

Summary Report Optimization in ARA (Aggregate API)

WICG Review on Blog Post

aksu@google.com,

10 December 2023

Outline

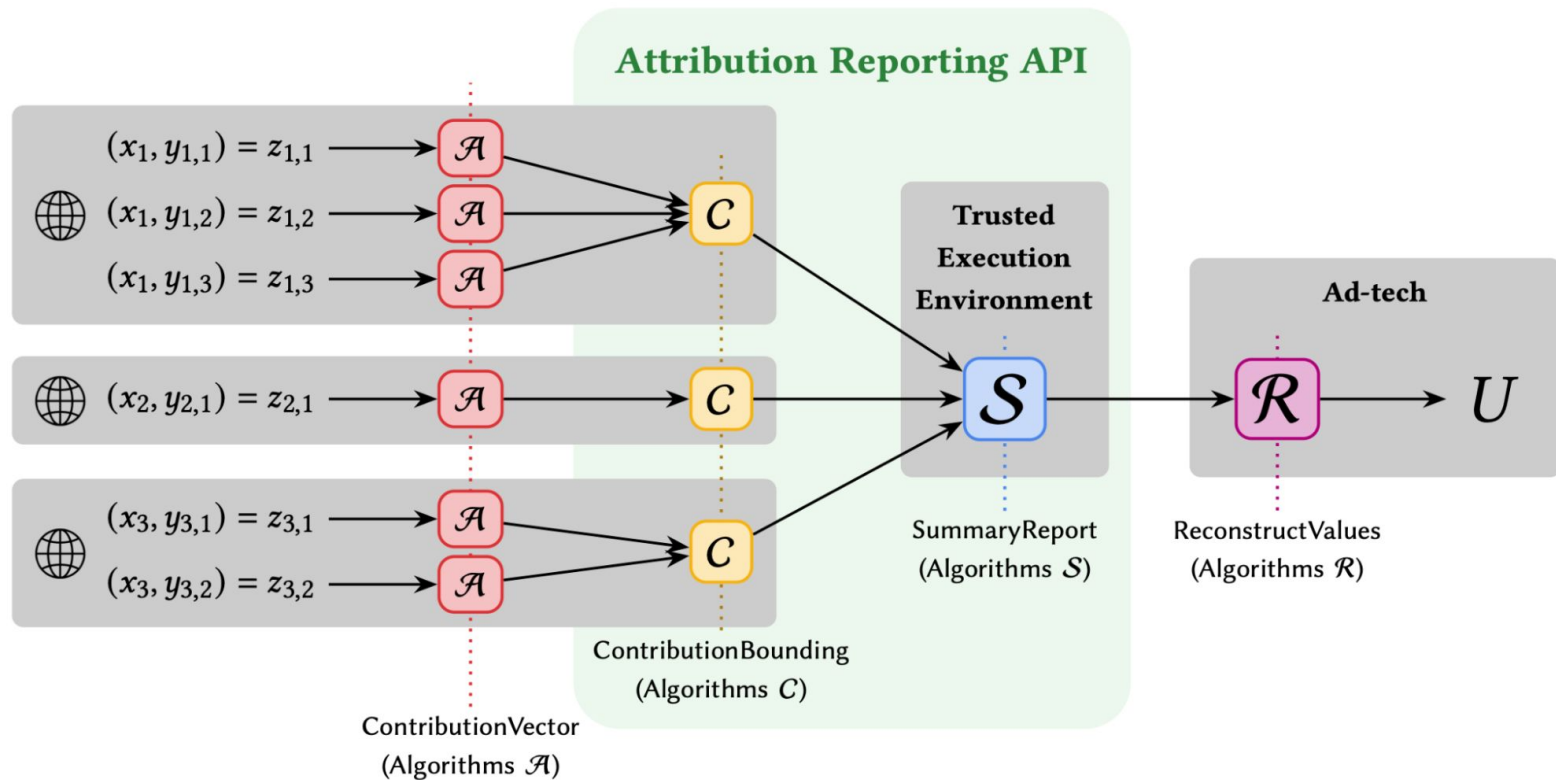
- Context & Problem
 - ARA Summary Reports
 - Maximizing the utility of summary reports
- Approach:
 - Mathematical model
 - Error metric
 - Optimization
- Synthetic Data
- Experimental Evaluations & Results
- Conclusions

ARA Summary Reports

	Impression features x			Conversion features y	
	Impression ID	Campaign	City	#items	value (\$)
$z_1 \rightarrow$	123	Thanksgiving	New York	3	21
$z_2 \rightarrow$	123	Thanksgiving	New York	1	5
$z_3 \rightarrow$	456	Thanksgiving	Boston	1	99
$z_4 \rightarrow$	123	Thanksgiving	New York	2	23
$z_5 \rightarrow$	101	Christmas	Boston	2	50
$z_6 \rightarrow$	789	Christmas	New York	3	15
$z_7 \rightarrow$	101	Christmas	Boston	1	5
...

Impression and conversion feature logs for a fictional gift shop called Du & Penc

Mathematical Model



Error Metrics

We have chosen:

- [\$\tau\$ -truncated root mean square relative error](#)

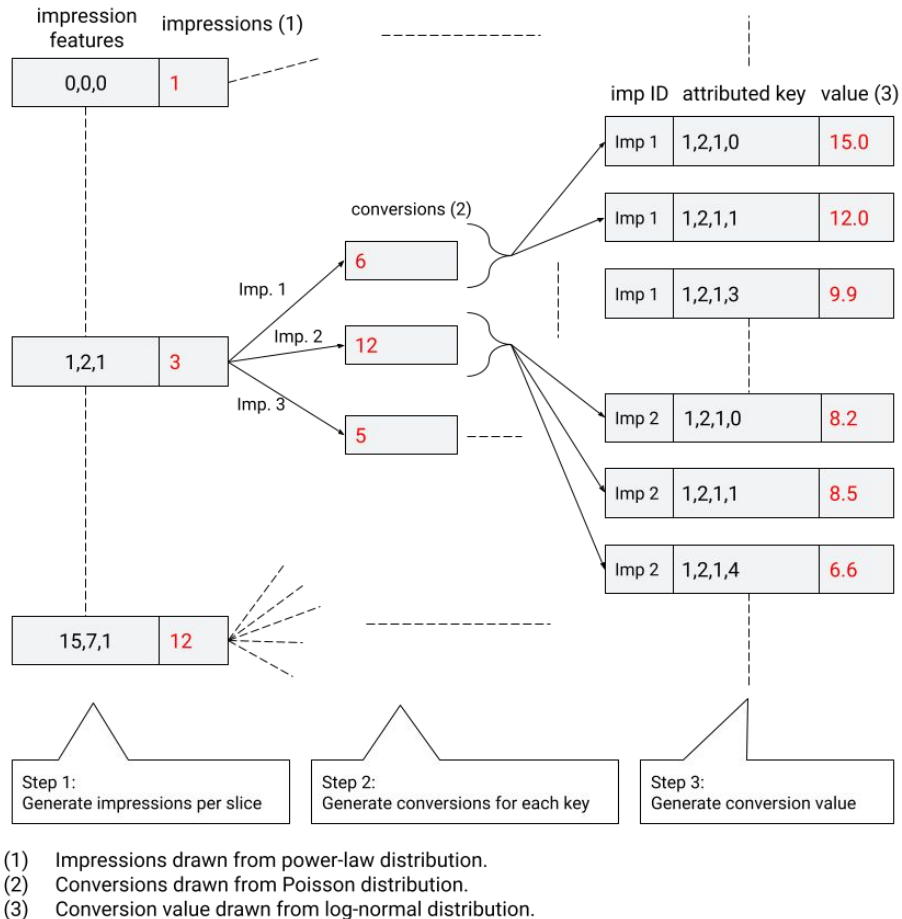
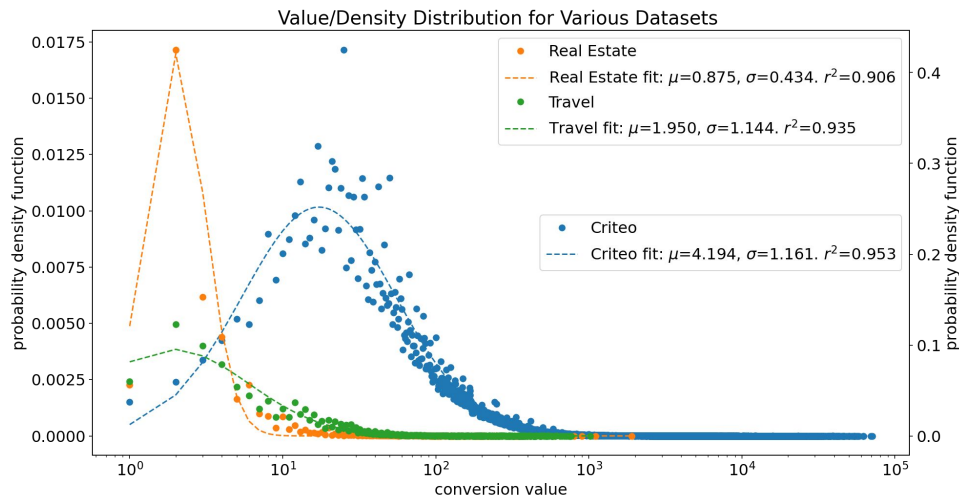
$$RMSRE_T(noise, true) = \sqrt{E\left(\left[\frac{|noise|}{\max\{T, true\}}\right]^2\right)}$$

Optimization

- The problem
 - maximize the utility of summary reports
- Reduce to optimization problem
 - optimize utility as measured by RMSRE_τ ,
 - capping parameter, C ,
 - privacy budget, α , for each slice
- Observations
 - RMSRE_τ = the bias from clipping and the variance of the noise distribution.
 - for for a fixed privacy budget, α , or a capping parameter, C ,
 - finding the optimal params is convex
 - for joint variables (C, α) it becomes non-convex
 - off-the-shelf optimizers works

Synthetic Data

- Generated synthetic data using
 - power law,
 - Poisson
 - log normal distributions



Experimental Evaluation

Datasets:

- Criteo: 15M clicks
- Real Estate: 100K conversions
- Travel: 30K conversions
- 3 Synthetic ones

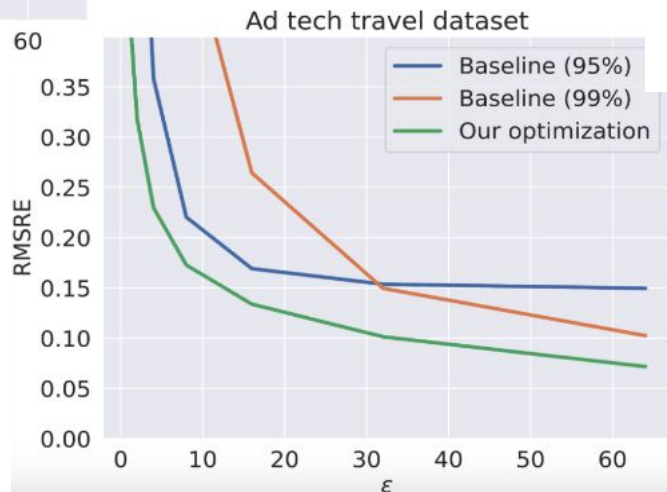
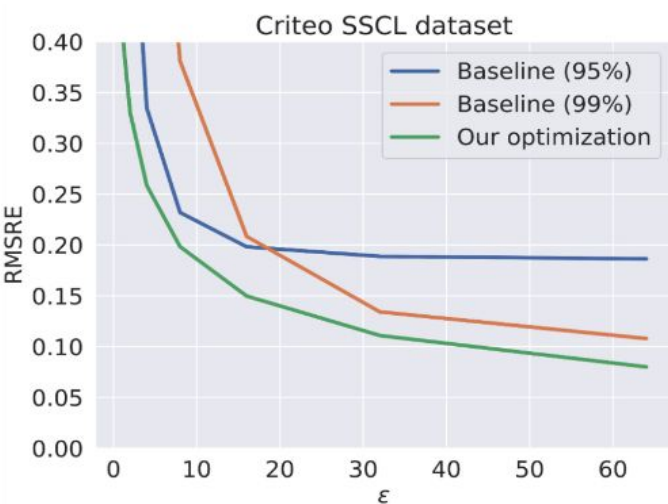
Dataset Partitioning:

- Training set: Choose budgets, thresholds, and limits
- Test set: Evaluate errors

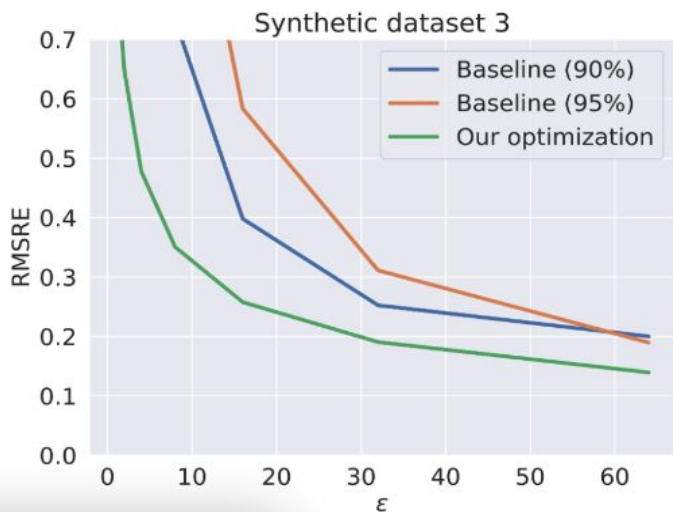
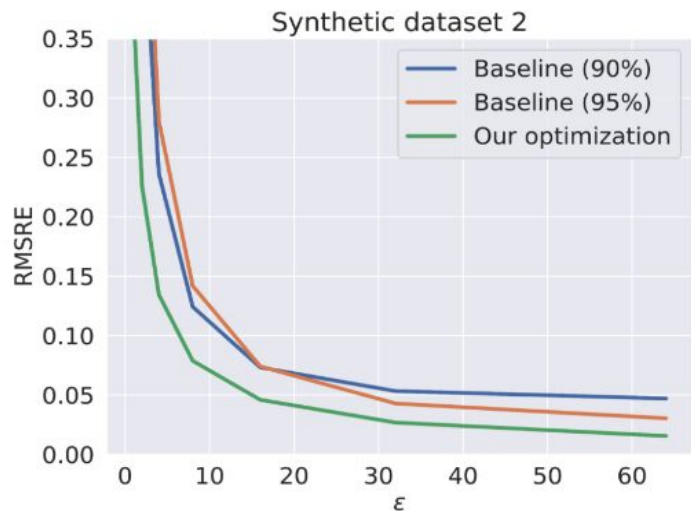
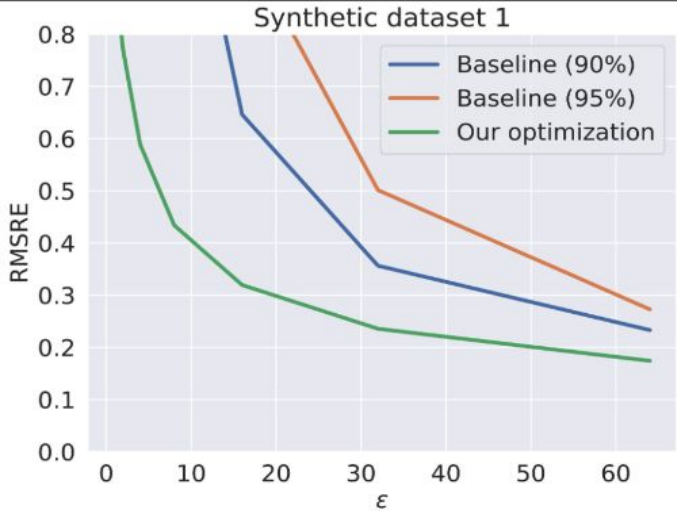
Error Metric: RMSRE_{τ} :

- Chosen to be 5 times the median value on training data
- Invariant to data rescaling
- Allows combining errors across different scales

Results



Baselines that use a fixed quantile for the clipping threshold and split the contribution budget equally among the queries.



Conclusion

- Leverages historical data
 - bound and scale the contributions of future data
 - use synthetic data
- Paper provide **Generalization bounds**
 - we're not overfitting to the historical data
- [Blog post](#), the [paper](#) and [accompanying code](#) are public.