Advanced Data Management (CSCI 680/490)

Time Series Data

Dr. David Koop





Spatial Data

Measure vegetation density



Measure snow melt



D. Koop, CSCI 680/490, Spring 2022

Track hurricanes



Track phytoplankton populations













Interactive Exploration of Spatial Data













Interactive Exploration of Spatial Data













Two Inputs to Exploratory Browsing Input Compute Respond Compute Respond Input Prepare data User Fetch results User Update Create in DBMS submits from DBMS visualization visualization pan/zoom query (Pre-comp. Structures)



Cold start time

D. Koop, CSCI 680/490, Spring 2022

interaction latency < 500ms









Systems for Interactive Exploration

		(Offline) Pre-computed structures	(Before interaction) Predictive	(After interaction) Progressive/Incremental
ormat	Sampling		DICE (ICI A-WARE (H	SampleAction (CHI 2012) Vizdom (VLDB 2015) DE 2014) HILDA 2016)
ut but f	gation	Nanocubes (Infovis 2013) imMens (Eurovis 2013)	ATLAS (VAST 2008) XmdvTool (<i>DASFAA</i> 2003)	
õ	Aggre	ForeCa	che	

D. Koop, CSCI 680/490, Spring 2022

Time





Nanocubes



Linked view of tweets in San Diego, US







From Tables and Spreadsheets to Data Cubes

- data in the form of a data cube
- A data cube, such as sales, allows data to be modeled and viewed in multiple dimensions
 - Dimension tables, such as item (item_name, brand, type), or time(day, week, month, quarter, year)
 - Fact table contains measures (such as dollars_sold) and keys to each of the related dimension tables
- called the apex cuboid. The lattice of cuboids forms a data cube.

• A data warehouse is based on a multidimensional data model which views

 In data warehousing literature, an n-D base cube is called a base cuboid. The top most 0-D cuboid, which holds the highest-level of summarization, is





Data Cube: A Lattice of Cuboids



D. Koop, CSCI 680/490, Spring 2022

0-D (*apex*) cuboid

4-D (base) cuboid





8

Data Cube Measures: Three Categories

- without partitioning
 - E.g., count(), sum(), min(), max()
- distributive aggregate function
 - E.g., avg(), min_N(), standard_deviation()
- Holistic: if there is no constant bound on the storage size needed to describe a subaggregate.
 - E.g., median(), mode(), rank()

• **Distributive**: if the result derived by applying the function to n aggregate values is the same as that derived by applying the function on all the data

• Algebraic: if it can be computed by an algebraic function with M arguments (where M is a bounded integer), each of which is obtained by applying a





9

Multidimensional Data

Sales volume as a function of product, month, and region



D. Koop, CSCI 680/490, Spring 2022

Dimensions: *Product, Location, Time* **Hierarchical summarization paths**









A Sample Data Cube







OLAP Operations







Efficient Processing of OLAP Queries

- - e.g., dice = selection + projection
- Determine which materialized cuboid(s) for OLAP operation:
 - Query: {brand, province or state} With "year = 2004"
 - 4 materialized cuboids available:
 - 1. {year, item name, city}
 - 2. {year, brand, country}
 - 3. {year, brand, province or state}
 - 4. {item name, province or state} Where year = 2004
 - Which should be selected to process the query?

D. Koop, CSCI 680/490, Spring 2022

• Determine which operations should be performed on the available cuboids - Transform drill, roll, etc. into corresponding SQL and/or OLAP operations,





	ult. If we	anow en	tries in the rec	(oxtato)	••			U	
	d repre	sent this a	gereation as t		Country	v Device	Language	Count	
				2 WOU	All	All	All	5	
Dala Ul	JDECounty		Lange SFO		All	Android	All	2	
	the repart	on above	Fould be to sel	S follo	All	iPhone	All	3	
	se using o		h ^a this case, fi	NO WEI	All	All	eu	4	
	tds in avre	allowtha	tiontanthass	Schall	All	All	ru	1	
	Bauta Tatil		BALAROBAS		All	iPhone	ru	1	
			Cations for a	e wol	All	Android	en	2	
	Gountry	Device	Language and	given a	All	iPhone	en	2	
	All_1	· All		5 S	1 100		UII	_	
	seffelat	10n given	a list of attribu	utes a					
	X CONTRACT						Γ		peration where
				C					Il un on Davies
	Country	Device	Language (Count Sa					ll up on Device
	All	Android	en en en	uage 2		-quivalent to	o Group By	on	up by s on: (1)
	A11	iPhone	nö attribu	itęs; (all possible s	subsets of		e. Note that the
	A 11	iPhone	group by	on Lin		Device, Lan	guage}		As the results
	1111	II HOHe	of group	by's				VII 40 IV	lations, we can
		teifugetio		nes atta		threations	As we wi	11 descri	be <i>nanocubes</i> is
				verbed at	blotto	tube the store	and query	cubeso	f roll uns
	COUNTRY	IS HOLALIOP	'Langelage'	HALEIST	ane me	A BAR BEILLE		cubes 0	, <i>1011 ups</i> .
			and brouge		Nipicic			98.	
btributes and	ro lbase Andraia	ninger	odeona ka seneccia	2 Millical	dedlad	biogses ble s	subsets dP	IEWPU	AL NULL-UP
les derived from	a hase relation		ist of attribute	stand a	n agor	Device, Lan	guage}		
Wentfonal-ways	Of descripting		ions for a grv	Plazatta	tione.eg	<i>Logant</i> uter a	are necessa	rilv bou	nded by display
by, Cube, and to	DELEP. PAGIOU	orephopera	HUH HS ONE IN	WINCH	a regata	offas trahe a	hle to quic	klycolle	ect subspaces of
D. Koop, CSCI 680/	490, Spring 202		storatti 10ertes	stand a	naggie	gale in th			orthern Illinois Universi
The full of the that to	r-ALARVI dittar	Anticomhi	hatid to the	the clideo	Wath thin		ie same DD		C SUICCII. IIUW-



Building a Nanocube









<u>Assignment 5</u>

- Chicago Bike Sharing Data
 - Spatial Analysis
 - Temporal Analysis
 - Graph Database (neo4j)





Teaching Evaluations

• This Wednesday (April 20) in class





TopKube: A Rank-Aware Data Cube for Real-Time Exploration of Spatiotemporal Data

F. Miranda, L. Lins, J. T. Klosowski, and C. T. Silva





TopKube: What about Top-k and Rankings?

- Aggregates are interesting
- Also, often interested in top-k answers given particular criteria
- ... or rankings
- Search over time and space but find specific examples
- TopKube is a rank-aware data structure that computes top-k queries with low latency so interactive exploration is possible





Example: Basketball

Shots by time, number of points scored, and location on the court

team	player	time	pts	X		У
CLE	L. James		5	0	13	28
BOS	R. Rondo		5	2	38	26
CLE	L. James		7	3	42	35

- Query: Ranked list of the 50 players who took the most shots
 - SELECT player, count (*) AS shots FROM table GROUP BY player ORDER BY shots DESC LIMIT 50
- Query: Rank the top 50 players by points made:
 - SELECT player, sum (pts) AS points FROM table GROUP BY player ORDER BY points DESC LIMIT 50









Ranking by Shot Location



1.	Anthony Parker	93 1.	Jamal Crawford	113
2.	Rasual Butler	84 2.	Arron Afflalo	105
3.	Mickael Pietrus	84 3.	Rashard Lewis	98
4.	Jeff Green	82 4.	Martell Webster	96
5.	Jared Dudley	80 5.	Joe Johnson	86
6.	Jason Richardson	80 6.	Rasual Butler	85
7.	Carlos Delfino	76 7.	Jason Terry	81
8.	George Hill	75 8.	Anthony Parker	80
9.	Shane Battier	72 9.	Danilo Gallinari	78
10.	Joe Johnson	71 10.	George Hill	75
11.	Matt Barnes	68 11.	Ray Allen	69
12.	Brandon Rush	67 12.	Steve Blake	68
13.	Mo Williams	65 13.	Mickael Pietrus	68
14.	Steve Blake	63 14.	Mo Williams	63
15.	Arron Afflalo	63 15.	Keith Bogans	63
16.	Charlie Bell	63 16.	Anthony Morrow	62
17.	Courtney Lee	62 17.	Mike Bibby	62
18.	Stephen Jackson	61 18.	Al Harrington	62
19.	Marvin Williams	61 19.	Shane Battier	61
20.	Ray Allen	60 20.	Carlos Delfino	61

D. Koop, CSCI 680/490, Spring 2022



[F. Miranda et al., 2017]





TopKube vs. Nanocubes

- Product bin: the combination of selections from dimensions
- Nanocubes maps each product bin ((01,10), iPhone) to a time series $\beta \mapsto ((t_1, v_1), (t_2, v_1 + v_2), \dots, (t_m, v_1 + \dots + v_m))$
- TopKube maps each product bin to rank-aware multi-set $\beta \mapsto \left\{ \operatorname{lst} = ((q_1, v_1, \sigma_1), \dots, (q_j, v_j, \sigma_j)), \operatorname{su} \right\}$
- q_i is the ith smallest key that appears in product bin
- v_i is the value of the measure for key q_i in the product bin
- σ_i is the index of the key with its largest value

$$\operatorname{um} = \sum_{i=1}^{j} v_i \bigg\}$$









Example: One Spatial Dim. and A,B,C events



D. Koop, CSCI 680/490, Spring 2022



[F. Miranda et al., 2017]





Problem: Lots of Bins!





Three Algorithms to Merge Bins

- Threshold: don't do a full scan, use extra information about ranking Sweep: Use a priority queue where the product bin with the current smallest
- key is on the top
- Hybrid:
 - Threshold has best theoretical guarantee but some sparse cases can be faster
 - Use Sweep on small input lists, Threshold on denser problem









Top-edited Wikipages in Nevada and Mississippi



A Ø /m A
Baton_Rou
University_
Mississipp
Jackson,_
Louisiana_
Mississipp
WVLA-TV
Ole_Miss_
List_of_Sta
Louisiana
New_Orlea

D. Koop, CSCI 680/490, Spring 2022

⇔		
ouge,_Louisiana		323
/_of_Mississippi		230
pi		216
_Mississippi		208
_State_University		189
pi_State_University		169
/		158
_Rebels_football		155
tar_Wars_books		131
		122
eans_Saints	1	107

압☆ø/m⇔

Reno,_Nevada Early_Christianity Comparison_of_the_AK-Las_Vegas_Academy Timeline_of_Christianity. Las_Vegas Council_of_Jerusalem Paul_the_Apostle University_of_Nevada,_L Nevada Antinomianism



North

Briste



Geolocated Flickr tags in Africa



Top Hashtags in Paris related to Charlie Hebdo

D. Koop, CSCI 680/490, Spring 2022

Northern Illinois University

D. Koop, CSCI 680/490, Spring 2022

te

Sun Jan 18 04:00:00 2015

Mon Jan 26 16:00:0

	얍 ☆ Ø /m ⇔			얍 ☆ ∅ /m ⇔		
4102	dotfiles		10022	dotfiles		890
1780	rust	•	2341	postgresql	•	212
1297	llvm-project		2321	mongo	•	210
1155	llvm-project		2306	pubsub	8	176
1105	titanium_mob		2266	docs	1	124
1057	git	•	2031	neon-rtk-gps	1	122
961	node	1	1837	pandas	1	117
843	gaia	1	1830	homebrew	1	116
793	wireshark-ht	1	1745	julia	1	110
790	spark	1	1741	sample_app	I.	973
758	phabricator	1	1720	astrometry.n	I	948
738	mozilla-cent	1	1709	.emacs.d	I	947
728	zamboni	1	1631	VTK		921
717	gecko-dev	1	1531	website		858
707	docs	1	1522	Cinder		749

Evaluation

D. Koop, CSCI 680/490, Spring 2022

[F. Miranda et al., 2017]

D. Koop, CSCI 680/490, Spring 2022

Aggregation

Split-Apply-Combine

- Coined by H. Wickham, 2011
- Similar to Map (split+apply) Reduce (combine) paradigm
- The Pattern:
 - 1. **Split** the data by some grouping variable
 - 2. Apply some function to each group independently
 - 3. Combine the data into some output dataset
- The apply step is usually one of :
 - Aggregate
 - Transform
 - Filter

Split-Apply-Combine

D. Koop, CSCI 680/490, Spring 2022

[W. McKinney, Python for Data Analysis]

Splitting by Variables

name	age	sex
John	13	Male
Mary	15	Female
Alice	14	Female
Peter	13	Male
Roger	14	Male
Phyllis	13	Female

age	sex
13	Male
13	Male
14	Male
	age 13 13 14

name	age	sex
Mary	15	Female
Alice	14	Female
Phyllis	13	Female

D. Koop, CSCI 680/490, Spring 2022

	(sex	K)
_	\	_/

name	age	sex
John	13	Male
Peter	13	Male
Phyllis	13	Female

.(age)

name	age	sex
Alice	14	Female
Roger	14	Male

name	age	sex
Mary	15	Female

Apply+Combine: Counting

.(sex)

.(age)

sex	value	
Male	3	
Female	3	

age	
13	
14	
15	

D. Koop, CSCI 680/490, Spring 2022

.(sex, age)

value	sex	age	value
3	Male	13	2
2	Male	14	1
1	Female	13	1
	Female	14	1
	Female	15	1

[H. Wickham, 2011]

In Pandas

- groupby method creates a GroupBy object
- groupby doesn't actually compute anything until there is an apply/aggregate step or we wish to examine the groups
- Choose keys (columns) to group by
- size() is the count of each group









Aggregation

- Operations:
 - count()
 - mean()
 - sum()
- May also wish to aggregate only certain subsets
 - Use square brackets with column names
- Can also write your own functions for aggregation and pass then to agg function
 - def peak_to_peak(arr):
 return arr.max() arr.min()
 grouped.agg(peak_to_peak)





Optimized groupby methods

Function name	Description
count	Number of non-NA
SUM	Sum of non-NA val
mean	Mean of non-NA va
median	Arithmetic median
std, var	Unbiased (n – 1 de
min, max	Minimum and max
prod	Product of non-NA
first, last	First and last non-l

- A values in the group
- ues
- alues
- of non-NA values
- enominator) standard deviation and variance
- ximum of non-NA values
- values
- NA values

[W. McKinney, Python for Data Analysis]









Iterating over groups

- for name, group in df.groupby('key1'): print (name) print (group)
- Can also .describe() groups







Apply: Generalized methods

<pre>In [74]: def top(df, n=5, column='tip_pct'):</pre>								
• • •	.: ret	urn di.	SOIL_V	alues	(by=cotu	יווין בוויו	:]	
In [75]: top(tip	s, n= <mark>6</mark>)						
Out[75]:							
t	otal_bill	tip s	moker	day	time	size	tip_pc	t
109	14.31	4.00	Yes	Sat	Dinner	2	0.27952	5
183	23.17	6.50	Yes	Sun	Dinner	4	0.28053	5
232	11.61	3.39	No	Sat	Dinner	2	0.29199	Θ
67	3.07	1.00	Yes	Sat	Dinner	1	0.32573	3
178	9.60	4.00	Yes	Sun	Dinner	2	0.41666	7
172	7.25	5.15	Yes	Sun	Dinner	2	0.71034	5
	l. tips or		smokar	חב (''				
Out[76]: ctps.gr	σαρυγ	SHUKET).ah	ριγ(ιορ)			
	tota	l_bill	tip	smoke	r day	time	e size	tip_pct
smoker								
No	88	24.71	5.85	N	o Thur	Lunch	ר 2	0.236746
	185	20.69	5.00	N	o Sun	Dinner	- 5	0.241663
	51	10.29	2.60	N	o Sun	Dinner	- 2	0.252672
	149	7.51	2.00	N	o Thur	Lunch	ר 2	0.266312
	232	11.61	3.39	N	o Sat	Dinner	- 2	0.291990
Yes	109	14.31	4.00	Ye	s Sat	Dinner	- 2	0.279525
	183	23.17	6.50	Ye	s Sun	Dinner	- 4	0.280535
	67	3.07	1.00	Ye	s Sat	Dinner	- 1	0.325733
	178	9.60	4.00	Ye	s Sun	Dinner	- 2	0.416667
	172	7.25	5.15	Ye	s Sun	Dinner	- 2	0.710345

size	tip_pct
2	0.279525
4	0.280535
2	0.291990
1	0.325733
2	0.416667
2	0 710345





- tips.groupby('smoker').apply(top)
- Function is an **argument**
- Function applied on each row group
- All row groups glued together using concat







Types of GroupBy

- Aggregation: agg
 - n:1 n group values become one value
 - Examples: mean, min, median
- Apply: apply
 - n:m n group values become m values
 - Most general (could do aggregation or transform with apply)
 - Example: top 5 in each group, filter
- Transform: transform
 - n:n n group values become n values
 - Cannot mutate the input







Transform Example

In	[76]:	df
Out	:[<mark>76</mark>]:	
	key	value
0	а	0.0
1	b	1.0
2	С	2.0
3	а	3.0
4	b	4.0
5	С	5.0
6	а	6.0
7	b	7.0
8	С	8.0
9	а	9.0
10	b	10.0
11	С	11.0

D. Koop, CSCI 680/490, Spring 2022

```
In [77]: g = df.groupby('key').value
In [78]: g.mean()
Out[78]:
key
     4.5
а
     5.5
b
     6.5
С
Name: value, dtype: float64
In [79]: g.transform(lambda x: x.mean())
Out[79]:
      4.5
0
      5.5
1
     6.5
2
     4.5
3
      5.5
4
     6.5
5
      4.5
6
      5.5
7
     6.5
8
      4.5
9
      5.5
10
11 6.5
Name: value, dtype: float64
```

[W. McKinney, Python for Data Analysis]









Transform Example

In	[76]:	df
Out	:[<mark>76</mark>]:	
	key	value
0	а	0.0
1	b	1.0
2	С	2.0
3	а	3.0
4	b	4.0
5	С	5.0
6	а	6.0
7	b	7.0
8	С	8.0
9	а	9.0
10	b	10.0
11	С	11.0

D. Koop, CSCI 680/490, Spring 2022

```
In [77]: g = df.groupby('key').value
In [78]: g.mean()
Out[78]:
key
    4.5
а
    5.5
b
    6.5
С
Name: value, dtype: float64
In [79]: g.transform(lambda x: x.mean())
Out[79]:
     4.5
0
              Of g.transform('mean')
     5.5
1
     6.5
2
     4.5
3
     5.5
4
     6.5
5
     4.5
6
     5.5
7
     6.5
8
     4.5
9
10
     5.5
11 6.5
Name: value, dtype: float64
```

[W. McKinney, Python for Data Analysis]









Normalization

```
def normalize(x):
    return (x - x.mean()) / x.std()
In [84]: g.transform(normalize)
Out[84]:
     -1.161895
0
     -1.161895
1
2
     -1.161895
3
     -0.387298
     -0.387298
4
     -0.387298
5
      0.387298
6
7
      0.387298
      0.387298
8
      1.161895
9
     1.161895
10
11
      1.161895
Name: value, dtype: float64
```

D. Koop, CSCI 680/490, Spring 2022

```
In [85]: g.apply(normalize)
Out[85]:
     -1.161895
0
     -1.161895
1
     -1.161895
2
     -0.387298
3
     -0.387298
4
5
     -0.387298
      0.387298
6
      0.387298
7
      0.387298
8
      1.161895
9
      1.161895
10
11
      1.161895
Name: value, dtype: float64
```





44

Normalization

```
def normalize(x):
    return (x - x.mean()) / x.std()
In [84]: g.transform(normalize)
Out[84]:
     -1.161895
0
     -1.161895
1
     -1.161895
2
3
     -0.387298
     -0.387298
4
     -0.387298
5
      0.387298
6
      0.387298
      0.387298
8
      1.161895
9
10
      1.161895
11
      1.161895
Name: value, dtype: float64
```

In [87]: normalized = (df['value'] - g.transform('mean')) / g.transform('std')

Fastest: "Unwrapped" group operation

D. Koop, CSCI 680/490, Spring 2022

```
In [85]: g.apply(normalize)
Out[85]:
     -1.161895
0
     -1.161895
     -1.161895
2
     -0.387298
3
     -0.387298
4
5
     -0.387298
      0.387298
6
      0.387298
      0.387298
8
      1.161895
9
      1.161895
10
11
      1.161895
Name: value, dtype: float64
```







44

Other Operations

- Quantiles: return values at particular splits
 - Median is a 0.5-quantile
 - df.quantile(0.1)
 - also works on groups
- Can return data from group-by without having the keys in the index (as index=False) or use reset index after computing
- Grouped weighted average via apply





Pivot Tables

- Data summarization tool in many spreadsheet programs
- Aggregates a table of data by one or more keys with some keys arranged on rows (index), others as columns (columns)
- Pandas supports via pivot table method
- margins=True gives partial totals
- Can use different aggregation functions via aggfunc kwarg

Function name	Description
values	Column name or names to aggregate. By default aggregates
rows	Column names or other group keys to group on the rows of th
cols	Column names or other group keys to group on the columns of
aggfunc	Aggregation function or list of functions; 'mean' by default.
fill_value	Replace missing values in result table
margins	Add row/column subtotals and grand total, False by default

D. Koop, CSCI 680/490, Spring 2022

- all numeric columns
- he resulting pivot table
- of the resulting pivot table
- Can be any function valid in a groupby context

[W. McKinney, Python for Data Analysis]











can also unstack this series into a dataframe
 tip

I	I		I				1	
	total_bill	tip	sex	smoker	day	time	size	tip_pct
0	16.99	1.01	Female	No	Sun	Dinner	2	0.059447
1	10.34	1.66	Male	No	Sun	Dinner	3	0.160542
2	21.01	3.50	Male	No	Sun	Dinner	3	0.166587
3	23.68	3.31	Male	No	Sun	Dinner	2	0.139780
4	24.59	3.61	Female	No	Sun	Dinner	4	0.146808
5	25.29	4.71	Male	No	Sun	Dinner	4	0.186240
6	8.77	2.00	Male	No	Sun	Dinner	2	0.228050

• tips_pivot_table(index=['sex', 'smoker'])

		size	tip	tip_pct	total_bill
sex	smoker				
Female	No	2.592593	2.773519	0.156921	18.105185
	Yes	2.242424	2.931515	0.182150	17.977879
Male	No	2.711340	3.113402	0.160669	19.791237
	Yes	2.500000	3.051167	0.152771	22.284500







Pivot Tables with Margins and Aggfunc

• tips.pivot table(['size'], index=['sex', 'day'], columns='smoker', aggfup.Gevisum([/simargieR-S-Fruge)], columns='smoker', aggfunc



D. Koop, CSCI 680/490, Spring 2022

	size		
oker	Νο	Yes	All
	2.0	7.0	9.0
	13.0	15.0	28.0
	14.0	4.0	18.0
r	25.0	7.0	32.0
	2.0	8.0	10.0
	32.0	27.0	59.0
	43.0	15.0	58.0
r	20.0	10.0	30.0
	151.0	93.0	244.0

tips.pivot table('size', index=['time', 'sex',





smoker]

NIU

Crosstabs

• crosstab is a special case for group frequencies (aggfunc='count')

<pre>In [293]: pd.cross</pre>	tab(data.Ge
Out[293]:	
Handedness Left-h	anded Righ
Gender	
Female	1
Male	2
All	3

- Tipping example
- Also see the Federal Election Database example in the book

D. Koop, CSCI 680/490, Spring 2022

- ender, data.Handedness, margins=True)
- it-handed All
 - 4 53 5 10





49

Crosstabs

margins=True)



fill value=0) smo

D. Koop, CSCI 680/490, Spring 2022

day time

• pd.crosstab([tips.timed.crosspace([diges.time.trips.sdaysmolice.rmpker, margins=True)

oker	No	Yes	AII
7			
	3	9	12
	45	42	87
ו	57	19	76
ır	1	0	1
	1	6	7
ır	44	17	61
	151	93	244

• Or... tips.pivot_table('#ton minic crosstab using axpivot_table ', 'day'], # doesn't-matter what the data (first argument) is Columns=['smoker'], aggps1p1x6E_ta61eUHota1_biMa,1641axAStime Ueday'], columns=['smoker

oker	No	Yes	All	
7				







Time Series Data







What is time series data?

- Technically, it's normal tabular data with a timestamp attached
- This allows more analysis
- Example: Web site database that tracks the last time a user logged in

 - 2: Add a new row with login information every time the user logs in
 - Option 2 takes more storage, but we can also do a lot more analysis!

But... we have observations of the same values over time, usually in order

- 1: Keep an attribute lastLogin that is overwritten every time user logs in







Time Series Data

- Metrics: measurements at regular intervals
- Events: measurements that are not gathered at regular intervals













Types of Time Series Data

- time series: observations for a **single** entity at **different** time intervals - one patient's heart rate every minute
- cross-section: observations for **multiple** entities at the **same** point in time
- heart rates of 100 patients at 8:01pm
- panel data: observations for **multiple** entities at **different** time intervals - heart rates of 100 patients every minute over the past hour











Time Series Databases

- Most time series data is heavy **inserts**, few updates
- Also analysis tends to be on ordered data with trends, prediction, etc.
- Can also consider stream processing
- Focus on time series allows databases to specialize
- Examples:
 - InfluxDB (noSQL)
 - TimescaleDB (SQL-based)









Features of Time Series Data

- Trend: long-term increase or decrease in the data
- Seasonal Pattern: time series is affected by seasonal factors such as the time of the year or the day of the week (fixed and of known frequency)
- Cyclic Pattern: rises and falls that are not of a fixed frequency
- Stationary: no predictable patterns (roughly horizontal with constant variance)
 - White noise series is stationary
 - Will look the basically the same whenever you observe it



































































Types of Time Data

- Timestamps: specific instants in time (e.g. 2018 11 27 14:15:00)
- Periods: have a standard start and length (e.g. the month November 2018)
- Intervals: have a start and end timestamp
 - Periods are special case
 - Example: 2018-11-21 14:15:00 2018-12-01 05:15:00
- Elapsed time: measure of time relative to a start time (15 minutes)









Dates and Times

- What is time to a computer?
 - Can be stored as seconds since Unix Epoch (January 1st, 1970)
- Often useful to break down into minutes, hours, days, months, years...
- Lots of different ways to write time:
 - How could you write "November 29, 2016"?
 - European vs. American ordering...
- What about time zones?









Python Support for Time

- The datetime package
 - Has date, time, and datetime classes
 - .now() method: the current datetime
- Can access properties of the time (year, month, seconds, etc.) Converting from strings to datetimes:
 - datetime.strptime: good for known formats
 - dateutil.parser.parse: good for unknown formats
- Converting to strings
 - str(dt) Or dt.strftime(<format>)







Datetime format specification

- Look it up:
 - <u>http://strftime.org</u>
- Generally, can create whatever format you need using these format strings

Code	Meaning	Example
%a	Weekday as locale's abbreviated name.	Mon
۶A	Weekday as locale's full name.	Monday
88	Weekday as a decimal number, where 0 is Sunday and 6 is Saturday.	1
%d	Day of the month as a zero-padded decimal number.	30
%−d	Day of the month as a decimal number. (Platform specific)	30
%b	Month as locale's abbreviated name.	Sep
۶B	Month as locale's full name.	September
%m	Month as a zero-padded decimal number.	09
%-m	Month as a decimal number. (Platform specific)	9
۶y	Year without century as a zero-padded decimal number.	13
8Y	Year with century as a decimal number.	2013
%H	Hour (24-hour clock) as a zero-padded decimal number.	07
%−H	Hour (24-hour clock) as a decimal number. (Platform specific)	7
%I	Hour (12-hour clock) as a zero-padded decimal number.	07
%-I	Hour (12-hour clock) as a decimal number. (Platform specific)	7
%p	Locale's equivalent of either AM or PM.	AM
۶M	Minute as a zero-padded decimal number.	06
%-M	Minute as a decimal number. (Platform specific)	6
°S	Second as a zero-padded decimal number.	05
%-S	Second as a decimal number. (Platform specific)	5









Pandas Support for Datetime

- pd.to datetime:
 - convenience method
 - can convert an entire column to datetime
- Has a Nat to indicate a missing time value
- Stores in a numpy.datetime64 format
- pd.Timestamp: a wrapper for the datetime 64 objects







More Pandas Support

- can be interpreted as a date:
 - ts['1/10/2011'] Or ts['20110110']
- Date ranges: pd.date range('4/1/2012', '6/1/2012', freq='4h')
- Slicing works as expected
- Can do operations (add, subtract) on data indexed by datetime and the indexes will match up
- As with strings, to treat a column as datetime, you can use the .dt accessor

D. Koop, CSCI 680/490, Spring 2022

Accessing a particular time or checking equivalence allows any string that









Generating Date Ranges

- index = pd.date range('4/1/2012', '6/1/2012')
- Can generate based on a number of periods as well - index = pd.date range('4/1/2012', periods=20)
- Frequency (freq) controls how the range is divided
 - Codes for specifying this (e.g. 4h, D, M)
 - In [90]: pd.date range('1/1/2000', '1/3/2000 23:59', freq='4h') Out[90]: <class 'pandas.tseries.index.DatetimeIndex'> $[2000-01-01 \ 00:00:00, \ldots, 2000-01-03 \ 20:00:00]$ Length: 18, Freq: 4H, Timezone: None
 - Can also mix them: '2h30m'











Time Series Frequencies

Alias	Offset Type
D	Day
В	BusinessDay
Н	Hour
T or min	Minute
S	Second
L or ms	Milli
U	Micro
Μ	MonthEnd
BM	BusinessMonthEnd
MS	MonthBegin
BMS	BusinessMonthBegin
W-MON, W-TUE,	Week
WOM-1MON, WOM-2MON,	WeekOfMonth

	Description	
	Calendar daily	
	Business daily	
	Hourly	
	Minutely	
	Secondly	
	Millisecond (1/1000th of 1 second)	
	Microsecond (1/1000000th of 1 second)
	Last calendar day of month	
	Last business day (weekday) of month	
	First calendar day of month	
1	First weekday of month	
	Weekly on given day of week: MON, TU or SUN.	IE, WED, THU, FRI, SAT,
	Generate weekly dates in the first, secor of the month. For example, WOM-3FR	nd, third, or fourth week E for the 3rd Friday of
	each month.	[W. McKinney, F








DatetimeIndex

- Can use time as an **index**
- data = [('2017 11 30', 48)]('2017 - 12 - 02', 45),(2017 - 12 - 03', 44),(2017 - 12 - 04', 48)dates, temps = zip(*data)

s = pd.Series(temps, pd.to datetime(dates))

- Accessing a particular time or checking equivalence allows any string that can be interpreted as a date:
 - s['12/04/2017'] Or s['20171204']
- Using a less specific string will get all matching data:
 - s['2017-12'] returns the three December entries







DatetimeIndex

• Time slices do not need to exist: - s['2017-12-01':'2017-12-31']

D. Koop, CSCI 680/490, Spring 2022









Shifting Data

Leading or Lagging Data

```
In [95]: ts = Series(np.random.randn(4),
                    index=pd.date_range('1/1/2000', periods=4, freq='M'))
   • • • • •
In [96]: ts
                            In [97]: ts.shift(2)
                                                         In [98]: ts.shift(-2)
Out[96]:
                            Out[97]:
                                                         Out[98]:
2000-01-31
            -0.066748
                                                         2000-01-31
                            2000-01-31
                                               NaN
                                                                      -0.117388
            0.838639
                                               NaN
2000-02-29
                            2000-02-29
                                                         2000-02-29
                                                                      -0.517795
2000-03-31 -0.117388
                            2000-03-31 -0.066748
                                                         2000-03-31
                                                                            NaN
                                                                            NaN
2000-04-30 -0.517795
                            2000-04-30
                                          0.838639
                                                         2000-04-30
Freq: M, dtype: float64
                            Freq: M, dtype: float64
                                                         Freq: M, dtype: float64
```

• Shifting by time:

```
In [99]: ts.shift(2, freq='M')
Out[99]:
            -0.066748
2000-03-31
2000-04-30
           0.838639
            -0.117388
2000-05-31
2000-06-30 -0.517795
Freq: M, dtype: float64
```







Shifting Time Series

• Data:

[(2017-11-30', 48), (2017-12-02', 45),('2017 - 12 - 03', 44), ('2017 - 12 - 04', 48)]

• Compute day-to-day difference in high temperature:

- s s.shift(1) (same as s.di
- 2017-11-30 NaN 2017 - 12 - 02 - 3.02017 - 12 - 03 - 1.04.0 2017 - 12 - 04

D. Koop, CSCI 680/490, Spring 2022





Timedelta

- Compute differences between dates
- Lives in datetime module
- diff = parse date("1 Jan 2017") datetime.now().date() diff.days
- Also a pd. Timedelta object that take strings:
 - datetime.now().date() + pd.Timedelta("4 days")
- Also, Roll dates using anchored offsets from pandas.tseries.offsets import Day, MonthEnd

now = datetime(2011, 11, 17)In [107]: now + MonthEnd(2) Out[107]: Timestamp('2011-12-31 00:00:00')





70