Advanced Data Management (CSCI 640/490)

Scalable Dataframes

Dr. David Koop





History of Dataframes

- R, open-source alternative to S, developed in 2000 (with dataframes)
- Pandas, 2009
- Spark, 2010 (resilient distributed dataset [RDD], Dataset API)

D. Koop, CSCI 640/490, Spring 2024

• Originally in Statistical Models in S, [J. M. Chambers & T. J. Hastie, 1992]











Pandas Workflow: Ingest, Cleaning, Analysis

R1. Read HTML

4.6

4.7

import pandas as pd
products = pd.read_html(...) products

 iPhone 11	iPhone Pro Max	iPhone 11 Pro	
 6.1-inch	6.5-inch	5.8-inch	Display
 Dual 12MP	Triple 12MP	Triple 12MP	Camera
 7MP	12MP	120MP	Front Camera

		(C4. Read	Exce
pric pric	es = pd .read_ es	_excel()	
		Price	Rating	
	iPhone 11 Pro	999.00	4.5	
	iPhone Pro Max	1099.00	5.0	

iPhone 11 699.99

iPhone XS 999.99

	C1. Ord	lered point	updates			C2.	Matrix-lik	e t	ranspose
product product	s .iloc [2, s	0] = "12	MP"		products = products	produc	ts .T		
i							0		Wireless
	Phone 11 Pro I	Phone Pro Max	iPhone 11			Display	Camera		Charging
Display	5.8-inch	Phone Pro Max 6.5-inch	iPhone 11 6.1-inch	 	iPhone 11 Pro	Display	Triple 12MP		Charging Yes
Display Camera	5.8-inch Triple 12MP	6.5-inch Triple 12MP	iPhone 11 6.1-inch Dual 12MP	 	iPhone 11 Pro	Display 5.8-inch	Triple 12MP		Charging Yes Yes
Display Camera Front Camera	5.8-inch Triple 12MP 12MP	6.5-inch Triple 12MP 12MP	iPhone 11 6.1-inch Dual 12MP 7MP	 	iPhone 11 Pro iPhone Pro Max iPhone 11	Display 5.8-inch 6.5-inch 6.1-inch	Triple 12MP Triple 12MP Dual 12MP	••• •••	Charging Yes Yes Yes

```
one_hot_df = pd.get_dummies(prod
iphone_df = prices.merge(
      one_hot_df,
left_index=True, right_index
iphone_df
```

D. Koop, CSCI 640/490, Spring 2024

A1. One-to-many column mapping A2. Joins

lucts)		Price	Rating	Wireless Charging	Display_\ 5.8-inch	
	iPhone 11 Pro	999.00	4.5	1	1	
(=True	iPhone Pro Max	1099.00	5.0	1	0	
	iPhone 11	699.99	4.6	1	0	
	iPhone XS	999.99	4.7	0	1	

		C3.	Column t	ran	sfoi	rma	tio
pr [pr	oducts = p "Wireless" lambda x: oducts	oroduc Charg 1 if	ts\ ing"] .ma x is "Ye	p (s"	el	se	0)
		Display	Camera		Wire Cha	eless rging	9
	iPhone 11 Pro	5.8-inch	Triple 12MP			-	1

iPhone Pro Max 6.5-inch	Triple 12MP	 1
iPhone 11 6.1-inch	Dual 12MP	 1
iPhone XS 5.8-inch	Dual 12MP	 0

A3. Matrix Covariance

iphone_df.cov()
iphone_df

29868.3	19.967	
19.967	0.0466667	
16.8317	-7.40149e-17	
33.3333	-0.0666667	
-	29868.3 19.967 16.8317 33.3333 	29868.3 19.967 19.967 0.0466667 16.8317 -7.40149e-17 33.3333 -0.06666667













Problems Scaling: From Pandas to Other Solutions







Modin as a Solution





Modin Positioning













Dataframe Data Model



- Combines parts of matrices, databases, and spreadsheets
- Ordered, but not necessarily sorted
 - Rows and columns
- No predefined schema necessary
 - Types can be induced at runtime
- Typed Row/column labels
 - Labels can become data
- Indexing by label or row/column number
 - "Named notation" or "Positional notation"







Comparing Dataframes and Relational Stores

- Dataframe Characteristics
 - Ordered table
 - Named rows labels
 - A lazily-induced schema
 - Column names from $d \in Dom$
 - Column/row symmetry
 - Support for linear alg. operators

- Relational Characteristics
 - Unordered table
 - No naming of rows
 - Rigid schema
 - Column names from att
 - Columns and rows are distinct
 - No native support











Comparing Dataframes and Matrices

- Dataframe Characteristics
 - Heterogeneously typed
 - Numeric & non-numeric types
 - Explicit row and column labels
 - Support for rel. algebra operators

- Matrix Characteristics
 - Homogeneously typed
 - Only numeric types
 - No row or column labels
 - No native support











Dataframe Algebra

Operator	(Met	a)data	Schema	Origin	Order	Description
SELECTION		×	static	REL	Parent	Eliminate rows
PROJECTION		×	static	REL	Parent	Eliminate columns
UNION		×	static	REL	Parent [†]	Set union of two dataframes
DIFFERENCE		×	static	REL	Parent [†]	Set difference of two dataframes
CROSS PRODUCT / JOIN		×	static	REL	Parent [†]	Combine two dataframes by element
DROP DUPLICATES		×	static	REL	Parent	Remove duplicate rows
GROUPBY		×	static	REL	New	Group identical values for a given (set of) attribute(s)
SORT		×	static	REL	New	Lexicographically order rows
RENAME	(\times)		static	REL	Parent	Change the name of a column
WINDOW		×	static	SQL	Parent	Apply a function via a sliding-window (either direction
TRANSPOSE	(\times)	×	dynamic	DF	Parent [◊]	Swap data and metadata between rows and columns
MAP	(\times)	×	dynamic	DF	Parent	Apply a function uniformly to every row
TOLABELS	(\times)	×	dynamic	DF	Parent	Set a data column as the row labels column
FROMLABELS	(\times)	×	dynamic	DF	Parent	Convert the row labels column into a data column



















Pivot Example

Wide Table of MONT

ЪТ	T 1 1 / C			Month	2001	2002	
Narro	w Table (S	ALES)	, Γ	Jan	100	150	
Year	Month	Sales		Feb	110	200	
2001	Jan	100	-	Mar	120	250	
2001	Feb	110					
2001	Mar	120		$\mathbf{Pivot} \longrightarrow$			
2002	Jan	150		← Unpivot			
2002	Feb	200					
2002	Mar	250		Year	Jan	Feb	Γ
2003	Jan	300		2001	100	110	
2003	Feb	310		2002	150	200	
		<u> </u>	l	2003	300	310	N

Wide Table of YEARs

D. Koop, CSCI 640/490, Spring 2024



[D. Petersohn et al., 2020]



Northern Illinois University



11

Modin Challenges

- Massive API: 240+ operators, but with a lot of redundancy Parallel Execution: row-based, column-based, and block-based Data Model Challenges: Schema induction, reusing type info

- Order is important
- Supporting billions of columns: Row/Column equivalence (transpose)
- Metadata is data (and vice versa)
- Users want immediate feedback
- Users want to create queries incrementally





<u>Assignment 4</u>

- Work on Data Integration and Data Fusion
- Integrate artist datasets from different institutions (Met, NGA, AIC, CMA)
 - Integrate information based on ids and matching
- Record Matching:
 - Which artists are the same?
- Data Fusion:
 - Names
 - Dates
 - Nationalities





Test 2

- Upcoming... April 8
- Similar format, but more emphasis research papers

• Similar format, but more emphasis on topics we have covered including the





Dataframes, Databases, and the Cloud

- How do we take advantage of different architectures? Lots of work in scaling databases and specialized computational engines • What is the code that people actually write?





Magpie: Python at Speed and Scale using Cloud Backends

A. Jindal





Data Science Jungle



D. Koop, CSCI 640/490, Spring 2024



Northern Illinois University

Magpie Goals



D. Koop, CSCI 640/490, Spring 2024









ConnectorX: Databases to Dataframes

- Write read_sql queries but write SQL
- Written in Rust
- Returns a dataframe

```
query = f"""
SELECT 1 orderkey,
 SUM(l_extendedprice * (1 - l_discount)) AS revenue,
 o orderdate,
 o shippriority
FROM customer,
 orders,
 lineitem
WHERE c mktsegment = 'BUILDING'
AND c custkey = o custkey
 AND 1 orderkey = o orderkey
 AND o orderdate < DATE '1995-03-15'
 AND 1 shipdate > DATE '1995-03-15'
GROUP BY 1 orderkey,
 o orderdate,
 o shippriority
\\ // //
df = read sql("postgresql://postgres:postgres@localhost:5432/tpch", query,
              partition on="l orderkey", partition num=4)
```





ConnectorX Speed & Memory











Improvements in ConnectorX

- Written in native language (Rust)
- Copy exactly once (even during parallel computations)
- CPU cache-friendly: process in a streaming fashion









An Opinionated Introduction to Polars

Nico Kreiling





Handling Large Data with Polars

Etienne Bacher





Discussion

- Data in the cloud and local exploration
- Languages: SQL or Pandas or Ibis or....?







