



COCOA: COrrrelation COefficient-Aware Data Augmentation

Mahdi Esmailoghli¹

Jorge-Arnulfo Quiané-Ruiz²

Ziawasch Abedjan^{1, 3}

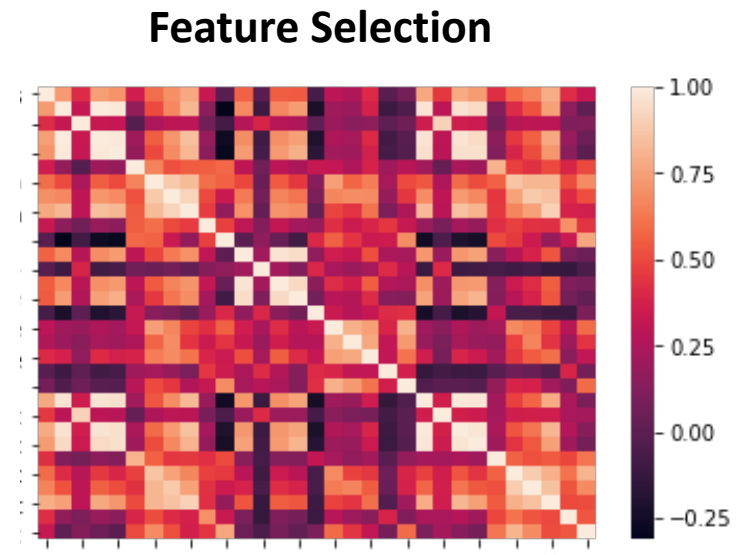
1 Leibniz Universität Hannover

2 TU Berlin

3 L3S Research Center

Correlating Coefficient

- Dependencies in the data
- Relevance of the feature
- Remove redundancy
- Feature importance
- Data enrichment



Data Enrichment



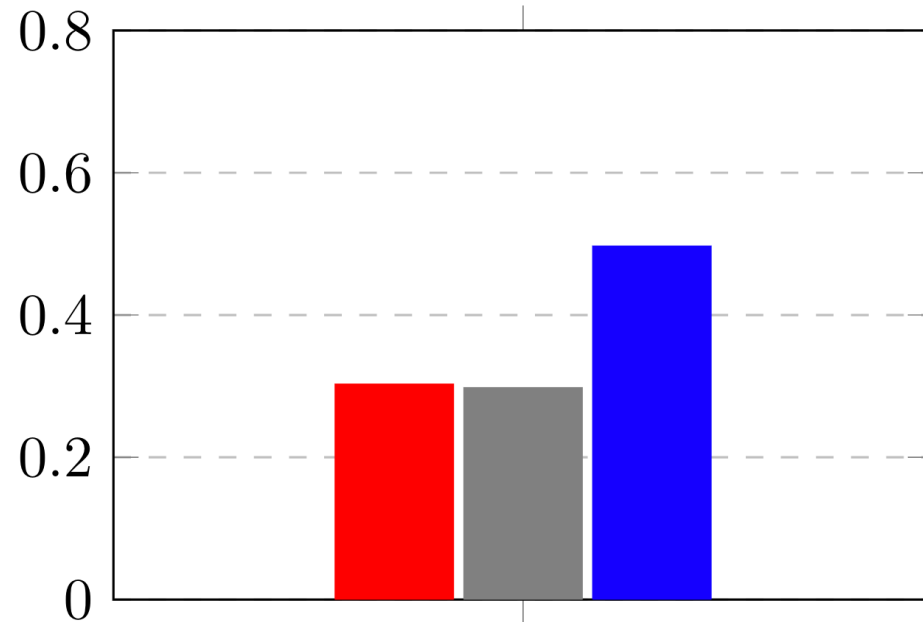
Correlating Coefficient

Poor Dataset

Country	RO
Germany	.7
U.S.A.	1.2
U.K.	1.0
China	.4



Average R^2 score



JOSIE¹ TR² Correlating features

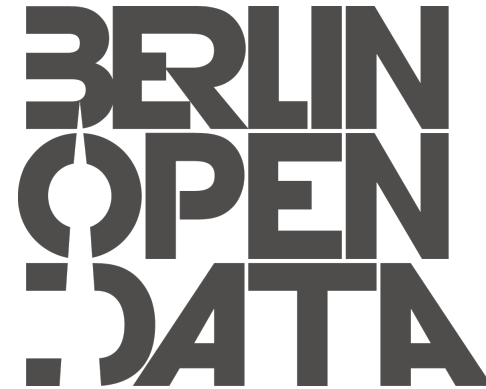
Enriched Dataset

Country	Education	Health	RO
Germany	95	Great	.7
U.S.A.	70	OK	1.2
U.K.	80	OK	1.0
China	90	Great	.4

1- Zhu, Erkang, et al. "Josie: Overlap set similarity search for finding joinable tables in data lakes." Proceedings of the 2019 International Conference on Management of Data. 2019.

2- Kumar, Arun, et al. "To join or not to join? Thinking twice about joins before feature selection." Proceedings of the 2016 International Conference on Management of Data. 2016.

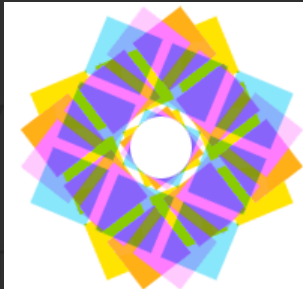
Available Data Lakes



Available Data Lakes



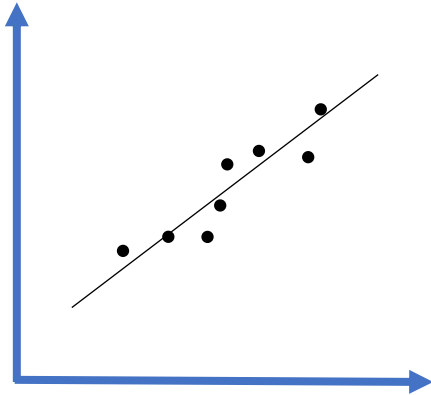
EU Open
Data Portal



DWTC
Dresden Web Table Corpus

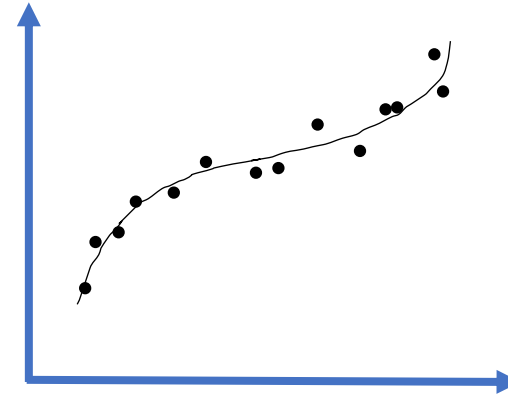
145 Million Tables
870 Million Columns

Time Complexity



Linear Correlation
e.g. Pearson

$O(n)$



Non-Linear Correlation
e.g. Spearman's

$O(n \cdot \log n)$

Categorical data: $O(n^2)$

Robustness

<i>Area (Million sq. miles)</i>	<i>Calling Code</i>
0.29	56
0.3	90
3.8	1
0.5	51
600	9800

Pearson = 1.0

Spearman's = 0.1

Goal

<i>Area (Million sq. miles)</i>	<i>Calling Code</i>
---------------------------------	---------------------

Calculate the **non-linear Spearman's correlation** in linear time.

0.5	51
600	9800

Pearson = 1.0

Spearman's = 0.1

Spearman's Correlation

Area	Code
0.3	90
0.29	56
3.8	1
0.5	51
600	9800



Area	Code
② 0.3	90 ④
① 0.29	56 ③
④ 3.8	1 ①
③ 0.5	51 ②
⑤ 600	9800 ⑤



$$s_{cc} = \frac{\sum_{i=1}^m (R(x_i) - \overline{R(x)})(R(y_i) - \overline{R(y)})}{\sqrt{\sum_{i=1}^m (R(x_i) - \overline{R(x)})^2 (R(y_i) - \overline{R(y)})^2}}$$

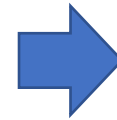
SCC= 0.1

Spearman's Correlation

Area	Code
0.29	56
0.3	90
3.8	1
0.5	51
600	9800



Area	Code
② 0.3	90 ④
① 0.29	56 ③
④ 3.8	1 ①
③ 0.5	51 ②
⑤ 600	9800 ⑤



$$S_{cc} = \frac{\sum_{i=1}^m (R(x_i) - \overline{R(x)})(R(y_i) - \overline{R(y)})}{\sqrt{\sum_{i=1}^m (R(x_i) - \overline{R(x)})^2 (R(y_i) - \overline{R(y)})^2}}$$

$O(n \cdot \log n)$

$O(n^2)$

Spearman's Correlation

Area	Code
0.29	56

Area	Code
① 0.29	56 ③

$$S_{cc} = \frac{\sum_{i=1}^m (R(x_i) - \overline{R(x)})(R(y_i) - \overline{R(y)})}{\sqrt{\sum_{i=1}^m (R(x_i) - \overline{R(x)})^2 (R(y_i) - \overline{R(y)})^2}}$$



COCOA leverages a **light-weight index** to calculate **non-linear correlation** in a **linear time**.

0.5	51
600	9800

③ 0.5	51 ②
⑤ 600	9800 ⑤

$O(n^2)$

Order Index (Offline)

Country	Area	Code
Turkey	0.3	90
Chile	0.29	56
Canada	3.8	1
Peru	0.5	51
Iran	600	9800



Area
2 0.3
1 0.29
4 3.8
3 0.5
5 600



Min

Area
0.3
0.29
3.8
0.5
600



Min

Row	Area
0	3
1	0
2	4
3	2
4	∅

COCOA (Online)

1. Detecting Tables with the highest overlap.

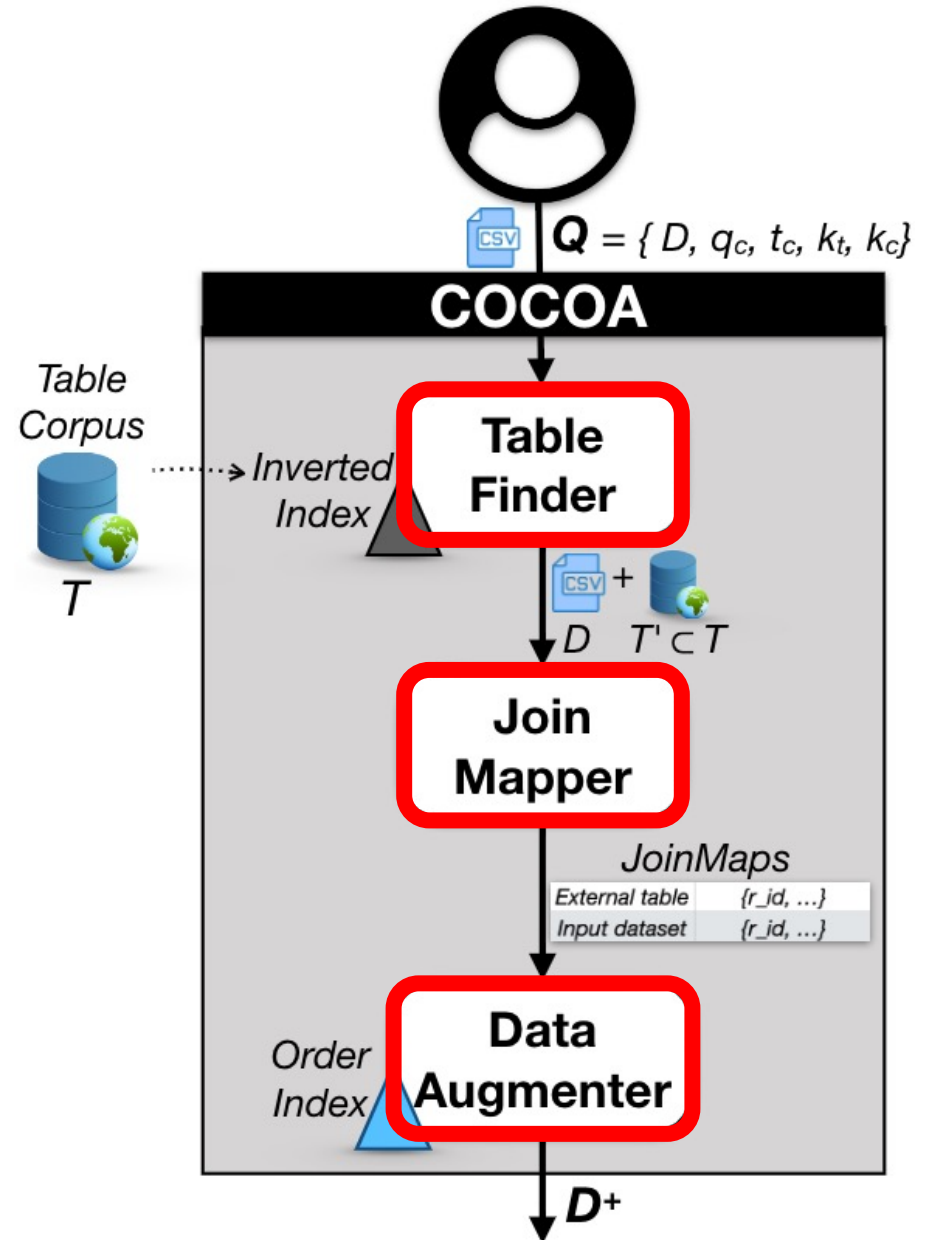
SELECT $T, C, \text{count}(*)$ **AS** *overlap*

FROM *tbl inverted index*

<i>row in C₁</i>	1	2	3	4	5	6	7	8	9
<i>row in q</i>	∅	3	∅	∅	∅	5	1	4	2

ORDER BY *overlap* **DESC**

LIMIT k_t



Correlation Calculation (Online)

Input Dataset

Row	Name (q)	Pop. (t)	Pop. (Rank)
1	Russia	144	5
2	Turkey	81	4
3	Switzerland	9	1
4	Sweden	10	2
5	Canada	37	3



Joinable Table

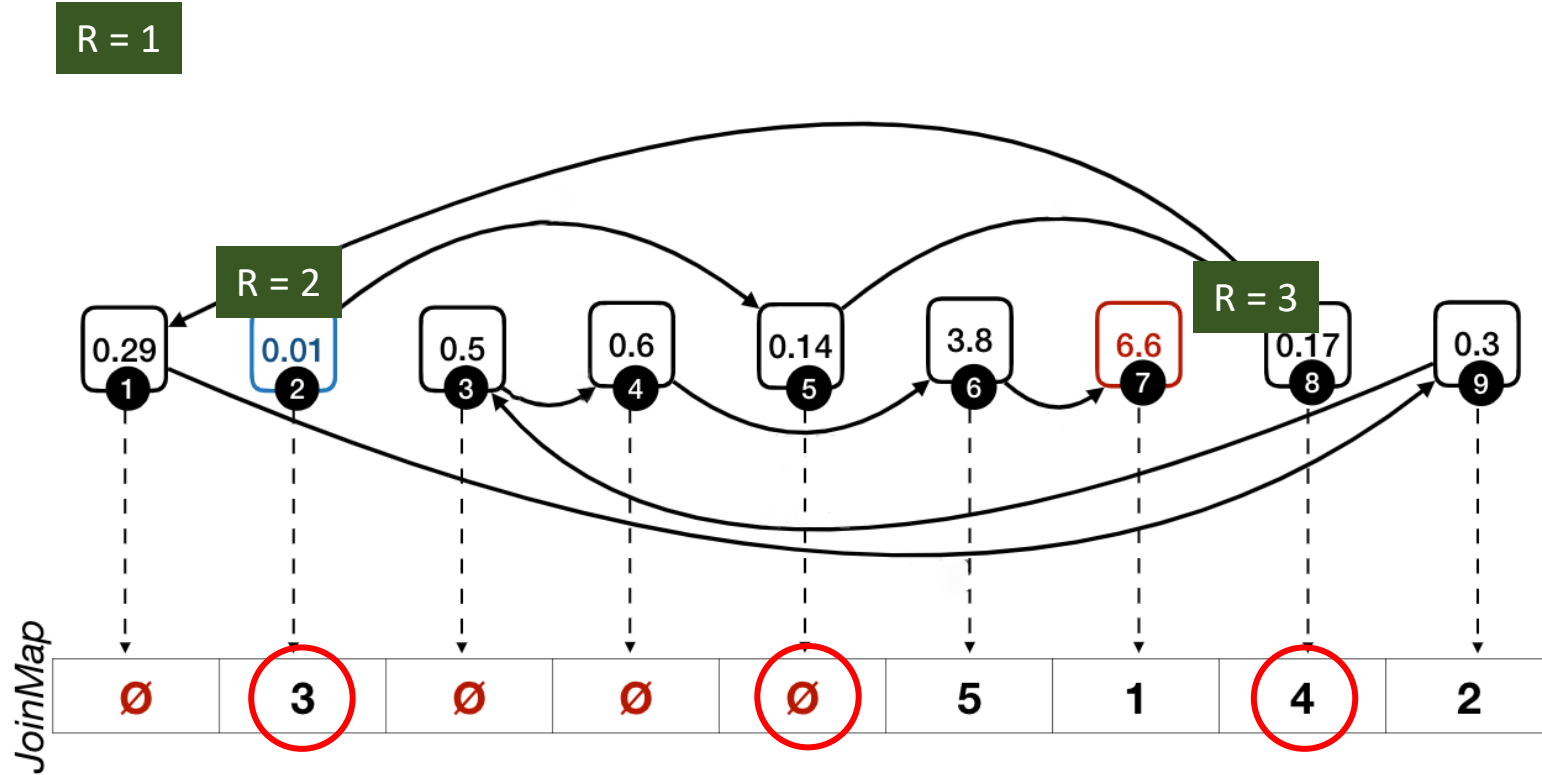
Row	Country	Area
1	Chile	0.29
2	Switzerland	0.01
3	Peru	0.5
4	Iran	0.6
5	Germany	0.14
6	Canada	3.8
7	Russia	6.6
8	Sweden	0.17
9	Turkey	0.3

Join Map

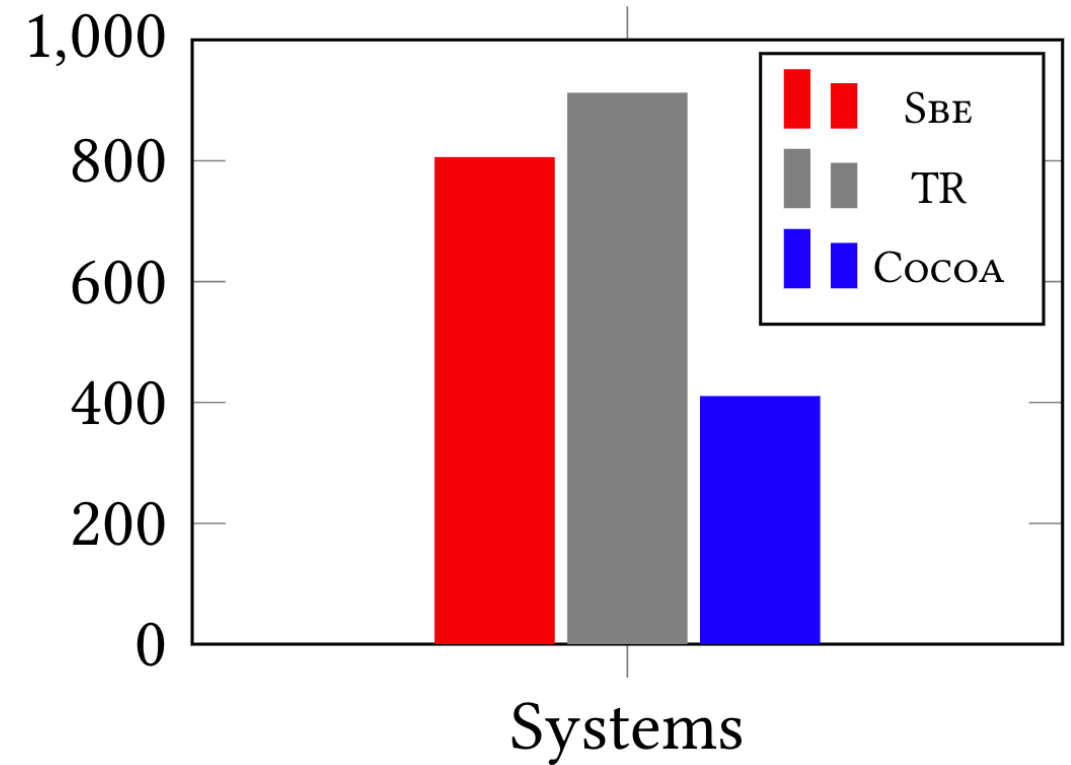
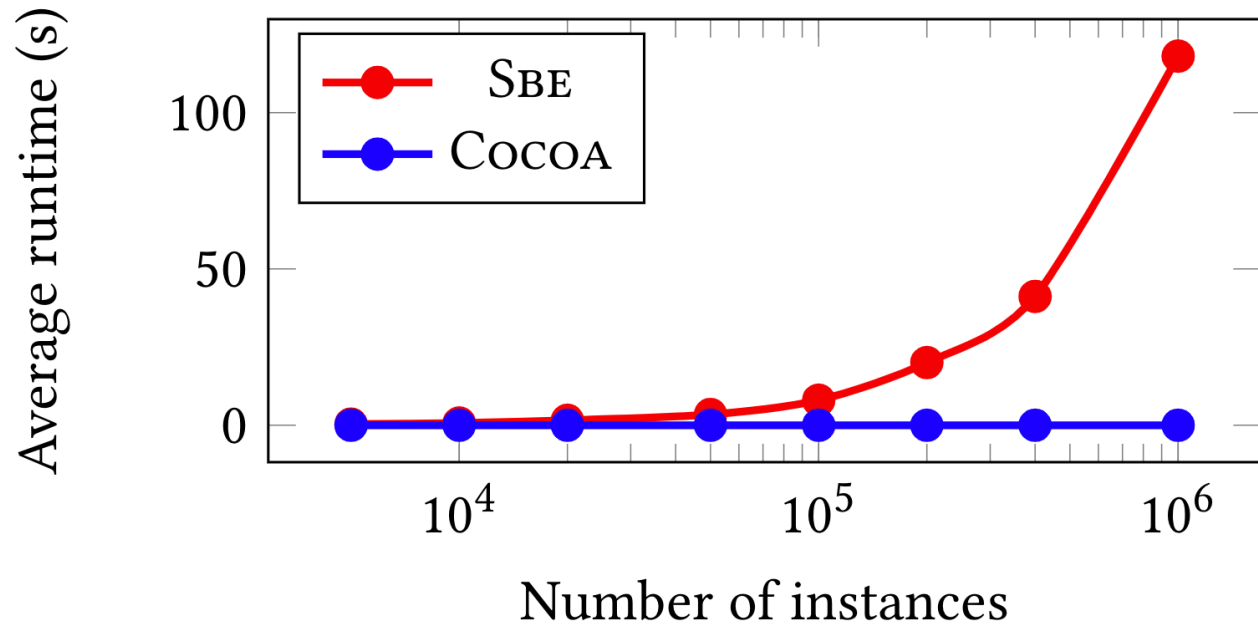
Row in Country	1	2	3	4	5	6	7	8	9
Row in q	∅	3	∅	∅	∅	5	1	4	2

Correlation Calculation (Online)

Row	Country	Pop.	Pop. (Rank)	Area (Rank)
1	Russia	144	5	5
2	Turkey	81	4	3
3	Switzerland	9	1	1
4	Sweden	10	2	2
5	Canada	37	3	4



Results





COCOA: COrrrelation COefficient-Aware Data Augmentation

Mahdi Esmailoghli¹

Jorge-Arnulfo Quiané-Ruiz²

Ziawasch Abedjan^{1, 3}

1 Leibniz Universität Hannover

2 TU Berlin

3 L3S Research Center