



Simple Binary Hypothesis Testing: Locally Private and Communication-Efficient

Ankit Pensia

ITA 2023



Joint Work With



Amir Asadi



Varun Jog



Po-Ling Loh

Outline

- ▶ Motivation
- ▶ Problem Statement
- ▶ Our Results
 - ▶ Statistical
 - ▶ Computational
- ▶ Proof Sketch
- ▶ Conclusion

Simple Hypothesis Testing: Centralized

- Let p and q be two known distributions over $\{1, \dots, k\}$

Problem (Simple Hypothesis Testing):

Input: i.i.d. samples from either p or q



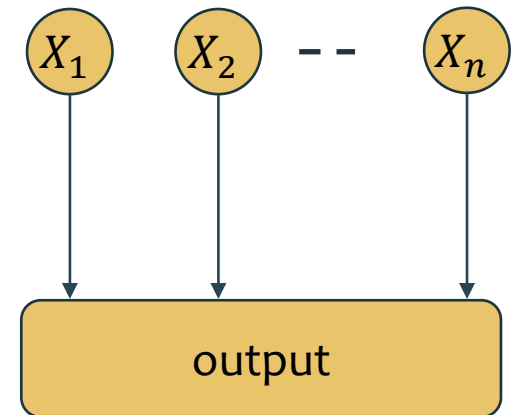
Simple Hypothesis Testing: Centralized

- Let p and q be two known distributions over $\{1, \dots, k\}$

Problem (Simple Hypothesis Testing):

Input: i.i.d. samples from either p or q

Output: whether they came from p or q



Simple Hypothesis Testing: Centralized

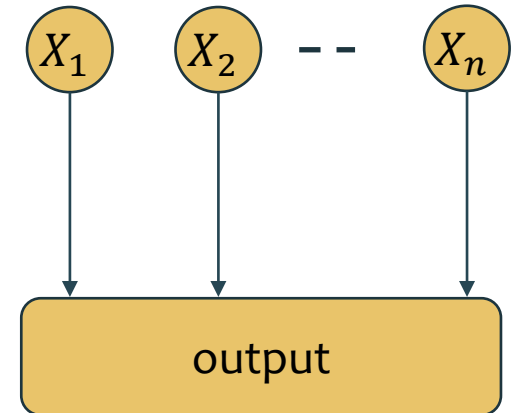
- Let p and q be two known distributions over $\{1, \dots, k\}$

Problem (Simple Hypothesis Testing):

Input: i.i.d. samples from either p or q

Output: whether they came from p or q

- Arguably, the simplest statistical problem



Simple Hypothesis Testing: Centralized

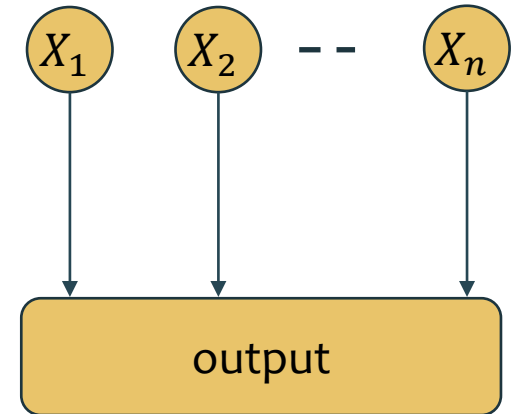
- Let p and q be two known distributions over $\{1, \dots, k\}$

Problem (Simple Hypothesis Testing):

Input: i.i.d. samples from either p or q

Output: whether they came from p or q

- Arguably, the simplest statistical problem
 - Optimal test: Likelihood ratio test



Simple Hypothesis Testing: Centralized

- Let p and q be two known distributions over $\{1, \dots, k\}$

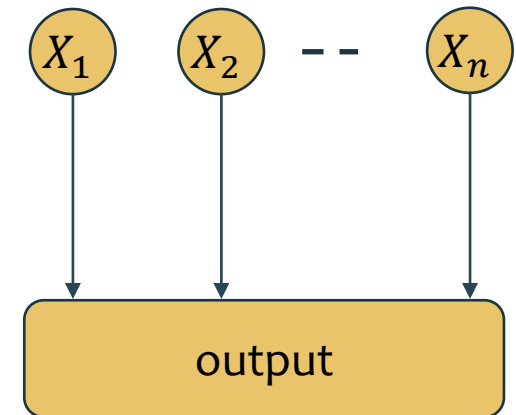
Problem (Simple Hypothesis Testing):

Input: i.i.d. samples from either p or q

Output: whether they came from p or q

- Arguably, the simplest statistical problem
 - Optimal test: Likelihood ratio test

Requires access to X_i 's



Simple Hypothesis Testing: Centralized

- Let p and q be two known distributions over $\{1, \dots, k\}$

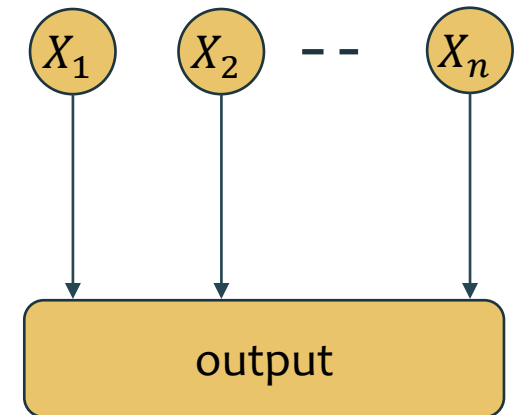
Problem (Simple Hypothesis Testing):

Input: i.i.d. samples from either p or q

Output: whether they came from p or q

- Arguably, the simplest statistical problem
 - Optimal test: Likelihood ratio test
- Data is distributed these days
 - Limited communication bandwidth
 - Privacy concerns

Requires access to X_i 's



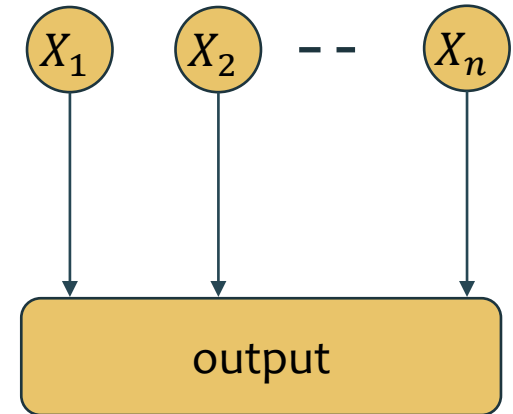
Simple Hypothesis Testing: Centralized

- Let p and q be two known distributions over $\{1, \dots, k\}$

Problem (Simple Hypothesis Testing):

Input: i.i.d. samples from either p or q

Output: whether they came from p or q



- Arguably, the simplest statistical problem

- Optimal test: Likelihood ratio test

Requires access to X_i 's

- Data is distributed these days

- Limited communication bandwidth

- Privacy concerns

Requires quantizing/privatizing X_i 's

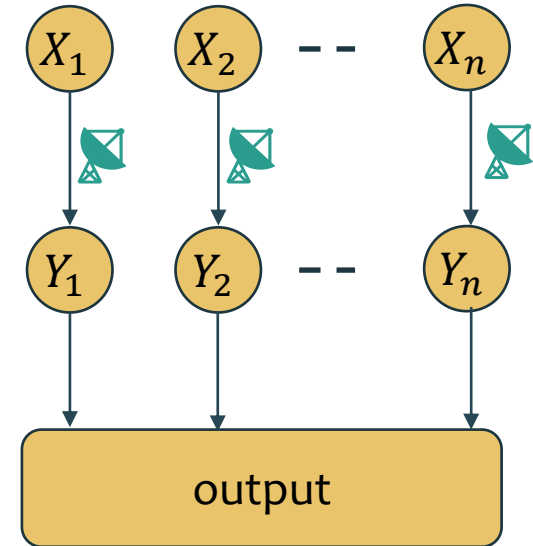
Simple Hypothesis Testing: Decentralized

- Let p and q be two known distributions over $\{1, \dots, k\}$

Problem (Simple Hypothesis Testing):

Input: i.i.d. samples from either p or q

Output: whether they came from p or q



Simple Hypothesis Testing: Decentralized

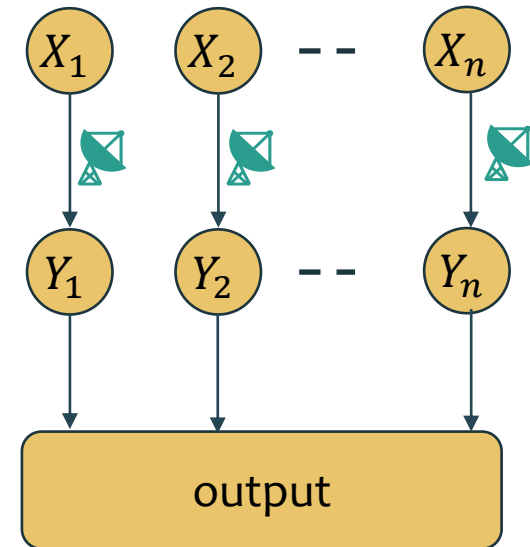
- Let p and q be two known distributions over $\{1, \dots, k\}$

Problem (Simple Hypothesis Testing):

Input: i.i.d. samples from either p or q

Output: whether they came from p or q

- : captures communication and/or privacy



Simple Hypothesis Testing: Decentralized

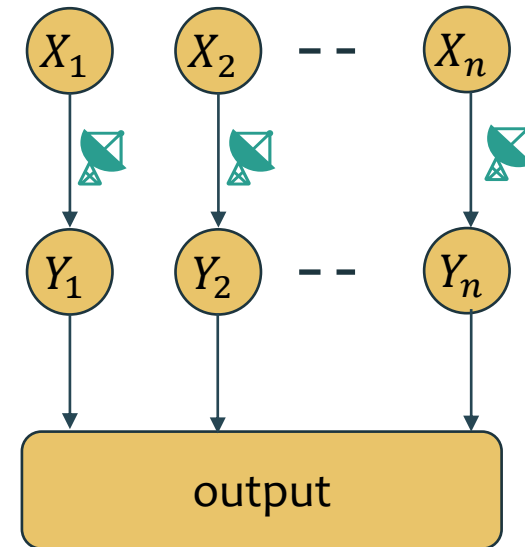
- Let p and q be two known distributions over $\{1, \dots, k\}$

Problem (Decentralized Simple Hypothesis Testing):

Input: **modified** samples from either p or q

Output: whether they came from p or q

- : captures communication and/or privacy



Simple Hypothesis Testing: Decentralized

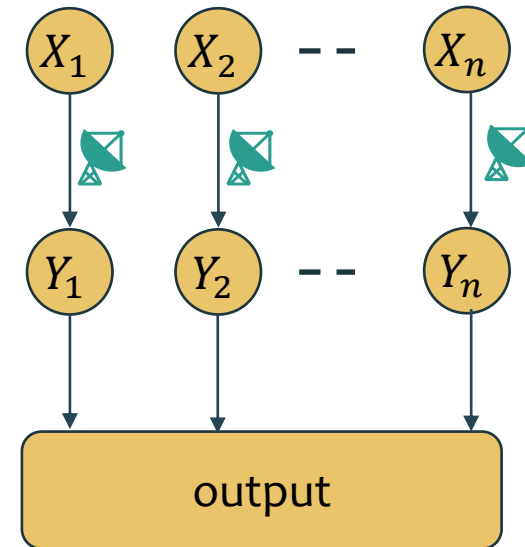
- Let p and q be two known distributions over $\{1, \dots, k\}$

Problem (Decentralized Simple Hypothesis Testing):

Input: **modified** samples from either p or q

Output: whether they came from p or q

- : captures communication and/or privacy



Simple Hypothesis Testing: Decentralized

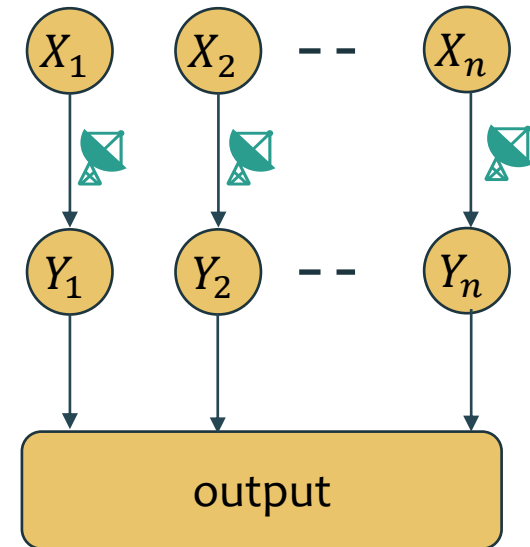
- Let p and q be two known distributions over $\{1, \dots, k\}$

Problem (Decentralized Simple Hypothesis Testing):

Input: **modified** samples from either p or q

Output: whether they came from p or q

- : captures communication and/or privacy



How do we perform decentralized hypothesis testing?

Outline

- ▶ Motivation
- ▶ **Problem Statement**
- ▶ Our Results
 - ▶ Statistical
 - ▶ Computational
- ▶ Proof Sketch
- ▶ Conclusion

Privacy Model and Communication Constraints

- Local Differential Privacy (LDP)

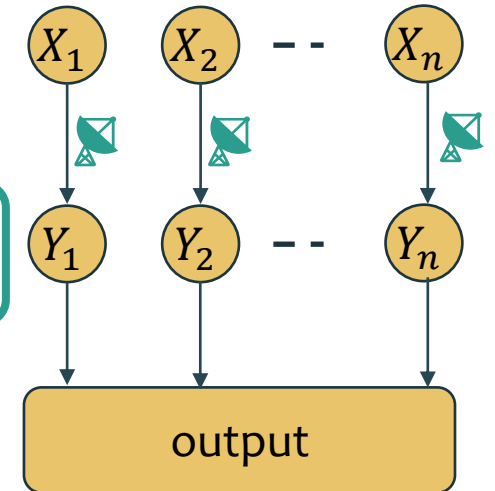
- Everyone releases a randomized version of data

- Channel  is ϵ -LDP if:

$$\frac{\mathbb{P}(Y_i=y | X_i=x)}{\mathbb{P}(Y_i=y | X_i=x')} \leq e^\epsilon \text{ for all } x, x', y$$

Can't reliably distinguish between x and x' using values of Y_i

- Non-interactive (private-coin): Y_i 's are independent



Privacy Model and Communication Constraints

- Local Differential Privacy (LDP)

- Everyone releases a randomized version of data

- Channel  is ϵ -LDP if:

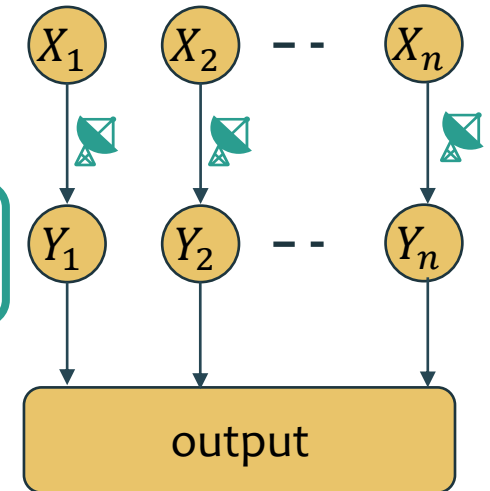
$$\frac{\mathbb{P}(Y_i=y | X_i=x)}{\mathbb{P}(Y_i=y | X_i=x')} \leq e^\epsilon \text{ for all } x, x', y$$

Can't reliably distinguish between x and x' using values of Y_i

- Non-interactive (private-coin): Y_i 's are independent


- Communication-constraints

- $Y_i \in \{1, \dots, \ell\}$ for some $\ell \ll k$



Privacy Model and Communication Constraints

- Local Differential Privacy (LDP)

- Everyone releases a randomized version of data
- Channel  is ϵ -LDP if:

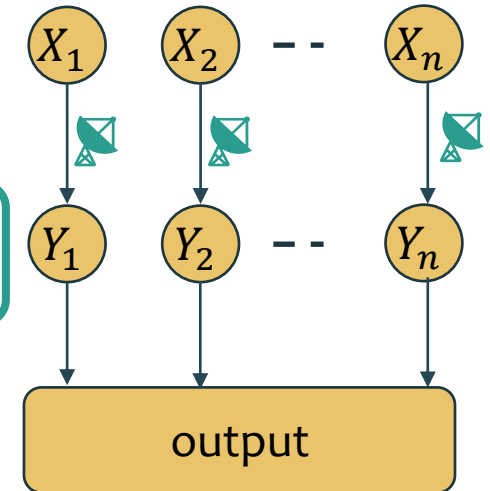
$$\frac{\mathbb{P}(Y_i=y | X_i=x)}{\mathbb{P}(Y_i=y | X_i=x')} \leq e^\epsilon \text{ for all } x, x', y$$

Can't reliably distinguish between x and x' using values of Y_i

- Non-interactive (private-coin): Y_i 's are independent

- Communication-constraints

- $Y_i \in \{1, \dots, \ell\}$ for some $\ell \ll k$



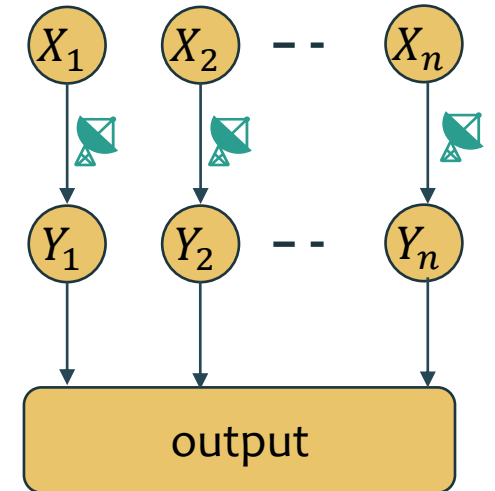
Today's focus: Privacy (LDP)

Questions of Interest

Problem (Decentralized Simple Hypothesis Testing):

Input: modified samples from either p or q

Output: whether they came from p or q



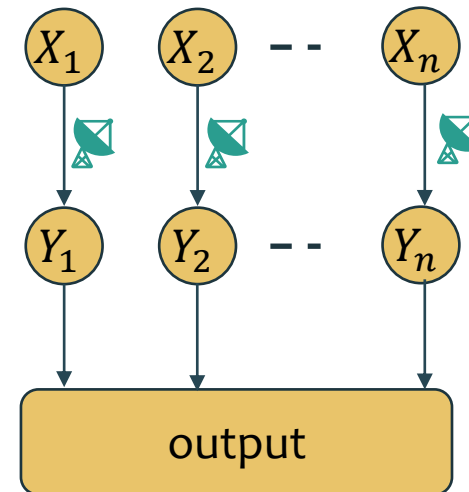
Questions of Interest

Problem (Decentralized Simple Hypothesis Testing):

Input: modified samples from either p or q

Output: whether they came from p or q

Goal: Design the test and channels  so that the probability of error ≤ 0.1



Questions of Interest

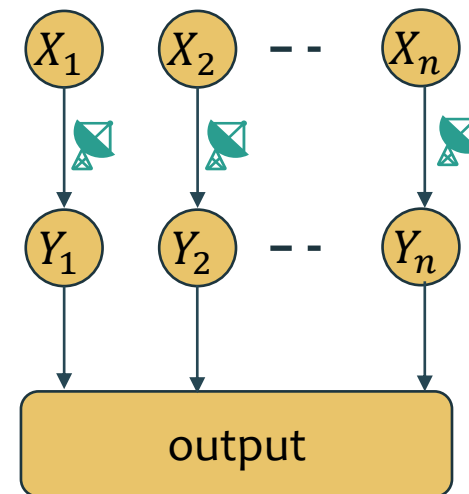
Problem (Decentralized Simple Hypothesis Testing):

Input: modified samples from either p or q

Output: whether they came from p or q

Goal: Design the test and channels  so that the probability of error ≤ 0.1

Sample Complexity: Minimum n to achieve above goal



Problem (Decentralized Simple Hypothesis Testing):

Input: modified samples from either p or q

Output: whether they came from p or q

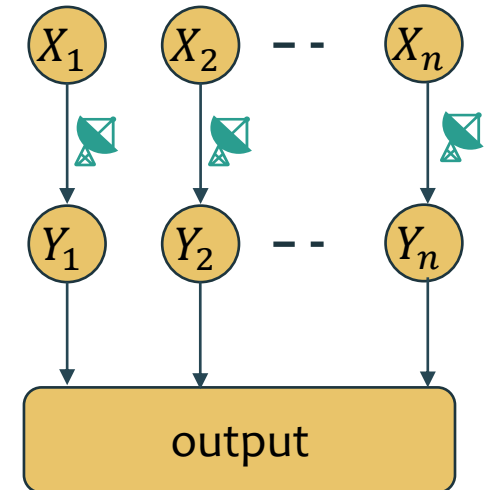
Questions of Interest

Goal: Design the test and channels  so that the probability of error ≤ 0.1

Sample Complexity: Minimum n to achieve above goal

n^* := Sample complexity (no constraints)

$n^*(\epsilon)$:= Sample complexity with channels satisfying ϵ -LDP



Problem (Decentralized Simple Hypothesis Testing):

Input: modified samples from either p or q

Output: whether they came from p or q

Questions of Interest

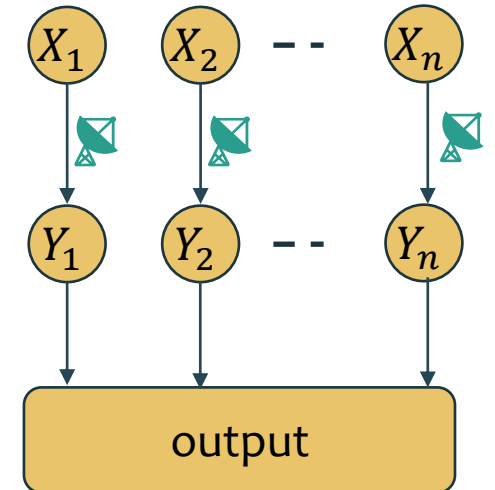
Goal: Design the test and channels  so that the probability of error ≤ 0.1

Sample Complexity: Minimum n to achieve above goal

n^* := Sample complexity (no constraints)

$n^*(\epsilon)$:= Sample complexity with channels satisfying ϵ -LDP

Questions:



Problem (Decentralized Simple Hypothesis Testing):

Input: modified samples from either p or q

Output: whether they came from p or q

Questions of Interest

Goal: Design the test and channels  so that the probability of error ≤ 0.1

Sample Complexity: Minimum n to achieve above goal

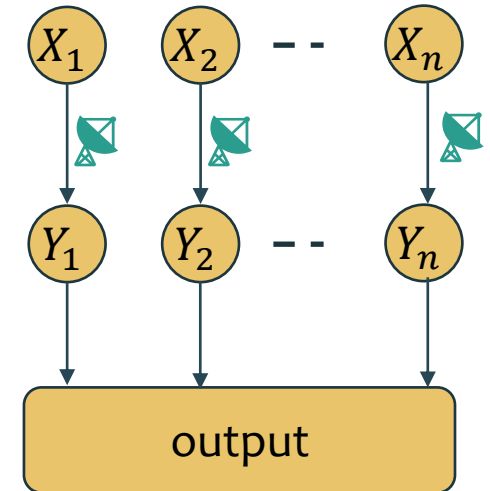
n^* := Sample complexity (no constraints)

$n^*(\epsilon)$:= Sample complexity with channels satisfying ϵ -LDP

Questions:

1. (Statistical) How much does sample complexity change?

$n^*(\epsilon)$ vs. n^*



Problem (Decentralized Simple Hypothesis Testing):

Input: modified samples from either p or q

Output: whether they came from p or q

Questions of Interest

Goal: Design the test and channels  so that the probability of error ≤ 0.1

Sample Complexity: Minimum n to achieve above goal

n^* := Sample complexity (no constraints)

$n^*(\epsilon)$:= Sample complexity with channels satisfying ϵ -LDP

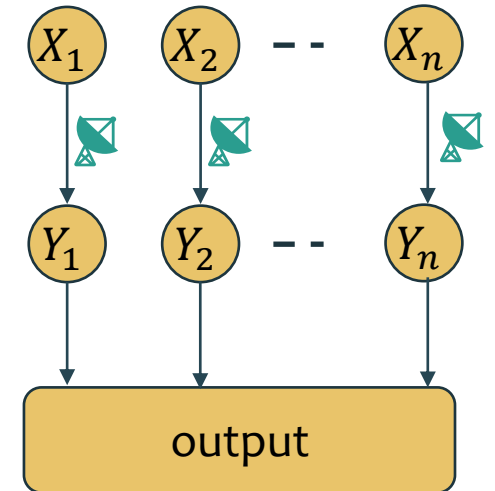
Questions:

1. (Statistical) How much does sample complexity change?

$n^*(\epsilon)$ vs. n^*

2. (Computational) How to find (near)-optimal channels fast?

polynomial in support size

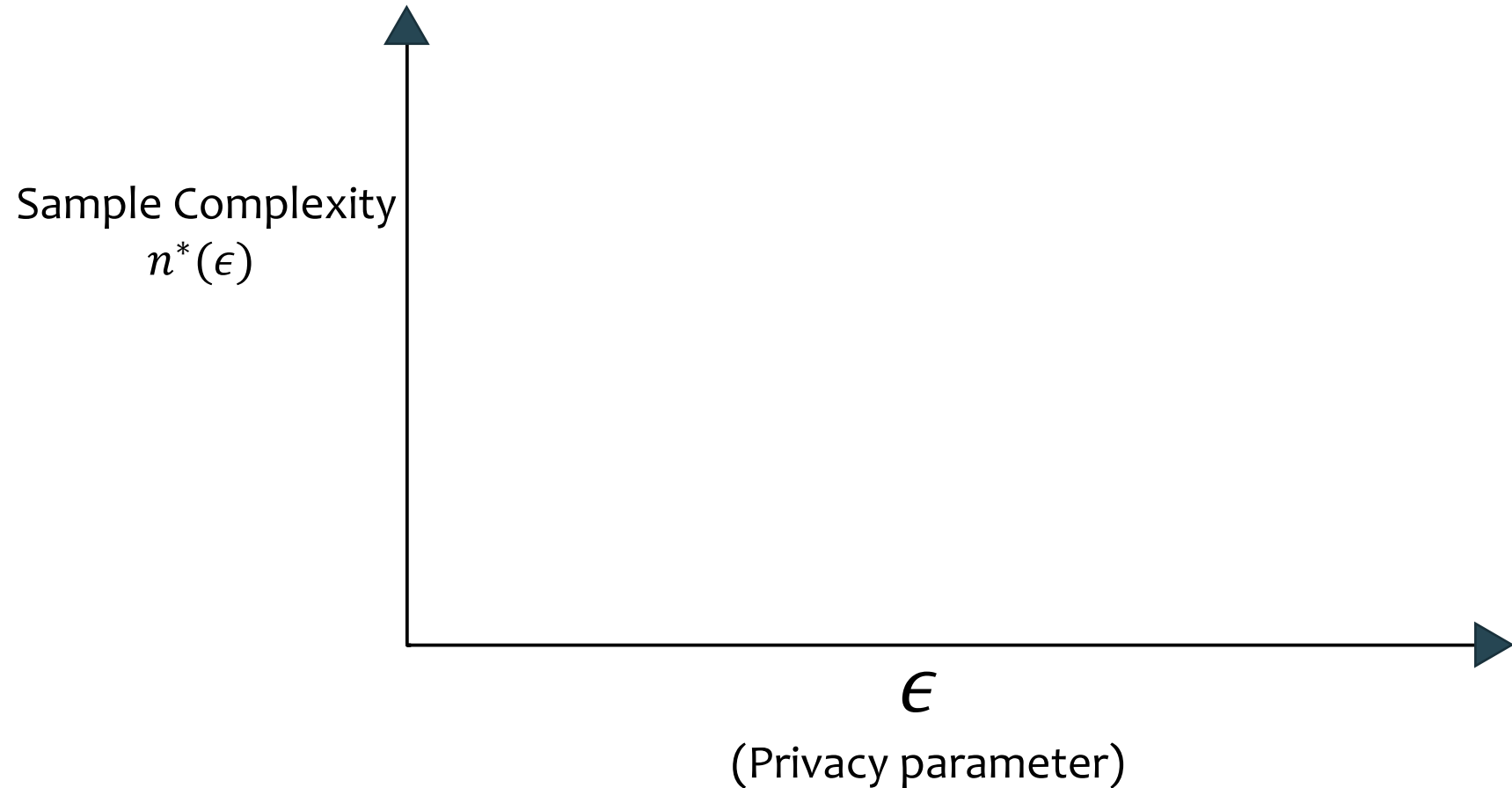


Outline

- ▶ Motivation
- ▶ Problem Statement
- ▶ **Our Results**
 - ▶ **Statistical**
 - ▶ Computational
- ▶ Proof Sketch
- ▶ Conclusion

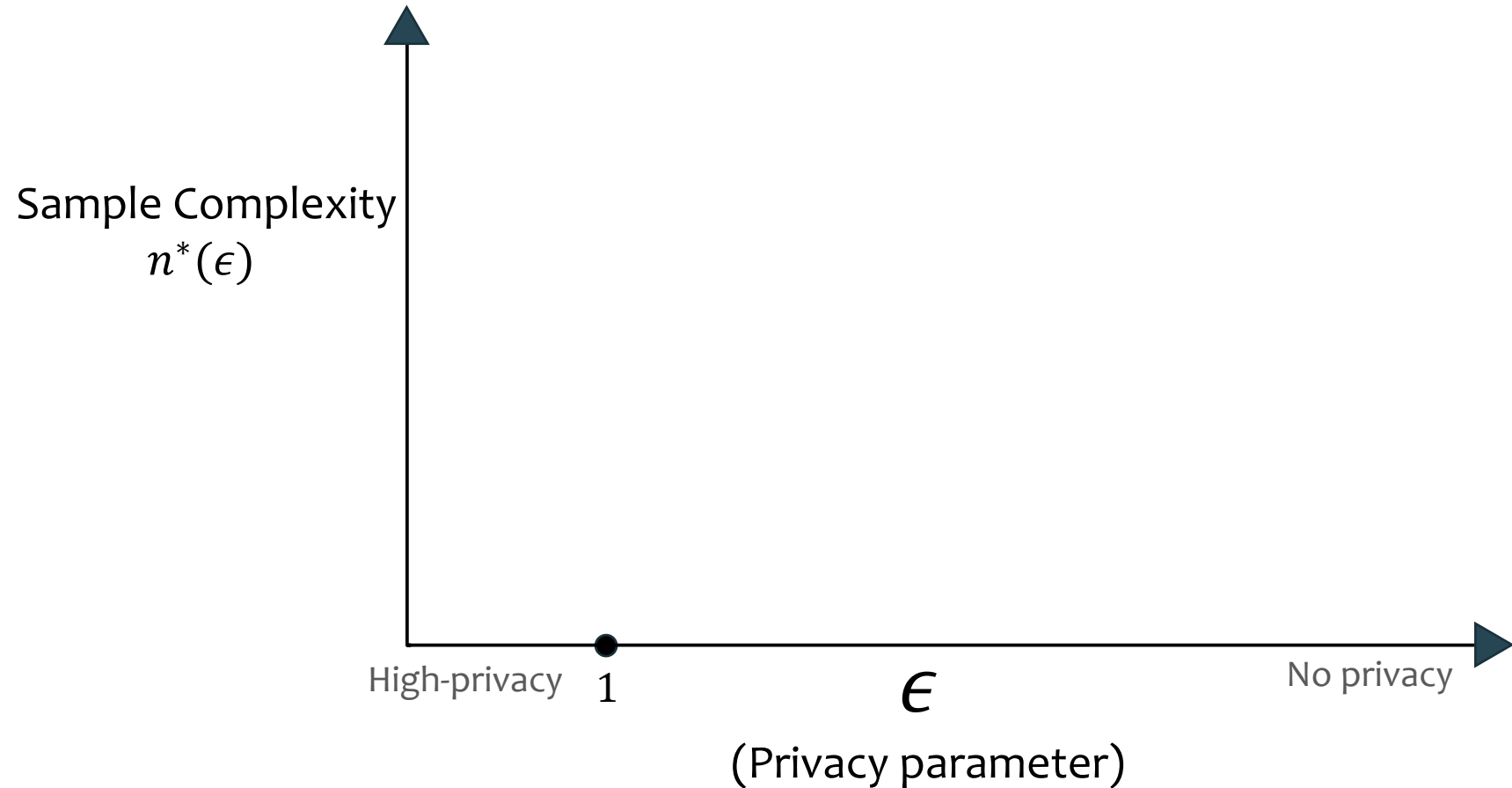
Statistical Cost of Privacy: Existing Results

- Sample Complexity



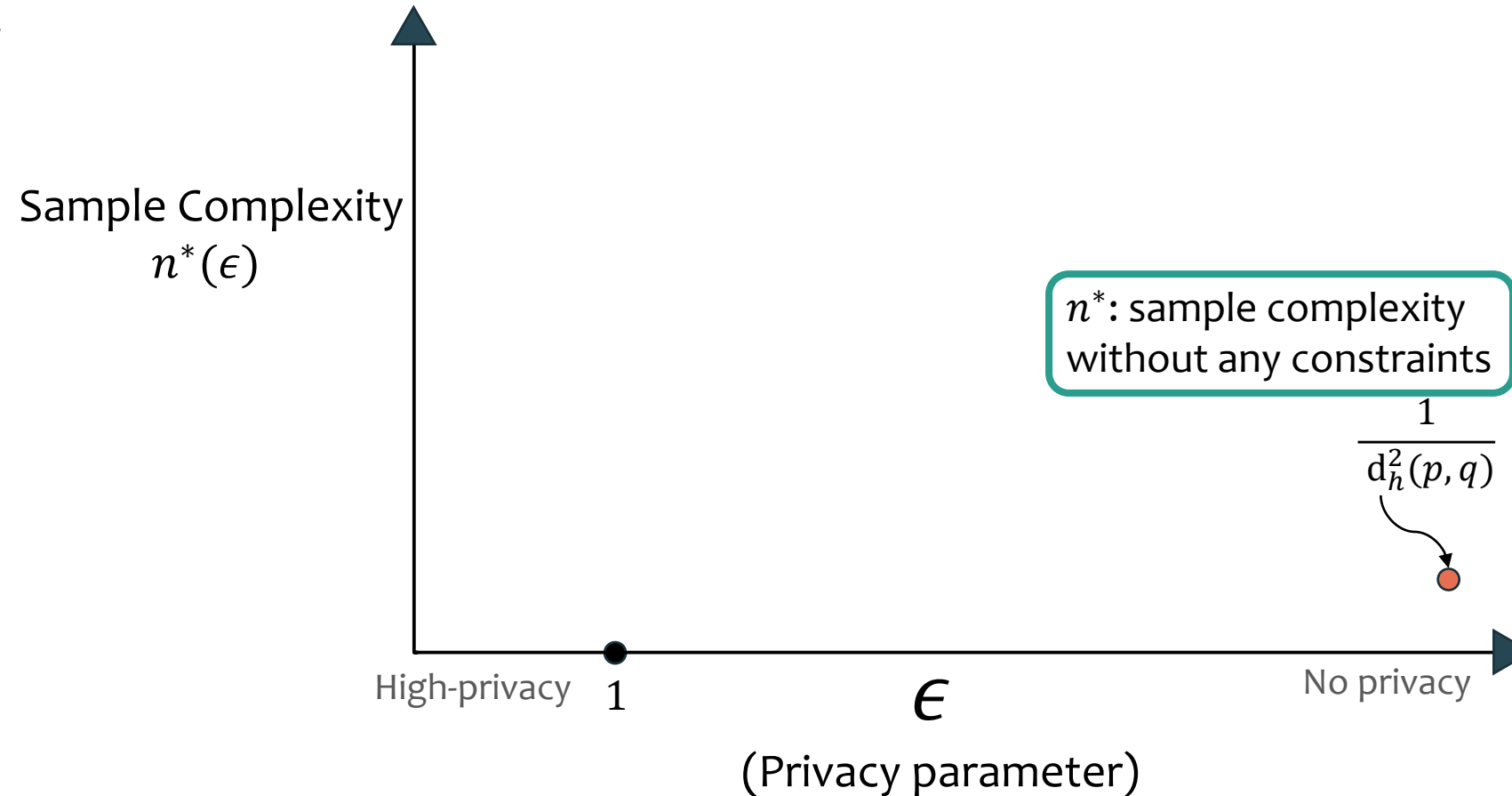
Statistical Cost of Privacy: Existing Results

- Sample Complexity



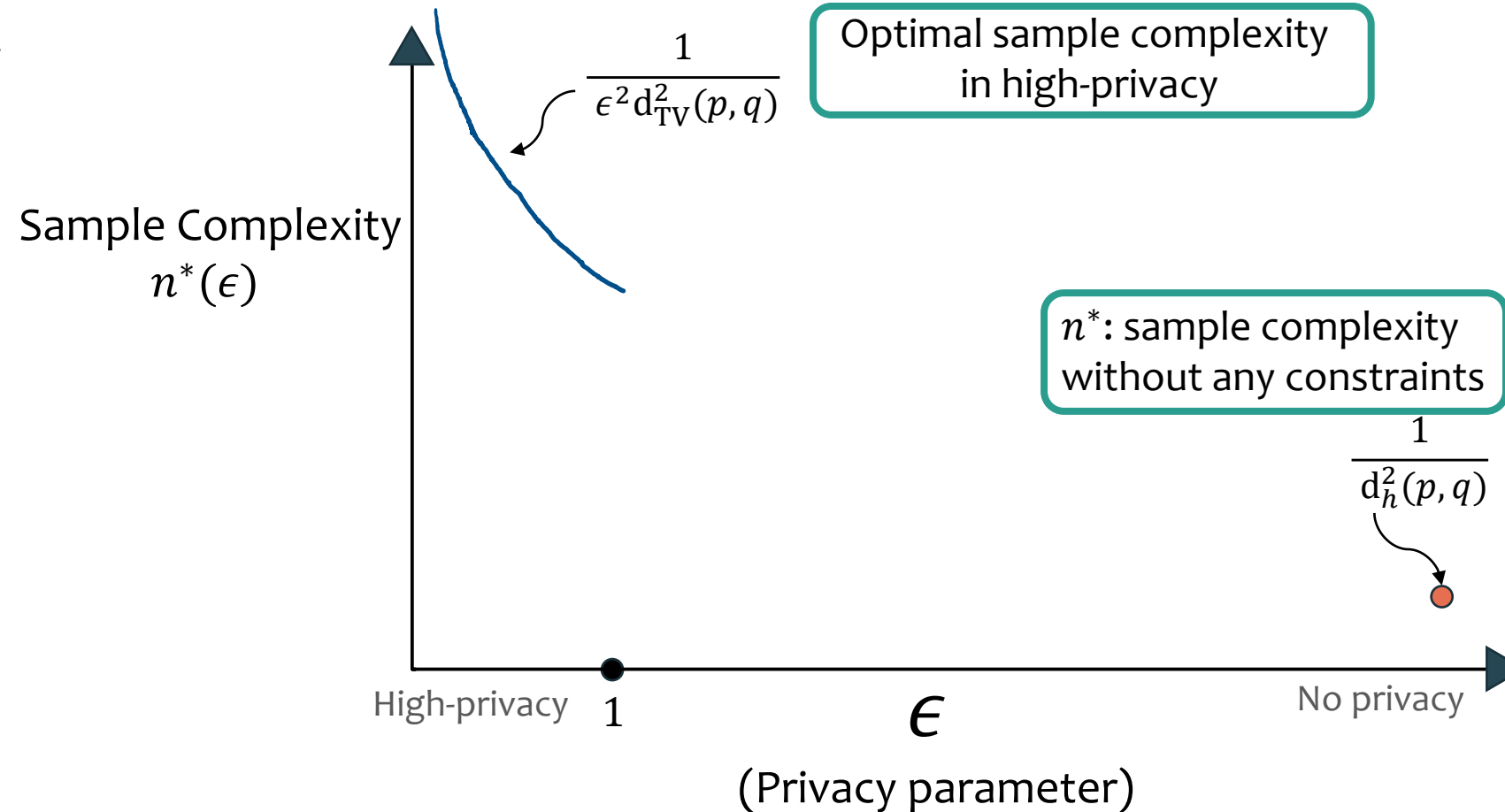
Statistical Cost of Privacy: Existing Results

- Sample Complexity



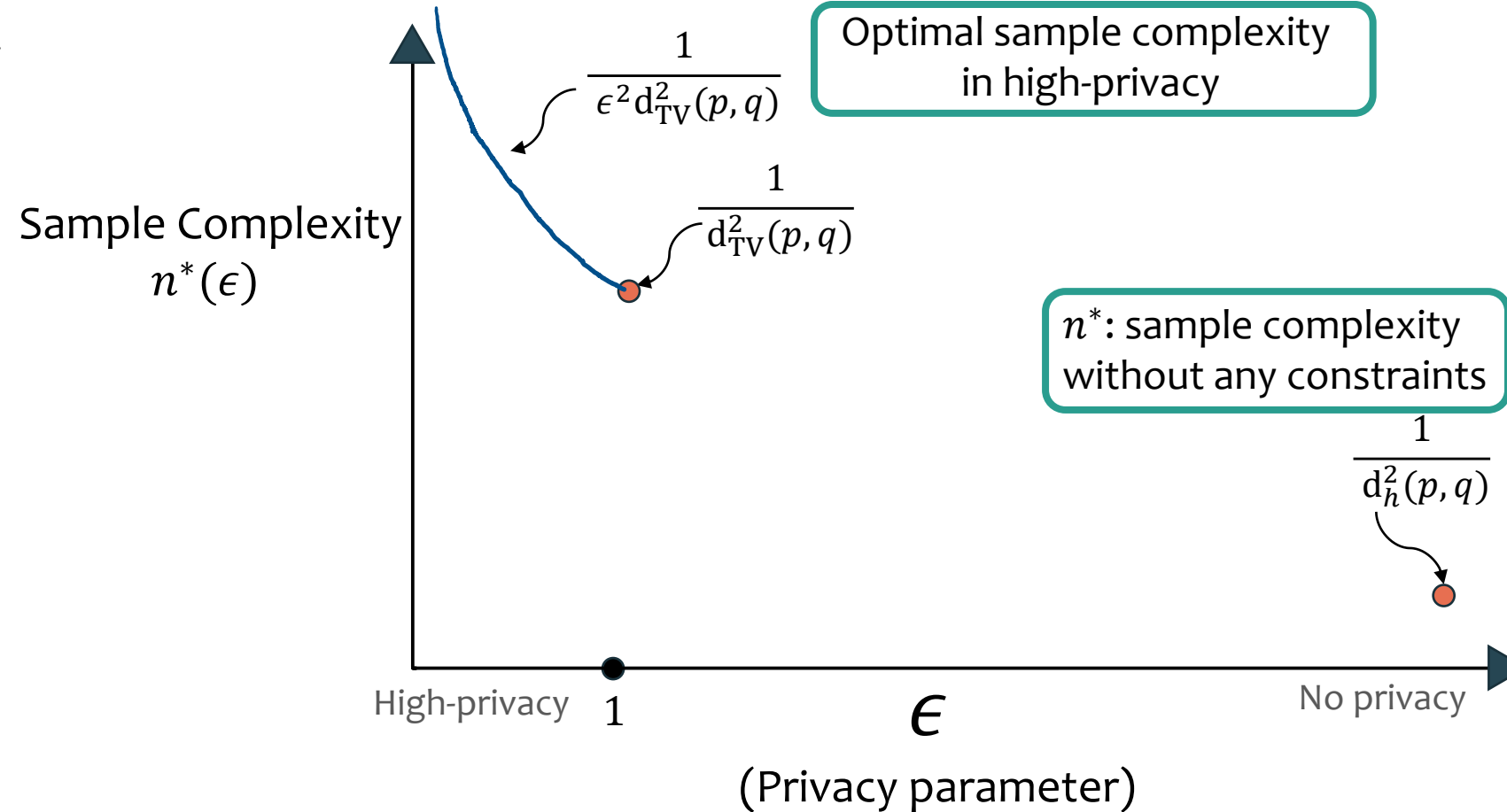
Statistical Cost of Privacy: Existing Results

- Sample Complexity



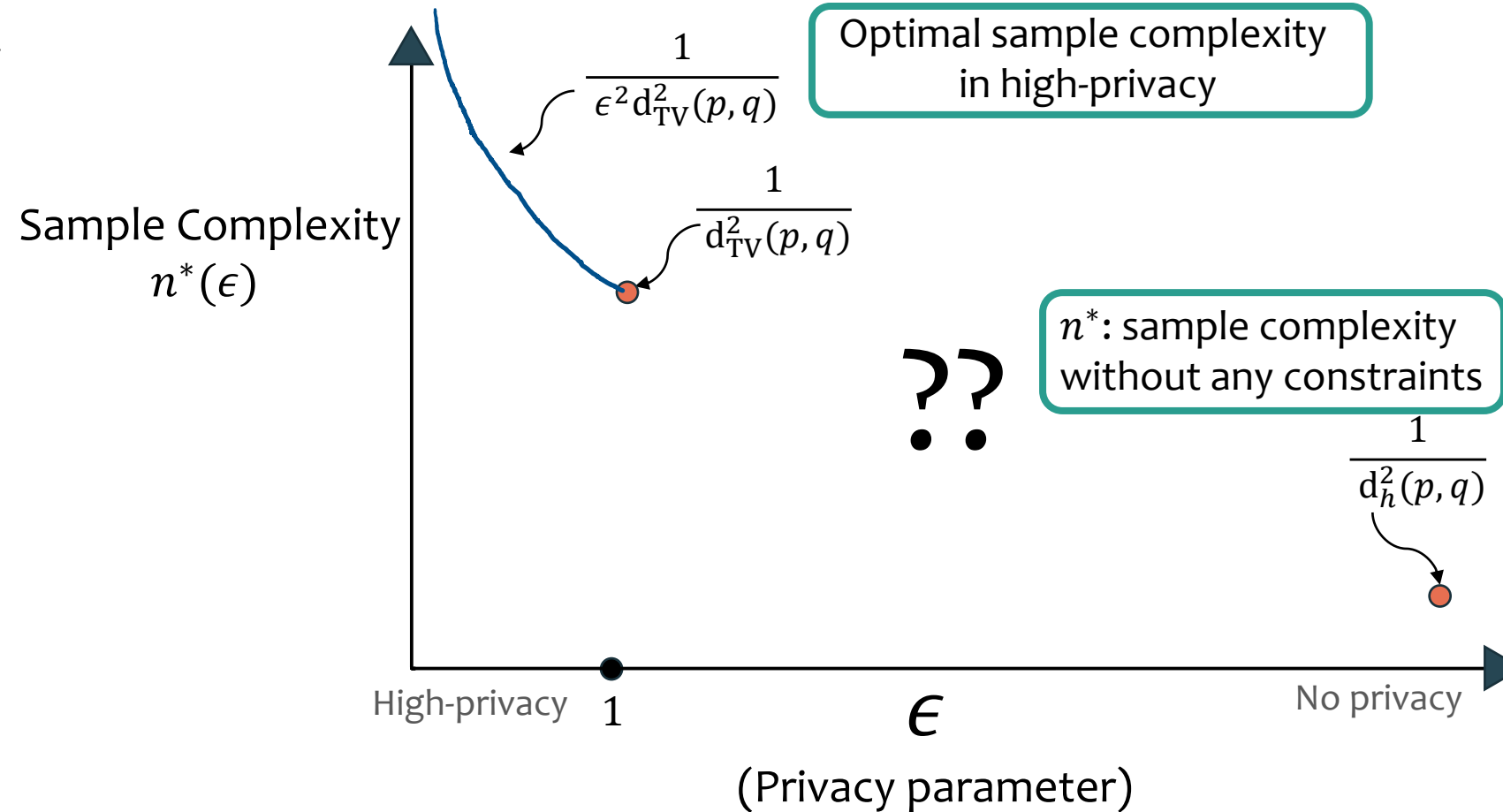
Statistical Cost of Privacy: Existing Results

- Sample Complexity



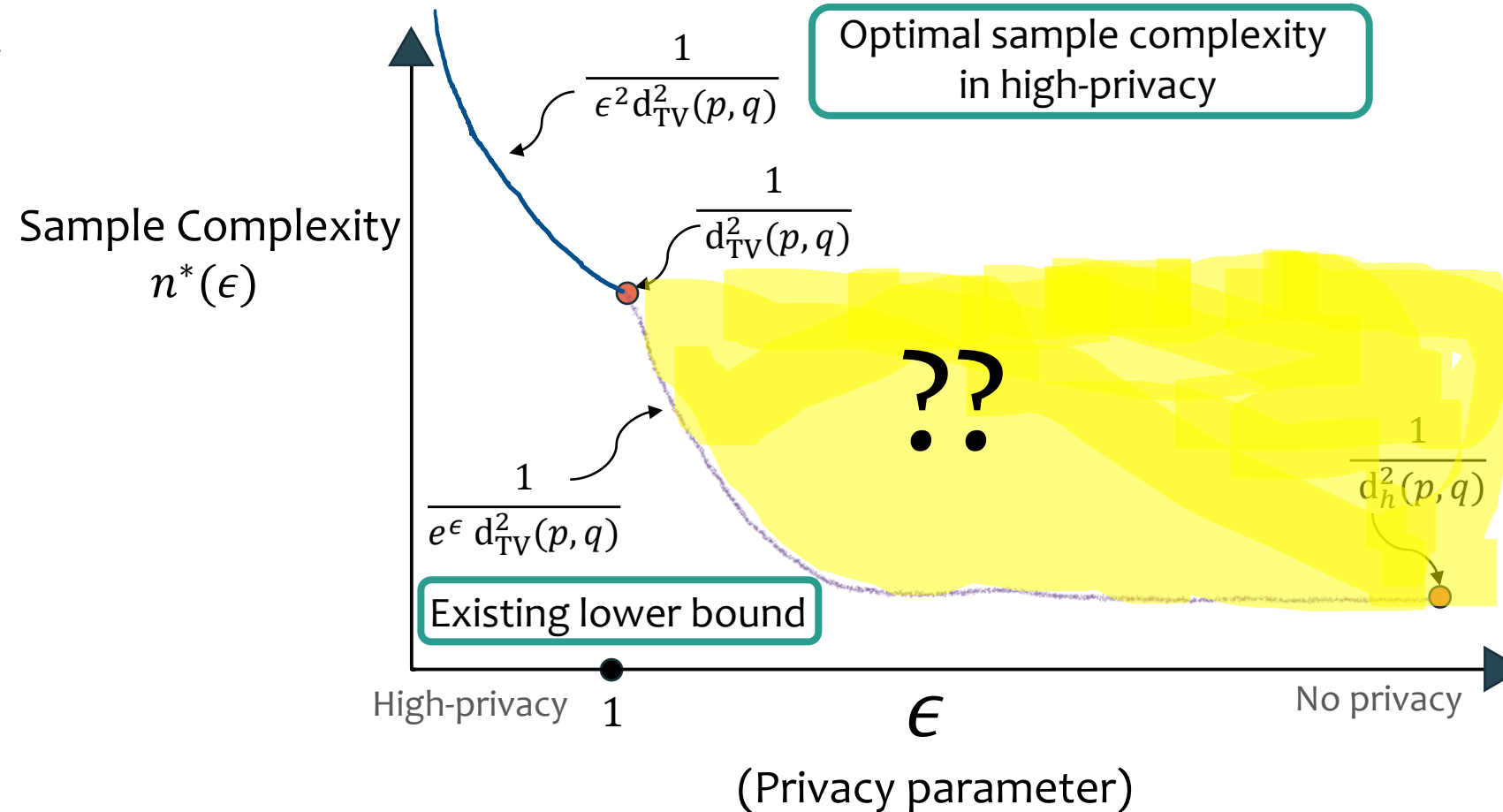
Statistical Cost of Privacy: Existing Results

- Sample Complexity



Statistical Cost of Privacy: Existing Results

- Sample Complexity



[DJW13] J. Duchi, M. Wainwright, M. Jordan. Minimax Optimal Procedures for Locally Private Estimation. 2013.

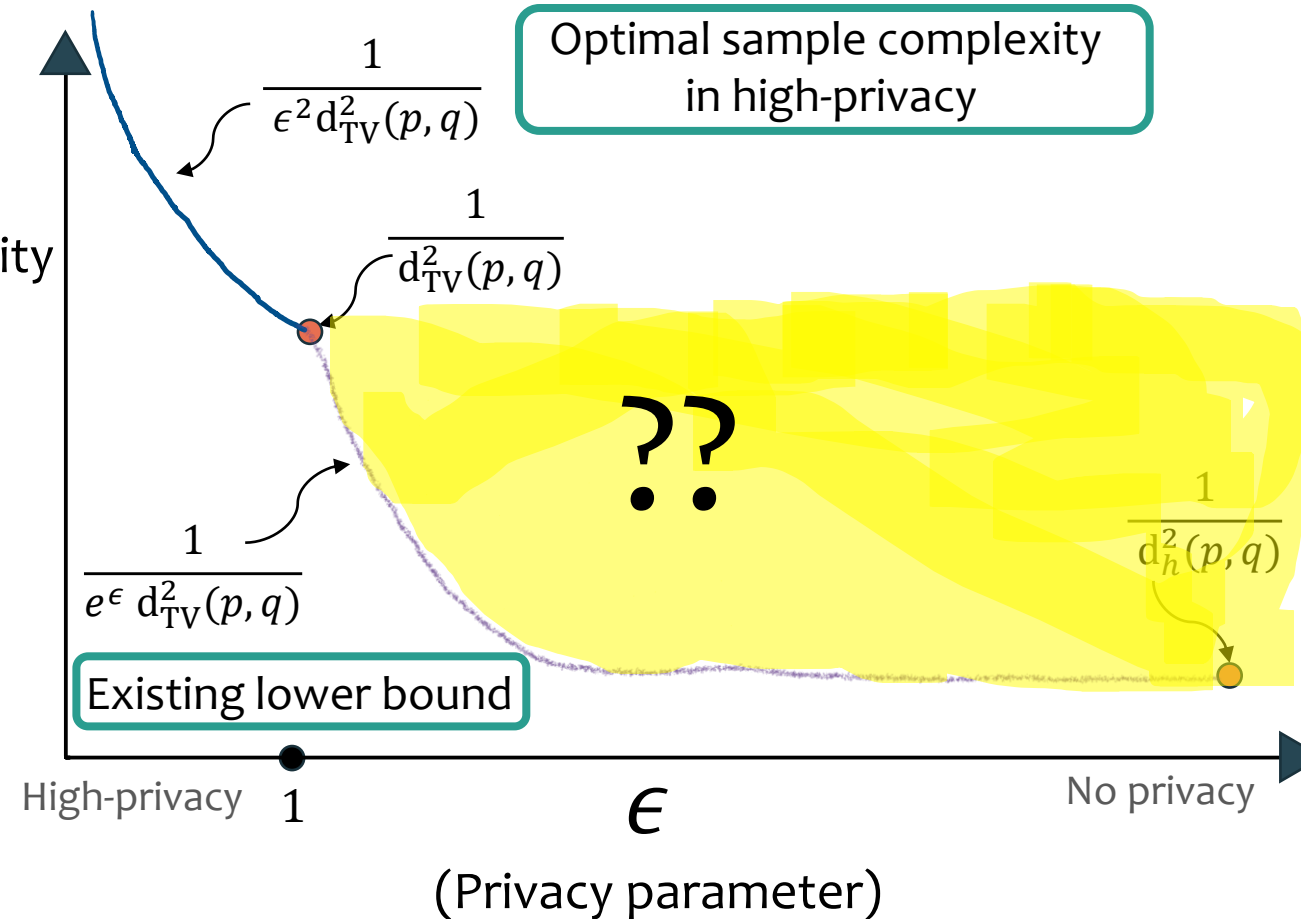
[AZ22] S. Asodeh, H. Zhang. Contraction of Locally Private Mechanisms. 2022.

Statistical Cost of Privacy: Existing Results

- Sample Complexity

Sample Complexity
 $n^*(\epsilon)$

[PAJL23]: Existing lower bound is tight for Bernoulli distributions



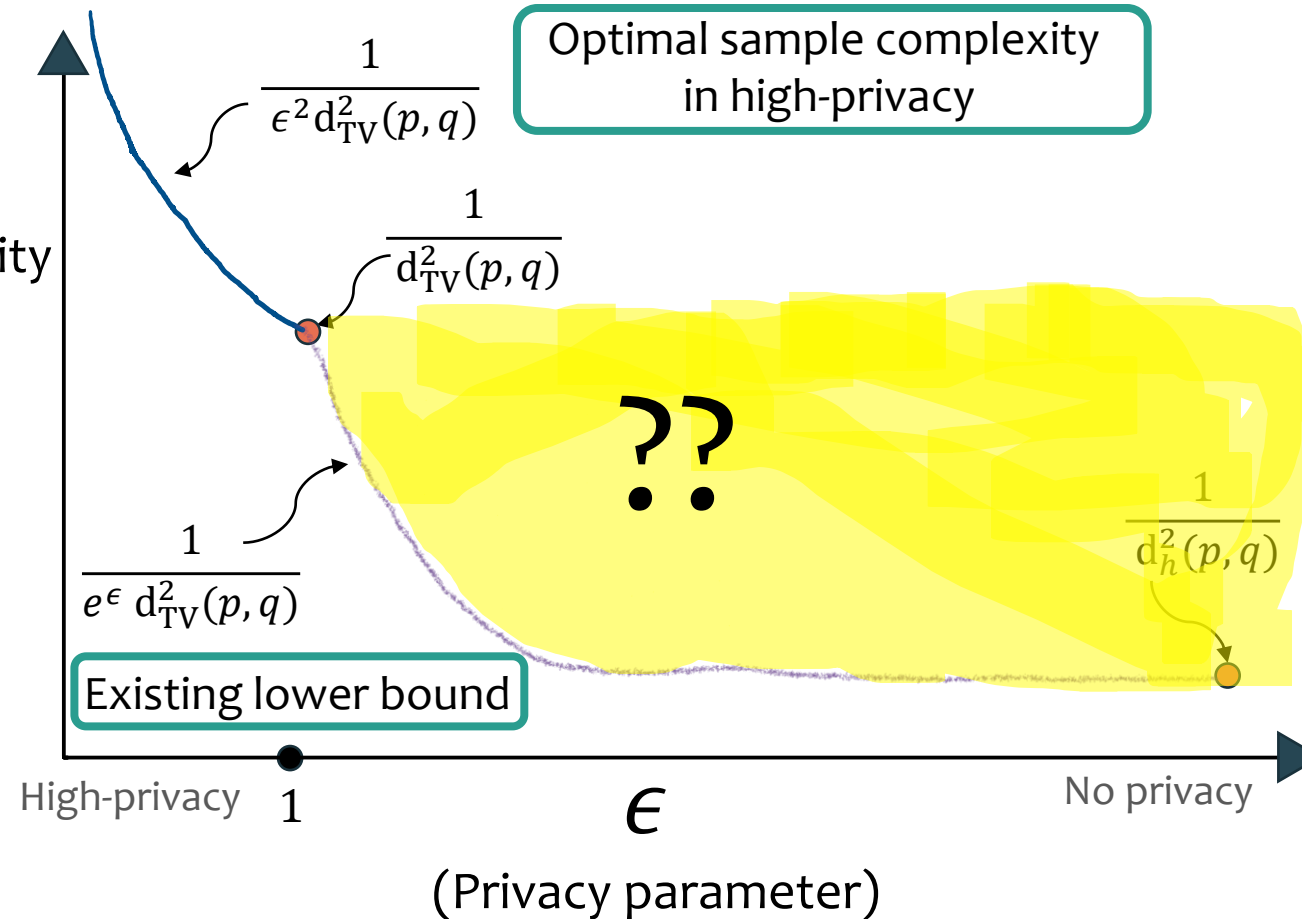
[DJW13] J. Duchi, M. Wainwright, M. Jordan. Minimax Optimal Procedures for Locally Private Estimation. 2013.

[AZ22] S. Asodeh, H. Zhang. Contraction of Locally Private Mechanisms. 2022.

Statistical Cost of Privacy: Existing Results

- Sample Complexity

Sample Complexity
 $n^*(\epsilon)$



[PAJL23]: Existing lower bound is tight for Bernoulli distributions

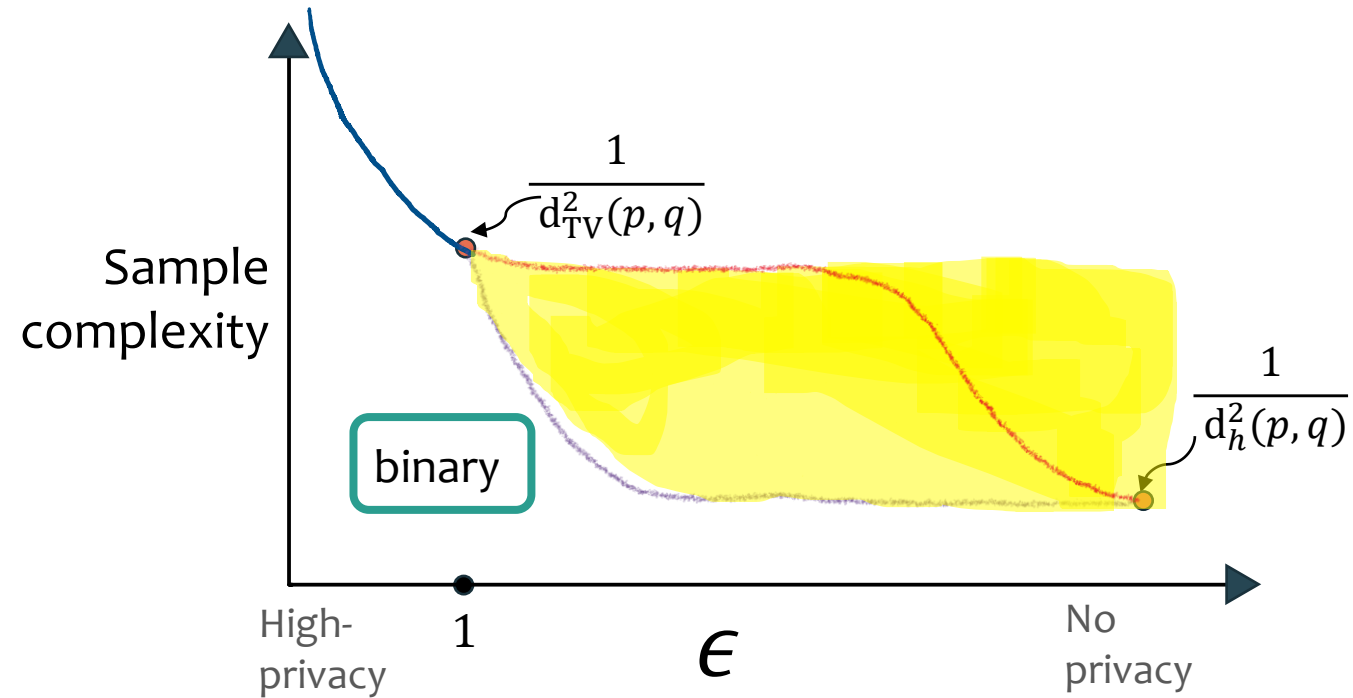
What about general distributions?

[DJW13] J. Duchi, M. Wainwright, M. Jordan. Minimax Optimal Procedures for Locally Private Estimation. 2013.

[AZ22] S. Asoodeh, H. Zhang. Contraction of Locally Private Mechanisms. 2022.

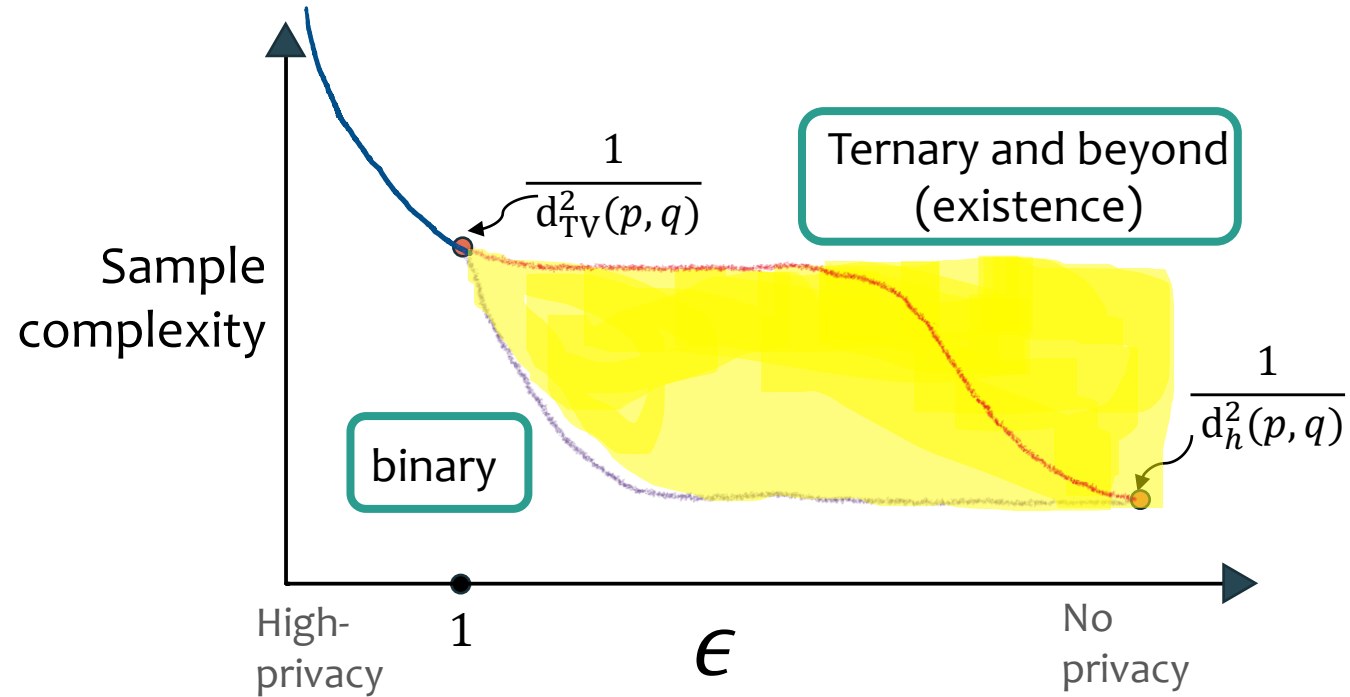
Our Results: Minimax Optimal Sample Complexity

Theorem[PAJL23] There exist ternary distributions p and q with larger sample complexities.



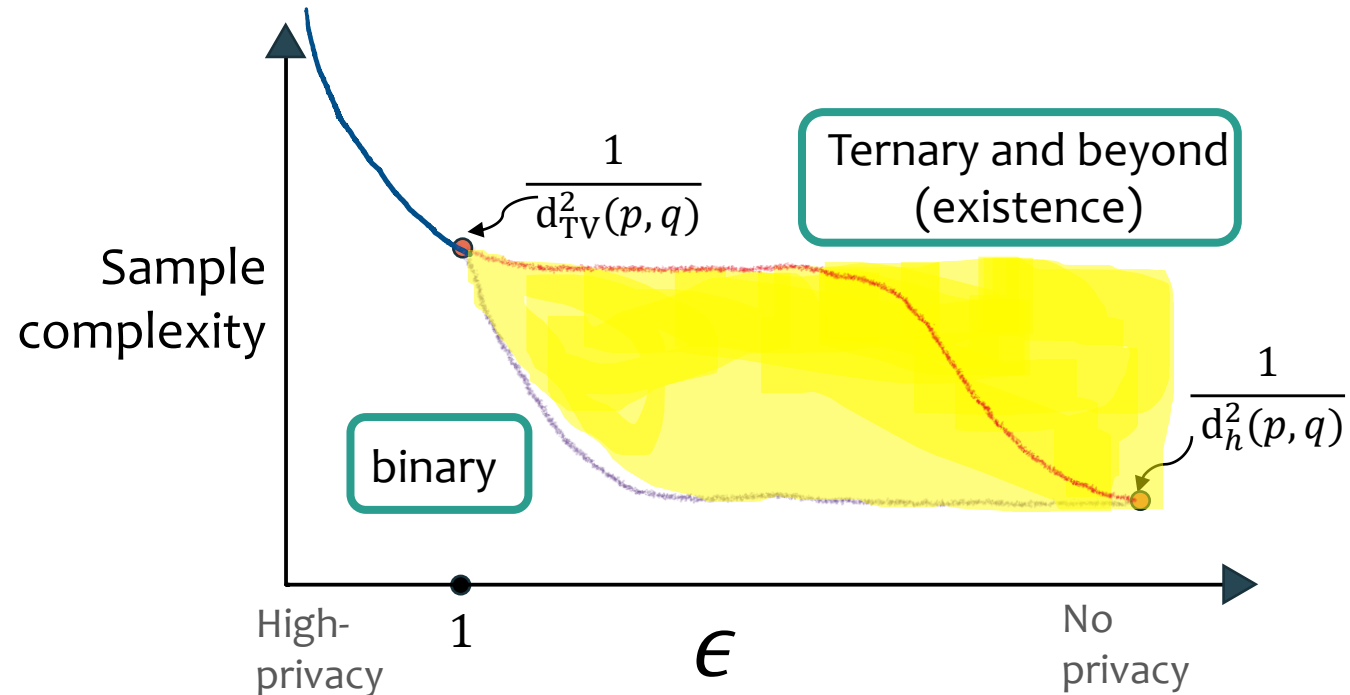
Our Results: Minimax Optimal Sample Complexity

Theorem[PAJL23] There exist ternary distributions p and q with larger sample complexities.



Our Results: Minimax Optimal Sample Complexity

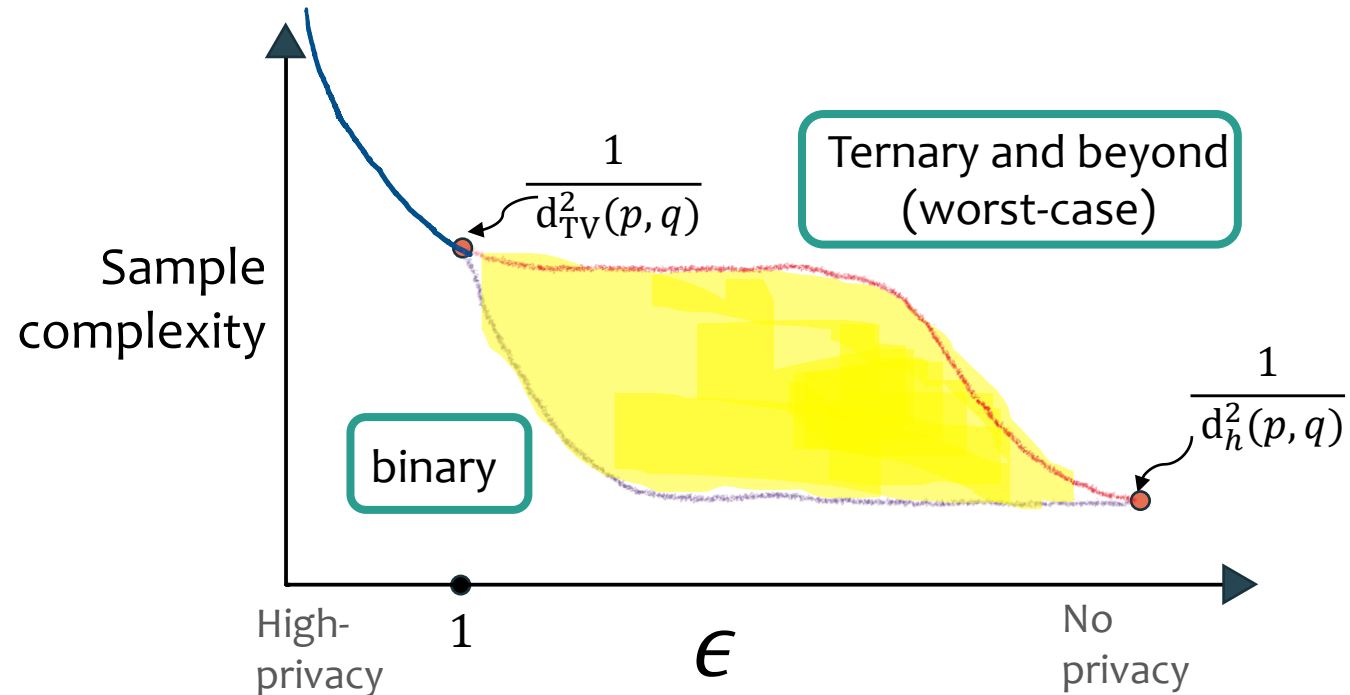
Theorem[PAJL23] There exist ternary distributions p and q with larger sample complexities.



Theorem[PAJL23] There is an efficient algorithm with nearly-matching upper bounds for all distributions.

Our Results: Minimax Optimal Sample Complexity

Theorem[PAJL23] There exist ternary distributions p and q with larger sample complexities.

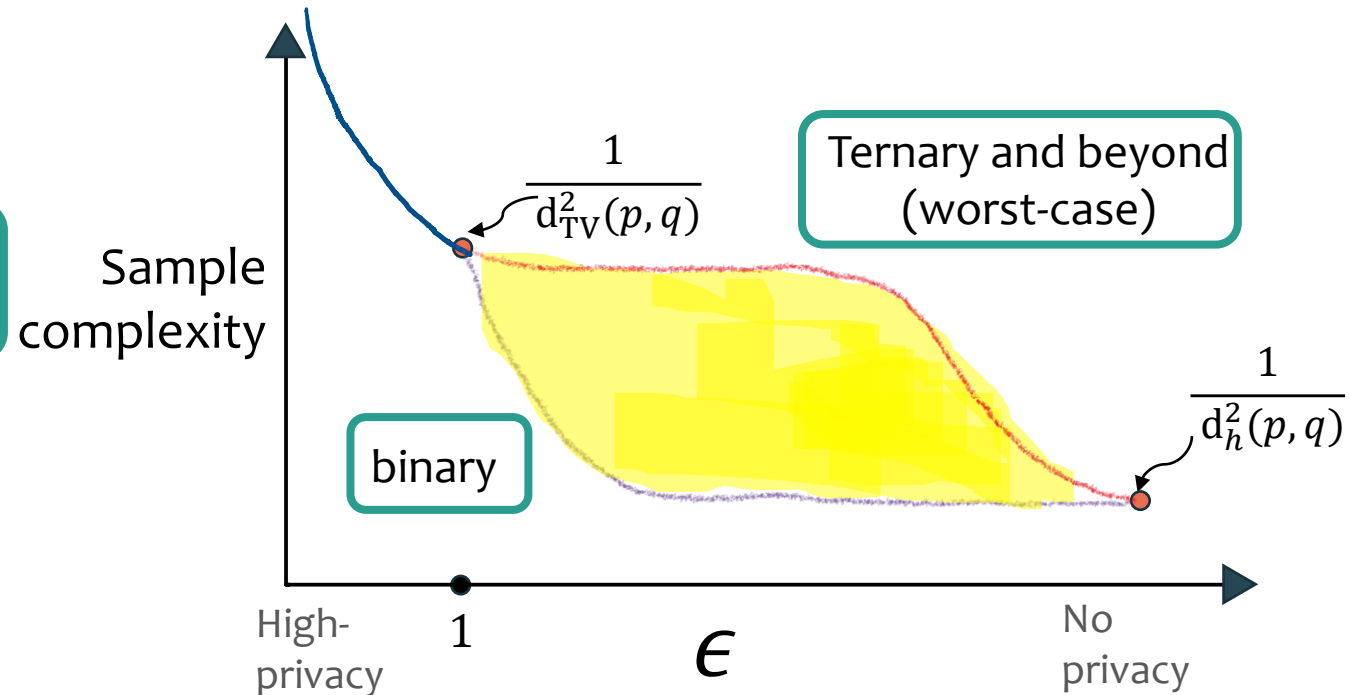


Theorem[PAJL23] There is an efficient algorithm with nearly-matching upper bounds for all distributions.

Our Results: Minimax Optimal Sample Complexity

Theorem[PAJL23] There exist ternary distributions p and q with larger sample complexities.

Overall, a satisfying story for minimax optimality



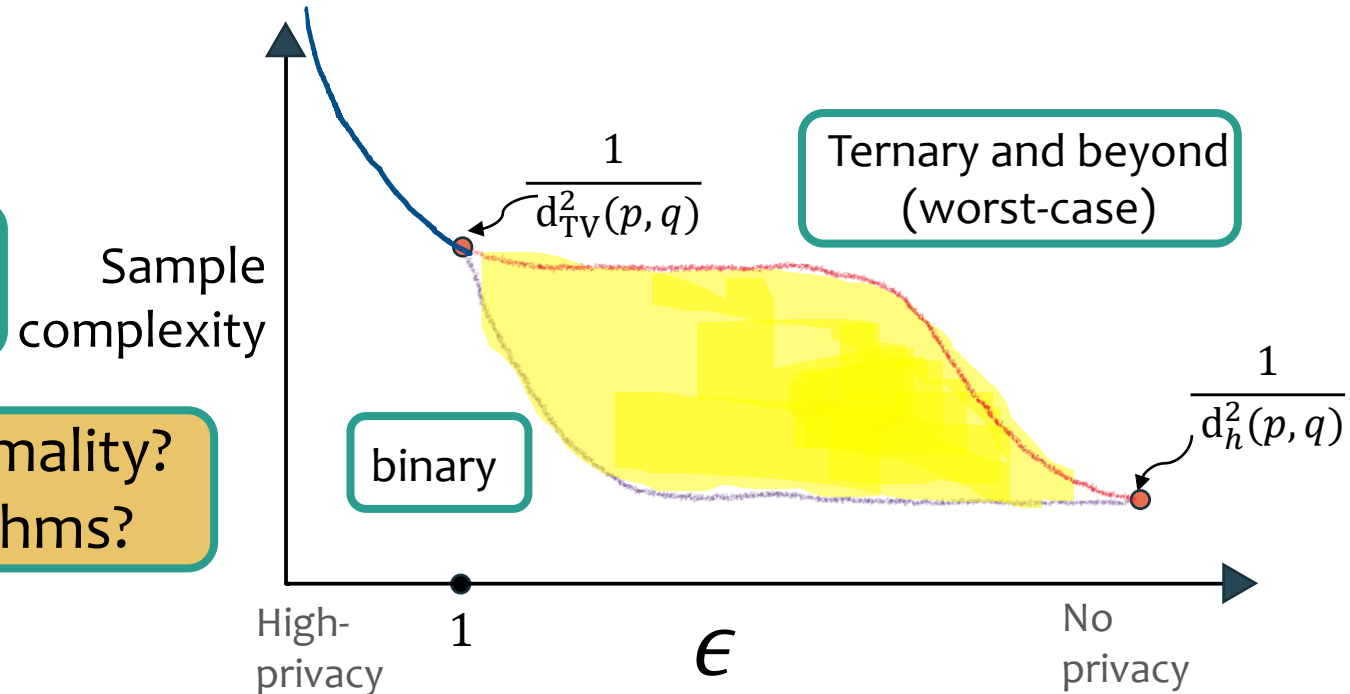
Theorem[PAJL23] There is an efficient algorithm with nearly-matching upper bounds for all distributions.

Our Results: Minimax Optimal Sample Complexity

Theorem[PAJL23] There exist ternary distributions p and q with larger sample complexities.

Overall, a satisfying story for minimax optimality

What about instance-optimality?
Are there efficient algorithms?



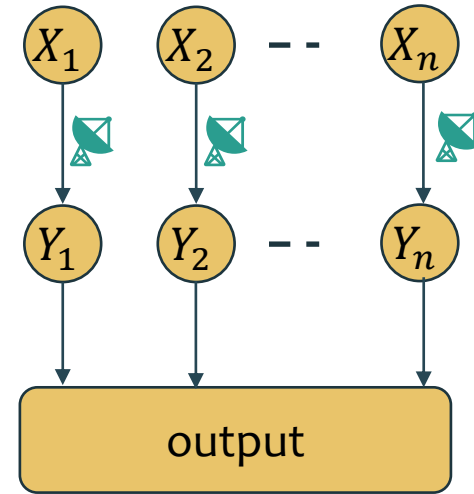
Theorem[PAJL23] There is an efficient algorithm with nearly-matching upper bounds for all distributions.

Outline


- ▶ Motivation
- ▶ Problem Statement
- ▶ Our Results
 - ▶ Statistical
 - ▶ Computational
- ▶ Proof Sketch
- ▶ Conclusion

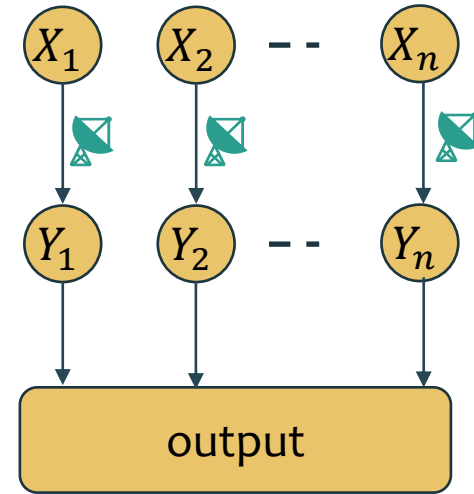
Computational Cost of Privacy

- Recall we need to map the original data $X_i \rightarrow Y_i$




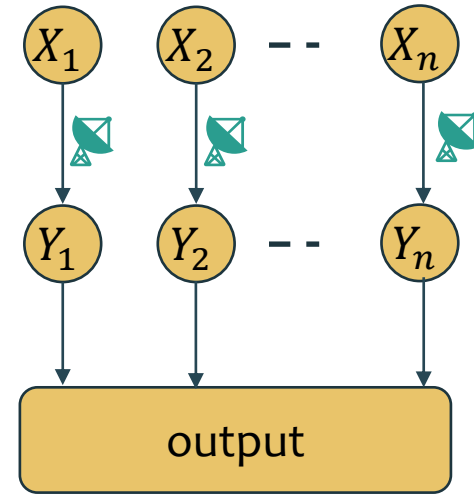
Computational Cost of Privacy

- Recall we need to map the original data $X_i \rightarrow Y_i$
- Performance depends on the channel 




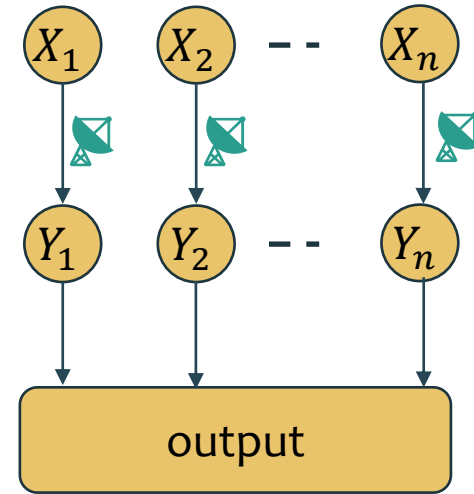
Computational Cost of Privacy

- Recall we need to map the original data $X_i \rightarrow Y_i$
- Performance depends on the channel 
 - Once the channel is fixed, perform likelihood ratio test




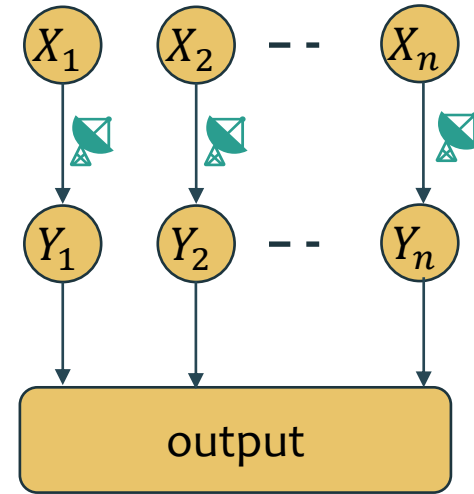
Computational Cost of Privacy

- Recall we need to map the original data $X_i \rightarrow Y_i$
- Performance depends on the channel 
 - Once the channel is fixed, perform likelihood ratio test
- Prior work on finding the optimal channel




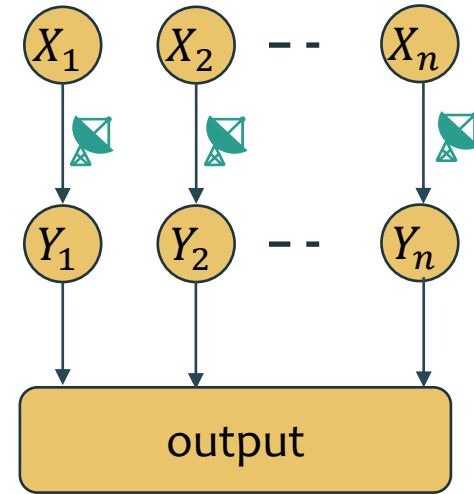
Computational Cost of Privacy

- Recall we need to map the original data $X_i \rightarrow Y_i$
- Performance depends on the channel 
 - Once the channel is fixed, perform likelihood ratio test
- Prior work on finding the optimal channel
 - $\epsilon \ll 1$: Well-understood




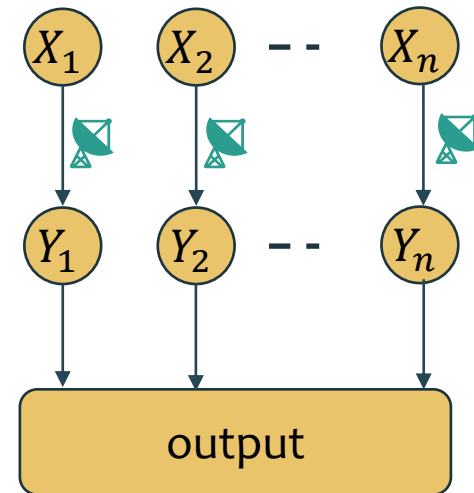
Computational Cost of Privacy

- Recall we need to map the original data $X_i \rightarrow Y_i$
- Performance depends on the channel 
 - Once the channel is fixed, perform likelihood ratio test
- Prior work on finding the optimal channel
 - $\epsilon \ll 1$: Well-understood
 - $\epsilon \gg 1$: No polynomial-time algorithm




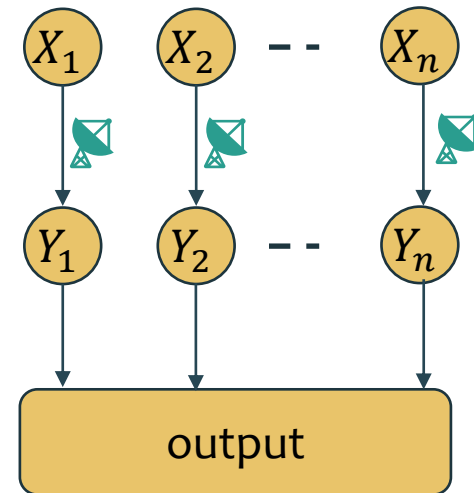
Computational Cost of Privacy

- Recall we need to map the original data $X_i \rightarrow Y_i$
- Performance depends on the channel 
 - Once the channel is fixed, perform likelihood ratio test
- Prior work on finding the optimal channel
 - $\epsilon \ll 1$: Well-understood
 - $\epsilon \gg 1$: No polynomial-time algorithm
 - [KOV14] gave an exponential-time algorithm



Computational Cost of Privacy

- Recall we need to map the original data $X_i \rightarrow Y_i$
- Performance depends on the channel 
 - Once the channel is fixed, perform likelihood ratio test
- Prior work on finding the optimal channel
 - $\epsilon \ll 1$: Well-understood
 - $\epsilon \gg 1$: No polynomial-time algorithm
 - [KOV14] gave an exponential-time algorithm




Can we efficiently find the (near)-optimal channel?


Our Results: Computational Cost of Privacy

Theorem[PAJL23] Given any two distributions p and q on $[k]$ and ϵ ,


Our Results: Computational Cost of Privacy

Theorem[PAJL23] Given any two distributions p and q on $[k]$ and ϵ , there is a **linear-time algorithm** to find an ϵ -LDP channel 

Our Results: Computational Cost of Privacy


Theorem[PAJL23] Given any two distributions p and q on $[k]$ and ϵ ,
there is a **linear-time algorithm** to find an ϵ -LDP channel 
whose sample complexity is **near-optimal**.

Our Results: Computational Cost of Privacy

Theorem[PAJL23] Given any two distributions p and q on $[k]$ and ϵ , there is a **linear-time algorithm** to find an ϵ -LDP channel  whose sample complexity is **near-optimal**.

- More broadly, consider the optimization problem

Our Results: Computational Cost of Privacy


Theorem[PAJL23] Given any two distributions p and q on $[k]$ and ϵ , there is a **linear-time algorithm** to find an ϵ -LDP channel  whose sample complexity is **near-optimal**.

- More broadly, consider the optimization problem

$$g(\text{📡 } p, \text{📡 } q)$$

g : a (quasi)-convex objective

Our Results: Computational Cost of Privacy

Theorem[PAJL23] Given any two distributions p and q on $[k]$ and ϵ , there is a **linear-time algorithm** to find an ϵ -LDP channel  whose sample complexity is **near-optimal**.


- More broadly, consider the optimization problem

$\mathcal{P}(\epsilon, \ell)$: All ϵ -LDP channels of output size ℓ

$$\max_{\mathcal{P}(\epsilon, \ell)} g(\mathcal{P}, p, q)$$

g : a (quasi)-convex objective

Our Results: Computational Cost of Privacy

Theorem[PAJL23] Given any two distributions p and q on $[k]$ and ϵ , there is a **linear-time algorithm** to find an ϵ -LDP channel  whose sample complexity is **near-optimal**.

- More broadly, consider the optimization problem


$\mathcal{P}(\epsilon, \ell)$: All ϵ -LDP channels of output size ℓ

$$\max_{\mathcal{P}(\epsilon, \ell)} g(\text{satellite } p, \text{satellite } q)$$

g : a (quasi)-convex objective

Recall: maximizing a convex objective is usually hard!

Our Results: Computational Cost of Privacy

Theorem[PAJL23] Given any two distributions p and q on $[k]$ and ϵ , there is a **linear-time algorithm** to find an ϵ -LDP channel  whose sample complexity is **near-optimal**.

- More broadly, consider the optimization problem

$\mathcal{P}(\epsilon, \ell)$: All ϵ -LDP channels of output size ℓ

$$\max_{\mathcal{P}(\epsilon, \ell)} g(\text{satellite } p, \text{satellite } q)$$

g : a (quasi)-convex objective

Recall: maximizing a convex objective is usually hard!

Theorem[PAJL23] There is a $\text{poly}(k^{\ell^2})$ -time algorithm to find the optimum.

Outline

- ▶ Motivation
- ▶ Problem Statement
- ▶ Our Results
 - ▶ Statistical
 - ▶ Computational
- ▶ **Proof Sketch**
- ▶ Conclusion

Proof Sketch: Exponential Search to Linear

- Say, we want to find the optimal binary channel \mathbf{T}^* $\max_{\mathbf{T} \in \mathcal{P}(\epsilon, 2)} g(\mathbf{T}p, \mathbf{T}q)$

Proof Sketch: Exponential Search to Linear

- Say, we want to find the optimal binary channel \mathbf{T}^*
- Can show that optimal \mathbf{T}^* is of the form:

$$\max_{\mathbf{T} \in \mathcal{P}(\epsilon, 2)} g(\mathbf{T}p, \mathbf{T}q)$$

Proof Sketch: Exponential Search to Linear

- Say, we want to find the optimal binary channel \mathbf{T}^* $\max_{\mathbf{T} \in \mathcal{P}(\epsilon, 2)} g(\mathbf{T}_p, \mathbf{T}_q)$
- Can show that optimal \mathbf{T}^* is of the form:
 - First, use a binary deterministic channel \mathbf{T}' to partition $[k]$ into two sets

Proof Sketch: Exponential Search to Linear

- Say, we want to find the optimal binary channel \mathbf{T}^* $\max_{\mathbf{T} \in \mathcal{P}(\epsilon, 2)} g(\mathbf{T}p, \mathbf{T}q)$
- Can show that optimal \mathbf{T}^* is of the form:
 - First, use a binary deterministic channel \mathbf{T}' to partition $[k]$ into two sets
 - Ensure privacy using the randomized response channel (BSC)

Proof Sketch: Exponential Search to Linear

- Say, we want to find the optimal binary channel \mathbf{T}^* $\max_{\mathbf{T} \in \mathcal{P}(\epsilon, 2)} g(\mathbf{T}_p, \mathbf{T}_q)$
- Can show that optimal \mathbf{T}^* is of the form:
 - First, use a binary deterministic channel \mathbf{T}' to partition $[k]$ into two sets
 - Ensure privacy using the randomized response channel (BSC)
- But the number of possible partitions: 2^k

Proof Sketch: Exponential Search to Linear

- Say, we want to find the optimal binary channel \mathbf{T}^* $\max_{\mathbf{T} \in \mathcal{P}(\epsilon, 2)} g(\mathbf{T}p, \mathbf{T}q)$
- Can show that optimal \mathbf{T}^* is of the form:
 - First, use a binary deterministic channel \mathbf{T}' to partition $[k]$ into two sets
 - Ensure privacy using the randomized response channel (BSC)
- But the number of possible partitions: 2^k
- Can we use p and q to reduce our search space?

Proof Sketch: Exponential Search to Linear

- Say, we want to find the optimal binary channel \mathbf{T}^* $\max_{\mathbf{T} \in \mathcal{P}(\epsilon, 2)} g(\mathbf{T}p, \mathbf{T}q)$
- Can show that optimal \mathbf{T}^* is of the form:
 - First, use a binary deterministic channel \mathbf{T}' to partition $[k]$ into two sets
 - Ensure privacy using the randomized response channel (BSC)
- But the number of possible partitions: 2^k
- Can we use p and q to reduce our search space?
- Our answer: yes!

Proof Sketch: Exponential Search to Linear

- Say, we want to find the optimal binary channel \mathbf{T}^* $\max_{\mathbf{T} \in \mathcal{P}(\epsilon, 2)} g(\mathbf{T}p, \mathbf{T}q)$
- Can show that optimal \mathbf{T}^* is of the form:
 - First, use a binary deterministic channel \mathbf{T}' to partition $[k]$ into two sets
 - Ensure privacy using the randomized response channel (BSC)
- But the number of possible partitions: 2^k
- Can we use p and q to reduce our search space?
- Our answer: yes!
 - Optimal partition must respect the likelihood ratios of p and q

Outline

- ▶ Motivation
- ▶ Problem Statement
- ▶ Our Results
 - ▶ Statistical
 - ▶ Computational
- ▶ Proof Sketch
- ▶ Conclusion

Conclusion and Future Directions

- Derived minmax-optimal sample complexities under privacy
 - No longer depends only on TV distance and Hellinger

Conclusion and Future Directions

- Derived minmax-optimal sample complexities under privacy
 - No longer depends only on TV distance and Hellinger
- Computationally and Communication-efficient algorithms

Conclusion and Future Directions

- Derived minmax-optimal sample complexities under privacy
 - No longer depends only on TV distance and Hellinger
- Computationally and Communication-efficient algorithms
- Open problems:

Conclusion and Future Directions

- Derived minmax-optimal sample complexities under privacy
 - No longer depends only on TV distance and Hellinger
- Computationally and Communication-efficient algorithms
- Open problems:
 - Role of interactivity

Conclusion and Future Directions

- Derived minmax-optimal sample complexities under privacy
 - No longer depends only on TV distance and Hellinger
- Computationally and Communication-efficient algorithms
- Open problems:
 - Role of interactivity
 - Characterization of instance-optimal sample complexity
 - Looking beyond TV distance and Hellinger divergence

Conclusion and Future Directions

- Derived minmax-optimal sample complexities under privacy
 - No longer depends only on TV distance and Hellinger
- Computationally and Communication-efficient algorithms
- Open problems:
 - Role of interactivity
 - Characterization of instance-optimal sample complexity
 - Looking beyond TV distance and Hellinger divergence
 - M-ary hypothesis testing, optimally

Conclusion and Future Directions

- Derived minmax-optimal sample complexities under privacy
 - No longer depends only on TV distance and Hellinger
- Computationally and Communication-efficient algorithms
- Open problems:
 - Role of interactivity
 - Characterization of instance-optimal sample complexity
 - Looking beyond TV distance and Hellinger divergence
 - M-ary hypothesis testing, optimally

Thank you!