# Introduction to Natural Language Processing

<span style="color:red">Out of Vocabulary Forms</span>

<span style="color:red">Spelling Error correction</span>

**Jean-Cédric Chappelier**

Jean-Cedric.Chappelier@epfl.ch

and

**Martin Rajman**

Martin.Rajman@epfl.ch

Artificial Intelligence Laboratory

LIA
I&C

Introduction to Natural Language Processing (CS-431)

M. Rajman
J.-C. Chappelier

1/32

# Contents

➥ Out of Vocabulary Forms

➥ Spelling Error Correction

  ✈ Edit distance

  ✈ Spelling error correction with FSA

  ✈ Weighted edit distance

*LIA*
*I&C*
ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

Introduction to Natural Language Processing (CS-431)

M. Rajman
J.-C. Chappelier

2/32

# Out of Vocabulary forms

- Out of Vocabulary (OoV) forms matter: they occur quite frequently (e.g. $\simeq 10\%$ in newspapers)

  What do they consist of?

  - spelling errors: *foget*, *summmary*, *usqge*, ...

  - neologisms: *Internetization*, *Tacherism*, ...

  - borrowings: *gestalt*, *rendez-vous*, ...

  - forms difficult to exhaustively lexicalize: (numbers,) proper names, abbreviations, ...

- identification based on patterns is not well-adapted for all OoV forms

  ☞ We will focus here on spelling errors, neologisms and borrowings

# Spelling errors and neologisms

- for **spelling errors** (resp. **neologisms**), distortions (resp. derivations) are modelled by *transformations*, i.e. rewriting rules (sometimes weighted)

  *Example*:

  – Transposition (distortion): `XY → YX [1.0]`

     where `X` and `Y` stands for variables

  – tripling (distortion): `XX → XXX [1.0]`

  – name derivation: `ize:INF → ization:N [1.0]`

- a given lexicon (regular language) and a set of transformations define the edit space to be explored

  ☞ The aim is to find the position of the OoV forms in the edit space with respect to known (lexicalized) forms (*neighbourhoods*, *similarity*, *distance*)

## Spelling errors and neologisms (2)

- if the transformation set is simple enough: automatic (or semi-automatic) learning of the transformation set is possible

  Examples:

  – morphological rules for Spanish

  – transformations for spelling error correction after OCR

LIA
I&C
ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

Introduction to Natural Language Processing (CS-431)

M. Rajman
J.-C. Chappelier

5/32

## Borrowings

For **borrowings** ☞ **identification of the source language**

$\Big[$ when no large coverage lexica are avalaible for the other languages, but only

representative texts $\Big]$

Decomposition into $n$-grams of characters: *Example*: for trigrams

$$\text{dribble} \rightarrow (\text{dri,rib,ibb,bbl,ble})$$

In practice: $n$ varies from $2$ to $4$

From reference corpora, computation of a frequency matrix ($n$-gram $\times$ language)

☞ approximation of likelihood of a word to belong to a given language

*Example* for trigrams:

$$V(\text{dribble}, L) = P(\text{dribble}|L) = P(\text{dri}|L) \cdot \frac{P(\text{rib}|L)}{P(\text{ri}|L)} \cdot ... \cdot \frac{P(\text{ble}|L)}{P(\text{bl}|L)}$$

Trigrams for French, English, German and Spanish: $87\%$ discrimination accuracy

LIA I&C
ÉCOLE POLYTECHNIQUE FÉDÉRALE DE LAUSANNE

Introduction to Natural Language Processing (CS-431)

M. Rajman
J.-C. Chappelier

6/32

# $n$-**gram approach**

A sequence of $x$s (letters, words, ...; was letters in the former case)

$(n-1)$-order Markov assumption:

$$P(x_1 \cdots x_N) = P(x_1 \cdots x_{n-1}) \cdot \prod_{i=n}^{N} P(x_i | x_{i-n+1} \cdots x_{i-1})$$

i.e. use this as a score to compare sequences ($n \geq 2$):

$$\frac{\displaystyle\prod_{i=1}^{N-n+1} P(x_i \cdots x_{i+n-1})}{\displaystyle\prod_{i=2}^{N-n+1} P(x_i \cdots x_{i+n-2})} \qquad P(x_i \cdots x_{i+n-1}) \text{ being estimated on some corpus}$$

Reminder: $P(x_i \cdots x_{i+n-2}) = \displaystyle\sum_{x} P(x_i \cdots x_{i+n-2}\, x)$

LIA
I&C
ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

Introduction to Natural Language Processing (CS-431)

M. Rajman
J.-C. Chappelier

7/32

# Contents

- Out of Vocabulary Forms

➤ Spelling Error Correction

  ✤ Edit distance

  ✤ Spelling error correction with FSA

  ✤ Weighted edit distance

*LIA*
*I&C*
ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

Introduction to Natural Language Processing (CS-431)

M. Rajman
J.-C. Chappelier

8/32

# Spelling error correction

all strings

input string

solutions

max. distance

correct strings

Two approaches:

|  | Exact lexicon-based | Statistical |
|---|---|---|
| correct forms: | lexicon | any string |
| metric: | edit distance | probability |

In this lecture: exact, lexicon-based

For statistical: see http://norvig.com/spell-correct.html

# Edit distance

also called Levenshtein distance

☞ distance between 2 forms

= minimal number of transformations to change one into the other

Example of transformations:

❋ insertion: *exmple → example*

❋ deletion: *example → exmple*

❋ substitution: *exemple → example*

❋ transposition: *exmaple → example*

LIA
I&C

ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

# Computation of edit distance (1)

Notations:

$X_i$: $i$th char of string $X$

$X_i^j$: if $i \leq j$: substring $X_i,...,X_j$; empty string otherwise

Example: $X$ = castle

$X_3$ = s $\qquad\qquad$ $X_4^6$ = tle $\qquad\qquad$ $X_1^4$ = cast $\qquad\qquad$ $X_1^0 = \varepsilon$

Computation of the distance $D(X, Y)$ by dynamic programming:

☞ step by step in a chart $m$ where each cell $m_{ij}$ contains the distance between the two substrings $X_1^i$ and $Y_1^j$ :
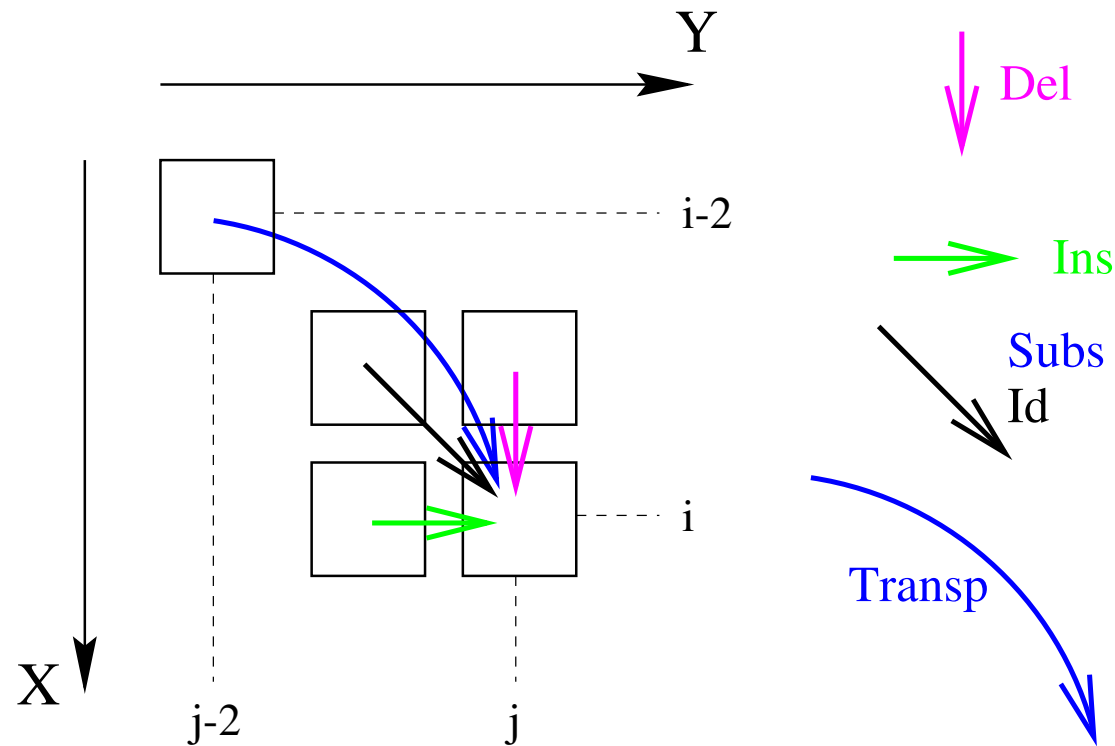
$$m_{ij} = D(X_1^i, Y_1^j)$$

LIA
I&C
ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

Introduction to Natural Language Processing (CS-431)

M. Rajman
J.-C. Chappelier

11/32

# Computation of edit distance (2)

$$D(X_1^0; Y_1^j) = \text{j} \qquad\qquad\qquad \text{initialization}$$

$$D(X_1^i; Y_1^0) = \text{i}$$

---

$$D(X_1^i; Y_1^j) = D(X_1^{i-1}; Y_1^{j-1}) \qquad\qquad \textbf{if } X_i = Y_j \text{ (equality)}$$

$$= 1 + \min \left\{ D(X_1^{i-2}; Y_1^{j-2}), \right. \qquad \textbf{else if } i \geq 2 \text{ and } j \geq 2$$

$$\left. D(X_1^{i-1}; Y_1^j), D(X_1^i; Y_1^{j-1}) \right\} \qquad \text{and } X_{i-1} = Y_j \text{ and } X_i = Y_{j-1} \text{ (transposition, deletion, insertion)}$$

$$= 1 + \min \left\{ D(X_1^{i-1}; Y_1^{j-1}), \right. \qquad \textbf{else } \text{(substitution, deletion, insertion)}$$

$$\left. D(X_1^{i-1}; Y_1^j), D(X_1^i; Y_1^{j-1}) \right\}$$

LIA
I&C
ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

Introduction to Natural Language Processing (CS-431)

M. Rajman
J.-C. Chappelier

12/32

# Computation order

several possible ways of computing: rowwise, columnwise or diagonal

## Computation of edit distance (3)

Example, columnwise:

> **for all** $i$ from $0$ to $|X|$ (size of $X$) **do**
>
> $\quad m_{i0} = i$
>
> **for all** $j$ from $1$ to $|Y|$ **do**
>
> $\quad m_{0j} = j$
>
> $\quad$ **for all** $i$ from $1$ to $|X|$ **do**
>
> $\quad\quad$ **if** $X_i = Y_j$ **then**
>
> $\quad\quad\quad m_{ij} = m_{i-1,j-1}$
>
> $\quad\quad$ **else if** $i \geq 2$ and $j \geq 2$ and $X_{i-1} = Y_j$ and $X_i = Y_{j-1}$ **then**
>
> $\quad\quad\quad m_{ij} = 1 + \min\Big\{ m_{i-2,j-2};\ m_{i,j-1};\ m_{i-1,j} \Big\}$
>
> $\quad\quad$ **else**
>
> $\quad\quad\quad m_{ij} = 1 + \min\Big\{ m_{i-1,j-1};\ m_{i,j-1};\ m_{i-1,j} \Big\}$
>
> **Return** $m_{|X|,|Y|}$

LIA
I&C
ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

M. Rajman
J.-C. Chappelier

14/32

# Edit Distance (example)

D(exmple;exemple)

|   |   | e | x | e | m | p | l | e |
|---|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| e | 1 | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| x | 2 | 1 | 0 | 1 | 2 | 3 | 4 | 5 |
| m | 3 | 2 | 1 | 1 | 1 | 2 | 3 | 4 |
| p | 4 | 3 | 2 | 2 | 2 | 1 | 2 | 3 |
| l | 5 | 4 | 3 | 3 | 3 | 2 | 1 | 2 |
| e | 6 | 5 | 4 | 3 | 4 | 3 | 2 | 1 |

D(exmaple;example)

|   |   | e | x | a | m | p | l | e |
|---|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| e | 1 | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| x | 2 | 1 | 0 | 1 | 2 | 3 | 4 | 5 |
| m | 3 | 2 | 1 | 1 | 1 | 2 | 3 | 4 |
| a | 4 | 3 | 2 | 1 | 1 | 2 | 3 | 4 |
| p | 5 | 4 | 3 | 2 | 2 | 1 | 2 | 3 |
| l | 6 | 5 | 4 | 3 | 3 | 2 | 1 | 2 |
| e | 7 | 6 | 5 | 4 | 4 | 3 | 2 | 1 |

LIA
I&C
ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

Introduction to Natural Language Processing (CS-431)

M. Rajman
J.-C. Chappelier

15/32

# Spelling error correction using a FSA

Problem: approximative search of lexicalized (surface) forms

= within a max. distance range

i.e. Fault-tolerant recognition (within a regular language):

Find **all** ending paths such that the corresponding string is within a distance range less than $\theta$ of the given input string.

Remark: a trie is a special case of FSA

*LIA I&C*

ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

## Pruning criteria: cut-off edit distance

To make it useful in practice $\Rightarrow$ Fast $\Rightarrow$ good pruning

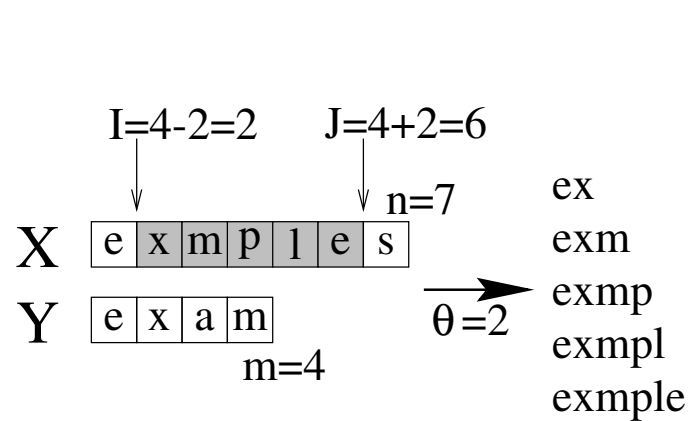☞ cut-off edit distance          [Oflazer 1996]

$$D_c(X_1^n, Y_1^m) = \min_{I(m) \leq i \leq J(m)} D(X_1^i; Y_1^m)$$

$$I(m) = \min(n, \max(1, m - \theta)) \qquad J(m) = \min(n, \max(1, m + \theta))$$

Important property:

$$D_c(X, Y) > \theta \implies \forall Z \ D(X, Y + Z) > \theta$$

LIA
I&C
ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

# Cut-off Edit Distance: example

I=4-2=2    J=4+2=6

n=7

X  | e | x | m | p | l | e | s |

Y  | e | x | a | m |

$\theta=2$

m=4

ex
exm
exmp
exmpl
exmple

|   |   | Y |   |   |   |
|---|---|---|---|---|---|
|   |   | e | x | a | m |
|   | 0 | 1 | 2 | 3 | 4 |
| e | 1 | 0 | 1 | 2 | 3 |
| x | 2 | 1 | 0 | 1 | 2 |
| m | 3 | 2 | 1 | 1 | 1 |
| p | 4 | 3 | 2 | 2 | 2 |
| l | 5 | 4 | 3 | 3 | 3 |
| e | 6 | 5 | 4 | 4 | 4 |
| s | 7 | 6 | 5 | 5 | 5 |

X

$$D_c(X,Y) = \min\{2,1,2,3,4\} = 1$$

## Walk through a FSA within a $\theta$ distance range
## Original algorithm

Input: a string to be corrected $(X)$, a lexicon in the form of a FSA and a maximal error threshold $(\theta)$

Push$(\varepsilon, q_0)$                                                 $\varepsilon$: empty string, $q_0$: inital state

**while** Stack is not empty **do**

  Pop$(Z, p)$

  **for all** $a \in \Sigma$ **do**

    **for all** $q$ such that $\delta(p, a) = q$ **do**
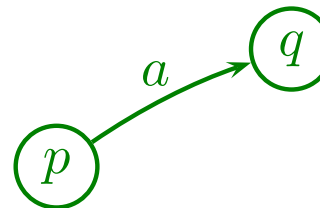
      $Y \leftarrow Z + a$

      **if** $D_c(X, Y) \leq \theta$ **then**

        Push$(Y, q)$

        **if** $q \in F$ **and** $D(X, Y) \leq \theta$        $F$: set of final states **then**

          Add $Y$ to solutions

## Walk through a FSA within a $\theta$ distance range

## Prefix-compatible Depth-first version

Push$(\varepsilon, \varepsilon, q_0)$

**while** Stack is not empty **do**

　Pop$(Z, c, p)$

　$(q, a) = \mathsf{nextAfter}(p, c)$

　**if** $(q, a) \neq \emptyset$ **then**
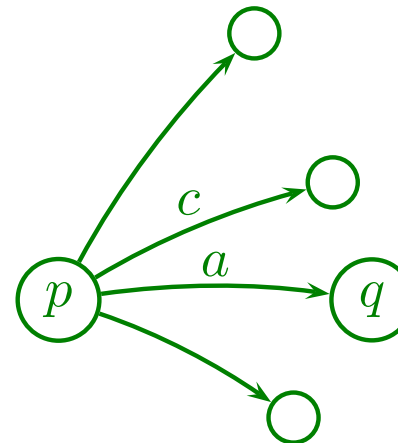
　　Push$(Z, a, p)$

　　$Y \leftarrow Z + a$

　　**if** $D_c(X, Y) \leq \theta$ **then**

　　　Push$(Y, \varepsilon, q)$

　　　**if** $q \in F$ **and** $D(X, Y) \leq \theta$ **then**

　　　　Add $Y$ to solutions



where:

nextAfter$(p, c) = \mathrm{Argmin}_\alpha \{(q, \alpha) \in Q \times \Sigma$ such that $\alpha > c$ and $\delta(p, \alpha) = q\}$

LIA
I&C
ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

# Implementation issues

① original version (simpler) .vs. Prefix-compatible Depth first (faster)

② Efficient computation of $D_c$ with the previously described chart :

☞ recomputation of the last column $(m)$ **only**

(pay attention to the backtrack case! Original version $\rightarrow$ the last *two* columns)

☞ Computation of $D$ and $D_c$ in the **same** loop

③ $Y \leftarrow Z + a$: beware (local copies, pointers etc...).

Similarly, do not naively implement "Push$(Y, q)$".

④ In some languages (especially POO): it could be worth transposing the algorithm: Y (which is changing) for rows and X for columns

LIA
I&C
ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

M. Rajman
J.-C. Chappelier

|   |   | a | b | a | a | b | a |
|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| a | 1 | 0 | 1 | 2 | 3 | 4 | 5 |
| b | 2 | 1 | 0 | 1 | 2 | 3 | 4 |
| a | 3 | 2 | 1 | 0 | 1 | 2 | 3 |
| b | 4 | 3 | 2 | 1 | 1 | 1 | 2 |
| a | 5 | 4 | 3 | 2 | 1 | 1 | 1 |

X=ababa .vs. (aba|bab)* for $\theta = 1$

Solutions ☞ abaaba, ababab, bababa

LIA I&C
ÉCOLE POLYTECHNIQUE FÉDÉRALE DE LAUSANNE

Introduction to Natural Language Processing (CS-431)

M. Rajman
J.-C. Chappelier

22/32

# Contents

- Out of Vocabulary Forms

- Spelling Error Correction

  ✈ Edit distance

  ✈ Spelling error correction with FSA

  ➤ Weighted edit distance

Introduction to Natural Language Processing (CS-431)

*LIA*
*I&C*

ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

M. Rajman
J.-C. Chappelier

23/32

# Limitations?

➥ weighting

Example: diacritics, uppercase

*eleves* → *élèves*          *aloves* → *élèves*

➥ specific transformations

Example: typing errors

*tupe* → *type*          more generally: **de**uit → *fruit*

*us**q**ge* → *usage*

➥ whitespaces

*theothers*→ *the others*          *othe rs*→ *others*

☞ 3 aspects of the **same** problem

Solution: generalization of the edit distance: weighted edit distance

## Weighted Edit Distance

**weighted** transformations such that :

➦ $C(\mathrm{Id}) = 0$

➦ $C(\mathrm{f}) > 0 \qquad \mathrm{f} \neq \mathrm{Id}$

➦ $C(\mathrm{f}^{-1}) = C(\mathrm{f})$

➦ $C(\mathrm{f} \circ \mathrm{g}) = C(\mathrm{f}) + C(\mathrm{g})$

$$D(X; Y) = \min_{\mathrm{f}:Y=\mathrm{f}(X)} C(f)$$

☞ It is actually a distance on $\Sigma^*$

Difference with the preceding distance: $C(\mathrm{f})$ is not necessarily the same (= 1).

# Remarks

❶ Distance on $\Sigma^* \Rightarrow \forall X\, Y, \quad \exists f : Y = f(X)$

True if $\mathrm{Ins}$ and $\mathrm{Del}$ are in the transformation set

❷ non overlapping transformations

i.e. cannot apply a transformation to the result of the previous transformation

Counter-Example: $\mathrm{ba} \overset{\mathrm{Transp}}{\to} \mathrm{ab} \overset{\mathrm{Sub}}{\to} \mathrm{ac}$

## Coherence Constraints

"Semantic Integrity":

- ❑ $C(\mathrm{Del}) + C(\mathrm{Ins}(x)) > C(\mathrm{Sub}(x))$
- ❑ $C(Split) < C(\mathrm{Ins}(x))$ ( $\Rightarrow C(Merge) < C(\mathrm{Del})$)
- ❑ $C(\mathrm{Transp}) < C(\mathrm{Ins}(x)) + C(\mathrm{Del})$

☞ Introduction a new $f$ such that $f = \circ_i f_i$, is useful if and only if

$$C(f) < \sum_i C(f_i)$$

LIA
I&C
ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

Introduction to Natural Language Processing (CS-431)

M. Rajman
J.-C. Chappelier

27/32

## Weighted Edit Distance: computation

$$\text{supp}(f)$$

$$X : \quad \text{xxxx} \boxed{\text{xxx}} \text{xxxx}$$

$$Y : \quad \text{yyy} \boxed{\text{yyy}} \text{yyyyy}$$

$$\text{min1} \qquad \qquad \text{min2}$$

$$\text{min2} = \min_{f} [\ \text{min1}(f) + C(f)\ ]$$

(min1 and min2 are the values stored in the chart)

## Weighted Edit Distance: computation (2)

$$D(X_1^0; Y_1^j) \quad = \quad j \qquad\qquad\qquad\qquad\qquad\qquad \text{initialization}$$

$$D(X_1^i; Y_1^0) \quad = \quad i$$

increasing $C(f)$

$$D(X_1^i; Y_1^j) \quad = \quad D(X_1^{i-1}; Y_1^{j-1}) \qquad\qquad \text{if } X_i = Y_j \text{ (equality)}$$

$$= \quad C(f) + \min\{\min_1(f)\} \qquad \text{for all applicable transfor-mations } f \text{ of the same weight}$$

$$= \quad \ldots \qquad\qquad\qquad\qquad\qquad \text{for all possible weights.}$$

☞ The optimization lies in the grouping of similar cases: same weight and compatible transformations (Example: previously $\mathrm{Transp}$ and $\mathrm{Sub}$ were incompatible because $C(\mathrm{Transp}) < 2\,C(\mathrm{Sub})$. But each of them is compatible with $\mathrm{Del}$ and $\mathrm{Ins}$.)

Note: $\{\min_1(f)\}$ is the set of all the minimal values for all possible $f$ at this point; they shall, of course, already be computed at this point (loop condition)

LIA I&C
ÉCOLE POLYTECHNIQUE FÉDÉRALE DE LAUSANNE

Introduction to Natural Language Processing (CS-431)

M. Rajman
J.-C. Chappelier

29/32

## Example

D(example;exemple)

|   |   | e | x | e | m | p | l | e |
|---|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| e | 1 | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| x | 2 | 1 | 0 | 1 | 2 | 3 | 4 | 5 |
| a | 3 | 2 | 1 | 1 | 2 | 3 | 4 | 5 |
| m | 4 | 3 | 2 | 2 | 1 | 2 | 3 | 4 |
| p | 5 | 4 | 3 | 3 | 2 | 1 | 2 | 3 |
| l | 6 | 5 | 4 | 4 | 3 | 2 | 1 | 2 |
| e | 7 | 6 | 5 | 4 | 4 | 3 | 2 | 1 |

D(exémple;exemple)

|   |   | e | x | e | m | p | l | e |
|---|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| e | 1 | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| x | 2 | 1 | 0 | 1 | 2 | 3 | 4 | 5 |
| é | 3 | 2 | 1 | 0.1 | 1.1 | 2.1 | 3.1 | 4.1 |
| m | 4 | 3 | 2 | 1.1 | 0.1 | 1.1 | 2.1 | 3.1 |
| p | 5 | 4 | 3 | 2.1 | 1.1 | 0.1 | 1.1 | 2.1 |
| l | 6 | 5 | 4 | 3.1 | 2.1 | 1.1 | 0.1 | 1.1 |
| e | 7 | 6 | 5 | 4 | 3.1 | 2.1 | 1.1 | 0.1 |

$$C(\text{é}\leftrightarrow\text{e})=0.1$$

## Keypoints

➡ One has to handle out of vocabulary forms

➡ Edit (Levenshtein) distance, weighted edit distance

➡ Spelling error correction with FSA

*LIA*
*I&C*
ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

Introduction to Natural Language Processing (CS-431)

M. Rajman
J.-C. Chappelier

31/32

# References

K. Oflazer, *Error-tolerant Finite State Recognition with Applications to Morphological Analysis and Spelling Correction*, Computational Linguistics, Volume 22, Number 1, 1996.

Section 8.2 in M. Rajman editor, "Speech and Language Engineering", EPFL Press, 2006.

Sections 3.10 and 3.11 in D. Jurafsky and J. H. Martin, "Speech and Language Processing", Prentice Hall, 2008 (2nd edition).

Section 3.3 in C. D. Manning, P. Raghavan and H. Schütze, "Introduction to Information Retrieval", Cambridge University Press. 2008

LIA
I&C
ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

Introduction to Natural Language Processing (CS-431)

M. Rajman
J.-C. Chappelier

32/32