt;s&gt; Sam I am &lt;/s&gt;
&lt;s&gt;&lt;s&gt; I do not like green eggs and ham &lt;/s&gt;

$P(I | <s><s>) = 2/3$
$P(am | <s> I) = 1/2$

(Q 3.2) $P(i\ want\ chinese\ food) = P(i|<s>)\ P(want|i)$
$P(chinese|want)\ P(food|chinese)$
$P(</s>|food)$

$= 0.25 \times .33 \times 0\ 0065 \times .52 \times 0.68$

$= 0.0001896$

$P(i\ want\ chinese\ food) = P(i|<s>)\ P(want|i)\ P(chinese|want)$
$P(food|chinese)\ P(</s>|food)$

$= .19 \times 0.21 \times 0.0029 \times 0.052 \times .4$

$= 0.000002406$

(Q 3.3) The unsmoothed probability is higher because the bigrams used in the sentences are very common and has probablities. However, in the smoothed case, their probablities are distributed among not-so-common bigrams which are not used in our test statement.

| | <s> | I | am | Sam | do | not | like | green | eggs | and | </s> |
|---|---|---|---|---|---|---|---|---|---|---|---|
| <s> | 1 | 4 | 1 | 2 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| I | 1 | 1 | 4 | 1 | 2 | 1 | 1 | 1 | 1 | 1 | 1 |
| am | 1 | 1 | 1 | 3 | 1 | 1 | 1 | 1 | " | 1 | 2 |
| Sam | 1 | 2 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 4 |
| do | " | 1 | " | " | " | 2 | " | " | " | " | 1 |
| not | " | " | " | " | " | 1 | 2 | " | " | " | " |
| like | " | " | " | " | " | " | 1 | 2 | " | " | " |
| green | " | " | " | " | " | " | " | 1 | 2 | " | " |
| eggs | " | " | " | " | " | " | " | " | 1 | 2 | " |
| and | " | " | " | 2 | " | " | " | " | " | 1 | " |
| </s> | " | " | " | 1 | " | " | " | " | " | " | " |

$$P(\text{Sam} \mid \text{am}) = \frac{C(\text{am Sam})}{C(\text{am})} = \frac{3}{14} = .214$$

| | <s> | a | b |
|---|---|---|---|
| <s> | 0 | 2 | 2 |
| a | 0 | 1 | 1 |
| b | 0 | 1 | 1 |

- $P(a\,a) = P(a \mid <s>)\,P(a \mid a) = 0.5 \times 0.5 = 0.25$
  $P(a\,b) = P(a \mid <s>)\,P(b \mid a) = 0.5 \times 0.5 = 0.25$
  $P(b\,b) = P(b \mid <s>)\,P(b \mid b) = 0.5 \times 0.5 = 0.25$
  $P(b\,a) = P(b \mid <s>)\,P(a \mid b) = 0.5 \times 0.5 = 0.25$

  $P(s \in \{a,b\}^2) = 1$

(Q3.6)

$$P(w_3 | w_1 w_2) = \frac{C(w_1 w_2 w_3) + 1}{\sum_w (C(w_1 w_2 w) + 1)}$$

$$= \frac{C(w_1 w_2 w_3) + 1}{C(w_1 w_2) + v}$$

(Q3.7) $= P(sam | am) = d_1 P(sam) + d_2 (sam | am)$

$$= \frac{1}{2} \times \frac{42}{25} + \frac{1}{2} \times \frac{2}{3}$$

$$= \frac{2}{25} + \frac{1}{3} = 0.41$$

(Q3.12) $PP(w) = \sqrt[N]{\prod_{i=1}^{N} \frac{1}{P(w_i)}}$

$P(0) = \frac{91}{100}$

$P(1) = P(2) \dots P(9) = \frac{1}{100}$

$$= \sqrt[10]{\frac{(100)^{10}}{(91)^9}}$$

$$= 1.726$$

(Q 4.1)    S =    "I always like foreign films"

$$P(neg \mid S) = \frac{P(S \mid neg)\ P(neg)}{P(S)}$$

$$P(pos \mid S) = \frac{P(S \mid pos)\ P(pos)}{P(S)}$$

→ ignore common base.

- $P(neg \mid S) = (0.16 \times 0.06^2 \times 0.15 \times 0.11) \times 0.5$

  $= 0.000004752$

- $P(pos \mid S) = (0.09 \times 0.07 \times 0.29 \times 0.04 \times 0.08) \times 0.5$

  $= 0.000002923$

The naive bayes will assign "neg" class to the sentence because $P(neg \mid S) > P(pos \mid S)$

(Q4.2)   $P(\text{comedy}) = 2/5$          $|V| = 7$
         $P(\text{action}) = 3/5$

$P(\text{fast} \mid \text{comedy}) = \dfrac{\text{Count}(\text{fast, comedy}) + 1}{\sum\limits_{w \in V}(\text{Count}(w, \text{comedy}) + 1)}$

$= \dfrac{2}{9+7} = \dfrac{2}{16}$

$P(\text{fast} \mid \text{action}) = \dfrac{3}{11+7} = \dfrac{3}{18}$

$P(\text{couple} \mid \text{comedy}) = \dfrac{3}{16}$      $P(\text{shoot} \mid \text{comedy}) = \dfrac{1}{16}$

$P(\text{couple} \mid \text{action}) = \dfrac{1}{18}$      $P(\text{shoot} \mid \text{action}) = \dfrac{5}{18}$

$P(\text{fly} \mid \text{comedy}) = \dfrac{2}{16}$

$P(\text{fly} \mid \text{action}) = \dfrac{2}{18}$

$P(\text{comedy} \mid D) = P(D \mid \text{comedy})\, P(\text{comedy})$

$= \dfrac{2}{16} \times \dfrac{3}{16} \times \dfrac{2}{16} \times \dfrac{1}{16} \times \dfrac{2}{5}$

$= 0.0000732$

$P(\text{action} \mid D) = 3/18 \times 1/18 \times 2/18 \times 5/18 \times 3/5$

$= 0.000171$

D will be classified as "action".

(Q 4:3) • Binarized naive Bayes

$P(\text{neg}) = 0.6$ $\qquad P(\text{pos}) = 0.4$

$P(\text{good} | \text{neg}) = 3/9$ $\qquad P(\text{good} | \text{pos}) = 2/7$

$P(\text{poor} | \text{neg}) = 4/9$ $\qquad P(\text{poor} | \text{pos}) = 2/7$

$P(\text{great} | \text{neg}) = 2/9$ $\qquad P(\text{great} | \text{pos}) = 3/7$

$P(\text{neg} | D) = \frac{3}{9} \times \frac{4}{9} \times \frac{2}{9} \times 0.6 = 0.0197$

$P(\text{pos} | D) = \frac{2}{7} \times \frac{2}{7} \times \frac{3}{7} \times 0.4 = 0.0139$

Classified as "neg" by BNB.

• Multinomial Naive Bayes

$P(\text{good} | \text{pos}) = 4/12$ $\qquad P(\text{good} | \text{neg}) = 3/17$

$P(\text{poor} | \text{pos}) = 2/12$ $\qquad P(\text{poor} | \text{neg}) = 11/17$

$P(\text{great} | \text{pos}) = 6/12$ $\qquad P(\text{great} | \text{neg}) = 3/17$

$P(\text{pos} | D) = \left(\frac{4}{12}\right)^2 \times \frac{2}{12} \times \frac{6}{12} \times 0.6$

$\qquad = 0.0055$

$P(\text{neg} | D) = \left(\frac{3}{17}\right)^2 \times \left(\frac{11}{17}\right) \times \left(\frac{3}{17}\right) \times 0.4$

$\qquad = 0.0014$

Classified as "pos" by MNB.
Both models disagree.

# Chapter 10:

**CQ 10·1)**    REF : witness the past     ( 18 unigram, 17 bigrams)

           HYP2 : past witness            ( 11 unig.,   10 bigrams )

Unigrams that match : past witness     ( 11 unigrams)

bigrams that match : pa as st   wi it tn ne es ss   ( 9 bigrams)

Unigram P =   11/11          Unigram R: 11/18

Bigram P =   9/10          Bigram R : 9/17

$$chr\, P = \left( 11/11 + 9/10 \right) / 2 = \cdot 95$$

$$chr\, R = \left( 11/18 + 9/17 \right) / 2 = 0\cdot 57$$

$$Chr\, F_{2,2} = 5\,\frac{chr\,P * chr\,R}{4\,chr\,P + chr\,R} = \frac{5 \times \cdot 95 \times 0\cdot 57}{4 \times \cdot 95 + 0\cdot 57}$$

$$= 0\cdot 6199$$