

# Recap of the last lecture

# Regular Expression (RE)

- A standard notation of characterizing a text sequence
- How can we search for any of the following:
  - woodchuck
  - woodchucks
  - Woodchuck
  - Woodchucks



- RE search requires a pattern and a **corpus** of texts to search through.

# Morphology: Definition

The study of words, how they are formed, and their relationship to other words in the same language.

# The Porter Stemmer (Porter, 1980)

- A simple rule-based algorithm for stemming
- An example of a HEURISTIC method
- Based on rules like:
  - ATIONAL -> ATE (e.g., *relational* -> *relate*)
- The algorithm consists of seven sets of rules, applied in order

# Spelling Error: Minimum Edit Distance



# How similar are two strings?

- Spell correction

- The user typed “Appl”

Which is closest?

- App
    - Appeal
    - Apple

- Computational Biology

- Align two sequences of nucleotides

```
AGGCTATCACCTGACCTCCAGGCCGATGCCC  
TAGCTATCACGACCGCGGGTCGATTTGCCCGAC
```

- Resulting alignment:

```
-AGGCTATCACCTGACCTCCAGGCCGA--TGCCC--  
TAG-CTATCAC--GACCGC--GGTCGATTTGCCCGAC
```

- 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 (**I**)
  - Deletion (**D**)
  - Substitution (**S**)
- Need to transform one into the other

# Minimum Edit Distance

- Two strings and their **alignment**: TRIAL vs ZEIL

T	R	I	A	L
Z	E	I	*	L
s	s		d	

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 3
- If substitutions cost 2 (Levenshtein)
  - Distance between them is 5



# Other uses of Edit Distance in NLP

- Evaluating Machine Translation and speech recognition

Spokesman confirms      senior government adviser was shot

Spokesman said      the senior      adviser was shot dead

S

I

D

I

- **Named Entity Extraction and Entity Coreference**

- IBM Inc. announced today
- IBM profits
- US President Donald Trump announced yesterday
- for United States President Donald Trump

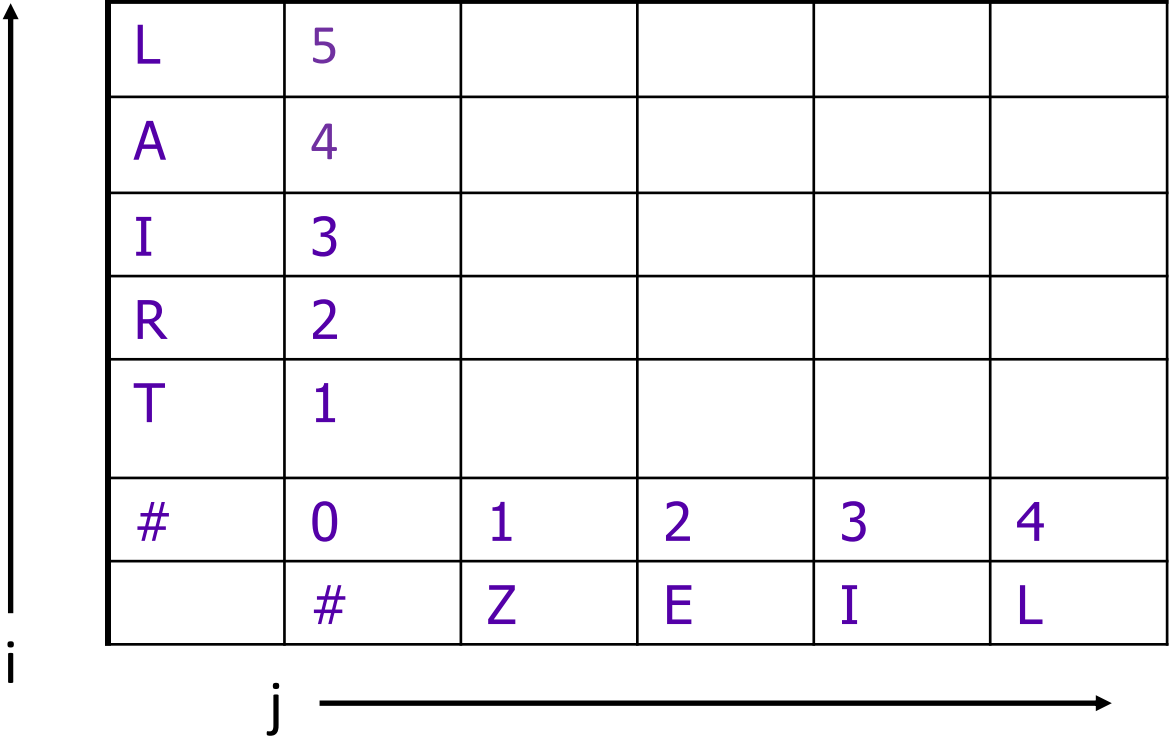
# Minimum Edit as Search

- 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

# 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  $D(n,m)$

# The Edit Distance Table



The diagram shows an edit distance table. To the left of the table is a vertical arrow pointing upwards, labeled with the letter 'i' at its base. Below the table is a horizontal arrow pointing to the right, labeled with the letter 'j' at its start. The table itself is a 7x6 grid. The first column contains the characters 'L', 'A', 'I', 'R', 'T', '#', and an empty cell. The first row contains the values '5', '4', '3', '2', '1', '0', and '#'. The remaining cells in the first row and first column are empty. The remaining cells in the table (from row 2 to row 7 and column 3 to column 6) are empty.

L	5				
A	4				
I	3				
R	2				
T	1				
#	0	1	2	3	4
	#	Z	E	I	L

# The Edit Distance Table

<div style="display: flex; align-items: center;"> <div style="width: 10px; height: 100px; border-left: 1px solid black; margin-right: 5px;"></div> <div style="writing-mode: vertical-rl; transform: rotate(180deg);">i</div> </div>	5	L	5	6	7	6	<b>5</b>
	4	A	4	5	6	5	6
	3	I	3	4	5	4	5
	2	R	2	3	4	5	6
	1	T	1	2	3	4	5
	0	#	0	1	2	3	4
			#	Z	E	I	L
			0	1	2	3	4
			j →				

$$D(i,j) = \min \begin{cases} D(i-1,j) + 1 \\ D(i,j-1) + 1 \\ D(i-1,j-1) + \begin{cases} 2; & \text{if } S_1(i) \neq S_2(j) \\ 0; & \text{if } S_1(i) = S_2(j) \end{cases} \end{cases}$$

# The Edit Distance Table

(INTENSION vs EXECUTION)

N	9	8	9	10	11	12	11	10	9	8
O	8	7	8	9	10	11	10	9	8	9
I	7	6	7	8	9	10	9	8	9	10
T	6	5	6	7	8	9	8	9	10	11
N	5	4	5	6	7	8	9	10	11	10
E	4	3	4	5	6	7	8	9	10	9
T	3	4	5	6	7	8	7	8	9	8
N	2	3	4	5	6	7	8	7	8	7
I	1	2	3	4	5	6	7	6	7	8
#	0	1	2	3	4	5	6	7	8	9
	#	E	X	E	C	U	T	I	O	N

# 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

# Edit Distance with Backtrace

T R I A L  
| | | | |  
Z E I \* L  
s s d

L	5	6	7	6	5
A	4	5	6	5	6
I	3	4	5	4	5
R	2	3	4	5	6
T	1	2	3	4	5
#	0	1	2	3	4
	#	Z	E	I	L

- ← Substitute
- ← No Change
- ← Delete

$$D(i,j) = \min \begin{cases} D(i-1,j) + 1 \\ D(i,j-1) + 1 \\ D(i-1,j-1) + \begin{cases} 2; & \text{if } S_1(i) \neq S_2(j) \\ 0; & \text{if } S_1(i) = S_2(j) \end{cases} \end{cases}$$



# Edit Distance with Backtrace (Another Path)

L	5	6	7	6	<b>5</b>
A	4	5	6	5	6
I	3	4	5	4	5
R	2	3	4	5	6
T	1	2	3	4	5
#	0	1	2	3	4
	#	Z	E	I	L

- ← Substitute
- ← No Change
- ← Delete
- ← Insert

Cost is same, i.e., 5

$$D(i,j) = \min \begin{cases} D(i-1,j) + 1 \\ D(i,j-1) + 1 \\ D(i-1,j-1) + \begin{cases} 2; & \text{if } S_1(i) \neq S_2(j) \\ 0; & \text{if } S_1(i) = S_2(j) \end{cases} \end{cases}$$

# MinEdit with Backtrace

<b>n</b>	9	↓ 8	↙←↓ 9	↙←↓ 10	↙←↓ 11	↙←↓ 12	↓ 11	↓ 10	↓ 9	↙ 8	
<b>o</b>	8	↓ 7	↙←↓ 8	↙←↓ 9	↙←↓ 10	↙←↓ 11	↓ 10	↓ 9	↙ 8	← 9	
<b>i</b>	7	↓ 6	↙←↓ 7	↙←↓ 8	↙←↓ 9	↙←↓ 10	↓ 9	↙ 8	← 9	← 10	
<b>t</b>	6	↓ 5	↙←↓ 6	↙←↓ 7	↙←↓ 8	↙←↓ 9	↙ 8	← 9	← 10	←↓ 11	
<b>n</b>	5	↓ 4	↙←↓ 5	↙←↓ 6	↙←↓ 7	↙←↓ 8	↙←↓ 9	↙←↓ 10	↙←↓ 11	↙↓ 10	
<b>e</b>	4	↙ 3	← 4	↙← 5	← 6	← 7	←↓ 8	↙←↓ 9	↙←↓ 10	↓ 9	
<b>t</b>	3	↙←↓ 4	↙←↓ 5	↙←↓ 6	↙←↓ 7	↙←↓ 8	↙ 7	←↓ 8	↙←↓ 9	↓ 8	
<b>n</b>	2	↙←↓ 3	↙←↓ 4	↙←↓ 5	↙←↓ 6	↙←↓ 7	↙←↓ 8	↓ 7	↙←↓ 8	↙ 7	
<b>i</b>	1	↙←↓ 2	↙←↓ 3	↙←↓ 4	↙←↓ 5	↙←↓ 6	↙←↓ 7	↙ 6	← 7	← 8	
<b>#</b>	0	1	2	3	4	5	6	7	8	9	
	<b>#</b>	<b>e</b>	<b>x</b>	<b>e</b>	<b>c</b>	<b>u</b>	<b>t</b>	<b>i</b>	<b>o</b>	<b>n</b>	

# Adding Backtrace to Minimum Edit Distance

- Base conditions:

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

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

Termination:

$D(N, M)$  is distance

- Recurrence Relation:

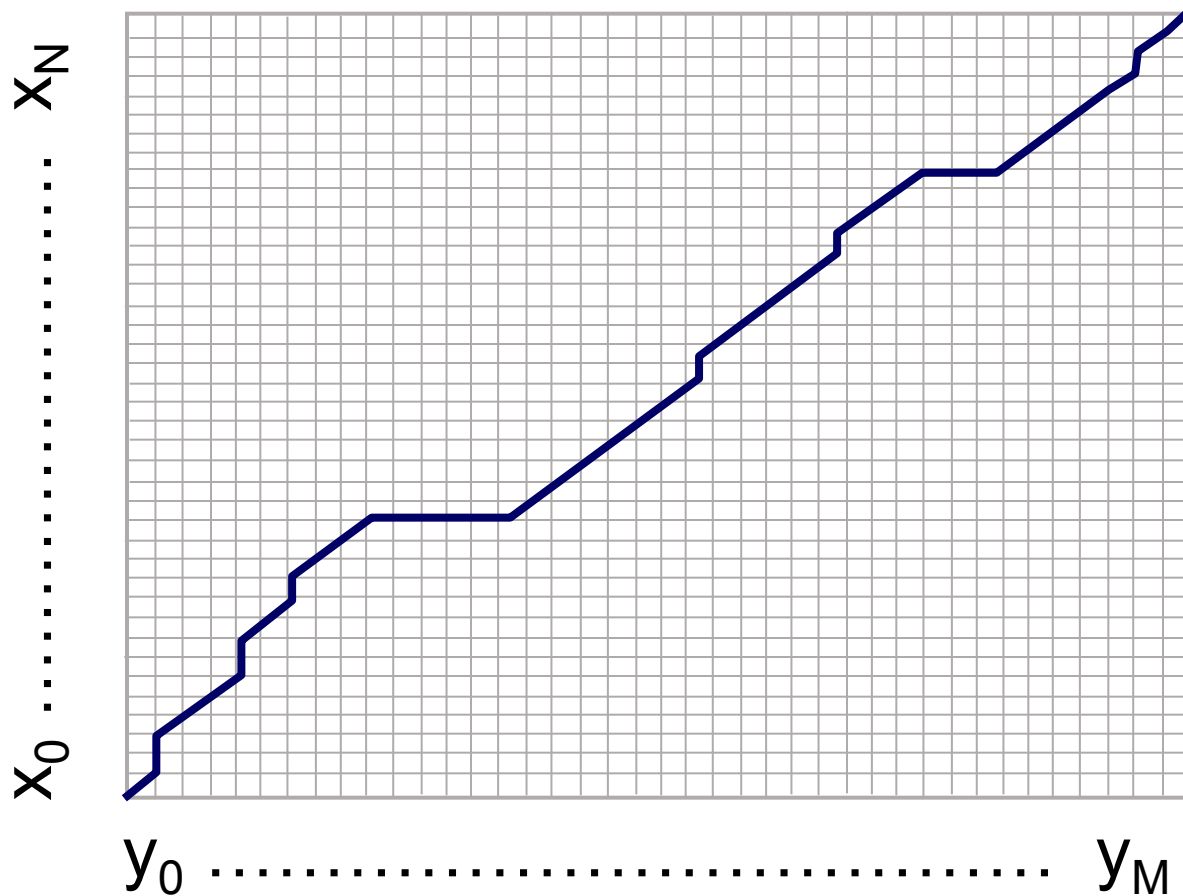
For each  $i = 1 \dots M$

For each  $j = 1 \dots N$

$$D(i, j) = \min \begin{cases} D(i-1, j) + 1 & \text{deletion} \\ D(i, j-1) + 1 & \text{insertion} \\ D(i-1, j-1) + \begin{cases} 2; & \text{if } X(i) \neq Y(j) \\ 0; & \text{if } X(i) = Y(j) \end{cases} & \text{substitution} \end{cases}$$

$$\text{ptr}(i, j) = \begin{cases} \text{LEFT} & \text{insertion} \\ \text{DOWN} & \text{deletion} \\ \text{DIAG} & \text{substitution} \end{cases}$$

# 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

# Performance

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