

# **String Matching**

**Pedro Szekely & Craig Knoblock**  
**University of Southern California**

# Isn't the Problem Solved?

String.equalsIgnoreCase(String x)

# How Many Publications?



# Why So Many?

The screenshot shows a web browser window titled "string matching – Google Scholar". The address bar indicates the search query is "string matching" on scholar.google.com. The search interface includes tabs for "Web", "Images", and "More...". The user's email, "szekely1401@gmail.com", is visible in the top right. Below the search bar, the Google logo is on the left, and the search term "string matching" is entered. To the right of the search bar is a blue "Search" button with a magnifying glass icon. Further down, the word "Scholar" is highlighted in red, followed by the text "About 614,000 results (0.03 sec)". On the right side of this row are filters for "Any time" and a notification bell icon showing "0" notifications. The main content area displays two search results. The first result is for a paper by AV Aho and MJ Corasick, titled "Efficient string matching: an aid to bibliographic search". It includes a link to "isrl.kr [PDF]", the authors' names, and the journal "Communications of the ACM, 1975 - dl.acm.org". The abstract describes a simple, efficient algorithm for string matching. The second result is for a paper by J Cleary and I Witten, titled "Data compression using adaptive coding and partial string matching". It includes a link to "ieeexplore.ieee.org", the authors' names, and the journal "Communications, IEEE Transactions on, 1984". The abstract discusses the use of arithmetic coding and a Markov model for data compression.

string matching – Google Scholar

scholar.google.com/scholar?hl=en&q=string+matching&btnG=

Web Images More... szekely1401@gmail.com

Google string matching

Scholar About 614,000 results (0.03 sec) Any time 0

[Efficient string matching: an aid to bibliographic search](#) isrl.kr [PDF]  
AV Aho, MJ Corasick - Communications of the ACM, 1975 - dl.acm.org  
Abstract This paper describes a simple, efficient algorithm to locate all occurrences of any of a finite number of keywords in a **string** of text. The algorithm consists of constructing a finite state pattern **matching** machine from the keywords and then using the pattern **matching** ...  
Cited by 2234 Related articles All 69 versions Import into BibTeX More ▾

[Data compression using adaptive coding and partial string matching](#)  
J Cleary, I Witten - Communications, IEEE Transactions on, 1984 - ieeexplore.ieee.org  
Abstract The recently developed technique of arithmetic coding, in conjunction with a Markov model of the source, is a powerful method of data compression in situations where a

# Multiple John Singer Sargents?

```
dallas:John_Singer_Sargent
  a foaf:Person;
  :dateOfBirth "1856" ;
  :dateOfDeath "1925" ;
  :name "John Singer Sargent" .
```

```
ima:John_Singer_Sargent
  a foaf:Person;
  :dateOfBirth "1856" ;
  :dateOfDeath "1925" ;
  :name "John S. Sargent" .
```

# Multiple John Singer Sargents?

```
dallas:John_Singer_Sargent
  a foaf:Person;
  :dateOfBirth "1856" ;
  :dateOfDeath "1925" ;
  :name "John Singer Sargent" .
```

string\_match("John Singer Sargent", "John S. Sargent") = ???

```
ima:John_Singer_Sargent
  a foaf:Person;
  :dateOfBirth "1856" ;
  :dateOfDeath "1925" ;
  :name "John S. Sargent" .
```

# String Matching Problem

myMatchFunction(x, y)

yourMatchFunction(x, y)

What does it mean that one is better than the other?

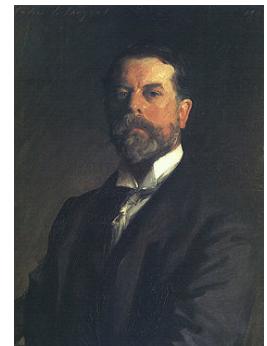
# Problem Definition

Given  $X$  and  $Y$  sets of strings

Find pairs  $(x, y)$   
such that both  $x$  and  $y$   
refer to the same real world entity

"John S. Sargent"

"John Singer Sargent"



# Problem Definition

Given  $X$  and  $Y$  sets of strings

Find pairs  $(x, y)$   
such that both  $x$  and  $y$   
refer to the same real world entity

We can use **precision** and **recall** to evaluate algorithms

# Problem Definition

Given  $X$  and  $Y$  sets of strings

Find pairs  $(x, y)$   
such that both  $x$  and  $y$   
refer to the same real world entity

fraction of pairs found that are correct



We can use **precision** and **recall** to evaluate algorithms

fraction of pairs found

# Why Strings Don't Match Perfectly?

typos     "Joh" vs "John"

OCR errors     "J0hn" vs "John"

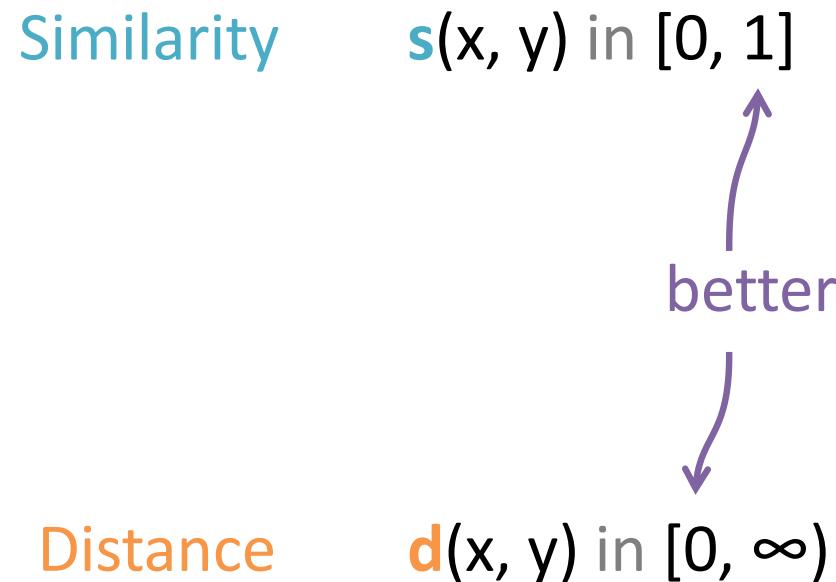
formatting conventions     "03/17" vs "March 17"

abbreviations     "J. S. Sargent" vs "John Singer Sargent"

nick names     "John" vs "Jock"

word order     "Sargent, John S." vs "John S. Sargent"

# Similiarity Measure



# Types of Similarity Metrics

- Sequence based
- Set based
- Hybrid
- Phonetic

# Sequence Based Metrics

# Edit Distance

"J0n Singer Sargent"



"John S. Sargent"

insert character

delete character

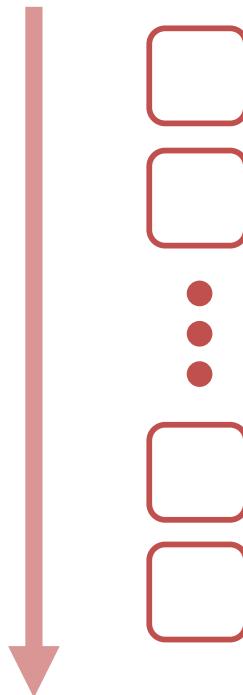
substitute character

transpose character

...

# Edit Distance

"J0n Singer Sargent"



"John S. Sargent"

costs

insert character  $c_1$

delete character  $c_2$

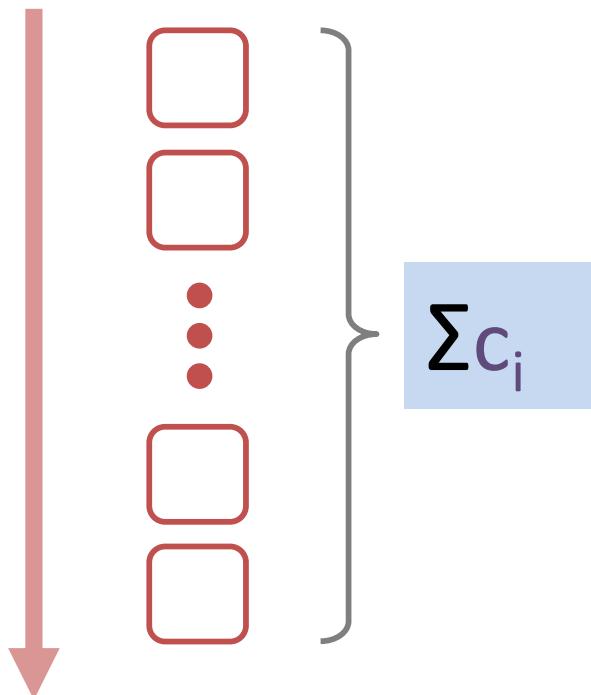
substitute character  $c_3$

transpose character  $c_4$

...  $c_i$

# Edit Distance

"J0n Singer Sargent"



"John S. Sargent"

costs

insert character  $C_1$

delete character  $C_2$

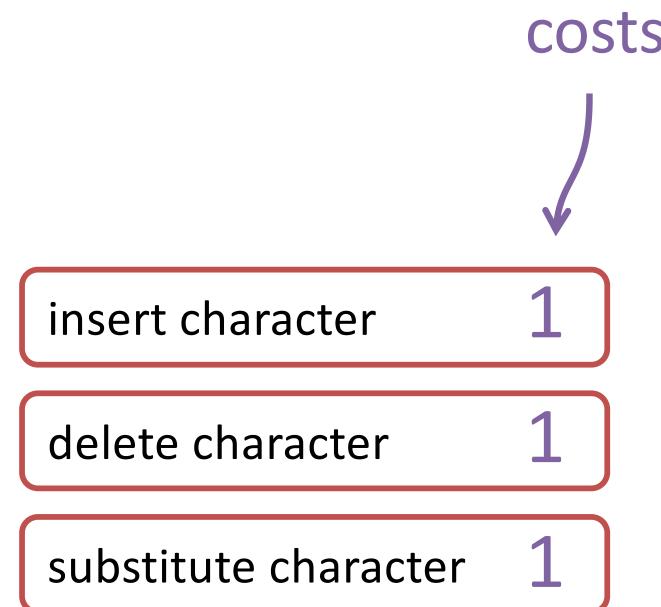
substitute character  $C_3$

transpose character  $C_4$

...  $C_i$

# Levenshtein Distance

Edit distance:



$\text{lev}(x, y)$  is the minimum cost to transform  $x$  to  $y$

Online calculator: <http://planetcalc.com/1721/>

# Computing Levenshtein Distance

$\text{lev}(x, y)$  is the minimum cost to transform  $x$  to  $y$

## Definitions

$$x = x_1 x_2 \dots x_n$$

$$y = y_1 y_2 \dots y_m$$

$$d(i, j) = \text{lev}(x_1 x_2 \dots \underline{x_i}, y_1 y_2 \dots \underline{y_j})$$

$$d(0, 0) = \text{lev}("", "") = 0$$

We want  $d(n, m)$

# Computing Levenshtein Distance

$$d(\emptyset, \emptyset) = 0$$

# Computing Levenshtein Distance

$$d(\emptyset, \emptyset) = 0$$

$x_1x_2 \dots x_{i-1}x_i$  is a suffix of  $x$        $y_1y_2 \dots y_{j-1}y_j$  is a suffix of  $y$

# Computing Levenshtein Distance

$$d(0, 0) = 0$$

$x_1x_2 \dots x_{i-1}x_i$  is a prefix of  $x$        $y_1y_2 \dots y_{j-1}y_j$  is a prefix of  $y$

Case	Distance	Operation
$x_i = y_j$	$d(i-1, j-1)$	keep $x_i$

# Computing Levenshtein Distance

$$d(\theta, \theta) = \theta$$

$x_1x_2 \dots x_{i-1}x_i$  is a prefix of  $x$        $y_1y_2 \dots y_{j-1}y_j$  is a prefix of  $y$

Case	Distance	Operation
$x_i = y_j$	$d(i-1, j-1)$	keep $x_i$
$x_i \neq y_j$		delete $x_i$ insert $y_j$ after $x_i$ replace $x_i$ with $y_j$

# Computing Levenshtein Distance

$$d(0, 0) = 0$$

$x_1x_2 \dots x_{i-1}x_i$  is a prefix of  $x$        $y_1y_2 \dots y_{j-1}y_j$  is a prefix of  $y$

Case	Distance	Operation
$x_i = y_j$	$d(i-1, j-1)$	keep $x_i$
$x_i \neq y_j$	$d(i-1, j) + 1$	delete $x_i$ insert $y_j$ after $x_i$ replace $x_i$ with $y_j$

# Computing Levenshtein Distance

$$d(0, 0) = 0$$

$x_1 x_2 \dots x_{i-1} x_i$  is a prefix of  $x$        $y_1 y_2 \dots y_{j-1} y_j$  is a prefix of  $y$

Case	Distance	Operation
$x_i = y_j$	$d(i-1, j-1)$	keep $x_i$
$x_i \neq y_j$	$d(i-1, j) + 1$	delete $x_i$
	$d(i, j-1) + 1$	insert $y_j$ after $x_i$
		replace $x_i$ with $y_j$

# Computing Levenshtein Distance

$$d(0, 0) = 0$$

$x_1 x_2 \dots x_{i-1} x_i$  is a prefix of  $x$        $y_1 y_2 \dots y_{j-1} y_j$  is a prefix of  $y$

Case	Distance	Operation
$x_i = y_j$	$d(i-1, j-1)$	keep $x_i$
$x_i \neq y_j$	$d(i-1, j) + 1$	delete $x_i$
	$d(i, j-1) + 1$	insert $y_j$ after $x_i$
	$d(i-1, j-1) + 1$	replace $x_i$ with $y_j$

# Computing Levenshtein Distance

$$d(0,0) = 0$$

$x_1x_2 \dots x_{i-1}x_i$  is a prefix of  $x$        $y_1y_2 \dots y_{j-1}y_j$  is a prefix of  $y$

Case	Distance	Operation
$x_i = y_j$	$d(i-1, j-1)$	keep $x_i$
$x_i \neq y_j$	$d(i-1, j) + 1$	delete $x_i$
	$d(i, j-1) + 1$	insert $y_j$ after $x_i$
	$d(i-1, j-1) + 1$	replace $x_i$ with $y_j$

$$d(i, j) = \text{minimum } \uparrow$$

# Computing Levenshtein Distance Using Dynamic Programming

- $x = \text{dva}$ ,  $y = \text{dave}$

		y0	y1	y2	y3	y4
		d	a	v	e	
x0		0	1	2	3	4
x1	d	1	0 ← 1			
x2	v	2				
x3	a	3				

		y0	y1	y2	y3	y4
		d	a	v	e	
x0		0	1	2	3	4
x1	d	1	0 ← 1	1 ← 2	2 ← 3	
x2	v	2	1	1	1 ← 2	
x3	a	3	2	1 ← 2	2	

$x = d - v a$   
| | | |  
 $y = d a v e$

substitute a with e  
insert a (after d)

- Cost of dynamic programming is  $O(|x| |y|)$

# Levenshtein Distance Complexity

Dynamic programming algorithm

Time Complexity =  $O(N * M)$ , Space Complexity =  $O(N * M)$

## Polylogarithmic Approximation for Edit Distance and the Asymmetric Query Complexity

Alexandr Andoni, Robert Krauthgamer, Krzysztof Onak

We present a near-linear time algorithm that approximates the edit distance between two strings within a polylogarithmic factor; specifically, for strings of length  $n$  and every fixed  $\epsilon > 0$ , it can compute a  $(\log n)^{O(1/\epsilon)}$  approximation in  $n^{(1+\epsilon)}$  time.

<http://arxiv.org/abs/1005.4033>

# Levenshtein Distance Examples

`lev(John Singer Sargent, John S. Sargent) =`

# Levenshtein Distance Examples

`lev(John Singer Sargent, John S. Sargent) = 5`

# Levenshtein Distance Examples

`lev(John Singer Sargent, John S. Sargent) = 5`

`lev(John Singer Sargent, Jane Klinger Sargent) =`

# Levenshtein Distance Examples

`lev(John Singer Sargent, John S. Sargent) = 5`

`lev(John Singer Sargent, Jane Klinger Sargent) = 5`

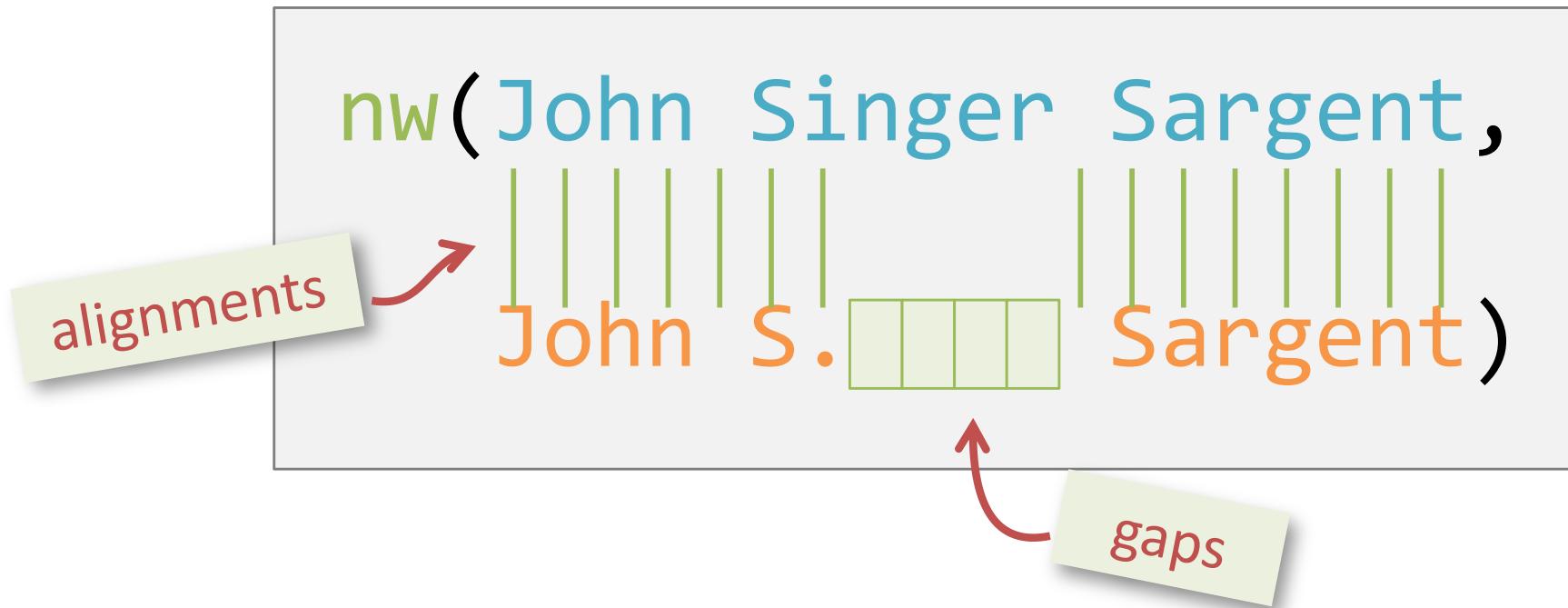
# Levenshtein Distance Examples

Too high a cost for deleting a sequence of characters


$$\text{lev}(\text{John Singer Sargent, John S.}, \text{Sargent}) = 5$$
$$\text{lev}(\text{John Singer Sargent, Jane Klinger Sargent}) = 5$$

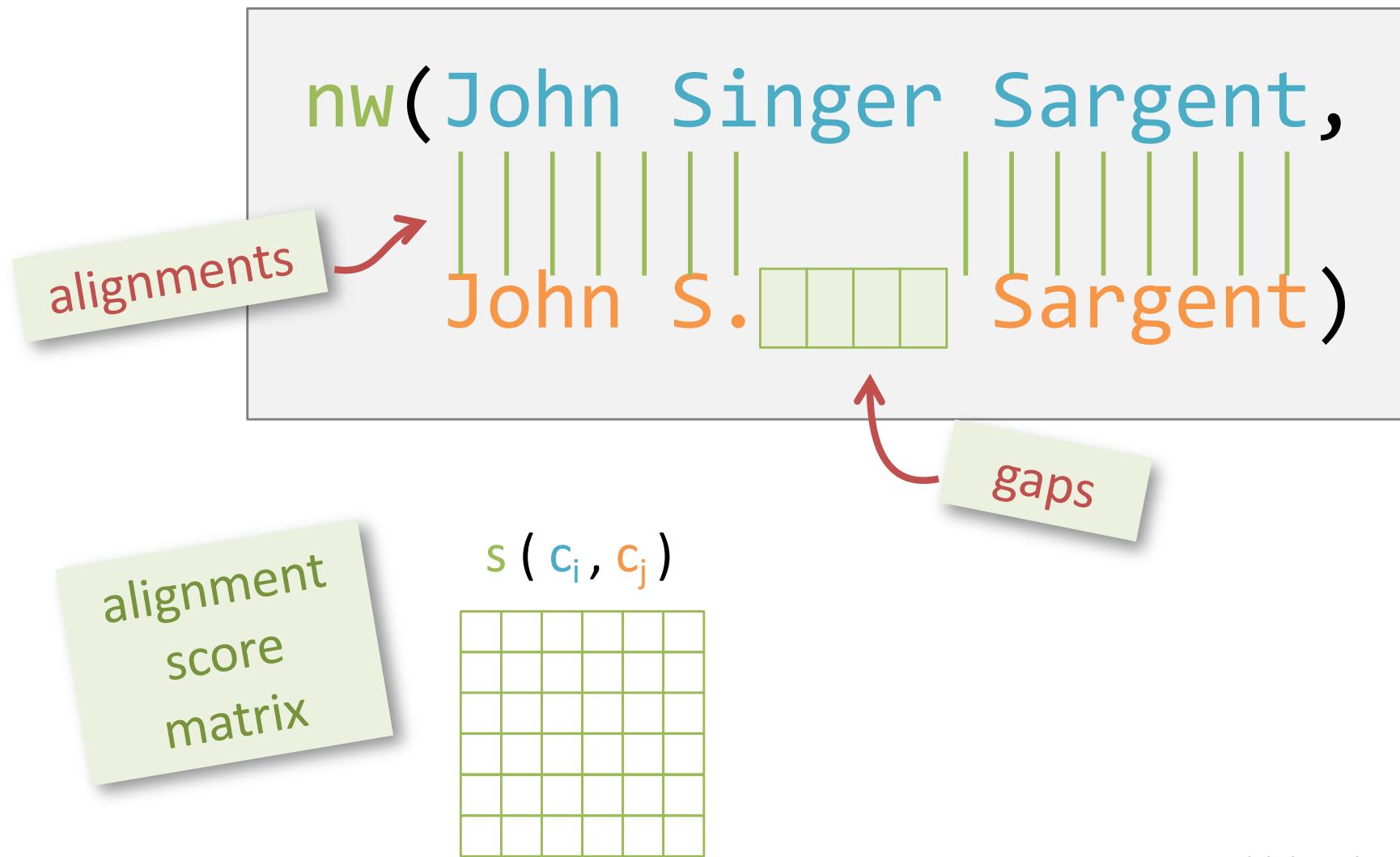
# Needleman-Wunch Measure

Generalization of  $\text{levenstein}(x, y)$



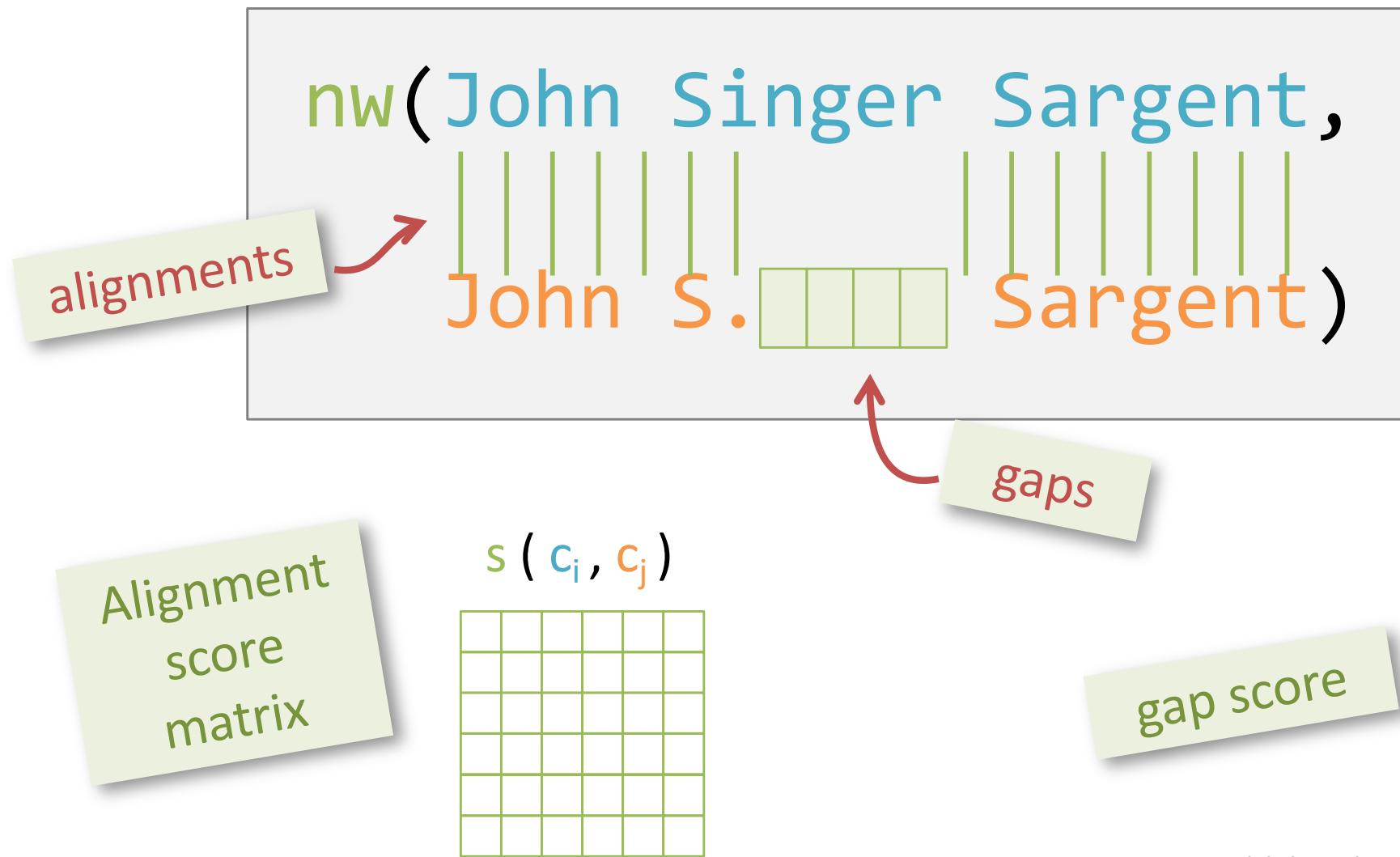
# Needleman-Wunch Measure

Generalization of levenstein( $x, y$ )



# Needleman-Wunch Measure

Generalization of levenstein( $x, y$ )



# Needleman-Wunch Measure

Generalization of levenstein( $x, y$ )

nw( John Singer Sargent,  
John S. [ ] Sargent)

$$s(c_i, c_j) = \begin{cases} 2 & \text{if } c_i = c_j \\ -1 & \text{if } c_i \neq c_j \end{cases}$$

gap-score = -0.5

# Needleman-Wunch Measure

Generalization of levenstein( $x, y$ )

$nw(\text{John Singer Sargent},$   
  
 $\text{John } S. \square\square\square\square \text{Sargent})$

$$2 * 14 + (-1) * 1 + (-0.5) * 4 = 25$$

$$s(c_i, c_j) = \begin{cases} 2 & \text{if } c_i = c_j \\ -1 & \text{if } c_i \neq c_j \end{cases}$$

gap-score = -0.5

slide by Pedro Szekely



# Comparison

	Levenshtein	Needleman-Wunch
Costs	1	matrix
Operations	insert/delete	gaps
Result	distance	similarity

OCR errors "J0hn" vs "John"

$$\text{score}(o, \theta) = -0.2$$

$$\text{score}(m, \theta) = -1.0$$

lower penalty

# Needleman-Wunch Example

$\text{nw}(\text{John Singer Sargent},$   
 $\text{John S.} \quad \text{Sargent})$

$$2 * 14 + (-1) * 1 + (-0.5) * 4 = 25$$

$$2 * 14 + (-1) * 1 + (-0.5) * 8 = 23$$

$\text{nw}(\text{John Stanislaus Sargent},$   
 $\text{John S.} \quad \text{Sargent})$

# Needleman-Wunch Example

$\text{nw}(\text{John Singer Sargent},$   
 $\text{John S.} \quad \text{Sargent})$

Longer gaps are  
penalized more

Bad for names

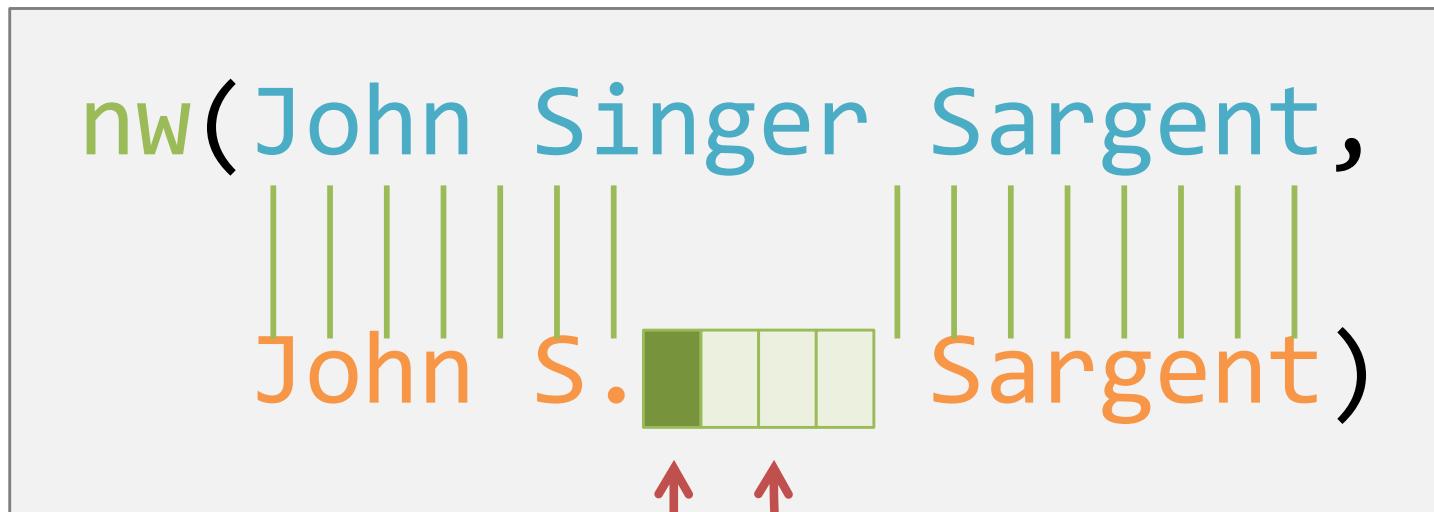
$$+ (-1) * 1 + (-0.5) * 4 = 25$$

$$+ (-1) * 1 + (-0.5) * 8 = 23$$

$\text{nw}(\text{John Stanislaus Sargent},$   
 $\text{John S.} \quad \text{Sargent})$

# Affine Gap Measure

Generalization of needleman-wunch( $x, y$ )

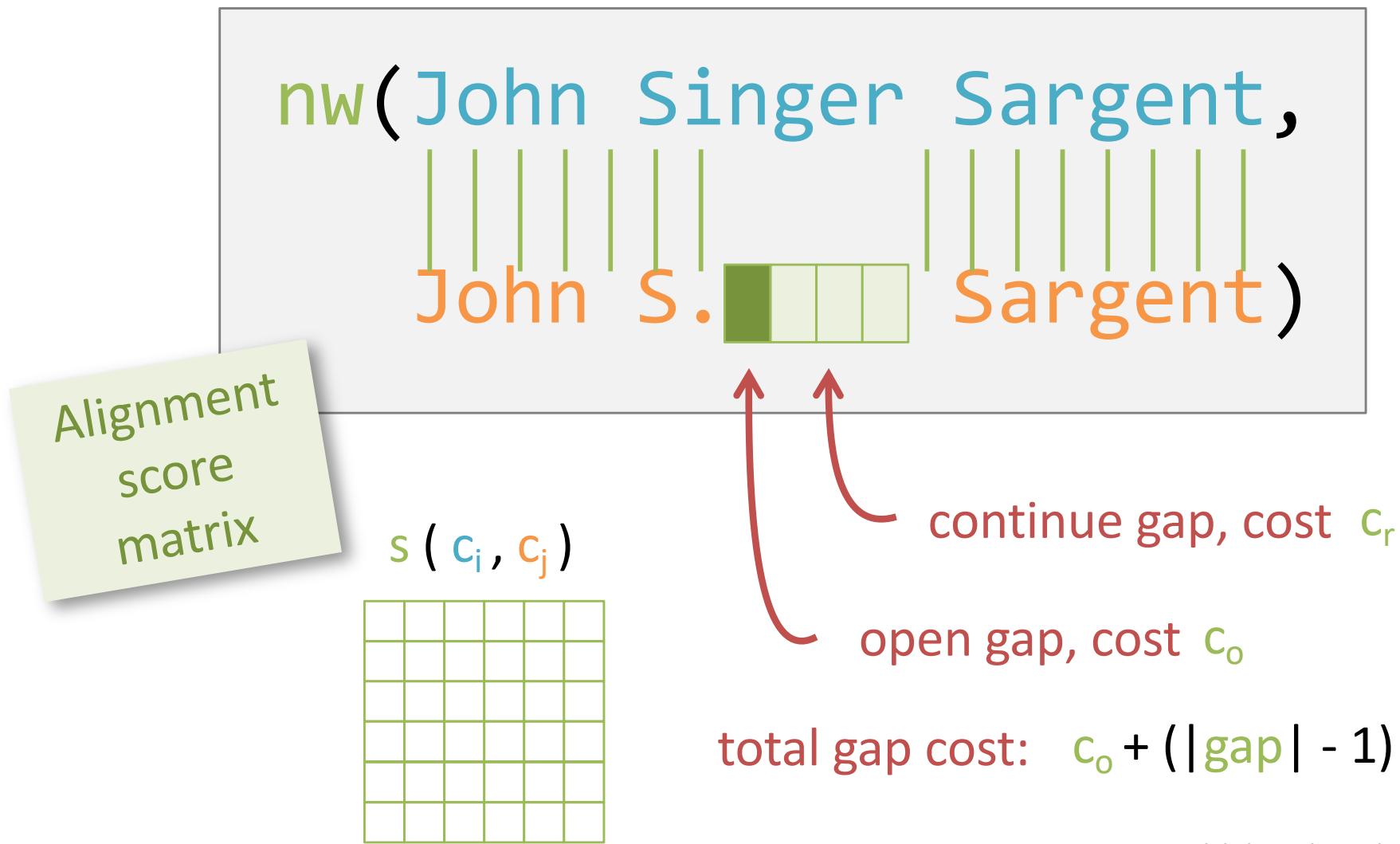


continue gap, cost  $c_r$

open gap, cost  $c_o$

# Affine Gap Measure

Generalization of needleman-wunch( $x, y$ )



# Smith-Waterman

## Global alignment

fully align sequences,  
opening gaps as needed

## Needleman-Wunch

## Local alignment

find best subsequences  
to align

Smith-Waterman: local alignment version of Needleman-Wunch

F	T	F	T	A	L	I	L	L	A	V	A	V
F	-	-	T	A	L	-	L	L	A	-	A	V

F	T	F	T	A	L	I	L	L	A	V	A	V
-	-	F	T	A	L	-	L	L	A	V	-	-

# Smith-Waterman Example

```
match(John Sargent, american painter,  
      American artist John S. Sargent)
```

```
nw(  
      John Sargent, american painter,  
      American artist John S. Sargent  
)
```

Needleman-Wunch: significant gap penalty

# Smith-Waterman Example

```
match(John Sargent, american painter,  
      American artist John S. Sargent)
```

```
nw()
```

Needleman-Wunch: significant gap penalty

# Smith-Waterman Example

```
match(John Sargent, american painter,  
      American artist John S. Sargent)
```

```
nw(John Sargent, american painter,  
      American artist John S. Sargent)
```

Needleman-Wunch: significant gap penalty

```
sw(John Sargent, american painter,  
      American artist John S. Sargent)
```

Smith-Waterman: identifies similar subsequences

# Jaro Similarity Measure

- Get points for having characters in common
  - but only if they are “close by”
- Get points for common characters in the same order
  - lose points for transpositions

# Jaro Similarity Measure $jaro(x, y)$

$$\text{max-distance} = \frac{\max(|x|, |y|)}{2} - 1$$

$x_i$  matches  $y_j$  if  $\begin{cases} x_i = y_j \\ |i - j| \leq \text{max-distance} \end{cases}$

$m$  = number of matching characters

$t$  = number of transpositions  
(of matching characters)

# Jaro Similarity Measure

$$\text{max-distance} = \frac{\max(|x|, |y|)}{2} - 1$$

$m$  = number of matching characters

$t$  = number of transpositions

$$\text{jaro}(x, y) = \begin{cases} 0 & \text{if } m = 0 \\ \frac{1}{3} \left( \frac{m}{|x|} + \frac{m}{|y|} + \frac{m-t}{m} \right) & \text{otherwise} \end{cases}$$

# Jaro Example

`lev(DIXON, DICKSONX) = 4 (4/8 = 0.5)`

`jaro(DIXON, DICKSONX) = ???`

# Jaro Example

	D	I	X	O	N
D	1	0	0	0	0
I	0	1	0	0	0
C	0	0	0	0	0
K	0	0	0	0	0
S	0	0	0	0	0
O	0	0	0	1	0
N	0	0	0	0	1
X	0	0	1	0	0

# Jaro Example

	D	I	X	O	N
D	1	0	0	0	0
I	0	1	0	0	0
C	0	0	0	0	0
K	0	0	0	0	0
S	0	0	0	0	0
O	0	0	0	1	0
N	0	0	0	0	1
X	0	0	0	0	0

$$|x| = 5$$

$$|y| = 8$$

max-distance

$$= (8/2) - 1$$

$$= 3$$

$$m = 4$$

$$t = 0$$

# Jaro Example

	D	I	X	O	N
D	1	0	0	0	0
I	0	1	0	0	0
C	0	0	0	0	0
K	0	0	0	0	0
S	0	0	0	0	0
O	0	0	0	1	0
N	0	0	0	0	1
X	0	0	0	0	0

$$|x| = 5$$

$$|y| = 8$$

$$m = 4$$

$$t = 0$$

$$\frac{1}{3} \left( \frac{m}{|x|} + \frac{m}{|y|} + \frac{m - t}{m} \right)$$

$$\frac{1}{3} \left( \frac{4}{5} + \frac{4}{8} + \frac{4 - 0}{4} \right)$$

$$= 0.767$$

# Jaro-Winkler Measure

Give a bonus if there is a common prefix

```
jwProximity(x,y,boostThreshold,prefixSize)
  = jaro(x,y) <= boostThreshold
  ? jaro (x,y)
  : jaro (x,y)
    + 0.1 * prefixMatch(x,y,prefixSize)
      * (1.0 - jaro(x,y))
```

```
prefixMatch(x,y,prefixSize) =
  min(prefixSize, common-prefix(x,y))
```

```
boostThreshold = 0.7
prefixSize = 4
```

# Jaro-Winkler Measure

```
jwProximity(x,y,boostThreshold,prefixSize)          boostThreshold = 0.7  
  = jaro(x,y) <= boostThreshold                      prefixSize = 4  
  ? jaro (x,y)  
  : jaro (x,y)  
    + 0.1 * prefixMatch(x,y,prefixSize)  
      * (1.0 - jaro(x,y))
```

jaro(DIXON, DICKSONX) = 0.767

$$\begin{aligned}\text{jwProximity(DIXON, DICKSONX)} &= 0.767 + 0.1 \cdot 2 \cdot (1 - 0.767) \\ &= 0.767 + 0.2 \cdot 0.233 \\ &= 0.767 + 0.0466 \\ &= 0.8136\end{aligned}$$

# Set-Based Metrics

# Set-Based Metrics

Generate set of tokens from the strings

Measure similarity between the sets of tokens

# Tokenizing a String

Words

# Tokenizing a String

Words

q-grams: substrings of length q

“david smith” 3-grams  
 $\{\text{##d, #da, dav, avi, ..., h##}\}$

# Jaccard Measure

$B_x = \text{tokens}(x)$

$B_y = \text{tokens}(y)$

$$\text{jaccard}(x, y) = \frac{|B_x \cap B_y|}{|B_x \cup B_y|}$$

`jaccard(dave, dav)`

$B_x = \{\#d, da, av, ve, e\#\}$

$B_y = \{\#d, da, av, v\#\}$

`jaccard(x,y) = 3/6`

# TF/IDF Measure

TF = term frequency

IDF = inverse document frequency

x = Apple Corporation, CA

y = IBM Corporation, CA

z = Apple Corp

...

blah blah Corporation

lev(x, y)



lev(x, z)

???

# TF/IDF Measure

TF = term frequency

IDF = inverse document frequency

x = Apple Corporation, CA

y = IBM Corporation, CA

z = Apple Corp

...

blah blah Corporation

$\text{lev}(x, y) < \text{lev}(x, z)$

... but intuitively (x, z) is a better match

# TF/IDF Measure

TF = term frequency

IDF = inverse document frequency

x = Apple Corporation, CA  
y = IBM Corporation, CA  
z = Apple Corp  
...  
blah blah Corporation



frequent term  
should not carry  
as much weight

$$\text{lev}(x, y) > \text{lev}(x, z)$$

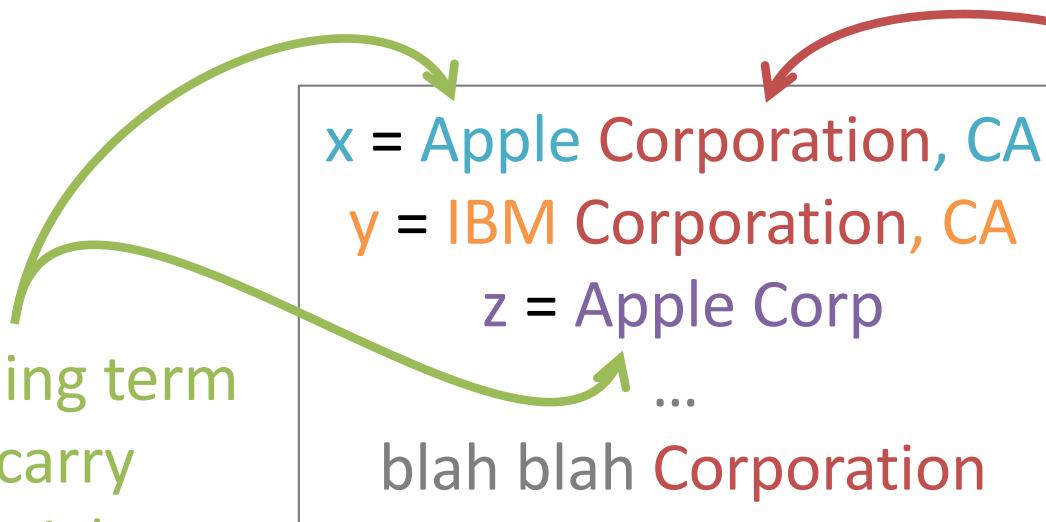
... but intuitively (x, z) is a better match

# TF/IDF Measure

TF = term frequency

IDF = inverse document frequency

distinguishing term  
should carry  
more weight



frequent term  
should not carry  
as much weight

$$\text{lev}(x, y) > \text{lev}(x, z)$$

... but intuitively  $(x, z)$  is a better match

# Term Frequencies and Inverse Document Frequencies

- Assume  $x$  and  $y$  are taken from a collection of strings
- Each string is converted into a bag of terms called a document
- term frequency  $tf(t,d) =$ 
  - number of times term  $t$  appears in document  $d$
- inverse document frequency  $idf(t) =$ 
  - $N / N_d$ , number of documents in collection divided by number of documents that contain  $t$
  - note: in practice,  $idf(t)$  is often defined as  $\log(N / N_d)$

# Example

$$x = aab \rightarrow B_x = \{a, a, b\}$$

$$y = ac \rightarrow B_y = \{a, c\}$$

$$z = a \rightarrow B_z = \{a\}$$

$$tf(a, x) = 2$$

$$tf(b, x) = 1$$

...

$$tf(c, z) = 0$$

$$idf(a) = 3/3 = 1$$

$$idf(b) = 3/1 = 3$$

$$idf(c) = 3/1 = 3$$

# Feature Vectors

- Each document  $d$  is converted into a feature vector  $\mathbf{v}_d$
- $\mathbf{v}_d$  has a feature  $v_d(t)$  for each term  $t$ 
  - value of  $v_d(t)$  is a function of TF and IDF scores
  - here we assume  $v_d(t) = \text{tf}(t,d) * \text{idf}(t)$

$$x = aab \Rightarrow B_x = \{a, a, b\}$$

$$y = ac \Rightarrow B_y = \{a, c\}$$

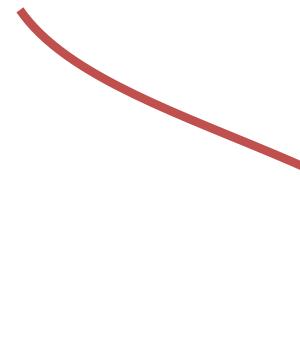
$$z = a \Rightarrow B_z = \{a\}$$

$$\text{tf}(a, x) = 2 \quad \text{idf}(a) = 3/3 = 1$$

$$\text{tf}(b, x) = 1 \quad \text{idf}(b) = 3/1 = 3$$

$$\dots \quad \text{idf}(c) = 3/1 = 3$$

$$\text{tf}(c, z) = 0$$



	a	b	c
$\mathbf{v}_x$	2	3	0
$\mathbf{v}_y$	3	0	3
$\mathbf{v}_z$	3	0	0

# TF/IDF Similarity Score

- Let  $p$  and  $q$  be two strings, and  $T$  be the set of all terms in the collection
- Feature vectors  $v_p$  and  $v_q$  are vectors in the  $|T|$ -dimensional space where each dimension corresponds to a term
- TF/IDF score of  $p$  and  $q$  is the cosine of the angle between  $v_p$  and  $v_q$ 
  - $s(p,q) = \sum_{t \in T} v_p(t) * v_q(t) / [\sqrt{\sum_{t \in T} v_p(t)^2} * \sqrt{\sum_{t \in T} v_q(t)^2}]$

# TF/IDF Similarity Score

- Score is high if strings share many frequent terms
  - terms with high TF scores
- Unless these terms are common in other strings
  - i.e., they have low IDF scores
- Dampening TF and IDF as commonly done in practice
  - use  $v_d(t) = \log(tf(t,d) + 1) * \log(idf(t))$  instead of  
 $v_d(t) = tf(t,d) * idf(t)$
- Normalizing feature vectors
  - $v_d(t) = v_d(t) / \sqrt{\sum_{\{t \in T\}} v_d(t)^2}$

# Hybrid Similarity Measures

# Hybrid Measures

Do the set-based thing  
but  
use a similarity metric for each element of the set

Aple mispelt

x = Apple Corporation, CA  
y = IBM Corporation, CA

z = Aple Corp



...

blah blah Corporation

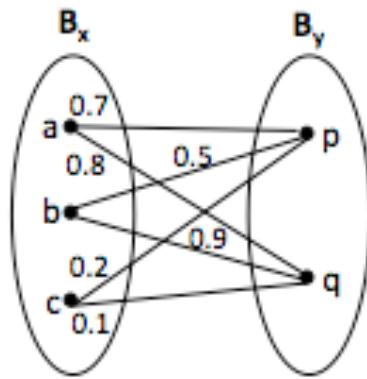
# Generalized Jaccard Measure

- Jaccard measure
  - considers overlapping tokens in both  $x$  and  $y$
  - a token from  $x$  and a token from  $y$  must be identical to be included in the set of overlapping tokens
  - this can be too restrictive in certain cases
- Example:
  - matching taxonomic nodes that describe companies
  - “Energy & Transportation” vs. “Transportation, Energy, & Gas”
  - in theory Jaccard is well suited here, in practice Jaccard may not work well if tokens are commonly misspelled
    - e.g., energy vs. energy
  - generalized Jaccard measure can help such cases

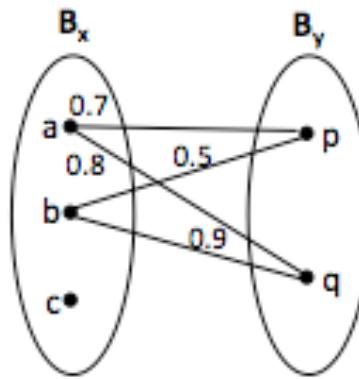
# Generalized Jaccard Measure

- Let  $B_x = \{x_1, \dots, x_n\}$ ,  $B_y = \{y_1, \dots, y_m\}$
- Step 1: find token pairs that will be in the “softened” overlap set
  - apply a similarity measure  $s$  to compute sim score for each pair  $(x_i, y_j)$
  - keep only those score  $>$  a given threshold  $\alpha$ , this forms a bipartite graph  $G$
  - find the maximum-weight matching  $M$  in  $G$
- Step 2: return normalized weight of  $M$  as generalized Jaccard score
  - $GJ(x,y) = \sum_{(x_i,y_j) \in M} s(x_i, y_j) / (|B_x| + |B_y| - |M|)$

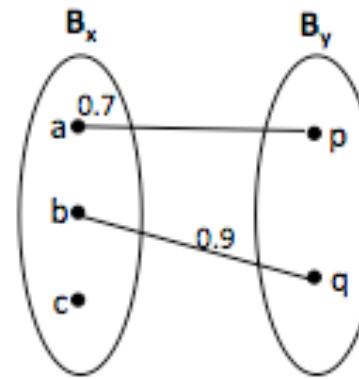
# Generalized Jaccard Example



(a)



(b)



(c)

$$\alpha = 0.5$$

- Generalized Jaccard score:  $(0.7 + 0.9)/(3 + 2 - 2) = 0.53$

# Monge-Elkan Measure

- Break strings  $x$  and  $y$  into multiple substrings
  - $x = A_1 \dots A_n, y = B_1 \dots B_m$
- Compute
  - $s(x,y) = 1/n * \sum_{i=1}^n \max_j \max_{j=1}^m s'(A_i, B_j)$
  - $s'$  is a secondary sim measure, such as Jaro-Winkler
  - Intuitively, we ignore the order of the matching of substrings and only consider the best match for substrings of  $x$  in  $y$

# Monge-Elkan Measure

$$s(x,y) = 1/n * \sum_{i=1}^n \max_j \min s'(A_i, B_j)$$

$$x = A_1 A_2 \quad y = B_1 B_2 B_3$$

$$\frac{\max(s'(A_1, B_1), s'(A_1, B_2), s'(A_1, B_3)) + \max(s'(A_2, B_1), s'(A_2, B_2), s'(A_2, B_3))}{2}$$

2

$s'$  could be any metric, e.g., levenshtein

# Monge-Elkan Measure

$$s(x,y) = 1/n * \sum_{i=1}^n \max_j \min_{s'}(A_i, B_j)$$

x = Comput. Sci. and Eng. Dept., University of California, San Diego

y = Department of Computer Science, Univ. of Calif., San Diego

what s' should we use?

levenshtein

needleman-wunch

affine-gap

smith-waterman

jaro

jaro-winkler

# Phonetic Similarity Measures

# Phonetic Similarity Measures

- Match strings based on their sound, instead of appearances
- Very effective in matching names, which often appear in different ways that sound the same
  - e.g., Meyer, Meier, and Mire; Smith, Smithe, and Smythe
- Soundex is most commonly used

# The Soundex Measure

- Used primarily to match surnames
  - maps a surname  $x$  into a 4-letter code
  - two surnames are judged similar if share the same code
- Algorithm to map  $x$  into a code:
  - Step 1: keep the first letter of  $x$ , subsequent steps are performed on the rest of  $x$
  - Step 2: remove all occurrences of W and H. Replace the remaining letters with digits as follows:
    - ❖ replace B, F, P, V with 1, C, G, J, K, Q, S, X, Z with 2, D, T with 3, L with 4, M, N with 5, R with 6
  - Step 3: replace sequence of identical digits by the digit itself
  - Step 4: Drop all non-digit letters, return the first four letters as the soundex code

# The Soundex Measure

- Example: x = Ashcraft
  - after Step 2: A226a13, after Step 3: A26a13, Step 4 converts this into A2613, then returns A261
  - Soundex code is padded with 0 if there is not enough digits
- Example: Robert and Rupert map into R163
- Soundex fails to map Gough and Goff, and Jawornicki and Yavornitzky
  - designed primarily for Caucasian names, but found to work well for names of many different origins
  - does not work well for names of East Asian origins
    - ❖ which uses vowels to discriminate, Soundex ignores vowels

# Other Readings

- [http://en.wikipedia.org/wiki/String\\_metric](http://en.wikipedia.org/wiki/String_metric)
- [http://en.wikipedia.org/wiki/Approximate\\_string\\_matching](http://en.wikipedia.org/wiki/Approximate_string_matching)
- [http://en.wikipedia.org/wiki/Edit\\_distance](http://en.wikipedia.org/wiki/Edit_distance)
- [http://en.wikipedia.org/wiki/Levenshtein\\_distance](http://en.wikipedia.org/wiki/Levenshtein_distance)
- [http://en.wikipedia.org/wiki/Jaro–Winkler\\_distance](http://en.wikipedia.org/wiki/Jaro–Winkler_distance)
- [http://en.wikipedia.org/wiki/Smith–Waterman\\_algorithm](http://en.wikipedia.org/wiki/Smith–Waterman_algorithm)
- [http://en.wikipedia.org/wiki/Jaccard\\_index](http://en.wikipedia.org/wiki/Jaccard_index)
  
- <http://alias-i.com/lingpipe/demos/tutorial/stringCompare/read-me.html>
- <https://web.archive.org/web/20160303200731/http://www.gettingcirrius.com/2011/01/calculating-similarity-part-2-jaccard.html>
- [http://en.wikipedia.org/wiki/Sequence\\_alignment#Pairwise\\_alignment](http://en.wikipedia.org/wiki/Sequence_alignment#Pairwise_alignment)

# Software

- <http://code.google.com/p/javasimilarity/source/browse/trunk/simmetrics/src/main/java/uk/ac/shef/wt/simmetrics/>
- <http://sourceforge.net/projects/secondstring/>
- <http://planetcalc.com/1721/> (Levenshtein calculator)

Source path: [svn/](#) [trunk/](#) [simmetrics/](#) [src/](#) [main/](#) [java/](#) [uk/](#) [ac/](#) [shef/](#) [wit/](#) [simmetrics](#)

Directories	Filename
simmetrics	<a href="#">AbstractStringMetric.java</a>
arbitrators	<a href="#">BlockDistance.java</a>
basiccontainers	<a href="#">ChapmanLengthDeviation.java</a>
math	<a href="#">ChapmanMatchingSoundex.java</a>
metrichandlers	<a href="#">ChapmanMeanLength.java</a>
similaritymetrics	<a href="#">ChapmanOrderedNameCompoundSimilarity.java</a>
costfunctions	<a href="#">CosineSimilarity.java</a>
task	<a href="#">DiceSimilarity.java</a>
tokenisers	<a href="#">EuclideanDistance.java</a>
utils	<a href="#">InterfaceStringMetric.java</a>
wordhandlers	<a href="#">JaccardSimilarity.java</a>
	<a href="#">Jaro.java</a>
	<a href="#">JaroWinkler.java</a>
	<a href="#">Levenshtein.java</a>
	<a href="#">MatchingCoefficient.java</a>
	<a href="#">MongeElkan.java</a>
	<a href="#">NeedlemanWunch.java</a>
	<a href="#">OverlapCoefficient.java</a>
	<a href="#">QGramsDistance.java</a>