

CS4248: Natural Language Processing

Lecture 2 — Strings & Words

Recap of Week 01



Outline

• Regular Expressions

- Basic Concepts
- Relationship to FSA
- Error Types

• Corpus Preprocessing

- Tokenization
- Normalization
- Stemming / Lemmatization
- Segmentation

• Word error handling

- Spelling Errors
- Minimum Edit Distance
- Noisy Channel Model

2

Regular Expressions

- Regular Expression Definition
 - Search pattern used to match character combinations in a string
 - Pattern = sequence of characters
- Common applications
 - Parse text documents to find specific character patterns
 - Validate text to ensure it matches predefined patterns
 - Extract, edit, replace, delete substrings matching a pattern
- Two basic search approaches
 - Default: match only <u>first</u> occurrence of pattern
 - Global search: match all occurrences of pattern (assumed in most following examples)

Example: password validation



Basic Patterns

• Fixed patterns



• Special characters (metacharacters)

	Character	Explanation
~	g . '	matches any character except line breaks
• >	~ ^	match the start of a string
	\$	match the end of a string
1 S , -	- > 1	matches RegEx either before or after the symbol (e.g., floor floors)
	\b	matches boundary between word and non-word

Character Classes

- Character class
 - Defines set of valid characters
 - Enclosed using "[...]"
 - Can be negated: "(^)..]"
 - [0-9][0-9] →
- My block has 15 floors, and I live on floor 5. (match all sequences of 2 digits)

01234 Aa

→ My block has 15 floors, and I live on floor 5. (match all sequences of length 1 that are either a period, comma, etc.)



[.,;:]

My block has 15 floors, and I live on floor 5.

(match all sequences of length 1 that are not a lowercase letter)

Predefined Character Classes

• Common character classes with their own shorthand notation (i.e., metacharacters)

	Class	Alternative	Explanation	
	\d	[0-9]	matches any digit $TP127456$	787]
•	\D	[^0-9]	matches any non-digit	
	\s	$[n\t]$	matches any whitespace character	
	\\s	$[^ \n\t]$	matches any non-whitespace character	
		[a-zA-Z0-9_]	matches any word character	
	\w	[^a-zA-Z0-9_]	matches any non-word character	

Repetition Patterns

- Very common: patterns with flexible lengths, e.g.:
 - All numbers with more than 2 digits
 - All words with less than 5 characters
- Repetition patterns metacharacters

	Pattern	Explanation
	b +	1 or more occurrences
\sim	-) (*)	0 or more occurrences
	?	0 or 1 occurrences
19	{n}	exactly n occurrences
	{1,u}	between 1 and u occurrences; can be unbounded: {1, } or {, u}

Repetition Patterns — Examples

- $d{2,}$
- My block has 15 floors, and I live on floor 5. (match all numbers with 2 or more digits)
- \d+ → My block has 15 floors, and I live on floor 5. (match all numbers with 1 or more digits)

- My block has 15 floors, and I live on floor 5. (match words with 2 to 4 characters)
- \b[Ff]loor[s]?\b →

My block has 15 floors, and I live on floor 5.

(match occurrences of "floor", either capitalized or not, either in singular or plural)

Groups

Quick quiz: In which case(s) would the RegEx below fail to correctly match an email address?

- Groups: Organizing patterns into parts
 - Groups are enclosed using "(...)"
 - While whole expression must match, groups are captures individually (a match is no longer a string but a tuple of strings, on for each group)
 - Groups can be nested, e.g., (...(...)...((...))...) (order of groups depends on the order in which the groups "open")

Send an email to alice@example.org for more information.



Match:	user@example.org
Group #1:	alice
Group #2:	example.org

Backreferences

Quick quiz: Can the same be achieved using only 1 group?

- Reference groups within a RegEx
 - Find repeated patterns (see example below)
 - Support only partial replacement of matches
- Example:
 - "My mom said I need to pass this test."
 - Goal: Find all words that start and end with the same letter



Match:	тот
Group #1:	тот
Group #2:	m

Match:	test
Group #1:	test
Group #2:	t

Lookarounds

- Special groups assertions
 - Match like any other group, but do not capture the match
 - 2 types: lookaheads and lookbehinds
 - 2 forms of assertion: positive and negative

	Туре	Example
(?=)	positive lookahead	A(?=B) \rightarrow finds expr. A but only when followed by expr. B
(?!)	negative lookahead	A(?!B) \rightarrow finds expr. A but only when not followed by expr. B
(?<=)	positive lookbehind	(?<=B) A → finds expr. A but only when preceded by expr. B
(?)</td <td>negative lookbehind</td> <td>(?<!--B) A → finds expr. A but only when not preceded by expr. B</td--></td>	negative lookbehind	(? B) A → finds expr. A but only when not preceded by expr. B</td

Lookarounds — Example

- Positive lookahead
 - "Paying 10 SGD for 1 kg of chicken seems fair."
 - Goal: Extract all kg values (numbers followed by the unit kg)

 \rightarrow



"Paying 10 SGD for 1 kg of chicken seems fair.

"Paying 10 SGD for 1.5 kg of chicken seems fair.

"Paying 10 SGD for 1,500.00 kg of chicken seems fair.

"Paying 10 SGD for <mark>1</mark> kg of chicken seems fair.

"Paying 10 SGD for 1,500.00 kg of chicken seems fair.

Outline

• Regular Expressions

- Basic Concepts
- Relationship to FSA
- Error Types

• Corpus Preprocessing

- Tokenization
- Normalization
- Stemming / Lemmatization
- Segmentation

• Word error handling

- Spelling Errors
- Minimum Edit Distance
- Noisy Channel Model

2

Pre-Lecture Activity from last week

Pre-Lecture Activity for Next Week

• Assigned Task (due before Jan 20)

Post a 1-2 sentence answer to the following question into the L1 Discussion forum (you will find the thread on Canvas > Discussion > Week 2 (L1))

"What is the relationship between a Finite State Machine and Regular Expressions?"

Side notes:

- This task is meant as a warm-up to provide some context for the next lecture
- No worries if you get lost; we will talk about this in the next lecture
- You can just copy-&-paste others' answers but this won't help you learn better





Google said they are the same. ChatGPT said there are totally different. Me myself think that regular expression could be drawn as state machine using criteria together with their groupings as state transitions. So deadend, spider trap any thing can happen in state graph in both deterministic and non. However, parsing a string to regex would be deterministic at some point.

Regular Expressions can be used to represent a Finite State Machine via a regular expression pattern. A regular expression pattern can be converted into a Finite State Machine directed graph that enable us to perform character matching and understanding a word's morphology.



References:

- <u>https://ljvmiranda921.github.io/notebook/2022/10/07/finite-</u> state-automata/ ▷

- https://www.cs.drexel.edu/~johnsojr/2006-

07/winter/cs360/lectures/lec2.html

YC

Regular expressions describe finite state machines, where instead of actions, we have individual letters/characters. However, FSM implementations of Regex cannot do capture groups, and in general, offer different conveniences/challenges in implementation as compared to direct Regex implementations.

References:

https://bakalian.cs.umd.edu/assets/notes/fa.pdf ⊟

FSM is a mathematical model of computation. According to the different transfer functions, it can be divided into Deterministic Finite Automaton(DFA) and Nondeterministic Finite Automaton(NFA).

LH

Regular expressions are based on the theory of automata. After the user writes the regular expression, the corresponding FSM can be constructed through the regular expression. (This FSM may be DFA or NFA)



Relationship to Finite State Automata

• Equivalence

 Regular Expressions describe Regular Languages (most restricted types of languages w.r.t Chomsky Hierarchy)

Regular Language = language accepted by a FSA

Example: FSA that accepts the Regular Language described by the Regular Expression **I(o+I)+**





Relationship to Finite State Automata

• Basic equivalences



🏃 🏃 🏃 In-Lecture Activity (5 mins)

- Task: Find the RegEx describing the FSA below
 - Post your RegEx to Canvas > Discussions
 (individually or as a group; include all group members' names in the post)
 - Optional: There is more than one correct answer → Why?



Outline

• Regular Expressions

- Basic Concepts
- Relationship to FSA
- Error Types

• Corpus Preprocessing

- Tokenization
- Normalization
- Stemming / Lemmatization
- Segmentation

• Word error handling

- Spelling Errors
- Minimum Edit Distance
- Noisy Channel Model

2

Error Types — What Can Go Wrong

Quick quiz: What would be a better RegEx for this task?

- Example: Find all occurrences of article "the"
 - Naive approach: "the" (fixed pattern)



Error Types

• 2 basic types of errors

Matching strings that we should <u>not</u> have matched (e.g., *other*, *theology*, *weather*, *bathe*, *mother*)

Not matching things that we should have matched (e.g., *THE*)



Error Types — Observations

- Many contexts deal with these 2 types of errors, e.g.:
 - Medical testing (e.g., ART test is positive but person is not infected with COVID → false positive)
- Information retrieval (e.g., a Web search is missing a relevant page → false negative)
 - Document classification (e.g., an abusive tweet has be classified as positive → false positive)

• Reducing errors

false negative

false positive

- Both error types not always equally bad (infected person tests negative vs. healthy person test positive)
- Reducing False Positives and False Negatives often in conflict (reducing False Positives often increases False Negatives, and vice versa)

Regular Expressions — Summary

- Know their powers
 - Extremely useful tool for many (low-level) text processing tasks (e.g., data preprocessing, tokenization, normalization)
 - Important skill for anyone working with strings or text

• Know their limitations

- Regular Expressions represent hard rules
- Higher-level text processing task generally require statistical models ("soft" rules)
- → Machine Learning classifiers

WHAT GIVES PEOPLE FEELINGS OF POWER



Outline

• Regular Expressions

- Basic Concepts
- Relationship to FSA
- Error Types

• Corpus Preprocessing



- Tokenization
- Normalization
- Stemming / Lemmatization
- Segmentation

• Word error handling

- Spelling Errors
- Minimum Edit Distance
- Noisy Channel Model

2

Tokenization

● Tokenization: splitting a string into tokens → vocabulary (set of all unique tokens)

- Token = character sequence with a semantic meaning (typically: words, numbers, punctuation — but may differ depending on applications)
- Very important for step for most NLP algorithms (tokenization errors quickly propagate up → "garbage in, garbage out")
- 3 basic approaches



Character-based tokenization trivial (e.g., using Regex: .)

Tokenization — Word-Based

Quick quiz: What is an important assumption for the 2 approaches?

→ \w+|\d+|[,.;:]

- 2 intuitive approaches (solved using RegEx)
 - Match all words, numbers and punctuation marks
 - Match boundaries between "words" and "non-words" → $(?=\setminus W) | (?<=\setminus W)$

 $w+|d+|[,.;:] \rightarrow NLP$ is fun, and there is so much to learn in 13 weeks.

 $(?=\setminus W) | (?<=\setminus W) \rightarrow NLP | is fun, and there | is so much to learn | in 13 weeks |$

Tokenization — It Quickly Gets Tricky

Multiword phrases

Common contractions

Hyphenations

Acronyms, names, etc.

Special tokens



- → I just came back from New York City.
 → I'm not home, so don't call.
- → NLP is a well-defined but non-trivial topic.
- → I watched a C++ documentary on T.V.
- → My email is chris@nus.comp.nus.sg :o)

RegEx used: \w+|\d+|[,.;:]

Example: spaCy Tokenizer



- (1) Split string on whitespace characters
- (2) From left to right, recursively check substrings:
- Does substring match an exception rule? (e.g., "don't" → "do", "n't", but keep "U.K.")
 - Can a prefix, suffix or infix be split of? (e.g., commas, periods, quotes, hyphens)

Substring checks based on

- Regular Expressions
- Hand-crafted rules / patterns

Example: Chris's Tokenizer

Sequential labeling of characters



→ Tokens = Substrings with adjacent characters with the same labels

Tokenization — Language Issues

- French
 - Different uses of apostrophes and hyphens (compared to English)



- German
 - Very common: compound nouns

Arbeiterunfallversicherungsgesetz "worker injury insurance act"

→ important: compound splitter

Tokenization — Language Issues

Languages without whitespaces separating words



Tokenization — Word Segmentation of Chinese Text

Baseline algorithm: Maximum Matching



Tokenization — Maximum Matching

• Surprisingly good performance on Chinese text (even better performance with probabilistic methods or extensions)

Generally does not work for English text



Tokenization — Subword-Based

- Subword-based tokenization
 - So far: a priori specification of rules (e.g., RegEx) what constitutes valid tokens
 - Now: use data to specify how to tokenize
- Why do we want to do this?

Out Of Vocabulary (OOV) words (word/token an NLP model has not seen before)

Very rare words in corpus

→ problematic when building statistical models



→ Goal: Split OOV and rare words into (some) known & frequent tokens

Tokenization — Subword-Based

- Different algorithms for subword tokenization
 - <u>Byte-Pair Encoding (BPE)</u>, Unigram Language Model Tokenization, WordPiece, etc.

• Different approaches, similar 2-parts setup

1) Token Learner

Takes raw training corpus and induces a vocabulary (i.e., set of tokens)

(2) Token Segmenter

Takes a raw text and tokenizes it according to vocabulary


corpus representation 6 n e w e s t _ 5 l o w _ 3 w i d e s t _ 2 l o w e r _ 1 l o n g e r _

vocabulary

d, e, g, i, l, n, o, r, s, t, w, _

merges

most frequent pair: e & s (9 occurrences)

corpus representation

6	n e w e	st_
5	low_	
3	wid e	st_
2	lowe	r _
1	long	er_

vocabulary

d, e, g, i, l, n, o, r, s, t, w, _, es

merges

(e, s)

most frequent pair: es & t (9 occurrences)

corpus representation

6	n	е	W	es	st	_			
5	1	0	W	_					
3	W	i	d	es	st	_			
2	1	0	W	е	r	_			
1	1	0	n	g	е	r	_		

vocabulary

d, e, g, i, l, n, o, r, s, t, w, _, es, est

merges

(e, s), (es, t)

_ most frequent pair: est & _ (9 occurrences)

corpus representation

6	n	е	W	es	st_	_			
5	1	0	W	_					
3	W	i	d	es	st_	_			
2	1	0	W	е	r	_			
1	1	0	n	g	е	r	_		

vocabulary

d, e, g, i, l, n, o, r, s, t, w, _, es, est, est_

merges

(e, s), (es, t), (est, _)

most frequent pair: 1 & o (8 occurrences)

corpus representation

6	n e w est_
5	lo w _
3	widest_
2	lower_
1	longer_

vocabulary

d, e, g, i, l, n, o, r, s, t, w, _, es, est, est_, **lo**

merges

(e, s), (es, t), (est, _), (1, o)

most frequent pair: lo & w (7 occurrences)

corpus representation

6	n e w est_
5	low _
3	w i d est_
2	low e r _
1	longer_

vocabulary

d, e, g, i, l, n, o, r, s, t, w, _, es, est, est_, lo, low

merges

(e, s), (es, t), (est, _), (l, o), (lo, w)

most frequent pair: n & e (6 occurrences)

- vocabulary d, e, g, i, l, n, o, r, s, t, w, _, es, est, est_, lo, low, ne
 - merges (e, s), (es, t), (est, _), (l, o), (lo, w), (n, e)



vocabulary d, e, g, i, l, n, o, r, s, t, w, _, es, est, est_, lo, low, ne, new
merges (e, s), (es, t), (est, _), (l, o), (lo, w), (n, e), (ne, w)



vocabulary d, e, g, i, l, n, o, r, s, t, w, _, es, est, est_, lo, low, ne, new, newest_
merges (e, s), (es, t), (est, _), (l, o), (lo, w), (n, e), (ne, w), (new, est_)

Tokenization — BPE Token Segmenter

vocabulary	d, e, g, i, l, n, o, r, s, t, w, _, es, est, est_, lo, low, ne, new, newest_,
	low_,er,er_,wi,wid,widest_,lower_,lon,long,longer_
merges	(e, s), (es, t), (est, _), (l, o), (lo, w), (<u>n, e</u>), (ne, w), (new, est_), (low, _), (e, r),
	(er, _), (w, i), (wi, d), (wid, est_), (low, er_), (lo, n), (lon, g), (long, er_)



Tokenization — Summary

- Tokenization as low-level NLP task
 - Challenges: important, non-trivial, language-dependent
 - Particularly tricky for informal language (e.g., social media)
- 3 basic approaches
 - Character-based (trivial to do but often not suitable individual characters generally carry no semantic meaning)
 - Word-based (a priori specification of rules; language-dependent; problem: OOV/rare words)
 - Subword-based (tokenization learned from data tokens are often morphemes!)
- Practical consideration (when using off-the-shell word-based tokenizers)
 - What is my type of text (e.g., formal or informal)? Are there special tokens (e.g., URLs, hashtags)?
 - Try and assess different tokenizers very, very last resort: write your own tokenizer

Outline

• Regular Expressions

- Basic Concepts
- Relationship to FSA
- Error Types

Corpus Preprocessing

- Tokenization
- Normalization
- Stemming / Lemmatization
- Segmentation

• Word error handling

- Spelling Errors
- Minimum Edit Distance
- Noisy Channel Model

2

Normalization

 Goal: Convert text into a canonical (s 	standard)	form
--	-----------	------

- Remove noise / "randomness" from text
- Affects characters, <u>words</u>, sentences, documents
- Implicit definition of equivalence classes
 - Suitable normalization steps depend on task/application

Alternative to equivalence classes: **asymmetric expansion** Example: Web Search (utilize case of search terms)

Entered term		Searched terms
window	→	window, windows
windows	→	Windows, windows, window
Windows	→	Windows

		Raw	Normalized
)	~	Germany GERMANY	germany
	>	USA U.S.A US of A	USA
		tonight tonite 2N8	tonight
		connects connected connecting connection	connect
		:) :-) :0)	smile

Normalization — Case Folding

- When to fold?
 - Common application: Information Retrieval (e.g., Web search where must users type only in lowercase anyway)
 - Potential problems: Bush vs. bush, MOM vs. mom, Cloud vs. cloud, etc. (potential exception: upper case word in mid sentence?)

• When NOT to fold?

- NLP tasks where case of letters or words are important features
- Examples: Named Entity Recognition, Machine Translation

They sent **us** a card from the **US** during their vacation. Distinction important for NER and MT!

Outline

• Regular Expressions

- Basic Concepts
- Relationship to FSA
- Error Types

Corpus Preprocessing

- Tokenization
- Normalization
- Stemming / Lemmatization
- Segmentation

• Word error handling

- Spelling Errors
- Minimum Edit Distance
- Noisy Channel Model

2

Normalization — Stemming & Lemmatization

• Motivating example:

"dogs make the best friends" vs. "a dog makes a good friend"

→ Very similar semantics but (very) different syntax

Common reasons for variations of the same word

- Singular vs. plural form (mainly of nouns)
- Different tenses of verbs
- Comparative/superlative of adjectives

Can we normalize words to abstract from such variations?



- Idea of Stemming
 - Reduce words to their stem
 - Approach: crude chopping of affixes based on rules (→ language dependent)
 - Different stemmers apply different rules

• Characteristics

- Pro: fast + no lexicon required
- Con: stemmed word not necessarily
 a proper word (i.e., not in dictionary)

Examples

(alternatives reflect results from different stemmers)

Raw	Stemmed
cats	cat
running	run
phones	phon(e)
presumably	presum
crying	cry/cri
went	went
worse	wors
best	best
mice	mic(e)

Normalization — Stemming: Porter Stemmer

- Porter Stemmer most common stemmer for English text
 - Simple, efficient + very good results in practice
- Series of rewrite rules that run in a cascade
 - Output of each pass is fed is input to the next pass
 - Stemming steps if a pass yields no more changes



	sses → ss	e.g.: possesses \rightarrow possess, classes \rightarrow class
	tional → tion	e.g., optional \rightarrow option, fictional \rightarrow function
	ies → i	e.g., cries \rightarrow cri, tries \rightarrow tri
stem must contain vowel	(*v*)ng → ε	e.g.: sing \rightarrow sing, singing \rightarrow sing, talking \rightarrow talk
stem must contain >1 chars —	(m>1)ement → ε	e.g., replacement \rightarrow replac, cement \rightarrow cement

Normalization *Lemmatization*

- Idea of Lemmatization
 - Reduce inflections or variant forms to base form
 - Find the correct dictionary headword form
 - Differentiates between word forms: nouns (N), verbs (V), adjectives (A)

	Raw	Lemmatized (N)	Lemmatized (V)	Lemmatized (A)
X	running	running	run	running
	phones	phone	phone	phones
	went	went	go	went
	worse	worse	worse	bad
	mice	mouse	mice	mice

Normalization — Lemmatization: Characteristics

• Pros

- Lemmatized words are proper words (i.e., dictionary words)
- Can normalize irregular forms (e.g., went \rightarrow go, worst \rightarrow bad)

• Cons

- Requires curated lexicons / lookup tables + rules (typically)
- Requires Part-of-Speech tags for correct results
- Generally slower as stemming

Normalization — Stemming & Lemmatization

• Back to our motivating example

Raw:	"dogs make the best friends"	"a dog makes a good friend"
Stemmed:	"dog make the best friend"	"a dog make a good friend"
Lemmatized:	"dog make the good friend"	"a dog make a good friend"

Normalization — Final Words

- Canonical form also effects tokenization, e.g.: Penn Treebank Tokenizer
 - Separate out clitics (e.g., $doesn't \rightarrow does n't$; $John's \rightarrow John 's$)
 - Keep hyphenated words together
 - Separate out all punctuation symbols
- Other common normalization steps
 - Removal of stopwords (e.g., a, an, the, not, and, or, but, to, from, at)
 - Removal of non-standard tokens (e.g., URs, emojis, emoticons)

In-Lecture Activity





Which preprocessing step would arguably affect **sentiment analysis** negatively?

Outline

• Regular Expressions

- Basic Concepts
- Relationship to FSA
- Error Types

Corpus Preprocessing

- Tokenization
- Normalization
- Stemming / Lemmatization
- Segmentation

• Word error handling

- Spelling Errors
- Minimum Edit Distance
- Noisy Channel Model

2

Sentence Segmentation

- Sound like a simple task but...
 - Period "." can be quite ambiguous (e.g., "1.25", "U.S.A.", "Dr.") "?", "!" relatively unambiguous
 - Poor punctuation in informal text (common: missing whitespaces, missing capitalization)
 - → RegEx for segmenting sentences quickly become very complex
 Example RegEx: (?<!\w\.\w.) (?<! [A-Z] [a-Z] \.) (?<=\.|\?) \s</p>

• Alternative: binary classifier

(Source: Stackoverflow)

- Consider each period "." in a text
- Classify: EndOfSentence or NotEndOfSentence
- → Possible approaches: handwritten rules, set of RegEx, machine learning

Example: Simple Rules (represented as a binary Decision Tree)

Quick quiz: What are some common cases where this classifier would fail?



Many Other Features Conceivable

- Example: numerical features
 - length of word before / after period "."
 - Distance (in #chars) to next punctuation mark
 - Probabilities derived from a dataset
 (e.g., probability of with "." occurs at the end of sentence)

Side note: In informal text (e.g., social media) people often use emoticons or emojis to separate sentences, making this task even more complicated.

Break

CI	pboard T ₂	Font ra	Alignme
82	* : X v	f_x Febuary	
	A	В	С
1	JAN	January	
2	FEB	Febuary	
3	MAR	Maruary	
4	APR	Apruary	
5	MAY	Mayuary	
6	JUN	Junuary	
7	JUL	Juluary	
8	AUG	Auguary	

Meme Credits: The Language Nerds @ Facebook



60

Why are you taking this module?

It's a core requirement of my programme.	23 respondents	9 %	\checkmark
It's an elective course (one of a basket) requirement for my programme.	103 respondents	40 %	T .
I need this course to graduate this coming semester.	55 respondents	21 %	
I'm taking a related focus area, minor or track that has this as part of the possible fulfilling courses.	139 respondents	54 [%]	
I'm a student doing research, and NLP is my research area or related to my research area.	24 respondents	9 %	
It looked interesting.	214 respondents	83 %	
I have friends that are taking this course.	70 respondents	27 %	
It fits well in my timetable.	52 respondents	20 %	
Someone referred me to take this course because of its content.	29 respondents	11 %	
Someone referred me to take this course because of its instruction staff.	26 respondents	10 %	
I like project courses.	40 respondents	16 %	
I saw that it was popular.	30 respondents	12 %	
I have fulfilled all the prerequisites for the course, so why not?	38 respondents	15 %	
Others, I've mentioned in the next question.	5 respondents	2 %	





Major or Programme

Arts		0 %	\checkmark
Business	9 respondents	3 %	
Computer Engineering	20 respondents	8 %	
Computer Sciences	175 respondents	68 %	
Data Science	51 respondents	20 %	
English Language, Foreign Language		0 %	
Linguistics		0 %	
Math / Statistics	22 respondents	9 %	
Masters	42 respondents	16 %	
Ph.D. Candidate		0 %	
SCALE / Continuing Education Learner		0 %	
Other	6 respondents	2 %	

Overseas?

Yes, abroad on exchange or NOC.	1 respondent	0%
Yes, on internship.	25 respondents	10 %
Considering one of the two above.	13 respondents	5 %
No, full-time local student (most students).	219 respondents	85 %

NLP Career Path

No idea.	16 respondents	6 %
Unlikely.	9 respondents	3 %
Somewhat Unlikely.	29 respondents	11 %
So-So.	72 respondents	28 %
Somewhat Likely.	100 respondents	39 [%]
Likely.	31 respondents	12 %
No Answer	1 respondent	0 %

What do you want to learn?





Outline

• Regular Expressions

- Basic Concepts
- Relationship to FSA
- Error Types

• Corpus Preprocessing

- Tokenization
- Normalization
- Stemming / Lemmatization
- Segmentation

• Word error handling

- Spelling Errors
- Minimum Edit Distance
- Noisy Channel Model

2

Spelling Errors

1. Non-word error detections

- Basically, word is not found in dictionary
 - Example: detecting graffe (misspelling of giraffe)

2. Isolated-word error correction

- Consider word in isolation (i.e., without surrounding words)
- Example: correcting *graffe* to *giraffe*

3. **Context-sensitive error detection & correction**

- Consider surrounding words to detect and correct errors
- Important for "wrong" words that a spelled correctly
- Examples: there vs. three, dessert vs. desert, son vs. song



avr

Spelling Errors — Common Patterns

- Observation
 - Most misspelled words in typewritten text are single-error
 - Damerau (1964): 80%, Peterson (1986): 93-95%
- Single-error misspellings
 - Insertion (e.g., *acress* vs. *acres*)
 - Deletion (e.g., *acress* vs. *actress*)
 - Substitution (e.g., *ac<u>r</u>ess* vs. *ac<u>c</u>ess)*
 - Transposition (e.g., <u>acress</u> vs. <u>caress</u>)

For non-word errors:

- → Good candidates are orthographically similar
- → Minimum Edit Distance

Outline

• Regular Expressions

- Basic Concepts
- Relationship to FSA
- Error Types

Corpus Preprocessing

- Tokenization
- Normalization
- Stemming / Lemmatization
- Segmentation

• Word error handling

- Spelling Errors
- Minimum Edit Distance
- Noisy Channel Model

2

Minimum Edit Distance (MED)

- Minimum Edit Distance between 2 strings s_1 and s_2
 - Minimum number of allowed edit operations to transform s_1 into s_2
 - Allowed edit operations: Insertion, Deletion, Substitution, Transposition -

Not covered here to keep examples simple

- Example
 - $s_1 = "LANGUAGE"$

→ Alignment of MED:

- $s_2 = "SAUSAGE"$
- L A N G U * A G E | | | | | | | | | S A * * U S A G E

- MED if all operations cost 1 \rightarrow 4
- MED if Substitution costs 2, Insertion 1, Deletion 1 → 5
- Problem formulation: Find a path (i.e., sequence of edits) from start string to final string
 - Initial state: the word being transformed (e.g., "LANGUAGE")
 - Target state: the word being transformed into (e.g., "SAUSAGE")
 - Operators: insert, delete, substitute
 - Path cost: aggregated costs of all edits



- → Potentially huge search space
- → Naive navigation of all path impractical

- Observations
 - Many distinct paths end up in the same state



- → No need to keep track of all paths
- → Only important: "cheapest" path to each revisited state (best in terms of costs, not just number of operations!)
- → Solve using Dynamic Programming solving problems by combining solutions to subproblems

- Input: 2 strings
 - $\bullet \quad \text{Source string } X \text{ of length } n \\$



- Bottom-up approach of Dynamic Programming
 - Compute D(i, j) for small i, j (base cases)
 - Compute D(i,j) for larger *i*, *j* based on previously computes D(i,j) for smaller *i*, *j* \checkmark

Initialization of bases cases



Assumptions for costs
Insert: 1
Delete: 1
Substitute: 2
→ Levenshtein MED

• For $0 < i \le n$ and $0 < j \le m$

$$D(i,j) = \min \begin{cases} D(i-1,j) + 1 & \text{Delete} \\ D(i,j-1) + 1 & \text{Insert} \\ D(i-1,j-1) + \begin{cases} 2, & if X[i] \neq Y[j] \\ 0, & if X[i] = Y[j] \end{cases} \text{Substitute} \end{cases}$$

Complexity	analysis
Space:	O(nm)
Time:	O(nm)



Minimum Edit Distance — Calculation Example

Ε	8	9	8	7	8	7	6	5
G	7	8	7	6	7	6	5	6
Α	6	7	6	5	6	5	6	7
U	5	6	5	4	5	6	7	8
G	4	5	4	5	6	7	6	7
Ν	3	4 -	7,3	4	5	6	7	8
Α	2	3	2	3	4	5	6	.7
\mathbf{L}	1	2	3	4	5	6	7	8
#	0	1	2	3	4	5	6	7
	#	\mathbf{S}	Α	U	S	Α	G	\mathbf{E}

Minimum Edit Distance — Backtrace & Alignments

- Current limitation
 - Base algorithm only returns the MED
 - Often important: alignment between strings



How do we get this?

• Keep track of backtrace

- Remember from which "direction" we entered a new cell
- At the end, trace path from upper right corner to read of alignment

Keep set of pointers for each i, j

Small extension to base algorithm:

$$PTR(i,j) = \begin{cases} \texttt{LEFT} & \texttt{Insert} \\ \texttt{DOWN} & \texttt{Delete} \\ \texttt{DIAG} & \texttt{Substitute} \end{cases}$$

Note: Backtraces are generally not unique → different alignments for the same MED possible

Minimum Edit Distance — Backtrace & Alignments

SA*

*

USAGE



Quick quiz: Why do we choose the diagonal path here?

Minimum Edit Distance — More Examples

• Biology: Align 2 sequences of nucleotides

AGGCTATCACCTGACCTCCAGGCCGATGCCC TAGCTATCACGACCGCGGTCGATTTGCCCGAC

C	31	$\downarrow 30$	1.29	↓ 28	1. 27	↓ 26	1 25	1.24	/ 1 23	$\checkmark \leftarrow \downarrow 24$	/ 1 23	1 22	↓ 21	1 20	19	↓ 18	117	$\downarrow 16$	↓ 15	$\checkmark \leftarrow \downarrow 16$	√↓15	√ ←↓ 16	15	↓ 14	./←↓15	√ ←., 16	15	14	1 18	/ 12	$\leftarrow 13$	$\leftarrow 14$	$\checkmark \leftarrow 15$
C	30	1 29	↓ 28	1 27	1 26	1 25	↓ 24	↓ 23	1 22	$\checkmark \leftarrow \downarrow 23$	1 22	↓ 21	1 20	/119	√↓18	1 17 L	∠ ↓ 16	↓ 15	14	√←↓ 15	114		↓ 14	1 13	./←↓14	√ ← ↓ 15	14	1 13	/ 12	$\swarrow \leftarrow 13$	$\leftarrow 14$	← 15	$\swarrow \leftarrow 16$
C	29	$\downarrow 28$	↓ 27	1.26	2. 25	↓ 24	1 23	↓ 22	21	$\langle \downarrow 22$	/1 21	4 20	19	118	2.17	↓ 16	√↓ 15	↓ 14	↓ 13	1 14	/ 13	< ↓14	13	↓ 12	./ + 13	1 . 14	13	/ 12	1 13	14 14	< 15	< 16	< 17
G	28	1 27	↓ 26	√↓ 25	↓ 24	1 23	↓ 22	↓ 21	$\checkmark \leftarrow \downarrow 22$	↓ 21	↓ 20	/ 19	18	117	$\downarrow 16$	/1 15	14	1.13	/ 12	$\leftarrow \downarrow 13$	√←↓14	113	$\downarrow 12$	↓ 11	./←↓ 12	√←↓13	/ 12	$\leftarrow 13$	$\leftarrow 14$	$\leftarrow 15$	$\checkmark \leftarrow 16$	$\leftarrow \downarrow 17$./←↓18
т	27 ,	/1.26	4.25	↓ 24	$\downarrow 23$	/ 1 22	↓ 21	20	<↓21	1, 20	↓ 19	↓ 18	↓ 17	↓ 16	$\downarrow 15$	↓ 14	$\downarrow 13$	↓ 12	$\checkmark \leftarrow \downarrow 13$	/ 12	$\leftarrow \downarrow 13$	$\downarrow 12$	↓ 11	10	14-11	× < 12	← 13	<−14	< 15	← 16	<↓17	$\downarrow 16$	1 + 117
A	26	1 25	1.24	1 23	↓ 22	↓ 21	√↓ 20	√ ← 21	1 20	√↓19	↓ 18	↓ 17	1 16	1 15	$\downarrow 14$	↓ 13	. 12	↓ 11	$\checkmark \leftarrow \downarrow 12$	$\checkmark \leftarrow \downarrow 13$	$\downarrow 12$	11	10	← 11	← 12	← 13	$\leftarrow 14$	← 15	$\leftarrow \downarrow 16$./←↓17	↓ 16	/ 15	$\leftarrow 16$
G	25	$\downarrow 24$	1.23	1 22	1.21	± 20	↓ 19	1 4. 20	↓ 19	↓ 18	↓ 17	14 16	15	114	13	12		10	14411	1 4 ↓ 12	↓ 11	10	← 11	$\leftarrow 12$	← 13	← 14		\leftarrow , 16	↓ 15	16 €	/ 15	$\leftarrow 16$	← 17
C	24	4 23	↓ 22	↓ 21	1.20	↓ 19	↓ 18	1 . 19	118	$\downarrow 17$	/116	4 15	14	113	√↓12	↓ 11	√↓10	14.11	↓ 10	1+11	√↓ 10	1 ←↓ 11	1+ 12	$\checkmark \leftarrow \downarrow 13$	1 4 14	$\checkmark \leftarrow \downarrow 15$	√←↓16	√↓ 15	1				18
C	23	$\downarrow 22$	1.21	↓ 20	1. 19	↓ 18	↓ 17	14. 18	117	$\downarrow 16$	/115	↓ 14	113	112	11	↓ 10	29	←, 10	1.0	10	19	$\leftarrow 10$	← 11	← 12	← 13	← 14	←↓ 15	14	14				17
G	22	↓ 21	1 20	√↓ 19	↓ 18	17	↓ 16	1 17	116	4 15	114	/ 13	112	↓11	↓ 10	V.1.9	1 4 10	1.9	18	< 9	< 10	1 4 11	1 . 12	1 13	1 4 14	1 . 15	114	1 + 15	1141		- N	- 16	18
G	21	↓ 20	↓ 19	118	$\downarrow 17$	↓ 16	↓ 15	√←↓ 16	↓ 15	$\downarrow 14$	$\downarrow 13$	112	111	1 10	$\downarrow 9$	18	$e \downarrow \rightarrow$	18	$\sqrt{\leftarrow 9}$	$\leftarrow 10$	+↓ 11	\$ 10	÷↓∏	√←↓12	√+↓13	14 → √	/ 13	← 14	÷ .		- U -	- Ii	19
Α	20	↓ 19	1.18	↓ 17	↓ 16	↓ 15	/ 14	←15	↓ 14	13	$\downarrow 12$	↓11	10	19	184		18		$\checkmark \leftarrow \downarrow 10$	14411	↓ 10	14411	10	← 11	← 12	← 13	← 14	←.↓ 15	144				18
C	19	1 18	+ 17	16	1.15	1 14	1 4 + 1 15	↓ 14	/1 13	↓ 12	/111	1 10 L	19	118	217	./+↓8	17			$\checkmark \leftarrow \downarrow 10$	2.9	√+110	√←↓11	√ ←↓ 12	./←↓13	$\checkmark \leftarrow \downarrow 14$	√+115	14					17
C	18	↓ 17	4.16	↓ 15	2.14	↓ 13	14 €	↓ 13	112	11	√↓10	U 9	18	147	16	1417	16	← 7	← 8	$e \downarrow \rightarrow$	18	$e \rightarrow$	← 10	← 11	← 12	← 13	← 14	$\checkmark \leftarrow 15$	$\swarrow \leftarrow \downarrow 16$	/ 15	← 16	$\leftarrow 17$	$\swarrow \leftarrow 18$
т	17 .	/16	115	14	↓ 13	12	$\leftarrow \downarrow 13$	1. 12	111	↓ 10	19	18	17	16	15	./+↓6	.∠+↓7	/←18		/8	$\leftarrow 9$	$\leftarrow 10$	$\leftarrow 11$	$\checkmark \leftarrow 12$	$\checkmark \leftarrow 13$	$\checkmark \leftarrow 14$	← 15	$\leftarrow \downarrow 16$	↓ 15	√+116	√←↓17		./ ←↓ 19
C	16	+ 15	↓ 14	↓ 13	1. 12	$\checkmark \leftarrow \downarrow 13$	↓ 12	11	V 10	<u>+ 9</u>	18	1,7	16	15	14	← 5	$\checkmark \leftarrow 6$	$\leftarrow 7$	← 8	$e \rightarrow$	√ ← 10	← 11	$\leftarrow \downarrow 12$	√←↓13	√←↓14	√←↓15	√←↓16	15	14	$\checkmark \leftarrow 15$	← 16	$\leftarrow 17$	$\checkmark \leftarrow 18$
C	15	1 14	↓ 13	↓ 12	11	$\leftarrow \downarrow 12$	↓ 11	$\downarrow 10$	19	↓8	17	16	15	14	$\swarrow \leftarrow 5$	$\leftarrow 6$	2 ← 7	$\leftarrow 8$	$\leftarrow 9$	$\leftarrow 10$		$\leftarrow \downarrow 12$	$\downarrow 11$	∠ ←↓ 12	./ ←↓ 13			14	√ ← 15	$\swarrow \leftarrow 16$	$\leftarrow 17$	$\leftarrow 18$	$\swarrow \leftarrow 19$
A	14	$\downarrow 13$	√↓ 12	111	$\checkmark \leftarrow \downarrow 12$	111	∠↓ 10	$\downarrow 9$	18	217	16	5	14	← 5	$\leftarrow 6$	$\leftarrow \downarrow 7$			$\checkmark \leftarrow \downarrow 10$		$\checkmark \leftarrow \downarrow 12$	\downarrow 11	10	$\leftarrow 11 \rightarrow$	$\leftarrow 12$	$\leftarrow 13$	$\leftarrow 14$	$\leftarrow 15$	$\leftarrow 16$	$\leftarrow 17$	$\leftarrow 18$	$\swarrow \leftarrow 19$	← 20
G	13	$\downarrow 12$	↓ 11	/ 10	←↓ 11	1 10	19	18	17	$\downarrow 6$	15	14	$\leftarrow 5$	$\leftarrow \downarrow 6$	1+17	16	$\leftarrow 7$	14-8		$\leftarrow \downarrow 10$	$\checkmark \leftarrow \downarrow 11$	10	← 11	$\leftarrow 12$	← 13	$\leftarrow 14$	$\checkmark \leftarrow 15$	$\leftarrow 16$	$\leftarrow 17$	← 18	$\checkmark \leftarrow 19$	$\leftarrow 20$	$\leftarrow 21$
Т	12 ,	/11	1 10 L	.∠+↓11	↓ 10	19	18	1.7	16	↓5	$\downarrow 4$./+↓5		15		./+↓7	1418		√ ←↓ 10	/ 9	$\leftarrow 10$	$\leftarrow 11$	$\leftarrow 12$	$\checkmark \leftarrow 13$	$\checkmark \leftarrow 14$	√ ← 15	$\leftarrow 16$	$\leftarrow 17$	$\leftarrow 18$	$\leftarrow 19$	← 20	$\leftarrow 21$	$\leftarrow 22$
C	11	↓ 10	19	√<↓10	/19	4.8	17	- 6	×45	1.4	/43	1.44	X + 5	14	1.5	< 6	1:7	< 8	< 9	< 10	/ < 11	< 12	< 13	< 14	< 15	< 16	< 17	√ < 18	1 19	∠< 20	< 21	< 22	$\checkmark \downarrow 23$
С	10	$\downarrow 9$	$\downarrow 8$	∠ ←↓ 9	218	$\downarrow 7$	↓ 6	$\downarrow 5$	244	$\downarrow 3$	12	$\leftarrow 3$	$\leftarrow 4$	$\swarrow \leftarrow 5$	$\swarrow \leftarrow 6$	← 7	2 ← 8	$e \rightarrow$	$\leftarrow 10$	← 11	$\swarrow \leftarrow 12$	$\leftarrow 13$	$\leftarrow 14$	$\leftarrow 15$	$\leftarrow 16$	$\leftarrow 17$	$\leftarrow 18$	$\swarrow \leftarrow 19$	$\swarrow \leftarrow 20$	$\swarrow \leftarrow 21$	$\leftarrow 22$	$\leftarrow \downarrow 23$	2 22
Α	9	4.8	1.47	√←↓8		.↓6	145	4	↓ 3	/ 2	$\leftarrow 3$	← 4	14-5	<- 6	$\leftarrow 7$	< <u>−</u> 8	$\leftarrow 9$	← 10	$\leftarrow 11$	← 12	← 13	← 14	$\sqrt{\leftarrow 15}$	← 16	← 17	← 18	$\leftarrow 19$	$\leftarrow \downarrow 20$	$\checkmark \leftarrow \downarrow 21$	1 - 1 22	$/ \leftarrow 23$	/ 22	$\leftarrow 23$
C	8	↓7	↓ 6	∠+↓7	216	↓5	14	↓ 3	12	$\leftarrow 3$	$\swarrow 4$	$\leftarrow 5$	$\leftarrow 6$	√ ← 7	$\swarrow \leftarrow 8$	$\leftarrow 9$	$\swarrow \leftarrow 10$	← 11	$\leftarrow 12$	$\leftarrow 13$	$\swarrow \leftarrow 14$	$\leftarrow 15$	$\leftarrow 16$	$\leftarrow \downarrow 17$		√←↓19	∠←↓ 20	2 19	$\swarrow \leftarrow 20$	$\swarrow \leftarrow 21$	← 22	$\leftarrow 23$	$\swarrow \leftarrow 24$
Т	7	1.6	ψð	√ ←↓ 6	45	14	13	12	$\leftarrow 3$	$\leftarrow 4$	$\leftarrow 5$	$\leftarrow 6$	$\leftarrow 7$	$\leftarrow 8$	$\leftarrow 9$	$\leftarrow 10$	← 11	$\leftarrow 12$	$\leftarrow 13$	$\checkmark \leftarrow 14$	← 15	$\leftarrow 16$	$\leftarrow \downarrow 17$	/ 16	$\sqrt{\leftarrow 17}$	√ ← 18	$\leftarrow 19$	$\leftarrow 20$	$\leftarrow 21$	$\leftarrow 22$	$\leftarrow 23$	$\leftarrow 24$	$\leftarrow 25$
Λ	6	15	14	←↓ 5	J 4	13	12	$\leftarrow 3$	$\leftarrow 4$	$\checkmark \leftarrow 5$	$\leftarrow 6$	$\leftarrow 7$	$\swarrow \leftarrow 8$	$e \rightarrow$	$\leftarrow 10$	$\leftarrow 11$	$\leftarrow 12$	$\leftarrow 13$	$\leftarrow 14$	$\leftrightarrow \downarrow 15$	√ ← ↓ 16	1 + 17	16	$\leftarrow 17$	$\leftarrow 18$	$\leftarrow 19$	$\leftarrow 20$	$\leftarrow 21$	$\leftarrow 22$	$\leftarrow 23$	$\leftarrow 24$	$\swarrow \leftarrow 25$	$\leftarrow 26$
т	5	11	←, 5	$\downarrow 4$., 3	1 2	$\leftarrow 3$	$\checkmark \leftarrow 4$	← 5	$\leftarrow 6$	$\leftarrow 7$	$\leftarrow 8$	$e \rightarrow$	$\leftarrow \downarrow 10$	11 + + 11	$\checkmark \leftarrow \downarrow 12$	√←↓13	14 - 14	√ ←↓ 15	14	← 15	$\leftarrow 16$	← 17	$\swarrow \leftarrow 18$	$\checkmark \leftarrow 19$	$\swarrow \leftarrow 20$	$\leftarrow 21$	$\leftarrow 22$	$\leftarrow 23$	$\leftarrow 24$	$\leftarrow 25$	$\leftarrow 26$	$\leftarrow 27$
C	4 1	+↓5	14	↓ 3	12	$\leftarrow 3$	$\leftarrow 4$	$\leftarrow 5$	$\checkmark \leftarrow 6$	$\leftarrow 7$./ ← 8	÷, 9	.∕ ←↓ 10	19	$\swarrow \leftarrow 10$	$\leftarrow 11$	$\checkmark \leftarrow 12$	$\leftarrow 13$	$\leftarrow 14$	$\leftarrow 15$	$\checkmark \leftarrow 16$	$\leftarrow 17$	$\leftarrow 18$	$\leftarrow 19$	$\leftarrow 20$	$\leftarrow 21$	$\leftarrow 22$	$\checkmark \leftarrow 23$	$\swarrow \leftarrow 24$	$\checkmark \leftarrow 25$	← 26	$\leftarrow 27$	$\checkmark \leftarrow 28$
G	3 2	$\leftarrow \downarrow 4$	13	12	√←13		2 ←1 5	2 ← 1 6	2+↓7	2 ← 1 8	∠+↓9	18	$e \rightarrow$	$\leftarrow 10$	← 11	√ ← 12	$\leftarrow 13$	$\swarrow \leftarrow 14$	$\swarrow \leftarrow 15$	$\leftarrow 16$	$\leftarrow 17$	$\checkmark \leftarrow 18$	← 19	$\leftarrow 20$	$\leftarrow 21$	← 22	$\swarrow \leftarrow 23$	$\leftarrow 24$	$\leftarrow 25$	$\leftarrow 26$	$\swarrow \leftarrow 27$	$\leftarrow 28$	$\leftarrow 29$
G	2 ./	(←↓ 3	↓ 2	1	$\leftarrow 2$	$\leftarrow 3$	$\leftarrow 4$	$\leftarrow 5$	$\leftarrow 6$	$\leftarrow 7$	$\leftarrow 8$./ ← 9	$\leftarrow 10$	$\leftarrow 11$	$\leftarrow 12$./ ← 13	$\leftarrow 14$	$\swarrow \leftarrow 15$	$\checkmark \leftarrow 16$	← 17	$\leftarrow 18$	$\swarrow \leftarrow 19$	$\leftarrow 20$	$\leftarrow 21$	← 22	$\leftarrow 23$	$\swarrow \leftarrow 24$	$\leftarrow 25$	$\leftarrow 26$	$\leftarrow 27$	$\checkmark \leftarrow 28$	$\leftarrow 29$	$\leftarrow 30$
Α	1 2	$\leftarrow \downarrow 2$	1	$\leftarrow 2$	$\leftarrow 3$	$\leftarrow 4$	$\swarrow \leftarrow 5$	$\leftarrow 6$	← 7	2 ← 8	$\leftarrow 9$	$\leftarrow 10$		$\leftarrow 12$	$\leftarrow 13$	$\leftarrow 14$	$\leftarrow 15$	$\leftarrow 16$	$\leftarrow 17$	$\leftarrow 18$	$\leftarrow 19$	$\leftarrow 20$	$\swarrow \leftarrow 21$	$\leftarrow 22$	$\leftarrow 23$	$\leftarrow 24$	$\leftarrow 25$	$\leftarrow 26$	$\leftarrow 27$	$\leftarrow 28$	$\leftarrow 29$	$\swarrow \leftarrow 30$	$\leftarrow 31$
#	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32
	11	m		0	0	m		00	0		0	0		0	0	0	0	0	0	ar i	0	0		m	ar .	m	0	0	C1	0	0		0



🏃 🏃 🏃 In-Lecture Activity (5 mins)

- Task: Compute the MED and alignment between "NUS" and "TRUST"
 - Post your MED (Levenshtein) and alignment to Canvas > Discussions (individually or as a group – add all group members' names to the post)



- Try to complete the table for this task (probably not needed as the words are very short)
- Some of you can share their solution

Example alignment (but bad one!)
NUS****
***TRUST

In-Lecture Activity



• Solution

S	3	$\swarrow \leftarrow \downarrow 4$	$\swarrow \leftarrow \downarrow 5$	$\downarrow 4$	3	$\leftarrow 4$
\mathbf{U}	2	$\swarrow \leftarrow \downarrow 3$	$\swarrow \leftarrow \downarrow 4$	23	$\leftarrow 4$	$\leftarrow 5$
Ν	1	$\swarrow \leftarrow \downarrow 2$	$\checkmark \leftarrow \checkmark 3$	$\swarrow \leftarrow \downarrow 4$	$\swarrow \leftarrow \downarrow 5$	$\swarrow \leftarrow \downarrow 6$
#	0	1	2	3	4	5
	#	Т	\mathbf{R}	\mathbf{U}	S	Т

* NUS* ____ | S) T

Minimum Edit Distance — Other Uses in NLP

• Evaluating Machine Translation and speech recognition

e.g., How similar are 2 translations?

Reference:	Spokesman	confirms	*	senior	government	adviser	was	shot	*
		- I			1				
Prediction:	Spokesman	said	the	senior	*	adviser	was	shot	dead

Named Entity Extraction and Entity Coreference

"We stayed at the * Merchant Court prior to a cruise" "The Swissotel Merchant Court is a great place to stay in Singapore" Referring to the same entity?



Minimum Edit Distance — Efficiency

• Time: *O(nm)*

• Space: *O(nm)*

• Backtrace: O(n+m)

Minimum Edit Distance — Extensions

- Weighted Minimum Edit Distance, e.g.:
 - Spell Correction: some letters are more likely to be mistyped than others
 - Biology: certain kinds of deletions or insertions are more likely than others

→ Generalization of algorithm

Application-dependent weights (i.e., costs for edit operations)

Initialization of base cases:Recurrence relation:
$$D(0,0) = 0$$
 $D(i,0) = D(i-1,0) + del(X[i]), \text{ for } 1 < i \le n$ $D(i,j) = min \begin{cases} D(i-1,j) \\ D(i,j-1) \\ D(i-1,j-1) \end{cases} + ins(Y[j]) + ins(Y[j]) \\ sub(X[i],Y[i]) \end{cases}$

Minimum Edit Distance — Extensions

• Needleman-Wunsch

- No penalty for gaps (*) at the beginning or the end of an alignment
- Good if strings have very different lengths

• Smith-Wasserman

- Ignore badly aligned regions
- Find optimal <u>local</u> alignments within substrings (Levenshtein finds the best global distance and alignment)

Common application: Alignment of nucleotides sequences

Outline

• Regular Expressions

- Basic Concepts
- Relationship to FSA
- Error Types

• Corpus Preprocessing

- Tokenization
- Normalization
- Stemming / Lemmatization
- Segmentation

• Word error handling

- Spelling Errors
- Minimum Edit Distance
- Noisy Channel Model

2

Where We are Right Now

- Given a misspelled word, generate suitable candidates for error correction
 - 80% of errors are within minimum edit distance 1
 - Almost all errors within minimum edit distance 2
 - Covers also missing spaces and hyphens (e.g., thisidea vs. this idea; inlaw vs. in-law)
- Still missing: Which is the most likely candidate?
 - Ranking of candidates to show top candidates first
 - Support for automated spelling correction

→ Noisy Channel Model

Idea: Assign each candidate a probability





Decoding: Observing error x, can we predict correct word w?

Noisy Channel Model — Bayesian Inferencing



Noisy Channel Model — Calculating/Estimating P(w)

- Approach using Maximum Likelihood Estimate (MLE)
 - $\blacksquare \;$ Required: Large text corpus with N words
 - Calculate/estimate P(w) with $P(w) = \frac{freq(w)}{N}$
- Example
 - 100 MB Wikipedia dump
 - Total of 14.4M+ words

			\mathcal{T}
	W	freq(w)	P(w)
->	actress	1,135	0.0000784
7	cress	1	0.00000
_	caress	3	0.00000
~	access	1,670	0.0001153
	across	1,756	0.0001213
	acres	177	0.0000122

Note: The frequencies can widely different across different corpora (e.g. Wikipedia articles vs. English Literature).

Noisy Channel Model — Calculating/Estimating P(x|w)

- In general, P(x|w) almost impossible to predict
 - Predictions depends on arbitrary factors

(e.g., proficiency of typist, lighting conditions, input device)

- Estimate P(x|w) based on simplifying assumptions (Kernighan et al., 1990)
 - Most misspelled words in typewritten text are single-error
 - Consider only single-error misspellings: Insertion, Deletion, Substitution, Transposition

Noisy Channel Model — Calculating/Estimating $P(\boldsymbol{x}|\boldsymbol{w})$

- Definition of 4 confusion matrices (1 for each single-error type)
 - Each confusion matrix lists the number of times one "thing" was confused with another
 - e.g., for substitution, an entry represents the number of times one letter was incorrectly used
- Underlying definitions for generate confusion matrices

	ins[x, y]	number of times x was typed as xy
	del[x,y]	number of times xy was typed as x
	sub[x, y]	number of times x is substituted for y
	trans[x, y]	number of times xy was typed as yx
7	count[x]	number of times that \boldsymbol{x} appeared in the training set
	count[x, y]	number of times that xy appeared in the training set



Noisy Channel Model — Calculating/Estimating P(x|w)

$$P(x|w) = \begin{cases} \frac{ins[w_{i-1}, x_i]}{count[w_i]} &, \text{ if insertion} \\ \frac{del[w_{i-1}, w_i]}{count[w_{i-1}, w_i]} &, \text{ if deletion} \\ \frac{sub[x_i, w_i]}{count[w_i]} &, \text{ if substitution} \\ \frac{trans[w_i, w_{i+1}]}{count[w_i, w_{i+1}]} &, \text{ if transposition} \end{cases}$$

 w_i = i-th character in the correct word w

 x_i = i-th character in the misspelled word x

Noisy Channel Model — Calculating/Estimating P(x|w)

X						ubl	,	-1-	Jul	Sere	a ci ci	Y	(col	rect	, iii	<i>cc)</i> 1		(0.	,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,							
	a	b	с	d	e	f	g	h	i	j	k	1	m	n	0	р	q	r	S	t	u	v	w	X	у	Z
a	0	0	7	1	342	0	0	2	118	0	1	0	0	3	76	0	0	1	35	9	9	0	1	0	5	0
b	0	0	9	9	2	2	3	1	0	0	0	5	11	5	0	10	0	0	2	1	0	0	8	0	0	0
c	6	5	0	16	0	9	5	0	0	0	1	0	7	9	1	10	2	5	39	40	1	3	7	1	1	0
d	1	10	13	0	12	0	5	5	0	0	2	3	7	3	0	1	0	43	30	22	0	0	4	0	2	0
e	388	0	3	11	0	2	2	0	89	0	0	3	0	5	93	0	0	14	12	6	15	0	1	0	18	0
f	0	15	0	3	1	0	5	2	0	0	0	3	4	1	0	0	0	6	4	12	0	0	2	0	0	0
g	4	1	11	11	9	2	0	0	0	1	1	3	0	0	2	1	3	5	13	21	0	0	1	0	3	0
h	1	8	0	3	0	0	0	0	0	0	2	0	12	14	2	3	0	3	1	11	0	0	2	0	0	0
i	103	0	0	0	146	0	1	0	0	0	0	6	0	0	49	0	0	0	2	1	47	0	2	1	15	0
i	0	1	1	9	0	0	1	0	0	0	0	2	1	0	0	0	0	0	5	0	0	0	0	0	0	0
k	1	2	8	4	1	1	2	5	0	0	0	0	5	0	2	0	0	0	6	0	0	0	4	0	0	3
1	2	10	1	4	0	4	5	6	13	0	1	0	0	14	2	5	0	11	10	2	0	0	0	0	0	0
m	1	3	7	8	0	2	0	6	0	0	4	4	0	180	0	6	0	0	9	15	13	3	2	2	3	0
n	2	7	6	5	3	0	1	19	1	0	4	35	78	0	0	7	0	28	5	7	0	0	1	2	0	2
0	91	1	1	3	116	0	0	0	25	0	2	0	0	0	0	14	0	2	4	14	39	0	0	0	18	0
p	0	11	1	2	0	6	5	0	2	9	0	2	7	6	15	0	0	1	3	6	0	4	1	0	0	0
q	0	0	1	0	0	0	27	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
r	0	14	0	30	12	2	2	8	2	0	5	8	4	20	1	14	0	0	12	22	4	0	0	1	0	0
s	11	8	27	33	35	4	0	1	0	1	0	27	0	6	1	7	0	14	0	15	0	0	5	3	20	1
t	3	4	9	42	7	5	19	5	0	1	0	14	9	5	5	6	0	11	37	0	0	2	19	0	7	6
u	20	0	0	0	44	0	0	0	64	0	0	0	0	2	43	0	0	4	0	0	0	0	2	0	8	0
v	0	0	7	0	0	3	0	0	0	0	0	1	0	0	1	0	0	0	8	3	0	0	0	0	0	0
w	2	2	1	0	1	0	0	2	0	0	1	0	0	0	0	7	0	6	3	3	1	0	0	0	0	0
x	0	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	9	0	0	0	0	0	0	0
y	0	0	2	0	15	0	1	7	15	0	0	0	2	0	6	1	0	7	36	8	5	0	0	1	0	0
Z	0	0	0	7	0	0	0	0	0	0	0	7	5	0	0	0	0	2	21	3	0	0	0	0	3	0

sub[X, V] = Substitution of X (incorrect) for V (correct)

Source: <u>A Spelling Correction Program Based on a Noisy Channel Model</u> (Kernighan et al., 1990)

Noisy Channel Model — Example

• Noisy channel probabilities for "acress"

Candidate Correction	Correct Letter	Error Letter	x w	P(x w)	P(w)	10 ⁹ *P(x w)P(w)	%	
actress	t		c ct	.000117	0117 .0000231 2.7			
cress		а	a #	.00000144	.00000054	.00078	~0	
caress	са	ac	ac ca	.00000164	.00000170	.0028	~0	
access	С	r	r c	.00000021	.0000916	.019	~0	
across	ο	е	eļo	.0000093	.000299	2.8	37.2	
acres		S	es e	.0000321	.0000318	1.0	13.3	
acres		S	ss s	.0000342	.0000318	1.0	13.3	

→ Choice of candidate for correction: *across*

Noisy Channel Model — Discussion

- Basic limitation: No consideration of additional context
 - Model only applicable for non-word errors
 - Basic model will always suggest "across" to correct "acress"

"The role was played by an acress famous for her comedic timing."

→ Language Models (next lecture)

Summary

- RegEx fundamental and useful tool
- Text Preprocessing getting your data ready for analysis
 - Tokenization
 - Stemming / Lemmatization
 - Normalization

typical very task-dependent!

- Error Handling (so far)
 - Focus on single-error misspellings
 - Focus on isolated-word error correction

already very non-trivial!

Student Learning Outcomes

Outlook for Next Week: Language Models

Image from Da Nina @ Unsplash

Pre-Lecture Activity for Next Week

- Assigned Task
 - Post a 1–2 sentence answer to the following question into the L1 Discussion (you will find the thread on Canvas > Discussions)

"What do we mean when we talk about the probability of a sentence?"

Side notes:

- This task is meant as a warm-up to provide some context for the next lecture
- No worries if you get lost; we will talk about this in the next lecture
- You can just copy-&-paste others' answers, but his won't help you learn better