

Edit distances, string alignments and dynamic programming

Marc-Olivier Buob, Maxime Raynal (Nokia Bell Labs/ AAAID)

LINCS, Network theory, February 3th, 2021

The human problem

Find the differences between two input strings

pg_original.txt	pg_modif.txt
1Voici un texte original :	1Voici le texte modifié :
2	2
3Solsbury hill	3Solsbury hill
4	4
5Climbing up on solsbury hill	5Climbing up on solsbury hill
6I could see the city light	6I could see the city light
7Wind was blowing, time stood still	7Wind was blowing, time stood still
8Eagle flew out of the night	8Eagle flew out of the night
9	9
10He was something to observe	10So I went from day to day
11Came in close, I heard a voice	11Tho my life was in a rut
12Standing stretching every nerve	12till I thought of what Id say
13I had to listen had no choice	13Which connection I should cut
14	14
15I did not believe the information	15He was something to observe
16Just had to trust imagination	16Came in close, I heard a voice
17My heart was going boom boom, boom	17Standing stretching every nerve
18Son, he said, grab your things, Ive come to take you home.	18I had to listen had no choice
19	19
20To keeping silence I resigned	20I did not believe the information
21My friends would think I was a nut	21Just had to trust imagination
22Turning water into wine	22My heart was going boom boom boom
23Open doors would soon be shut	23Son, he said, grab your things, I've come to take you home.
24	24
25So I went from day to day	25To keeping silence I resigned
26Tho my life was in a rut	26My friends would think I was a nut
27till I thought of what Id say	27Turning water into wine
28Which connection I should cut	28Open doors would soon be shut
29	29
30I was feeling part of the scenery	30I was feeling part of the scenery
31I walked right out of the machinery	31I walked right out of the machinery
32My heart was going boom boom boom	32My heart was going boom boom boom
33Hey, he said, grab your things, Ive come to take you home.	33Hey, he said, grab your things, Ive come to take you home.

Use cases

Wide range of applications

- **Computer science**
 - File comparison (diff, git, ...)
 - Approximate string matching
 - spell checkers,
 - fuzzy string search,
 - fraud detection...
 - Optical character recognition
- **Bioinformatic**
 - Nucleic acid sequence and protein alignment
- **Linguistic**
 - Distance between two languages



Edit operations

Used to "count" the number of differences

- **Insertion**

- ABCDE vs ABXCDE

- **Deletion**

- ABXCDE vs ABCDE

- **Substitution**

- ABCDE vs ABXDE

- **Transposition**

- Cyclic permutation
- 1234 vs 2413
- ABCD vs CBAD

pg_original.txt	pg_modif.txt
1 Voici un texte original :	1 Voici le texte modifié :
2	2
3 Solsbury hill	3 Solsbury hill
4	4
5 Climbing up on solsbury hill	5 Climbing up on solsbury hill
6 I could see the city light	6 I could see the city light
7 Wind was blowing, time stood still	7 Wind was blowing, time stood still
8 Eagle flew out of the night	8 Eagle flew out of the night
9	9
10 He was something to observe	10 So I went from day to day
11 Came in close, I heard a voice	11 Tho my life was in a rut
12 Standing stretching every nerve	12 till I thought of what Id say
13 I had to listen had no choice	13 Which connection I should cut
14	14
15 I did not believe the information	15 He was something to observe
16 Just had to trust imagination	16 Came in close, I heard a voice
17 My heart was going boom boom, boom	17 Standing stretching every nerve
18 Son, he said, grab your things, Ive come to take you home.	18 I had to listen had no choice
19	19
20 To keeping silence I resigned	20 I did not believe the information
21 My friends would think I was a nut	21 Just had to trust imagination
22 Turning water into wine	22 My heart was going boom boom boom
23 Open doors would soon be shut	23 Son, he said, grab your things, I've come to take you home.
24	24
25 So I went from day to day	25 To keeping silence I resigned
26 Tho my life was in a rut	26 My friends would think I was a nut
27 till I thought of what Id say	27 Turning water into wine
28 Which connection I should cut	28 Open doors would soon be shut
29	29
30 I was feeling part of the scenery	30 I was feeling part of the scenery
31 I walked right out of the machinery	31 I walked right out of the machinery
32 My heart was going boom boom boom	32 My heart was going boom boom boom
33 Hey, he said, grab your things, Ive come to take you home.	33 Hey, he said, grab your things, Ive come to take you home.

Popular edit distances

Support/count a subset of edit operations

Distance	Insertion	Deletion	Substitution	Transposition
Levenstein	✓	✓	✓	
Damerau-Levenstein	✓	✓	✓	... of consecutive char pairs
LCS	✓	✓		
Hamming			✓	
Jaro-Winkler				... of "closed" chars

- As Hamming distance only supports substitution, it can only compare two strings of same length.
- To compute distances involving insertion and deletion, we typically rely on **dynamic programming**.

Longest common subsequence (LCS)

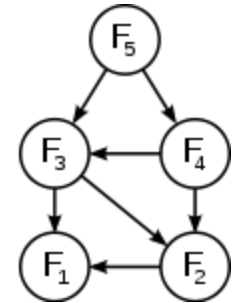
This part presents:

- the LCS problem,
- the underlying model graph (called edit graph)
- the resulting dynamic programming model.
- the main algorithms related to LCS problem.

Dynamic programming (Richard Bellman, 1950)

An optimization method and a computer programming method

- **Key idea:** break down a complex problem into a simpler subproblems in a **recursive** manner.
- **Scope:** it applies on any problem having an **optimal substructure** and **overlapping subproblems**.
 - *Optimal substructure* means that the solution can be obtained by the combination of optimal solutions to its sub-problems.
 - *Overlapping* sub-problems means that sub-problems dependencies form a DAG, by contrast to divide and conquer (D&C) where this graph is a tree.
- **DP:** Fibonacci series ($F_i = F_{i-1} + F_{i-2}$), Dijkstra algorithm, LCS.
- **D&C:** merge sort, quick sort.



The LCS (Longest Common Subsequence) problem, Maier, 1978

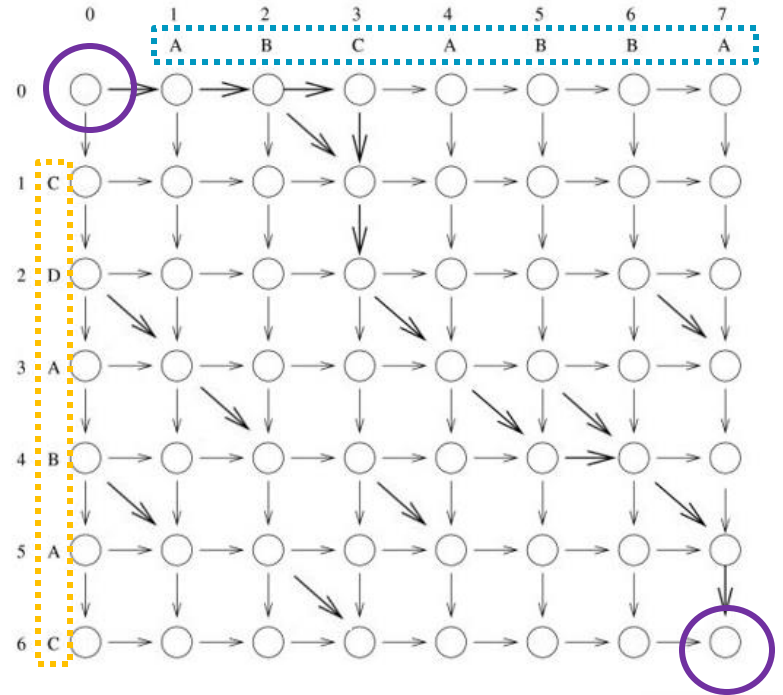
Definitions

- Consider an arbitrary word w . Any sub-word obtained by selecting a subset of characters with distinct index and sorted by increasing index is said to be a **subsequence**.
 - *Example:* **ACEF** is a subsequence of **ABCDEFGF**
- Consider two strings X and Y . If s is a subsequence of X and Y , it is said to be a **common subsequence** of X and Y
 - *Example:* **ABC** is a common subsequence of **ABCABBA** and **CDABAC**
- The LCS problem aims at finding a **longest common subsequence** of two words.
 - *Example:* **CABA** is the LCSs of **ABCABBA** and **CDABAC**.
 - In general, two words may have several LCSs.
- How to determine them?

Edit graph

Model

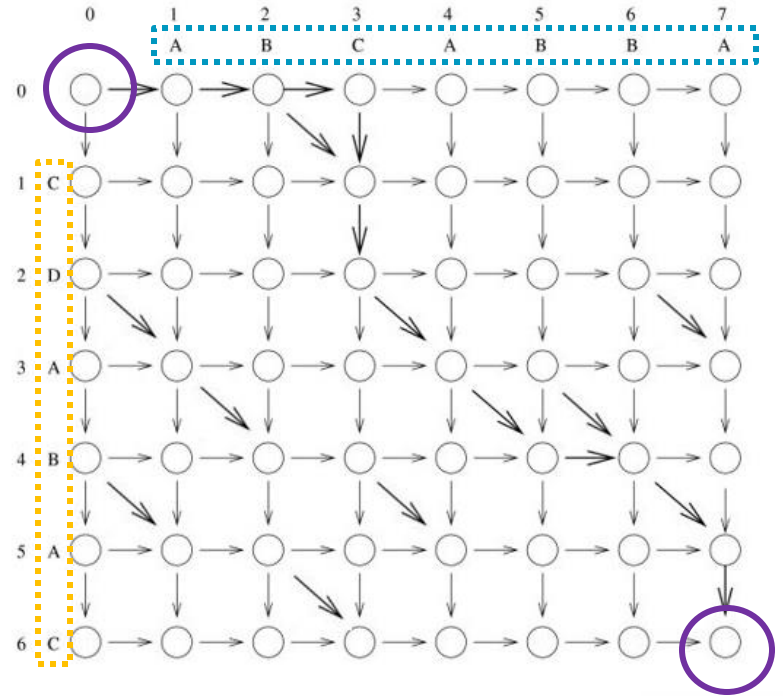
- Consider two words:
 - X, of length m (e.g. **A**BCABBA)
 - Y, of length n (e.g. C**D**ABAC)
- The **edit graph** is defined as follows:
 - **Vertices:**
 - Each $(i, j) \in \{0 \dots m\} \times \{0 \dots n\}$
 - **Arcs:**
 - **Horizontal:** from (i, j) to $(i, j+1)$ (insertion)
 - **Vertical:** from (i, j) to $(i+1, j)$ (deletion)
 - **Diagonal:** from (i, j) to $(i+1, j+1)$ if and only if $X[i] == Y[j]$ (match)



Edit graph

Find optimal string alignments

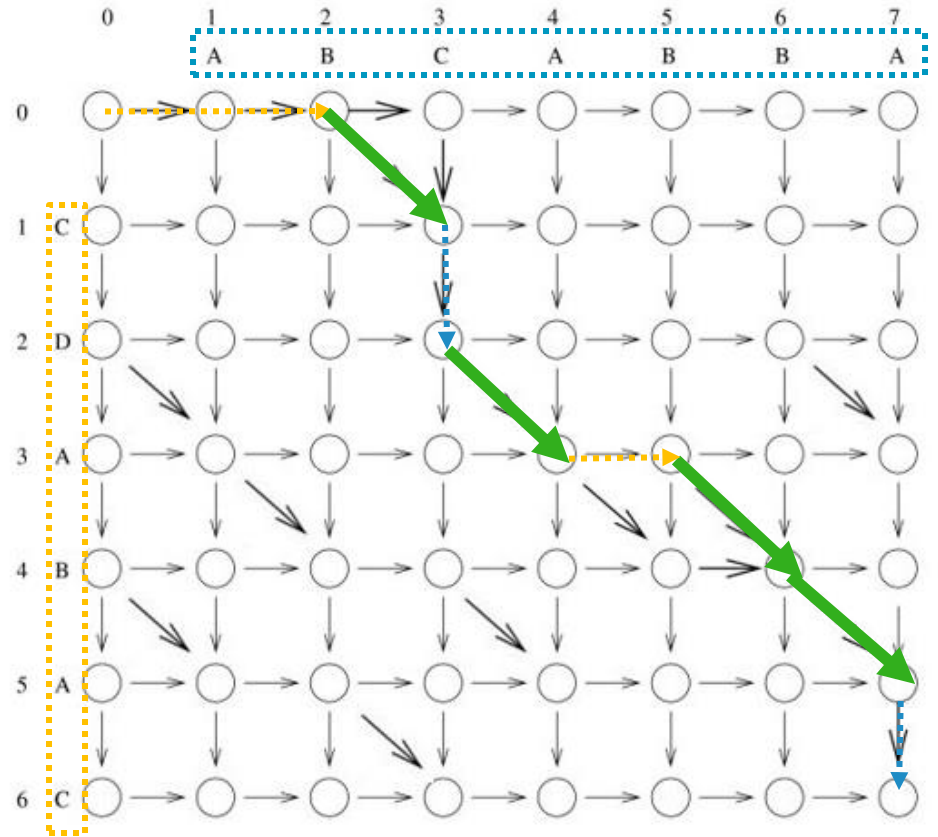
- The LCS must entirely consider X and Y \Rightarrow Path from $(0, 0)$ to (m, n)
- Horizontal and diagonal arc = edit operation
- Path = set of editions to move from X to Y (and conversely) called **alignments**
- The LCS maximizes the #matching characters \Rightarrow The path maximizes the #diagonal edges.



Edit graph

LCS extraction

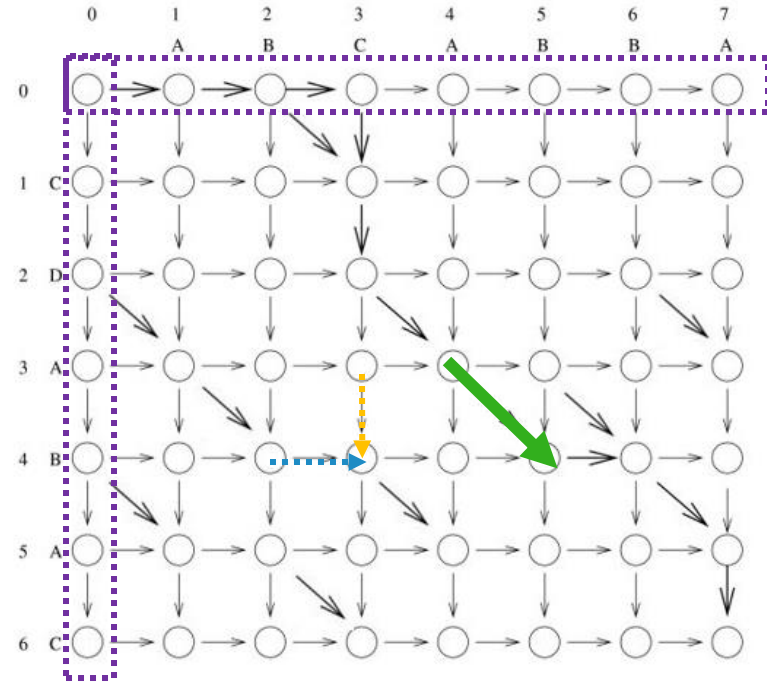
- Paths maximizing #diagonal arcs reveal LCS of $X=ABCABBA$ and $Y=CDABAC$.
- Here: $ABCDABBAC$.
- To transform X into Y:
 - 3 insertions
 - 2 deletions
 - 4 matches (LCS = CABA)
- $LCS(X, Y) = 4$



LCS and dynamic programming

Edit graph ~ DP model

- Initialization: $i = 0$ or $j = 0$
 - Empty LCS.
- Recursion (from (n, n')): $i > 0$ and $j > 0$
 - Prefer **diagonal** arc (if any):
 - Intuition: It's always better to go through a diagonal arc
 - Score : +1
 - Otherwise: **horizontal** and **vertical** arcs.
 - Intuition: Best effort fallback
 - Score : +0



$$LCS(X_i, Y_j) = \begin{cases} \emptyset & \text{if } i = 0 \text{ or } j = 0 \\ LCS(X_{i-1}, Y_{j-1}) \hat{x}_i & \text{if } i, j > 0 \text{ and } x_i = y_j \\ \max\{LCS(X_i, Y_{j-1}), LCS(X_{i-1}, Y_j)\} & \text{if } i, j > 0 \text{ and } x_i \neq y_j. \end{cases}$$

Concatenation

LCS and dynamic programming

Algorithm

- Compute recursively B and C where:
 - C[i, j] stores the **score** obtained for LCS(X[:i], Y[:j])
 - B[i, j] stores an optimal **predecessor** of (i, j) predecessor chosen to obtain C[i, j]
- **Optimization:** you could only store the two last rows of C.

```
function LCSLength(X[1..m], Y[1..n])
    C = array(0..m, 0..n)
    for i := 0..m C[i,0] = 0
    for j := 0..n C[0,j] = 0
    for i := 1..m
        for j := 1..n
            if X[i] = Y[j]
                C[i,j] := C[i-1,j-1] + 1
            else
                C[i,j] := max(C[i,j-1], C[i-1,j])
    return C[m,n]
```

$$LCS(X_i, Y_j) = \begin{cases} \emptyset & \text{if } i = 0 \text{ or } j = 0 \\ LCS(X_{i-1}, Y_{j-1}) \hat{x}_i & \text{if } i, j > 0 \text{ and } x_i = y_j \\ \max\{LCS(X_i, Y_{j-1}), LCS(X_{i-1}, Y_j)\} & \text{if } i, j > 0 \text{ and } x_i \neq y_j. \end{cases}$$

Needleman Wunsch, 1970

Custom score function

- **Problem:** In LCS: match: +1 ; mismatch: +0. What if two characters are **similar**?
- **Key idea:** build diagonal arcs (from (i, j) to (i+1, j+1)), and weight them using a **score function** s:
 - **Match:** $s(a, b) > 0$ if a and b are **equal** or **similar**
 - **Mismatch:** $s(a, b) = D$, where D is a (negative) constant.
- **Trick:** update DP model as follows:
 - **Initialization:**
 - $C[i, 0]$ are initialized to $-D \cdot i$
 - $C[0, j]$ are initialized to $-D \cdot j$
 - **Recursion:**
 - $C[i, j] = \max(C[i-1, j-1] + s(a, b), C[i, j-1] + D, C[i-1, j] + D)$






Needleman-Wunsch

match = 1 mismatch = -1 gap = -1

		G	C	A	T	G	C	U
	0	-1	-2	-3	-4	-5	-6	-7
G	-1	1	0	-1	-2	-3	-4	-5
A	-2	0	0	1	0	-1	-2	-3
T	-3	-1	-1	0	2	1	0	-1
T	-4	-2	-2	-1	1	1	0	-1
A	-5	-3	-3	-1	0	0	0	-1
C	-6	-4	-2	-2	-1	-1	1	0
A	-7	-5	-3	-1	-2	-2	0	0

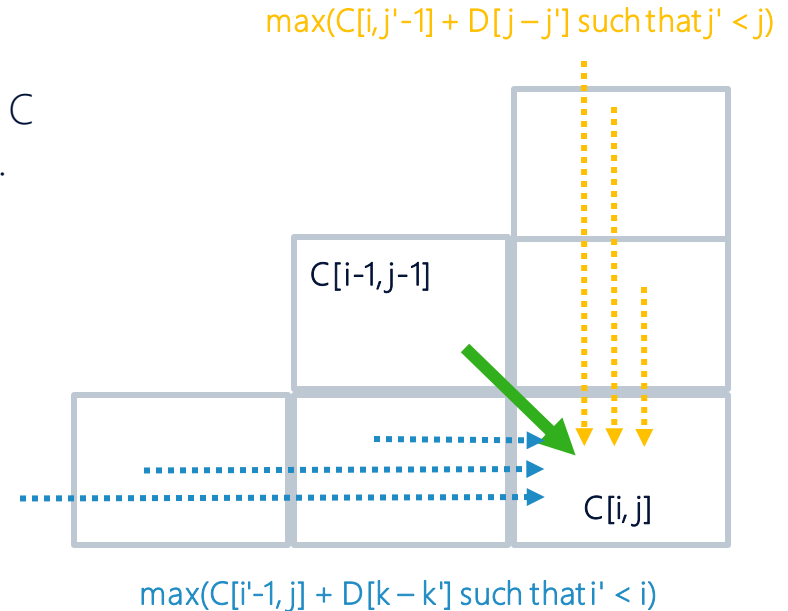
Example taken from [\[wikipedia\]](#)

- $D = -1$ if indel 
- If $a = b$: $s(a, b) = 1$ 
- If $a \sim b$: $s(a, b) = -1$ 
- We only pay -1 (instead of $2 \cdot D$) when aligning "similar" characters.

Smith-Waterman algorithm, 1981

Local alignment

- **Problem:** how to find the best **local alignment** (between any pair of vertices of the edit graph)
- **Trick:** uses an extended neighborhood to compute C and define a penalty depending on the gap length.
- Initialization:
 - $C[i, 0]$ and $C[0, j]$ are initialized to 0.
- Recursion:
 - $C[i, j] = \max(\begin{array}{l} C[i-1, j-1] + s(a, b), \\ \max(C[i, j'-1] + D[j-j'] \text{ such that } j' < j), \\ \max(C[i'-1, j] + D[k-k'] \text{ such that } i' < i) \end{array})$



Back to the other edit distances

This part presents:

- the Hamming distance
- the Levenshtein distance
- the Damerau-Levenshtein distance
- I skip the Jaro-Winkler distance, see wikipedia if you want further details...

Hamming distance, 1950

Count matching characters at fixed indices

- Consider two words X and Y of length m and n, such that m = n. The **Hamming distance** is defined by:

$$\text{hamming}(a, b) = \sum \mathbf{0}(X[i], Y[i]) \text{ where } \mathbf{0}(a, b) = 0 \text{ if } a = b, 1 \text{ otherwise}$$

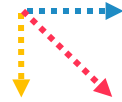
- **Interpretation:** The Hamming distance counts the #mismatching characters).
- *Examples:*
 - "karolin" and "kathrin" is 3.
 - "karolin" and "kerstin" is 3.
- The Hamming distance can be computed in $O(n)$ with a footprint in $O(1)$. It's significantly cheaper than computing $|a| - |\text{LCS}(X, Y)|$ which is done in $O(m.n)$ with a footprint in $O(n)$

Levenshtein distance (Владимир Иосифович Левенштéйн, 1956)

Count #operations to transform w into w' (insertion, deletion, copy)

- Consider two words a and b . The Levenstein distance **lev** is defined by:

$$\text{lev}(a, b) = \begin{cases} |a| & \text{if } |b| = 0, \\ |b| & \text{if } |a| = 0, \\ \text{lev}(\text{tail}(a), \text{tail}(b)) & \text{if } a[0] = b[0] \\ 1 + \min \begin{cases} \text{lev}(\text{tail}(a), b) \\ \text{lev}(a, \text{tail}(b)) \\ \text{lev}(\text{tail}(a), \text{tail}(b)) \end{cases} & \text{otherwise.} \end{cases}$$



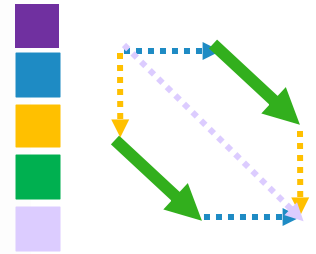
- ... where **tail** be the primitive returning the suffix starting from index 1 of a string w
 - Example:* $\text{tail}(\text{"abcde"}) = \text{"bcde"}$
- It computes the shortest path length from $(0, 0)$ to (m, n) in the edit graph by minimizing **#insertion** + **#deletion** + **#substitution**, while LCS maximizes **#matches**.
- lev** can be computed using [Wagner Fisher algorithm](#) (which is slight modification of [Needleman Wunsh algorithm](#))

Damerau-Levenshtein distance, 1964

Same as Levenshtein distance, with some transposition operations

- Consider two words a and b . The **Damerau-Levenshtein distance** of two words a and b is recursively defined by:

$$d_{a,b}(i, j) = \min \begin{cases} 0 & \text{if } i = j = 0 \\ d_{a,b}(i - 1, j) + 1 & \text{if } i > 0 \\ d_{a,b}(i, j - 1) + 1 & \text{if } j > 0 \\ d_{a,b}(i - 1, j - 1) + 1_{(a_i \neq b_j)} & \text{if } i, j > 0 \\ d_{a,b}(i - 2, j - 2) + 1 & \text{if } i, j > 1 \text{ and } a[i] = b[j - 1] \text{ and } a[i - 1] = b[j] \end{cases}$$



- Compared to Levenshtein distance, **swapped characters** are rewarded by slightly modifying the neighborhood definition.
 - Intuitively, if $a[i]a[i-1] = b[j-1]b[j]$, there is a corresponding arc that only costs 1 instead, and which is cheaper than the indel paths of cost 2.

NOKIA

Wrap-up

Complexity, footprint and properties

Distances are not necessarily metrics!

Distance	Metric	Compl.	Footp.
Levenshtein	✓	$O(m.n)$	$O(m.n)$
Damerau-Levenshtein	!Δ	$O(m.n)$	$O(\min(m,n))$
LCS	✓	$O(m.n)$	$O(n)$
Hamming	✓	$O(n)$	$O(1)$
Jaro-Winkler	!Δ, !id		

- d is a **metric** iff the following properties hold:
 - **Symmetry**: $d(a,b) = d(b,a)$
 - **Identity**: $d(a,b) = 0 \Leftrightarrow a = b$
 - **Triangle inequality**: $d(a,c) \leq d(a,b) + d(b,c)$
- Check requirements to compute correct results!
 - Some algorithms require metric (e.g. Dijkstra).
 - Some algorithms only require quasi-metric (e.g. some clustering algorithms).

Conclusion

Edit distances and string alignments

- Edit distances are **widely used** to check if two input strings are "similar"
 - The edit distance results from **string alignments**, that are ruled by a set of **edit operations** (insertion, deletion, substitution and transposition) and a **cost function**.
 - Edit distances are not always **metrics**.
- Edit distances are typically computed using **dynamic programming**.
 - Dynamic programming is a technique applies on any problem having an **optimal substructure** and **overlapping subproblems**. As these dependencies form a DAG (not a tree) this is DP (not D&C).
- When computing edit distances, the DP model is closely related to the **edit graph**.
 - DP model ~ edit graph (topology + weights)
 - DP solution ~ best path ~ best alignment.
 - Edit distance ~ best path length ~ best alignment score

NOKIA