# Natural Language Processing

## Regular Expressions

### Regular Expressions: Disjunctions

- Letters inside square brackets []

| Pattern | Matches |
|---|---|
| [wW]oodchuck | Woodchuck, woodchuck |
| [1234567890] | Any digit |

- Ranges [A-Z]

| Pattern | Matches | |
|---|---|---|
| [A-Z] | An upper case letter | Drenched Blossoms |
| [a-z] | A lower case letter | my beans were impatient |
| [0-9] | A single digit | Chapter 1: Down the Rabbit Hole |

### Regular Expressions: Negation in Disjunction

- Negations [^Ss]
  - Carat means negation only when first in []

| Pattern | Matches | |
|---|---|---|
| [^A-Z] | Not an upper case letter | Oyfn pripetchik |
| [^Ss] | Neither 'S' nor 's' | I have no exquisite reason" |
| [^e^] | Neither e nor ^ | Look here |
| a^b | The pattern a carat b | Look up a^b now |

1/34

# Note:

- "^" inside the bracket and outside the bracket acts differently
- Caret outside the bracket means start of the string



Dan Jurafsky

## Regular Expressions: More Disjunction

- Woodchucks is another name for groundhog!
- The pipe | for disjunction

| Pattern | Matches |
|---|---|
| groundhog\|woodchuck | |
| yours\|mine | yours<br>mine |
| a\|b\|c | = [abc] |
| [gG]roundhog\|[Ww]oodchuck | |

- Pipe "|" stands for or



Dan Jurafsky

## Regular Expressions: ?  *  +  .

| Pattern | Matches | | | | |
|---|---|---|---|---|---|
| colou?r | Optional previous char | color | colour | | |
| oo*h! | 0 or more of previous char | oh! | ooh! | oooh! | ooooh! |
| o+h! | 1 or more of previous char | oh! | ooh! | oooh! | ooooh! |
| baa+ | | baa | baaa | baaaa | baaaaa |
| beg.n | | begin | begun | begun | beg3n |

Stephen C Kleene

Kleene *,  Kleene +

## Note:

• The character in the last example in place of "." can't be empty



• Caret stands for start of the string
• Dollar stands for end of string

## Cheat Sheet

**Character classes**

| | |
|---|---|
| . | any character except newline |
| \w \d \s | word, digit, whitespace |
| \W \D \S | not word, digit, whitespace |
| [abc] | any of a, b, or c |
| [^abc] | not a, b, or c |
| [a-g] | character between a & g |

**Anchors**

| | |
|---|---|
| ^abc$ | start / end of the string |
| \b | word boundary |

**Escaped characters**

| | |
|---|---|
| \. \* \\ | escaped special characters |
| \t \n \r | tab, linefeed, carriage return |
| \u00A9 | unicode escaped © |

**Groups & Lookaround**

| | |
|---|---|
| (abc) | capture group |
| \1 | backreference to group #1 |
| (?:abc) | non-capturing group |
| (?=abc) | positive lookahead |
| (?!abc) | negative lookahead |

**Quantifiers & Alternation**

| | |
|---|---|
| a* a+ a? | 0 or more, 1 or more, 0 or 1 |
| a{5} a{2,} | exactly five, two or more |
| a{1,3} | between one & three |
| a+? a{2,}? | match as few as possible |
| ab|cd | match ab or cd |

## *Tokenization*

- Fragments ⇒ main- mainly
- Filled pauses ⇒ uh

Dan Jurafsky

# How many words?

- I do uh main- mainly business data processing
  - Fragments, filled pauses
- Seuss's cat in the hat is different from other cats!
  - **Lemma**: same stem, part of speech, rough word sense
    - cat and cats = same lemma
  - **Wordform**: the full inflected surface form
    - cat and cats = different wordforms

01:21 ━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 14:25

---

Dan Jurafsky

# How many words?

$N$ = number of tokens

$V$ = vocabulary = set of types

  $|V|$ is the size of the vocabulary

Church and Gale (1990): $|V| > O(N^{\frac{1}{2}})$

| | Tokens = N | Types = \|V\| |
|---|---|---|
| Switchboard phone conversations | 2.4 million | 20 thousand |
| Shakespeare | 884,000 | 31 thousand |
| Google N-grams | 1 trillion | 13 million |

04:41 ━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 14:25

Media  Playback  Audio  Video  Subtitle  Tools  View  Help

Dan Jurafsky

# Issues in Tokenization

- Finland's capital          →   Finland Finlands Finland's  ?
- what're, I'm, isn't        →   What are, I am, is not
- Hewlett-Packard            →   Hewlett Packard ?
- state-of-the-art           →   state of the art ?
- Lowercase              →   lower-case lowercase lower case ?
- San Francisco          →   one token or two?
- m.p.h., PhD.            →   ??

09:31                                                                    14:25

---

Media  Playback  Audio  Video  Subtitle  Tools  View  Help

Dan Jurafsky

# Maximum Matching
# Word Segmentation Algorithm

- Given a wordlist of Chinese, and a string.
1) Start a pointer at the beginning of the string
2) Find the longest word in dictionary that matches the string starting at pointer
3) Move the pointer over the word in string
4) Go to 2

12:39                                                                    14:25

This algorithm doesn't work for English

## *Normalization*

Dan Jurafsky

# Normalization

- Need to "normalize" terms
  - Information Retrieval: indexed text & query terms must have same form.
    - We want to match *U.S.A.* and *USA*
- We implicitly define equivalence classes of terms
  - e.g., deleting periods in a term
- Alternative: asymmetric expansion:
  - Enter: *window*          Search: *window, windows*
  - Enter: *windows*         Search: *Windows, windows, window*
  - Enter: *Windows*         Search: *Windows*
- Potentially more powerful, but less efficient

01:07 ————————————————————— 11:47

Generally it is done to make SEARCHING easy

## Case Folding

Dan Jurafsky

# Case folding

- Applications like IR: reduce all letters to lower case
  - Since users tend to use lower case
  - Possible exception: upper case in mid-sentence?
    - e.g., *General Motors*
    - *Fed* vs. *fed*
    - *SAIL* vs. *sail*
- For sentiment analysis, MT, Information extraction
  - Case is helpful (*US* versus *us* is important)

01:56 ──────────────────────────────────────────── 11:47

*Lemmatization*

Dan Jurafsky

# Lemmatization

- Reduce inflections or variant forms to base form
  - *am, are, is → be*
  - *car, cars, car's, cars' → car*
- *the boy's cars are different colors → the boy car be different color*
- Lemmatization: have to find correct dictionary headword form
- Machine translation
  - Spanish quiero ('I want'), quieres ('you want') same lemma as querer 'want'

02:50                                                                  11:47

---

## Morphology

Dan Jurafsky

# Morphology

- **Morphemes**:
  - The small meaningful units that make up words
  - **Stems**: The core meaning-bearing units
  - **Affixes**: Bits and pieces that adhere to stems
    - Often with grammatical functions

03:37                                                                  11:47

For example:

In the string "Stems" - "stem" is a stem and "s" is an affix

## *Stemming*



Stemming is a simplified version of Lemmatization

## Porter's Algorithm

Dan Jurafsky

# Porter's algorithm
# The most common English stemmer

**Step 1a**

| | | | | |
|---|---|---|---|---|
| sses | → ss | caresses | → | caress |
| ies | → i | ponies | → | poni |
| ss | → ss | caress | → | caress |
| s | → ø | cats | → | cat |

**Step 1b**

| | | | |
|---|---|---|---|
| (*v*)ing | → ø | walking | → walk |
| | | sing | → sing |
| (*v*)ed | → ø | plastered | → plaster |

...

**Step 2 (for long stems)**

ational→ ate    relational→ relate
izer→ ize       digitizer → digitize
ator→ ate       operator → operate

...

**Step 3 (for longer stems)**

| | | | |
|---|---|---|---|
| al | → ø | revival | → reviv |
| able | → ø | adjustable | → adjust |
| ate | → ø | activate | → activ |

...

---

Dan Jurafsky

# Viewing morphology in a corpus
# Why only strip –ing if there is a vowel?

$$(*v*)ing \rightarrow ø \quad walking \quad \rightarrow walk$$
$$sing \quad \rightarrow sing$$

33

*Sentence Segmentation*



Dan Jurafsky
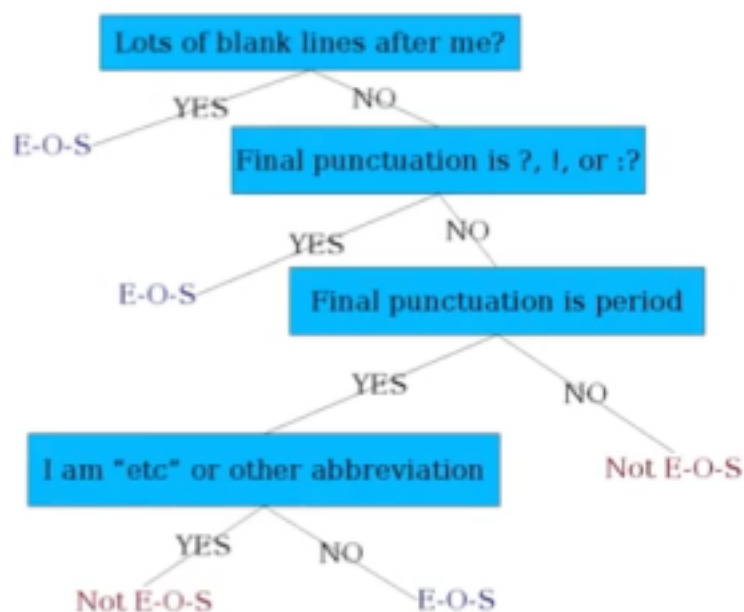
# Sentence Segmentation

- !, ? are relatively unambiguous
- Period "." is quite ambiguous
  - Sentence boundary
  - Abbreviations like Inc. or Dr.
  - Numbers like .02% or 4.3
- Build a binary classifier
  - Looks at a "."
  - Decides EndOfSentence/NotEndOfSentence
  - Classifiers: hand-written rules, regular expressions, or machine-learning

Dan Jurafsky

# Determining if a word is end-of-sentence: a Decision Tree

Lots of blank lines after me?

YES → E-O-S

NO → Final punctuation is ?, !, or :?

YES → E-O-S

NO → Final punctuation is period

YES → I am "etc" or other abbreviation

NO → Not E-O-S

YES → Not E-O-S

NO → E-O-S

---

Dan Jurafsky

# More sophisticated decision tree features

- Case of word with ".": Upper, Lower, Cap, Number
- Case of word after ".": Upper, Lower, Cap, Number

- Numeric features
  - Length of word with "."
  - Probability(word with "." occurs at end-of-s)
  - Probability(word after "." occurs at beginning-of-s)

# *Minimum Edit Distance*



**Edit Distance**

Dan Jurafsky

- The minimum edit distance between two strings
- Is the minimum number of editing operations
  - Insertion
  - Deletion
  - Substitution
- Needed to transform one into the other

Dan Jurafsky

# Defining Min Edit Distance (Levenshtein)

- Initialization

  $D(i,0) = i$

  $D(0,j) = j$

- Recurrence Relation:

  For each   i = 1...M

       For each   j = 1...N

$$D(i,j)= \min \begin{cases} D(i-1,j) + 1 \\ D(i,j-1) + 1 \\ D(i-1,j-1) + \begin{array}{l} 2; & \text{if } X(i) \neq Y(j) \\ 0; & \text{if } X(i) = Y(j) \end{array} \end{cases}$$

- Termination:

  $D(N,M)$ is distance

01:04                                                                              05:54

125%

16/34

# Minimum Edit Distance



|   |   | a | b | c | d | e | f |
|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| a | 1 | 0 | 1 | 2 | 3 | 4 | 5 |
| j | 2 | 1 | 1 | 2 | 3 | 4 | 5 |
| c | 3 | 2 | 2 | 1 | 2 | 3 | 4 |
| e | 4 | 3 | 3 | 2 | 2 | 2 | 3 |
| d | 5 | 4 | 4 | 3 | 2 | 3 | 3 |

abcdef
ajced

f → d

dx

b → 3

## Note:

• If the comparing characters are different then the value will be minimum of (left value or top value or left diagonal value) plus 1
• if the comparing characters are same then the value will be equal to the left diagonal value

## *Language Modelling*

WHY?

Dan Jurafsky

# Probabilistic Language Models

- Today's goal: assign a probability to a sentence
    - Machine Translation:
        - P(**high** winds tonite) > P(**large** winds tonite)

Why?
    - Spell Correction
        - The office is about fifteen **minuets** from my house
            - P(about fifteen **minutes** from) > P(about fifteen **minuets** from)
    - Speech Recognition
        - P(I saw a van) >> P(eyes awe of an)
    - + Summarization, question-answering, etc., etc.!!

01:02 ──────────────────────────── 08:41

---

Dan Jurafsky

# Probabilistic Language Modeling

- Goal: compute the probability of a sentence or sequence of words:

$$P(W) = P(w_1, w_2, w_3, w_4, w_5 \ldots w_n)$$

- Related task: probability of an upcoming word:

$$P(w_5 | w_1, w_2, w_3, w_4) \qquad P(w_1 w_2 v_3 w_4 w_5 -)$$

- A model that computes either of these:

$$P(W) \quad \text{or} \quad P(w_n | w_1, w_2 \ldots w_{n-1})$$    is called a **language model**.

01:50 ──────────────────────────── 08:41

# Compute the Joint Probability of a sentence or a sequence of words

P (its, water, is, so, transparent, that)

This can be calculated using **Chain Rule of Probability**



## Reminder: The Chain Rule

- Recall the definition of conditional probabilities

$$P(A|B) = \frac{P(A,B)}{P(B)}$$

Rewriting: $P(A|B)P(B) = P(A,B)$

$$P(A,B) = P(A|B)P(B)$$

- More variables:

$$P(A,B,C,D) = P(A)P(B|A)P(C|A,B)P(D|A,B,C)$$

- The Chain Rule in General

$$P(x_1,x_2,x_3,...,x_n) = P(x_1)P(x_2|x_1)P(x_3|x_1,x_2)...P(x_n|x_1,...,x_{n-1})$$

Dan Jurafsky

# The Chain Rule applied to compute joint probability of words in sentence

$$P(w_1 w_2 \ldots w_n) = \prod_i P(w_i \mid w_1 w_2 \ldots w_{i-1})$$

P("its water is so transparent") =

P(its) × P(water|its) × P(is|its water)

× P(so|its water is) × P(transparent|its water is so)

04:09 ──────────────────────────────── 08:41

---

How to calculate this probability?

---

Dan Jurafsky

# How to estimate these probabilities

- Could we just count and divide?

$$P(\text{the} \mid \text{its water is so transparent that}) =$$
$$\frac{Count(\text{its water is so transparent that the})}{Count(\text{its water is so transparent that})}$$

- No!  Too many possible sentences!
- We'll never see enough data for estimating these

04:57 ──────────────────────────────── 08:41

Markov Assumption is used to calculate this probability

## *Markov Assumption*

Dan Jurafsky

## Markov Assumption

Andrei Markov

- Simplifying assumption:

$$P(\text{the} \mid \text{its water is so transparent that}) \approx P(\text{the} \mid \text{that})$$

- Or maybe

$$P(\text{the} \mid \text{its water is so transparent that}) \approx P(\text{the} \mid \text{transparent that})$$

---

Dan Jurafsky

## Markov Assumption

$$P(w_1 w_2 \ldots w_n) \approx \prod_i P(w_i \mid w_{i-k} \ldots w_{i-1})$$

- In other words, we approximate each component in the product

$$P(w_i \mid w_1 w_2 \ldots w_{i-1}) \approx P(w_i \mid w_{i-k} \ldots w_{i-1})$$

Dan Jurafsky

# Simplest case: Unigram model

$$P(w_1 w_2 \ldots w_n) \approx \prod_i P(w_i)$$

Some automatically generated sentences from a unigram model

```
fifth, an, of, futures, the, an, incorporated, a,
a, the, inflation, most, dollars, quarter, in, is,
mass

thrift, did, eighty, said, hard, 'm, july, bullish

that, or, limited, the
```

In Unigram Model, the words are independent of each other. So the sentence apparently doesn't make sense.

Dan Jurafsky

# Bigram model

- Condition on the previous word:

$$P(w_i \mid w_1 w_2 \ldots w_{i-1}) \approx P(w_i \mid w_{i-1})$$

```
texaco, rose, one, in, this, issue, is, pursuing, growth, in,
a, boiler, house, said, mr., gurria, mexico, 's, motion,
control, proposal, without, permission, from, five, hundred,
fifty, five, yen

outside, new, car, parking, lot, of, the, agreement, reached

this, would, be, a, record, november
```

In Bigram Model the word is dependent on its previous word.



Like in the above example, N-gram model can fail with sentences having long distance dependencies.
Because the word "crashed" is not related to floor. Here the actual subject is computer. Hence the N-gram model will fail here.
But more or less it works OK for most of the sentences.

# An example

Estimating Bigram Probabilities

$$P(w_i \mid w_{i-1}) = \frac{c(w_{i-1}, w_i)}{c(w_{i-1})}$$

<s> I am Sam </s>
<s> Sam I am </s>
<s> I do not like green eggs and ham </s>

$P(\text{I} \mid <s>) = \frac{2}{3} = .67$    $P(\text{Sam} \mid <s>) = \frac{1}{3} = .33$    $P(\text{am} \mid \text{I}) = \frac{2}{3} = .67$

$P(</s> \mid \text{Sam}) = \frac{1}{2} = 0.5$    $P(\text{Sam} \mid \text{am}) = \frac{1}{2} = .5$    $P(\text{do} \mid \text{I}) = \frac{1}{3} = .33$

$$\frac{\text{Count ( am, Sam )}}{\text{Count ( am )}} = \frac{1}{2}$$

02:05 ———————————————————————— 09:38

# Practical Issues

- We do everything in log space
  - Avoid underflow    Multiplying Numbers like 0.002 * 0.00001
  - (also adding is faster than multiplying)

$$p_1 \times p_2 \times p_3 \times p_4 = \log p_1 + \log p_2 + \log p_3 + \log p_4$$

07:19 ———————————————————————— 09:38

## *Evaluation of Model*

## Model can be evaluated in two ways

1. Extrinsic
2. Intrinsic

Extrinsic evaluation is done by using external data and comparing the accuracy. Like using the unseen test data and checking the accuracy on it.



 Intrinsic evaluation is the evaluation of the model itself rather than any external application

Dan Jurafsky

# Difficulty of extrinsic (in-vivo) evaluation of N-gram models

- Extrinsic evaluation
  - Time-consuming; can take days or weeks
- So
  - Sometimes use **intrinsic** evaluation: **perplexity**
  - Bad approximation
    - unless the test data looks **just** like the training data
    - So **generally only useful in pilot experiments**

A type of intrinsic evaluation is Perplexity
Perplexity is the probability of the test set, normalized by the number of words.

Dan Jurafsky

# Perplexity

The best language model is one that best predicts an unseen test set

- Gives the highest P(sentence)

Perplexity is the probability of the test set, normalized by the number of words:

$$PP(W) = P(w_1 w_2 \ldots w_N)^{-\frac{1}{N}}$$

$$= \sqrt[N]{\frac{1}{P(w_1 w_2 \ldots w_N)}}$$

Chain rule: $\quad PP(W) = \sqrt[N]{\prod_{i=1}^{N} \frac{1}{P(w_i | w_1 \ldots w_{i-1})}}$

For bigrams: $\quad PP(W) = \sqrt[N]{\prod_{i=1}^{N} \frac{1}{P(w_i | w_{i-1})}}$

Minimizing perplexity is the same as maximizing probability

06:56 ————————————————————————— 11:09

---

Dan Jurafsky

# The Shannon Game intuition for perplexity

- From Josh Goodman
- How hard is the task of recognizing digits '0,1,2,3,4,5,6,7,8,9'
  - Perplexity 10
- How hard is recognizing (30,000) names at Microsoft.
  - Perplexity = 30,000
- If a system has to recognize
  - Operator (1 in 4)
  - Sales (1 in 4)
  - Technical Support (1 in 4)
  - 30,000 names (1 in 120,000 each)     *The perplexity is actually 53, not 54*
  - Perplexity is 54
- Perplexity is weighted equivalent branching factor

*average branching factor*

08:45 ————————————————————————— 11:09

Dan Jurafsky

# Perplexity as branching factor

- Let's suppose a sentence consisting of random digits
- What is the perplexity of this sentence according to a model that assign P=1/10 to each digit?

$$PP(W) = P(w_1 w_2 \ldots w_N)^{-\frac{1}{N}}$$
$$= \left(\frac{1}{10}^N\right)^{-\frac{1}{N}}$$
$$= \frac{1}{10}^{-1}$$
$$= 10$$

$$P(w_1 w_2 \cdot w_N)^{-\frac{1}{N}}$$
$$P(\tfrac{1}{10} \cdot \tfrac{1}{10} \cdot \tfrac{1}{10} \cdot \tfrac{1}{10}$$
$$\left(\tfrac{1}{10}^N\right)^{-\frac{1}{N}}$$

10:03                                                                                    11:09

---

Dan Jurafsky

# Lower perplexity = better model

- Training 38 million words, test 1.5 million words, WSJ

| N-gram Order | Unigram | Bigram | Trigram |
|---|---|---|---|
| Perplexity | 962 | 170 | 109 |

10:29                                                                                    11:09

# Generalization

## Zeros

**Dan Jurafsky**

- Training set:
  ... denied the allegations
  ... denied the reports
  ... denied the claims
  ... denied the request

  $P(\text{"offer"} \mid \text{denied the}) = 0$

- Test set
  ... denied the offer
  ... denied the loan

## Zero Probability Problem

When the sentences in the test set have never appeared in the training set, the probability assigned to it is zero.
This produces Zero Probability Problem.

## Smoothing (Add One Estimate)

To avoid Zero Probability Problem, Smoothing or Add One Estimate technique is used.

Dan Jurafsky

## The intuition of smoothing (from Dan Klein)

- When we have sparse statistics:

  P(w | denied the)
  - 3 allegations
  - 2 reports
  - 1 claims
  - 1 request
  - 7 total

- Steal probability mass to generalize better
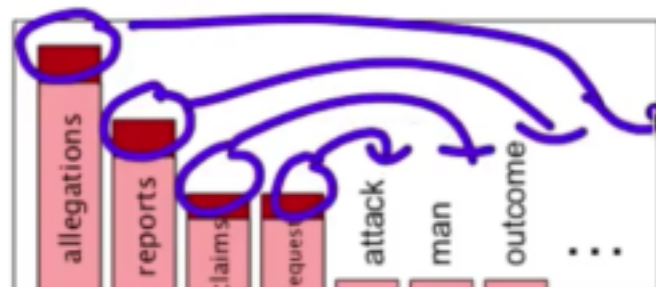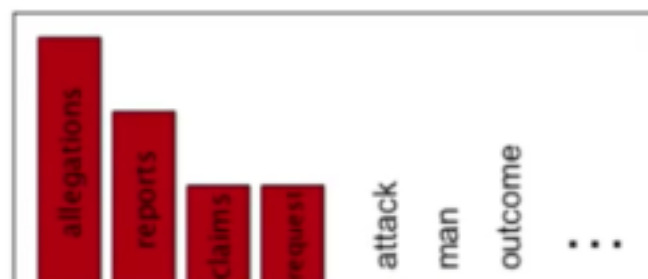
  P(w | denied the)
  - 2.5 allegations
  - 1.5 reports
  - 0.5 claims
  - 0.5 request
  - 2 other
  - 7 total



01:04 ──────────────────────── 06:30

---

Dan Jurafsky

## Add-one estimation

- Also called Laplace smoothing
- Pretend we saw each word one more time than we did
- Just add one to all the counts!

- MLE estimate:

$$P_{MLE}(w_i \mid w_{i-1}) = \frac{c(w_{i-1}, w_i)}{c(w_{i-1})}$$

- Add-1 estimate:

$$P_{Add-1}(w_i \mid w_{i-1}) = \frac{c(w_{i-1}, w_i) + 1}{c(w_{i-1}) + V}$$

02:06 ──────────────────────── 06:30

# Add-1 estimation is a blunt instrument

Dan Jurafsky

- So add-1 isn't used for N-grams:
  - We'll see better methods
- But add-1 is used to smooth other NLP models
  - For text classification
  - In domains where the number of zeros isn't so huge.