


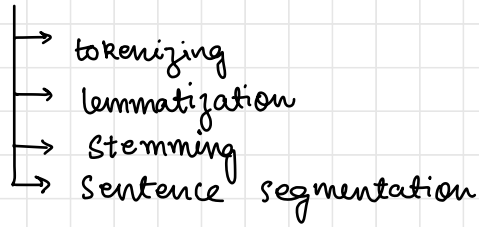
Speech and Text Processing

- Dan Jurafsky
- James Martin



Chapter 2 : Regular Expressions, Text Normalization, Edit Distance

- Text Normalization



Exercises

(Q2.1)

(a) $\backslash b [a-zA-Z]^+ \backslash b$

(b) $\backslash b [a-z]^* b \backslash b$

(c) $\backslash b b + (ab^+)^+ \backslash b$

(Q2.2)

(a) $(\cdot +) _ _ 1$

(b) $^ [0-9]^+ \cdot [A-Za-z]^* \cdot$

(c) $\backslash b \text{grotto} \backslash b \cdot \backslash b \text{raven} \backslash b \mid$
 $\backslash b \text{raven} \backslash b \cdot \backslash b \text{grotto} \backslash b$

(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

$$(\textcircled{3.1}) \quad P(w_n | w_{n-2} w_{n-1}) = \frac{C(w_{n-2} w_{n-1} w_n)}{C(w_{n-2} w_{n-1})}$$

<S><S> I am Sam </S>

<S><S> Sam I am </S>

<S><S> I do not like green eggs and ham </S>

$$P(I | <S><S>) = 2/3$$

$$P(\text{am} | <S>I) = 1/2$$

$$\begin{aligned} (\textcircled{3.2}) \quad P(i \text{ want chinese food}) &= P(i | <S>) P(\text{want} | i) \\ &\quad P(\text{chinese} | \text{want}) P(\text{food} | \text{chinese}) \\ &\quad P(</S> | \text{food}) \\ &= 0.25 \times 0.33 \times 0.0065 \times 0.52 \times 0.68 \\ &= 0.0001896 \end{aligned}$$

$$\begin{aligned} P(i \text{ want chinese food}) &= P(i | <S>) P(\text{want} | i) P(\text{chinese} | \text{want}) \\ &\quad P(\text{food} | \text{chinese}) P(</S> | \text{food}) \\ &= 0.19 \times 0.21 \times 0.0029 \times 0.052 \times 0.4 \\ &= 0.00002406 \end{aligned}$$

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

(Q3.4)

| | <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 | 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 | " | " | " | " | " | " | " |

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

(Q3.5)

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

- $$P(a a) = P(a | <s>) P(a | a) = 0.5 \times 0.5 = 0.25$$

$$P(a b) = P(a | <s>) P(b | a) = 0.5 \times 0.5 = 0.25$$

$$P(b b) = P(b | <s>) P(b | b) = 0.5 \times 0.5 = 0.25$$

$$P(b a) = P(b | <s>) P(a | b) = 0.5 \times 0.5 = 0.25$$

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

$$\begin{aligned}
 (\text{Q3.6}) \quad P(\omega_3 | \omega_1 \omega_2) &= \frac{C(\omega_1 \omega_2 \omega_3) + 1}{\sum_{\omega} (C(\omega_1 \omega_2 \omega) + 1)} \\
 &= \frac{C(\omega_1 \omega_2 \omega_3) + 1}{C(\omega_1 \omega_2) + 9}
 \end{aligned}$$

$$\begin{aligned}
 (\text{Q3.7}) \quad P(\text{sam} | \text{am}) &= d_1 P(\text{sam}) + d_2 P(\text{sam} | \text{am}) \\
 &= \frac{1}{2} \times \frac{42}{25} + \frac{1}{2} \times \frac{2}{3} \\
 &= \frac{2}{25} + \frac{1}{3} = 0.41
 \end{aligned}$$

$$\begin{aligned}
 (\text{Q3.12}) \quad PP(\omega) &= \sqrt[N]{\prod_{i=1}^N \frac{1}{P(\omega_i)}} \\
 &= \sqrt[10]{\frac{(100)^{10}}{(91)^9}} \\
 &= 1.0726
 \end{aligned}$$

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

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

Chapter 4: Naive Bayes and Sentiment Classification

(Q4.1) $S =$ "I always like foreign films"

$$P(\text{neg}|S) = \frac{P(S|\text{neg}) P(\text{neg})}{P(S)}$$

$$P(\text{pos}|S) = \frac{P(S|\text{pos}) P(\text{pos})}{P(S)}$$

→ ignore common base.

$$\begin{aligned} \bullet P(\text{neg}|S) &= (0.16 \times 0.06^2 \times 0.15 \times 0.11) \times 0.5 \\ &= 0.000004752 \end{aligned}$$

$$\begin{aligned} \bullet P(\text{pos}|S) &= (0.09 \times 0.07 \times 0.29 \times 0.04 \times 0.08) \times 0.5 \\ &= 0.000002923 \end{aligned}$$

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

$$\textcircled{34.2} \quad P(\text{comedy}) = 2/5 \quad |V| = 7$$

$$P(\text{action}) = 3/5$$

$$P(\text{fast} | \text{comedy}) = \frac{\text{count}(\text{fast}, \text{comedy}) + 1}{\sum_{w \in V} (\text{count}(w, \text{comedy}) + 1)}$$

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

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

$$P(\text{couple} | \text{comedy}) = \frac{3}{16} \quad P(\text{shoot} | \text{comedy}) = \frac{1}{16}$$

$$P(\text{couple} | \text{action}) = \frac{1}{18} \quad P(\text{shoot} | \text{action}) = \frac{5}{18}$$

$$P(\text{fly} | \text{comedy}) = \frac{2}{16}$$

$$P(\text{fly} | \text{action}) = \frac{2}{18}$$

$$P(\text{comedy} | D) = P(D | \text{comedy}) P(\text{comedy})$$

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

$$= 0.0000732$$

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

$$= 0.000171$$

D will be classified as "action".

(Q 4.3) • Binarized naive Bayes

$$P(\text{neg}) = 0.6$$

$$P(\text{pos}) = 0.4$$

$$P(\text{good} | \text{neg}) = 3/9$$

$$P(\text{good} | \text{pos}) = 2/7$$

$$P(\text{poor} | \text{neg}) = 4/9$$

$$P(\text{poor} | \text{pos}) = 2/7$$

$$P(\text{great} | \text{neg}) = 2/9$$

$$P(\text{great} | \text{pos}) = 3/7$$

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

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

Classified as "neg" by BNB.

• Multinomial Naive Bayes

$$P(\text{good} | \text{pos}) = 4/12$$

$$P(\text{good} | \text{neg}) = 3/17$$

$$P(\text{poor} | \text{pos}) = 2/12$$

$$P(\text{poor} | \text{neg}) = 11/17$$

$$P(\text{great} | \text{pos}) = 6/12$$

$$P(\text{great} | \text{neg}) = 3/17$$

$$\begin{aligned} P(\text{pos} | D) &= \left(\frac{4}{12}\right)^2 \times \frac{2}{12} \times \frac{6}{12} \times 0.6 \\ &= 0.0055 \end{aligned}$$

$$\begin{aligned} P(\text{neg} | D) &= \left(\frac{3}{17}\right)^2 \times \left(\frac{11}{17}\right) \times \left(\frac{3}{17}\right) \times 0.4 \\ &= 0.0014 \end{aligned}$$

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

Chapter 10:

(Q 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 tu ne es ss (9 bigrams)

$$\text{Unigram } P = 11/11$$

$$\text{Unigram } R = 11/18$$

$$\text{Bigram } P = 9/10$$

$$\text{Bigram } R = 9/17$$

$$\text{chr}P = (11/11 + 9/10) / 2 = .95$$

$$\text{chr}R = (11/18 + 9/17) / 2 = 0.57$$

$$\begin{aligned}\text{chr}_{f2,2} &= \frac{5 \text{ chr}P * \text{chr}R}{4 \text{ chr}P + \text{chr}R} = \frac{5 \times .95 \times 0.57}{4 \times .95 + 0.57} \\ &= 0.6199\end{aligned}$$

Chapter 12: Constituency Grammar

(Q 12.1)

Grammar used:

$S \rightarrow NP \mid VP \mid NP VP \mid NP PP$
 $NP \rightarrow Pronoun \mid Proper-Noun \mid Det Nominal \mid Nominal \mid Adj$
 $Nominal \rightarrow Nominal N \mid Adj N \mid N$
 $VP \rightarrow V \mid V NP \mid V NP PP \mid V PP$
 $PP \rightarrow Preposition NP$

Lexicals used:

Adj \rightarrow early | all | one-way | any

N \rightarrow p.m. | flights | redeye | fare | delays | five

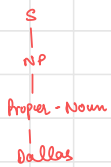
Det \rightarrow a

V \rightarrow arriving

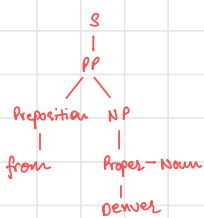
Preposition \rightarrow from | in | after | on

Proper-Noun \rightarrow Denver | Dallas | Washington | Thursday

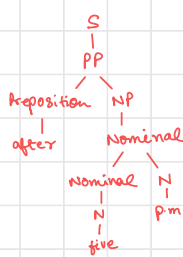
(a) Dallas



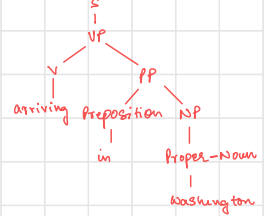
(b) from Denver



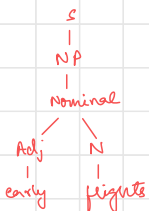
(c) after five p.m



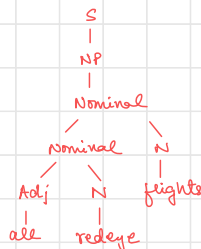
(d) arriving in Washington



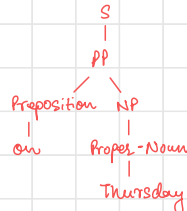
(e) early flights



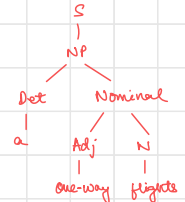
(f) all redeye flights



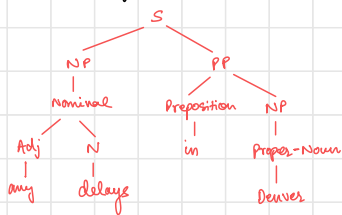
(g) on Thursday



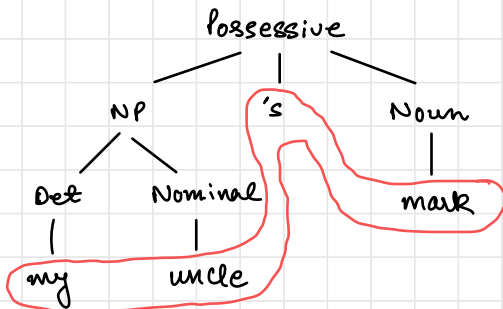
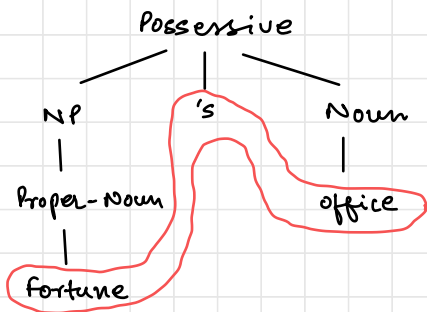
(h) a one-way-flight



(i) my delays in denver



(Q12.6) Possessive \rightarrow NP 's Noun



(Q12.8)

- Step 1: Remove all terminals from right hand side if they occur with a non-terminal.

Example:

$A \rightarrow bC \Rightarrow A \rightarrow BC$
 $B \rightarrow b$

- Step 2: on rules with more than 2 non-terminals on RHS, create new rules and substitute.

Example: $A \rightarrow BCD \Rightarrow A \rightarrow BE$
 $E \rightarrow CD$

chapter 13 : Constituency Parsing

(Q13.3)

function CKY-PARSE (words, grammar) returns table

```

for j ← from 1 to LENGTH(words) do
  for all {A | A → words[j] ∈ grammar}
    table[j-1, j] ← CHECK-UNITS(table[j-1, j] ∪ A)

  for i ← from j-2 downto 0 do
    for k ← i+1 to j-1 do
      for all {A | A → BC ∈ grammar and B ∈ table[i, k]
        and C ∈ table[k, j]}
        table[i, j] ← CHECK-UNITS(table[i, j] ∪ A)
  
```

function CHECK-UNITS (cell) return updated-cell

```

updated-cell = cell
for i ← from 1 to LENGTH(cell)
  B ← cell[i]
  for all {A | A →* B ∈ grammar}
    updated-cell ← cell ∪ A
  
```

S → NP VP
 S → X1 VP
 X1 → Aux NP
 S → VP
 NP → Pronoun
 NP → Proper-Noun
 NP → Det Nominal
 Nominal → Noun
 Nominal → Nominal Noun
 Nominal → Nominal PP
 VP → Verb
 VP → Verb NP
 VP → XL PP
 X2 → Verb NP
 VP → Verb PP
 VP → VP PP
 PP → Preposition NP
 Det → the / this / that
 noun → book / flight / meal / money
 verb → bark / include / prefer
 Pronoun → I / she / me
 Proper-noun → Houston / nowa
 Aux → does
 Preposition → from / to / on / near / through

| BoPR | the | flight | through | Houston |
|----------------------------------|-----|----------------------|---------|------------------------|
| noun, verb ↓ nominal VP, S | | VP, X2 ↓ S | | VP VP X2 ↓ ↓ S S |
| 0,1 | 0,2 | 0,3 | 0,4 | 0,5 |
| | Det | NP | | NP |
| | 1,2 | 1,3 | 1,4 | 1,5 |
| | | Noun ↓ Nominal | | Nominal |
| | | 2,3 | 2,4 | 2,5 |
| | | | Prep | PP |
| | | | 3,4 | 3,5 |
| | | | | Proper-Noun ↓ NP |
| | | | | 4,5 |

Fig: Running above Algo.

(Q13.4) Partial parsing is mostly used in information retrieval systems and its main advantage is the accelerated processing speed. Since you are only parsing chunks of data instead of each word, partial parsing is much faster. Also, as mentioned in section 13.6, eliminating post-head modifiers obviates the need to resolve attachment ambiguities.

But a major disadvantage is probably the fact that you end up losing a lot of valuable information.

- (Q13.5) We can do the following things:
- (i) In case of spelling mistakes, we can flag the incorrect words and use string matching algorithm to correct the word.
 - (ii) Or, we can learn the n -gram probabilities, and use them to predict the closest word to replace the incorrect word.
 - (iii) We can extend (ii) to produce a set of candidate sentences, and choose one with the highest probability.