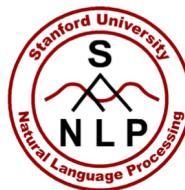


Minimum Edit Distance

Definition of Minimum
Edit Distance



How similar are two strings?

- Spell correction
 - The user typed “graffe”
 Which is closest?
 - graf
 - graft
 - grail
 - giraffe
- Computational Biology
 - Align two sequences of nucleotides


```
AGGCTATCACCTGACCTCCAGGCCGATGCC  
TAGCTATCACGACCAGCGGTGATTTGCCCGAC
```
 - Resulting alignment:


```
—AGGCTATCACCTGACCTTCAGGCCGA--TGCCC---  
TAG—CTATCAC--GACCAGC--GGTCGAATTGCCCGAC
```
- Also for Machine Translation, Information Extraction, Speech Recognition



Edit Distance

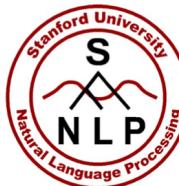
- The minimum edit distance between two strings
- Is the minimum number of editing operations
 - Insertion
 - Deletion
 - Substitution
- Needed to transform one into the other



Minimum Edit Distance

- Two strings and their **alignment**:

I	N	T	E	*	N	T	I	O	N
*	E	X	E	C	U	T	I	O	N



Minimum Edit Distance

I N T E * N T I O N

| | | | | | | | |

* E X E C U T I O N

d s s i s

- If each operation has cost of 1
 - Distance between these is 5
- If substitutions cost 2 (Levenshtein)
 - Distance between them is 8



Alignment in Computational Biology

- Given a sequence of bases

AGGCTATCACCTGACCTCCAGGCCGATGCC
TAGCTATCACGACCGCGGTGATTGCCCGAC

- An alignment:

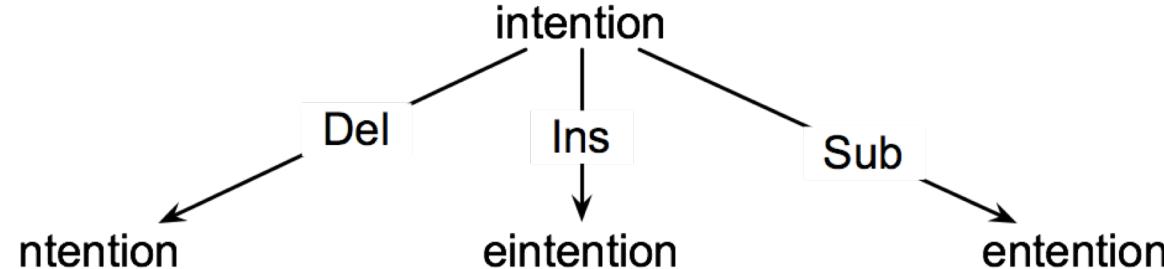
-**A**GG**C**TAT**C**AC**C**T**G**ACC**T**CCA**GG**C**CG**A--TG**CCC**--
TAG-CTAT**C**AC--**G**ACC**G**C--GGT**CG**A**TT**TGCC**C**GAC

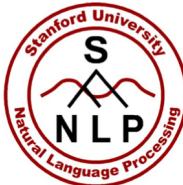
- Given two sequences, align each letter to a letter or gap



How to find the Min Edit Distance?

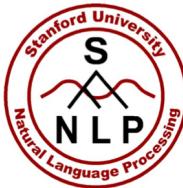
- Searching for a path (sequence of edits) from the start string to the final string:
 - **Initial state:** the word we're transforming
 - **Operators:** insert, delete, substitute
 - **Goal state:** the word we're trying to get to
 - **Path cost:** what we want to minimize: the number of edits





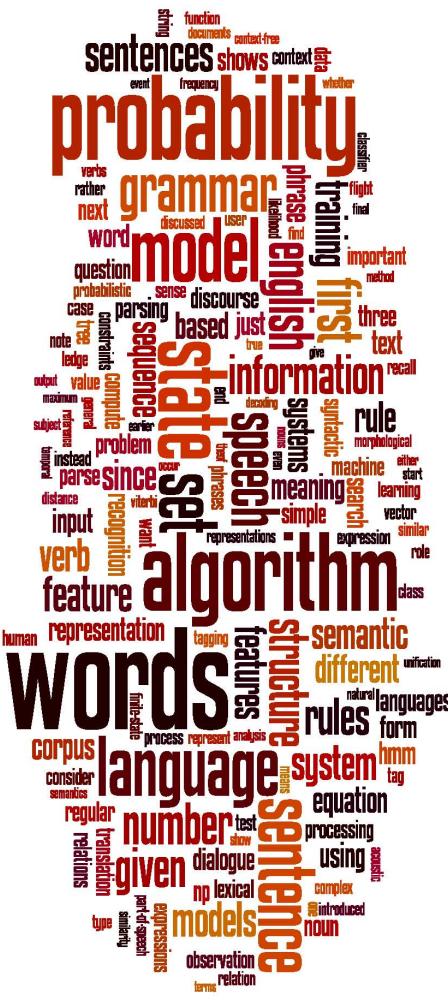
Minimum Edit as Search

- But the space of all edit sequences is huge!
 - We can't afford to navigate naïvely
 - Lots of distinct paths wind up at the same state.
 - We don't have to keep track of all of them
 - Just the shortest path to each of those revisited states.



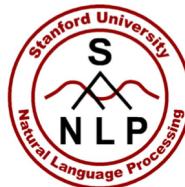
Defining Min Edit Distance

- For two strings
 - X of length n
 - Y of length m
- We define $D(i,j)$
 - the edit distance between $X[1..i]$ and $Y[1..j]$
 - i.e., the first i characters of X and the first j characters of Y
 - The edit distance between X and Y is thus $D(n,m)$



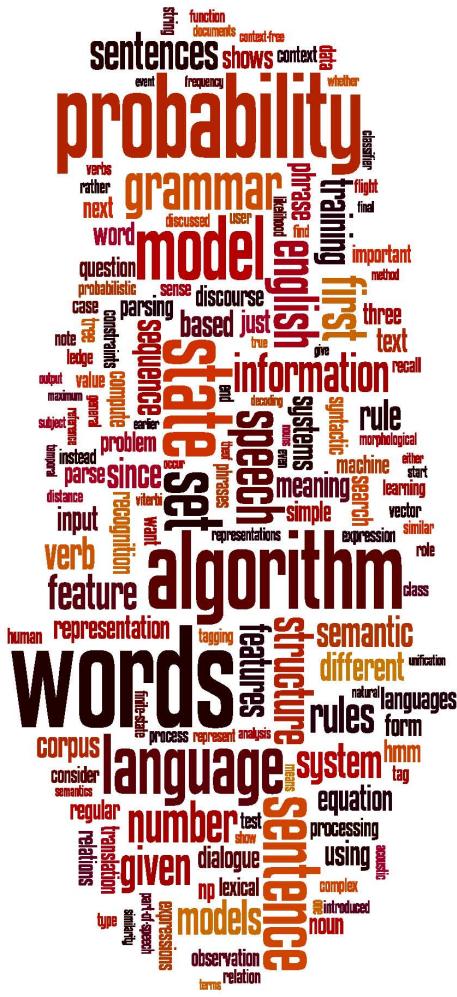
Minimum Edit Distance

Computing Minimum
Edit Distance



Dynamic Programming for Minimum Edit Distance

- **Dynamic programming:** A tabular computation of $D(n,m)$
- Solving problems by combining solutions to subproblems.
- Bottom-up
 - We compute $D(i,j)$ for small i,j
 - And compute larger $D(i,j)$ based on previously computed smaller values
 - i.e., compute $D(i,j)$ for all i ($0 < i < n$) and j ($0 < j < m$)



Minimum Edit Distance

Backtrace for
Computing Alignments

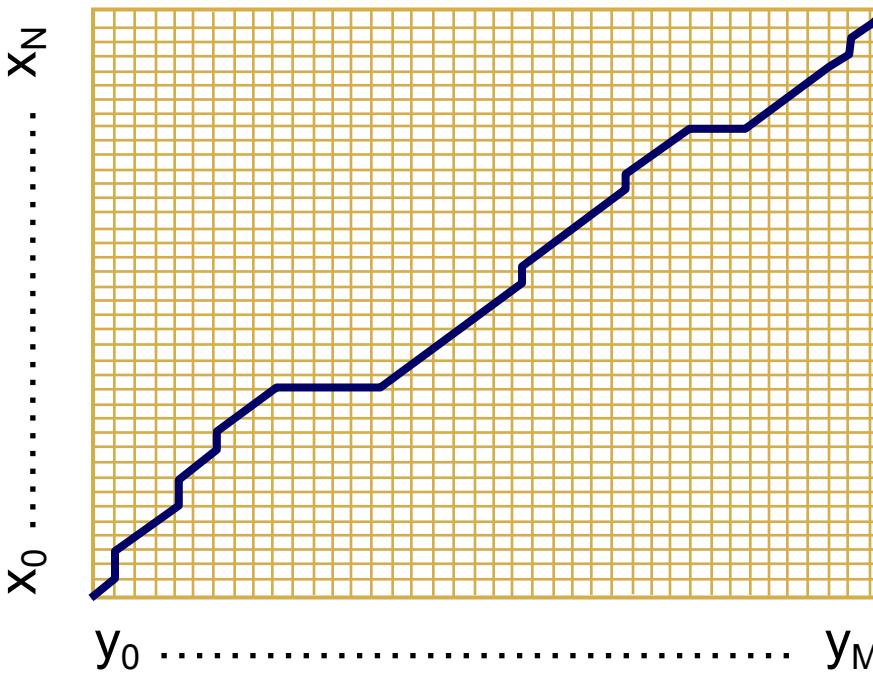


Computing alignments

- Edit distance isn't sufficient
 - We often need to **align** each character of the two strings to each other
- We do this by keeping a “backtrace”
- Every time we enter a cell, remember where we came from
- When we reach the end,
 - Trace back the path from the upper right corner to read off the alignment



The Distance Matrix



Every non-decreasing path
from $(0,0)$ to (M, N)

corresponds to
an alignment
of the two sequences

An optimal alignment is composed
of optimal subalignments



Result of Backtrace

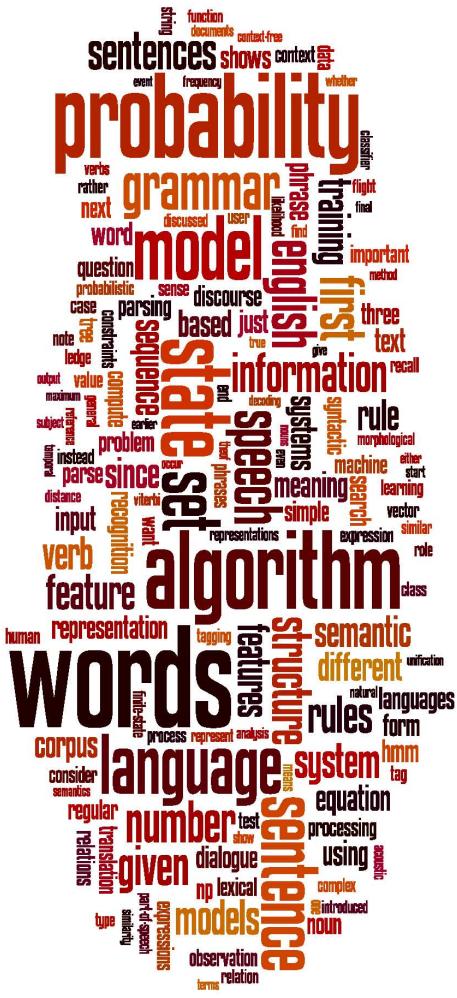
- Two strings and their **alignment**:

I	N	T	E	*	N	T	I	O	N
*	E	X	E	C	U	T	I	O	N



Performance

- Time:
 $O(nm)$
- Space:
 $O(nm)$
- Backtrace
 $O(n+m)$



Minimum Edit Distance

Weighted Minimum Edit Distance



Weighted Edit Distance

- Why would we add weights to the computation?
 - Spell Correction: some letters are more likely to be mistyped than others
 - Biology: certain kinds of deletions or insertions are more likely than others





Weighted Min Edit Distance

- Initialization:

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

$$D(i, 0) = D(i-1, 0) + \text{del}[x(i)]; \quad 1 < i \leq N$$

$$D(0, j) = D(0, j-1) + \text{ins}[y(j)]; \quad 1 < j \leq M$$

- Recurrence Relation:

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

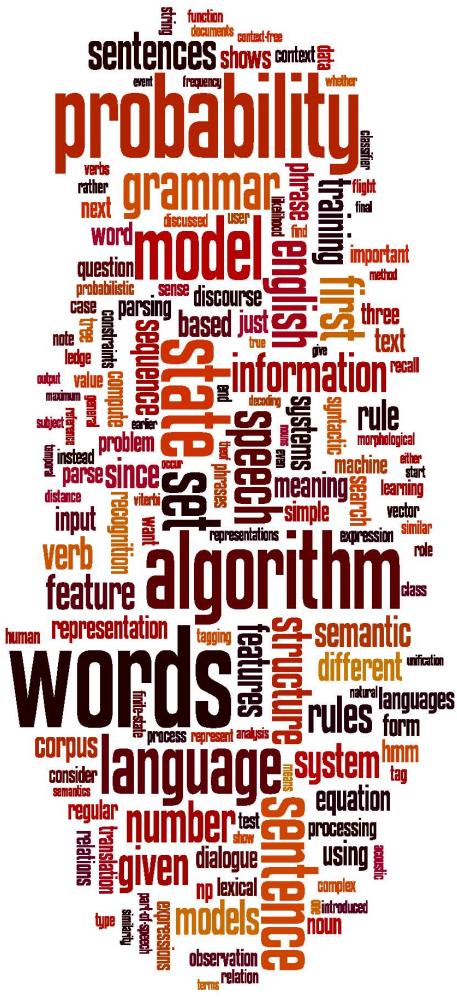
- Termination:

$D(N, M)$ is distance



- Show the edit distance between “giraffe” and “coffee” assuming insertion and deletion cost is 1 while substitution cost is 2. Fill up the table given below. Suggest one possible alignment along with the operation applied at each position (i for insertion/ d for deletion/ s for substitution) for achieving optimized edit distance.

e						
f						
f						
a						
r						
i						
g						
		c	o	f	f	e
						e



Minimum Edit Distance

Weighted Minimum Edit Distance