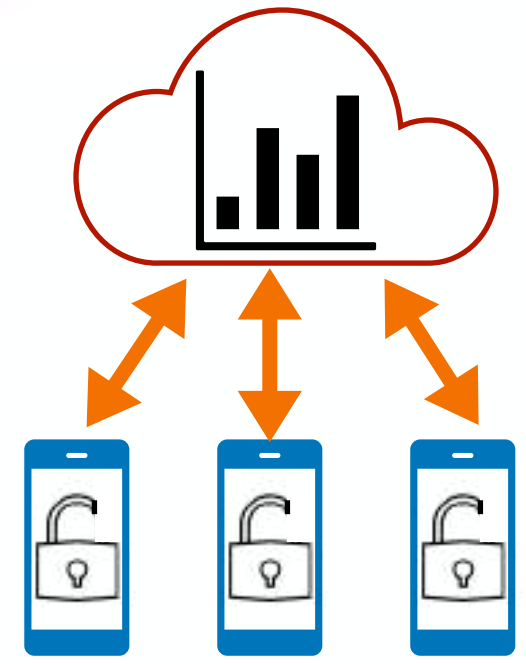


Data Acquisition Problem

- Advances in AI → demand for data is increasing!
- Companies collect data in various ways:
 - offering payments, e.g., Nielsen.
 - collecting it as a byproduct of their services, e.g., Netflix, Google, or Facebook.



- A common concern: **Privacy!**
- As a user's data is harnessed, more & more information about her behavior & preferences are uncovered.
- Differential Privacy (DP)** [Dwork et al. 2006]: a popular framework for characterizing privacy losses:
 - Widely used by tech companies, including Apple, Google, and Microsoft.



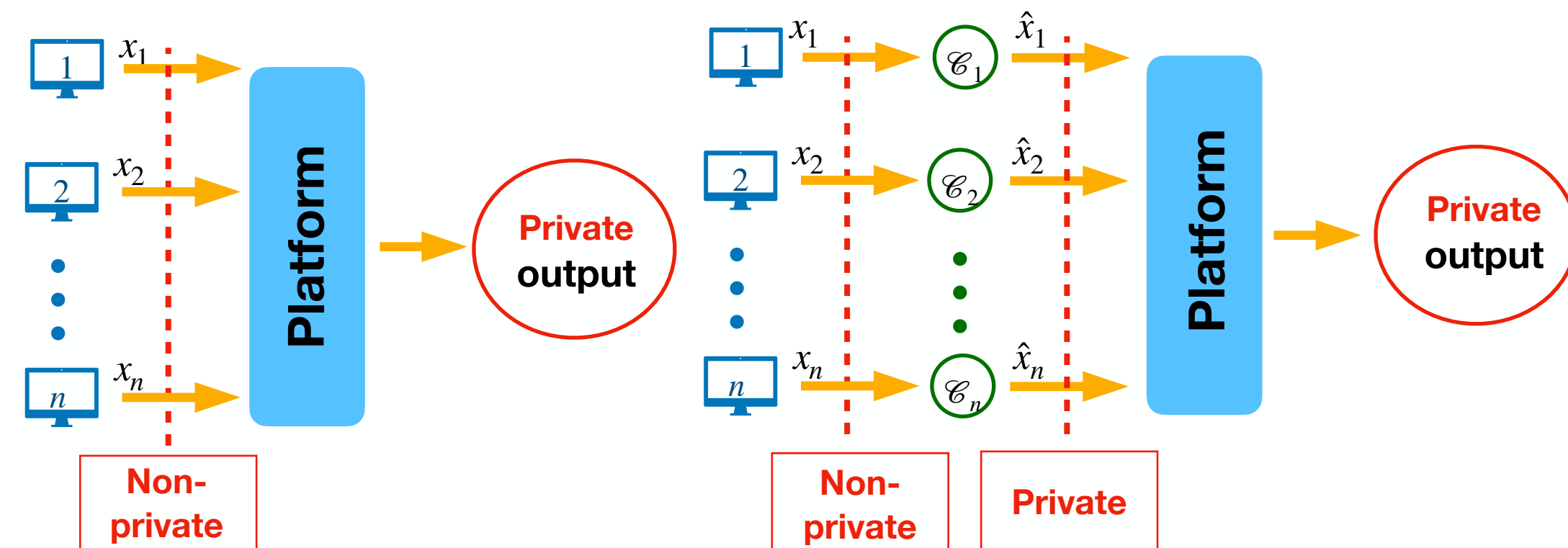
How much privacy an individual will obtain?

- Two individuals have similar data, but one is more privacy conscious.
- We want to provide different privacy levels to them.
- How should we do this? How should we elicit their privacy sensitivities?

- Model:**
- A platform is interested in estimating a parameter θ by collecting data from users $1, \dots, n$.
- Data of user i : $x_i = \theta + Z_i \in \mathcal{X}$.
- Z_i 's: i.i.d. zero mean with variance VAR and $|Z_i| \leq 1/2$.

Central & Local Differential Privacy

- In the **central setting**, users trust the platform and the platform releases a private estimate.
- In the **local setting**, users make their data private from the beginning: user i shares $\mathcal{G}_i(x_i)$ with the platform:



Central Differential Privacy Local Differential Privacy

Central & Local Differential Privacy

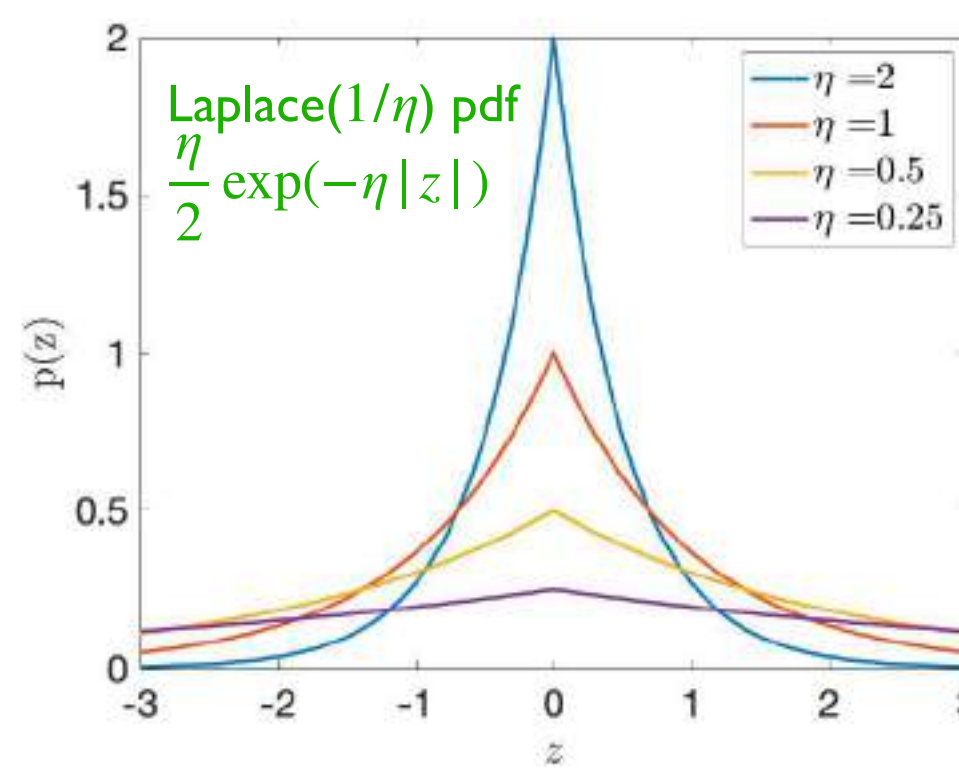
- An algorithm \mathcal{A} is $(\epsilon_i)_{i=1}^n$ -**centrally DP** if:

$$\mathbb{P}\{\mathcal{A}(\mathcal{S}) \in \mathcal{W}\} \leq e^{\epsilon_i} \mathbb{P}\{\mathcal{A}(\mathcal{S}') \in \mathcal{W}\}$$
- ϵ_i : maximum **privacy loss level** of user i .
- Smaller ϵ_i = user i 's data has less impact on the output of Algorithm → higher privacy guarantee for user i .

- An algorithm \mathcal{A} is $(\epsilon_i)_{i=1}^n$ -**locally DP** if:

$$\mathbb{P}\{\mathcal{G}_i(x) \in \mathcal{W}\} \leq e^{\epsilon_i} \mathbb{P}\{\mathcal{G}_i(x') \in \mathcal{W}\}, \quad \forall x, x' \in \mathcal{X}$$

- Laplace mechanism:**
 - Adding Laplace noise to ensure privacy.
 - Lemma:** Linear estimator $\sum_{i=1}^n w_i x_i + \text{Laplace}(1/\eta)$ is $(w_i \eta)_{i=1}^n$ centrally DP.



Question

- Assume $(\epsilon_i)_{i=1}^n$ are given.
- Which estimator is **optimal** w.r.t. mean square error?

Minimax error $\mathcal{L}(\Theta, \mathcal{P}) := \inf_{\hat{\theta} \in \Theta} \sup_{P \in \mathcal{P}} \mathbb{E}_{(x_i)_{i=1}^n \sim P, \hat{\theta}} \left[\left| \hat{\theta}(x_{1:n}) - \theta(P) \right|^2 \right]$

Theorem 1

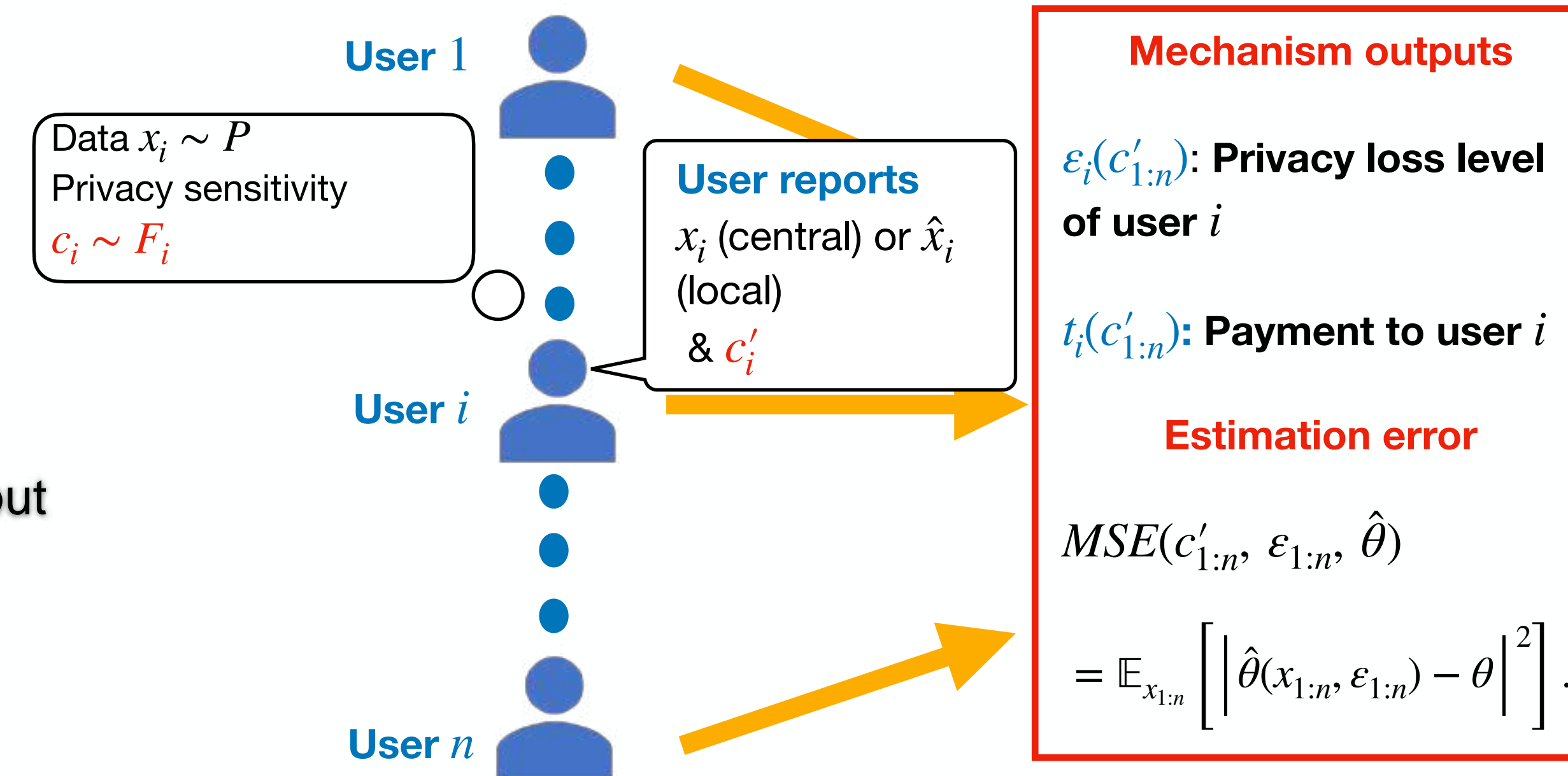
- Assume $\epsilon_1 \leq \dots \leq \epsilon_n \leq 1$. Let \mathcal{P} be family of distributions such that $|X| \leq 1/2$ a.s.
- For $\Theta =$ family of $(\epsilon_i)_{i=1}^n$ -centrally DP estimators

$$\mathcal{L}(\Theta, \mathcal{P}) \geq \mathcal{O}(1) \left(\max_{k \in \{0, 1, \dots, n\}} \frac{1}{n - k + \left(\sum_{i=1}^k \epsilon_i \right)^2} \wedge 1 \right)$$
- For $\Theta =$ family of $(\epsilon_i)_{i=1}^n$ -locally DP estimators

$$\mathcal{L}(\Theta, \mathcal{P}) \geq \mathcal{O}(1) \left(\frac{1}{\sum_{i=1}^n \epsilon_i^2} \wedge 1 \right)$$
- Achievable (up to log factor) by **linear estimators**.

Private Data Acquisition Mechanism

- How $(\epsilon_i)_{i=1}^n$ are endogenized?
- A user's **privacy sensitivity**: per unit cost of privacy loss.



- The cost of user i with privacy sensitive c_i who reports c_i' :

$$\text{Cost}_i(c_i', c_i; \epsilon_{1:n}, t_i, \hat{\theta}) = \mathbb{E} \left[\text{MSE}(c_i', c_{-i}, \epsilon_{1:n}, \hat{\theta}) + c_i \epsilon_i(c_i', c_{-i}) - t_i(c_i', c_{-i}) \right]$$
- The cost of a nonparticipating user: VAR (her best estimate given her data alone).

Platform's problem

$$\min_{\hat{\theta}, \epsilon_{1:n}(\cdot), t_{1:n}(\cdot)} \mathbb{E}_{c_{1:n}} \left[\text{MSE}(c_{1:n}, \epsilon_{1:n}, \hat{\theta}) + \sum_{i=1}^n t_i(c_{1:n}) \right]$$

- Incentive Compatibility (IC):** user i has no incentive to misrepresent her privacy sensitivity when others report truthfully:

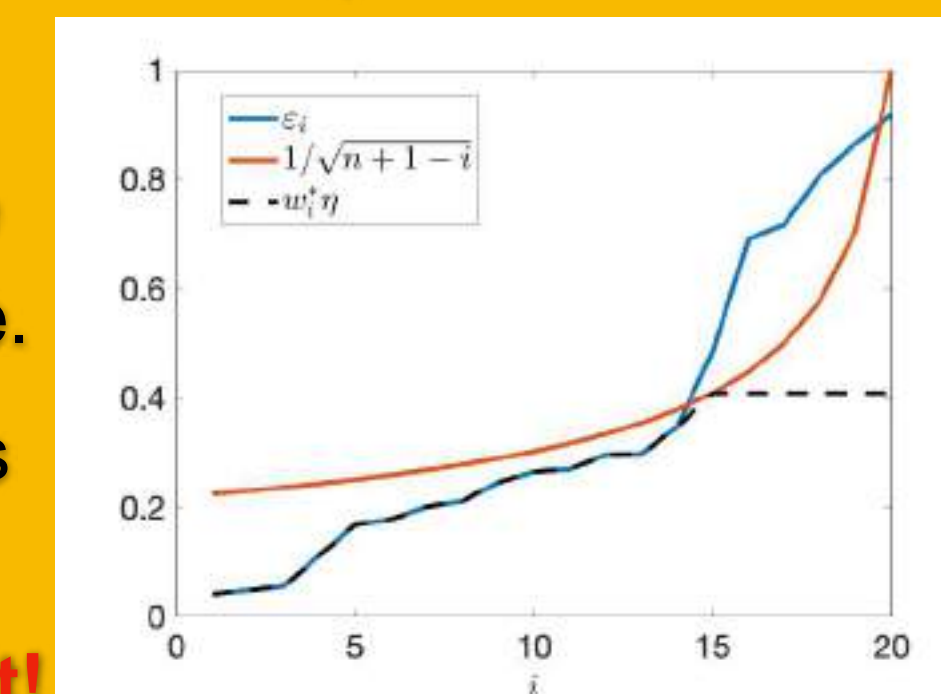
$$\text{Cost}_i(c_i, c_i; \epsilon_{1:n}, t_i, \hat{\theta}) \leq \text{Cost}_i(c_i', c_i; \epsilon_{1:n}, t_i, \hat{\theta})$$
- Individual Rationality (IR):** the platform does not make users worse off by participating in the mechanism:

$$\text{Cost}_i(c_i, c_i; \epsilon_{1:n}, t_i, \hat{\theta}) \leq VAR$$

Proof Ideas [for the centrally DP case]

- Lower bound: Le Cam's method**
 - Replace sup by average over arbitrary $P_1, P_2 \in \mathcal{P}$.
 - We show $\mathcal{L}(\Theta, \mathcal{P}) \geq \frac{(\mu_{P_1} - \mu_{P_2})^2}{8} (1 - \|Q_1 - Q_2\|_{TV})$ where Q_i is the distribution of $\hat{\theta}(x_{1:n})$ when $x_i \sim P_i$.
 - Then, we show that for any k

$$\|Q_1 - Q_2\|_{TV} \leq 2 \|P_1 - P_2\|_{TV} \sum_{i=1}^k (e^{\epsilon_i} - 1) + \sqrt{\frac{n-k}{2}} D_{KL}(P_1, P_2)$$
 - Finally, choose P_1 & P_2 as two Bernoullis and optimize over their means' difference.
- Upper bound idea:** w_i grows proportional to ϵ_i up to some k and then **remains constant!**



From Mechanism Design to Optimization

- We characterize the payment that satisfies IC & IR.
- Plugging it back, we simplify the optimization problem:

$$\min_{\epsilon_{1:n}(c)} \mathbb{E}_{c_{1:n}} \left[(n+1) \text{MSE}(c_{1:n}, \epsilon_{1:n}, \hat{\theta}) + \sum_{i=1}^n \epsilon_i(c_{1:n}) \psi_i(c_i) \right] - nVAR$$

s.t. $\mathbb{E}_{c_{-i}} [\epsilon_i(z, c_{-i})]$ is decreasing in z for all i
- Assume **virtual cost** $\psi_i(c) = c + \frac{F_i(c)}{f_i(c)}$ is increasing.

Theorem II

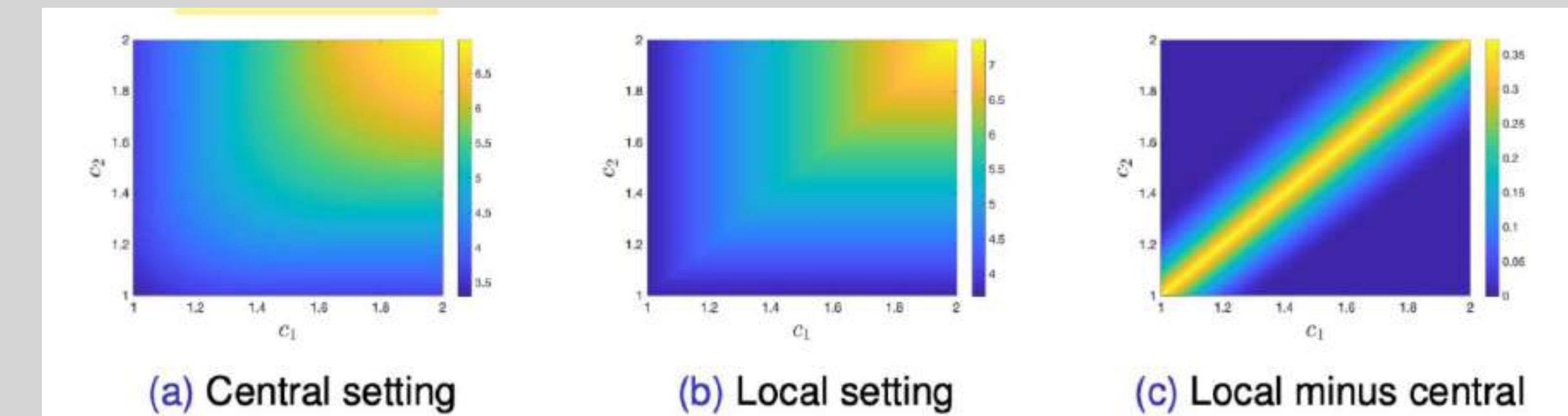
- Let us focus on **linear estimators**.
- For central DP, the optimization problem reduces to

$$\min_{\epsilon_{1:n}} \frac{n+1}{\left(\sum_{j=1}^n \epsilon_j \right)^2} \left(2 + \sum_{i=1}^n \epsilon_i^2 VAR \right) + \sum_{i=1}^n \psi_i(c_i) \epsilon_i$$
- We develop a score-based algorithm to solve this problem in $\mathcal{O}(n \log(n))$.
- For local DP, the optimization problem reduces to

$$\min_{\epsilon_{1:n}} \frac{n+1}{\sum_{i=1}^n \frac{1}{VAR + 2/\epsilon_i^2}} + \sum_{i=1}^n \psi_i(c_i) \epsilon_i$$
- We develop an algorithm to find a δ -approximate solution in time **poly($n, 1/\delta$)**.
- Proofs rely on using KKT conditions, defining auxiliary variables, and characterizing structures of the solution.

Central vs. Local: Platform's objective

- We establish that the platform's optimal cost under central differential privacy setting is always **weakly smaller**.
- Example with two users with c_i uniform over $[1, 2]$:



- However, user's privacy loss can be smaller in the local setting.

