# Natural Language Processing (NLP)

## Basic Text Processing
### Based on Jurafsky & Martin (2018), Chapter 2

Katrien Beuls

Artificial Intelligence Laboratory
Vrije Universiteit Brussel

VRIJE
UNIVERSITEIT
BRUSSEL

ELIZA

```
User:  You are like my father in some ways.
ELIZA: WHAT RESEMBLANCE DO YOU SEE
User:  You are not very aggressive but I think you don't want me to
       notice that.
ELIZA: WHAT MAKES YOU THINK I AM NOT VERY AGGRESSIVE
User:  You don't argue with me.
ELIZA: WHY DO YOU THINK I DON'T ARGUE WITH YOU
User:  You are afraid of me.
ELIZA: DOES IT PLEASE YOU TO BELIEVE I AM AFRAID OF YOU
```
Weizenbaum (1966)

► A formal language for specifying text strings

► How can we search for any of these?
   ► woodchuck
   ► woodchucks
   ► Woodchuck
   ► Woodchucks

► Letters inside square brackets: `[]`

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

► Ranges: `[A-Z]`

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

▶ Negations: `[^Ss]`

| Pattern | Matches | |
|---------|---------|---|
| `[^A-Z]` | Not an upper case letter | `O`<u>`y`</u>`fn pripetchik` |
| `[^Ss]` | Neither 'S' nor 's' | <u>`I`</u>` have no reason` |
| `[^\.]` | Not a period | <u>`o`</u>`ur resident Djinn` |
| `a^b` | The pattern 'a^b' | `look up a^b now` |

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

| Pattern | Matches |
|---|---|
| `groundhog|woodchuck` | groundhog<br>woodchuck |
| `gupp(y|ies)` | guppy<br>guppies |
| `a|b|c` | = `[abc]` |
| `[gG]roundhog|[Ww]oodchuck` | |

# REGULAR EXPRESSIONS

`? * + .`

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

## ANCHORS ^ $

| Pattern | Matches |
|---|---|
| `^[A-Z]` | <u>P</u>alo Alto |
| `^[^A-Za-z]` | <u>1</u>    <u>"</u>Hello" |
| `\.$` | The end<u>.</u> |
| `.$` | The end<u>?</u>    The end<u>!</u> |

- Find all instances of the word "the" in a text.

  - `/the/`
    Misses capitalised examples

  - `/[tT]he/`
    Incorrectly returns `other` or `theology`

  - `/[^a-zA-Z][tT]he[^a-zA-Z]/`
    Does not return "the" when it begins a line

  - `/(^|[^a-zA-Z])[tT]he([^a-zA-Z]|$)/`

- From highest to lowest precedence:

| | |
|---|---|
| **Parenthesis** | `()` |
| **Counters** | `*  +  ?    {}` |
| **Sequences and anchors** | `the   ^my   end$` |
| **Disjunction** | `|` |

- **Counters > Sequences**: `/the*/` matches *theeeee* but not *thethe*

- **Sequences > Disjunction**: `/the|any/` matches *the* or *any* but not *theny*

▶ Substitution operator `s/regexp1/pattern/` used in Unix commands like `vim` or `sed` allows a string characterized by a regular expression to be replaced by another string:

$$s/colour/color$$

▶ Referring to a subpart of the string matching the first pattern: e.g. put angle brackets around all integers in a text:

$$s/[0-9]+/<\backslash1>/$$

```
/the (.*)er they were, the \1er they will be/
```

► Parenthesis used for storing a pattern in memory
   = **a capture group**
   ► Resulting match is stored in a numbered **register**
   ► If you match two different sets of parentheses, \2 means
     whatever matched the *second* capture group
     ```
     /the (.*)er they (.*), the \1er we \2 /
     ```
► Use a **non-capturing group** if you don't want to capture the
  resulting pattern in a register:
     ```
     /((?:some|a few) (people|cats) like some \1/
     ```

- ► Works by having a cascade of (ranked) regular expression substitutions
    - ► Input lines are first uppercased
    - ► First substitutions change all instances of *MY* to *YOUR* and *I'M* to *YOU ARE*
    - ► The next set of substitutions matches and replaces other patterns in the input

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

- How many words are in the following Brown sentence?
  - *He stepped out into the hall, was delighted to encounter a water brother.*
  - 13 words (15 with punctuation marks)
- How many words are in the following utterance from the Switchboard corpus?
  - *I do uh main- mainly business data processing*
  - Two disfluencies: a fragment (*main-*) and a filled pause (*uh*)

- How about inflected forms like *cats* versus *cat*?
  - They have the same lemma *cat* but different wordforms
  - A **lemma** is a set of lexical forms having the same stem, the same major part-of-speech, and the same word sense
  - The **wordform** is the fully inflected or derived form of the word

- **Types** are the number of distinct words in a corpus
  - If the set of words in the vocabulary is *V*, the number of types is the vocabulary size |*V*|
- **Tokens** are the total number *N* of running words
- If we ignore punctuation, the following Brown sentence has 16 tokens and 14 types:
  - *They picknicked by the pool, then lay back on the grass and looked at the stars*.

| Corpus | Tokens = $N$ | Types = $|V|$ |
|---|---:|---:|
| Shakespeare | 884 thousand | 31 thousand |
| Brown corpus | 1 million | 38 thousand |
| Switchboard telephone conversations | 2.4 million | 20 thousand |
| COCA | 440 million | 2 million |
| Google N-grams | 1 trillion | 13 million |

The relationship between the number of types and number of tokens is called **Herdan's Law** (Herdan, 1960) or **Heap's Law** (Heaps, 1978), where *k* and $\beta$ are positive constants and $0 < \beta < 1$:

$$|V| = kN^{\beta}$$

- Look at the **number of lemmas** instead of wordform types
- **Dictionaries** can help in giving lemma counts
    - Dictionary entries or boldface forms are a very rough upper bound on the number of lemmas
    - The 1989 edition of the Oxford English Dictionary had 615,000 entries

At least three tasks are commonly applied as part of any normalization process:

1. Segmenting/tokenizing words from running text
2. Normalizing word formats
3. Segmenting sentences in running text

Main challenges:

- ▶ Break off **punctuation** as a separate token but preserve it when it occurs word internally (*Ph.D.*, *AT&T*, ...)
- ▶ Keep **special characters** and **numbers** in prices ($45.55) and dates (15/02/2019)
- ▶ Expand **clitic** contractions that are marked by apostrophes (*we're → we are*)
- ▶ Tokenize **multiword expressions** like *New York* or *rock 'n' roll* as a single token

**Penn Treebank tokenization** standard separates out clitics (*doesn't* becomes *does* plus *n't*), keeps hyphenated words together, and separates out all punctuation:

**Input**: "The San Francisco-based restaurant," they said, "doesn't charge $10".

**Output**:

| " | The | San | Francisco-based | restaurant | , | " | they |
|---|-----|-----|-----------------|------------|---|---|------|

| said | , | " | does | n't | charge | $ | 10 | " | . |
|------|---|---|------|-----|--------|---|----|---|---|

Choosing a **single normalized form** for words with multiple forms such as *USA* and *US*.

- ▶ Valuable for **information retrieval** if you want to query for *US* to match a document that has *USA*
- ▶ In **information extraction** we might want to extract coherent information that is consistent across differently-spelled instances

**Case folding** is another kind of normalization

- ► For tasks like speech recognition and information retrieval, everything is mapped to lower case
- ► For sentiment analysis and other text classification tasks, information extraction and machine translation, case is quite helpful and case folding is generally not done

Because tokenization needs to be run before any other
language processing, it is important to be **very fast**.

- **Deterministic algorithms based on regular expressions**
  compiled into very efficient finite state automata
- Carefully designed deterministic algorithms can deal with
  the **ambiguities** that arise (e.g. the apostrophe)

- **Lemmatization** is the task of determining that two words have the same root, despite their surface differences
- How is lemmatization done?
    - Most sophisticated methods for lemmatization involve complete **morphological parsing** of the word
    - Two broad classes of morphemes can be distinguished:
        1. **Stems** - the central morpheme of the word, supplying the main meaning
        2. **Affixes** - adding "additional" meanings of various kinds

- **Stemming** is a simpler but cruder method, which mainly consists of chopping off word-final affixes
- One of the most widely used stemming algorithms for English is **the Porter stemmer** (1980)

```
This was not the map we found in Billy Bones's chest, but
an accurate copy, complete in all things-names and heights
and soundings-with the single exception of the red crosses
and the written notes.

Thi wa not the map we found in Billi Bone s chest but an
accur copi complet in all thing name and height and sound
with the singl except of the red cross and the written note
```

▶ The Porter stemmer algorithm is based on a series of rewrite rules in series, as a **cascade**

ATIONAL $\rightarrow$ ATE  (e.g. relational $\rightarrow$ relate)
　　ING $\rightarrow$ $\epsilon$  if stem contains vowel (e.g., motoring $\rightarrow$ motor)
　SSES $\rightarrow$ SS  (e.g., grasses $\rightarrow$ grass)

▶ Detailed rule lists for the Porter stemmer can be found on Martin Porter's homepage

**Sentence segmentation** is another important step in text processing

- ▶ The most useful cues are **punctuation** (periods, question marks, exclamation points)
- ▶ Periods can be ambiguous: *Mr.* or *Inc.*

In general, **sentence tokenization methods** work by building a **binary classifier** that decides if a period is a part of the word or is a sentence-boundary marker

# MINIMUM EDIT DISTANCE

## STRING SIMILARITY

- Calculating the similarity between two strings is useful in many NLP tasks, such as spelling correction or coreference resolution

- The **minimum edit distance** between two strings is defined as the minimum number of editing operations (operations like insertion, deletion, substitution) needed to transform one string into another

```
I N T E * N T I O N
| | | | | | | | | |
* E X E C U T I O N
d s s   i s
```

► The gap between *intention* and *execution* is 5

- ► We can also assign a particular **cost** or **weight** to each of these operations
- ► **Levenshtein distance**: Each insertion or deletion has a cost of 1 and substitutions are not allowed
  - ► (This is equivalent to allowing substitution, but giving each substitution a cost of 2 since any substitution can be represented by one insertion and one deletion)
- ► Using this metric, the Levenshtein distance between *intention* and *execution* is 8

- **How do we find the minimum edit distance?**
    - We can think of this as a search task, in which we are searching for **the shortest path - a sequence of edits - from one string to another**
- The space of all possible edits is enormous, so we can't search naively
- However, lots of distinct edit paths will end up in the same state (string), so rather than recomputing all those paths, we could just remember the shortest path to a state each time we saw it

- **Dynamic programming** is the name for a class of algorithms, first introduced by Bellman (1957), that apply **a table-driven method to solve problems by combining solutions to sub-problems**
- Some of the most commonly used algorithms in NLP make use of dynamic programming, such as the Viterbi algorithm and the CKY algorithm for parsing

```
i n t e n t i o n
                        ←—— delete i
n t e n t i o n
                        ←—— substitute n by e
e t e n t i o n
                        ←—— substitute t by x
e x e n t i o n
                        ←—— insert u
e x e n u t i o n
                        ←—— substitute n by c
e x e c u t i o n
```

## ALGORITHM

**function** MIN-EDIT-DISTANCE(*source*, *target*) **returns** *min-distance*

$n \leftarrow$ LENGTH(*source*)
$m \leftarrow$ LENGTH(*target*)
Create a distance matrix *distance[n+1,m+1]*

# *Initialization: the zeroth row and column is the distance from the empty string*
    $D[0,0] = 0$
    **for** each row $i$ **from** 1 **to** $n$ **do**
        $D[i,0] \leftarrow D[i-1,0] + del\text{-}cost(source[i])$
    **for** each column $j$ **from** 1 **to** $m$ **do**
        $D[0,j] \leftarrow D[0,j-1] + ins\text{-}cost(target[j])$

# *Recurrence relation:*
**for** each row $i$ **from** 1 **to** $n$ **do**
    **for** each column $j$ **from** 1 **to** $m$ **do**
      $D[i, j] \leftarrow$ MIN( $D[i-1,j] + del\text{-}cost(source[i])$,
                       $D[i-1,j-1] + sub\text{-}cost(source[i], target[j])$,
                       $D[i,j-1] + ins\text{-}cost(target[j]))$
# *Termination*
**return** $D[n,m]$

# MINIMUM EDIT DISTANCE

## THE EDIT DISTANCE MATRIX

| Src\Tar | # | e | x | e | c | u | t | i | o | n |
|---|---|---|---|---|---|---|---|---|---|---|
| # | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| i | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 6 | 7 | 8 |
| n | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 7 | 8 | 7 |
| t | 3 | 4 | 5 | 6 | 7 | 8 | 7 | 8 | 9 | 8 |
| e | 4 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 9 |
| n | 5 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 10 |
| t | 6 | 5 | 6 | 7 | 8 | 9 | 8 | 9 | 10 | 11 |
| i | 7 | 6 | 7 | 8 | 9 | 10 | 9 | 8 | 9 | 10 |
| o | 8 | 7 | 8 | 9 | 10 | 11 | 10 | 9 | 8 | 9 |
| n | 9 | 8 | 9 | 10 | 11 | 12 | 11 | 10 | 9 | 8 |

# MINIMUM EDIT DISTANCE

## PRODUCING AN ALIGNMENT

|   | # | e | x | e | c | u | t | i | o | n |
|---|---|---|---|---|---|---|---|---|---|---|
| # | 0 | ← 1 | ← 2 | ← 3 | ← 4 | ← 5 | ← 6 | ← 7 | ← 8 | ← 9 |
| i | ↑ **1** | ↖←↑ 2 | ↖←↑ 3 | ↖←↑ 4 | ↖←↑ 5 | ↖←↑ 6 | ↖←↑ 7 | ↖ 6 | ← 7 | ← 8 |
| n | ↑ 2 | ↖←↑ **3** | ↖←↑ 4 | ↖←↑ 5 | ↖←↑ 6 | ↖←↑ 7 | ↖←↑ 8 | ↑ 7 | ↖←↑ 8 | ↖ 7 |
| t | ↑ 3 | ↖←↑ 4 | ↖←↑ **5** | ↖←↑ 6 | ↖←↑ 7 | ↖←↑ 8 | ↖ 7 | ←↑ 8 | ↖←↑ 9 | ↑ 8 |
| e | ↑ 4 | ↖ 3 | ← 4 | ↖← **5** | ← **6** | ← 7 | ←↑ 8 | ↖←↑ 9 | ↖←↑ 10 | ↑ 9 |
| n | ↑ 5 | ↑ 4 | ↖←↑ 5 | ↖←↑ 6 | ↖←↑ 7 | ↖←↑ **8** | ↖←↑ 9 | ↖←↑ 10 | ↖←↑ 11 | ↖↑ 10 |
| t | ↑ 6 | ↑ 5 | ↖←↑ 6 | ↖←↑ 7 | ↖←↑ 8 | ↖←↑ 9 | ↖ **8** | ← 9 | ← 10 | ←↑ 11 |
| i | ↑ 7 | ↑ 6 | ↖←↑ 7 | ↖←↑ 8 | ↖←↑ 9 | ↖←↑ 10 | ↑ 9 | ↖ **8** | ← 9 | ← 10 |
| o | ↑ 8 | ↑ 7 | ↖←↑ 8 | ↖←↑ 9 | ↖←↑ 10 | ↖←↑ 11 | ↑ 10 | ↑ 9 | ↖ **8** | ← 9 |
| n | ↑ 9 | ↑ 8 | ↖←↑ 9 | ↖←↑ 10 | ↖←↑ 11 | ↖←↑ 12 | ↑ 11 | ↑ 10 | ↑ 9 | ↖ **8** |

- The algorithm **allows arbitrary weights** on the operations
  - For spelling correction, substitutions are more likely to happen between letters that are next to each other on the keyboard
- The **Viterbi algorithm** is a probabilistic extension of minimum edit distance
  - Viterbi computes the "maximum probability alignment" of one string with another (cf. Ch. 8 on POS tagging)

## SUMMARY

- The **regular expression** language is a powerful tool for pattern-matching.

- Basic operations in regular expressions include **concatenation** of symbols, **disjunction** of symbols ([], |, and ̀), **counters** (*, +, and {n,m}), anchors (ˆ, $) and **precedence** operators ((,)).

- **Word tokenization and normalization** are generally done by cascades of simple regular expressions substitutions or finite automata.

- The **Porter algorithm** is a simple and efficient way to do **stemming**, stripping off affixes. It does not have high accuracy but may be useful for some tasks.

- The **minimum edit distance** between two strings is the minimum number of operations it takes to edit one into the other. Minimum edit distance can be computed by **dynamic programming**, which also results in an **alignment** of the two strings.