# Minimum Cost Edit Distance

- Edit a source string into a target string
- Each edit has a cost
- Find the minimum cost edit(s)



## Minimum Cost Edit Distance

а	C	t	r	е	S		S	target
	I							
-	C	_	r	e	S	t	—	source

rce



minimum cost edit distance can be accomplished in multiple ways

ress a С t е S





Only 4 ways to edit source to target for this pair

## Levenshtein Distance

- Cost is fixed across characters
  - Insertion cost is 1
  - Deletion cost is 1
- Two different costs for substitutions
  - Substitution cost is 1 (transformation)
  - Substitution cost is 2 (one deletion + one insertion)



What's the edit distance?

Левенштейн Владимир Vladimir Levenshtein

## Minimum Cost Edit Distance

• An alignment between target and source

$$t_{1}, t_{2}, \dots, t_{n}$$
Find  $D(n,m)$  recursively  

$$s_{1}, s_{2}, \dots, s_{m}$$

$$D(i, j) = min \begin{cases} D(i-1, j) + \cos(t_{i}, \emptyset) \text{ insertion into target} \\ D(i-1, j-1) + \cos(t_{i}, s_{j}) \text{ substitution/identity} \\ D(i, j-1) + \cos((\emptyset, s_{j}) \text{ deletion from source} \end{cases}$$

$$D(0,0) = 0 \qquad D(i,0) = D(i-1,0) + \cos(t_{i}, \emptyset) \\ D(0, j) = D(0, j-1) + \cos((\emptyset, s_{j}))$$

Function MinEditDistance (target, source)

```
n = length(target)
m = length(source)
Create matrix D of size (n+1,m+1)
D[0,0] = 0
for i = 1 to n
  D[i,0] = D[i-1,0] + insert-cost
for j = 1 to m
 D[0,j] = D[0,j-1] + delete-cost
for i = 1 to n
  for j = 1 to m
    D[i,j] = MIN(D[i-1,j] + insert-cost,
                 D[i-1,j-1] + subst/eq-cost,
                 D[i,j-1] + delete-cost)
return D[n,m]
```

Consider two strings:

$$target = g_1 a_2 m_3 b_4 l_5 e_6$$

source=  $g_1 u_2 m_3 b_4 o_5$ 

- We want to find D(6,5)
- We find this recursively using values of D(i,j) where  $i \le 6 \le 5$
- For example, consider how to compute D(4,3)

 $target = g_1 a_2 m_3 b_4$ source=  $g_1 u_2 m_3$ D(3,2) D(4,2) D(3,3) \leftarrow D(4,3)

- Case 1: SUBSTITUTE  $b_4$  for  $m_3$
- Use previously stored value for D(3,2)
- $Cost(g_1a_2m_3b \text{ and } g_1u_2m) = D(3,2) + cost(b \approx m)$
- For substitution: D(i,j) = D(i-1,j-1) + cost(subst)
  - Case 2: INSERT b<sub>4</sub>
- Use previously stored value for D(3,3)
- $Cost(g_1a_2m_3b \text{ and } g_1u_2m_3) = D(3,3) + cost(ins b)$
- For substitution: D(i,j) = D(i-1,j) + cost(ins)
- Case 3: DELETE m<sub>3</sub>
- Use previously stored value for D(4,2)
- $Cost(g_1a_2m_3b_4 \text{ and } g_1u_2m) = D(4,2) + cost(del m)$
- For substitution: D(i,j) = D(i,j-1) + cost(del)

b e g a m  $\mathbf{0}$  $\mathbf{0}$ g D u 2) m D b i е b 1 e m а g b u m g \_

source

target

- Algorithm using a Finite-state transducer:
  - construct a finite-state transducer with all possible ways to transduce source into target
  - We do this transduction one char at a time
  - A transition x:x gets zero cost and a transition on ε:x (insertion) or x:ε (deletion) for any char x gets cost 1
  - Finding minimum cost edit distance == Finding the shortest path from start state to final state

- Lets assume we want to edit source string 1010 into the target string 1110
- The alphabet is just 1 and 0



• Construct a FST that allows strings to be edited



• Compose SOURCE and EDITS and TARGET



• The shortest path is the minimum edit FST from SOURCE (1010) to TARGET (1110)

$$6 \xrightarrow{1:1} 5 \xrightarrow{0:} 4 \xrightarrow{1:1} 3 \xrightarrow{0:} 2 \xrightarrow{:1} 1 \xrightarrow{:0} 0$$

## Edit distance

- Useful in many NLP applications
- In some cases, we need edits with multiple characters, e.g. 2 chars deleted for one cost
- Comparing system output with human output, e.g. <u>input:</u> ibm <u>output:</u> IBM vs. Ibm (TrueCasing of speech recognition output)
- Error correction
- Defined over character edits or word edits, e.g. MT evaluation:
  - Foreign investment in Jiangsu 's agriculture on the increase
  - Foreign investment in Jiangsu agricultural investment increased

Pronunciation dialect map of the Netherlands based on phonetic edit-distance (W. Heeringa Phd thesis, 2004)



14

## Variable Cost Edit Distance

- So far, we have seen edit distance with uniform insert/ delete cost
- In different applications, we might want different insert/ delete costs for different items
- For example, consider the simple application of spelling correction
- Users typing on a qwerty keyboard will make certain errors more frequently than others
- So we can consider insert/delete costs in terms of a probability that a certain alignment occurs between the *correct* word and the *typo* word

# Spelling Correction

• Types of spelling correction – non-word error detection

e.g. hte for the

- isolated word error detection

e.g. *acres* vs. *access* (cannot decide if it is the right word for the context)

– context-dependent error detection (real world errors)

e.g. she is a talented acres vs. she is a talented actress

• For simplicity, we will consider the case with exactly 1 error



### Bayes Rule: computing P(orig | noisy)

• let *x* = original input, *y* = noisy observation  $p(x | y) = \frac{p(x, y)}{p(y)}$   $p(y | x) = \frac{p(y, x)}{p(x)}$ p(x, y) = p(y, x) $p(x \mid y) \times p(y) = p(y \mid x) \times p(x)$  $p(x \mid y) = \frac{p(y \mid x) \times p(x)}{p(y)} \qquad \underline{Bayes \ Rule}$ 

#### Chain Rule

$$p(a,b,c \mid d) = p(a \mid b,c,d) \times$$
$$p(b \mid c,d) \times$$
$$p(c \mid d)$$

Approximations: Bias vs. Variance

 $p(a \mid b, c, d) \approx p(a \mid b, c)$  less bias  $p(a \mid b)$ p(a) less variance

# Single Error Spelling Correction

• Insertion (addition)

acress vs. cress

- Deletion
  - acress vs. actress
- Substitution
  - acress vs. access
- Transposition (reversal)

– acress vs. caress

Noisy Channel Model for Spelling Correction (Kernighan, Church and Gale, 1990)

• *t* is the word with a single typo and *c* is the correct word

$$P(c \mid t) = p(t \mid c) \times p(c)$$
 Bayes Rule

• Find the best candidate for the correct word

$$\hat{c} = \frac{\arg \max}{c \in C} P(t \mid c) \times P(c)$$
$$P(t \mid c) = ?? \qquad P(c) = \frac{f(c)}{N}$$

C is all the words in the vocabulary; |C| = N

#### Noisy Channel Model for Spelling Correction (Kernighan, Church and Gale, 1990) single error, condition on previous letter



P(poton | potion)

$$P(t \mid c) =$$

P(poton | piton)



$$\frac{del[c_{p-1},c_p]}{chars[c_{p-1},c_p]} (xy)_c \text{ typed as } (x)_t$$

$$t = poton$$

$$c = potion$$

$$del[t,i]=427$$

$$chars[t,i]=575$$

$$P = .7426$$

$$\frac{sub[t_p,c_p]}{chars[c_p-1]} (x)_c \text{ typed as } (xy)_t$$

$$t = poton$$

$$c = poton$$

$$del[t,i]=427$$

$$chars[t,i]=575$$

$$P = .7426$$

$$t = poton$$

$$c = piton$$

$$c = piton$$

$$sub[o,i]=568$$

$$chars[i]=1406$$

$$P = .4039$$

#### Noisy Channel model for Spelling Correction

• The *del*, *ins*, *sub*, *rev* matrix values need data in which contain known errors (training data)

#### (training data)

e.g. Birbeck spelling error corpus (from 1984!)

 Accuracy on single errors on unseen data (<u>test data</u>)

#### Noisy Channel model for Spelling Correction

- Easily extended to multiple spelling errors in a word using edit distance algorithm (however, using learned costs for ins, del, replace)
- Experiments: 87% accuracy for machine vs. 98% average human accuracy
- What are the limitations of this model?

. . .

... was called a "stellar and versatile **acress** whose combination of sass and glamour has defined her

KCG model best guess is acres