

14.2

Edit Distance and Sequence Alignment

14.2.1

Problem definition and background

Spell Checking Problem

Given a string “exponen” that is not in the dictionary, how should a spell checker suggest a nearby string?

What does nearness mean?

Question: Given two strings $x_1x_2 \dots x_n$ and $y_1y_2 \dots y_m$ what is a distance between them?

Edit Distance: minimum number of “edits” to transform x into y .

Spell Checking Problem

Given a string “exponen” that is not in the dictionary, how should a spell checker suggest a nearby string?

What does nearness mean?

Question: Given two strings $x_1x_2 \dots x_n$ and $y_1y_2 \dots y_m$ what is a distance between them?

Edit Distance: minimum number of “edits” to transform x into y .

Spell Checking Problem

Given a string “exponen” that is not in the dictionary, how should a spell checker suggest a nearby string?

What does nearness mean?

Question: Given two strings $x_1x_2 \dots x_n$ and $y_1y_2 \dots y_m$ what is a distance between them?

Edit Distance: minimum number of “edits” to transform x into y .

Edit Distance

Definition 14.1.

Edit distance between two words X and Y is the number of letter insertions, letter deletions and letter substitutions required to obtain Y from X .

Example 14.2.

The edit distance between FOOD and MONEY is at most **4**:

FOOD \rightarrow MOOD \rightarrow MONOD \rightarrow MONED \rightarrow MONEY

Edit Distance: Alternate View

Alignment

Place words one on top of the other, with gaps in the first word indicating insertions, and gaps in the second word indicating deletions.

F	O	O		D
M	O	N	E	Y

Formally, an **alignment** is a set M of pairs (i, j) such that each index appears at most once, and there is no “crossing”: $i < i'$ and i is matched to j implies i' is matched to $j' > j$. In the above example, this is $M = \{(1, 1), (2, 2), (3, 3), (4, 5)\}$. Cost of an alignment is the number of mismatched columns plus number of unmatched indices in both strings.

Edit Distance: Alternate View

Alignment

Place words one on top of the other, with gaps in the first word indicating insertions, and gaps in the second word indicating deletions.

F	O	O		D
M	O	N	E	Y

Formally, an **alignment** is a set M of pairs (i, j) such that each index appears at most once, and there is no “crossing”: $i < i'$ and i is matched to j implies i' is matched to $j' > j$. In the above example, this is $M = \{(1, 1), (2, 2), (3, 3), (4, 5)\}$. Cost of an alignment is the number of mismatched columns plus number of unmatched indices in both strings.

Edit Distance: Alternate View

Alignment

Place words one on top of the other, with gaps in the first word indicating insertions, and gaps in the second word indicating deletions.

F	O	O		D
M	O	N	E	Y

Formally, an **alignment** is a set M of pairs (i, j) such that each index appears at most once, and there is no “crossing”: $i < i'$ and i is matched to j implies i' is matched to $j' > j$. In the above example, this is $M = \{(1, 1), (2, 2), (3, 3), (4, 5)\}$. Cost of an alignment is the number of mismatched columns plus number of unmatched indices in both strings.

Edit Distance Problem

Problem

Given two words, find the edit distance between them, i.e., an alignment of smallest cost.

Applications

- ① Spell-checkers and Dictionaries
- ② Unix `diff`
- ③ DNA sequence alignment ... but, we need a new metric

Similarity Metric

Definition 14.3.

For two strings X and Y , the cost of alignment M is

- 1 [Gap penalty] For each gap in the alignment, we incur a cost δ .
- 2 [Mismatch cost] For each pair p and q that have been matched in M , we incur cost α_{pq} ; typically $\alpha_{pp} = 0$.

Edit distance is special case when $\delta = \alpha_{pq} = 1$.

Similarity Metric

Definition 14.3.

For two strings X and Y , the cost of alignment M is

- 1 [Gap penalty] For each gap in the alignment, we incur a cost δ .
- 2 [Mismatch cost] For each pair p and q that have been matched in M , we incur cost α_{pq} ; typically $\alpha_{pp} = 0$.

Edit distance is special case when $\delta = \alpha_{pq} = 1$.

THE END

...

(for now)