



**NUS**  
National University  
of Singapore

| **Computing**

# **CS4248: Natural Language Processing**

## **Lecture 2 — Strings & Words**

# 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

# 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 first occurrence of pattern
- Global search: match all occurrences of pattern (assumed in most following examples)

**Example: password validation**

- \* Must have a minimum of 8 characters
- \* Must not contain username
- \* Must include at least 1 uppercase
- \* Must include at least 1 lowercase
- \* Must include at least 1 digit or 1 special character:  
~ ! @ # \$ % ^ & \* \_ - + = ` | \ ( ) { } [ ] ; : ' < > , . ? /

# Basic Patterns

- Fixed patterns

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

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

`blocks` → *My block has 15 floors, and I live on floor 5.*

- Special characters (metacharacters)

Character	Explanation
.	matches any character except line breaks
^	match the start of a string
\$	match the end of a string
	matches RegEx either before or after the symbol (e.g., <code>floor floors</code> )
\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)

`[.,;:]` → *My block has 15 floors, and I live on floor 5.*

(match all sequences of length 1 that are either a period, comma, etc.)

`[^a-z]` → *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
<b>\d</b>	[0-9]	matches any digit
<b>\D</b>	[^0-9]	matches any non-digit
<b>\s</b>	[ \n\r\t\f]	matches any whitespace character
<b>\S</b>	[^ \n\r\t\f]	matches any non-whitespace character
<b>\w</b>	[a-zA-Z0-9_]	matches any word character
<b>\W</b>	[^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
+	1 or more occurrences
*	0 or more occurrences
?	0 or 1 occurrences
{n}	exactly n occurrences
{l,u}	between l and u occurrences; can be unbounded: {l,} 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)

`\b\w{2,4}\b` → *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)



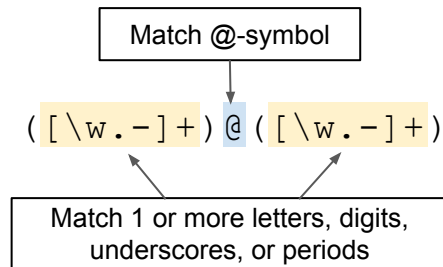
# Groups

- 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, one 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.



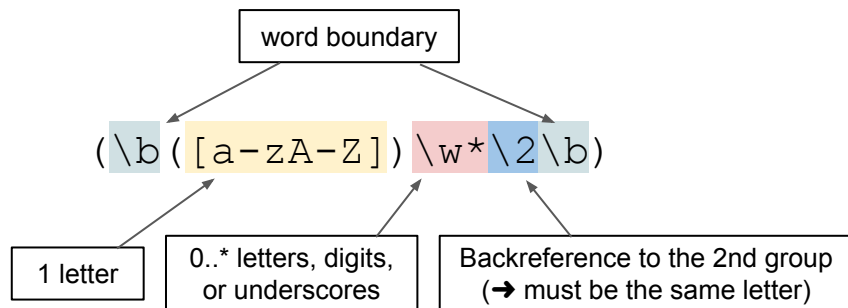
<b>Match:</b>	<code>user@example.org</code>
<b>Group #1:</b>	<code>alice</code>
<b>Group #2:</b>	<code>example.org</code>

# Backreferences

- 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:	<i>mom</i>
Group #1:	<i>mom</i>
Group #2:	<i>m</i>

Match:	<i>test</i>
Group #1:	<i>test</i>
Group #2:	<i>t</i>



# 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

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

# 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*)

`\d+ ( ?=\s*kg)`



*"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."*

`[0-9., ]* [0-9]+ ( ?=\s*kg)`



*"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."*

# Outline

- **Regular Expressions**

- Basic Concepts
- **Relationship to FSA**
- Error Types

- **Corpus Preprocessing**

- Tokenization
- Normalization
- Stemming / Lemmatization
- Segmentation

- **Word error handling**

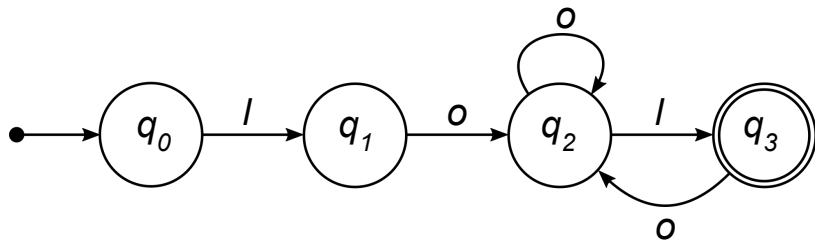
- Spelling Errors
- Minimum Edit Distance
- Noisy Channel Model

# 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  $I(o+I)^+$



Regular Expression

$I(o+I)^+$

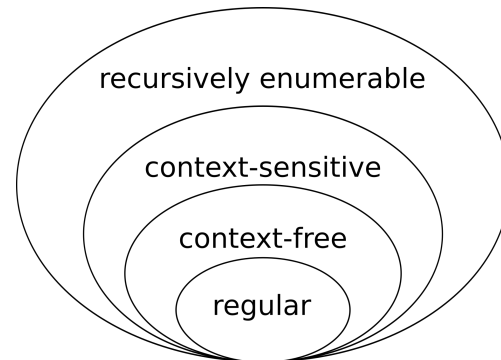


Regular Language

$\{I, Io, IIo, IIoIo, \dots\}$

## Chomsky Hierarchy

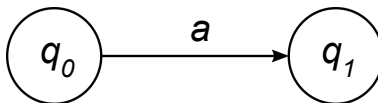
(Source: [Wikipedia](https://en.wikipedia.org/wiki/Chomsky_hierarchy))



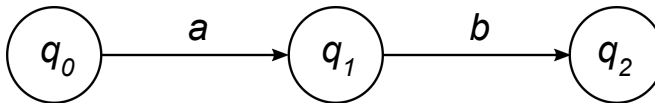
# Relationship to Finite State Automata

- Basic equivalences

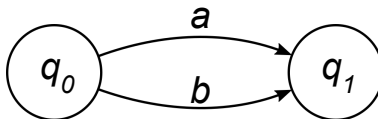
**a**



**ab**



**a | b**



**a\***



# In-Lecture Activity (10 mins)





# Outline

- **Regular Expressions**

- Basic Concepts
- Relationship to FSA
- **Error Types**

- **Corpus Preprocessing**

- Tokenization
- Normalization
- Stemming / Lemmatization
- Segmentation

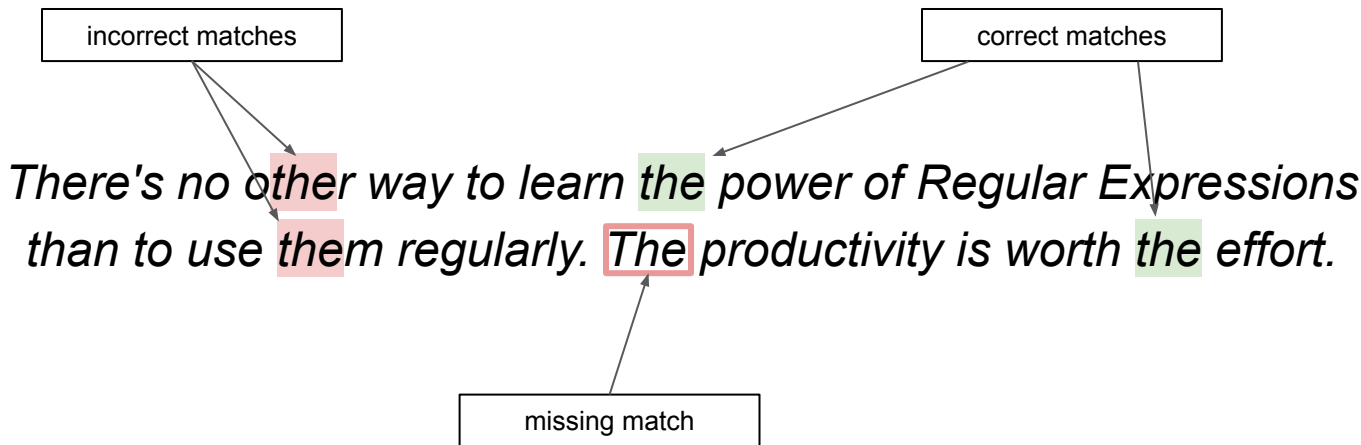
- **Word error handling**

- Spelling Errors
- Minimum Edit Distance
- Noisy Channel Model

# Error Types — What Can Go Wrong



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



# Error Types

- 2 basic types of errors

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



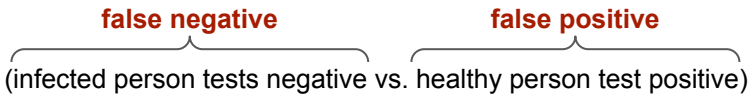
**False Positives**  
(Type I Errors)

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



**False Negatives**  
(Type II Errors)

# 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
  - 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

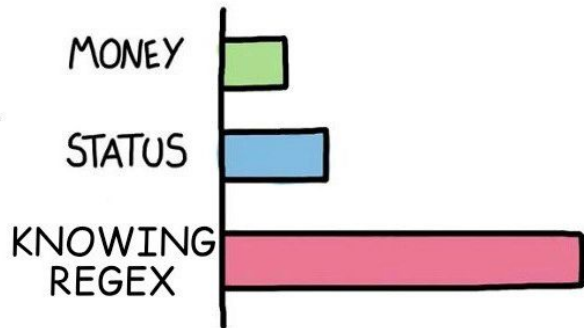
- Know their powers

- 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

# 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")

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

- 3 basic approaches

character-based	S	h	e	'	s	d	r	i	v	i	n	g	f	a	s	t	e	r	t	h	a	n	a	l	l	o	w	e	d	.
subword-based	She	's	driv	ing	fast	er	than	allow	ed	.																				
word-based	She	's	driving	faster	than	allowed	.																							

# Tokenization — Word-Based



- 2 intuitive approaches (solved using RegEx)

- Match all words, numbers and punctuation marks → `\w+ | \d+ | [ , . ; : ]`
- Match boundaries between "words" and "non-words" → `(?=\W) | (?<=\W)`

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

`(?=\W) | (?<=\W)` → NLP|is|fun|and|there|is|so|much|to|learn|in|13|weeks|



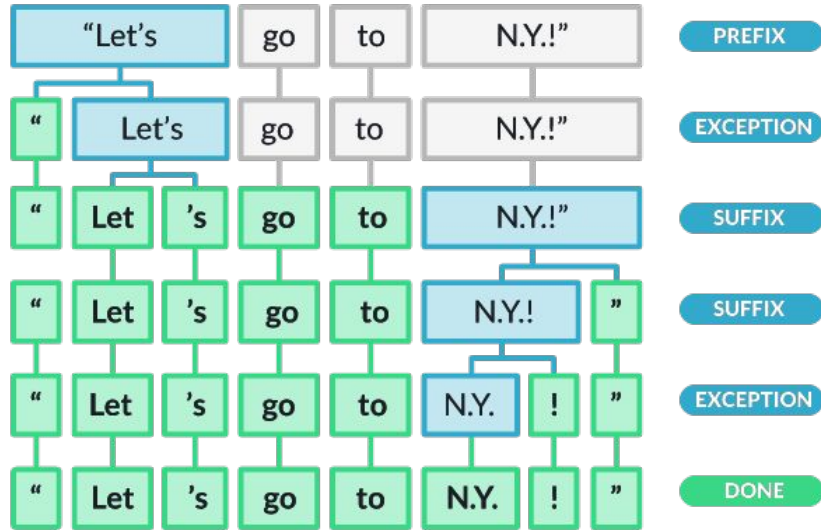
# Tokenization — It Quickly Gets Tricky

- Multiword phrases → *I just came back from New York City.*
- Common contractions → *I'm not home, so don't call.*
- Hyphenations → *NLP is a well-defined but non-trivial topic.*
- Acronyms, names, etc. → *I watched a C++ documentary on T.V.*
- Special tokens → *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


@ B o b . . . I  # D u n e M o v i e : o ) ) )

@ B o b . . . I  # D u n e M o v i e : o ) ) )

@ B o b . . . I  # D u n e M o v i e : o ) ) )

@ B o b . . . I  # D u n e M o v i e : o ) ) )

@ B o b . . . I  # D u n e M o v i e : o ) ) )

@ B o b . . . I  # D u n e M o v i e : o ) ) )

Sequential labeling of characters

Label all whitespace characters



Label all unicode characters



Label all emoticons



Label all special token types



Label all punctuation marks



Label all all alphanumeric characters

→ Tokens = Substrings with adjacent characters with the same labels

# Tokenization — Language Issues

- French

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

direct article

*l'ensemble*

"the whole" / "all"

indicates imperative

*donne-moi*

"give me!"

→ 1 token or 2 tokens?

- German

- Very common: compound nouns

*Arbeiterunfallversicherungsgesetz*

"worker injury insurance act"

→ important: **compound splitter**

# Tokenization — Language Issues

- Languages without whitespaces separating words

## Chinese

莎 拉 波 娃 | 现 在 | 居 住 | 在 | 美 国 | 东 南 部 | 的 | 佛 罗 里 达  
" Sharapova now lives in US southeastern Florida "

## Japanese

- multiple syllabaries
- multiple formats for dates and amounts

フォーチュン500社は情報不足のため時間あた\$500K(約6,000万円)

Katakana      Hiragana      Kanji      Romanji

# Tokenization — Word Segmentation of Chinese Text

- Baseline algorithm: **Maximum Matching**



莎拉波娃现在居住在美国东南部的佛罗里达



**莎拉波娃**现在居住在美国东南部的佛罗里达  
*Sharapova*



莎拉波娃|现在居住在美国东南部的佛罗里达

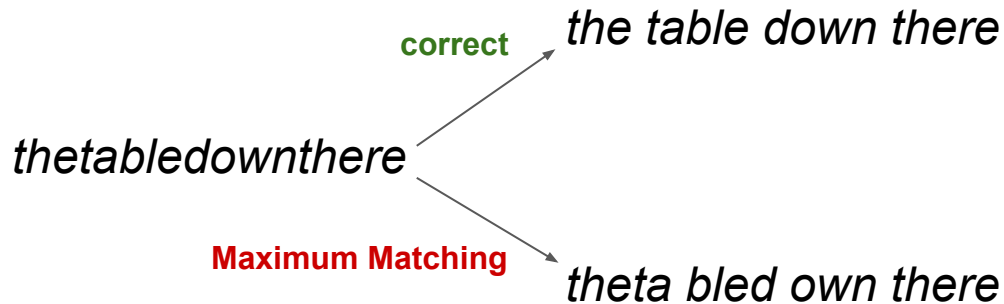


莎拉波娃|**现在**居住在美国东南部的佛罗里达  
*now*

- (1) Place a pointer at the beginning of the string
- (2) Find longest word in dictionary that matches string starting the pointer
- (3) Mover the pointer over the word in the string
- (4) Goto #2 to process the whole string

# Tokenization — Maximum Matching

- Surprisingly good performance on Chinese text  
(even better performance with probabilistic methods or extensions)
- Generally does not work for English text



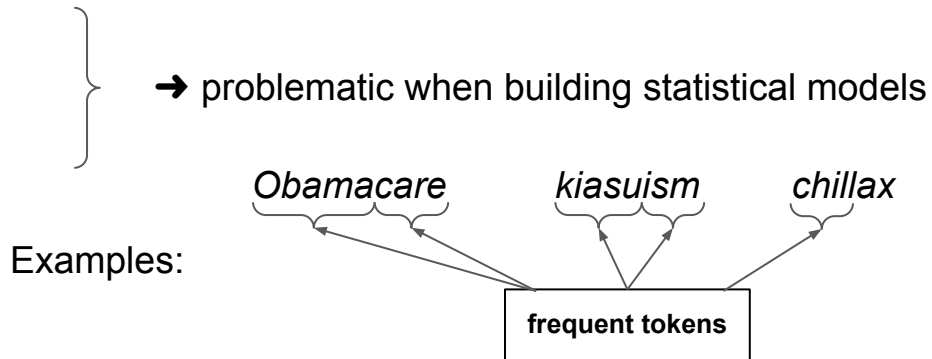
# Tokenization — Subword-Based

- 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



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



# Tokenization — Subword-Based

- Different algorithms for subword tokenization
  - Byte-Pair Encoding (BPE), 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

# Tokenization — BPE Token Learner



**Corpus:** *"low low low low low lower lower newest newest newest  
newest newest newest widest widest widest longer"*

special end-of-word token

Initialize vocabulary (e.g., {'d', 'e', 'g', 'i', 'l', 'n', 'o', 'r', 's', 't', 'w', \_})

## REPEAT

Find the 2 tokens most frequently adjacent to each other (e.g., 'e', 's')

Add a new merged token 'es' to vocabulary

Replace every adjacent 'e' 's' in corpus with 'es'

**UNTIL** k merges have been done

parameter of algorithm

# Tokenization — BPE Token Learner

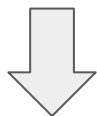
## corpus representation

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

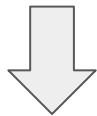
6	n e w <b>es</b> t _
5	l o w _
3	w i d <b>es</b> t _
2	l o w e r _
1	l o n g e r _

## vocabulary

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

## merges

**(e, s)**



most frequent pair: **es** & **t** (9 occurrences)

# Tokenization — BPE Token Learner

## corpus representation

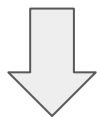
6	n e w <b>est</b> _
5	l o w _
3	w i d <b>est</b> _
2	l o w e r _
1	l o n g e r _

## vocabulary

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

## merges

(e, s), (**es**, t)



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

## corpus representation

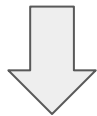
6	n e w <b>est</b> _
5	l o w _
3	w i d <b>est</b> _
2	l o w e r _
1	l o n g e r _

## vocabulary

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

## merges

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



most frequent pair: **l** & **o** (8 occurrences)

# Tokenization — BPE Token Learner

## corpus representation

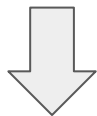
6	n e w est_
5	lo w _
3	w i d est_
2	lo w e r _
1	lo n g e r _

## 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

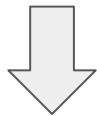
6	n e w est_
5	low _
3	w i d est_
2	low e r _
1	lo n g e r _

## vocabulary

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

## merges

(e, s), (es, t), (est, \_), (1, o), (lo, w)



most frequent pair: n & e (6 occurrences)

# Tokenization — BPE Token Learner

**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), (n, e), (ne, w), (new, est\_), (low, \_), (e, r),  
(er, \_), (w, i), (wi, d), (wid, est\_), (low, er\_), (lo, n), (lon, g), (long, er\_)

Tokenize/segment  
*"newer"*

Run each merge in order  
they have been learned

n e w e r \_  
ne w e r \_  
new e r \_  
new er \_  
new er \_

(n, e)  
(ne, w)  
(e, r)  
(er, \_)

→ tokens: "new", "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

# Normalization

- **Goal: Convert text into a canonical (standard) form**
  - Remove noise / "randomness" from text
  - Affects characters, words, 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
connect connects connected connecting connection	connect
:) :-) :o)	smile

# Normalization — Case Folding

- When to fold?

- Common application: Information Retrieval  
(e.g., Web search where most 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

# 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?

# Normalization — Stemming

- 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
<i>cats</i>	<i>cat</i>
<i>running</i>	<i>run</i>
<i>phones</i>	<i>phon(e)</i>
<i>presumably</i>	<i>presum</i>
<i>crying</i>	<i>cry/cri</i>
<i>went</i>	<i>went</i>
<i>worse</i>	<i>wors</i>
<i>best</i>	<i>best</i>
<i>mice</i>	<i>mic(e)</i>

# 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 stops if a pass yields no more changes

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

# 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)
<i>running</i>	<i>running</i>	<i>run</i>	<i>running</i>
<i>phones</i>	<i>phone</i>	<i>phone</i>	<i>phones</i>
<i>went</i>	<i>went</i>	<i>go</i>	<i>went</i>
<i>worse</i>	<i>worse</i>	<i>worse</i>	<i>bad</i>
<i>mice</i>	<i>mouse</i>	<i>mice</i>	<i>mice</i>



# Normalization — Lemmatization: Characteristics

- Pros

- Lemmatized words are proper words (i.e., dictionary words)
- Can normalize irregular forms (e.g., *went* → *go*, *worst* → *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* → *does n't*; *John's* → *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)
  - ...

# Quick Quiz



# 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

# 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`

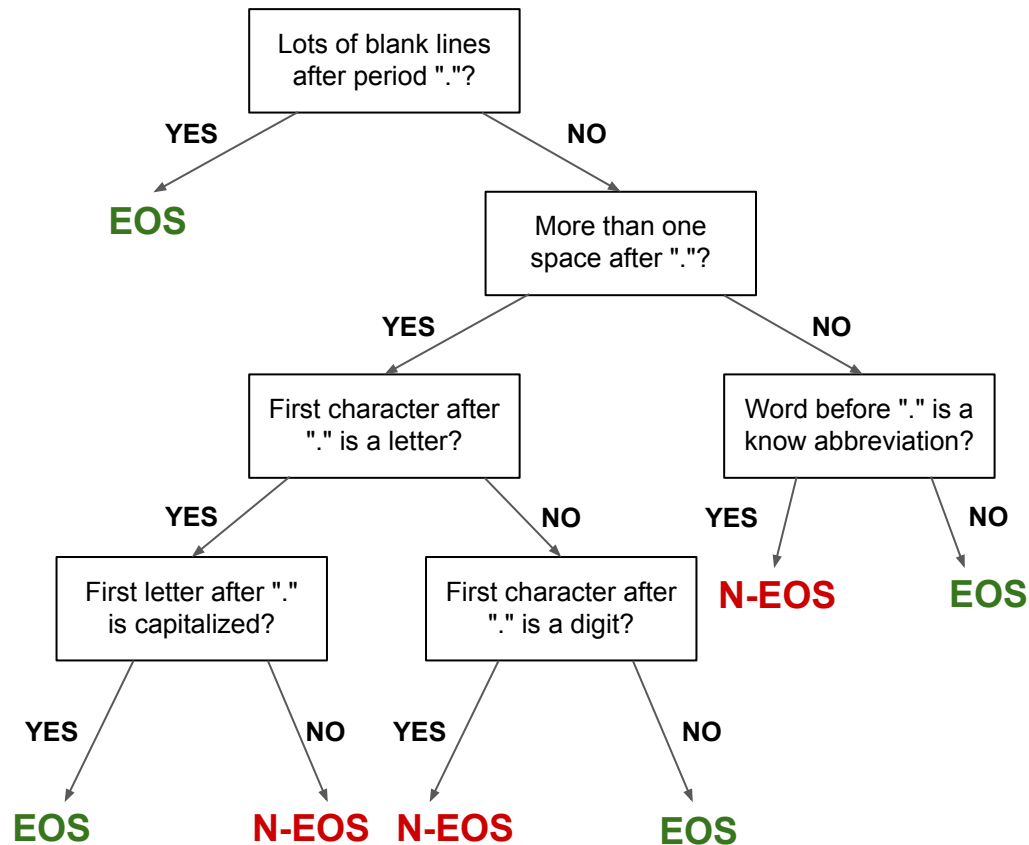
(Source: [Stackoverflow](#))

- Alternative: binary classifier

- 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)



# 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.



# 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

# Spelling Errors

Increasing Complexity

## 1. Non-word error detections

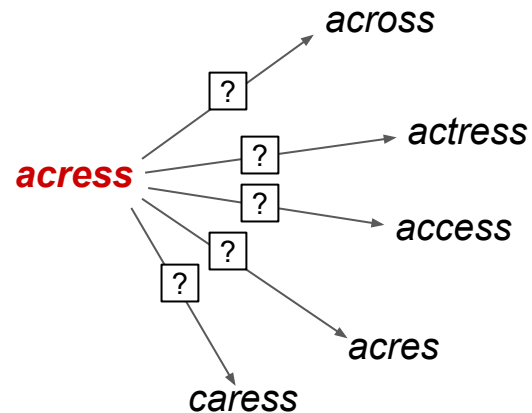
- Basically, word is not found in dictionary
- Example: detecting *graffe* (misspelling of *giraffe*)

## 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 are spelled correctly
- Examples: *there* vs. *three*, *dessert* vs. *desert*, *son* vs. *song*



# 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., *acress* vs. *access*)
- Transposition (e.g., *acress* vs. *caress*)

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

# 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

- $s_1 = \text{"LANGUAGE"}$

- $s_2 = \text{"SAUSAGE"}$

→ **Alignment** of MED:

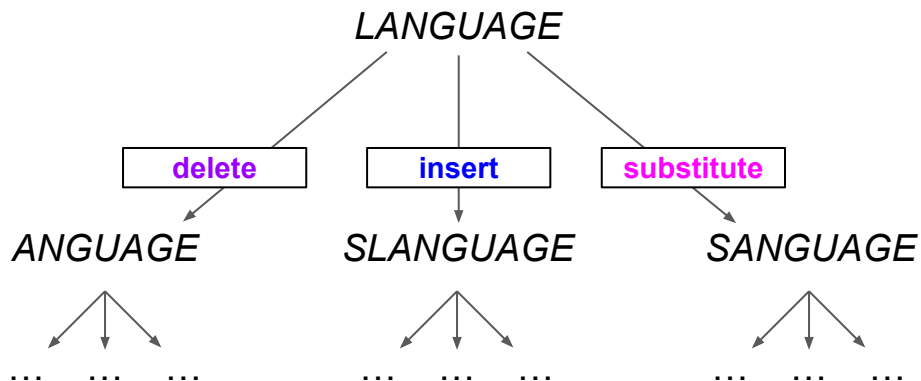
L	A	N	G	U	*	A	G	E
S	A	*	*	U	S	A	G	E

MED if all operations cost 1 → 4

MED if Substitution costs 2,  
Insertion 1, Deletion 1 → 5

# Minimum Edit Distance — Calculation

- 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



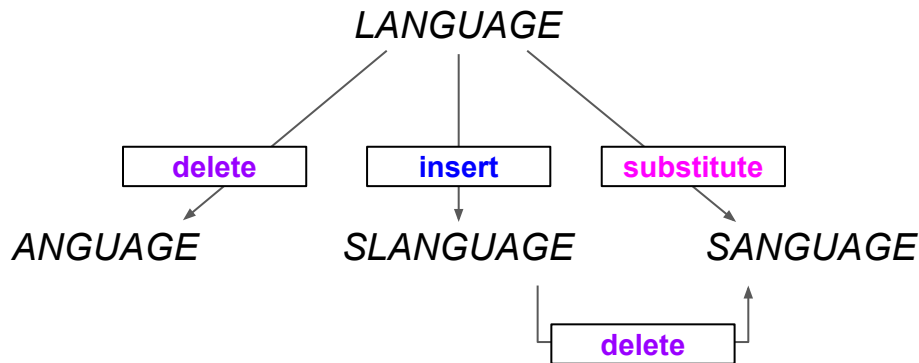
→ Potentially huge search space

→ Naive navigation of all path impractical

# Minimum Edit Distance — Calculation

- 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!)

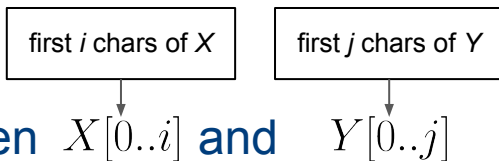
→ Solve using **Dynamic Programming**

solving problems by combining solutions to subproblems

# Minimum Edit Distance — Calculation

- Input: 2 strings

- Source string  $X$  of length  $n$
- Target string  $Y$  of length  $m$



- Define  $D(i, j)$  as MED between  $X[0..i]$  and  $Y[0..j]$

→ MED between  $X$  and  $Y$  is thus  $D(n, m)$

- 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  $i, j$



# Minimum Edit Distance — Calculation

- Initialization of bases cases

- $D(i, 0) = i$  (getting from  $X[0..i]$  to empty target string requires  $i$  deletions)
- $D(0, j) = j$  (getting from empty source string to  $Y[0..j]$  requires  $j$  insertions)

- For  $0 < i \leq n$  and  $0 < j \leq m$

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

Assumptions for costs

Insert: 1

Delete: 1

Substitute: 2

→ Levenshtein MED

Complexity analysis

Space:  $O(nm)$

Time:  $O(nm)$

# Minimum Edit Distance — Calculation Example

<b>E</b>	8								
<b>G</b>	7								
<b>A</b>	6								
<b>U</b>	5								
<b>G</b>	4								
<b>N</b>	3								
<b>A</b>	2								
<b>L</b>	1								
<b>#</b>	0	1	2	3	4	5	6	7	
	<b>#</b>	<b>S</b>	<b>A</b>	<b>U</b>	<b>S</b>	<b>A</b>	<b>G</b>	<b>E</b>	

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

## Minimum Edit Distance — Calculation Example

<b>E</b>	8	9	8	7	8	7	6	5
<b>G</b>	7	8	7	6	7	6	5	6
<b>A</b>	6	7	6	5	6	5	6	7
<b>U</b>	5	6	5	4	5	6	7	8
<b>G</b>	4	5	4	5	6	7	6	7
<b>N</b>	3	4	3	4	5	6	7	8
<b>A</b>	2	3	2	3	4	5	6	7
<b>L</b>	1	2	3	4	5	6	7	8
<b>#</b>	0	1	2	3	4	5	6	7
	<b>#</b>	<b>S</b>	<b>A</b>	<b>U</b>	<b>S</b>	<b>A</b>	<b>G</b>	<b>E</b>

# Minimum Edit Distance — Backtrace & Alignments

- Current limitation

- Base algorithms only returns the MED
- Often important: alignment between strings

L	A	N	G	U	*	A	G	E
S	A	*	*	U	S	A	G	E

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:

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

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

# Minimum Edit Distance — Backtrace & Alignments

E	8	↖←↓ 9	↓ 8	↓ 7	↖←↓ 8	↓ 7	↓ 6	↖ 5
G	7	↖←↓ 8	↓ 7	↓ 6	↖←↓ 7	↓ 6	↖ 5	← 6
A	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
N	3	↖←↓ 4	↓ 3	↖←↓ 4	↖←↓ 5	↖←↓ 6	↖←↓ 7	↖←↓ 8
A	2	↖←↓ 3	↖ 2	← 3	← 4	↖← 5	← 6	← 7
L	1	↖←↓ 2	↖←↓ 3	↖←↓ 4	↖←↓ 5	↖←↓ 6	↖←↓ 7	↖←↓ 8
#	0	1	2	3	4	5	6	7
	#	S	A	U	S	A	G	E



L A N G U \* A G E  
 | | | | | | | |  
 S A \* \* U S A G E

Complexity analysis

Time:  $O(n+m)$

## Minimum Edit Distance — More Examples

- Biology: Align 2 sequences of nucleotides

AGGCTATCACCTGACCTCCAGGCCGATGCCC  
TAGCTATCACGACCGCGGTCGATTTGCCCGAC

[illegible]

\* A G G C T A T C A C C T G A C C T C C A G G C C G A \* \* T G \* C C \* \* C  
 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |  
 T A \* G C T A T C A \* C \* G A C C \* G C \* G G T C G A T T T G C C C G A C

# In-Lecture Activity (10 mins)



# In-Lecture Activity (10 mins)





# Minimum Edit Distance — Other Uses in NLP

- Evaluating Machine Translation and speech recognition

e.g., How similar are 2 translations?

<b>Reference:</b>	Spokesman	confirms	*	senior	government	adviser	was	shot	*
<b>Prediction:</b>	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 — 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:**

$$D(0, 0) = 0$$

$$D(i, 0) = D(i - 1, 0) + \text{del}(X[i]), \quad \text{for } 1 < i \leq n$$

$$D(0, j) = D(0, j - 1) + \text{ins}(Y[j]), \quad \text{for } 1 < j \leq m$$

**Recurrence relation:**

$$D(i, j) = \min \begin{cases} D(i - 1, j) & + \text{del}(X[i]) \\ D(i, j - 1) & + \text{ins}(Y[j]) \\ D(i - 1, j - 1) & + \text{sub}(X[i], Y[j]) \end{cases}$$

# 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 local 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**

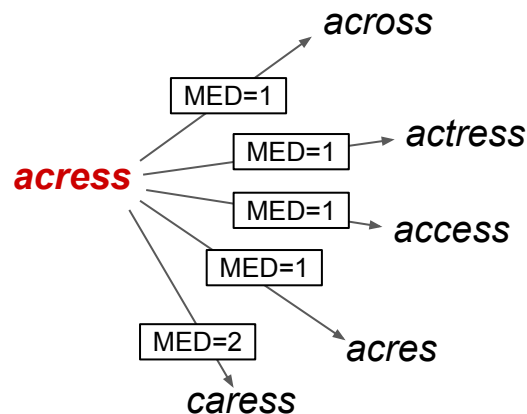
# Where We are Right Now

- 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?

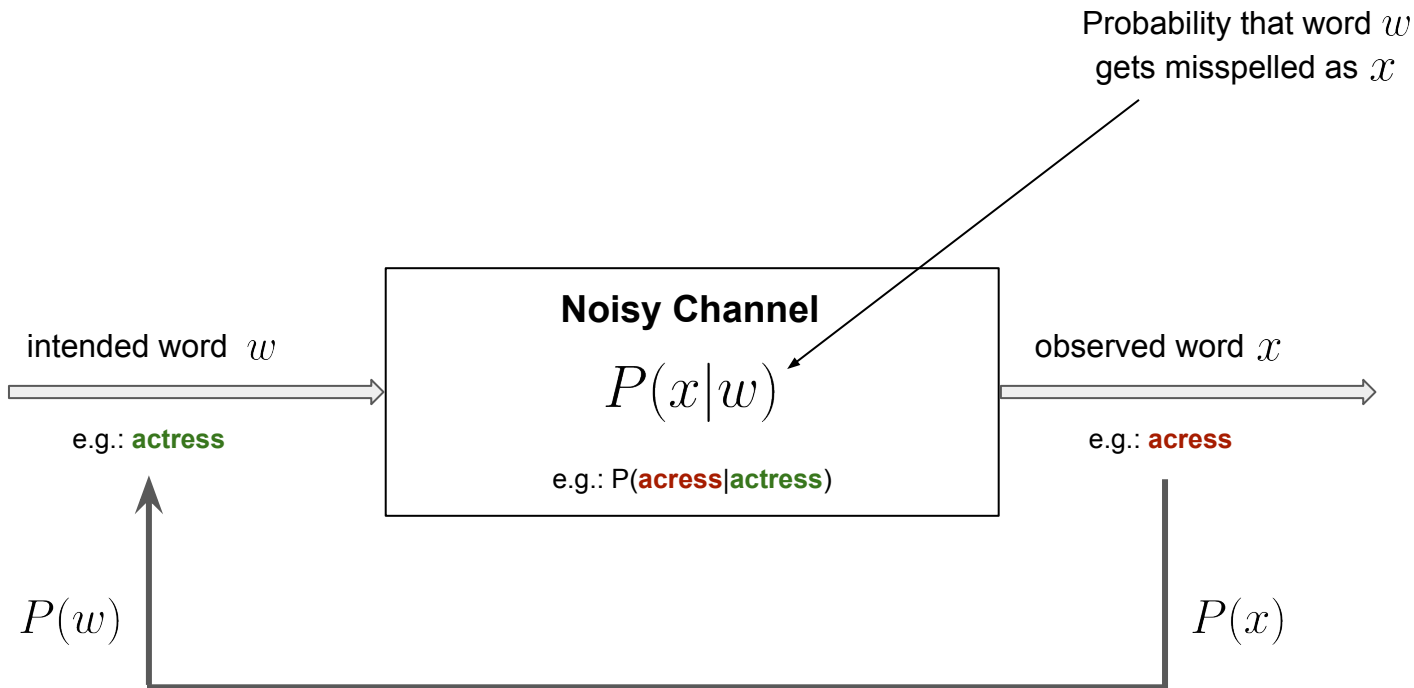
- Ranking of candidates to show top candidates first
- Support for automated spelling correction



## → Noisy Channel Model

Idea: Assign each candidate a probability

# Noisy Channel Model — Intuition



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

# Noisy Channel Model — Bayesian Inferencing

Given an observation  $x$  of a misspelled word,  
find the correct word  $w$ :

$$\hat{w} = \operatorname{argmax}_{w \in V} P(w|x)$$

$$\hat{w} = \operatorname{argmax}_{w \in V} \frac{P(x|w)P(w)}{P(x)}$$

$$\hat{w} = \operatorname{argmax}_{w \in V} P(x|w)P(w)$$

**Quick refresher: Bayes' Theorem**

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

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

$$\rightarrow P(A|B)P(B) = P(B|A)P(A)$$

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

→ How to calculate  $P(x|w)$  and  $P(w)$  ?

# Noisy Channel Model — Calculating/Estimating $P(w)$

- Approach using Maximum Likelihood Estimate (MLE)

- 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

$w$	$freq(w)$	$P(w)$
<i>actress</i>	1,135	0.0000784
<i>cress</i>	1	0.00000...
<i>caress</i>	3	0.00000...
<i>access</i>	1,670	0.0001153
<i>across</i>	1,756	0.0001213
<i>acres</i>	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 $P(x|w)$

- 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 $yx$
$count[x]$	number of times that $x$ appeared in the training set
$count[x, y]$	number of times that $xy$ appeared in the training set

$$x, y \in \{a, b, c, \dots, z\}$$

# 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 $P(x|w)$

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

X	Y (correct)																									
	a	b	c	d	e	f	g	h	i	j	k	l	m	n	o	p	q	r	s	t	u	v	w	x	y	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
j	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
l	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	9	15	13	3	2	2	3	0	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
o	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

# Noisy Channel Model — Example

- Noisy channel probabilities for "acress"

Candidate Correction	Correct Letter	Error Letter	$x w$	$P(x w)$	$P(w)$	$10^9 \cdot P(x w)P(w)$	%
<i>actress</i>	t		c ct	.000117	.0000231	2.7	35.9
<i>cress</i>		a	a #	.00000144	.00000054	.00078	~0
<i>caress</i>	ca	ac	ac ca	.00000164	.00000170	.0028	~0
<i>access</i>	c	r	r c	.00000021	.0000916	.019	~0
<b><i>across</i></b>	<b>o</b>	<b>e</b>	<b>e o</b>	<b>.0000093</b>	<b>.000299</b>	<b>2.8</b>	<b>37.2</b>
<i>acres</i>		s	es e	.0000321	.0000318	1.0	13.3
<i>acres</i>		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."*



"actress" here the better candidate

→ **Language Models** (next lecture)

# 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

# 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!**



# Pre-Lecture Activity for Next Week



# Solutions to Quick Quizzes



# Solutions to Quick Quizzes

