

Extracting Aggregate Answer Statistics for Integration

Zainab Zolaktaf ¹ Jian Xu ² Rachel Pottinger ¹

¹Department of Computer Science
University of British Columbia

²Microsoft Corporation

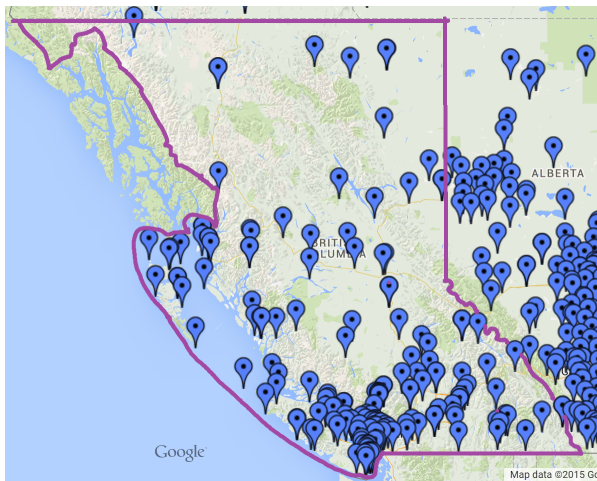


Average high temperature across British Columbia

- Climate change?
 - What is the average high temperature in British Columbia for each year?
 - Averaging across the temperature over the entire province seems reasonable?
 - Query the weather stations to find the average



Weather stations not distributed uniformly across BC



Sources have inconsistent values for same data points

- Weather of Vancouver on 11-June-2006?



Location	Avg Temp	Date
Burnaby	21	10-June-06
Vancouver	19	11-June-06
Surrey	18	11-June-06
...

Sources have inconsistent values for same data points

- Weather of Vancouver on 11-June-2006?



Location	Avg Temp	Date
Burnaby	21	10-June-06
Vancouver	19	11-June-06
Surrey	18	11-June-06
...



City	Temp	Date
Burnaby	21	06/10/06
Vancouver	22	06/11/06
Richmond	18	06/12/06
Richmond	18	06/13/06
...

Sources have different coverage and quality

- Coverage: single source contains information about a subset of objects and a subset of object attributes
- Quality: inconsistent or even conflicting values for the same object

The diagram illustrates four data sources, each represented by a colored cylinder, pointing to a table of data. The sources are: 1. An orange cylinder pointing to a table with columns Location, Avg Temp, and Date. 2. A purple cylinder pointing to a table with columns City, Temp, and Date. 3. A green cylinder pointing to a table with columns City, Temp, Date, and Total Rain. 4. A blue cylinder pointing to a table with columns Location, Temp, Date, Total Snow, and Total Rain.

Location	Avg Temp	Date
Burnaby	21	10-June-06
Vancouver	19	11-June-06
...

City	Temp	Date
Burnaby	21	06/10/06
Vancouver	22	06/11/06
Richmond	18	06/12/06
Richmond	18	06/13/06
...

City	Temp	Date	...	Total Rain
Burnaby	19	10-June-06	...	0.2
Vancouver	17	11-June-06	...	0.0
Surrey	15	11-June-06	...	0.0
Vancouver	20	12-June-06	...	1.4
...

Location	Temp	Date	Total Snow	Total Rain
Surrey	15	06/11/06	0.0	0.0
Surrey	19	06/12/06	0.0	1.2
...

Aggregate queries

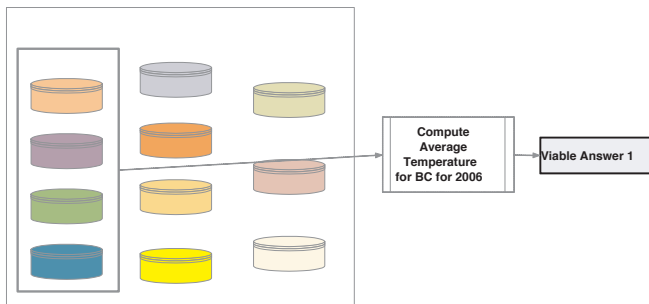
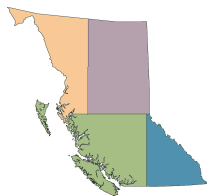
- Group set of data values and calculate informative statistics
 - Sum, Median, Avg ...
- Answering them in integration contexts
 - Requires combining sets of data that are segmented across multiple data sources
- Standard aggregation averages over all the points
 - It is incorrect!
 - Some data points have duplicates across the sources
 - The duplicates can have different values in the sources

Viable answer

- Correct aggregation requires using one value per data point
- Choosing the values from different sources will result in different answers
- Each possible answer called a viable answer
- Which set of sources and value combinations to use?

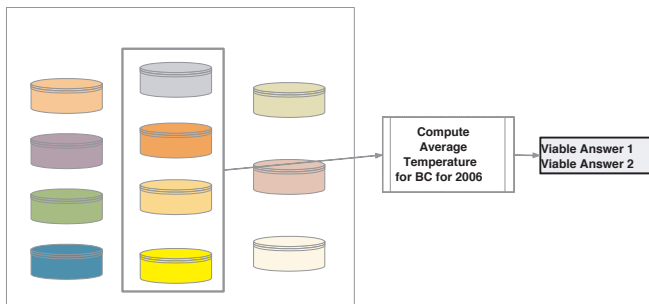
Different combinations of sources and values are possible

- Which set of sources and value combinations to use?



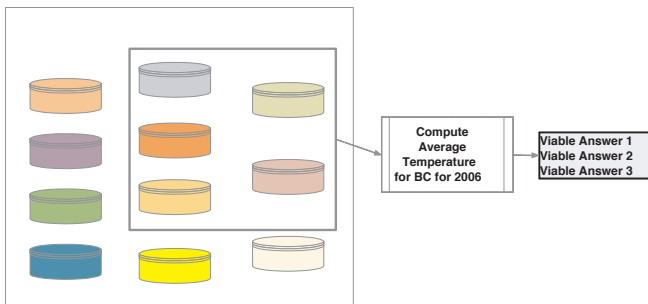
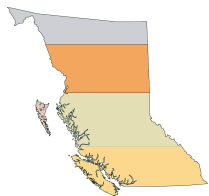
Different combinations of sources and values are possible

- Which set of sources and value combinations to use?



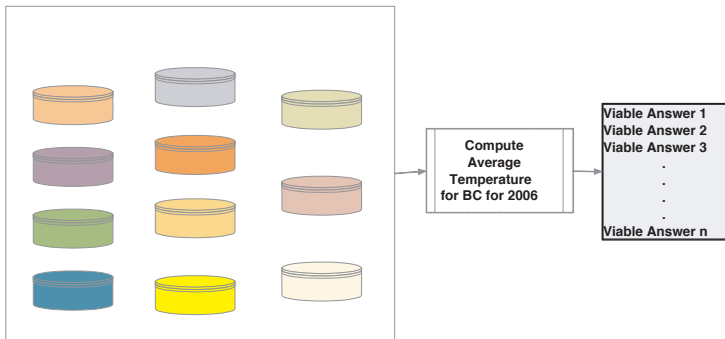
Different combinations of sources and values are possible

- Which set of sources and value combinations to use?



Different combinations of sources and values are possible

- Depending on the choice of sources and value combinations, there can be a whole range of viable answers
- Aggregate query answer is a distribution rather than a single scalar value



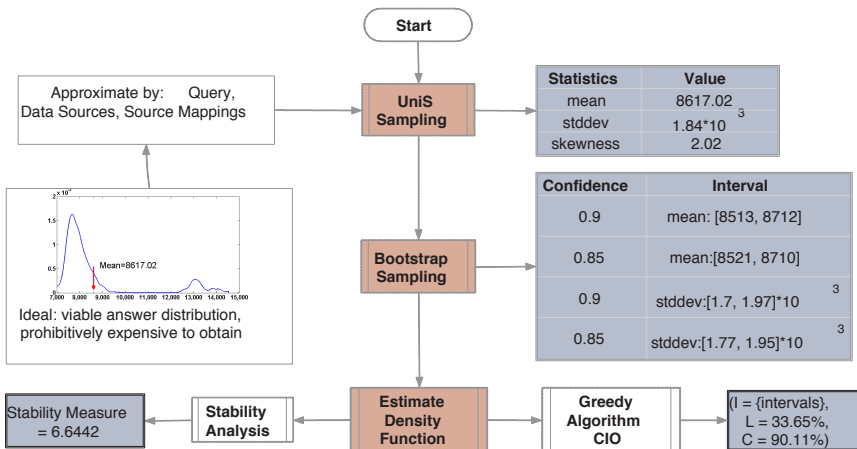
Problem formalization

- Assume meta-information regarding mapping and bindings is available
- Question: What is the viable answer distribution?
 - Enumerating all the possible value combinations is impractical
 - Estimating the exact distribution is infeasible
 - Scalability issues
 - User still has to interpret and analyze

Contributions

- We define aggregate answers as a distribution of viable answers
- We provide summary statistics for the viable answer distribution
 - Key point statistics
 - High coverage intervals
 - Stability score
- We provide algorithms for the efficient extraction of above statistics
- We verify the effectiveness of our methods using real-life and synthetic data

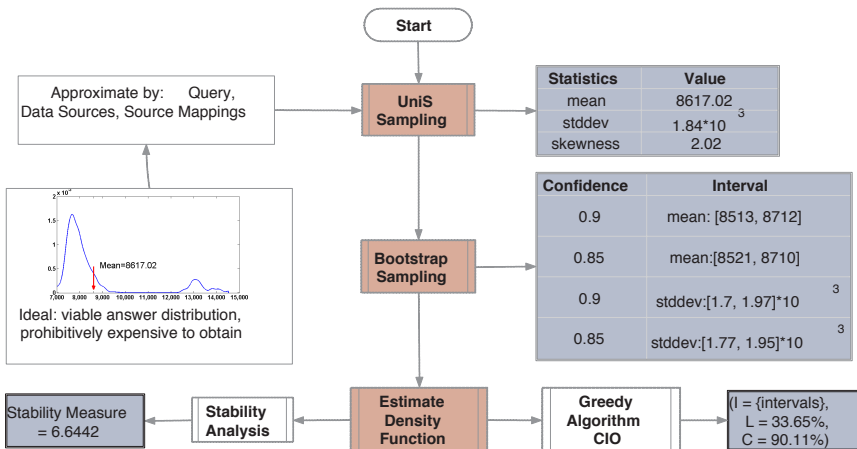
Overview



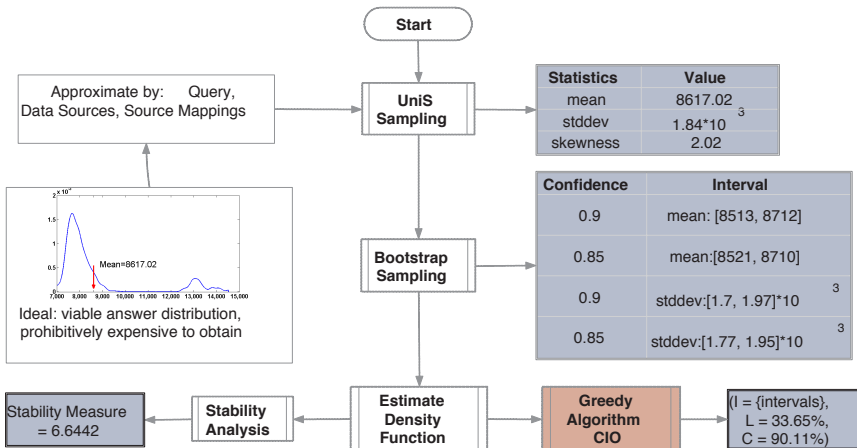
Sampling and point statistics

- Goal: Efficiently approximate viable answer distribution
- Sample a set of viable answers
 - No prior knowledge regarding coverage, accuracy and quality
- Sampling scheme? Uniform sampling
 - Choose sources uniformly at random
 - Stay at source until source is exhausted (all relevant components used)
- Apply bootstrap sampling and bagging
- Apply kernel density estimation

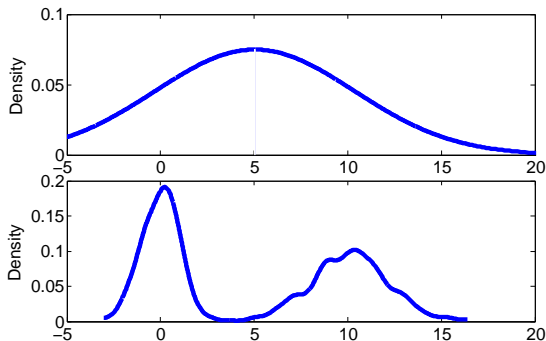
Overview



Overview



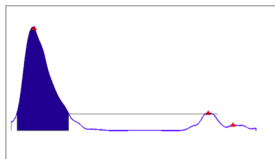
High coverage intervals and optimization



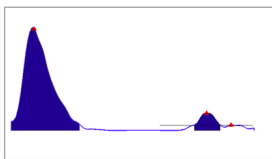
- Point statistics such as mean and variance are insufficient
- Statistics that convey shape information are needed!

High coverage intervals and optimization

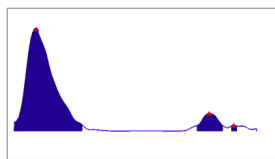
- Goal: Communicate shape information about viable answer distribution
- Greedy algorithm CIO
 - Minimize interval length so that coverage of viable answers is above a certain threshold



(a) initial high coverage interval finding

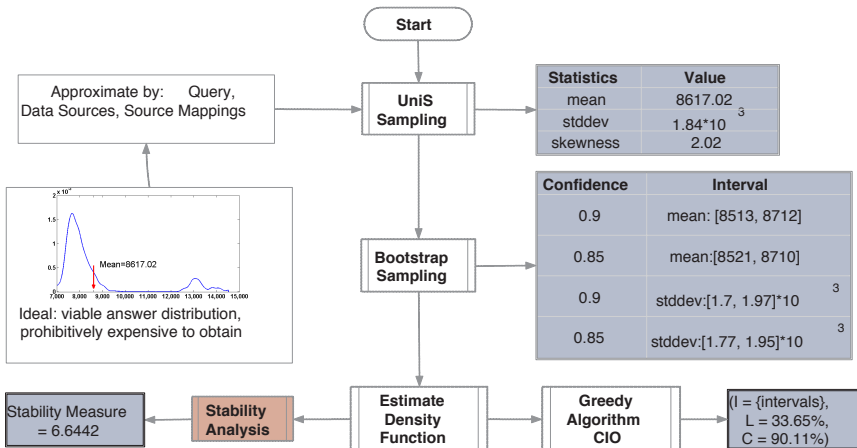


(b) intermediate



(c) high coverage intervals that are above our threshold ($\theta\%$)

Overview



Stability score

- Goal: How much change is caused in the estimated distribution when some sources leave?
- Quantify the change as the distance between two probability distributions
 - The original viable answer distribution
 - The viable answer distribution when some of the sources are removed
- Obtain stability analytically for the L_2 measure
- Helps prioritize re-evaluation and updating of queries need updating when sources are updated

Empirical study

- Dataset
 - Synthetic data
 - D2 - Mixtures of four Gaussians
 - D3 - Mixture of Gaussians, Cauchy and Gamma
 - Real-life data
 - C - Monthly Canadian climate data for the year 2006, from 1672 stations for 104 districts
- Aggregate query: Sum temperature data over 500 components from datasets

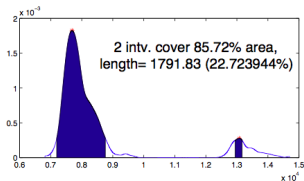
Bootstrapping vs Direct inference

- Bootstrapping helps derive tighter confidence intervals for point statistics
 - Smaller confidence intervals represent more reliable estimates
 - Improvement ratio $i_r = \frac{\text{len}(CI_{di})}{\text{len}(CI_{boot})}$

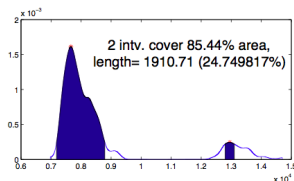
$ \mathcal{S}_{uniS} $	$1 - \alpha$	$\max i_r$	$\text{avg } i_r$
200	0.8	4.248	2.556
200	0.9	3.309	2.119
400	0.8	2.896	2.001
400	0.9	2.293	1.655

Greedy algorithm output

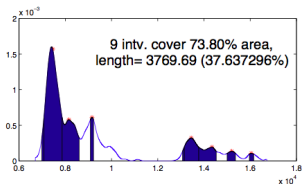
- By returning dense areas, the intervals cover a small percentage of the range of data



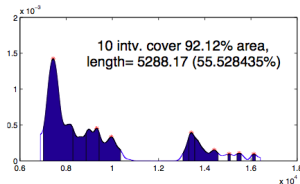
(a) S_1 (Climate Data C)



(b) S_2 (Climate Data C)



(c) S_3 (Synthetic data D3)



(d) S_4 (Synthetic data D3)

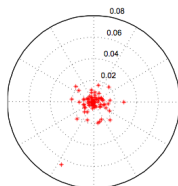
High coverage intervals

- Optimal Slicing: slice the area under the distribution into 4028 slices
- Optimal Slicing vs Greedy Algorithm CIO
 - Optimal Slicing returns tighter intervals, but does not guarantee the continuity of returned intervals

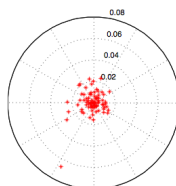
Fig	Greedy	Optimal	Cover	Greedy/Optimal
a	0.2272	0.2272	85.72%	1.0
b	0.2475	0.2475	85.44%	1.0
c	0.3764	0.2724	73.82%	1.38
d	0.5552	0.5150	92.12%	1.08

Stability score

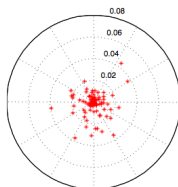
- Closer to the center, and the more dense around the center, the more stable the result



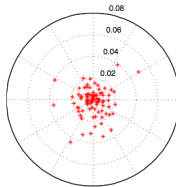
(a) S_1 L_2 score=6.5882



(b) S_2 L_2 score=6.4139



(c) S_3 L_2 score=6.4217



(d) S_4 L_2 score=6.3204

Selected related work

- Data Fusion - assumes a single true answer
 - Bleiholder, Jens, and Felix Naumann. "Data fusion." ACM Computing Surveys (CSUR) 41.1 (2008): 1.
 - Dong, Xin Luna, and Felix Naumann. "Data fusion: resolving data conflicts for integration." Proceedings of the VLDB Endowment 2.2 (2009): 1654-1655.
- Value-level heterogeneity in the Flight and Stock domain, due to sources applying different semantics
- No single true answer
 - Li, Xian, et al. "Truth finding on the deep web: is the problem solved?." Proceedings of the VLDB Endowment 6.2 (2012): 97-108.
- Applications in probabilistic databases

Future work

- Investigate the stability analysis
- Improve uniS sampling to consider quality and coverage
- Make inferences regarding data and sources based on non-normality of estimated viable distribution
 - Multi-modal distributions can indicate mapping problems
 - Find homogeneous sources that apply similar semantics

Contributions

- We define aggregate answers as a distribution of viable answers
- We provide summary statistics for the viable answer distribution
 - Key point statistics
 - High coverage intervals
 - Stability score
- We provide algorithms for the efficient extraction of above statistics
- We verify the effectiveness of our methods using real-life and synthetic data

Dataset

Home > About the Data > Resources

Monthly Data Report for 2006

Station Name

Metadata including Station Name, Province, Latitude, Longitude, Elevation, Climate ID, 1950 ID, 11 ID

1962 LAMLEY VILLAGE/ST. JOHN'S WESTERN COLUMBIA

Latitude	49°04'41.800" N	Longitude	122°51'03.600" W	Elevation	76.00 m
Climate ID	1108913	1950 ID		11 ID	

Home > About the Data > Resources

Monthly Data Report for 2006

Station Name

Metadata including Station Name, Province, Latitude, Longitude, Elevation, Climate ID, 1950 ID, 11 ID

WHITTE ROCK CAMP/ST. SCIENTIFIC WESTERN COLUMBIA

Latitude	49°04'05.000" N	Longitude	122°47'03.000" W	Elevation	13.00 m
Climate ID	1108910	1950 ID	71780	11 ID	WVK

Home > About the Data > Resources

Monthly Data Report for 2006

Station Name

Metadata including Station Name, Province, Latitude, Longitude, Elevation, Climate ID, 1950 ID, 11 ID

1111 HANCOCK CO WESTERN COLUMBIA

Latitude	49°17'20.040" N	Longitude	122°41'24.870" W	Elevation	5.00 m
Climate ID	1108179	1950 ID	71778	11 ID	WVK

Monthly Data Report for 2006

Month	Mean Max Temp	Mean Min Temp	Mean Snow Depth	Mean Rain	Total Rain	Total Snow	Max on Ground	Max on Last Day	Max on First Day	Max on Last Day	Max on First Day
Jan	8.1	4.6	6.0	13.0	1.02	481.0	0.0	481.0	0		
Feb	8.2	0.8	4.4	13.0	-0.5	79.0	2.4	81.4	0		
Mar	10.78	2.8	4.7	16.9	-0.89	90.8	10.6	101.4	0		
Apr	13.40	4.40	9.30	17.40	4.00	16.40	0.0	16.40	0		
May											
Jun	19.30	12.40	19.10	26.40	9.40	4.10	0.0	4.10	0		
Jul	22.40	13.40	17.40	23.40	11.40	0.40	0.0	6.40	15		
Aug	21.8	9.9	18.9	21.8	9.0	101.0	0.0	101.0	0		
Sep	15.0	5.9	10.4	21.50	-4.0	92.4	0.0	92.4	0		
Oct	7.4	3.7	9.1	18.0	-0.9	402.0	49.0	451.0	36		
Nov	4.90	1.40	3.90	13.90	-0.90	133.90	0.00	133.90	0		
Dec											
Sum											
Avg	14.4	7.2	10.9								
Wmax											

Summary, average and extreme values are based on the data above.

Monthly Data Report for 2006

Month	Mean Max Temp	Mean Min Temp	Mean Snow Depth	Mean Rain	Total Rain	Total Snow	Max on Ground	Max on Last Day	Max on First Day	Max on Last Day	Max on First Day
Jan	9.2	5.1	7.2	13.4	2.3		206.8				
Feb	9.1	1.9	5.2	14.5	-4.5		94.4				
Mar	11.4	2.8	7.6	16.9	-0.2		96.2				
Apr	13.8	6.2	10.0	21.3	2.3		76.8				
May	17.6	8.6	13.2	26.8	3.9		44.4				
Jun	20.1	12.2	16.4	27.9	9.1		23.9				
Jul	23.2	12.7	18.0	26.6	10.3		9.8				
Aug	21.2	12.5	16.9	26.6	10.3		7.8				
Sep	16.9	10.2	14.6	27.9	8.8		69.0				
Oct	14.2	6.4	10.3	23.0	-1.6		236.9				
Nov	8.6	3.7	6.3	15.3	-0.2		178.6				
Dec	7.9	3.6	5.3	13.9	-0.0		129.4				
Sum											
Avg	14.4	7.2	10.9								
Wmax											

Summary, average and extreme values are based on the data above.

Monthly Data Report for 2006

Month	Mean Max Temp	Mean Min Temp	Mean Snow Depth	Mean Rain	Total Rain	Total Snow	Max on Ground	Max on Last Day	Max on First Day	Max on Last Day	Max on First Day
Jan	8.3	3.6	4.2	12.8	-0.2		375.4				
Feb	8.7	-0.1	4.3	15.2	-7.1		63.6				
Mar	11.0	1.9	6.9	18.2	-3.7		112.0				
Apr	14.7	4.9	9.6	23.2	0.2		62.0				
May	19.4	7.4	13.0	29.2	0.9		82.4				
Jun	22.3	10.3	16.8	32.9	7.2		51.90				
Jul	25.1	12.9	19.0	26.7	8.7		22.6				
Aug	24.7	11.3	18.0	25.0	7.7		20.2				
Sep	22.3	9.4	16.9	22.0	6.9		61.2				
Oct	15.8	4.9	10.4	21.7	-4.1		61.9				
Nov	7.9	2.1	5.0	16.3	-13.8		417.4				
Dec	6.6	0.9	3.7	13.3	-4.9		188.6				
Sum											
Avg	13.6	5.9	10.7								
Wmax											

Summary, average and extreme values are based on the data above.

Home > Data Quality

Legend

- = How often an occurrence and estimated
- = Estimated
- = Missing
- = How often an occurrence

* = How often an occurrence and estimated
 * = Estimated
 * = Missing
 * = How often an occurrence

Home > Data Quality

Legend

- = How often an occurrence and estimated
- = Estimated
- = Missing
- = How often an occurrence

* = How often an occurrence and estimated
 * = Estimated
 * = Missing
 * = How often an occurrence

Home > Data Quality

Legend

- = How often an occurrence and estimated
- = Estimated
- = Missing
- = How often an occurrence

* = How often an occurrence and estimated
 * = Estimated
 * = Missing
 * = How often an occurrence

Preliminaries

- Bootstrap sampling and bagging
- Kernel density estimation
- Distance measures for distributions


Finding high coverage intervals - optimization approach

Given a density function f_X^D for a distribution defined on a finite range, a coverage threshold $0 \leq \theta \leq 1$, and a constant t representing the number of modes, the CIO problem finds k intervals I_1, I_2, \dots, I_k where $k \leq t$, to minimize $\sum_{i=1}^k |I_i|$ subject to $\sum_{i=1}^k \int_{I_i} f_X^D(x) dx \geq \theta$.

$$\begin{aligned} & \underset{k, I_1, \dots, I_k}{\text{minimize}} && \sum_{i=1}^k |I_i| \\ & \text{subject to} && \sum_{i=1}^k \int_{I_i} f_X^D(x) dx \geq \theta. \end{aligned}$$


Heterogeneity

- Heterogeneity at three levels
 - Schema-level




→

Location	Avg Temp	Date
Burnaby	21	10-June-06
Vancouver	19	11-June-06
...




→

City	Temp	Date
Burnaby	21	06/10/06
Vancouver	22	06/11/06
Richmond	18	06/12/06
Richmond	18	06/13/06
...



→

City	Temp	Date	...	Total Rain
Burnaby	19	10-June-06	...	0.2
Vancouver	17	11-June-06	...	0.0
Surrey	15	11-June-06	...	0.0
Vancouver	20	12-June-06	...	1.4
...




→

Location	Temp	Date	Total Snow	Total Rain
Surrey	15	06/11/06	0.0	0.0
Surrey	19	06/12/06	0.0	1.2
...


Heterogeneity

- Heterogeneity at three levels


- Schema-level
- Instance-level




Location	Avg Temp	Date
Burnaby	21	10-June-06
Vancouver	19	11-June-06
...



City	Temp	Date
Burnaby	21	06/10/06
Vancouver	22	06/11/06
Richmond	18	06/12/06
Richmond	18	06/13/06
...



City	Temp	Date	...	Total Rain
Burnaby	19	10-June-06	...	0.2
Vancouver	17	11-June-06	...	0.0
Surrey	15	11-June-06	...	0.0
Vancouver	20	12-June-06	...	1.4
...




Location	Temp	Date	Total Snow	Total Rain
Surrey	15	06/11/06	0.0	0.0
Surrey	19	06/12/06	0.0	1.2
...


Heterogeneity

- Heterogeneity at three levels


- Schema-level
- Instance-level
- Value-level




Location	Avg Temp	Date
Burnaby	21	10-June-06
Vancouver	19	11-June-06
...



City	Temp	Date
Burnaby	21	06/10/06
Vancouver	22	06/11/06
Richmond	18	06/12/06
Richmond	18	06/13/06
...




City	Temp	Date	...	Total Rain
Burnaby	19	10-June-06	...	0.2
Vancouver	17	11-June-06	...	0.0
Surrey	15	11-June-06	...	0.0
Vancouver	20	12-June-06	...	1.4
...




Location	Temp	Date	Total Snow	Total Rain
Surrey	15	06/11/06	0.0	0.0
Surrey	19	06/12/06	0.0	1.2
...

Value-level Heterogeneity


- Focus of our work is on value-level heterogeneity
 - Problem exists in various domains, e.g., stock, flight, weather domain
- Prior work assumes a single true answer exists, which we do not




Location	Avg Temp	Date
Burnaby	21	10-June-06
Vancouver	19	11-June-06
...



City	Temp	Date
Burnaby	21	06/10/06
Vancouver	22	06/11/06
Richmond	18	06/12/06
Richmond	18	06/13/06
...



City	Temp	Date	...	Total Rain
Burnaby	19	10-June-06	...	0.2
Vancouver	17	11-June-06	...	0.0
Surrey	15	11-June-06	...	0.0
Vancouver	20	12-June-06	...	1.4
...



Location	Temp	Date	Total Snow	Total Rain
Surrey	15	06/11/06	0.0	0.0
Surrey	19	06/12/06	0.0	1.2
...

Processing overhead of operations

- KDE dominates the processing overhead for extracting statistics
- Sampling the viable answers dominates the overall time needed for sampling and extracting statistics

