

Automatic, Fine-Grained Algorithmic Choice for Differential Privacy

Jacob Imola (Advisor: Jean Yang)

Carnegie Mellon University

2017

Motivation

Netflix wants to anonymize database and release it:

	Titanic	Date	The Notebook	Date
Jordan	8.5	1/17	7.5	3/20
Jean	6	3/5	5	3/7
Scott	2	1/21	9	2/26
Serena	4.5	4/5	9	3/7

Motivation

Just cross out names!

	Titanic		The Notebook	
		Date		Date
User 1	8.5	1/17	7.5	3/20
User 2	6	3/5	5	3/7
User 3	2	1/21	9	2/26
User 4	4.5	4/5	9	3/7

Months later...

Motivation

Netflix

	Titanic	Date	The Notebook	Date
User 1	8.5	1/17	7.5	3/20
User 2	6	3/5	5	3/7
User 3	2	1/21	9	2/26
User 4	4.5	4/5	9	3/7

IMDB

Scott on 1/22:
(1.5/5) Titanic was terrible!

Jordan on 1/20:
(4.0/5) Enjoyed Titanic!

Jean on 3/6:
(2.0/5) The Notebook was
pretty overrated :/

Netflix promised that all movie ratings would be protected!

Differential Privacy

Definition

For all D and D' differing in 1 row: P is ϵ -DP if $\Pr(P(D) = O) < e^\epsilon \Pr(P(D') = O)$ for all O .

$$\underbrace{P \left(\begin{array}{ccccc} \text{Jordan} & 8.5 & 1/17 & 7.5 & 3/20 \\ \text{Jean} & 6 & 3/5 & 5 & 3/7 \\ \text{Scott} & 2 & 1/21 & 9 & 2/26 \\ \text{Serena} & 4.5 & 4/5 & 9 & 3/7 \end{array} \right)}_D = \underbrace{\begin{array}{ccccc} \text{User 1} & 8.5 & 1/17 & 7.5 & 3/20 \\ \text{User 2} & 6 & 3/5 & 5 & 3/7 \\ \text{User 3} & 2 & 1/21 & 9 & 2/26 \\ \text{User 4} & 4.5 & 4/5 & 9 & 3/7 \end{array}}_O$$
$$\underbrace{P \left(\begin{array}{ccccc} \text{Jordan} & 8.5 & 1/17 & 7.5 & 3/20 \\ \text{Jean} & 6 & 3/5 & 5 & 3/7 \\ \text{Scott} & 2 & 1/21 & 9 & 2/26 \\ \text{Serena} & 8.5 & 4/6 & 3 & 1/20 \end{array} \right)}_{D'} = \underbrace{\begin{array}{ccccc} \text{User 1} & 8.5 & 1/17 & 7.5 & 3/20 \\ \text{User 2} & 6 & 3/5 & 5 & 3/7 \\ \text{User 3} & 2 & 1/21 & 9 & 2/26 \\ \text{User 4} & 8.5 & 4/6 & 3 & 1/20 \end{array}}_{O'}$$

Violation: $\Pr(P(D) = O) = 1$ and $\Pr(P(D') = O) = 0$

Example

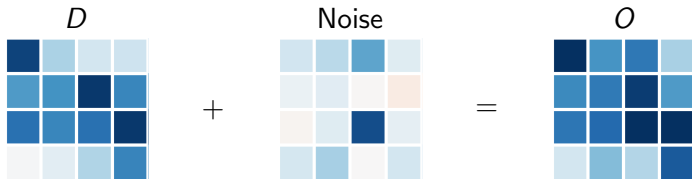
A representation change

User 1	8.5	$1/17$	7.5	$3/20$
User 2	6	$3/5$	5	$3/7$
User 3	2	$1/21$	9	$2/26$
User 4	4.5	$4/5$	9	$3/7$



Example

Method 1



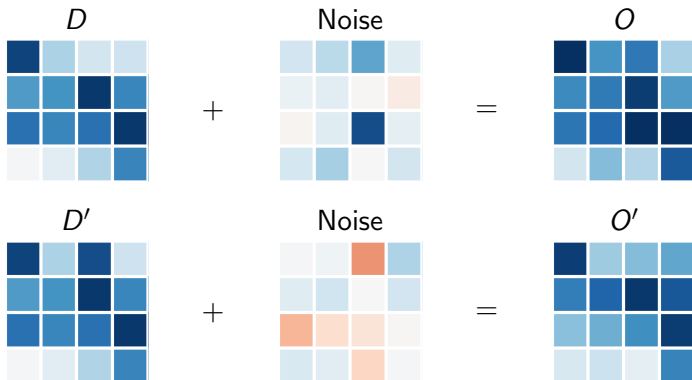
User 1	8.5	1/17	7.5	3/20
User 2	6	3/5	5	3/7
User 3	2	1/21	9	2/26
User 4	4.5	4/5	9	3/7

0.3	1.5	-1.1	2.3
-0.1	0.3	-0.2	-0.7
0.0	-2.4	0.9	-0.9
-0.2	1.2	0.4	1.1

User 1	8.0	1/19	6.5	3/22
User 2	6	3/5	5	3/6
User 3	2	1/19	10	2/25
User 4	4.5	4/6	9.5	3/8

Example

Method 1



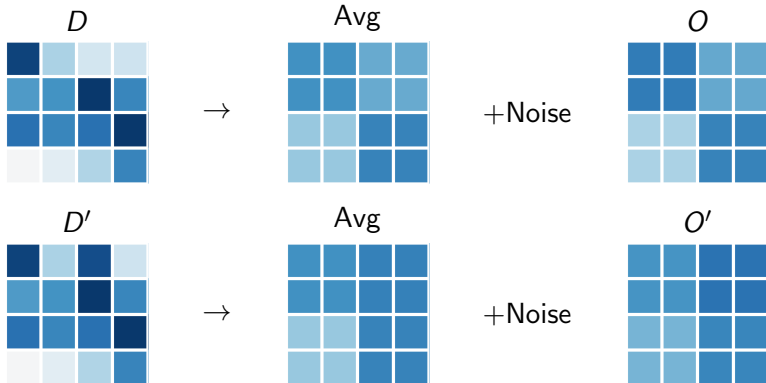
$$\Pr(P(D) = O) \approx 10^{-8} \quad \Pr(P(D') = O) \approx 2 \times 10^{-9}$$

Seeing O , attacker cannot distinguish D and D' .

Example

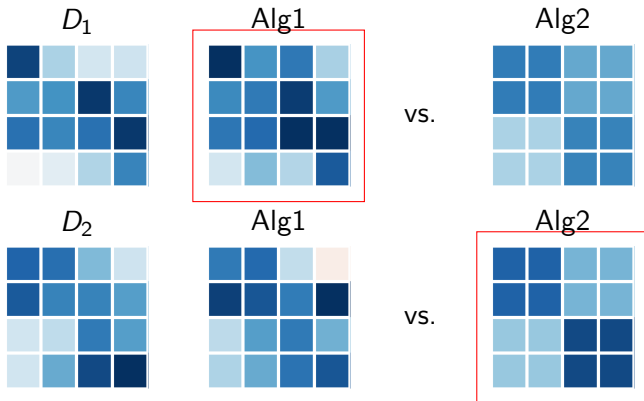
Method Two

- Sum into 4 buckets, add noise, divide by 4



Example

Which is better?



DP complicates algorithm analysis due to noise, makes algorithm deployment hard.

Vision

Task: Remove burden of DP algorithm analysis.

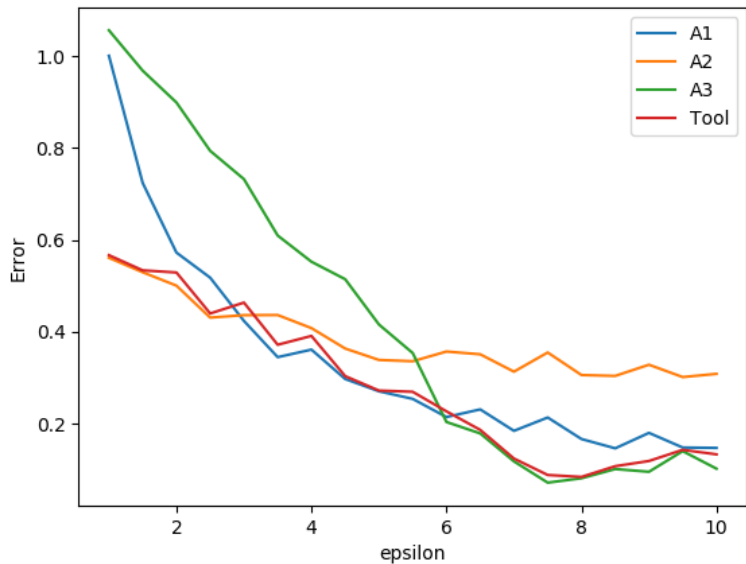
1. **Correctness** Differential privacy is never violated.
2. **Generalizability** Works on arbitrary code.
3. **Performance** Makes choice “close enough” to optimal.

Solution: A programming language (Jostle)!

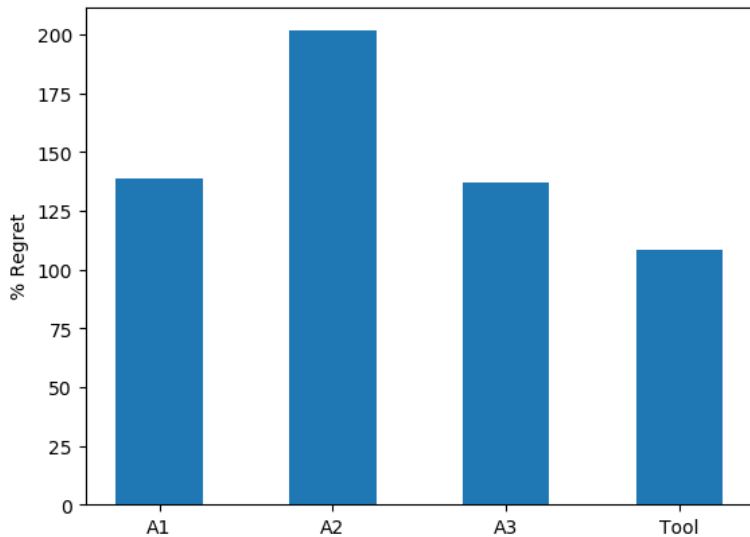
Represent any choice with `ChoiceMaker` construct.

```
1 answerHistQueries = MkChoiceMaker among {Alg1, Alg2}
2 answers = answerHistQueries(data, queries)
```

Note on performance



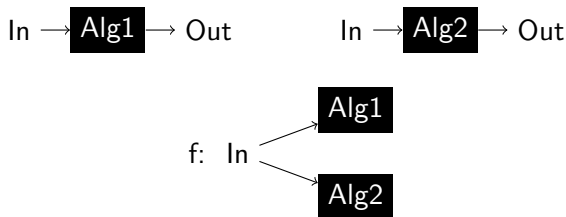
Note on performance



Challenges

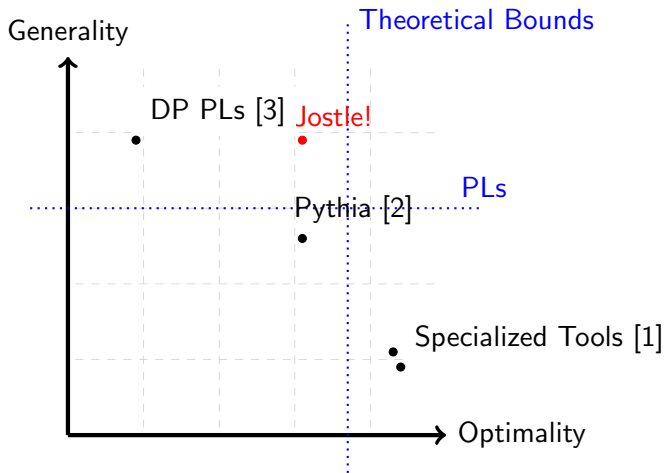
What's hard about writing this code?

```
1 answerHistQueries = MkChoiceMaker among {Alg1, Alg2}  
2 answers = answerHistQueries(data, queries)
```



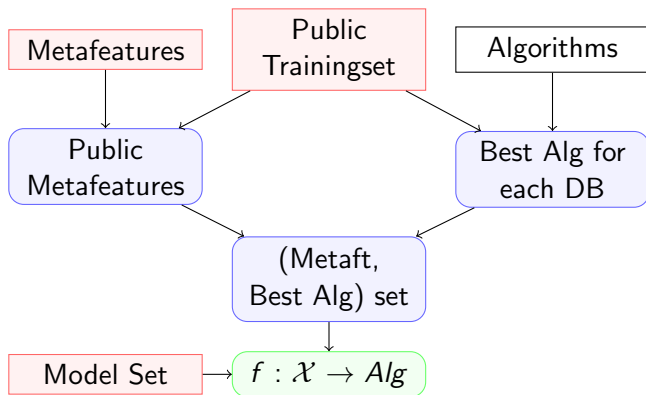
- ▶ Generality \implies Alg1, Alg2 are black boxes.
- ▶ Approach: Meta-machine learning: $f : DB \rightarrow Alg$
- ▶ Difficult—data science cannot be automated well.

Existing Work



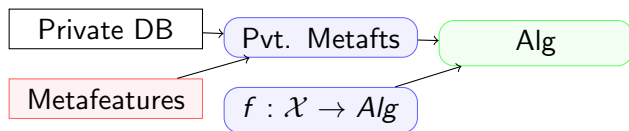
Solution Overview

- ▶ Metafeatures modeled after data science approach
- ▶ $f : DB \rightarrow Alg$ becomes $f : \mathcal{X} \rightarrow Alg$.



Important: Trainingset must have lots of DB's for training!

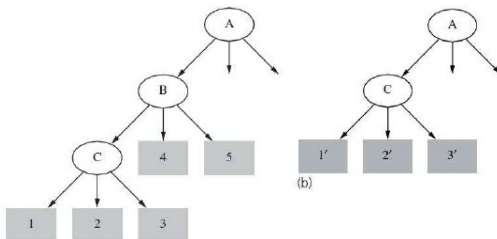
Solution Overview



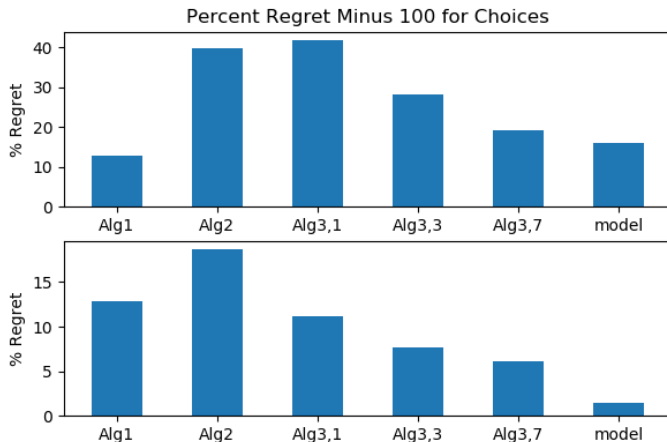
```
1 answerHistQueries =  
2   MkChoiceMaker among {Alg1, Alg2}  
3   informed by {dbSize, dbNumRows}  
4   modeled by LinearModel with ErrorFunc  
5   trained on TrainingSet}  
6  
7 answers = answerHistQueries(data, queries)
```

Experimental Setup

- ▶ **Algorithms** Stopping Criteria for Private Decision Trees.
- ▶ **Metafeatures** DB size, epsilon, domain size.
- ▶ **Classification** Linear Classifiers
- ▶ **Training Set**
 1. 300 real DB snapshots, 100 real DB snapshots
 2. 300 synth. DB snapshots, 100 synth. DB snapshots



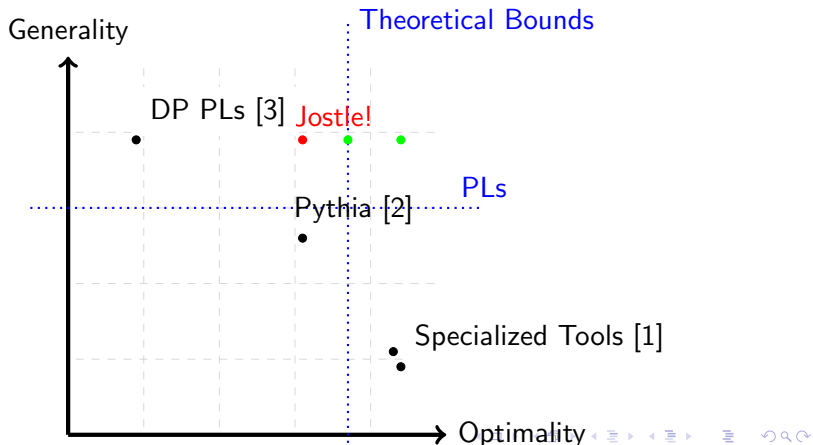
Results



Always possible to have similar test DB in training DB set?

Conclusion

- ▶ Performs as well as Pythia [2] with same expressiveness as PINQ [3].
- ▶ Only as good as how well the programmer frames the ML.
- ▶ Future work: Better optimization for multiple ChoiceMakers. Static analysis (open black boxes).



References



Kamalika Chaudhuri and Staal A. Vinterbo.

A stability-based validation procedure for differentially private machine learning.

In *NIPS*, 2013.



Ios Kotsogiannis, Ashwin Machanavajjhala, Michael Hay, and Gerome Miklau.

Pythia: Data dependent differentially private algorithm selection.

In *Proceedings of the 2017 ACM International Conference on Management of Data*, SIGMOD '17, pages 1323–1337, New York, NY, USA, 2017. ACM.



Frank McSherry.

Privacy integrated queries: An extensible platform for privacy-preserving data analysis.

Commun. ACM, 53(9):89–97, September 2010.