# WordNet Similarity & Edit Distance

Nathan Schneider
ANLP | 18 September 2017

*Most slides from Jurafsky & Martin SLP3 lectures*

1

# Computing with a thesaurus

Word Similarity:
Thesaurus Methods

# Word Similarity

- **Synonymy**: a binary relation
  - Two words are either synonymous or not
- **Similarity** (or **distance**): a looser metric
  - Two words are more similar if they share more features of meaning
- Similarity is properly a relation between **senses**
  - The word "bank" is not similar to the word "slope"
  - Bank[1] is similar to fund[3]
  - Bank[2] is similar to slope[5]
- But we'll compute similarity over both words and senses

# Why word similarity

- A practical component in lots of NLP tasks
  - Question answering
  - Natural language generation
  - Automatic essay grading
  - Plagiarism detection
- A theoretical component in many linguistic and cognitive tasks
  - Historical semantics
  - Models of human word learning
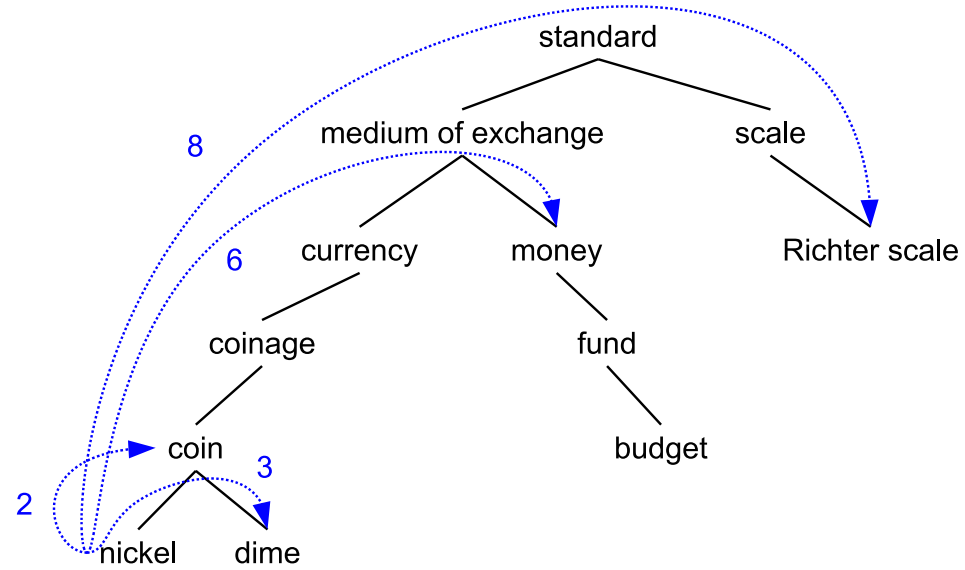  - Morphology and grammar induction

# **Word similarity and word relatedness**

- We often distinguish **word similarity** from **word relatedness**
  - **Similar words**: near-synonyms
  - **Related words**: can be related any way
    - `car, bicycle:` **similar**
    - `car, gasoline:` **related**, not similar

# Two classes of similarity algorithms

- Thesaurus-based algorithms
  - Are words "nearby" in hypernym hierarchy?
  - Do words have similar glosses (definitions)?
- Distributional algorithms
  - Do words have similar distributional contexts?
  - Distributional (Vector) semantics on Thursday!

# Path length



- Two concepts (senses/synsets) are similar if they are near each other in the thesaurus hierarchy
  - =have a short path between them
  - concepts have path 1 to themselves

# Path similarity

- pathlen($c_1, c_2$) = 1 + number of edges in the shortest path in the hypernym graph between sense nodes $c_1$ and $c_2$

- ranges from 0 to 1 (identity):

  - simpath($c_1, c_2$) = $\dfrac{1}{\text{pathlen}(c_1, c_2)}$

  - wordsim($w_1, w_2$) = $\max\limits_{c_1 \in \text{senses}(w_1), c_2 \in \text{senses}(w_2)} \text{sim}(c_1, c_2)$

# Example: path-based similarity

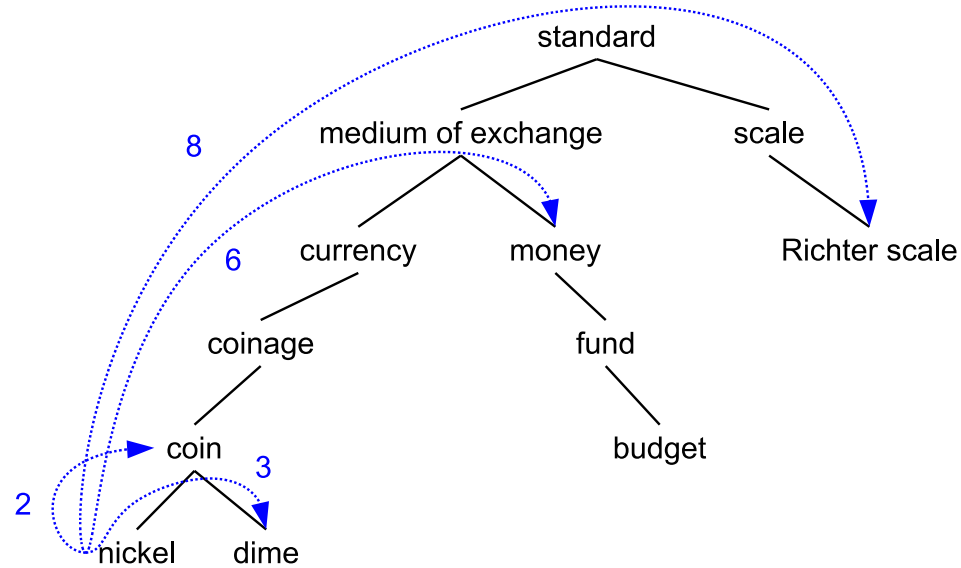$$\text{simpath}(c_1, c_2) = 1/\text{pathlen}(c_1, c_2)$$



simpath(*nickel,coin*) = 1/2 = .5

simpath(*fund,budget*) = 1/2 = .5

simpath(*nickel,currency*) = 1/4 = .25

simpath(*nickel,money*) = 1/6 = .17

simpath(*coinage,Richter scale*) = 1/6 = .17

# Problem with basic path-based similarity

- Assumes each link represents a uniform distance
  - But *nickel* to *money* seems to us to be closer than *nickel* to *standard*
  - Nodes high in the hierarchy are very abstract
- We instead want a metric that
  - Represents the cost of each edge independently
  - Words connected only through abstract nodes
    - are less similar

# Wu-Palmer similarity

- Let depth in the hierarchy inform semantic closeness

  - Least Common Subsumer (LCS) of two nodes: the most specific common ancestor in the hierarchy



simwup(nickel, dime) = 2*2/(1+1+2*2) = 2/3

simwup(nickel, medium of exchange) = 2*1/(2+0+2*1) = 1/2

simwup(nickel, nickel) = 2*3/(0+0+2*3) = 1

$$simwup(c_1, c_2) = \frac{2 \cdot d(root, LCS)}{d(LCS, c_1) + d(LCS, c_2) + 2 \cdot d(root, LCS)}$$

# Similarity & Distance: Form vs. Meaning

- Path similarity, Wu-Palmer similarity defined on [0,1]
  - Higher = more similar; 1 = identical
  - (There are other ways of measuring semantic similarity between words, as we will see later in the course!)
- With *distance* measures, higher = more different
- WordNet measures are about **meaning**. It is also useful measuring similarity/distance with respect to **form**.

# 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

    AGGCTATCACCTGACCTCCAGGCCGATGCCC
    TAGCTATCACGACCGCGGTCGATTTGCCCGAC

  - 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
  - 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 (Levenshtein)
  - Distance between these is 5
- If substitutions cost 2
  - Distance between them is 8

# Alignment in Computational Biology

- Given a sequence of bases

  AGGCTATCACCTGACCTCCAGGCCGATGCCC
  TAGCTATCACGACCGCGGTCGATTTGCCCGAC

- An alignment:

  -**AG**G**CTATCAC**CT**GACC**TC**CA**GG**C**CGA**--**TGCCC**---
  T**AG**-**CTATCAC**--**GACC**G**C**--**GG**T**CGA**TT**TGCCC**GAC

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

# Other uses of Edit Distance in NLP

- Evaluating Machine Translation and speech recognition

```
R Spokesman confirms    senior government adviser was shot
H Spokesman said    the senior          adviser was shot dead
             S       I                D                    I
```

- Named Entity Extraction and Entity Coreference

  - IBM Inc. announced today

  - IBM profits

  - Stanford President John Hennessy announced yesterday

  - for Stanford University President John Hennessy

# 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 revisted states.

21

# 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)

# Defining Min Edit Distance

- Initialization

  ```
  D(i,0) = i
  D(0,j) = j
  ```

- Recurrence Relation:

  ```
  For each  i = 1…M
     For each  j = 1…N
  ```

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

- Termination:

  ```
  D(N,M) is distance
  ```

# The Edit Distance Table

| N | 9 | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| O | 8 | | | | | | | | | |
| I | 7 | | | | | | | | | |
| T | 6 | | | | | | | | | |
| N | 5 | | | | | | | | | |
| E | 4 | | | | | | | | | |
| T | 3 | | | | | | | | | |
| N | 2 | | | | | | | | | |
| I | 1 | | | | | | | | | |
| # | **0** | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| | # | E | X | E | C | U | T | I | O | N |

# The Edit Distance Table

| N | 9 | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| O | 8 | | | | | | | | | |
| I | 7 | | | | | | | | | |
| T | 6 | | | | | | | | | |
| N | 5 | | | | | | | | | |
| E | 4 | | | | | | | | | |
| T | 3 | | | | | | | | | |
| N | 2 | | | | | | | | | |
| I | 1 | | | | | | | | | |
| # | **0** | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| | # | E | X | E | C | U | T | I | O | N |

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

# Edit Distance

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

| N | 9 | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| O | 8 | | | | | | | | | |
| I | 7 | | | | | | | | | |
| T | 6 | | | | | | | | | |
| N | 5 | | | | | | | | | |
| E | 4 | | | | | | | | | |
| T | 3 | | | | | | | | | |
| N | 2 | | | | | | | | | |
| I | 1 | | | | | | | | | |
| # | **0** | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|   | # | E | X | E | C | U | T | I | O | N |

# Edit Distance

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

| 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 |

# 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

# Edit Distance

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

| N | 9 | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| O | 8 | | | | | | | | | |
| I | 7 | | | | | | | | | |
| T | 6 | | | | | | | | | |
| N | 5 | | | | | | | | | |
| E | 4 | | | | | | | | | |
| T | 3 | | | | | | | | | |
| N | 2 | | | | | | | | | |
| I | 1 | | | | | | | | | |
| # | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| | # | E | X | E | C | U | T | I | O | N |

# MinEdit with Backtrace

| | | # | e | x | e | c | u | t | i | o | n |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **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** | |

# Adding Backtrace to Minimum Edit Distance

- Base conditions:

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

- Recurrence Relation:

  Termination:

  $D(N,M)$ is distance

```
For each  i = 1…M
  For each  j = 1…N
```

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

deletion

insertion

substitution

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

insertion

deletion

substitution

# 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

Slide adapted from Serafim Batzoglou

# 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

# Confusion matrix for spelling errors

**sub[X, Y] = Substitution of X (incorrect) for Y (correct)**

| X | a | b | c | d | e | f | g | h | i | j | k | l | m | n | o | p | q | r | s | t | u | v | w | x | y | z |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| a | 0 | 0 | 7 | 1 | 342 | 0 | 0 | 2 | 118 | 0 | 1 | 0 | 0 | 3 | 76 | 0 | 0 | 1 | 35 | 9 | 9 | 0 | 1 | 0 | 5 | 0 |
| b | 0 | 0 | 9 | 9 | 2 | 2 | 3 | 1 | 0 | 0 | 0 | 5 | 11 | 5 | 0 | 10 | 0 | 0 | 2 | 1 | 0 | 0 | 8 | 0 | 0 | 0 |
| c | 6 | 5 | 0 | 16 | 0 | 9 | 5 | 0 | 0 | 0 | 1 | 0 | 7 | 9 | 1 | 10 | 2 | 5 | 39 | 40 | 1 | 3 | 7 | 1 | 1 | 0 |
| d | 1 | 10 | 13 | 0 | 12 | 0 | 5 | 5 | 0 | 0 | 2 | 3 | 7 | 3 | 0 | 1 | 0 | 43 | 30 | 22 | 0 | 0 | 4 | 0 | 2 | 0 |
| e | 388 | 0 | 3 | 11 | 0 | 2 | 2 | 0 | 89 | 0 | 0 | 3 | 0 | 5 | 93 | 0 | 0 | 14 | 12 | 6 | 15 | 0 | 1 | 0 | 18 | 0 |
| f | 0 | 15 | 0 | 3 | 1 | 0 | 5 | 2 | 0 | 0 | 0 | 3 | 4 | 1 | 0 | 0 | 0 | 6 | 4 | 12 | 0 | 0 | 2 | 0 | 0 | 0 |
| g | 4 | 1 | 11 | 11 | 9 | 2 | 0 | 0 | 0 | 1 | 1 | 3 | 0 | 0 | 2 | 1 | 3 | 5 | 13 | 21 | 0 | 0 | 1 | 0 | 3 | 0 |
| h | 1 | 8 | 0 | 3 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 0 | 12 | 14 | 2 | 3 | 0 | 3 | 1 | 11 | 0 | 0 | 2 | 0 | 0 | 0 |
| i | 103 | 0 | 0 | 0 | 146 | 0 | 1 | 0 | 0 | 0 | 0 | 6 | 0 | 0 | 49 | 0 | 0 | 0 | 2 | 1 | 47 | 0 | 2 | 1 | 15 | 0 |
| j | 0 | 1 | 1 | 9 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 2 | 1 | 0 | 0 | 0 | 0 | 0 | 5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| k | 1 | 2 | 8 | 4 | 1 | 1 | 2 | 5 | 0 | 0 | 0 | 0 | 5 | 0 | 2 | 0 | 0 | 0 | 6 | 0 | 0 | 0 | 4 | 0 | 0 | 3 |
| l | 2 | 10 | 1 | 4 | 0 | 4 | 5 | 6 | 13 | 0 | 1 | 0 | 0 | 14 | 2 | 5 | 0 | 11 | 10 | 2 | 0 | 0 | 0 | 0 | 0 | 0 |
| m | 1 | 3 | 7 | 8 | 0 | 2 | 0 | 6 | 0 | 0 | 4 | 4 | 0 | 180 | 0 | 6 | 0 | 0 | 9 | 15 | 13 | 3 | 2 | 2 | 3 | 0 |
| n | 2 | 7 | 6 | 5 | 3 | 0 | 1 | 19 | 1 | 0 | 4 | 35 | 78 | 0 | 0 | 7 | 0 | 28 | 5 | 7 | 0 | 0 | 1 | 2 | 0 | 2 |
| o | 91 | 1 | 1 | 3 | 116 | 0 | 0 | 0 | 25 | 0 | 2 | 0 | 0 | 0 | 0 | 14 | 0 | 2 | 4 | 14 | 39 | 0 | 0 | 0 | 18 | 0 |
| p | 0 | 11 | 1 | 2 | 0 | 6 | 5 | 0 | 2 | 9 | 0 | 2 | 7 | 6 | 15 | 0 | 0 | 1 | 3 | 6 | 0 | 4 | 1 | 0 | 0 | 0 |
| q | 0 | 0 | 1 | 0 | 0 | 0 | 27 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| r | 0 | 14 | 0 | 30 | 12 | 2 | 2 | 8 | 2 | 0 | 5 | 8 | 4 | 20 | 1 | 14 | 0 | 0 | 12 | 22 | 4 | 0 | 0 | 1 | 0 | 0 |
| s | 11 | 8 | 27 | 33 | 35 | 4 | 0 | 1 | 0 | 1 | 0 | 27 | 0 | 6 | 1 | 7 | 0 | 14 | 0 | 15 | 0 | 0 | 5 | 3 | 20 | 1 |
| t | 3 | 4 | 9 | 42 | 7 | 5 | 19 | 5 | 0 | 1 | 0 | 14 | 9 | 5 | 5 | 6 | 0 | 11 | 37 | 0 | 0 | 2 | 19 | 0 | 7 | 6 |
| u | 20 | 0 | 0 | 0 | 44 | 0 | 0 | 0 | 64 | 0 | 0 | 0 | 0 | 2 | 43 | 0 | 0 | 4 | 0 | 0 | 0 | 0 | 2 | 0 | 8 | 0 |
| v | 0 | 0 | 7 | 0 | 0 | 3 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 8 | 3 | 0 | 0 | 0 | 0 | 0 | 0 |
| w | 2 | 2 | 1 | 0 | 1 | 0 | 0 | 2 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 7 | 0 | 6 | 3 | 3 | 1 | 0 | 0 | 0 | 0 | 0 |
| x | 0 | 0 | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 9 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| y | 0 | 0 | 2 | 0 | 15 | 0 | 1 | 7 | 15 | 0 | 0 | 0 | 2 | 0 | 6 | 1 | 0 | 7 | 36 | 8 | 5 | 0 | 0 | 1 | 0 | 0 |
| z | 0 | 0 | 0 | 7 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 7 | 5 | 0 | 0 | 0 | 0 | 2 | 21 | 3 | 0 | 0 | 0 | 0 | 3 | 0 |

# Weighted Min Edit Distance

- Initialization:

  ```
  D(0,0) = 0
  D(i,0) = D(i-1,0) + del[x(i)];    1 < i ≤ N
  D(0,j) = D(0,j-1) + ins[y(j)];    1 < j ≤ M
  ```

- Recurrence Relation:

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

- Termination:

  ```
  D(N,M) is distance
  ```

# Where did the name, dynamic programming, come from?

…The 1950s were not good years for mathematical research. [the] Secretary of Defense …had a pathological fear and hatred of the word, research…

 I decided therefore to use the word, "**programming**".

 I wanted to get across the idea that this was dynamic, this was multistage… I thought, let's … take a word that has an absolutely precise meaning, namely **dynamic**… it's impossible to use the word, **dynamic**, in a pejorative sense. Try thinking of some combination that will possibly give it a pejorative meaning. It's impossible.

Thus, I thought dynamic programming was a good name. It was something not even a Congressman could object to."

Richard Bellman, "Eye of the Hurricane: an autobiography" 1984.