

Setting

- large tables with outliers in aggregation columns
- approximate aggregation queries

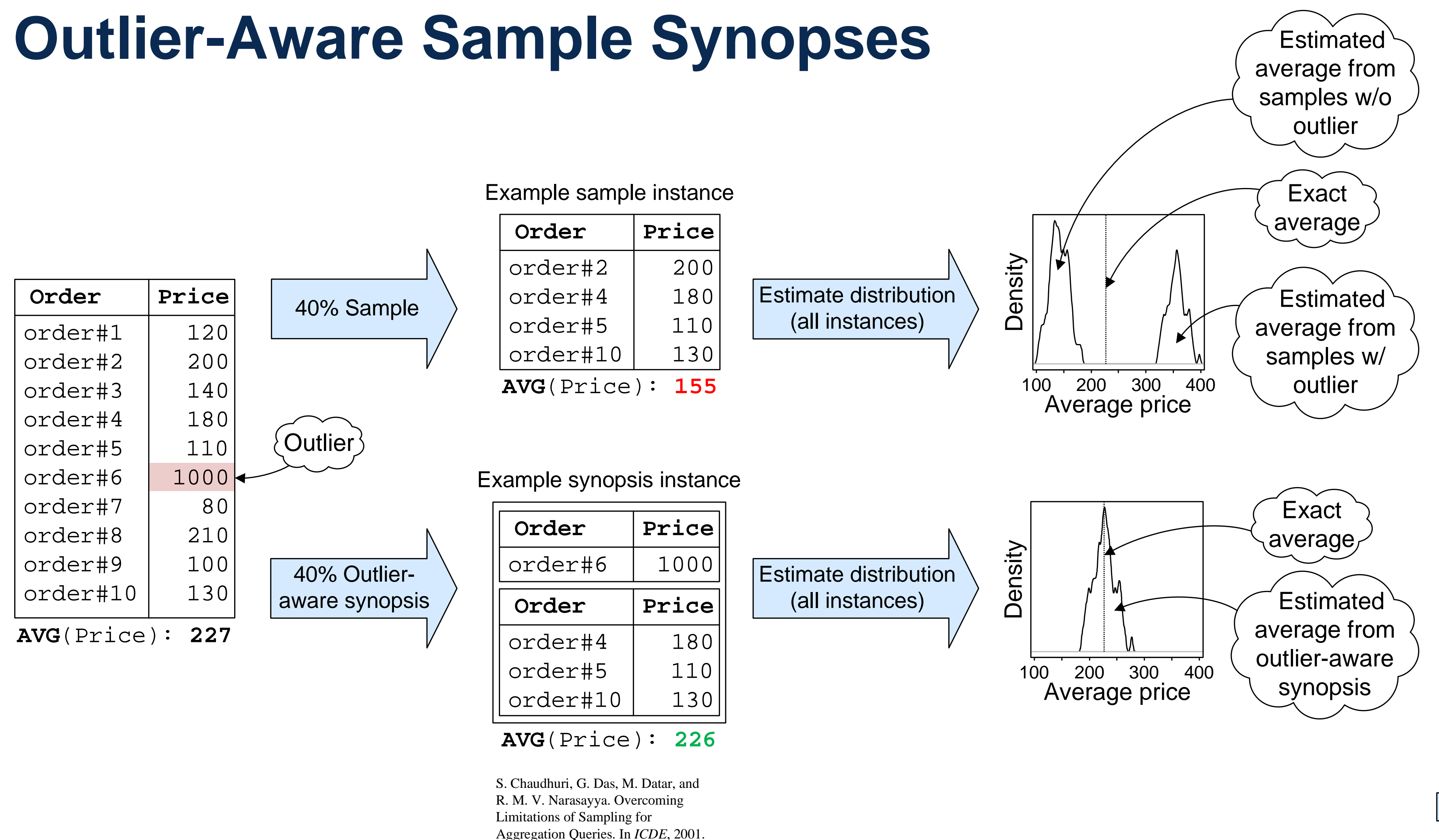
Goal

- memory-bounded outlier-aware sample synopses

Challenge

- efficient computation
- low estimation error for all aggregation columns

A Outlier-Aware Sample Synopses



B Extension to Multiple Columns

Problems

- What is an outlier?
 - Outliers of one column may not be outliers in other columns
 - Multiple single-column instances are not optimal under space constraint

Solution

- Quantification of a synopsis' quality with **measures** based on the relative standard error (RSE) of the estimates

- How to select outliers?
 - Number of possible outliers is prohibitively large

Solution

- Speed up computation by using **heuristics** and **greedy** proceeding

Measures

MAX-Measure

- effect: minimizes maximum RSE of the estimates + most intuitive
- fails if column with maximum RSE has no outliers

GEO-Measure

- effect: maximizes improvement in RSE compared to simple random sampling
- + eager outlier selection
- a lot of effort into columns with very low RSE

AVG-Measure

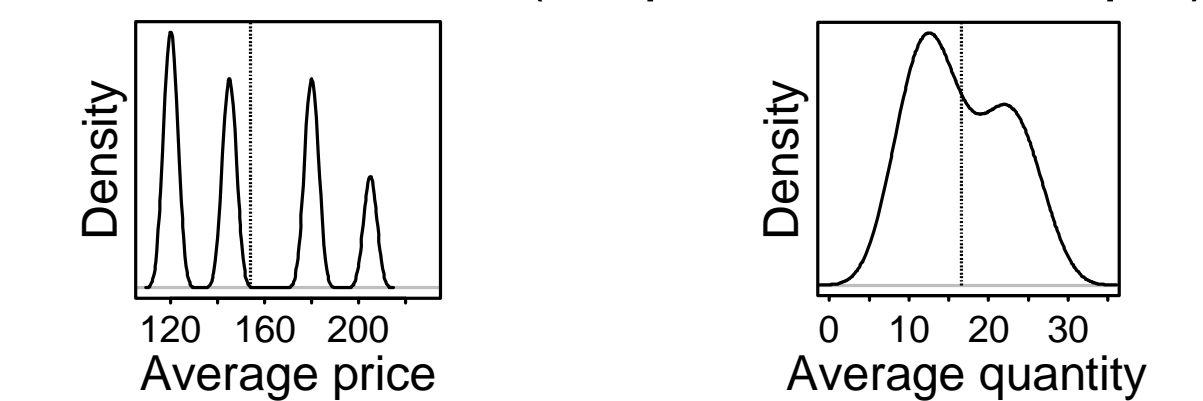
- effect: minimizes the average RSE of the estimates
- + prefers outliers from columns with high RSE
- + shows best overall performance

Order	Price	Qty.
order#1	120	9
order#2	125	4
order#3	115	12
order#4	119	54
order#5	130	1
order#6	110	5
order#7	360	15
order#8	123	21
order#9	118	18
order#10	220	27

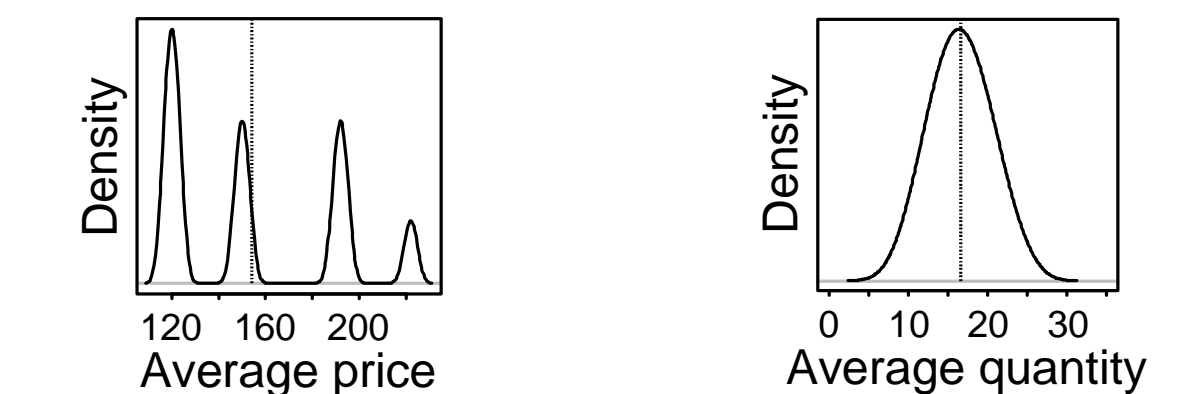
RSD: 48.7% 88.3%

Estimate distribution of synopses of size 40%

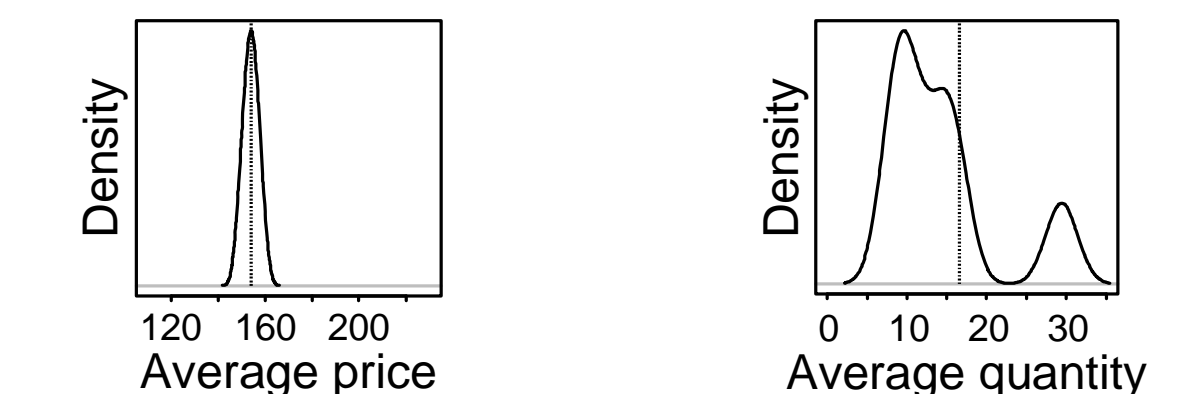
Selected outliers: - (simple random sample)



Selected outliers: [119|54]



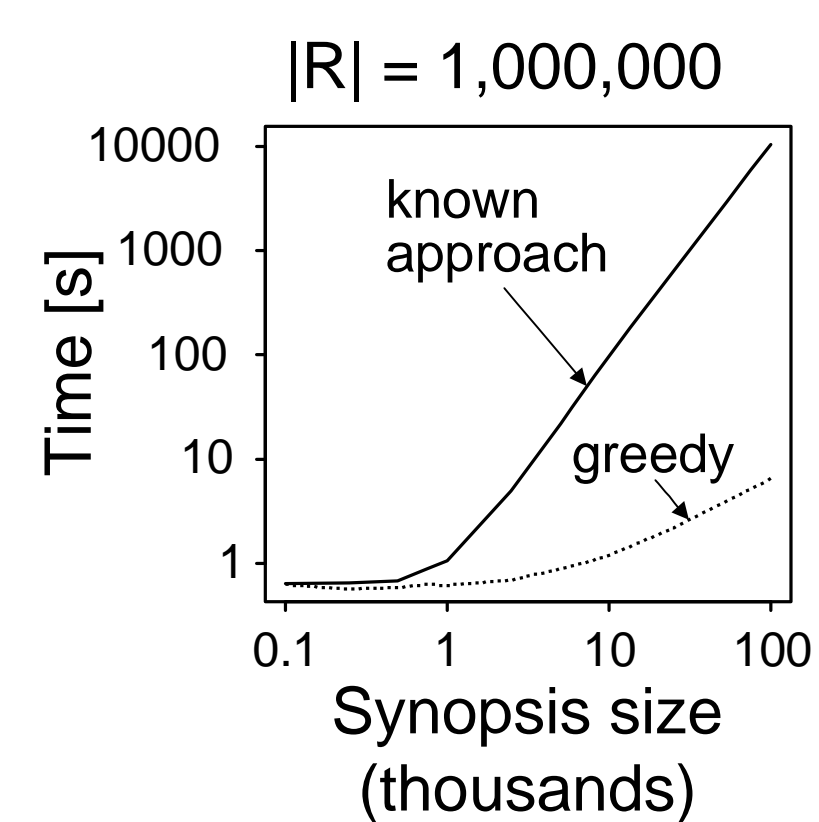
Selected outliers: [360|15], [220|27]



Selected outliers: [119|54], [360|15], [220|27]



C Impact on Single-Column Outlier-Aware Synopses



Computational Contribution

Search Space

known: $2M$ candidates from lower/upper value range with memory bound M

Outlier Selection

Problem: $O(M^2)$ combinations for M candidates

Solution: select outliers greedily from candidate set

Computational Contribution

Search Space

Problem: huge number of possible outlier sets $\sum_{i=0}^M \binom{|R|}{i}$ with memory bound M

Solution: select M outlier candidates during the single scan of the data selection based on weights specific to measure of choice

Outlier Selection

Problem: 2^M combinations for M candidates

Solution: select outliers greedily from candidate set based on weights

