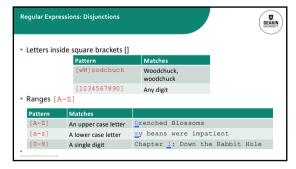
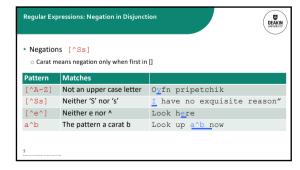


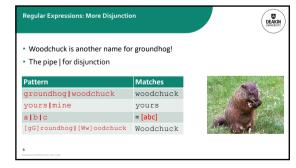
Regular expressions

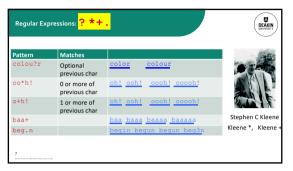
• A formal language for specifying text strings
• How can we search for any of these?
• woodchuck
• woodchucks
• Woodchuck
• Woodchucks

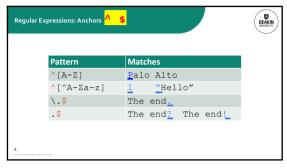
2











• Find me all instances of the word "the" in a text.

the

Misses capitalized examples

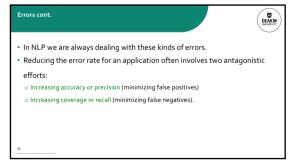
[tT]he

Incorrectly returns other or theology

[^a-zA-Z][tT]he[^a-zA-Z]

7 8

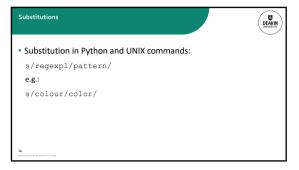




Regular expressions play a surprisingly large role
 Sophisticated sequences of regular expressions are often the first model for any text processing text
 For hard tasks, we use machine learning classifiers
 But regular expressions are still used for pre-processing, or as features in the classifiers
 Can be very useful in capturing generalizations

10 11 12





• Say we want to put angles around all numbers:

the 35 boxes → the <35> boxes

• Use parens () to "capture" a pattern into a numbered register (1, 2, 3...)

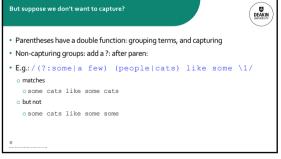
• Use \1 to refer to the contents of the register

s/([0-9]+)/<\1>/

13

14 15

• /the (.*)er they (.*), the \ler we \2/
• Matches
the faster they ran, the faster we ran
• But not
the faster they ran, the faster we ate



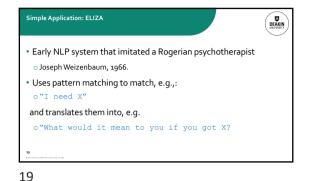
• (?= pattern) is true if pattern matches, but is zero-width;
doesn't advance character pointer

• (?! pattern) true if a pattern does not match

• How to match, at the beginning of a line, any single word that doesn't start with "Volcano":

• /^ (?!Volcano) [A-Za-z]+/

16 17 18





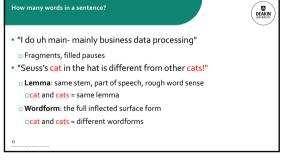
* s/.* I'M (depressed|sad) .*/I AM SORRY TO HEARYOU ARE \1/
* s/.* I AM (depressed|sad) .*/WHY DO YOU THINK YOU ARE \1/
* s/.* all .*/IN WHAT WAY?/
* s/.* always .*/CAN YOU THINK OF A SPECIFIC EXAMPLE?/

21

SIT330-770: Natural Language Processing
Week 3.3 - Words and Corpora

Dr. Mohamed Reda Bouadjenek
School of Information Technology,
Faculty of Sci Eng & Built Env

22



they lay back on the San Francisco grass and looked at the stars and their

Type: an element of the vocabulary.

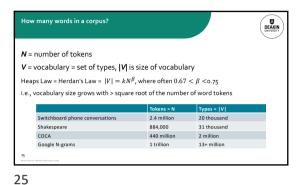
Token: an instance of that type in running text.

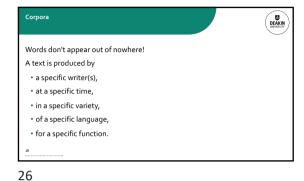
How many?

15 tokens (0r 14)

13 types (0r 12) (0r 11?)

23 24





Corpora vary along dimension like DEAKIN UNIVERSITY • Language: 7097 languages in the world · Variety, like African American Language varieties. AAE Twitter posts might include forms like "iont" (I don't) Code switching, e.g., Spanish/English, Hindi/English: S/E: Por primera vez veo a @username actually being hateful! It was beautiful:) [For the first time I get to see @username actually being hateful! it was beautiful:)] H/E: dost tha or ra- hega ... dont wory ... but dherya rakhe ["he was and will remain a friend ... don't worry ... but have faith"] · Genre: newswire, fiction, scientific articles, Wikipedia • Author Demographics: writer's age, gender, ethnicity, SES

27

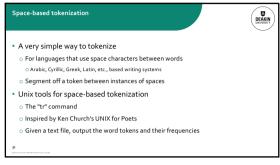
Corpus datasheets DEAKIN UNIVERSITY Gebru et al (2020), Bender and Friedman (2018) Motivation: · Why was the corpus collected? By whom?Who funded it? Situation: In what situation was the text written? Collection process: If it is a subsample how was it sampled? Was there consent? Pre-+Annotation process, language variety, demographics, etc.

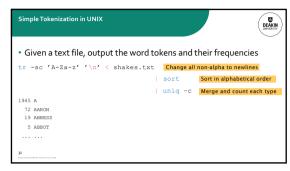
28



Text Normalization DEAKIN UNIVERSITY • Every NLP task requires text normalization: 1. Tokenizing (segmenting) words 2. Normalizing word formats 3. Segmenting sentences

29 30





The first step: tokenizing

EAKIN

EXAMPLE 1

THE

SOUNDETS

by

William

Shakespears

From

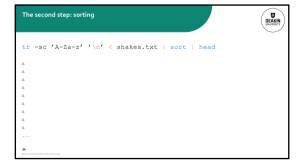
fairest

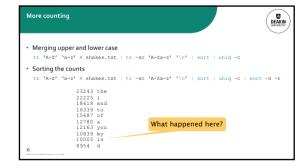
creatures

We

...

31 32 33



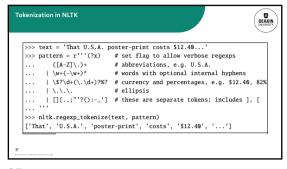


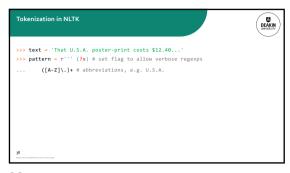
Can't just blindly remove punctuation:

o mp.h., Ph.D., AT&T, cap'n
o prices (14,5,55)
odates (solva2/o6)
o URLs (http://www.stanford.edu)
o hashtags (#nlproc)
o email addresses (someone@ics.colorado.edu)
Clitic: a word that doesn't stand on its own
o "are" in we're, French "je" in ji, "le" in Irhonneur

When should multiword expressions (MWE) be words?
o New York, rock 'n' roll

34 35 36





Many languages (like Chinese, Japanese, Thai) don't use spaces to separate words!

 How do we decide where the token boundaries should be?

37 38

Chinese words are composed of characters called "hanzi" (or sometimes just "zi")
 Each one represents a meaning unit called a morpheme.
 Each word has on average 2.4 of them.
 But deciding what counts as a word is complex and not agreed upon.

40



So, in Chinese it's common to just treat each character (zi) as a token.
So, the segmentation step is very simple
In other languages (like Thai and Japanese), more complex word segmentation is required.
The standard algorithms are neural sequence models trained by supervised machine learning.

41 42



Instead of
 white-space segmentation
 single-character segmentation
 single-character segmentation
 Subword tokenization (because tokens can be parts of words as well as whole words)

Three common algorithms:

Byte-Pair Encoding (BPE) (Sennrich et al., 2016)

Unigram language modeling tokenization (Kudo, 2018)

WordPiece (Schuster and Nakajima, 2012)

All have 2 parts:

A token learner that takes a raw training corpus and induces a vocabulary (a set of tokens).

A token segmenter that takes a raw test sentence and tokenizes it according to that vocabulary

45

43

44

• Let vocabulary be the set of all individual characters

= {A, B, C, D,..., a, b, c, d....}

• Repeat:

• Choose the two symbols that are most frequently adjacent in the training corpus (say 'A', 'B')

• Add a new merged symbol 'AB' to the vocabulary

• Replace every adjacent 'A' 'B' in the corpus with 'AB'.

• Until & merges have been done.

Function ByTE-PAIR ENCODING(strings C, number of merges k) returns vocab V $V \leftarrow$ all unique characters in C # initial set of tokens is characters for i = 1 to k do # merge tokens til k times t_L , $t_R \leftarrow$ Most frequent pair of adjacent tokens in C $t_{NEW} \leftarrow t_L + t_R$ # make new token by concatenating $V \leftarrow V + t_{NEW}$ # update the vocabulary

Replace each occurrence of t_L , t_R in C with t_{NEW} # and update the corpus return V

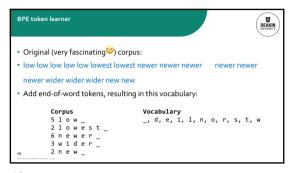
Byte Pair Encoding (BPE) Addendum

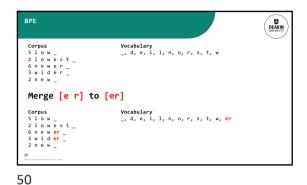
 Most subword algorithms are run inside space-separated tokens.

 So we commonly first add a special end-of-word symbol '__' before space in training corpus

 Next, separate into letters.

47 48





DEAKIN UNIVERSITY Vocabulary _, d, e, i, 1, n, o, r, s, t, w, er 5 1 o w _ 2 1 o w e s t _ 6 n e w er _ 3 w i d er _ 2 n e w _ Merge [er _] to [er_] Vocabulary
, d, e, i, l, n, o, r, s, t, w, er, er 5 1 o w _ 2 1 o w e s t _ 6 n e w er_ 3 w i d er_ 2 n e w _

51

DEAKIN UMVERSITY

49

Vocabulary
, d, e, i, 1, n, o, r, s, t, w, er, er

, d, e, i, l, n, o, r, s, t, w, er, er, ne

Vocabulary

DEAKIN UNIVERSITY

• The next merges are: Merge Current Vocabulary _, d, e, i, 1, n, o, r, s, t, w, er, er_, ne, new (ne, w) _, d, e, i, 1, n, o, r, s, t, w, er, er_, ne, new, lo (1, o) (lo, w) _, d, e, i, 1, n, o, r, s, t, w, er, er_, ne, new, lo, low (new, er_) _, d, e, i, 1, n, o, r, s, t, w, er, er_, ne, new, lo, low, newer_ (low, _) _, d, e, i, 1, n, o, r, s, t, w, er, er_, ne, new, lo, low, newer_, low_

BPE token segmenter algorithm DEAKIN UNIVERSITY On the test data, run each merge learned from the training data: o Greedily o In the order we learned them o (test frequencies don't play a role) So: merge every [e r] to [er], then merge [er _] to [er_], etc. • Result: o Test set "n e w e r _" would be tokenized as a full word o Test set "1 o w e r _" would be two tokens: "low er_"

52

Cornus 5 low_ 2 lowest_

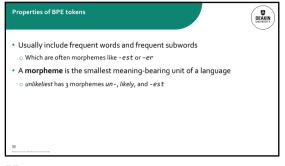
Corpus

5 1 o w _ 2 1 o w e s t _

6 ne w er_ 3 w i d er_ 2 ne w _

6 n e w er_ 3 w i d er_ 2 n e w _

Merge [n e] to [ne]





Putting words/tokens in a standard format

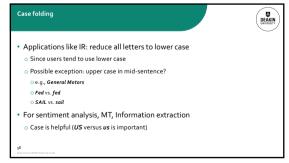
OU.S.A. or USA

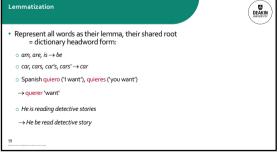
Ouhhuh or uh-huh

OFed or fed

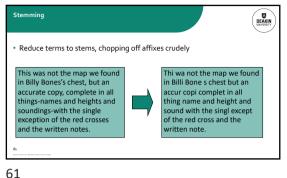
Oam, is, be, are

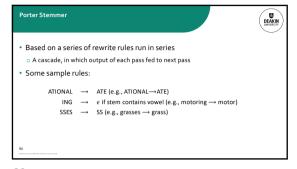
55 56 57





58 59 60





Dealing with complex morphology is necessary for many languages

• e.g., the Turkish word:

• Uygarlastiramadiklarimizdanmissinizcasina

• `(behaving) as if you are among those whom we could not civilize'

• Uygar `civilized' + las `become'

+ tir `cause' + ama `not able'

+ dik `past' + lar `plural'

+ imiz 'papl' + dan 'abl'

+ mis `past' + siniz '2pl' + casina `as if'

63

62

• 1, ? mostly unambiguous but period "." is very ambiguous

• Sentence boundary

• Abbreviations like Inc. or Dr.

• Numbers like .02% or 4.3

• Common algorithm: Tokenize first: use rules or ML to classify a period as either (a) part of the word or (b) a sentence-boundary.

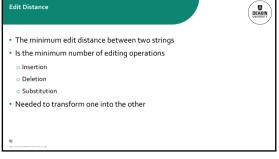
• An abbreviation dictionary can help

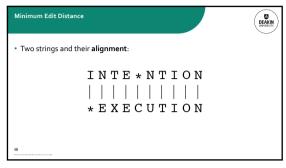
• Sentence segmentation can then often be done by rules based on this tokenization.

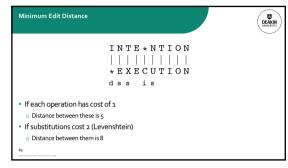
64



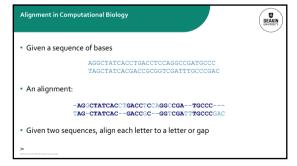
65 66



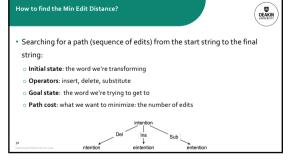




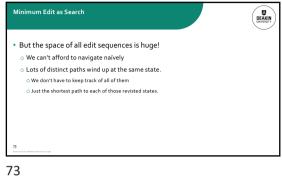
67 68 69







70 71 72



• For two strings

○ X of length n

○ Y of length m

• We define D(i,j)

○ the edit distance between X[1../] and Y[1../]

○ i.e., the first i characters of X and the first j characters of Y

○ The edit distance between X and Y is thus D(n,m)



3 74

Dynamic Programming for Minimum Edit Distance

• Dynamic programming: A tabular computation of D(n,m)

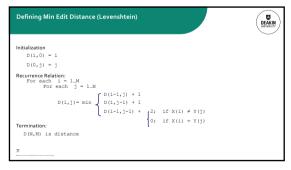
• Solving problems by combining solutions to subproblems.

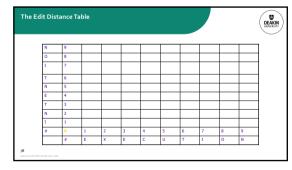
• Bottom-up

○ We compute D(i,j) for small i,j

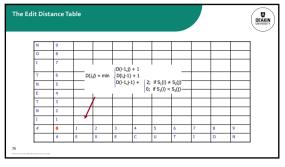
○ And compute larger D(i,j) based on previously computed smaller values

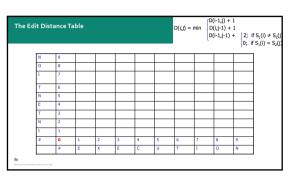
○ i.e., compute D(i,j) for all i (o < i < n) and j (o < j < m)

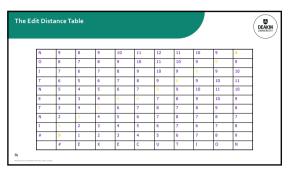




76 77 78

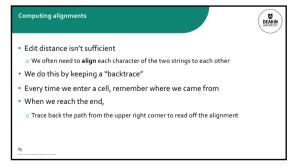


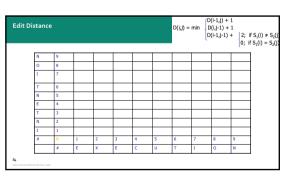




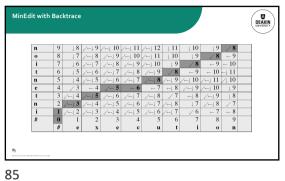
79 80 81

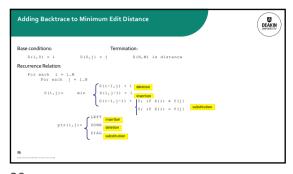


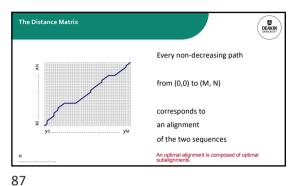




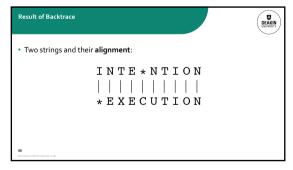
82 83 84







86







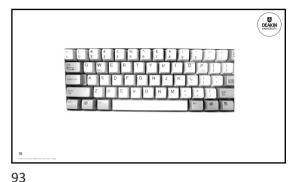
88 89 90

Weighted Edit Distance

Why would we add weights to the computation?

Spell Correction: some letters are more likely to be mistyped than others
Biology: certain kinds of deletions or insertions are more likely than others

| Confusion matrix for spelling errors | SubJix | SubJix



91 92 9

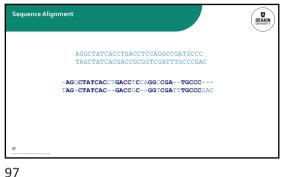
where did the name, dynamic programming, come from?

...The 1950s were not good years for mathematical research. [the] Secretary of Defense ...had a pathological fear and hatred of the word, research...
I decided therefore to use the word, "programming".
I wanted to get across the idea that this was dynamic, this was multistage... I thought, let's ... take a word that has an absolutely precise meaning, namely dynamic... it's impossible to use the word, dynamic, in a pejorative sense. Try thinking of some combination that will possibly give it a pejorative meaning. It's impossible.
Thus, I thought dynamic programming was a good name. It was something not even a Congressman could object to."

Richard Bellman, "Eye of the Hurricane: an autobiography" 1984.



94 95 96



DEAKIN UNIVERSITY Why sequence alignment? · Comparing genes or regions from different species o to find important regions o determine function o uncover evolutionary forces Assembling fragments to sequence DNA · Compare individuals to looking for mutations

DEAKIN UNIVERSITY Alignments in two fields • In Natural Language Processing We generally talk about distance (minimized) OAnd weights • In Computational Biology We generally talk about similarity (maximized) OAnd scores

99

102

DEAKIN UMAYORSITY

98

101

DEAKIN UNIVERSITY

The Needleman-Wunsch Matrix (Note that the origin is at the upper left.)

A variant of the basic algorithm: DEAKIN UNIVERSITY • Maybe it is OK to have an unlimited # of gaps in the beginning and end: -----CTATCACCTGACCTCCAGGCCGATGCCCCTTCCGGC GCGAGTTCATCTATCAC--GACCGC--GGTCG------• If so, we don't want to penalize gaps at the ends

D(i-1,j-1) + s[x(i),y(j)]• Termination: D(N,M) is distance

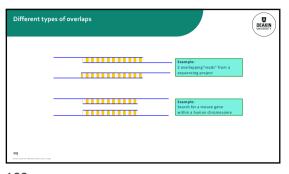
· Initialization:

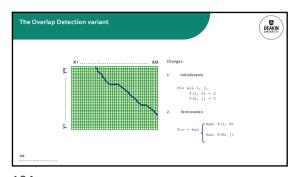
D(i,0) = -i * dD(0,j) = -j * d

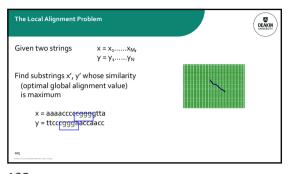
Recurrence Relation:

The Needleman-Wunsch Algorithm

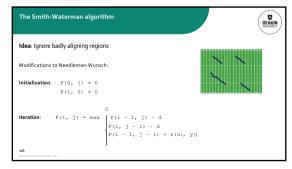
D(i-1,j) - dD(i,j) = min D(i,j-1) - d

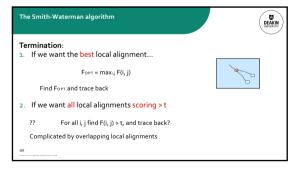




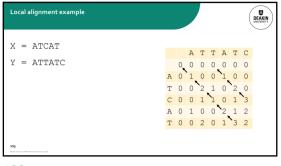


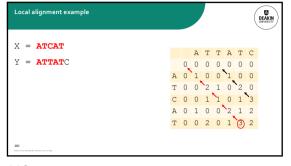
103 104 105

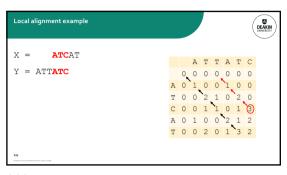




106 107 108







109 110 111