



Data deduplication

Robert Wrembel **Poznan University of Technology** Faculty of Computing and Telecommunications, Institute of Computing Science **Interdisciplinary Centre for Artificial Intelligence and Cybersecurity** Poznań, Poland



© R.Wrembel (PUT and CAICS, Poland)



IT landscape in a FI



3

- Dozens to over a hundred of data repositories to integrate
 - on premise databases from all software houses (mainly relational)
 - xls files
 - csv files
 - streaming sources
 - ...
- Thousands to dozens of thousands of integration (ETL) processes

© R.Wrembel (PUT and CAICS, Poland)



© R.Wrembel (PUT and CAICS, Poland)



Motivating facts



- Sources of duplicate data
 - FI acquisition
 - different bank products require separate customer records
 - imperfection of software used
 - data errors
- Duplicated data cause
 - loss of money
 - deterioration of FI reputation

© R.Wrembel (PUT and CAICS, Poland)





© R.Wrembel (PUT and CAICS, Poland)



8



Deduplication: introduction



- Before comparing, records should be cleaned
 - homogenizing values (abbreviations, units of measurement, symbols, ...)
 - no special signs, no punctuations
 - no abbreviations
- Problem: how to decide if 2 records represent the same entity?
 - [Wrembel, Robert, ul. Matejki, Poznań]
 - [Wrębel, Robert, ul. Matejki, Poznań]
- Case 1: natural identifiers (e.g., ID, SSN, PESEL, email, mobile#) available
 - but email, mobile# may change in time for the same person
- Case 2: no natural identifiers available
 - approximate/probabilistic decision based on a similarity measure

© R.Wrembel (PUT and CAICS, Poland)





How to compare records?

- the worst case: each entity compared to all the other entities
 - $O(n^2) \rightarrow deduplication$
 - O(n*m) → record linkeage
- problem: efficiency



Deduplication complexity



11

Deduplication complexity: O(n²)



Bank example: over 20mln customer records

- 2*10⁷ * 2*10⁷ = 4*10¹⁴ comparisons
- one comparison = 10⁻⁹s → comparing the whole set: 4*10⁵s = 4.6 days → must be optimized

© R.Wrembel (PUT and CAICS, Poland)





Deduplication pipeline



13

Improving performance

- avoiding N² comparisons of records
- finding groups of similar records
- Base-line deduplication pipeline (BLDDP)



- G. Papadakis, L. Tsekouras, E. Thanos, G. Giannakopoulos, T. Palpanas, M. Koubarakis: Domain- and Structure-Agnostic End-to-End Entity Resolution with JedAI. SIGMOD Record, Vol. 48, No. 4, 2019
- G. Papadakis, D. Skoutas, E. Thanos, T. Palpanas: Blocking and Filtering Techniques for Entity Resolution: A Survey. ACM Computing Surveys, 52:(3), 2021





© R.Wrembel (PUT and CAICS, Poland)



- values of all attributes are compared
- Challenge: selecting the right blocking key(s)

© R.Wrembel (PUT and CAICS, Poland)



Blocking



17

Canopy clustering

- computationally cheap algorithm for pre-clustering for more accurate algorithms (k-means, dbscan)
- two distances are used
 - T_L loose distance
 - T_T tight distance
 - T_L>T_T
- 1. randomly select a center point pⁱ_C and create canopy C_i where pⁱ_C is the center point of C_i
- 2. assign every point to C_i if the distance between a given point and $p^i{}_{\rm C} < T_L$
 - points within circles of radius $T_{\ensuremath{\mathsf{T}}}$ cannot be center points in the next iteration
- 3. repeat step 1-2 until there are no more data points

© R.Wrembel (PUT and CAICS, Poland)

<image><image><image><section-header><section-header>



Blocking



19

Hashing

- on some attributes (hash keys) → O(n)
- drawback: records must have identical values of hash keys to hash into the same bucket → applicable for exact matching

| © R.Wrembel (PUT and CAICS, Poland) | | |
|-------------------------------------|--|--|
| | | |



recordM

step2

© R.Wrembel (PUT and CAICS, Poland)

recordM

step1

20

recordM

step3



Blocking



21

■ Token-based → redundancy positive blocks

- r1: {name: Robert Wrembel, degree: prof.}
- r2: {name: Robert Wrębel, degree: professor}
- r3: {name: Witold Andrzejewski, degree: dr}
- r4: {name: Witold Andrzejewski, degree: dr inż.}
- r5: {name: Paweł Boiński, degree: dr}
- r6: {name: Bartosz Bębel, degree: dr}

block fname (Robert): {r1, r2} block fname (Witold): {r3, r4} block fname (W): {r4}

block Iname (Andrzejewski): {r3, r4} block degre (dr): {r3, r5, r6}

© R.Wrembel (PUT and CAICS, Poland)

Blocking bi-gram example Ρ i. e r r e e Pi Fe ie er -------> er rr ·····> rr ro re n-gram based blocking group similar records together based on n-grams they share records that share a parameterized minimum number of ngrams are grouped into the same block n-grams can be organized into an inverted index all records that contain the same token in their blocking key reside in the same inverted index list



- Goal: to minimize the number of comparisons, i.e., to discard unnecessary comparisons:
 - between blocks (in case of overlapping blocks)
 - between records in a given block

© R.Wrembel (PUT and CAICS, Poland)

© R.Wrembel (PUT and CAICS, Poland)







Duplicate propagation

- as blocks are overlapping, the same objects are compared more than once → avoiding this by propagating the identified matches to the subsequently processed blocks
- a central hash table contains all the matches detected so far
- before comparing a pair of objects, check whether any of them is registered in the hash table → if true for at least one of them then skip the comparison

| • | G. Papadakis, E. Ioannou, C. Niederée, P. Fankhauser: Efficient entity resolution for large heterogeneous informatio spaces. WSDM, 2011 |
|---|---|
| | |

© R.Wrembel (PUT and CAICS, Poland)







Block filtering

- importance of a block for a given object o_i = maximum number of blocks o_i resides in \rightarrow count block assignments per object
- remove o_i from the least important block(s)

| G. Papadakis, G. Papastefanatos, T. Palpanas, M. Koubarakis: Scaling Entity Resolution to Large, Het with Enhanced Meta-blocking. EDBT, 2016 | erogeneous Data |
|--|-----------------|
| © R.Wrembel (PUT and CAICS, Poland) | 27 |



- Split large blocks into smaller (below max allowed size)
 [1]
- Merge blocks with similar blocking keys (similarity threshold) [2]

| J. Fisher, P. Christen, Q. Wang, E. Rahm: A Clustering-Based Framework to Control Block Sizes for Entity Resolution KDD, 2015 | 1. | A. Das Sarma, A. Jain, A. Machanavajjhala, P. Bohannon: An automatic blocking mechanism for large-scale de- duplication tasks. CIKM, 2012 |
|---|----|--|
| | 2. | J. Fisher, P. Christen, Q. Wang, E. Rahm: A Clustering-Based Framework to Control Block Sizes for Entity Resolution. KDD, 2015 |





- Size-based block clustering
 - strategies
 - merging small blocks that correspond to similar blocking keys
 - splitting large blocks into smaller ones
 - to balance block sizes → balancing parallel processing of record matchings in blocks

| J. Fisher, P. Christen, Q. Wang, E. Rahm: A Clustering-Based Framework to Control Block Sizes for Entity Resolution KDD, 2015 | |
|---|----|
| © R.Wrembel (PUT and CAICS, Poland) | 29 |



Iterative blocking

- whenever a new pair of duplicates is detected ($r_{i_{\prime}}r_{m}$) their descriptions are merged $\rightarrow r_{im}$
- r_{im} replaces (r_i,r_m) in all blocks
- compare r_{im} to other records in its blocks and merge if matching
- repeat until no matches can be found

| - | S.E. Whang, D. Menestrina, G. Koutrika, M. Theobald, H. Garcia-Molina: Entity resolution with iterative blocking. SIGMOD, 2009 | |
|---|---|--|
| | | |





Meta-blocking

- uses a graph to represent comparisons
- eliminates the same comparisons in multiple blocks
- uses labels of graph edges to eliminate comparisons below certain threshold
 - methods for computing the values of the labels

| G. Papadakis, G. Koutrika, T. Palpanas, W. Nejdl: Meta-Blocking: Taking Entity Resolution to the Next Level. IEEE Trans. Knowl. Data Eng. 26(8), 2014 | |
|---|----|
| © R.Wrembel (PUT and CAICS, Poland) | 31 |

















Similarity measures for text data



- The longest common sub-sequence or the longest common sub-string
- Vector representation in an m-dimensional space
 e.g., TFIDF, Cosine
- Compression techniques
 - e.g., BZ2, Lempel-Ziv-Mrkov, ZLib
- Ensemble of measures
 - e.g., Monge-Elkan + Damerau-Levenshtein + n-gram, Cosine+n-gram

| © R.Wrembel (PUT | and CAICS, Poland) |
|------------------|--------------------|
|------------------|--------------------|



Example



39

 Different similarity measures produce different values for the same compared strings

| r1 | Robert | Wrembel | Wyspiańskiego | | |
|--------------------|------------------|---------|---------------|--|--|
| r2 | Robret | Wręble | Wyspiankiego | | |
| Jaro-Winkler 0.961 | | 0.826 | 0.953 | | |
| Levenshtein | evenshtein 0.667 | | 0.846 | | |
| Overlap | 1.000 | 0.833 | 0.917 | | |

 $\label{eq:avg} \begin{array}{l} \text{AVG}(\text{simJW}(r1,r2)) = (0.961 + 0.826 + 0.953)/3 = 0.913 \xleftarrow{} \text{similar} \\ \text{AVG}(\text{simLe}(r1,r2)) = (0.667 + 0.429 + 0.846)/3 = 0.647 \xleftarrow{} \text{probably similar} \\ \text{AVG}(\text{simOv}(r1,r2)) = (1 + 0.833 + 0.917)/3 = 0.916 \xleftarrow{} \text{similar} \\ \end{array}$



Example



41

| are they similar? | | | | | | |
|------------------------------------|---------------------------|--|--|--|--|--|
| B. CHROBREGO < | >BOLERSŁAWA CHROBREGO | | | | | |
| B CHROBREGO | >BOLESLAWA CHROBREGO | | | | | |
| B CHROBREGO 10 ←−−−−−−−−−− | >BOLESŁ. CHROBREGO | | | | | |
| B CHROBREGO 42B | BOLESŁAW CHROBREGO | | | | | |
| B,CHROBREGO | BOLESŁAWA CHROBREGO | | | | | |
| B.CHROBREGO | BOLESŁAWA CHROBREGO | | | | | |
| B.CHROBREGO 33/72 | BOLESŁ.CHROBREGO | | | | | |
| B.CHROBREGO OŁOBOK | BOLESŁWA CHROBREGO | | | | | |
| B.CHROBREGO SKR.60 | BOOESŁAWA CHROBREGO | | | | | |
| BOL. CHROBREGO | CHROBREGO | | | | | |
| BOL CHROBREGO | CHROBREGO 10 | | | | | |
| .BOL.CHROBREGO | CHROBREGO 22A/6 | | | | | |
| BOL.CHROBREGO | CHROBREGO KOWALEW | | | | | |
| BOL.CHROBREGO OŁOBOK | CHROBREGO PAŃSTWOWY DOM D | | | | | |

© R.Wrembel (PUT and CAICS, Poland)



© R.Wrembel (PUT and CAICS, Poland)



Semantic relationships augmenting similarities

- Context-based similarity (semantic relationships)
 - data source S1 (healthcare): object patient100
 - data source S2 (healthcare): object person131
 - data source S3 (banking): object customer456
 - patient100 more similar to person131 since both exist in the same healthcare context





Entity clustering



45

Multiple clustering algorithms

use the similarity measure between records



© R.Wrembel (PUT and CAICS, Poland)



Merge semantically identical records in a cluster into one final augmented record

| TaxpayerID | Fnan | ne | Lname | Counti | ry Titl | e | Domain | | | |
|------------|--------|---------|----------|-------------|-----------|-----------|------------|---------|-------|---------------|
| 132018 | Robe | ert \ | Wrembel | Polan | id Pro | f. data v | vrehouses | | | |
| | + | | | | | | | | | |
| TaxpayerID | 1sName | e 2no | dName | Last_Name | Education | | City B | Born | | |
| 132018 | Robert | t A | ndrew | Wrembel | univ. | Po | znań 1 | .968 | | |
| | | | | | | | | | | |
| TaxpayerID | 1sName | 2ndName | Last_Nam | e Education | City | Born | TaxpayerID | Country | Title | Doma |
| 132018 | Robert | Andrew | Wrembe | el univ. | Poznań | 1968 | 132018 | Poland | Prof. | data wrehouse |



- 35280+ combinations of these algorithms
- searching the full space of alg. combinations is impossible
- techniques for pruning search space
 - some methods may be better for deduplicating personal data, some for bibliographical data, some for addresses
 - knowledge of a domain expert is crucial
- an automatic approach does not exist

© R.Wrembel (PUT and CAICS, Poland)





ML for deduplication



49

Problem

- large set of training data is needed
- labeling record pairs is a time consuming task
- Standard solution → classification
 - classes: duplicates (T), non-duplicates (F)
 - classes: duplicates (T), probable-duplicates (P), nonduplicates (F)









ML for deduplication (1)



- HeALER
- Active learning approach
- Committee of heterogeneous classifiers
 - SVM, one-vs-rest logistic regression, logistic regression, decision tree
- Data sets
 - ACM-DBLP (4800+ records)
 - GScholar-DBLP (66000+ records)







ML for deduplication (2)



Characteristics

- active learning approach
- integrated matcher and blocker
- deep neural network used to build matcher and blocker
- Blocker creates blocks of possible duplicates → detects duplicates based on embedding similarity







- Matcher uses transformer-based pretrained language models (TPLM, e.g., BERT, RoBERTa)
- For each pair of records (r, s) Matcher assigns a probability of the pair representing a duplicate
- Matcher uses Transformer to get a joint embedding E(r, s) of pair (r, s)



ML for deduplication (2)



57

 Transformer receives concatenated tokens of record pairs

- [start], r₁...r_n, [separator], s₁...s_m, [separator]
- multiple encoders enc₁, . . . enc_m for creating embeddings
- transformer assigns for each token an embedding that represents its semantics in the context of the current record
- the contextual embeddings can eliminate problems with spelling mistakes and different abbreviations
- embeddings are indexed for faster retrieval





ML for deduplication (2)



59

Data sets in experiments

- Walmart, Amazon, Google
- DBLP, ACM, GScholar
- max. number of records: ~67000

© R.Wrembel (PUT and CAICS, Poland)



Goal

- to select the most suitable classification algorithms (classes T, N)
- to tune hyperparameters
 - KNN: the number of neighbors
 - SVM: kernel (linear, polynomial)
 - random forest: number of random features to sample at each split
 point







ML for deduplication (3)



Embedder

- encodes a token sequence into a multi-dimensional vector \rightarrow embedding
- typically, a pre-trained word embedding NN is used
 - in the paper the following embedders were applied: Bert, DistilBert, Albert, Roberta, XLNET

Combiner

 as multiple embeddings can be generated for the same pair of input records, they must be aggregated into a single multi-dimensional vector



© R.Wrembel (PUT and CAICS, Poland)



Data sets in experiments

- Walmart, Amazon, Google
- DBLP, ACM, GScholar
- max. number of records: ~28000



ML for deduplication: challenges



65

C5: how to construct a learning data set of a reasonable size for 20M customers' database?

- manual tagging of 1000 pairs of rows
 - 3 persons, including a bank expert
 - tagging throughput: 1.33 rec/min
 - 22.2h for tagging 1000 rows
- conclusion: impossible to construct a learning data set manually

| © R.Wrembel (PUT a | and CAICS, Poland) |
|--------------------|--------------------|
|--------------------|--------------------|





Requirements = reality



- A project must finish within a given monetary budget and time
- Processed data must be accessible either from a relational database or from files (csv, xls)
- Developed algorithms and models must be intuitive, understandable and easy to implement by a technical IT staff at a company
- Ideally the developed methods should be based on outof-the-box software components

| © R.Wrembel (Pl | JT and CAICS, Poland) |
|-----------------|-----------------------|
|-----------------|-----------------------|







- The developed methods must be efficient → they will be applied to several millions of rows
- The solution must be deployable in the IT architecture used by a company (specific hardware and software) → either in a database or in a standard data science environment
- Only licensed/certified software can be used (especially in financial institutions)

© R.Wrembel (PUT and CAICS, Poland)