CS 8803-MDS Human-in-the-loop Data Analytics

Lecture 12 10/03/22

Logistics

Grading

<= 300 words per answer for paper reviews project proposal grades will be released by next week Project timeline

project update presentation 10/31 Assignment 2-6: progress report (1% each) due Fridays 5PM at 10/21, 10/28, 11/4, 11/11, 11/18 Assignment 7: evaluation plan (5%) due Friday 5PM at 11/18

Next week:

how to make progress in research?

Today's class

Finding Related Tables in Data Lakes for Interactive Data Science

Author: Qiandong, Shen En

Reviewer: Vishnu

Archaeologist: Yanhao

Practioner: Haotian



1





ACM SIGMOD/PODS International Conference on Management of Data June 14 - June 19, 2020 Portland, OR, USA

Welcome Homepage	2020 ACM SIGMOD/PODS @ Portland, OR, USA SIGMOD/PODS 2020 Experience Report now available	Calls For Submissions
News		Calls for Submissions
Organization	<u>Coronavirus Updates</u>	PODS Program
Experience Report	Undates on SIGMOD/PODS 2020 Pagistration (12 May	Program Overview
Organization	2020)	Detailed Program
SIGMOD PC		Research Papers
PODS PC		Keynote Talk

- Many data science applications build on data lakes today
- Data Lakes
 - A central repository to store any types of data
 - Schema-on-Read
 - Easy to store data, low storage cost

Schema-on-Read



Npujwbujpo

- Data scientists do lots of repeated work every day
 - Importing and cleaning today's data
 - Data wrangling and exploration on public datasets

- Data scientists do lots of repeated work every day
 - Importing and cleaning today's data
 - Data wrangling and exploration on public datasets



- Data scientists do lots of repeated work every day
 - Importing and cleaning today's data
 - Data wrangling and exploration on public datasets
- Same data might already exist



- Data scientists do lots of repeated work every day
 - Importing and cleaning today's data
 - Data wrangling and exploration on public datasets
- Same data might already exist
- Data lakes do little to help user find related data
 - Users don't know what data is available, or unable to trust what they find
 - Users reinvent schema redundant work, data inconsistency



- What if we could help users find related data?
 - Promote reusable units of data and processing



- What if we could help users find related data?
 - Promote reusable units of data and processing
- Goal: Build a search framework to find related tables



- Based on Jupyter Notebook
- Considering
 - source data
 - cells
 - intermediate data
- Github Copilot for suggesting tables

Import the Data

Then use the read_csv() method in the cell below to import the data into a data frame - let's call it df. The data is in a file called new_data.csv

In [2]: df = pd.read_csv('new_data.csv')

Look at the Data

Take a little time in the next cell to explore the data. If you want to add new cells feel free to "Insert - Insert Cell Below" in the Jupyter Notebook.

[n [3]:	: df.head()													
Out[3]:		CompanyName	StateIncorporated	NumberofEmployees	EntityType									
	0	Ameren Corporation	Delaware	3000	C corp									
	1	America West Holdings Corporation	Colorado	100	lic									
	2	American Axle & Manufacturing Holdings, Inc.	California	30	S corp									
	3	American Eagle Outfitters, Inc.	California	200	c corp									
	4	American Electric Power Company, Inc.	Delaware	5000	c corp									

In [4]: df.info()

class 'pandas.core.frame.DataFrame'>
Rangeindex:55 entries,0 to 54
Data columns (total 4 columns):
CompanyName 55 non-null object
StateIncorporated 55 non-null object
NumberofEmployees 55 non-null int64
EntityType 55 non-null object
dtypes: int64(1), object(3)
memory usage: 1.84 KB

In [8]:	for i in df.a df['Enti	<pre>df.index: t[i, 'EntityType'] = df.at[i, 'EntityType'].upper() tyType'].value_counts()</pre>
Out[8]:	C CORP	27
	LLC	19
	S CORP	9
	Name: En	tityType, dtype: int64
A	re yc	ou looking for df_xxx?

- Augmenting training / validating data
 - Union of two tables
- Linking data
 - Join of two tables
- Extracting machine learning features
 - Adding new columns to the original table
- Cleaning data
 - Fill in missing values in the original table



- Augmenting training / validating data
 - Union of two tables
- Linking data
 - Join of two tables
- Extracting machine learning features
 - Adding new columns to the original table
- Cleaning data
 - Fill in missing values in the original table



- Augmenting training / validating data
 - Union of two tables
- Linking data
 - Join of two tables
- Extracting machine learning features
 - Adding new columns to the original table
- Cleaning data
 - Fill in missing values in the original table



- Augmenting training / validating data
 - Union of two tables
- Linking data
 - Join of two tables
- Extracting machine learning features
 - Adding new columns to the original table
- Cleaning data
 - Fill in missing values in the original table



Problem Statement



Data Lake

Problem Statement



Problem Statement



An Example - Feature Extraction

- Feature extraction
 - Produce tables that preserve inputs' keys and other columns, but add new columns

An Example - Feature Extraction

	Passengerld	Survived	Pclass	Name	Age	SibSp	Parch	Ticket	Fare	Cabin
0	1	0	3	Braund, Mr. Owen Harris	22.0	1	0	A/5 21171	7.2500	NaN
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th	38.0	1	0	PC 17599	71.2833	C85
2	3	1	3	Heikkinen, Miss. Laina	26.0	0	0	STON/02. 3101282	7.9250	NaN





	Passengerid	Survived	Pclass	Name	Sex	SibSp	Parch	Fare	Embarked	First Name	Last Name	AgeBand
0	1	0	3	Braund, Mr. Owen Harris	0	1	0	7.2500	S	Braund	Owen Harris	(16.0, 32.0]
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th	1	1	0	71.2833	C	Cumings	John Bradley (Florence Briggs Thayer)	(32.0, 48.0]
2	3	1	3	Heikkinen, Miss. Laina	1	0	0	7.9250	S	Heikkinen	Laina	(16.0, 32.0]

Feature Extraction - Table Overlap



Feature Extraction - Table Overlap

	Passengerid	Survived	Pclass	Name	Age	SibSp	Parch	Ticket	Fare	Cabin
0	1	0	3	Braund, Mr. Owen Harris	22.0	1	0	A/5 21171	7.2500	NaN
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th	38.0	1	0	PC 17599	71.2833	C85
2	3	1	3	Heikkinen, Miss. Laina	26.0	0	0	STON/02. 3101282	7.9250	NaN



Schema Overlap

	Passengerld	Survived	Pclass	Name	Sex	SibSp	Parch	Fare	Embarked	First Name	Last Name	AgeBand
0	1	0	3	Braund, Mr. Owen Harris	0	1	0	7.2500	S	Braund	Owen Harris	(16.0, 32.0]
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th	1	1	0	71.2833	C	Curnings	John Bradley (Florence Briggs Thayer)	(32.0, 48.0]
2	3	1	3	Heikkinen, Miss. Laina	1	0	0	7.9250	S	Heikkinen	Laina	(16.0, 32.0]



Feature Extraction - New Information

	Passengerid	Survived	Pclass	Name	Age	SibSp	Parch	Ticket	Fare	Cabin
0	1	0	3	Braund, Mr. Owen Harris	22.0	1	0	A/5 21171	7.2500	NaN
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th	38.0	1	0	PC 17599	71.2833	C85
2	3	1	3	Heikkinen, Miss. Laina	26.0	0	0	STON/02. 3101282	7.9250	NaN





New Columns

	Passengerid	Survived	Pclass	Name	Sex	SibSp	Parch	Fare	Embarked	First Name	Last Name	AgeBand
0	1	0	3	Braund, Mr. Owen Harris	0	1	0	7.2500	S	Braund	Owen Harris	(16.0, 32.0)
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th	1	1	0	71.2833	С	Curnings	John Bradley (Florence Briggs Thayer)	(32.0, 48.0)
2	3	1	3	Heikkinen, Miss. Laina	1	0	0	7.9250	S	Heikkinen	Laina	(16.0, 32.0)

Feature Extraction - New Information

	Passengerid	Survived	Pclass	Name	Age	SibSp	Parch	Ticket	Fare	Cabin
0	1	0	3	Braund, Mr. Owen Harris	22.0	1	0	A/5 21171	7.2500	NaN
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th	38.0	1	0	PC 17599	71.2833	C85
2	3	1	3	Heikkinen, Miss. Laina	26.0	0	0	STON/02. 3101282	7.9250	NaN





Provenance

	Passengerld	Survived	Pclass	Name	Sex	SibSp	Parch	Fare	Embarked	First Name	Last Name	AgeBand
0	1	0	3	Braund, Mr. Owen Harris	0	1	0	7.2500	S	Braund	Owen Harris	(16.0, 32.0)
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th	1	1	0	71.2833	С	Curnings	John Bradley (Florence Briggs Thayer)	(32.0, 48.0)
2	3	1	3	Heikkinen, Miss. Laina	1	0	0	7.9250	S	Heikkinen	Laina	(16.0, 32.0)

Full Suite of Basic Measures

- 1. Row/Column Match
 - Measures the row overlap and column overlap between two tables S and T
- 2. Row/Column Information Gain
 - Measures the new row rate and new column rate of T compared to query table S
- 3. Shared Provenance
 - Measures the purpose/workflow similarity of two tables S and T based on variable dependencies
- 4. Null Value Decrement
 - Measures the decrease in null value in T compared to query table S
- 5. Description Similarity
 - Measures the similarity of table metadata such as problem type, text about the workflow, etc.

Full Suite of Basic Measures

Use Cases	Row Match	Column Match	Row Info Gain	Column Info Gain	Shared Provenance	Null Value Decrement	Description Similarity
Data Augmentation		Î Î	ÎÎ		î		ſ
Feature Extraction	† †			† †	TT		Î
Data Cleaning	† †	† †			Î	Î	ſ
Linkable Data	î î						Î



Example: Assume our goal is to find the top-2 (k = 2) candidate tables T_i

Step 1: For each measure M_i, order the candidate tables T_i in descending order of scores

M ₁	Score	M ₂	Score	M ₃	Score
T ₁	1	T ₂	0.8	T ₄	0.8
T ₂	0.8	T ₃	0.7	T ₃	0.6
T_3	0.5	T ₁	0.3	T ₁	0.2
T ₄	0.3	T ₄	0.2	T ₅	0.1
T_5	0.1	T ₅	0.1	T ₂	0

Step 2: Perform sequential scan of the tables, set a threshold t = sum of the scores in the scan

M ₁	Score	M ₂	Score	M ₃	Score	
T ₁	1	T ₂	0.8	T ₄	0.8	→ t = 1 + 0.8 + 0.8 = 2.6
T ₂	0.8	T ₃	0.7	T ₃	0.6	
T ₃	0.5	T ₁	0.3	T ₁	0.2	
T ₄	0.3	T ₄	0.2	T ₅	0.1	
Т ₅	0.1	Т ₅	0.1	T ₂	0	

Step 3: At the same time, for each T_i appearing in the scan, calculate their total scores

M ₁	Score	_	M ₂	Score		M ₃	Score	
T ₁	1		T ₂	0.8		T ₄	0.8	→ t = 1 + 0.8 + 0.8 = 2.6
T ₂	0.8	+	T ₃	0.7		T ₃	0.6	s = 1 + 0.3 + 0.2 = 1.5
T ₃	0.5		T ₁	0.3	+	T ₁	0.2	$3_1 - 1 + 0.5 + 0.2 - 1.5$
T ₄	0.3		T ₄	0.2		T ₅	0.1	
T ₅	0.1		T ₅	0.1		T ₂	0	

Step 3: At the same time, for each T_i appearing in the scan, calculate their total scores

M ₁	Score		M ₂	Score	_	M ₃	Score	
T ₁	1		T ₂	0.8		T ₄	0.8	→ t = 1 + 0.8 + 0.8 = 2.6
T ₂	0.8	+	T ₃	0.7		T ₃	0.6	s = 1 + 0.3 + 0.2 = 1.5
T ₃	0.5		T ₁	0.3	+	T ₁	0.2	3 ₁ 1 0.0 0.2 1.5
T ₄	0.3		T ₄	0.2		T ₅	0.1	$s_2 = 0.8 + 0.8 + 0 = 1.6$
T ₅	0.1		T ₅	0.1		T ₂	0	

Step 3: At the same time, for each T_i appearing in the scan, calculate their total scores

M ₁	Score		M ₂	Score		M ₃	Score	
T ₁	1		T ₂	0.8		T ₄	0.8	→ t = 1 + 0.8 + 0.8 = 2.6
T_2	0.8		T ₃	0.7		T ₃	0.6	s = 1 + 0.3 + 0.2 = 1.5
T ₃	0.5		T ₁	0.3	+	T ₁	0.2	3 ₁ 1 0.0 0.2 1.3
T ₄	0.3	+	T ₄	0.2		T ₅	0.1	$s_2 = 0.8 + 0.8 + 0 = 1.6$
Т ₅	0.1		Т ₅	0.1		T ₂	0	s ₄ = 0.3 + 0.2 + 0.8 = 1.3

Step 3.5: We only need the top-2 tables, so we only maintain the top-2 total scores

M ₁	Score	M ₂	Score	M ₃	Score]
T ₁	1	T ₂	0.8	T ₄	0.8	→ t = 1 + 0.8 + 0.8 = 2.6
T ₂	0.8	T ₃	0.7	T ₃	0.6	s = 1.6
T ₃	0.5	T ₁	0.3	T ₁	0.2	s ₂ = 1.5
T ₄	0.3	T ₄	0.2	T ₅	0.1	
T ₅	0.1	T ₅	0.1	T ₂	0	s ₄ = 0.3 + 0.2 + 0.8 = 1.3

Step 4: Repeat Step 2 and Step 3 by:

(1) continuing the scan, (2) updating the threshold, (3) scale total scores of new tables

M ₁	Score		M ₂	Score		M ₃	Score	Updated threshold
T ₁	1		T ₂	0.8		T ₄	0.8	→ t = 0.8 + 0.7 + 0.6 = 2.1
T ₂	0.8		T ₃	0.7		T ₃	0.6	s = 18
T ₃	0.5		T ₁	0.3		T ₁	0.2	$s_3 = 1.6$
T ₄	0.3		T ₄	0.2		T ₅	0.1	Updated top-2 total scores:
T ₅	0.1		T ₅	0.1		T ₂	0	1. Calculate score for T_3
		4	-		4			2. $s_3 = 1.8$ is larger than $s_1 = 1.5$

Step 5: Stop when all the top-2 total scores maintained ≥ threshold

M ₁	Score		M ₂	Score	M ₃	Score	
T ₁	1		T ₂	0.8	T ₄	0.8	→ t = 0.5 + 0.3 + 0.2 = 1
T ₂	0.8		T ₃	0.7	T ₃	0.6	s = 1.8 We have found our top-2 tables
T ₃	0.5		T ₁	0.3	T ₁	0.2	s ₃ = 1.6
T ₄	0.3		T ₄	0.2	T ₅	0.1	2
T ₅	0.1	1	Т ₅	0.1	T ₂	0	
Step 5: Stop when all the top-2 total scores maintained ≥ threshold

M ₁	Score	M ₂	Score	M ₃	Score	
T ₁	1	T ₂	0.8	T ₄	0.8	\rightarrow t = 0.5 + 0.3 + 0.2 = 1
T ₂	0.8	T ₃	0.7	T ₃	0.6	s = 18 We have found our top-2 tables
T ₃	0.5	T ₁	0.3	T ₁	0.2	$s_{3} = 1.6$
T ₄	0.3	T ₄	0.2	T ₅	0.1	
T ₅	0.1	T ₅	0.1	T ₂	0	Why?

Step 5: Stop when all the top-2 total scores maintained ≥ threshold

M ₁	Score	M ₂	Score	M ₃	Score	
T ₁	1	T ₂	0.8	T ₄	0.8	t = 0.5 + 0.3 + 0.2 = 1
T ₂	0.8	T ₃	0.7	T ₃	0.6	s = 18 We have found our top-2 tables
T ₃	0.5	T ₁	0.3	T ₁	0.2	$s_{3} = 1.6$
T ₄	0.3	T ₄	0.2	T ₅	0.1	
T ₅	0.1	T ₅	0.1	T ₂	0	Why?
						been in the top-k list, we would have seen

it during the sequential scan

Challenge:

Scores requiring relation mapping are much more expensive to compute

M ₁	Score	M ₂	Score	M ₃	Score	
T ₁	1	T ₂	0.8	T ₄	0.8	\rightarrow t = 0.5 + 0.3 + 0.2 = 1
T ₂	0.8	T ₃	0.7	T ₃	0.6	s = 18
T ₃	0.5	T ₁	0.3	T ₁	0.2	$s_3 = 1.6$
T ₄	0.3	T ₄	0.2	T ₅	0.1	
T ₅	0.1	T ₅	0.1	Т ₂	0	

Challenge:

e.g. To calculate Row Overlap, we need to know how the table headers map to each other Scores requiring relation mapping are much more expensive to compute

M ₁	Score	M ₂	Score	M ₃	Score	
T ₁	1	T ₂	0.8	T ₄	0.8	→ t = 0.5 + 0.3 + 0.2 = 1
T ₂	0.8	T ₃	0.7	T ₃	0.6	s = 18
T ₃	0.5	T ₁	0.3	T ₁	0.2	$s_3 = 1.6$
T ₄	0.3	T ₄	0.2	T ₅	0.1	
T ₅	0.1	T ₅	0.1	T ₂	0	

Computational Challenges

- Relation mappings are expensive
 - \circ For two tables with size (R x C), computing table overlaps: O(RC²)
- Top-k search makes it worse
 - \circ $\,$ $\,$ For N tables, top-k search: O(NRC^2) $\,$

Example: Assume M_4 is much more expensive to compute than M_1 , M_2 , and M_3 .

M ₁	Score	M ₂	Score	M ₃	Score	M ₄	Score
T ₁	1	T ₂	0.8	T ₄	0.8		
T ₂	0.8	T ₃	0.7	T ₃	0.6		
T ₃	0.5	T ₁	0.3	T ₁	0.2		
T ₄	0.3	T ₄	0.2	T ₅	0.1		
T ₅	0.1	T ₅	0.1	T ₂	0		

Example: Assume M_4 is much more expensive to compute than M_1 , M_2 , and M_3 .

M ₁	Score	M ₂	Score	M ₃	Score	M ₄	Score	
T ₁	1	T ₂	0.8	T ₄	0.8			Instead of pre-compute all
T ₂	0.8	T ₃	0.7	T ₃	0.6			scores for M ₄ ,
T ₃	0.5	T ₁	0.3	T ₁	0.2			we do it on-demand
T ₄	0.3	T ₄	0.2	T ₅	0.1			
T ₅	0.1	T ₅	0.1	T ₂	0			

Trick: In each sequential scan, use the Upper Bound of M_4 instead of the individual M_4 scores.

M ₁	Score		M ₂	Score	M ₃	Score	M ₄	Score	
T ₁	1		T ₂	0.8	T ₄	0.8	UB	0.9	+ t = 0.1 + 0.8 + 0.8 + UB = 3.5
T ₂	0.8		T ₃	0.7	T ₃	0.6			
T ₃	0.5		T ₁	0.3	T ₁	0.2			
T ₄	0.3	-	T ₄	0.2	T ₅	0.1			
T ₅	0.1		T ₅	0.1	T ₂	0			

Trick: Calculate individual M₄ scores on-demand when covered by the scan

M ₁	Score	M ₂	Score	M ₃	Score	M ₄	Score	
T ₁	1	T ₂	0.8	T ₄	0.8	UB	0.9	★ t = 0.1 + 0.8 + 0.8 + UB = 3.5
T ₂	0.8	T ₃	0.7	T ₃	0.6	T ₁	0.6	$s_4 = 1 + 0.3 + 0.2 + 0.6 = 2.1$
T ₃	0.5	T ₁	0.3	T ₁	0.2	T ₂	0.5	s = 0.8 + 0.8 + 0 + 0.5 = 2.1
T ₄	0.3	T ₄	0.2	T ₅	0.1	T ₄	0.3	$s_2 = 0.3 \pm 0.2 \pm 0.8 \pm 0.3 \pm 1.4$
T ₅	0.1	T ₅	0.1	T ₂	0			$s_4 = 0.5 \pm 0.2 \pm 0.0 \pm 0.5 = 1.0$

	Passengerid	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket
590	591	0	3	Rintamaki, Mr. Matt	male	35.0	0	0	STON/O 2. 3101273
131	132	0	3	Coelho, Mr. Domingos Fernandeo	male	20.0	0	0	SOTON/O.Q. 3101307
628	629	0	3	Bostandyeff, Mr. Guentcho	male	26.0	0	0	349224

- 2. Compositional Profile
- 3. Workflow Graph



	Passengerld	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket
590	591	0	3	Rintamaki, Mr. Matt	male	35.0	0	0	STON/O 2. 3101273
131	132	0	3	Coelho, Mr. Domingos Fernandeo	male	20.0	0	0	SOTON/O.Q. 3101307
628	629	0	3	Bostandyeff, Mr. Guentcho	male	26.0	0	0	349224

Matcher tests whether a column belong to the domain

- 2. Compositional Profile
- 3. Workflow Graph





	Name	Sex_Code	Pclass	Embarked_Code	Title_Code	SibSp	Parch	Age	Fare
590	Rintamaki, Mr. Matti	1	3	2	3	0	0	35.0	7.1250
131	Coelho, Mr. Domingos Fernandeo	1	3	2	3	0	0	20.0	7.0500
628	Bostandyeff, Mr. Guentcho	1	3	2	3	0	0	26.0	7.8958



	Passengerid	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket
590	591	0	3	Rintamaki, Mr. Matti	male	35.0	0	0	STON/O 2. 3101273
131	132	0	3	Coelho, Mr. Domingos Fernandeo	male	20.0	0	0	SOTON/O.Q. 3101307
628	629	0	3	Bostandyeff, Mr. Guentcho	male	26.0	0	0	349224

1. Data Profiles

- 2. Compositional Profile
- 3. Workflow Graph

Data Profile Indices



Pas	sengerid	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	BirthDate
590	591	0	3	Rintamaki, Mr. Matti	male	35.0	0	0	STON/O 2. 3101273	1877/09/01
131	132	0	3	Coelho, Mr. Domingos Fernandeo	male	20.0	0	0	SOTON/O.Q. 3101307	1892/10/23
528	629	0	3	Bostandyeff, Mr. Guentcho	male	26.0	0	0	349224	1886/03/15
	ln co	dex co mmon	lumn se co-occi	ts based on urrences	C	ompo	sitional	Profile	Index	

- 1. Data Profiles
- 2. Compositional Profile
- 3. Workflow Graph

	Name	Sex_Code	Pclass	Embarked_Code	Title_Code	SibSp	Parch	Age	Fare	BirthDate
590	Rintamaki, Mr. Matti	1	3	2	3	0	0	35.0	7.1250	1877/09/01
131	Coelho, Mr. Domingos Fernandeo	1	3	2	3	0	0	20.0	7.0500	1892/10/23
628	Bostandyeff, Mr. Guentcho	1	3	2	3	0	0	26.0	7.8958	1886/03/15



- 1. Data Profiles
- 2. Compositional Profile
- 3. Workflow Graph



- 1. Data Profiles
- 2. Compositional Profile
- 3. Workflow Graph

- 2. Compositional Profile
- 3. Workflow Graph



- 2. Compositional Profile
- 3. Workflow Graph



- 2. Compositional Profile
- 3. Workflow Graph



System Implementation



System Implementation



Experiment Setup

- 5000+ indexed tables
- 102 notebooks from Kaggle, hand-labeled for use cases
- 14 different Kaggle tasks

top-k	Augm	nenting	Training	g Data	F	eature l	Extractio	n	Alter	native 1	Data Cle	aning	Linking Data			
Q_4	L64	TA	TA+P	SJ	L64	TA	TA+P	SJ	L64	TA	TA+P	SJ	L64	TA	TA+P	SJ
1	919.4	1017	9.776	1.481	875.3	2036	8.288	2.332	947.7	1454	5.833	3.133	22.64	2252	14.67	3.033
5	919.4	1086	9.944	3.254	875.3	2149	8.624	2.336	947.7	1461	5.999	3.145	22.64	2252	15.077	3.261
10	919.4	1106	12.17	5.105	875.3	2228	9.744	3.302	947.7	1488	7.029	4.364	22.64	2252	15.510	4.169
Q_5	L64	TA	TA+P	SJ	L64	TA	TA+P	SJ	L64	TA	TA+P	SJ	L64	TA	TA+P	SJ
1	1755	NF	23.51	5.119	1691	NF	21.30	5.188	1871	NF	20.46	4.318	127.9	NF	22.44	4.246
5	1755	NF	23.86	5.290	1691	NF	21.87	5.417	1871	NF	20.87	4.517	127.9	NF	22.51	4.343
10	1755	NF	24.15	5.532	1691	NF	23.29	5.666	1871	NF	20.95	4.892	127.9	NF	22.73	4.864

Table 3: Average running time (seconds) of returning top-k tables when searching for related tables

L64 = LSH Ensemble

TA = (Ours) Threshold Algorithm

(without Data Profiles and Indices)

TA+P = (Ours) Threshold Algorithm

(with Data Profiles and Indices)

top-k	Augm	enting	Training	g Data	F	eature l	Extractio	n	Alter	native l	Data Cle	aning	Linking Data			
Q_4	L64	TA	TA+P	SJ	L64	TA	TA+P	SJ	L64	TA	TA+P	SJ	L64	TA	TA+P	SJ
1	919.4	1017	9.776	1.481	875.3	2036	8.288	2.332	947.7	1454	5.833	3.133	22.64	2252	14.67	3.033
5	919.4	1086	9.944	3.254	875.3	2149	8.624	2.336	947.7	1461	5.999	3.145	22.64	2252	15.077	3.261
10	919.4	1106	12.17	5.105	875.3	2228	9.744	3.302	947.7	1488	7.029	4.364	22.64	2252	15.510	4.169
Q_5	L64	TA	TA+P	SJ	L64	TA	TA+P	SJ	L64	TA	TA+P	SJ	L64	TA	TA+P	SJ
1	1755	NF	23.51	5.119	1691	NF	21.30	5.188	1871	NF	20.46	4.318	127.9	NF	22.44	4.246
5	1755	NF	23.86	5.290	1691	NF	21.87	5.417	1871	NF	20.87	4.517	127.9	NF	22.51	4.343
10	1755	NF	24.15	5.532	1691	NF	23.29	5.666	1871	NF	20.95	4.892	127.9	NF	22.73	4.864

Table 3: Average running time (seconds) of returning top-k tables when searching for related tables

> 2300 seconds

L64 = LSH Ensemble

TA = (Ours) Threshold Algorithm

(without Data Profiles and Indices)

TA+P = (Ours) Threshold Algorithm

(with Data Profiles and Indices)

top-k	Augm	enting	Trainin	g Data	F	eature l	Extractio	n	Alter	native 1	Data Cle	aning	Linking Data			
Q_4	L64	TA	TA+P	SJ	L64	TA	TA+P	SJ	L64	TA	TA+P	SJ	L64	TA	TA+P	SJ
1	919.4	1017	9.776	1.481	875.3	2036	8.288	2.332	947.7	1454	5.833	3.133	22.64	2252	14.67	3.033
5	919.4	1086	9.944	3.254	875.3	2149	8.624	2.336	947.7	1461	5.999	3.145	22.64	2252	15.077	3.261
10	919.4	1106	12.17	5.105	875.3	2228	9.744	3.302	947.7	1488	7.029	4.364	22.64	2252	15.510	4.169
Q_5	L64	TA	TA+P	SJ	L64	TA	TA+P	SJ	L64	TA	TA+P	SJ	L64	TA	TA+P	SJ
1	1755	NF	23.51	5.119	1691	NF	21.30	5.188	1871	NF	20.46	4.318	127.9	NF	22.44	4.246
5	1755	NF	23.86	5.290	1691	NF	21.87	5.417	1871	NF	20.87	4.517	127.9	NF	22.51	4.343
10	1755	NF	24.15	5.532	1691	NF	23.29	5.666	1871	NF	20.95	4.892	127.9	NF	22.73	4.864

Table 3: Average running time (seconds) of returning top-k tables when searching for related tables

Adding data profiles and indices speeds up the search by 100x.

L64 = LSH Ensemble

TA = (Ours) Threshold Algorithm

(without Data Profiles and Indices)

TA+P = (Ours) Threshold Algorithm

(with Data Profiles and Indices)

top-k	Augm	nenting '	Training	g Data	F	eature l	Extractio	n	Alter	native 1	Data Cle	aning	Linking Data			
Q_4	L64	TA	TA+P	SJ	L64	TA	TA+P	SJ	L64	TA	TA+P	SJ	L64	TA	TA+P	SJ
1	919.4	1017	9.776	1.481	875.3	2036	8.288	2.332	947.7	1454	5.833	3.133	22.64	2252	14.67	3.033
5	919.4	1086	9.944	3.254	875.3	2149	8.624	2.336	947.7	1461	5.999	3.145	22.64	2252	15.077	3.261
10	919.4	1106	12.17	5.105	875.3	2228	9.744	3.302	947.7	1488	7.029	4.364	22.64	2252	15.510	4.169
Q_5	L64	TA	TA+P	SJ	L64	TA	TA+P	SJ	L64	TA	TA+P	SJ	L64	TA	TA+P	SJ
1	1755	NF	23.51	5.119	1691	NF	21.30	5.188	1871	NF	20.46	4.318	127.9	NF	22.44	4.246
5	1755	NF	23.86	5.290	1691	NF	21.87	5.417	1871	NF	20.87	4.517	127.9	NF	22.51	4.343
10	1755	NF	24.15	5.532	1691	NF	23.29	5.666	1871	NF	20.95	4.892	127.9	NF	22.73	4.864

Table 3: Average running time (seconds) of returning top-k tables when searching for related tables

Including workflow indices further decreases the compute time by 3-5x.

L64 = LSH Ensemble

TA = (Ours) Threshold Algorithm

(without Data Profiles and Indices)

TA+P = (Ours) Threshold Algorithm

(with Data Profiles and Indices)

top-k	Augm	nenting	Training	g Data	F	eature l	Extractio	n	Alter	native 1	Data Cle	aning	Linking Data			
Q_4	L64	TA	TA+P	SJ	L64	TA	TA+P	SJ	L64	TA	TA+P	SJ	L64	TA	TA+P	SJ
1	919.4	1017	9.776	1.481	875.3	2036	8.288	2.332	947.7	1454	5.833	3.133	22.64	2252	14.67	3.033
5	919.4	1086	9.944	3.254	875.3	2149	8.624	2.336	947.7	1461	5.999	3.145	22.64	2252	15.077	3.261
10	919.4	1106	12.17	5.105	875.3	2228	9.744	3.302	947.7	1488	7.029	4.364	22.64	2252	15.510	4.169
Q_5	L64	TA	TA+P	SJ	L64	TA	TA+P	SJ	L64	TA	TA+P	SJ	L64	TA	TA+P	SJ
1	1755	NF	23.51	5.119	1691	NF	21.30	5.188	1871	NF	20.46	4.318	127.9	NF	22.44	4.246
5	1755	NF	23.86	5.290	1691	NF	21.87	5.417	1871	NF	20.87	4.517	127.9	NF	22.51	4.343
10	1755	NF	24.15	5.532	1691	NF	23.29	5.666	1871	NF	20.95	4.892	127.9	NF	22.73	4.864

Table 3: Average running time (seconds) of returning top-k tables when searching for related tables

L64 = LSH Ensemble

TA = (Ours) Threshold Algorithm

(without Data Profiles and Indices)

TA+P = (Ours) Threshold Algorithm

(with Data Profiles and Indices)

SJ = (Ours) Full Juneau

Juneau is multiple orders of magnitude faster than the baseline LSH ensemble.

How do we measure the result quality?

- Data Augmentation
 - Precision of classifier trained on augmented data
- Feature Extraction
 - Precision of classifier trained on extracted features
- Data Cleaning
 - Precision of classifier trained on cleaned data
- Linkable Data
 - Check if the returned table R is joined with T

(T is some table joined with the query table S in the original notebook workflow)



Figure 3: MAP@K of tables returned by different search classes

- NB = Original performance in the notebook
- KS = Keyword Search
- L64 = LSH Ensemble
- **PS** = (Ours) Row/Column Similarity Only
- PS + J = (Ours) R/C Sim. + Provenance Sim.
- full SJ = (Ours) Full Juneau



Figure 3: MAP@K of tables returned by different search classes

1. **Data Augmentation:** Significantly better than the baseline and other strategies

- NB = Original performance in the notebook
- KS = Keyword Search
- L64 = LSH Ensemble
- PS = (Ours) Row/Column Similarity Only
- PS + J = (Ours) R/C Sim. + Provenance Sim.
- full SJ = (Ours) Full Juneau



Figure 3: MAP@K of tables returned by different search classes

- 1. **Data Augmentation**: Significantly better than the baseline and other strategies
- 2. Feature Extraction: Slightly improved results

- NB = Original performance in the notebook
- KS = Keyword Search
- L64 = LSH Ensemble
- PS = (Ours) Row/Column Similarity Only
- PS + J = (Ours) R/C Sim. + Provenance Sim.
- full SJ = (Ours) Full Juneau



Figure 3: MAP@K of tables returned by different search classes

- 1. **Data Augmentation**: Significantly better than the baseline and other strategies
- 2. Feature Extraction: Slightly improved results
- 3. Data Cleaning: Measurable but minor impacts

- NB = Original performance in the notebook
- KS = Keyword Search
- L64 = LSH Ensemble
- PS = (Ours) Row/Column Similarity Only
- PS + J = (Ours) R/C Sim. + Provenance Sim.
- full SJ = (Ours) Full Juneau



Figure 3: MAP@K of tables returned by different search classes

- 1. **Data Augmentation**: Significantly better than the baseline and other strategies
- 2. Feature Extraction: Slightly improved results
- 3. Data Cleaning: Measurable but minor impacts
- 4. Linkable Data: Provides meaningful tables

- NB = Original performance in the notebook
- KS = Keyword Search
- L64 = LSH Ensemble
- PS = (Ours) Row/Column Similarity Only
- PS + J = (Ours) R/C Sim. + Provenance Sim.
- full SJ = (Ours) Full Juneau

Conclusions & Future Work

We contributed to the reusability and modularity for data processing through

- 1. Building a query-by-table framework for data lakes
- 2. Developing basic measures for table relatenesses and specialized compositional measures for 4 main use cases.
- 3. Developing algorithmic and indexing strategies for efficient top-k search
- 4. Showing promising scalability and quality for the search

Next Steps: scaling to distributed settings + heterogeneous data types

Demo



Figure 2: Demo Step 1. The user clicks on the Show Notebook Datasets button in the tool bar, and JUNEAU lists all datasets (tables) in the current notebook.
Demo

	File	e Edit	View Insert Cell Kernel V				Trusted	Python 3 O		
	8	+ %	선 🖪 🛧 🔸 H Run 🔳 C 🕨	Code	+ • • •	Additional Tr	aining/Val	idation Da	atasets	č – 3
lotebo	ok Data	asets		N WITH		rtable9_df_train_14	IdMSSubC 0 1 60 1 2 20 2 3 60	lassSaleConditi Normal Normal Normal	onSalePrice 208500 181500 223500	
Name 🜩	Туре 💠	Shape	Value	Search			3 4 70 4 5 60	Abnormi Normal	250000	
cm	ndarray	(10, 10)	[[1. 0.7909816 0.70862448 0	\$ = *			24524620 246247190	Normal Normal	241500 137000	
corrmat	DataFrame	(38, 38)	Id MSSubClass	‡ ⊨ =			24724820 24824960	Normal Normal	140000 180000	
data	DataFrame	(1460, 2)	SalePrice YearBuilt 0 2	‡ ⊨⇒ ≈		rtable10_df_train_14	24925050 IdMSSubC	Normal	277000 on SalePrice	
df_train	DataFrame	(1460, 81)	Id MSSubClass MSZoning LotF	t 🛏 🖛	ardScaler		0 1 60 1 2 20	Normal Normal	0.355663 0.000637	
missing_data	DataFrame	(81, 2)	Total Percent PoolQC	‡ ⊨⇒ =			2 3 60 3 4 70 4 5 60	Normal Abnorml Normal	0.552899 -0.545050 0.901350	
percent	Series	(81,)	PoolQC 0.995205 MiscFeature	\$ ⊨⇒ =			245 24620	Normal	0.789583	
otal	Series	(81,)	PoolQC 1453 MiscFeature	‡ ⊨> =	Pycham Projects note		24724820 24824960 24925050	Normal Normal Normal	-0.545050 -0.019086 1.256375	
		In [3]:	<pre>#skewness and kurtosis print("Skewness: %f" % df_trai print("Kurtosis: %f" % df trai</pre>	n['SalePr: n['SalePr:	<pre>ice'].skew()) ice'].kurt())</pre>	rtable15_df_train_16	IdMSSubC 0 1 60 1 2 20	lassSaleConditi Normal Normal	on SalePrice 12.247694 12.109011	

Figure 3: Demo Step 2. The user selects a table with a search mode, and the system will rapidly return a ranked list of tables that are related. In this example, the user selects "D", which means the user is looking for augmenting training or validation data.

Demo

e Edit View	Inser	t Cell Ke	ernel Wid	gets He	elp	Additional Tra	aining/Va	alidation D	atasets	<u>c</u> –	x on 3
+ 🗶 🖄 🖪	• •	🖌 🕅 Run 🔳	C #	Code	\$	Table	Additional cor	ntent			
	FUUNAU	1400 0.000200					-				
Mise	cFeature	1406 0.963014				rtable9_df_train_14	IdMSSub	bClassSaleCondit	tionSalePrice		
1110	oreature						0 1 60	Normal	208500		
	Alley	1369 0.937671					1 2 20	Normal	181500		
	Ferrer	1170 0.907524					2 3 60	Normal	223500		
	rence	11/9 0.00/534					3 4 70	Abnorml	140000		
Fire	placeQu	690 0.472603					4 5 60	Normal	250000		
							24524620	Normal	241500		
							246247100	Normal	137000		
In [17]: #Imp	port Ne	w Table					24724820	Normal	140000		
	-						24824960	Normal	180000		
							F-101-1000		100000		
Tn [181. ong	= conn	ect2db()					24925050	Normal	277000		
In [18]: eng	= conn	ect2db() d read sol t	able/'rta	ble9 df	train 1	ong schams s	249 25050	Normal	277000		
In [18]: eng df_1	= conn new = p	ect2db() d.read_sql_t	able(<mark>'rta</mark>	ble9_df_	train_1	rtable10_df_train_14	24925050 IdMSSub	Normal	277000		
In [18]: eng df_1	= conn new = p	ect2db() d.read_sql_t	able('rta	ble9_df_	train_1	rtable10_df_train_14	24925050 IdMSSut 0 1 60	Normal bClassSaleCondit Normal	277000 tion SalePrice 0.355663		
In [18]: eng df_u In [19]: prin	= conn new = p nt(df_n	ect2db() d.read_sql_t ew)	able('rta	ble9_df_	train_1	rtable10_df_train_14	24925050 IdMSSut 0 1 60 1 2 20	Normal bClassSaleCondit Normal Normal	277000 tion SalePrice 0.355663 0.000637		
In [18]: eng df_1 In [19]: prin	= conn new = p nt(df_n	ect2db() d.read_sql_t ew)	able('rta	ble9_df_	train_1	rtable10_df_train_14	24925050 IdMSSut 0 1 60 1 2 20 2 3 60	Normal bClassSaleCondit Normal Normal Normal	277000 tion SalePrice 0.355663 0.000637 0.552899		
In [18]: eng df_1 In [19]: prin	= conn new = p nt(df_n Id	ect2db() d.read_sql_t ew) MSSubClass	able('rta	ble9_df_	train_1 Street	rtable10_df_train_14	24925050 IdMSSut 0 1 60 1 2 20 2 3 60 3 4 70	Normal bClassSaleCondit Normal Normal Normal Abnorml	277000 tion SalePrice 0.355663 0.000637 0.552899 -0.545050		
In [18]: eng df_1 In [19]: prin	= conn new = p nt(df_n Id 1	ect2db() d.read_sql_t ew) MSSubClass 60	able('rta MSZoning RL	ble9_df_ LotArea 8450	train_1 Street Pave	rtable10_df_train_14 LotShape LandCon Reg	24925050 IdMSSut 0 1 60 1 2 20 2 3 60 3 4 70 4 5 60	Normal	277000 tionSalePrice 0.355663 0.000637 0.552899 -0.545050 0.901350		
In [18]: eng df_ In [19]: prin 0 1	= conn new = p nt(df_n Id 1 2	ect2db() d.read_sql_t ew) MSSubClass 60 20	able('rta MSZoning RL RL	ble9_df_ LotArea 8450 9600	train_1 Street Pave Pave	rtable10_df_train_14	24925050 IdMSSut 0 1 60 1 2 20 2 3 60 3 4 70 4 5 60 	Normal bClassSaleCondit Normal Normal Normal Normal 1PubNormal 1PubNormal	277000 tion SalePrice 0.355663 0.000637 0.552899 -0.545050 0.901350 2.300500		
In [18]: eng df_1 In [19]: prin 0 1 2	= conn new = p nt(df_n Id 1 2 3	ect2db() d.read_sql_t ew) MSSubClass 60 20 60	able('rta MSZoning RL RL RL	ble9_df_ LotArea 8450 9600 11250	train_1 Street Pave Pave Pave	rtable10_df_train_14 Lot Shape LandCon Reg Reg IR1	24925050 IdMSSut 0 1 60 1 2 20 2 3 60 3 4 70 4 5 60 	Normal bolassSaleCondit Normal Normal Normal Normal Normal 	277000 tion SalePrice 0.355663 0.000637 0.552899 -0.545050 0.901350 0.789583 0.584407		
In [18]: eng df_1 In [19]: prin 0 1 2 3	= conn new = p nt(df_n Id 1 2 3 4	ect2db() d.read_sql_t ew) MSSubClass 60 20 60 70	able('rta MSZoning RL RL RL RL	ble9_df_ LotArea 8450 9600 11250 9550	train_1 Street Pave Pave Pave Pave	rtable10_df_train_14 Lot Shape LandCon Reg Reg IR1 IR1	24925050 IdMSSut 0 1 60 1 2 20 2 3 60 3 4 70 4 5 60 	Normal bClassSaleCondit Normal Normal Normal Normal IPHNormal IPHNormal Normal	277000 tion SalePrice 0.355663 0.000637 0.552899 -0.545050 0.901350 0.789583 -0.584497 0.548497 0.584497		
In [18]: eng df_ In [19]: prin 0 1 2 3 4	= conn new = p nt(df_n Id 1 2 3 4 5	ect2db() d.read_sql_t ew) MSSubClass 60 20 60 70 60	able('rta MSZoning RL RL RL RL RL RL	ble9_df_ LotArea 8450 9600 11250 9550 14260	Street Pave Pave Pave Pave Pave Pave	rtable10_df_train_14 LotShape LandCon Reg Reg IR1 IR1	24925050 IdMSSut 0 1 60 1 2 20 2 3 60 3 4 70 4 5 60 	Normal Normal Normal Normal Normal Normal Normal Normal Normal Normal Normal Normal Normal	277000 tion SalePrice 0.355663 0.000837 0.552899 -0.545050 0.901350 0.789583 -0.584497 -0.545050 0.901926		
In [18]: eng df_1 In [19]: prin 0 1 2 3 4 5	= conn new = p nt(df_n Id 1 2 3 4 5 6	ect2db() d.read_sql_t ew) MSSubClass 60 20 60 70 60 50	able('rta MSZoning RL RL RL RL RL RL	LotArea 8450 9600 11250 9550 14260 14115	Street Pave Pave Pave Pave Pave Pave Pave	rtable10_df_train_14	24925050 IdMSSut 0 1 60 1 2 20 2 3 60 3 4 70 4 5 60 	Normal bClassSaleCondit Normal Normal Normal Normal Normal Normal Normal Normal Normal Normal Normal	277000 tion SalePrice 0.355633 0.000637 0.552899 -0.545050 0.901350 0.789583 -0.584497 -0.545050 -0.019086 1.06975		
In [18]: eng df_1 In [19]: prin 0 1 2 3 4 5 6	= conn new = p nt(df_n Id 1 2 3 4 5 6 7	ect2db() d.read_sql_t ew) MSSubClass 60 20 60 70 60 50 20	able('rta MSZoning RL RL RL RL RL RL RL	ble9_df_ LotArea 8450 9600 11250 9550 14260 14115 10084	Street Pave Pave Pave Pave Pave Pave Pave Pave	rtable10_df_train_14	24925050 IdMSSut 0 1 60 1 2 20 2 3 60 3 4 70 4 5 60 	Normal bClassSaleCondit Normal Normal Normal Normal Normal Normal Normal Normal Normal Normal Normal	277000 tion SalePrice 0.355663 0.000637 0.552899 -0.545050 0.901350 0.789583 -0.584497 -0.545050 -0.019086 1.256375		
In [18]: eng df_1 In [19]: prin 0 1 2 3 4 5 6 7	= conn new = p nt(df_n Id 1 2 3 4 4 5 6 7 8	ect2db() d.read_sql_t ew) MSSubClass 60 20 60 70 60 50 20 60	able('rta MSZoning RL RL RL RL RL RL RL RL	LotArea 8450 9600 11250 9550 14260 14115 10084 10382	Street Pave Pave Pave Pave Pave Pave Pave Pave	rtable10_df_train_14	24925050 IdMSSut 0 1 60 1 2 20 2 3 60 3 4 70 4 5 60 	Normal Normal Normal Normal Normal Normal Normal Normal Normal Normal Normal Normal Normal Normal Normal	277000 tion SalePrice 0.355663 0.000637 0.552899 -0.545050 0.901350 0.789583 -0.584497 -0.545050 -0.019086 1.256375 tion SalePrice		
In [18]: eng df_1 In [19]: prin 0 1 2 3 4 5 6 7 8	= conn new = p nt(df_n Id 1 2 3 4 5 6 7 8 9	ect2db() d.read_sql_t ew) MSSubClass 60 20 60 70 60 50 20 60 50 20 60 50	able('rta MSZoning RL RL RL RL RL RL RL RL RL RL RL RL	ble9_df_ LotArea 8450 9600 1250 14260 14115 10084 10382 6120	train_1 Street Pave Pave Pave Pave Pave Pave Pave Pave	rtable10_df_train_14	24925050 IdMSSut 0 1 60 1 2 20 2 3 60 3 4 70 4 5 60 	Normal bClassSaleCondit Normal Normal Normal Normal Normal Normal Normal Normal Normal Normal Normal Normal	277000 tion SalePrice 0.355663 0.000637 0.552899 -0.545050 0.901350 0.799583 -0.584497 -0.584497 -0.545050 -0.019086 1.256375 tion SalePrice 12.247694		
In [18]: eng df_1 In [19]: prin 0 1 2 3 4 5 6 7 8 9	= conn new = p nt(df_n Id 1 2 3 4 5 6 7 8 9 10	ect2db() d.read_sql_t ew) MSSubClass 60 20 60 70 60 50 20 60 50 20 60 50	able('rta MSZoning RL RL RL RL RL RL RL RL RL RL	LotArea 8450 9600 11250 9550 14260 14115 10084 10382 6120 7420	train_1 Street Pave Pave Pave Pave Pave Pave Pave Pave	rtable10_df_train_14	24925050 IdMSSut 0 1 60 1 2 20 2 3 60 3 4 70 4 5 60 	Normal Normal Normal Normal Normal Normal Normal Normal Normal Normal Normal Normal Normal Normal Normal Normal Normal Normal Normal Normal	277000 tion SalePrice 0.355663 0.000637 0.552899 -0.545050 0.901350 0.789583 -0.584497 -0.545050 -0.019086 1.256375 tion SalePrice 12.247694 12.109011		
In [18]: eng df_1 In [19]: prin 0 1 2 3 4 5 6 7 8 9 10	= conn new = p nt(df_n Id 1 2 3 4 4 5 6 7 8 9 10 0 11	ect2db() d.read_sql_t ew) MSSubClass 60 20 60 70 60 50 20 60 50 20 60 50 20 60 50 20 60 50 20 60 50 20 60 50 20 20 20 20 20 20 20 20 20 20 20 20 20	able('rta MSZoning RL RL RL RL RL RL RL RL RL RL RL RL	ble9_df_ LotArea 8450 9600 11250 9550 14260 14215 10084 10382 6120 7420 11200	street Pave Pave Pave Pave Pave Pave Pave Pave	rtable10_df_train_14	24925050 IdMSSut 0 1 60 1 2 20 2 3 60 3 4 70 4 5 60 24524620 24524620 24524620 24524820 24524820 24524820 24524820 2452495050 IdMSSut 0 1 60 1 2 20	Normal bClassSaleCondit Normal Normal Normal Normal Normal Normal Normal Normal Normal Normal Normal Normal	277000 tion SalePrice 0.355663 0.000637 0.552899 -0.545050 0.901350 0.789583 -0.584497 -0.584497 -0.545050 -0.019086 1.256375 tion SalePrice 12.247694 12.109011		
In [18]: eng df_1 In [19]: prin 0 1 2 3 4 5 6 7 8 9 10 1 1	= conn new = p nt(df_n Id 1 2 3 4 5 6 6 7 7 8 9 10 11	ect2db() d.read_sql_t ew) MSSubClass 60 20 60 70 60 50 20 60 50 20 60 50 20 60 50 20 60 50 60 50 60 50 60 50 60 50 60 50 60 50 60 50 60 50 60 50 60 50 60 50 60 50 60 50 60 50 60 50 60 50 60 50 50 50 60 50 50 50 50 50 50 50 50 50 50 50 50 50	able('rta MSZoning RL RL RL RL RL RL RL RL RL RL RL	LotArea 8450 9600 11250 9550 14260 14115 10084 10382 6120 7420 11220	train_1 Street Pave Pave Pave Pave Pave Pave Pave Pave	rtable10_df_train_14	24925050 IdMSSut 0 1 60 1 2 20 2 3 60 3 4 70 4 5 60 	Normal Normal	277000 tion SalePrice 0.355663 0.000637 0.552899 -0.545050 0.901350 0.789583 -0.584497 -0.545050 -0.019086 1.256375 tion SalePrice 12.247694 12.109011		4

Figure 5: Demo Step 3. The user clicks on a table of interest, and JUNEAU adds a new cell to the notebook that loads the table.

Finding Related Tables in Data Lakes for Interactive Data Science

From the archaeologist's perspective

Yanhao Wang

The problem space: Data Lake Systems

schema-on-use



Hai, Rihan, Christoph Quix, and Matthias Jarke. "Data lake concept and systems: a survey." arXiv preprint arXiv:2106.09592 (2021).

The problem space: Related Dataset Discovery



JOSIE [Zhu, 2019]

- Problem: Joinable Table Discovery
 - Given (T, C), return top-K tables that can be joined with T on C
- Focus on data-value overlap
- New perspective: Consider the table columns as sets, and same tuple values as the set intersection, and formalize the problem as *overlap set similarity* problem.
- Apply existing overlap set similarity search techniques in data lakes, resolving the unique problems in this scenario:
 - Large number of tables in a massive data lake
 - Hard to decide the threshold for the intersection size threshold

Erkang Zhu, Dong Deng, Fatemeh Nargesian, and Renée J Miller. 2019. JOSIE: Overlap Set Similarity Search for Finding Joinable Tables in Data Lakes. In Proceedings of the 2019 International Conference on Management of Data. ACM, 847–864.

Aurum [Fernandez, 2018]

SRQL Query SRQL Engine 1 Builder 2 Profiler

- Problem: Joinable Table Discovery
- Indexing using data profiles and sketches
- Use hypergraphs to find similar datasets, where nodes are columns
 - Profiles each column with a signature f(card., data distr., min-hash, ...)
 - Index signatures in *Locality-Sensitive Hashing (LSH)*
 - Assign a weighted edge between two column (nodes) if their index fall into the same bucket, with the similarity value as the edge weight
 - Construct hyperedges based on these weighted edges
- Fast and robust against data value changes

Raul Castro Fernandez, Ziawasch Abedjan, Famien Koko, Gina Yuan, Samuel Madden, and Michael Stonebraker. 2018. Aurum: A data discovery system. In 2018 IEEE 34th International Conference on Data Engineering (ICDE). IEEE, 1001–1012.

hash functions with high collision probability for similar inputs *E.g. MinHash, random projections, etc.*

D3L [Bogatu, 2020]

- Problem: Related Table Discovery
 - Given T, return top-K tables that contain relevant attributes for populating T
- Extended types of similarity computed at schema- and instance-level
 - Representation patterns of instance values
 - Semantics of textual attributes
 - Data distribution for numerical values
- Map these features (after LSH) to a 5-d Euclidean space for calculating relatedness

Juneau [Zhang, 2020]

- Problem: (Task-specific) Related Table Discovery
 - \circ Given (T, τ), return top-K tables that are most relevant to T in terms of search type τ
- Additionally consider the **context** (table provenance) and **intent** (task) of the user
- Further extended the notion of relatedness, on top of the often-used value overlap and attribute overlap:
 - New attribute/instance rate
 - Provenance similarity
 - Description similarity
- Introduced workflow graph & variable dependency graph as auxiliary structures to support provenance similarity computation
- Proposed a practical framework to **integrate a subset of relatedness metrics** to address different data science tasks

		Query-driven]	Discover Related	ł		Dataset Preparation
Svstems	Target Pro	blem	Relatedness Criteria	Datasets	Similarity metrics	Auxilia	and Organization
Aurum	Joinable Ta Discovery	able	Instance value overlap Attribute name Primary-key/foreign-key cane	didate	Jaccard similarity (MinHash) Cosine similarity (TF-IDF)	Hypergraph	
JOSIE	Joinable Ta Discovery	able	Instance value overlap		Intersection size of sets	Inverte	ed Index
D3L	Related Ta	ble Discovery	Instance value overlap Attribute name Semantics Data value representation pa (Numerical) data distribution	ttern	Jaccard similarity (MinHash) Cosine similarity (Random projections)	5–dimensional Euclidean space	
Juneau	(Task-spec Table Disco	ific) Related overy	Instance value overlap Domain overlap Attribute name Key constraint New attributes rate New instance rate Variable dependency Descriptive metadata Null Values		Jaccard similarity	Workfl Variab graph	ow graph le dependency

Hai, Rihan, Christoph Quix, and Matthias Jarke. "Data lake concept and systems: a survey." arXiv preprint arXiv:2106.09592 (2021).

The problem space: Related Dataset Discovery



Edmint [Amsterdamer, 2022]

- Problem: Given a dataset T, automatically select multiple datasets for extension
 - The value of each potential dataset is dependent on the choice of other datasets
- One important metric is still relatedness, and the author uses the relatedness metrics proposed by [Zhang, 2020] to integrate relatedness from different dimensions
- Not only consider the relatedness from the search table, but also balancing:
 - The integration gain: number of acquired table cells that are truthful and correctly integrated
 - The integration cost: decrease in quality, increase of incompleteness

Finding Related Tables in Data Lakes for Interactive Data Science

Reviewer Role : Vishnu K Krishnan

Paper Summary

What is the paper about? –

- Developed framework Juneau to support search for semantically related tables.
- Focused towards data science use-cases cleaning, augmenting, etc.
- What does it solve?
 - It works to aid data scientists through promotion of reuse of data/workflows, consistency in processing and reduce redundant and inefficient work.
- How It functions? -
 - Usage of top-k, pruning and approximation strategies to return the most relatable tables

Strong Points

- Identified the bottleneck while doing top-k search relational mapping and pruned the candidate tables effectively to Improve the top-k performance.
- Threshold algorithm's on the fly calculation to save on precomputation costs.
- Novel method of profiling tables using various overlapping methods and definitions. Provided clear formal definitions.
- Proposed workaround for the NP-hard subset selection problem by limiting a variable to a small integer.(choosing key mapping to ensure maximum overlap – small key mapping size allows limit)
- Usage of Data profiles, compositional profiles and workflow graphs to speed up computation.
- Utilized Workflow graphs as directed bi-partite graphs for Juneau computational notebook software. Computational notebook platforms directly capture workflows and data exploration by interleaving source code cells with their outputs.

Opportunities of Improvement

- More representations, more data points and performance figures could have been provided.
- Can it be used in the industry? Privacy issues for example, if this were to be utilized by industry practitioners where some data could be considered confidential, can the current solution be augmented to ensure privacy of the data?
- Trade-off precision for speed using approximation in the threshold algorithm. (early stop used)

Accept

Data Lake vs Data Warehouse

Data lake:

raw data, can be unstructured low-cost storage, but no transactions, data quality checks data scientists and engineers

Data warehouse:

structured data (schema-on-write) expensive for large data volumes managers and business analysts

"Lakehouse"

Source: Lakehouse: A New Generation of Open Platforms that Unify Data Warehousing and Advanced Analytics

Contribution/Strengths

Novel problem formulation driven by use cases

- Query by table framework they hypothesized that interactive DS users search on data lake not only by keyword but also by table.
- More prior work on table search focus on the web
- Despite the difference, the searching for all four tasks can be captured under a general framework and solved by Juneau. This applicationdriven reformulation of the table searching problem is novel and inspiring.
- The authors exploited the characteristics of different use cases (e.g., data augmentation vs feature extraction) to establish measures that capture the relatedness of tables.

Contribution/Strengths

Implementation and evaluation

- Integrated with Jupyter Notebook (among others) makes adoption easy
- It leverages multiple indices to perform staged relation mapping detection to speed up the topk search: data profile index on columns, compositional profile index on co-occurring sets of columns, and workflow graph index.
- Juneau evaluates the results on objective measures, for example, the precision of the classifier on the original table vs the precision on the table returned by Juneau.
- Evaluated using real data science workflows from kaggle.com.

Lack of user perception

- Given that targeted usage scenarios are the highlight of the table relatedness measure (thus usefulness of the result), testing out if users actually define related tables in the same way as the relatedness measure is crucial
- No qualitative measures with data scientists How do actual data scientists interpret the search results? Do they feel comfortable integrating search results into their data pipelines? Are there times that they like the results but they don't perform well, and vice versa?

Unclear generalization

- All the metrics and searching methods are specifically designed for these mappings, it remains unclear how easy the proposed system generalizes to new settings.
- So, the query is based on a table and a specific task at hand, what if I wish to perform task 1 followed by task 3 maybe the query can combine two tasks and return the best-suited table for this.
- Assumption that cells in a notebook are in sequence is too optimistic. Notebooks are designed for random and unstructured runs, so assuming the order is a big assumption.

Weak baseline

- In the experiments, only a brute-force baseline is compared for efficiency. There is no state-of-the-art solution for finding similar joinable or unionable tables
- The experiments section does not sufficiently compare against other works - the authors mention several related approaches like Google Dataset search, q system, etc.
- The experiments do not conclusively show that Juneau is superior to keywork search. Although Juneau has slightly higher accuracy, its query times are much longer than keyword search.
- The corpus used in the evaluation consists of data ~5GB. This is pretty small in comparison to the sizes of data stored on data lakes in the real world, which might make the results biased.

Data quality issues

- Sometimes people leave trails of abandoned methods that are just not correct, outdated, or contain bugs. Storing and indexing these seems wasteful, and can also lead to poor/misleading search results.
- I wonder how they have handled different versions of the data are all available or just the latest one?
- Data cleaning and augmentation can be quite task specific, but the search doesn't explicitly encode the task, and incorrect methods can slip in.

Next class

Auto-Suggest: Learning-to-Recommend Data Preparation Steps Using Data Science Notebooks

- Author: Bojun, Siddhi
- Reviewer: Shubham, Shen En
- Archaeologist: Aniruddha
- Practioner: Jingfan
- Researcher: Ting