



INTEGRATION

Summer 2022
Data Integration

Thorsten Papenbrock

Duplicate: Two representations (e.g. records, objects, XML structures) that represent the same real-world entity or concept

- Example:

Name	Street	Number
Ernie	Sesamstr.	2
Fienchen	Sesamstr.	1
Bert	Sesamstr.	2
Samson	Sesamstr.	1
Tiffy	Sesamstr.	3
Kermit	Sesamstr.	6
Grobi	Sesamstr.	4
Krümelmonster	Sesamstr.	8
Mumpitz	Sesamstr.	7
Oscar	Sesamstr.	5
Bibo	Sesamstr.	9
Graf Zahl	Sesamstr.	9¾
Kermitt	Sesamstr.	6
Finchen	Sesamstr.	1
Rumpel	Sesamstr.	1.5

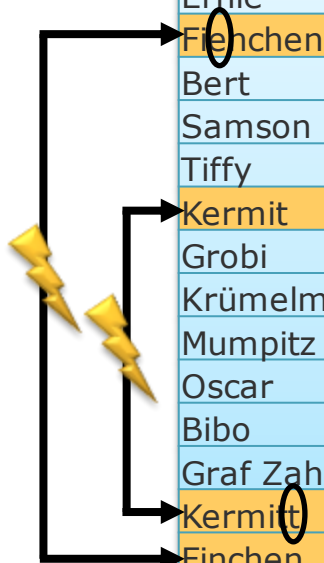
Data Integration

Entity Resolution

Thorsten Papenbrock
Slide 2

Duplicate: Two representations (e.g. records, objects, XML structures) that represent the same real-world entity or concept

- Example:



Name	Street	Number
Ernie	Sesamstr.	2
Finchen	Sesamstr.	1
Bert	Sesamstr.	2
Samson	Sesamstr.	1
Tiffy	Sesamstr.	3
Kermit	Sesamstr.	6
Grobi	Sesamstr.	4
Krümelmonster	Sesamstr.	8
Mumpitz	Sesamstr.	7
Oscar	Sesamstr.	5
Bibo	Sesamstr.	9
Graf Zahl	Sesamstr.	9¾
Kermitt	Sesamstr.	6
Finchen	Sesamstr.	1
Rumpel	Sesamstr.	1.5

Data Integration

Entity Resolution

ThorstenPapenbrock
Slide 3

Entity Resolution

Duplicates

Duplicate: Two representations (e.g. records, objects, XML structures) that represent the same real-world entity or concept

- Typically not equal but similar (w.r.t. values, attribute sets, references)
- Example:

QWMQ0071368
Dr Felix Naumann
72 A R.-Breitscheid-Str
Potsdam
14482
GERMANY

QWMX0071362
Felix Naumann
Rudolf-Breitscheid-Str 72A
Potsdam
14482
GERMANY

- Origin examples:
 - Data integration: data silos, data sharing, data discovery
 - Errors in data entry: typos, redundancies, transmission
 - Fraudulent actions: manipulation, human mistakes

Data Integration

Entity Resolution

Duplicate: Two representations (e.g. records, objects, XML structures) that represent the same real-world entity or concept

- Typically not equal but similar (w.r.t. values, attribute sets, references)
- Effects:
 - **Wrong decision making**
 - Actions are carried out multiple times per entity
 - **Inaccurate statistics**
 - Number of entities is counted too high
 - **Poorly trained machine learning models**
 - Duplicate entries introduce bias and misclassifications
 - **Software failures**
 - Unique constraints may be violated
 - ...

All this costs a lot of money in practice!

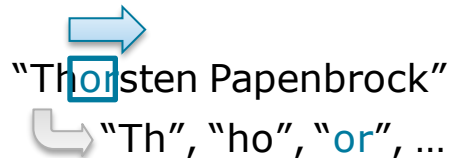
Data Integration

Entity Resolution

Thorsten Papenbrock
Slide **5**

k

- The function $to(x)$ splits the string x into short **substrings of length n** by sliding a window of size n over x ; every slide creates an n -gram.
 - $n=2$: Bigrams; $n=3$: Trigrams
 - Number of n -grams = $|x| - n + 1$
- Variation 1: **Pad with $n - 1$ special characters**
 - Emphasizes beginning and end of string
- Variation 2: **Include positional information**
 - Useful to weight tokens by their positions


 “Th” “or”, “ho”, “or”, ...

String	Bigrams	Padded bigrams	Positional bigrams	Trigrams
gail	ga, ai, il	\odot g, ga, ai, il, \otimes	(ga,1), (ai,2), (il,3)	gai, ail
gayle	ga, ay, yl, le	\odot g, ga, ay, yl, le, $e\otimes$	(ga,1), (ay,2), (yl,3), (le,4)	gay, ayl, yle
peter	pe, et, te, er	\odot p, pe, et, te, er, $r\otimes$	(pe,1), (et,2), (te,3), (er,4)	pet, ete, ter
pedro	pe, ed, dr, ro	\odot p, pe, ed, dr, ro, $o\otimes$	(pe,1), (ed,2), (dr,3), (ro,4)	ped, edr, dro

Data Integration

Data Matching

 ThorstenPapenbrock
 Slide **6**

String Similarity

Token-based Similarity: n-grams



n-grams "levenshtein"



NATURAL LANGUAGE



MATH INPUT



EXTENDED KEYBOARD



EXAMPLES



UPLOAD



RANDOM

Assuming the input is referring to string encodings | Use "n" as a [variable](#) instead

Input interpretation

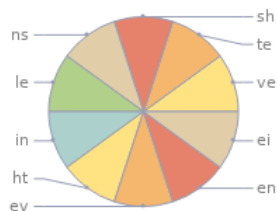
n-grams levenshtein

Character-level bigrams

le | ev | ve | en | ns | sh | ht | te | ei | in

Character-level bigrams frequency pie chart

Show bar chart



Data Integration
Data Matching

ThorstenPapenbrock
Slide **7**

String Similarity

Levenshtein Distance

This is the challenge!

Definition:

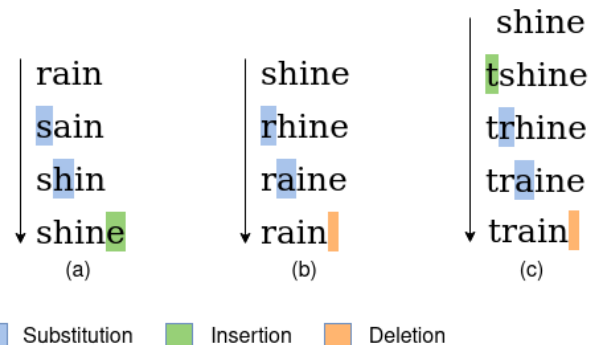
- $dist_{levenshtein}(x, y)$ = minimum number of edits (insert, delete, replace) that transform the string x into string y
- The most popular metric for describing the edit-distance of strings

Operations:

1. **insert** a character into the string
2. **delete** a character from the string
3. **replace** (substitute) a character with a different character

Examples:

- $dist_{levenshtein}(\text{'table'}, \text{'cable'}) = 1$
(1 replace)
- $dist_{levenshtein}(\text{'Thorsten Papenbrock'}, \text{' Papenbrock, Thorsten'}) = 19$
(9 deletes + 10 inserts)



String Similarity

Levenshtein Distance

Calculation:

- Calculating the minimum edit-distance is a case for **dynamic programming**
- Optimality principle:
 - Any minimum edit-distance of two substrings must be part of the best overall solution.
- Dynamic programming algorithm:
 1. Initialize a matrix M of size $(|x|+1) \times (|y|+1)$
 2. Fill matrix: $M_{i,0} = i$ and $M_{0,j} = j$
 3. Recursion:
$$M_{i,j} = \begin{cases} M_{i-1,j-1} & \text{if } x[i] = y[j] \\ 1 + \min(M_{i-1,j}, M_{i,j-1}, M_{i-1,j-1}) & \text{els} \end{cases}$$
 4. Distance: $dis_{levens}(x, y) = M_{|x|,|y|}$

Data Integration
Data Matching

Thorsten Papenbrock
Slide 9

String Similarity

Levenshtein Distance

Calculation example:

		J	O	N	E	S
	0	1	2	3	4	5
J	1					
O	2					
H	3					
N	4					
S	5					
O	6					
N	7					

		J	O	N	E	S
	0	1	2	3	4	5
J	1	0	1	2	3	4
O	2					
H	3					
N	4					
S	5					
O	6					
N	7					

		J	O	N	E	S
	0	1	2	3	4	5
J	1	0	1	2	3	4
O	2	1	0	1	2	3
H	3	2	1	1	2	3
N	4	3	2	1	2	3
S	5	4	3	2	2	2
O	6	5	4	3	3	3
N	7	6	5	4	4	4

Every path leads to the same optimal solution!

$$M_{i,0} = i \text{ and } M_{0,j} = j$$

$$M_{i,j} = \begin{cases} M_{i-1,j-1} & \text{if } x[i] = y[j] \\ 1 + \min(M_{i-1,j}, M_{i,j-1}, M_{i-1,j-1}) & \text{els} \end{cases}$$

Data Integration
Data Matching

ThorstenPapenbrock
Slide **10**

String Similarity

Levenshtein Distance

Similarity:

- $sim_{levenshtein}(x, y) = 1 - dist_{levenshtein}(x, y) / \max(|x|, |y|)$

x	y	Lev. Distance	Lev. Similarity
Jones	Johnson	4	0.43
Paul	Pual	2	0.5
Paul Jones	Jones, Paul	11	0

Discussion:

- Robust against different forms and positions of typos
- Easy to understand and implement
- Unsupervised similarity metric
- Sensitive to word order changes
- Expensive to calculate

Data Integration
Data Matching

ThorstenPapenbrock
Slide **11**

String Similarity

Levenshtein Distance

Complexity:

- Time: $O(|x| \cdot |y|)$ (calculate entire matrix)
- Space: $O(\min(|x|, |y|))$ (store only two rows of the matrix)

Properties:

- $0 \leq \text{dist}_{\text{levenshtein}}(x, y) \leq \max(|x|, |y|)$
- $||x| - |y|| \leq \text{dist}_{\text{levenshtein}}(x, y)$

Good lower-bound estimate
to possibly skip the exact
distance computation!

Cost models (extension): Assign different costs to edit-operations.

- By operation:
 - E.g. replace costs 0.5 but insert/delete cost 1.0 to punish string length changes.
- By character:
 - E.g. OCR ($m \simeq n$, $1 \simeq l$) or keyboard ($a \simeq s$) or brain ($6 \simeq 9$) or biology ($a \simeq t$)

Data Integration
Data Matching

ThorstenPapenbrock
Slide **12**

Sorted Neighborhood Method (SNM)

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
1																				
2																				
3																				
4																				
5																				
6																				
7																				
8																				
9																				
10																				
11																				
12																				
13																				
14																				
15																				
16																				
17																				
18																				
19																				
20																				

Sorted Neighborhood Method

- Sort tuples so that similar tuples are close to each other.
- Compare tuples only within a small neighborhood (= window).

Sorted Neighborhood Method (SNM)

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
1																				
2																				
3																				
4																				
5																				
6																				
7																				
8																				
9																				
10																				
11																				
12																				
13																				
14																				
15																				
16																				
17																				
18																				
19																				
20																				

Sorted Neighborhood Method

- Sort tuples so that similar tuples are close to each other.
 - Compare tuples only within a small neighborhood (= window).
1. **Generate a key**
 - E.g. SSN
 - E.g. Name[1-3] + Age + ...
 2. **Sort entire relation by the key**
 - Sort records physically or create a sorted index.
 3. **Slide a window over sorted tuples**
 - Compare all pairs of tuples within the sliding window.

Sorted Neighborhood Method (SNM)

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
1																				
2																				
3																				
4																				
5																				
6																				
7																				
8																				
9																				
10																				
11																				
12																				
13																				
14																				
15																				
16																				
17																				
18																				
19																				
20																				

Sorted Neighborhood algorithm:

// Sorting

```
Arrays.sort(records, (r1, r2) -> r1[6].compareTo(r2[6]));
```

Key here is simply
attribute 6

"w" is the
window size

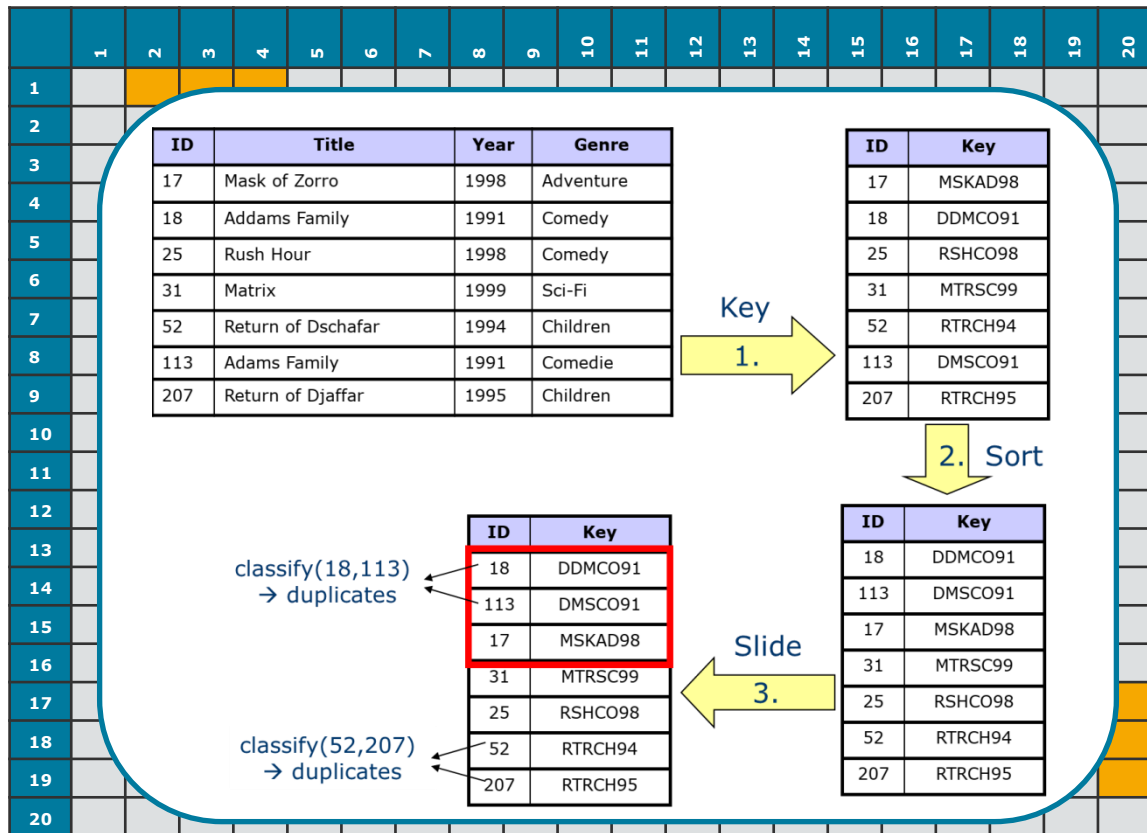
// Comparing

```
for (int i = 0; i < n; i++)
  for (int j = i+1; j < min(i+w, n); j++)
    match(records[i], records[j]);
```

Sorted Neighborhood Method

- Sort tuples so that similar tuples are close to each other.
 - Compare tuples only within a small neighborhood (= window).
1. **Generate a key**
 - E.g. SSN
 - E.g. Name[1-3] + Age + ...
 2. **Sort entire relation by the key**
 - Sort records physically or create a sorted index.
 3. **Slide a window over sorted tuples**
 - Compare all pairs of tuples within the sliding window.

Sorted Neighborhood Method (SNM)



Sorted Neighborhood Method

- Sort tuples so that similar tuples are close to each other.
 - Compare tuples only within a small neighborhood (= window).
- Generate a key**
 - E.g. SSN
 - E.g. Name[1-3] + Age + ...
 - Sort entire relation by the key**
 - Sort records physically or create a sorted index.
 - Slide a window over sorted tuples**
 - Compare all pairs of tuples within the sliding window.

Sorted Neighborhood Method (SNM)

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
1																				
2																				
3																				
4																				
5																				
6																				
7																				
8																				
9																				
10																				
11																				
12																				
13																				
14																				
15																				
16																				
17																				
18																				
19																				
20																				

Sorted Neighborhood complexity

N: Number of tuples

w: Window size

Three steps: create key, sort, window

Computational complexity:

- $O(N) + O(N \log N) + O(w N)$
 - $= O(N \log N)$ if $w < \log N$
 - $= O(w N)$ else

I/O complexity

- Linear in N
- Three passes over table on disk
- Sorting with e.g. TPMMS in $O(3N)$

Sorted Neighborhood Method

- Sort tuples so that similar tuples are close to each other.
 - Compare tuples only within a small neighborhood (= window).
1. **Generate a key**
 - E.g. SSN
 - E.g. Name[1-3] + Age + ...
 2. **Sort entire relation by the key**
 - Sort records physically or create a sorted index.
 3. **Slide a window over sorted tuples**
 - Compare all pairs of tuples within the sliding window.