

The Role of Pseudorandomness in (Computational) Differential Privacy

Ludmila Glinskih

Examination committee:

Mark Bun

Steven Homer

Sofya Raskhodnikova

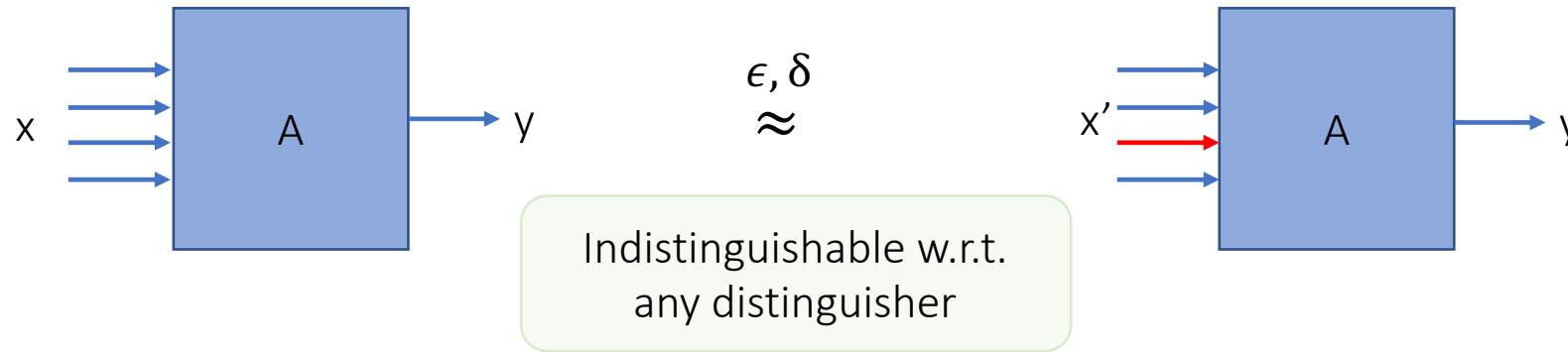
Boston University,

December 20, 2021

Plan of the talk

1. Overview of the papers
 1. Computational Differential Privacy
 1. Dense Model Theorem
 2. SV-sources for DP algorithms
 3. PRGs for Local DP algorithms
2. Proof of the equivalence result
 1. Simple Direction: hybrid argument
 2. Hard direction: dense model theorem
3. Possible extensions and future work

Differential Privacy

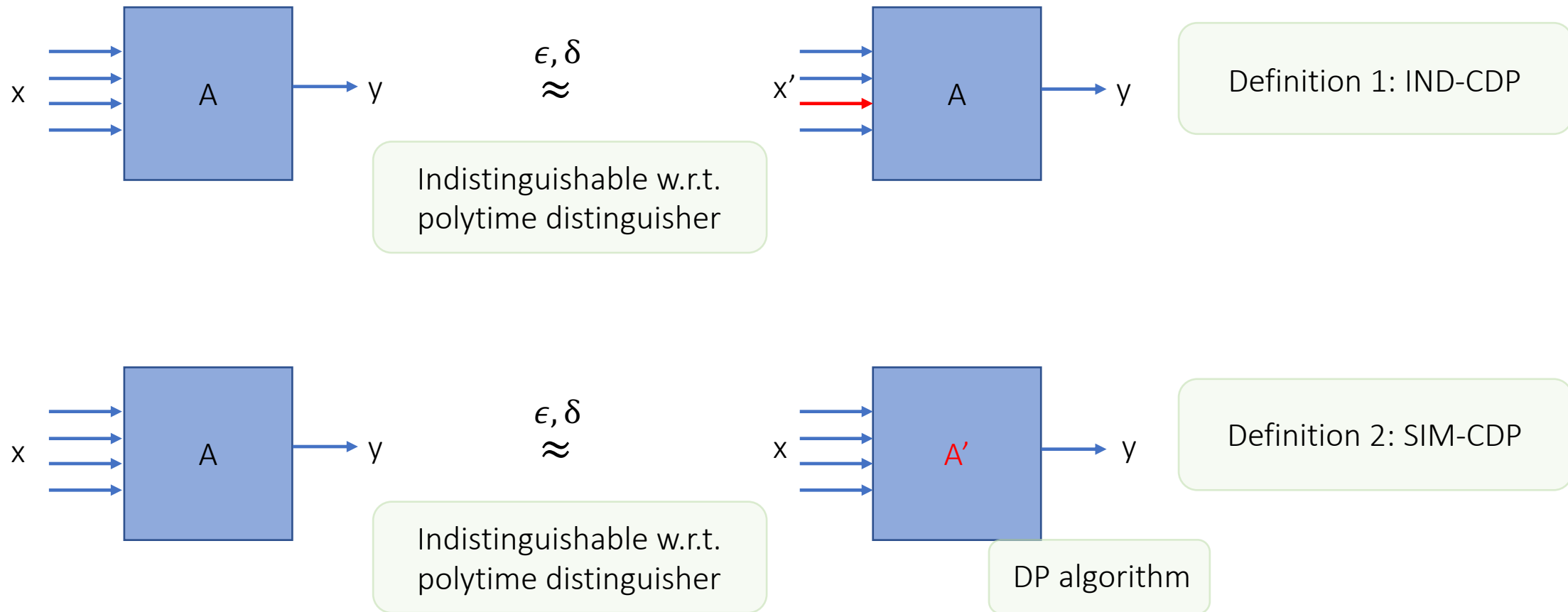


A is (ϵ, δ) -differentially private if for every set of possible outputs O , and for every neighboring x, x' :

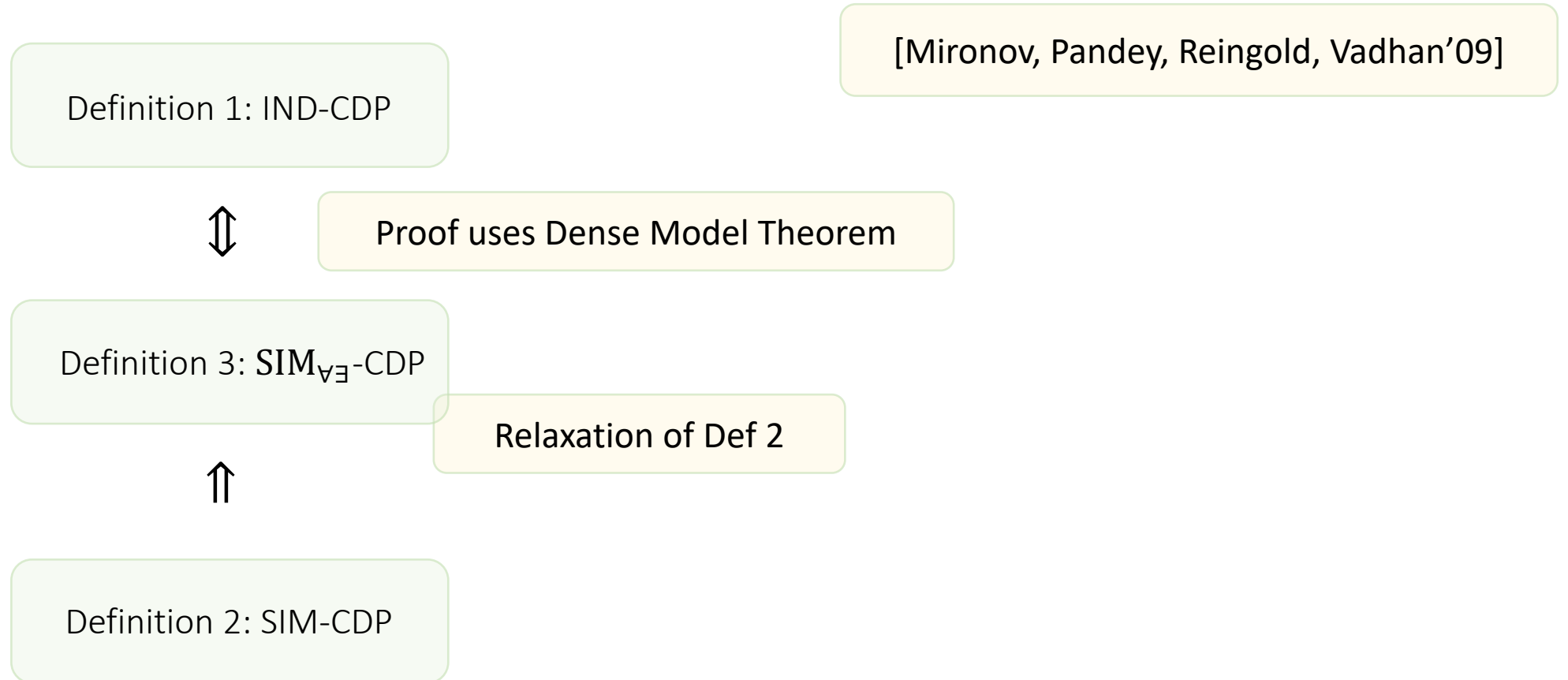
$$P[A(x) \in O] \leq e^\epsilon \cdot P[A(x') \in O] + \delta$$

Computational Differential Privacy

If a real adversary is a polysize probabilistic circuit, can we relax the definition of Differential Privacy?



Main Theorem of this talk



Differential Privacy and Randomness Sources

All non-trivial DP algorithms should be randomized

- Standard algorithms sample from the uniform distribution



Can we use Santha-Vazirani random sources instead of the Uniform distribution?

- Every i -th bit in the gamma-SV sequence has bias gamma
- Obstacles:
 - SV sources are non-extractable
 - Cannot construct signatures and other basic “privacy” protocols out of it

Differential Privacy with SV sources of Randomness

Can we build DP protocols that uses SV sources of randomness?

[Dodis, Lopez-Alt, Mironov, Vadhan'09]

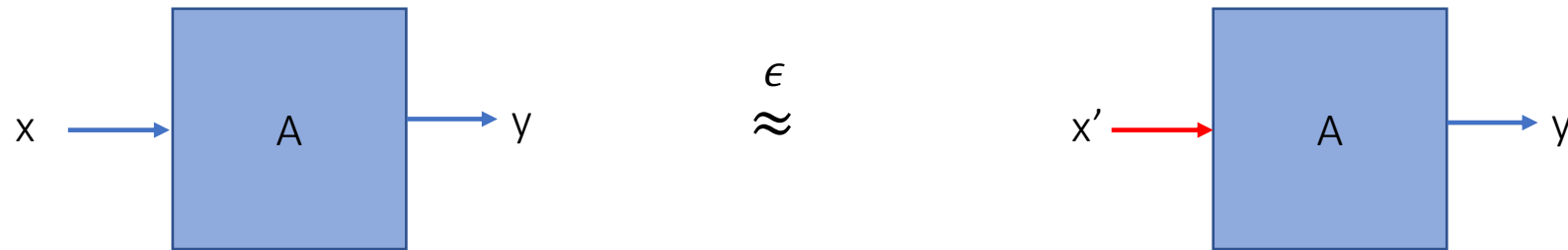


- No, if we use additive noise DP algorithms
 - $A(x) = f(x) + \text{random_noise}$, where f is a true, non-DP answer
- Yes, if we use additive noise DP algorithms with the discretization
 - The main intuition is that sets of probabilities of getting the same output values should be almost the same for all neighboring datasets

Local Differential Privacy

Users may not trust a centralized database,

- Then, they can add noise to their data guarantee their privacy



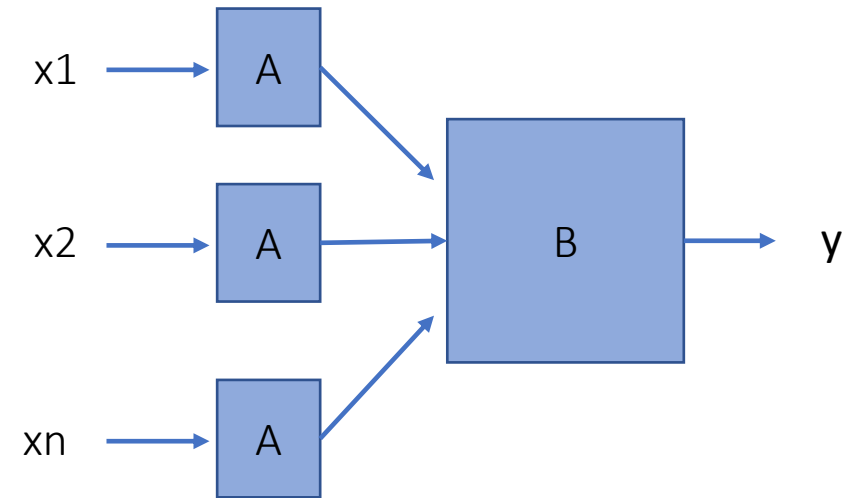
For every user, for every pair of values x, x' that this user may have, for every possible output o :

$$P[A(x) = o] \leq e^\epsilon \cdot P[A(x') = o]$$

Local DP: communication overhead

Each user adds noise, encodes, sends their data to the centralized server

- Many users (millions in case of Google, Apple, Facebook) send their data
- Amount of information in each message is small due to noise
- If encoding is not efficient, then communication overhead is huge



What if result of each algorithm A could be encoded using an output of a PRG, so we would need to send only a seed of such PRG?

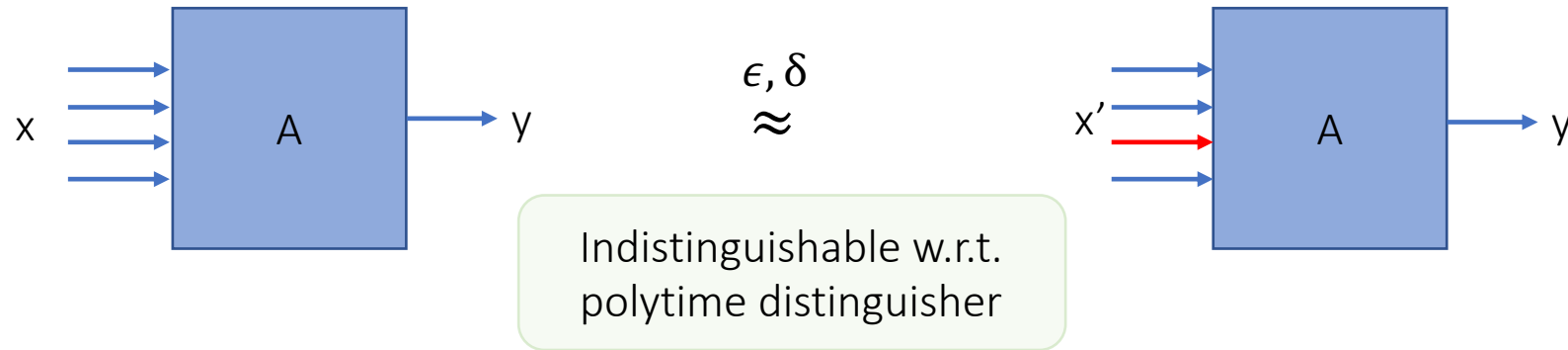
Can significantly decrease communication

- For element counting
- For mean estimation

[Feldman, Talwar'21]

Proof of $\text{SIM}_{\forall \exists}$ -CDP \Leftrightarrow IND-CDP

CDP: two definitions



Definition $\text{SIM}_{\forall\exists}$ -CDP:

Distribution on outputs of A is indistinguishable by polysized circuits from a distribution on outputs of a family of DP algorithm

Definition IND-CDP:

Distribution on outputs of A on different inputs is indistinguishable

Equivalence result

Theorem:

A mechanism B is IND-CDP if and only if B is $\text{SIM}_{\forall\exists}$ -CDP

Proof idea:

$\text{SIM}_{\forall\exists}$ -CDP \Rightarrow IND-CDP

Hybrid argument

$\text{SIM}_{\forall\exists}$ -CDP \Leftarrow IND-CDP

Dense Model Theorem

This talk

Dense, Pseudodense, and Indist Distributions

Distribution X is e^ϵ -dense in Y if

- X, Y are distribution over the same set R ,
- $\forall x \in R \Pr[X = x] \leq e^\epsilon \cdot \Pr[Y = x]$.

Distribution X is δ -indistinguishable from Y w.r.t. a family of predicates $\mathbf{A}: R \rightarrow \{0,1\}$ if

- $\forall A \in \mathbf{A} \quad | \Pr[A(X) = 1] - \Pr[A(Y) = 1] | \leq \delta$.

Distribution X is (e^ϵ, δ) -pseudodense in Y w.r.t. a family of predicates $\mathbf{A}: R \rightarrow \{0,1\}$ if

- $\forall A \in \mathbf{A} \quad \Pr[A(X) = 1] \leq e^\epsilon \cdot \Pr[A(Y) = 1] + \delta$.

DP Definitions via Pseudodensity

$f: \mathcal{D} \rightarrow R$ is ϵ -DP if and only if for all neighboring $D, D' \in \mathcal{D}$

- $f(D)$ is e^ϵ -dense in $f(D')$.

$\{f_k\}: \mathcal{D} \rightarrow R$ is ϵ_k -IND-CDP if and only if for all

- Exists $s(k) = k^{\omega(1)}$
- For all neighboring $D, D' \in \mathcal{D}$ of size $\leq s(k)$ $f_k(D)$ is $(e^\epsilon, \frac{1}{s(k)})$ -pseudodense in $f_k(D')$
 - w.r.t. circuits of size $\leq s(k)$

$\{f_k\}: \mathcal{D} \rightarrow R$ is ϵ_k -SIM _{$\forall\exists$} -CDP if and only if for all sequences of neighboring inputs $\{(D_k, D_k')\}$

- Exists $s(k) = k^{\omega(1)}$
- Exists $\{F_k\}$, such that every F_k is ϵ_k -DP
- $f_k(D)$ is $\frac{1}{s(k)}$ -indistinguishable from $F_k(D)$.
 - w.r.t. circuits of size $\leq s(k)$

Dense Model Theorem

Theorem:

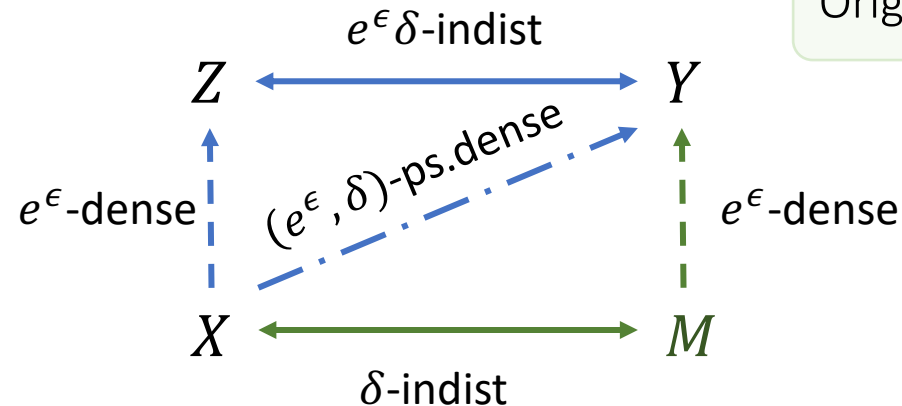
X, Y , and Z are distributions over finite set R

- X is e^ϵ -dense in Z
- $A = \{a_i\}$ is a set of predicates $a_i: R \rightarrow \{0,1\}$

If all distributions that are e^ϵ -dense in Y are δ -distinguishable from X ,

Then Z is $\Omega(e^\epsilon \delta)$ -distinguishable from Y .

[Reingold, Trevisan, Tulsiani, Vadhan'08]



Original proof for $Y = U_X$

New Dense Model Theorem

Theorem:

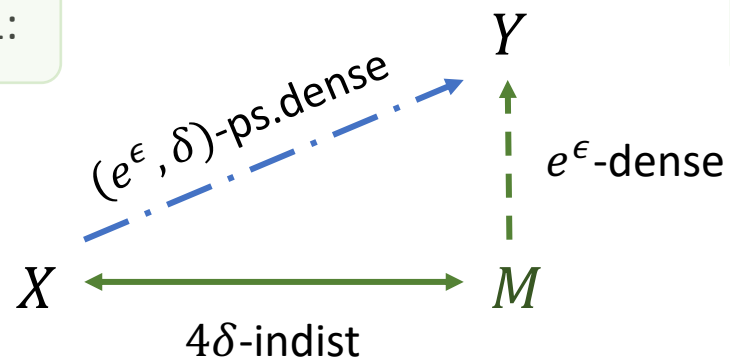
[Mironov, Pandey, Reingold, Vadhan'09]

X, Y are distributions over finite set R

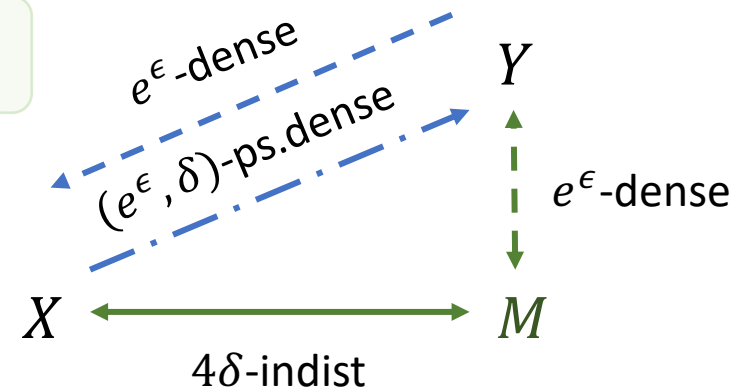
- X is (e^ϵ, δ) -pseudodense in Y
- $G(A)$ is a family of threshold predicates

1. Then there exists a distribution M that is e^ϵ -dense in Y and 4δ -distinguishable from X with respect to family of predicates A ;
2. If Y is e^ϵ -dense in X , then Y is e^ϵ -dense in M .

Statement 1:

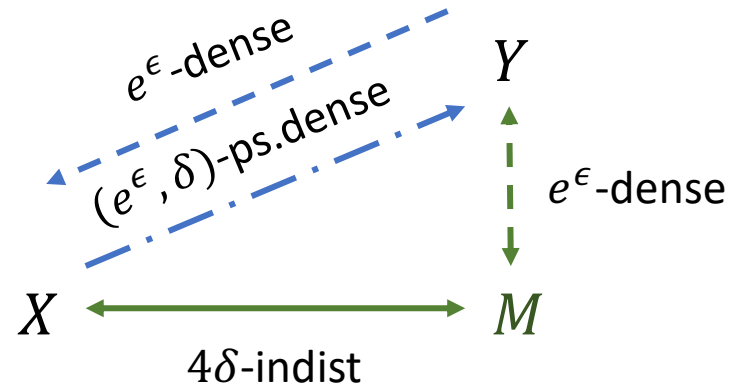


Statement 2:



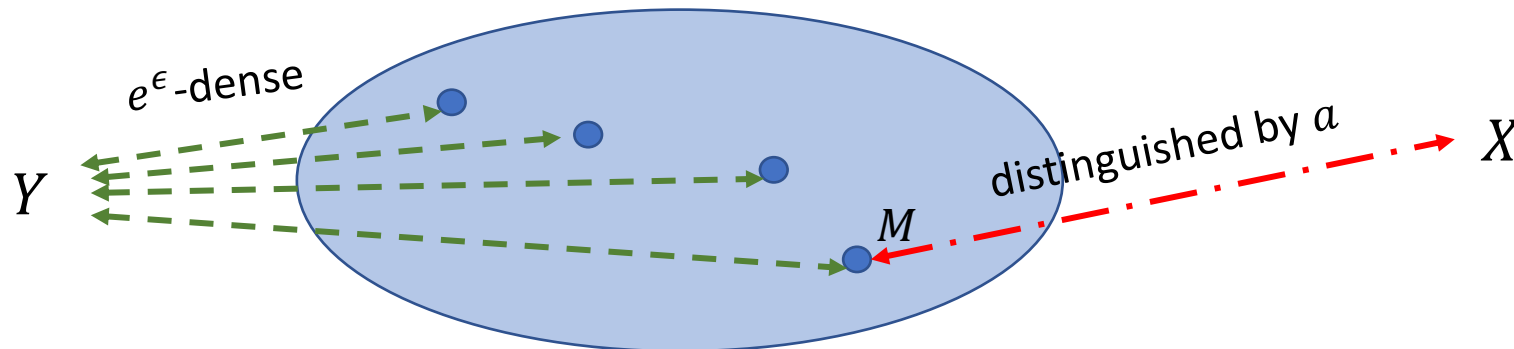
Proof of New Dense Model Theorem

Theorem:



Proof by contradiction:

Assume for any M that is e^ϵ -dense in Y exists $a \in \mathbf{A}$ that distinguishes M from X w.p. $\geq \mu = 4\delta$.



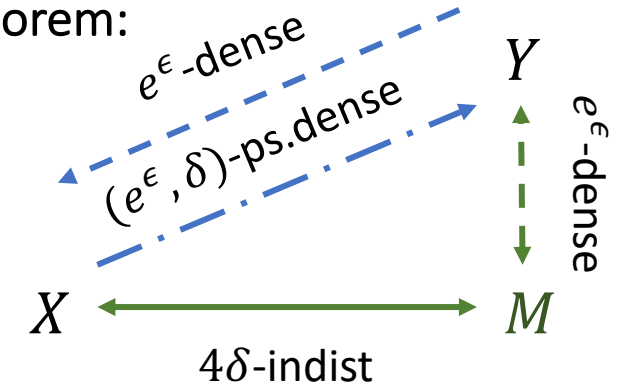
Proof of NDMT

Set of all such M is convex

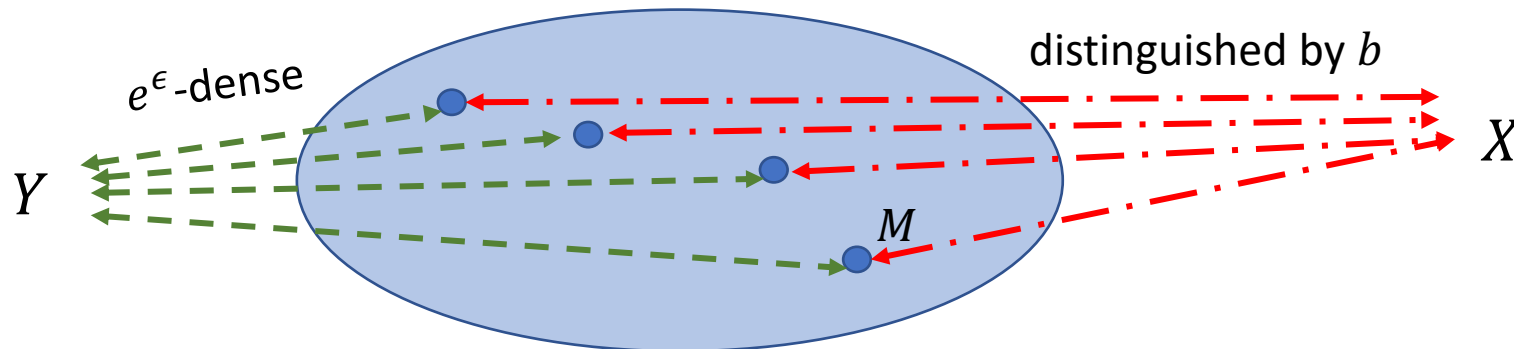
Idea: apply MinMax principle to get a new predicate b that distinguishes every M from X w.p. $\geq \mu = 4\delta$

- Player A picks a predicate a
- Player B picks a distribution D
- The payoff is $E[a(D)] - E[a(X)]$
- \exists mixed strategy $b: \forall M E[b(M)] - E[b(X)] \geq \gamma$
- \exists mixed strategy $M': \forall a E[a(M')] - E[a(X)] \leq \gamma$

Theorem:



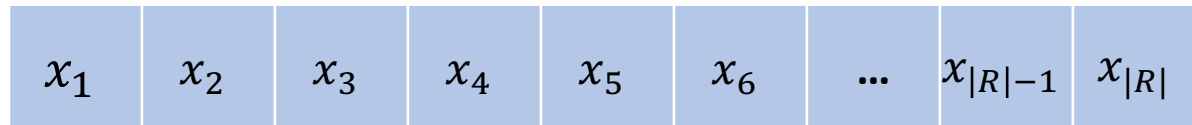
Distribution over dense distributions is a dense-distribution $\Rightarrow \gamma \geq \mu$



Proof of NDMT

Consider $b: \forall M \mathbb{E}[b(M)] - \mathbb{E}[b(X)] \geq \mu$

Sort all elements $x \in R$ in decreasing order of $\Pr[b(x) = 1]$



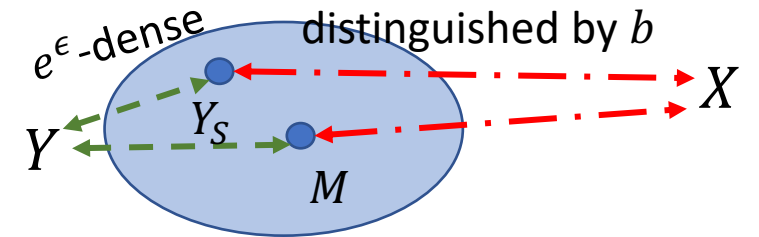
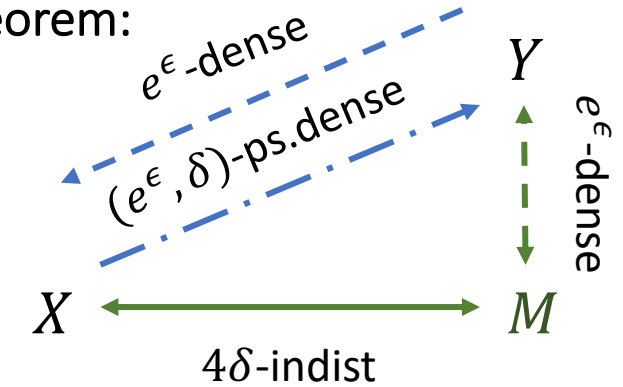
$$S: \Pr[Y \in S] = \frac{1}{1 + e^\epsilon}$$

Define a new distribution Y_S :

$$\Pr[Y_S = y] = \Pr[Y = y] \cdot \begin{cases} e^\epsilon, & \text{if } y \in S \\ e^{-\epsilon}, & \text{otherwise} \end{cases}$$

Y_S and Y are e^ϵ -dense $\Rightarrow b$ distinguishes Y_S from X

Theorem:



Proof of NDMT

Lemma:

$F: X \rightarrow [0,1]$

- Z and W are distributions, such that $E[F(Z)] \geq E[F(W)] + \mu$

Then exists $t \in [\frac{\mu}{2}, 1]$, such that

- $\Pr[F(Z) > t] \geq \Pr[F(W) \geq t - \frac{\mu}{2}] + \frac{\mu}{2}$

Applying this lemma for $F(x) = \Pr[b(x) = 1]$, X , and Y_S we get:

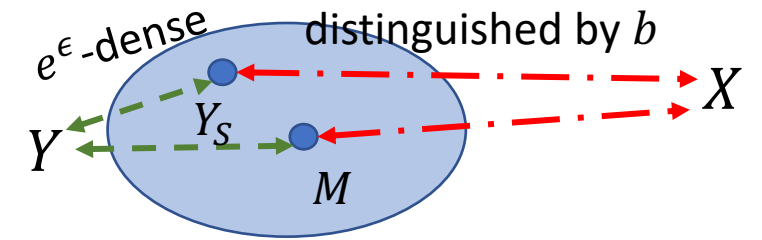
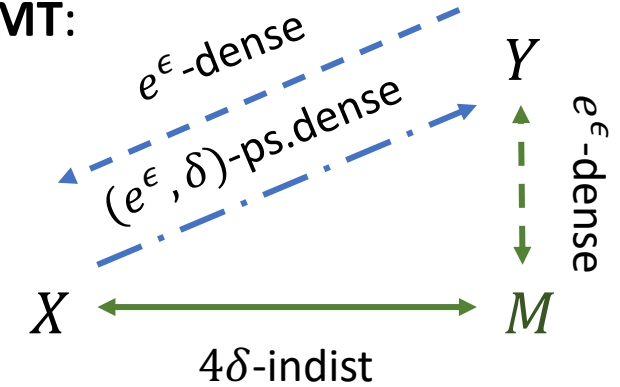
- $\Pr\left[\underbrace{\Pr[b(X) = 1]}_{\geq t + \frac{\mu}{2}} > t + \frac{\mu}{2}\right] \geq \Pr\left[\underbrace{\Pr[b(Y_S) = 1]}_{\geq t} \geq t\right] + \frac{\mu}{2}$

If $\geq t + \frac{\mu}{2}$: $b'(x) = 1$
 “looks like X ”

If $< t$: $b'(x) = 0$
 “looks like Y_S ”

new classifier, based on
 the output probability of b

NDMT:



Proof of NDMT

$$* \Pr \left[\underbrace{\Pr[b(X) = 1] > t + \frac{\mu}{2}}_{\text{If } \geq t + \frac{\mu}{2} : b'(x) = 1 \text{ "looks like X"}} \right] \geq \Pr \left[\underbrace{\Pr[b(Y_S) = 1] \geq t}_{\text{If } < t : b'(x) = 0 \text{ "looks like } Y_S"}} \right] + \frac{\mu}{2}$$

If $\geq t + \frac{\mu}{2} : b'(x) = 1$
"looks like X"

If $< t : b'(x) = 0$
"looks like Y_S "

Claim: $b'(y) = 0 \forall y \notin S$

Proof by contradiction:

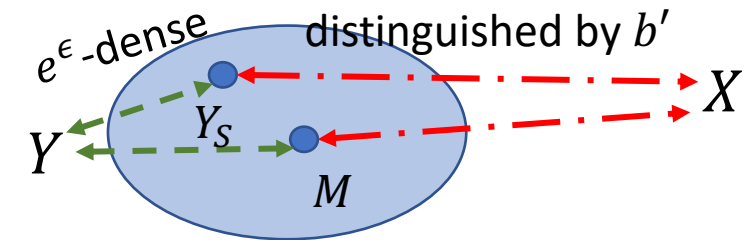
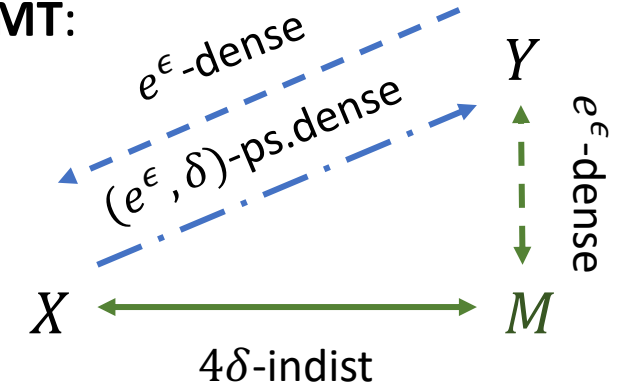
By construction $b'(y) \neq 0 \forall y \in S$

Y is dense in $X \Rightarrow$ for all $y \notin S$

$$\bullet \Pr[Y_S = y] = e^{-\epsilon} \cdot \Pr[Y = y] \leq e^{-\epsilon} e^{\epsilon} \cdot \Pr[X = y] = \Pr[X = y]$$

$$\Pr[b'(Y_S) = 0] = \sum_{\substack{y \notin S \\ b(y)=0}} \Pr[Y_S = y] \leq \sum_{\substack{y \notin S \\ b(y)=0}} \Pr[X = y] = \Pr[b'(X) = 0]$$

NDMT:



Contradiction with *:

$$\begin{aligned} \Pr[b'(X) = 0] &\leq 1 - \Pr[b'(X) = 1] \\ &\leq 1 - \Pr[b'(Y) \neq 0] - \frac{\mu}{2} = \Pr[b'(Y) = 0] - \frac{\mu}{2} \end{aligned}$$

Final part of the proof of NDMT

Using the fact that $b'(y) = 0 \ \forall y \notin S$:

- $\Pr[b'(Y) \neq 0] = e^{-\epsilon} \cdot \Pr[b'(Y_S) \neq 0] < e^{-\epsilon} (\Pr[b'(X) = 1] - \frac{\mu}{2})$
- $\Pr[b'(X) = 1] > e^{\epsilon} \cdot \Pr[b'(Y) \neq 0] + \frac{\mu}{2}$

This could be a contradiction with a pseudosensitivity of X in Y if b' was from the original family of predicates

Lemma:

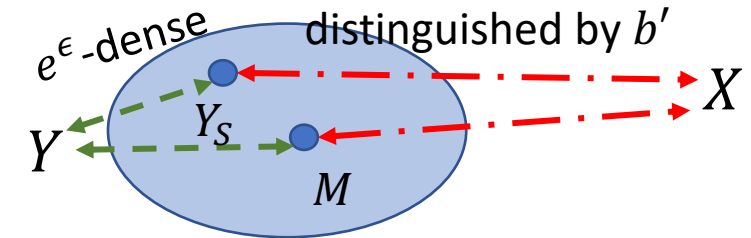
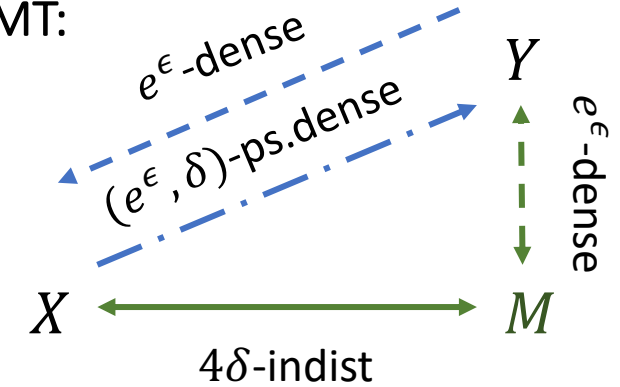
$F: R \rightarrow [0,1]$ is a convex combination of bounded functions from G

- Z_1, Z_2 distributions of G
- $\alpha, \beta > 0$

Then for $k = O(\frac{1}{\alpha^2} \cdot \log \frac{1}{\beta})$ there are $f_1, f_2, \dots, f_k \in G$ such that

- $\Pr \left[\left| F(Z_i) - \frac{1}{k} \cdot (f_1(Z_i) + f_2(Z_i) + \dots + f_k(Z_i)) \right| > \alpha \right] \leq \beta, \text{ for } i = 1, 2$

NDMT:



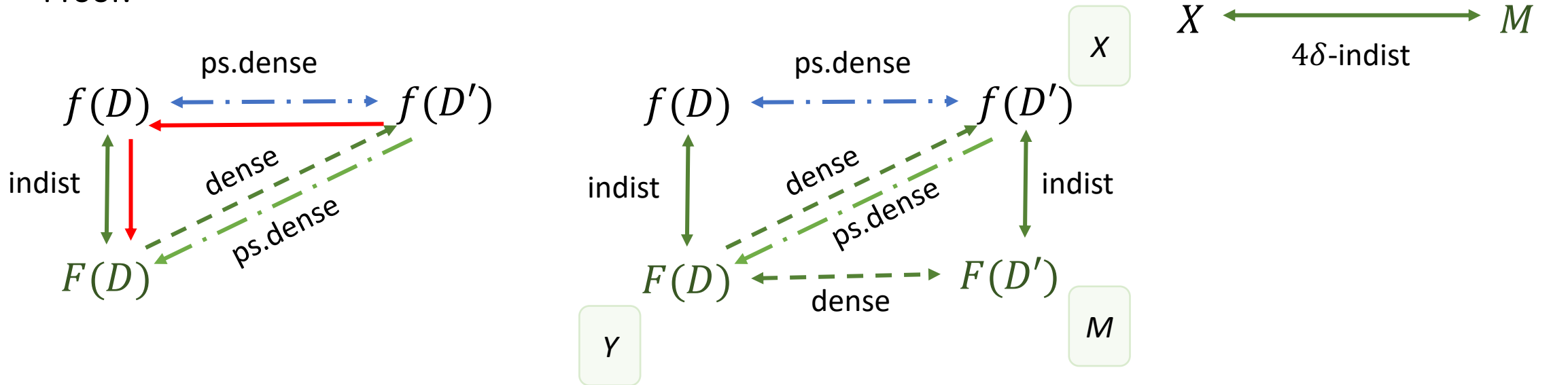
IND-CDP \Rightarrow $\text{SIM}_{\forall\exists}$ -CDP using NDMT

Theorem:

If a family $\{f_k\}: D \rightarrow R_k$ is ϵ_k -IND-CDP,

Then it is also ϵ_k - $\text{SIM}_{\forall\exists}$ -CDP.

Proof:



Conclusion

- Pseudorandomness can help to construct new DP algorithms
 - Can use PRGs to reduce communication complexity for Local-DP algorithms
 - Can use SV-sources with imperfect randomness to construct DP algorithms
 - Can use PRGs and computational indistinguishability to construct various Computationally DP algorithms
- Technical result that we discussed today:
 - Proof of the New Dense Model Theorem
 - Its application to showing that $\text{IND-CDP} \Rightarrow \text{SIM}_{\forall\exists}\text{-CDP}$

Open Problems

- The separation between $\text{SIM}_{\forall \exists}$ -CDP from SIM-CDP remains an open question
- How else, could we apply the new Dense Model Theorem?
- What other random sources (with imperfect randomness) could we use to build DP algorithms?
- Can we unconditionally compress communication in other Local-DP algorithms?

