

On the Computational and Statistical Interface and “Big Data”

Michael I. Jordan
University of California, Berkeley

February 17, 2015

What Is the Big Data Phenomenon?

- Science in confirmatory mode (e.g., particle physics)
- Science in exploratory mode (e.g., astronomy, genomics)
- Measurement of human activity, particularly online activity, is generating massive datasets that can be used (e.g.) for personalization and for creating markets
- Sensor networks are becoming pervasive

What Are the Conceptual/Mathematical Issues?

- The need to control **statistical risk** under constraints on algorithmic **runtime**

What Are the Conceptual/Mathematical Issues?

- The need to control **statistical risk** under constraints on algorithmic **runtime**
 - how do risk and runtime trade off as a function of the amount of data?

What Are the Conceptual/Mathematical Issues?

- The need to control **statistical risk** under constraints on algorithmic **runtime**
 - how do risk and runtime trade off as a function of the amount of data?
- Statistical with distributed and streaming data
 - how is inferential quality impacted by communication constraints?

What Are the Conceptual/Mathematical Issues?

- The need to control **statistical risk** under constraints on algorithmic **runtime**
 - how do risk and runtime trade off as a function of the amount of data?
- Statistical with distributed and streaming data
 - how is inferential quality impacted by communication constraints?
- The tradeoff between statistical risk and privacy (and other externalities)

What Are the Conceptual/Mathematical Issues?

- The need to control **statistical risk** under constraints on algorithmic **runtime**
 - how do risk and runtime trade off as a function of the amount of data?
- Statistical with distributed and streaming data
 - how is inferential quality impacted by communication constraints?
- The tradeoff between statistical risk and privacy (and other externalities)
- Many other issues that require a blend of **statistical thinking** (e.g., a focus on sampling, confidence intervals, evaluation, diagnostics, causal inference) and **computational thinking** (e.g., scalability, abstraction)

Our Approach

- Take (classical) statistical decision theory as a mathematical point of departure
- Treat computation, communication, privacy, etc as **constraints** on statistical risk
- This induces tradeoffs among these quantities and the number of data points

Our Approach

- Take (classical) statistical decision theory as a mathematical point of departure
- Treat computation, communication, privacy, etc as **constraints** on statistical risk
- This induces **tradeoffs** among these quantities and the number of data points
- Under the hood: geometry, information theory and optimization

Background

- In the 1930's, Wald laid the foundations of statistical decision theory
- Given a family of probability distributions \mathcal{P} , a **parameter** $\theta(P)$ for each $P \in \mathcal{P}$, an **estimator** $\hat{\theta}$, and a **loss** $l(\hat{\theta}, \theta(P))$, define the **risk**:

$$R_P(\hat{\theta}) := \mathbb{E}_P \left[l(\hat{\theta}, \theta(P)) \right]$$

- Minimax principle [Wald, '39, '43]: choose estimator minimizing worst-case risk:

$$\sup_{P \in \mathcal{P}} \mathbb{E}_P \left[l(\hat{\theta}, \theta(P)) \right]$$

Part I: Privacy and Minimax Risk

with John Duchi and Martin Wainwright
University of California, Berkeley

Privacy and Risk

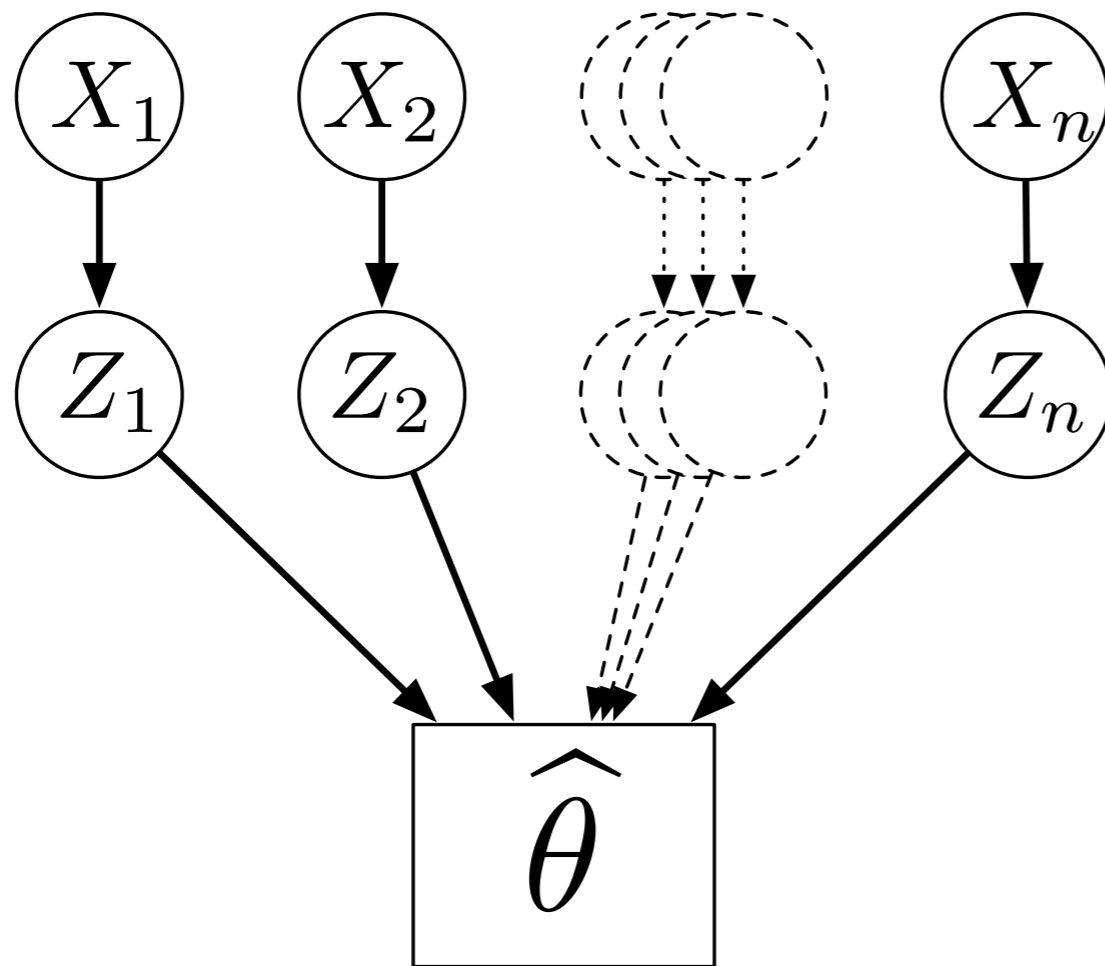
- Individuals are not generally willing to allow their personal data to be used without control on how it will be used and how much privacy loss they will incur
- We will quantify “privacy loss” via [differential privacy](#)
- We then treat differential privacy as a [constraint](#) on inference via statistical decision theory
- This yields (personal) tradeoffs between privacy loss and inferential gain

A model of privacy

Local privacy: providers do not trust collector [Warner 65, Evfimievski et al. 03]

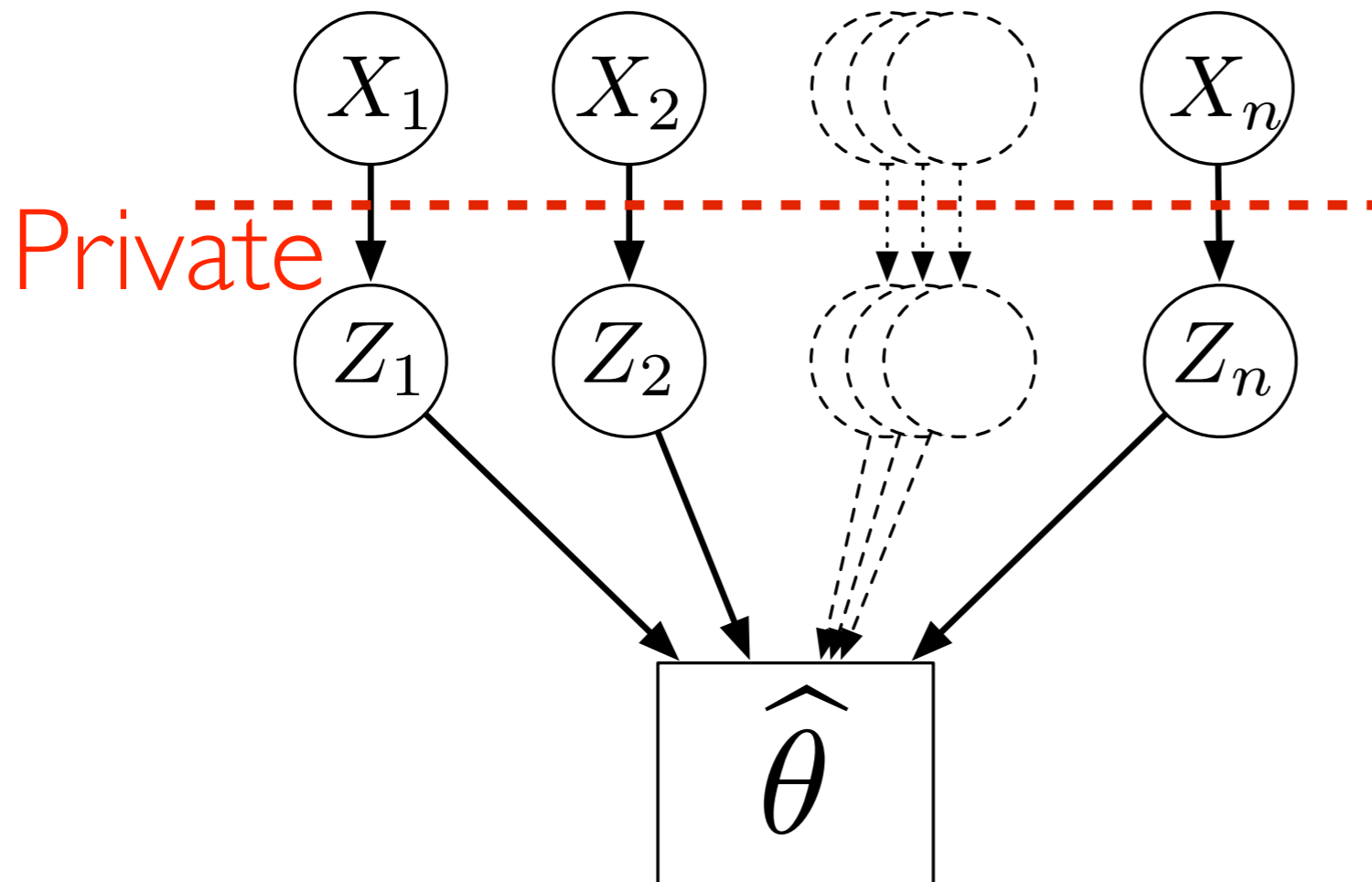
A model of privacy

Local privacy: providers do not trust collector [Warner 65, Evfimievski et al. 03]



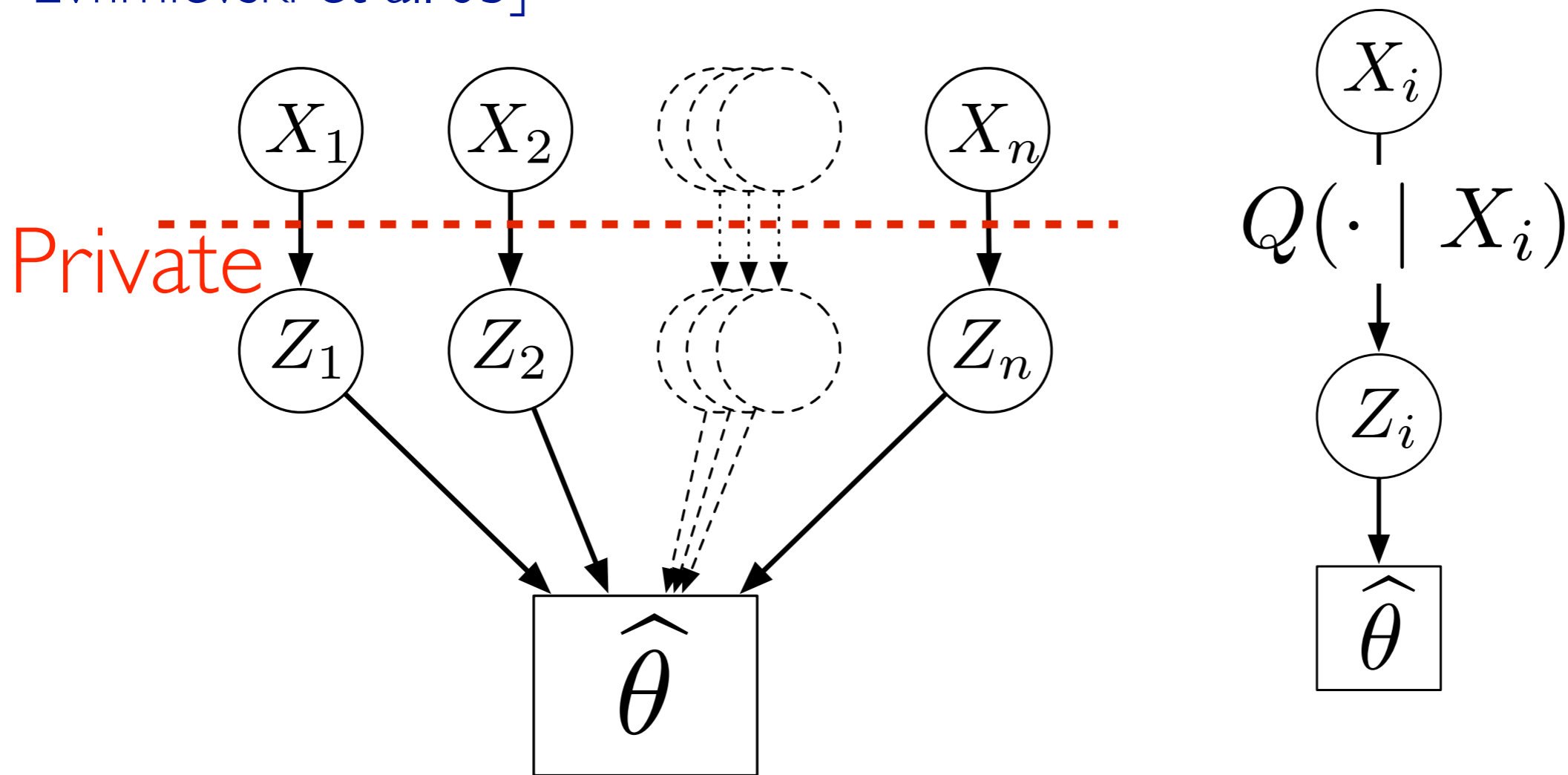
A model of privacy

Local privacy: providers do not trust collector [Warner 65, Evfimievski et al. 03]



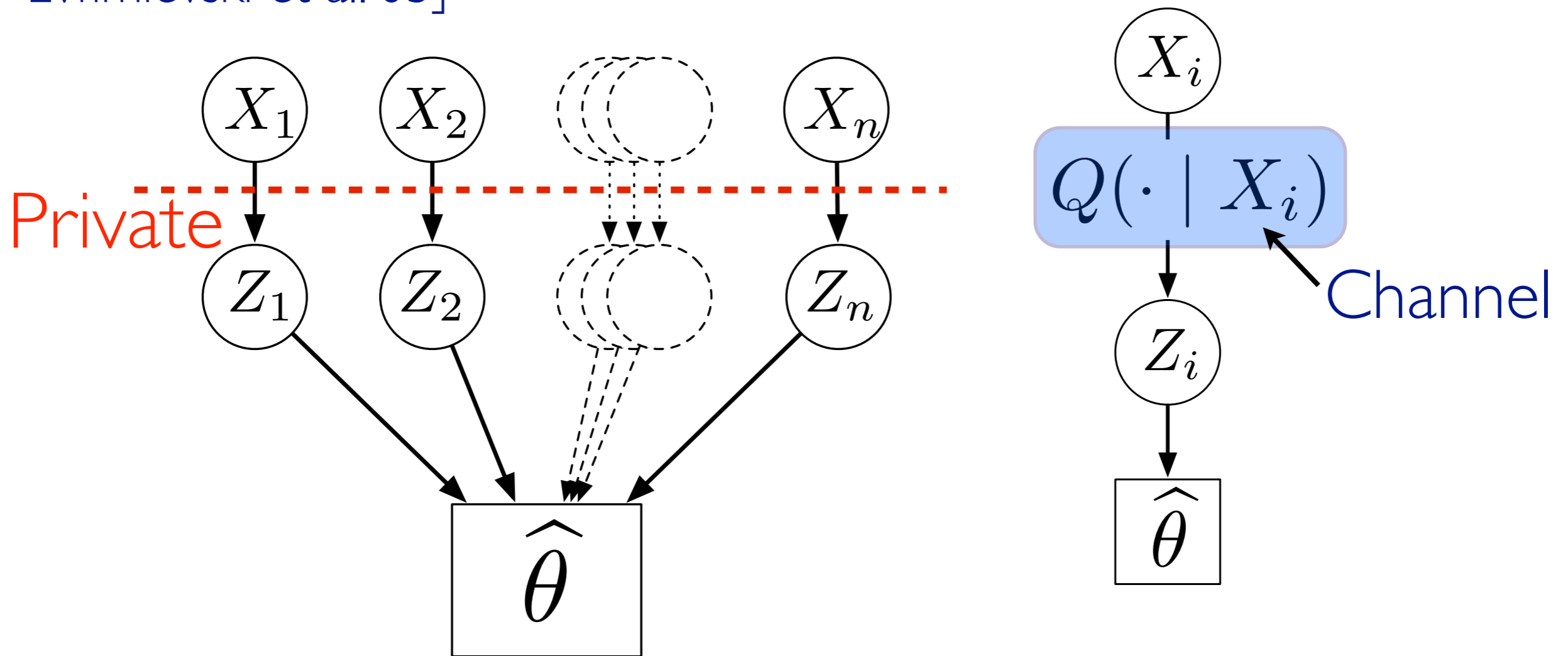
A model of privacy

Local privacy: providers do not trust collector [Warner 65, Evfimievski et al. 03]



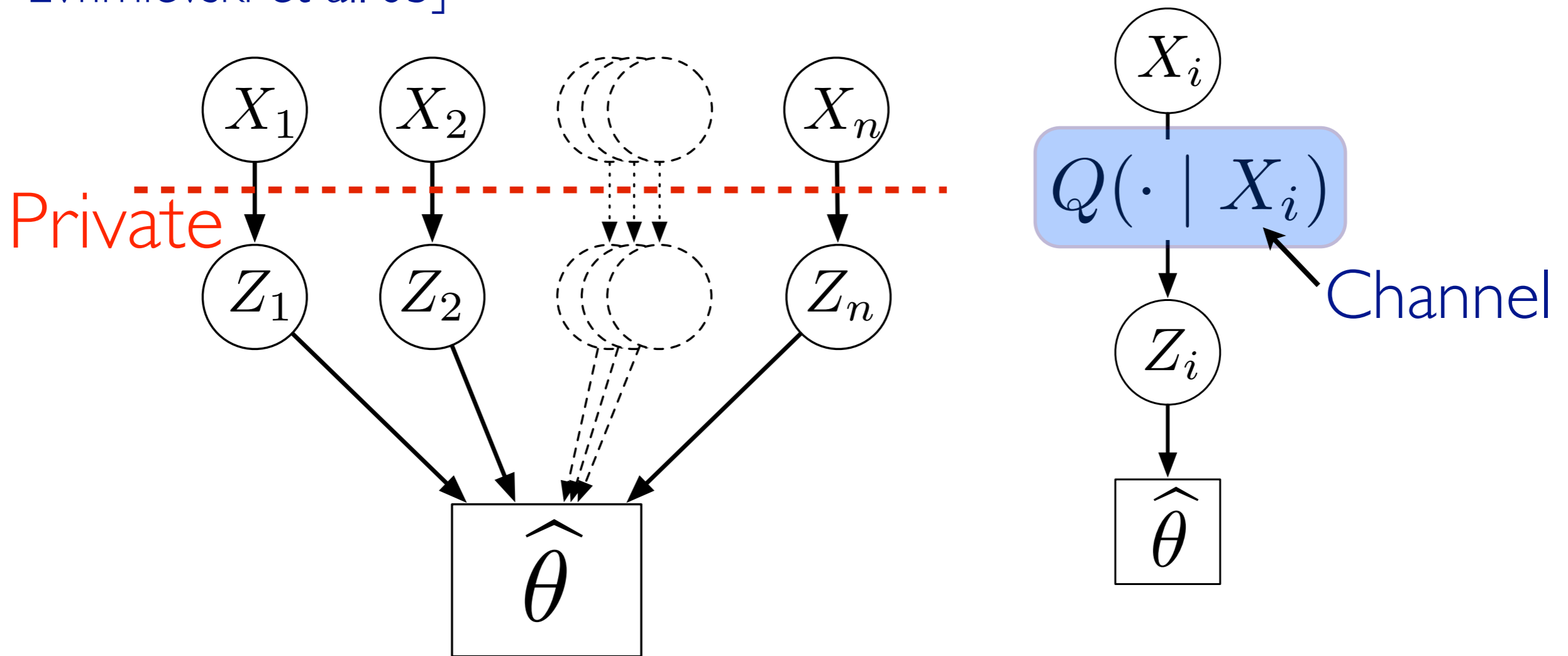
A model of privacy

Local privacy: providers do not trust collector [Warner 65, Evfimievski et al. 03]



A model of privacy

Local privacy: providers do not trust collector [Warner 65, Evfimievski et al. 03]

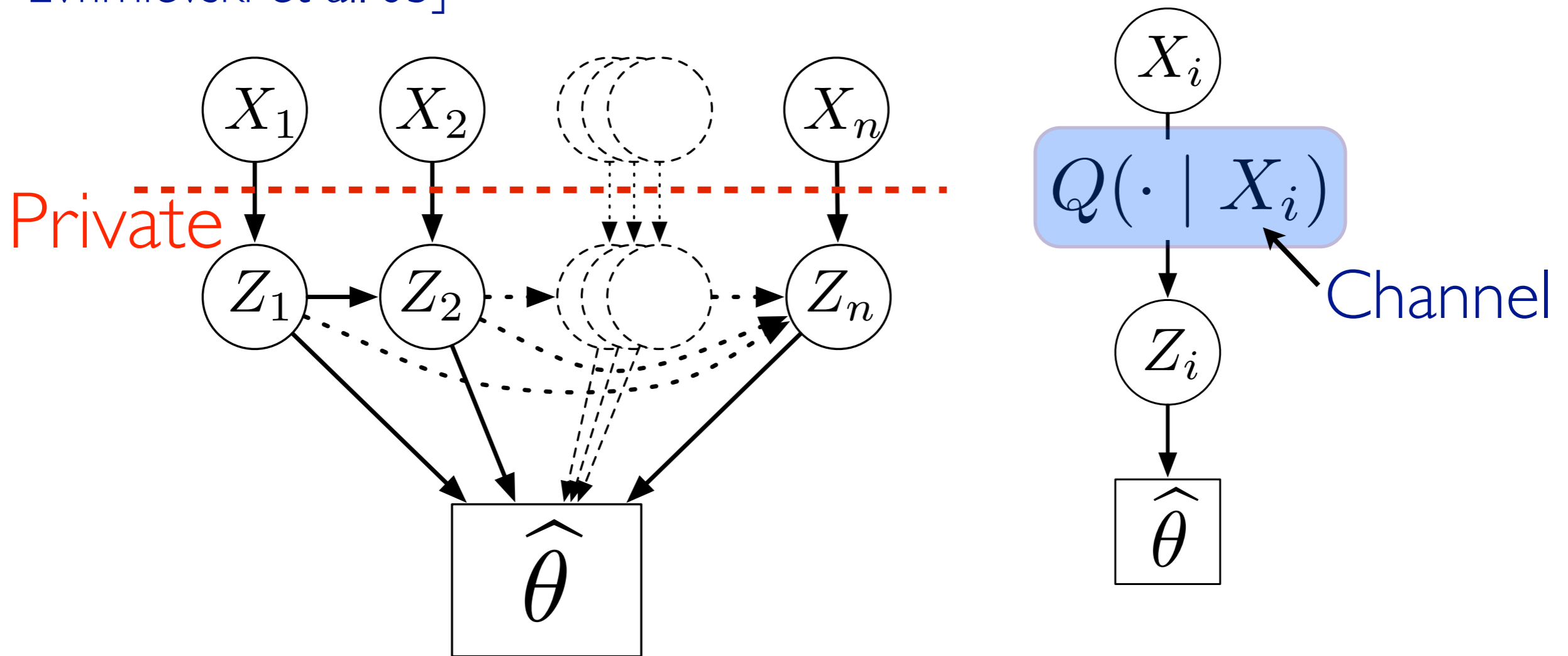


Individuals $i \in \{1, \dots, n\}$ with **private data** $X_i \stackrel{\text{iid}}{\sim} P$

Estimator $Z_1^n \mapsto \hat{\theta}(Z_1^n)$

A model of privacy

Local privacy: providers do not trust collector [Warner 65, Evfimievski et al. 03]



Individuals $i \in \{1, \dots, n\}$ with **private data** $X_i \stackrel{\text{iid}}{\sim} P$

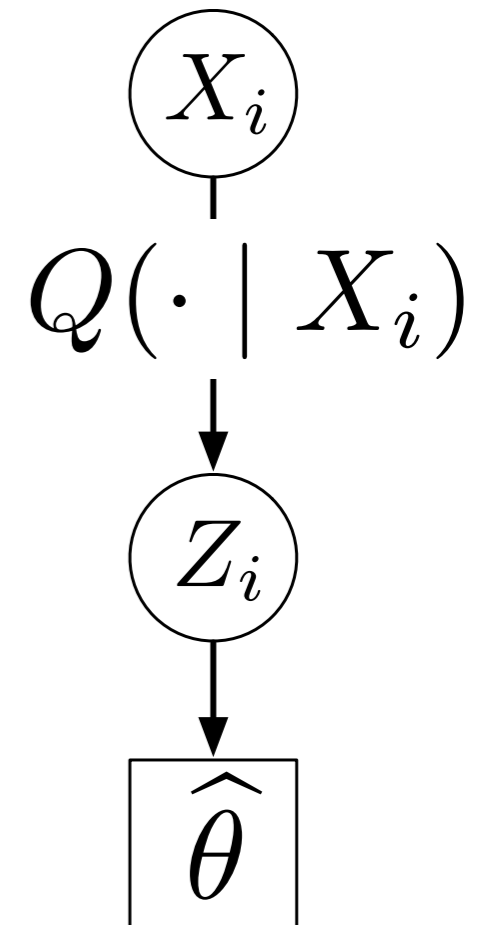
Estimator $Z_1^n \mapsto \hat{\theta}(Z_1^n)$

Definitions of privacy

Definition: channel Q is α -differentially private if

$$\sup_{S, x \in \mathcal{X}, x' \in \mathcal{X}} \frac{Q(Z \in S | x)}{Q(Z \in S | x')} \leq \exp(\alpha)$$

[Dwork, McSherry, Nissim, Smith 06]

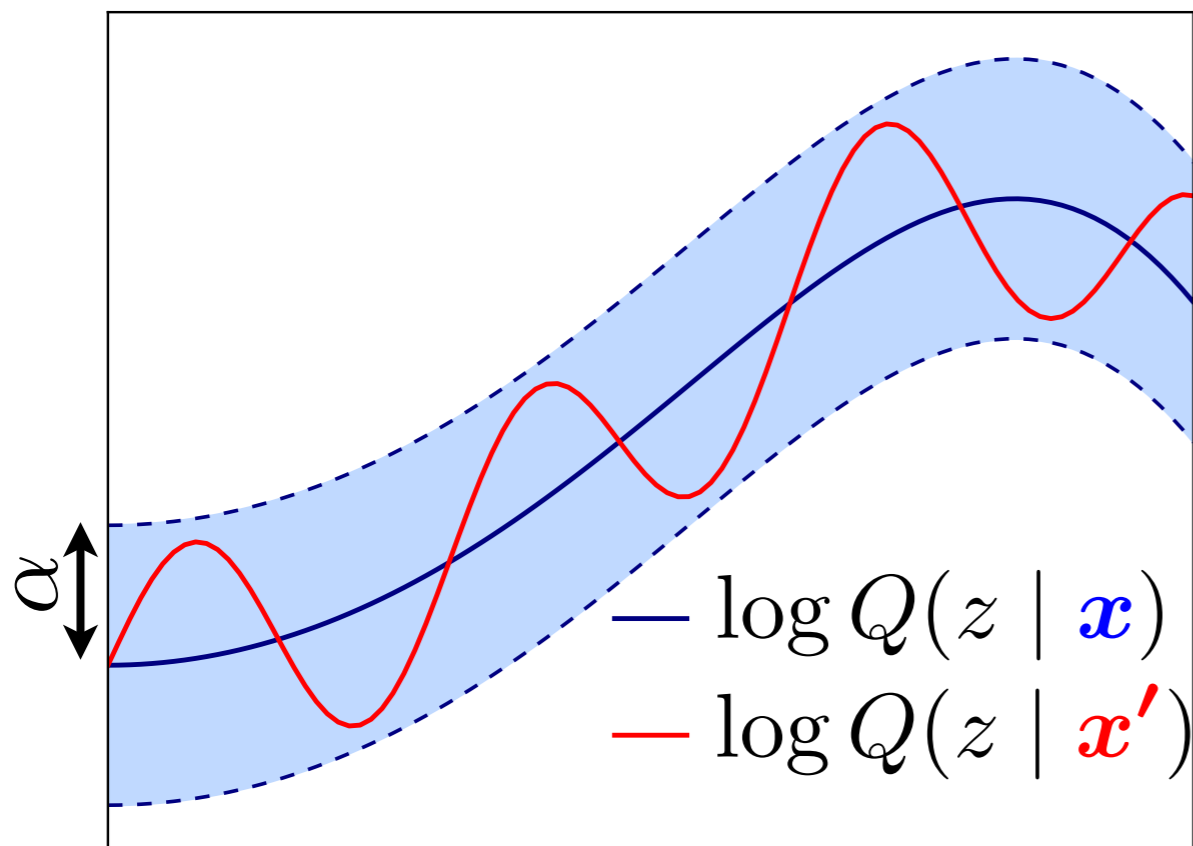
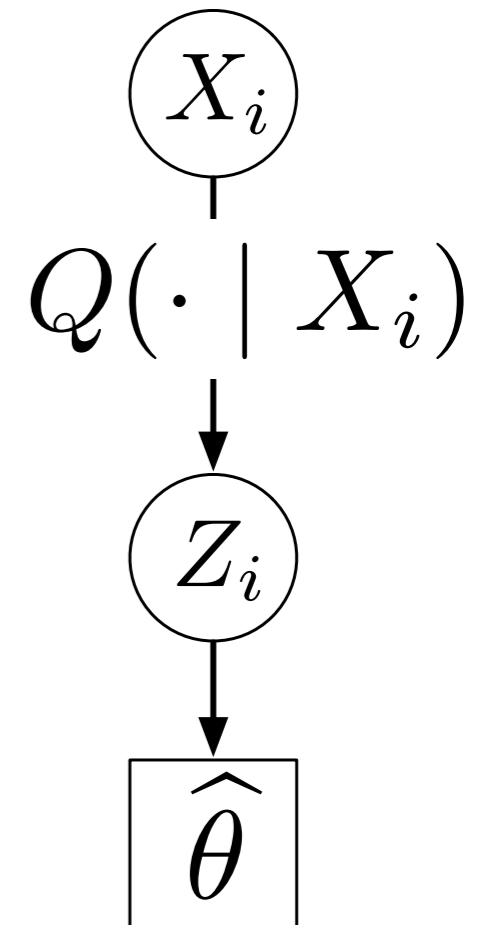


Definitions of privacy

Definition: channel Q is α -differentially private if

$$\sup_{S, x \in \mathcal{X}, x' \in \mathcal{X}} \frac{Q(Z \in S | x)}{Q(Z \in S | x')} \leq \exp(\alpha)$$

[Dwork, McSherry, Nissim, Smith 06]



Private Minimax Risk

Central object of study: minimax risk

- Parameter $\theta(P)$ of distribution
- Family of distributions \mathcal{P}
- Loss ℓ measuring error

Minimax risk

$$\mathfrak{M}_n(\theta(\mathcal{P}), \ell) := \inf_{\hat{\theta}} \sup_{P \in \mathcal{P}} \mathbb{E}_P \left[\ell(\hat{\theta}(X_1^n), \theta(P)) \right]$$

Private Minimax Risk

Central object of study: minimax risk

- Parameter $\theta(P)$ of distribution
- Family of distributions \mathcal{P}
- Loss ℓ measuring error
- Family \mathcal{Q}_α of **private** channels

α -private Minimax risk

$$\mathfrak{M}_n(\theta(\mathcal{P}), \ell, \alpha) := \inf_{Q \in \mathcal{Q}_\alpha} \inf_{\hat{\theta}} \sup_{P \in \mathcal{P}} \mathbb{E}_{P, Q} \left[\ell(\hat{\theta}(Z_1^n), \theta(P)) \right]$$

Private Minimax Risk

Central object of study: minimax risk

- Parameter $\theta(P)$ of distribution
- Family of distributions \mathcal{P}
- Loss ℓ measuring error
- Family \mathcal{Q}_α of **private** channels

α -private Minimax risk

$$\mathfrak{M}_n(\theta(\mathcal{P}), \ell, \alpha) := \inf_{Q \in \mathcal{Q}_\alpha} \inf_{\hat{\theta}} \sup_{P \in \mathcal{P}} \mathbb{E}_{P, Q} \left[\ell(\hat{\theta}(Z_1^n), \theta(P)) \right]$$

↑
Best α -private channel

Private Minimax Risk

Central object of study: minimax risk

- Parameter $\theta(P)$ of distribution
- Family of distributions \mathcal{P}
- Loss ℓ measuring error
- Family \mathcal{Q}_α of **private** channels

α -private Minimax risk

$$\mathfrak{M}_n(\theta(\mathcal{P}), \ell, \alpha) := \inf_{Q \in \mathcal{Q}_\alpha} \inf_{\hat{\theta}} \sup_{P \in \mathcal{P}} \mathbb{E}_{P, Q} \left[\ell(\hat{\theta}(Z_1^n), \theta(P)) \right]$$

Best α -private channel

Minimax risk under privacy constraint

Vignette: private mean (location) estimation

Example: estimate reasons for hospital visits
Patients admitted to hospital for substance abuse
Estimate prevalence of different substances

Vignette: private mean (location) estimation

Example: estimate reasons for hospital visits

Patients admitted to hospital for substance abuse

Estimate prevalence of different substances

<input checked="" type="checkbox"/> Alcohol	$\theta_1 = .45$
<input checked="" type="checkbox"/> Cocaine	$\theta_2 = .32$
<input type="checkbox"/> Heroin	$\theta_3 = .16$
<input type="checkbox"/> Cannabis	$\theta_4 = .20$
<input type="checkbox"/> LSD	$\theta_5 = .00$
<input type="checkbox"/> Amphetamines	$\theta_6 = .02$

Proportions $\theta =$

Vignette: mean estimation

Consider estimation of mean $\theta(P) := \mathbb{E}_P[X] \in \mathbb{R}^d$, errors measured in ℓ_∞ -norm, i.e. $\mathbb{E}[\|\hat{\theta} - \theta\|_\infty]$ for

$$\mathcal{P}_d := \left\{ \text{distributions } P \text{ supported on } [-1, 1]^d \right\}$$

Vignette: mean estimation

Consider estimation of mean $\theta(P) := \mathbb{E}_P[X] \in \mathbb{R}^d$, errors measured in ℓ_∞ -norm, i.e. $\mathbb{E}[\|\hat{\theta} - \theta\|_\infty]$ for

$$\mathcal{P}_d := \left\{ \text{distributions } P \text{ supported on } [-1, 1]^d \right\}$$

Proposition:

Minimax rate

$$\mathfrak{M}_n(\mathcal{P}_d, \|\cdot\|_\infty) \asymp \min \left\{ 1, \frac{\sqrt{\log d}}{\sqrt{n}} \right\}$$

(achieved by sample mean)

Vignette: mean estimation

Consider estimation of mean $\theta(P) := \mathbb{E}_P[X] \in \mathbb{R}^d$, errors measured in ℓ_∞ -norm, i.e. $\mathbb{E}[\|\hat{\theta} - \theta\|_\infty]$ for

$$\mathcal{P}_d := \left\{ \text{distributions } P \text{ supported on } [-1, 1]^d \right\}$$

Proposition:

Private minimax rate for $\alpha = O(1)$

$$\mathfrak{M}_n(\mathcal{P}_d, \|\cdot\|_\infty, \alpha) \asymp \min \left\{ 1, \frac{\sqrt{d \log d}}{\sqrt{n\alpha^2}} \right\}$$

Vignette: mean estimation

Consider estimation of mean $\theta(P) := \mathbb{E}_P[X] \in \mathbb{R}^d$, errors measured in ℓ_∞ -norm, i.e. $\mathbb{E}[\|\hat{\theta} - \theta\|_\infty]$ for

$$\mathcal{P}_d := \left\{ \text{distributions } P \text{ supported on } [-1, 1]^d \right\}$$

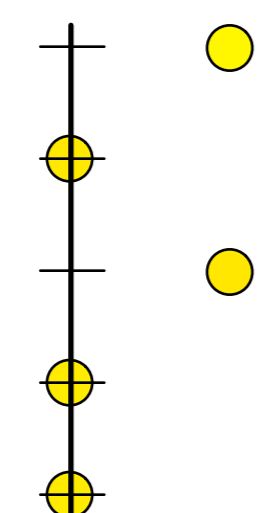
Proposition:

Private minimax rate for $\alpha = O(1)$

$$\mathfrak{M}_n(\mathcal{P}_d, \|\cdot\|_\infty, \alpha) \asymp \min \left\{ 1, \frac{\sqrt{d \log d}}{\sqrt{n\alpha^2}} \right\}$$

Note: Effective sample size $n \mapsto n\alpha^2/d$

Optimal mechanism?

$$X = \begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$


Non-private
observation

Optimal mechanism?

$$X = \begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \quad \begin{array}{c} + \\ \oplus \\ + \\ \oplus \\ \oplus \end{array} \quad \begin{array}{c} \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \end{array} \quad Z = X + W = \begin{bmatrix} 1 + W_1 \\ 0 + W_2 \\ 1 + W_3 \\ 0 + W_4 \\ 0 + W_5 \end{bmatrix} \quad \begin{array}{c} | \text{---} \oplus \text{---} \bullet | \\ | \text{---} \bullet \text{---} \oplus \text{---} | \\ | \text{---} \oplus \text{---} \bullet | \\ | \text{---} \oplus \text{---} \bullet | \\ | \text{---} \oplus \text{---} \bullet | \end{array}$$

Non-private
observation

Idea 1: add independent **noise**
(e.g. standard Laplace
mechanism)

Optimal mechanism?

$$X = \begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \quad \begin{array}{c} + \\ \oplus \\ + \\ \oplus \\ \oplus \end{array} \quad \begin{array}{c} \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \end{array} \quad Z = X + W = \begin{bmatrix} 1 + W_1 \\ 0 + W_2 \\ 1 + W_3 \\ 0 + W_4 \\ 0 + W_5 \end{bmatrix} \quad \begin{array}{c} | \text{---} \oplus \text{---} \bullet | \\ | \text{---} \bullet \text{---} \oplus \text{---} | \\ | \text{---} \oplus \text{---} \bullet | \\ | \text{---} \oplus \text{---} \bullet | \\ | \text{---} \oplus \text{---} \bullet | \end{array}$$

Non-private
observation

Idea 1: add independent **noise**
(e.g. standard Laplace
mechanism)

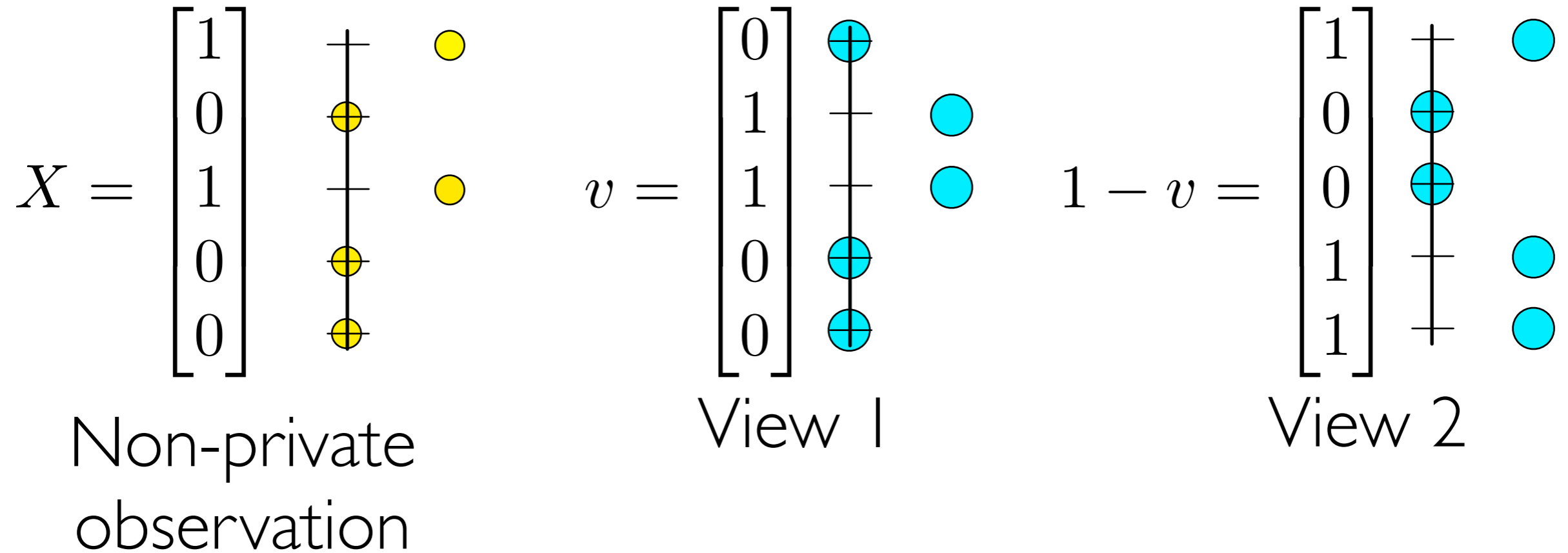
Problem: magnitude much too large
(this is unavoidable: *provably sub-optimal*)

Optimal mechanism

$$X = \begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \quad \begin{array}{c} + \\ \bullet \\ + \\ \bullet \\ \bullet \end{array} \quad \begin{array}{c} \bullet \\ \\ \bullet \\ \\ \end{array}$$

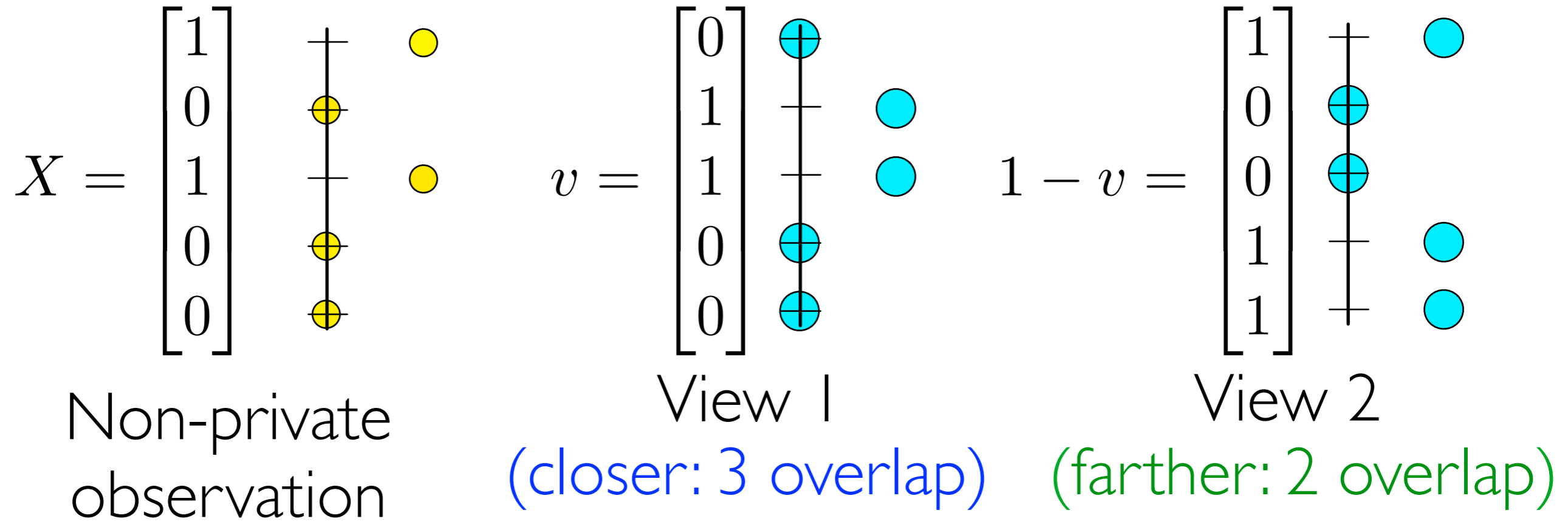
Non-private
observation

Optimal mechanism



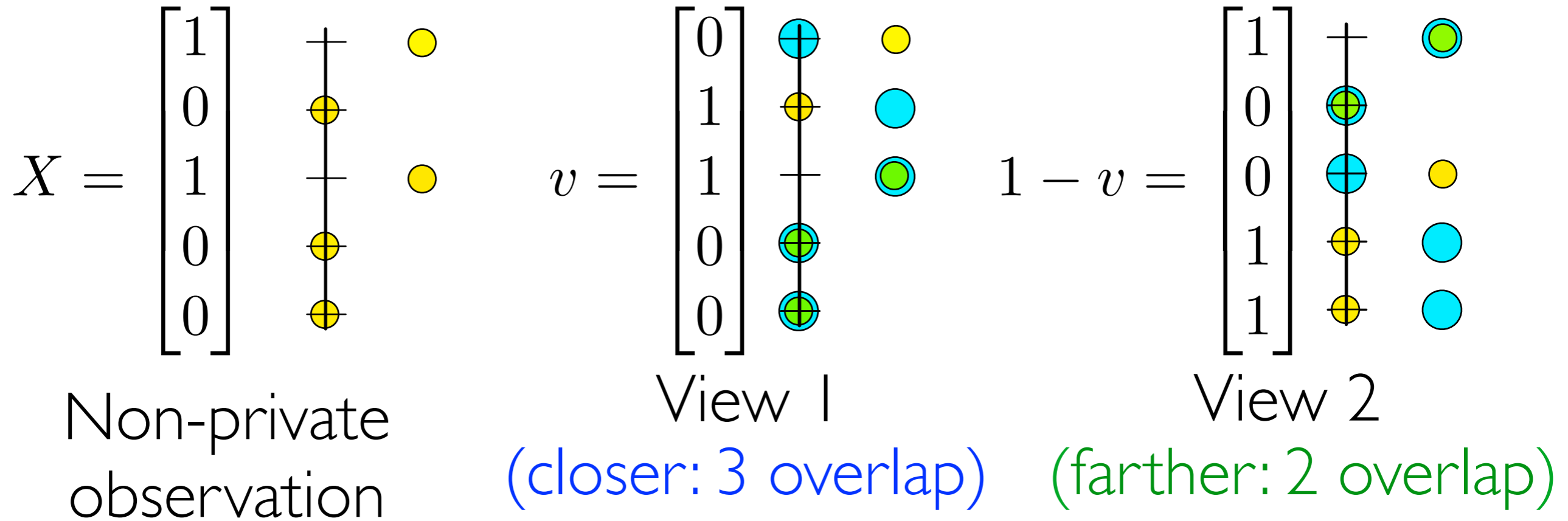
- Draw v uniformly in $\{0, 1\}^d$

Optimal mechanism



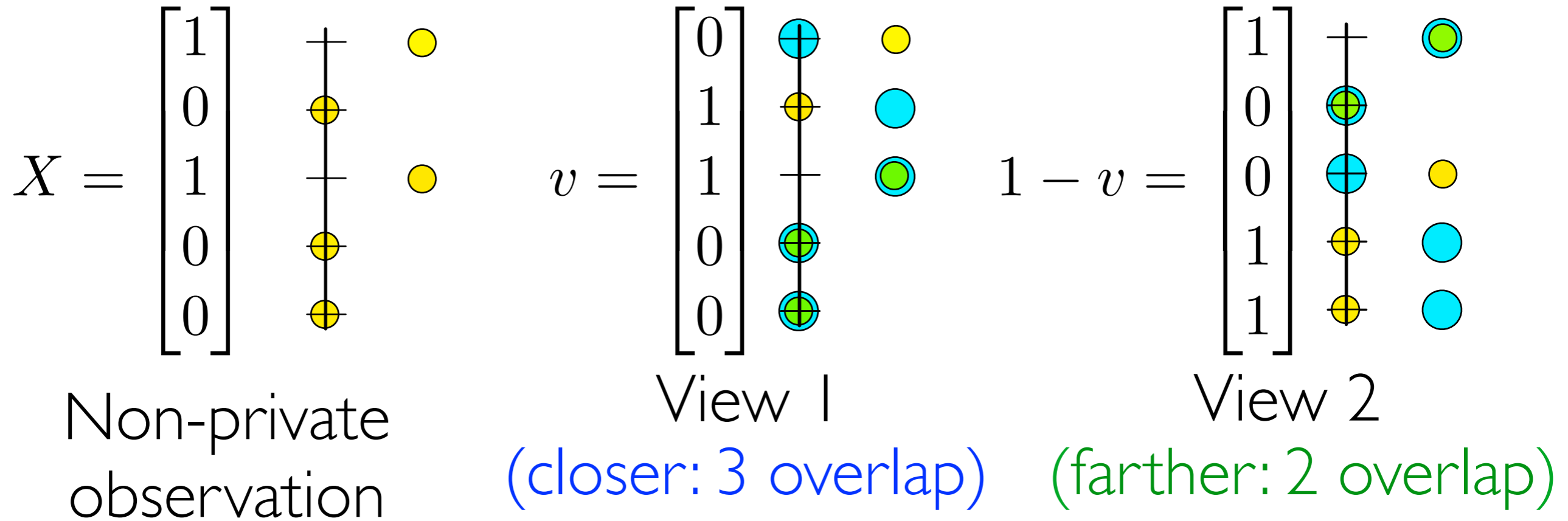
- Draw v uniformly in $\{0, 1\}^d$
- With probability $\frac{e^\alpha}{1 + e^\alpha}$ choose **closer** of v and $1 - v$ to X
- otherwise, choose **farther**

Optimal mechanism



- Draw v uniformly in $\{0, 1\}^d$
- With probability $\frac{e^\alpha}{1 + e^\alpha}$ choose closer of v and $1 - v$ to X
- otherwise, choose farther

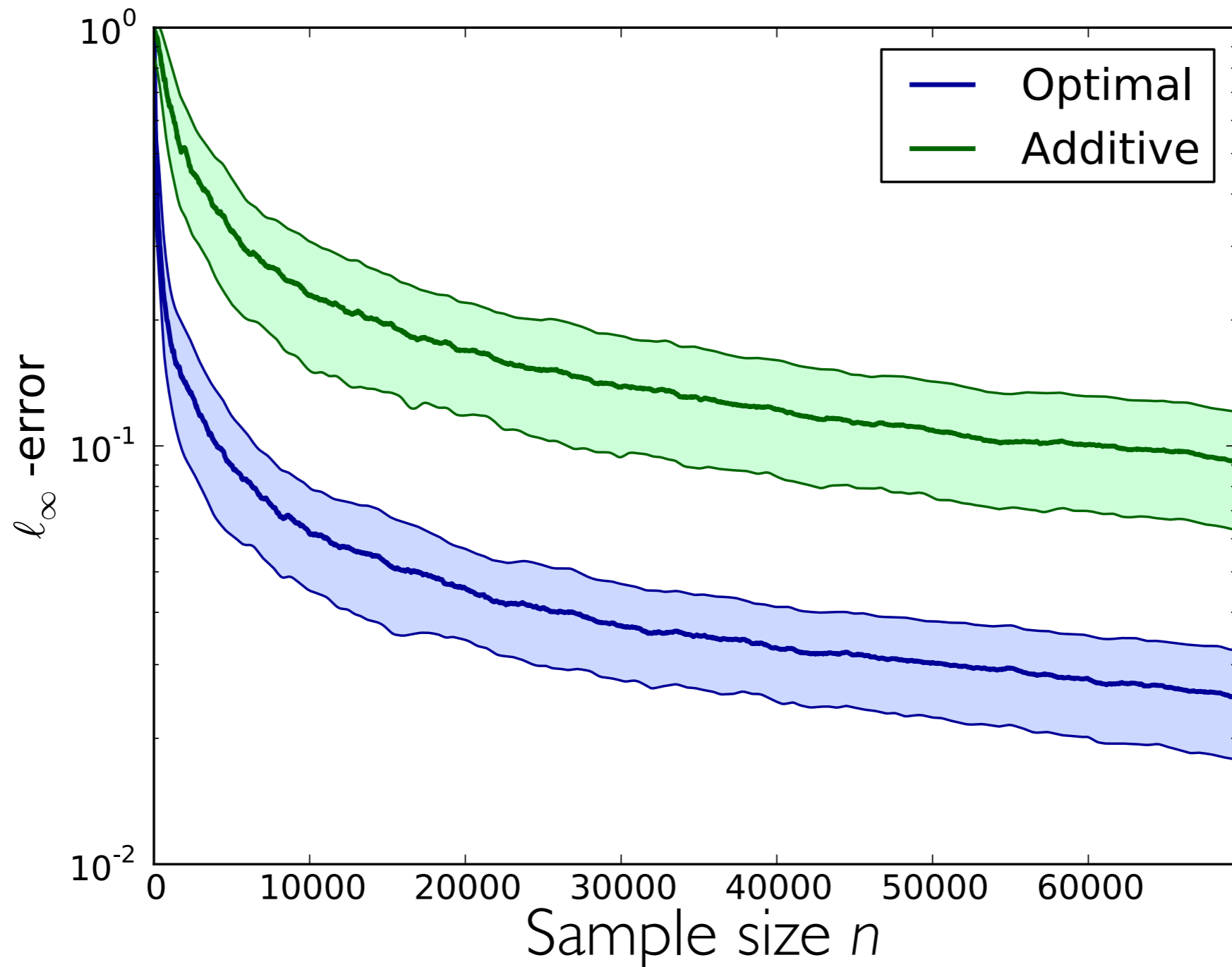
Optimal mechanism



- Draw v uniformly in $\{0, 1\}^d$
- With probability $\frac{e^\alpha}{1 + e^\alpha}$ choose **closer** of v and $1 - v$ to X
- otherwise, choose **farther**

At end:
Compute sample average and de-bias

Empirical evidence



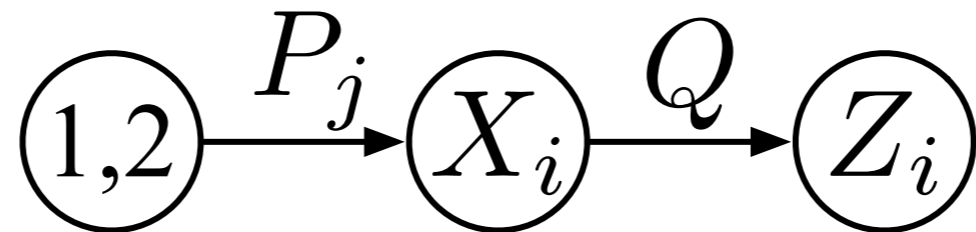
Data source: Drug Abuse Warning Network

Estimate proportion of emergency room visits involving different substances

Sample size reductions

Given α -private channel Q , pair $\{P_1, P_2\}$ induces marginal

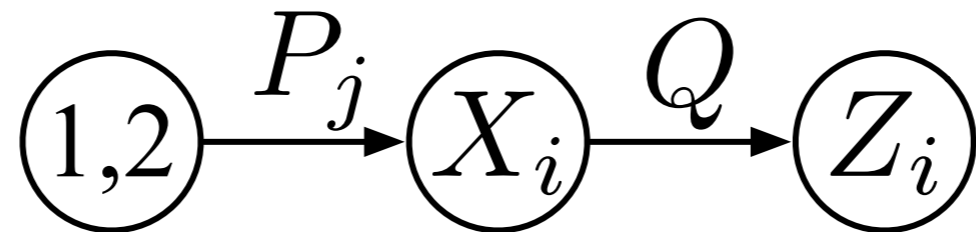
$$M_j(S) := \int Q(S | x_1, \dots, x_n) dP_j^n(x_1, \dots, x_n)$$



Sample size reductions

Given α -private channel Q , pair $\{P_1, P_2\}$ induces marginal

$$M_j(S) := \int Q(S | x_1, \dots, x_n) dP_j^n(x_1, \dots, x_n)$$

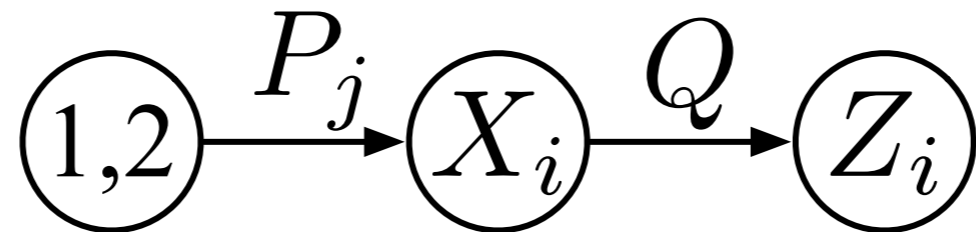


Question: How much “contraction” does privacy induce?

Sample size reductions

Given α -private channel Q , pair $\{P_1, P_2\}$ induces marginal

$$M_j(S) := \int Q(S | x_1, \dots, x_n) dP_j^n(x_1, \dots, x_n)$$



Question: How much “contraction” does privacy induce?

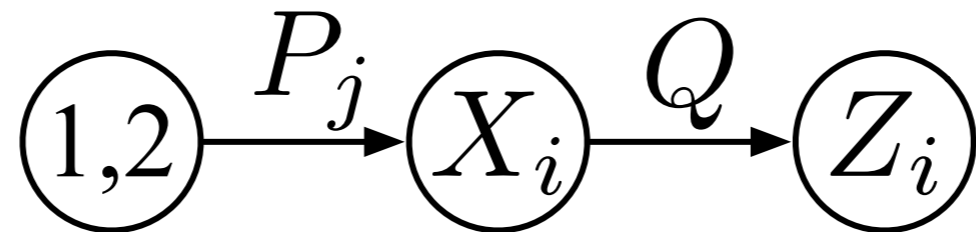
Theorem (data processing): for any α -private channel and i.i.d. sample of size n

$$D_{\text{kl}}(M_1 \| M_2) + D_{\text{kl}}(M_2 \| M_1) \leq 4n(e^\alpha - 1)^2 \|P_1 - P_2\|_{\text{TV}}^2$$

Sample size reductions

Given α -private channel Q , pair $\{P_1, P_2\}$ induces marginal

$$M_j(S) := \int Q(S | x_1, \dots, x_n) dP_j^n(x_1, \dots, x_n)$$



Question: How much “contraction” does privacy induce?

Theorem (data processing): for any α -private channel and i.i.d. sample of size n

$$D_{\text{kl}}(M_1 \| M_2) + D_{\text{kl}}(M_2 \| M_1) \leq 4n(e^\alpha - 1)^2 \|P_1 - P_2\|_{\text{TV}}^2$$

Note: $n \mapsto n\alpha^2$ for $\alpha \lesssim 1$

Final remarks: privacy

Rough technique: Reduction of estimation to testing, then apply information-theoretic testing lower bounds

[Le Cam, Hasminskii, Ibragimov, Assouad, Birge, Barron, Yu, ...]

Additional examples

- Fixed-design regression
- Convex risk minimization
- Multinomial estimation
- Nonparametric density estimation

Almost always: effective sample size reduction $n \mapsto n\alpha^2$

In d -dimensional problems: $n \mapsto \frac{n\alpha^2}{d}$

Key: Allows identification of new optimal mechanisms

Part II: Communication and Minimax Risk

with John Duchi, Martin Wainwright and Yuchen
Zhang

University of California, Berkeley

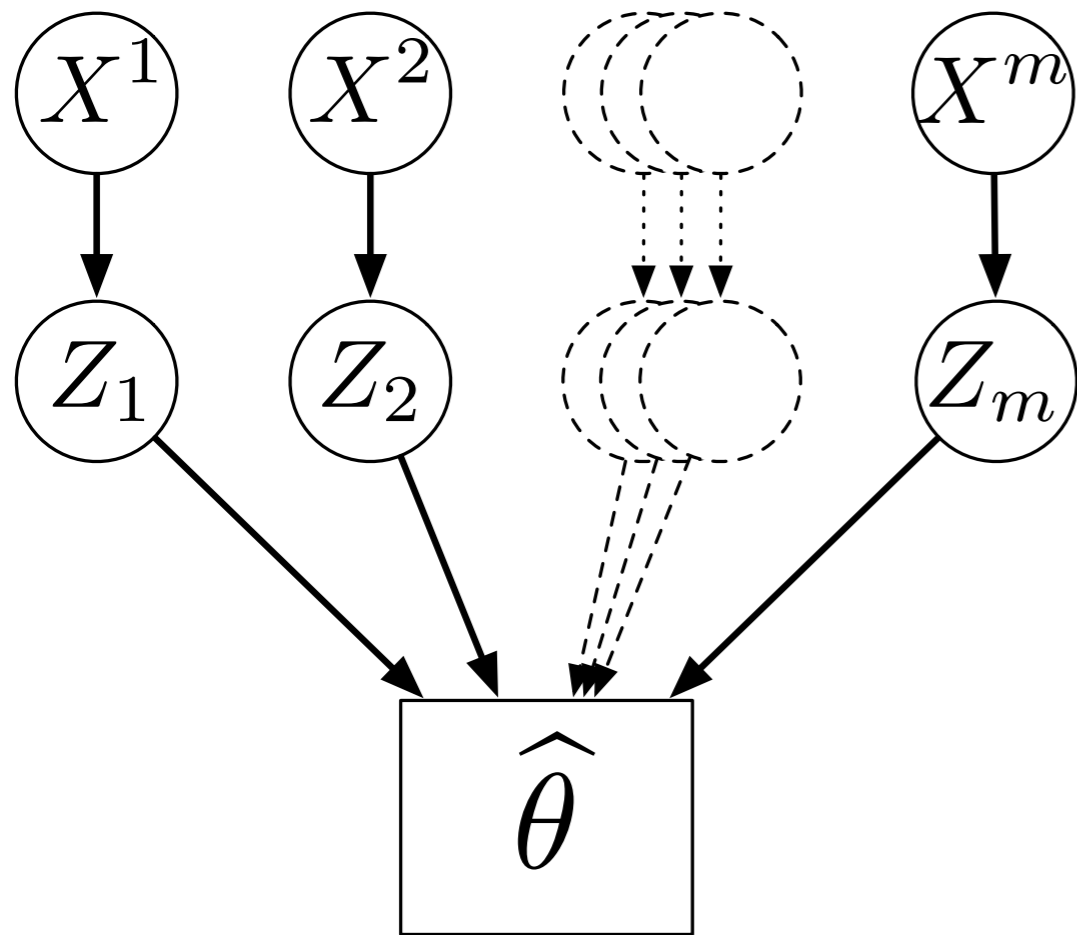
Communication-constraints

Communication-constraints

- Large data necessitates distributed storage
- Independent data collection (hospitals)
- Privacy?

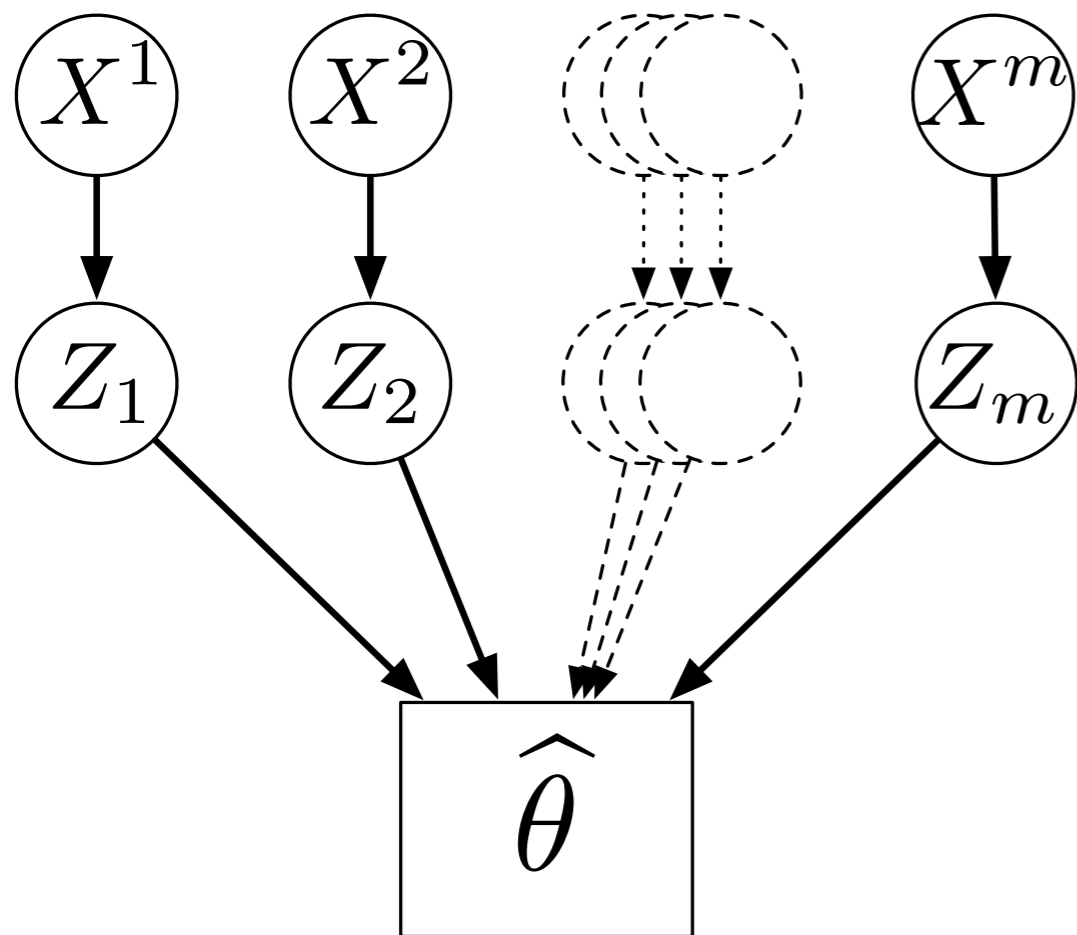
Communication-constraints

- Large data necessitates distributed storage
- Independent data collection (hospitals)
- Privacy?



Communication-constraints

- Large data necessitates distributed storage
- Independent data collection (hospitals)
- Privacy?



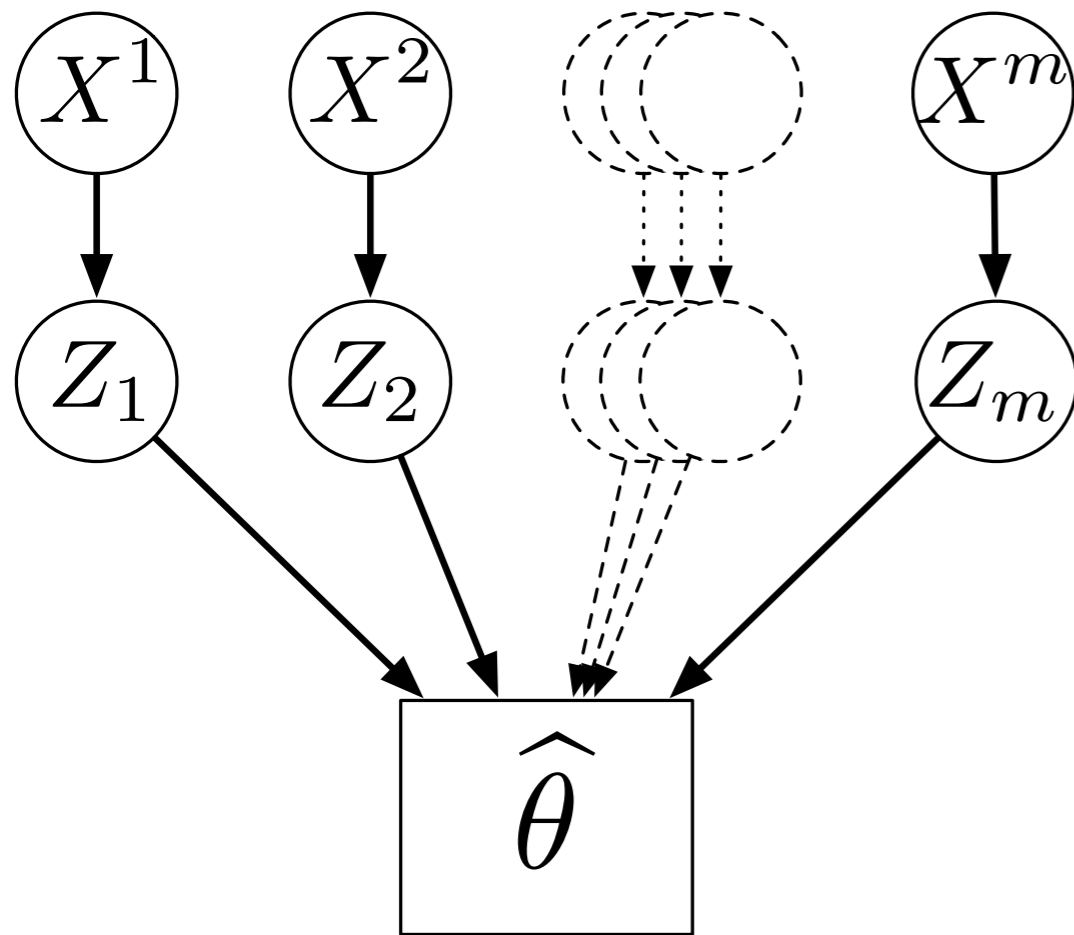
Setting: each of m agents has sample of size n

$$X^i = (X_1^i, X_2^i, \dots, X_n^i)$$

Messages Z_i to fusion center

Communication-constraints

- Large data necessitates distributed storage
- Independent data collection (hospitals)
- Privacy?



Setting: each of m agents has sample of size n

$$X^i = (X_1^i, X_2^i, \dots, X_n^i)$$

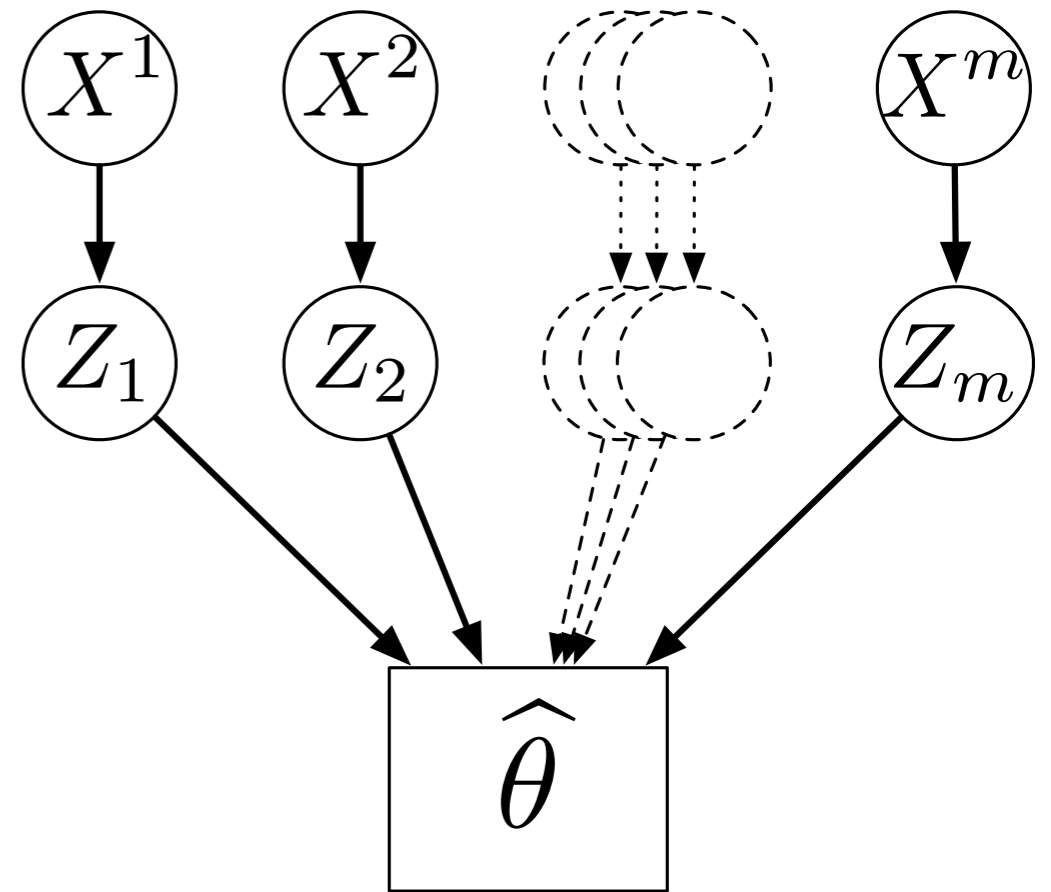
Messages Z_i to fusion center

Question: tradeoffs
between communication
and statistical utility?

Minimax communication

Central object of study:

- Parameter $\theta(P)$ of distribution
- Family of distributions \mathcal{P}
- Loss $\|\cdot\|_2^2$



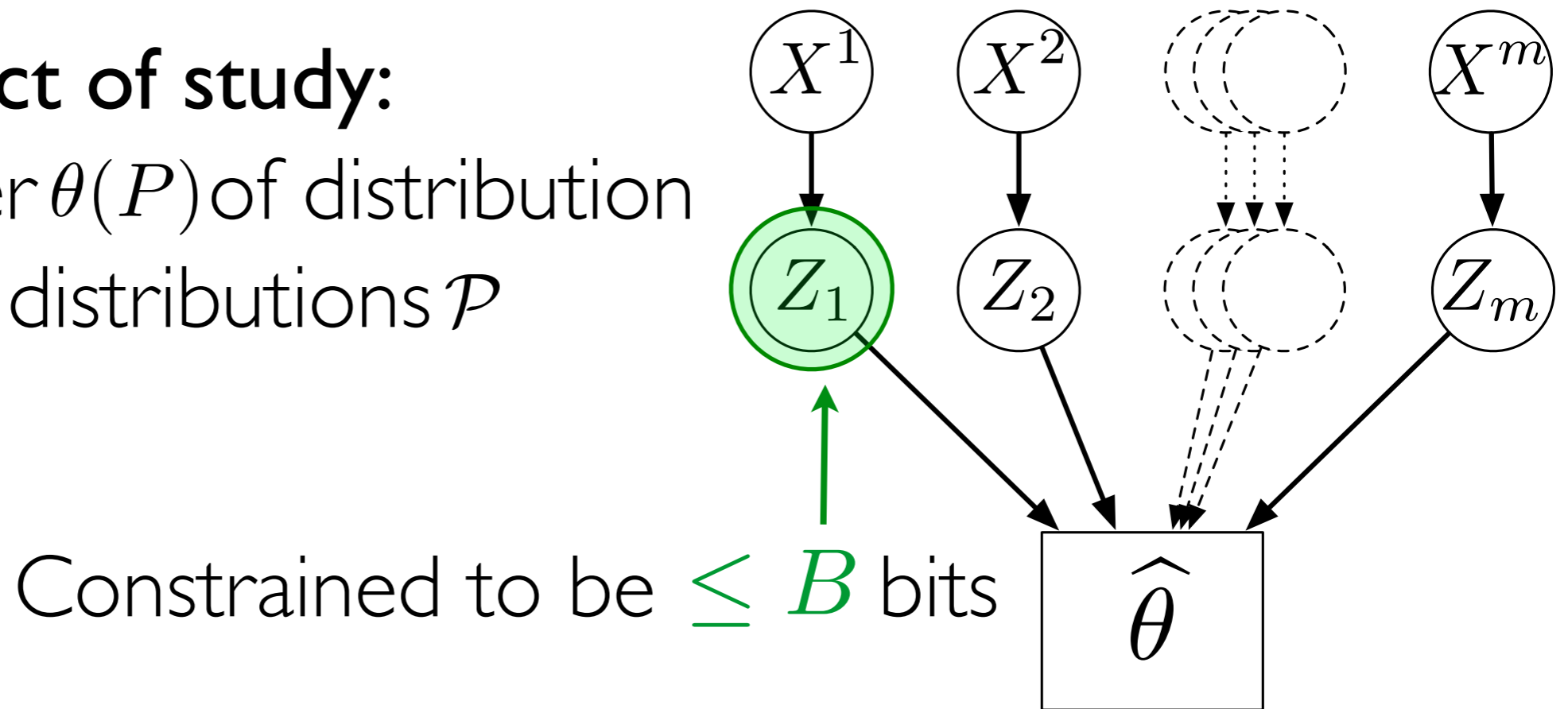
Minimax risk with B -bounded communication

$$\mathfrak{M}_n(\theta(\mathcal{P}), B) := \inf_{\pi \in \Pi_B} \inf_{\hat{\theta}} \sup_{P \in \mathcal{P}} \mathbb{E}_P \left[\|\hat{\theta}(Z_1^m) - \theta(P)\|_2^2 \right]$$

Minimax communication

Central object of study:

- Parameter $\theta(P)$ of distribution
- Family of distributions \mathcal{P}
- Loss $\|\cdot\|_2^2$



Minimax risk with B -bounded communication

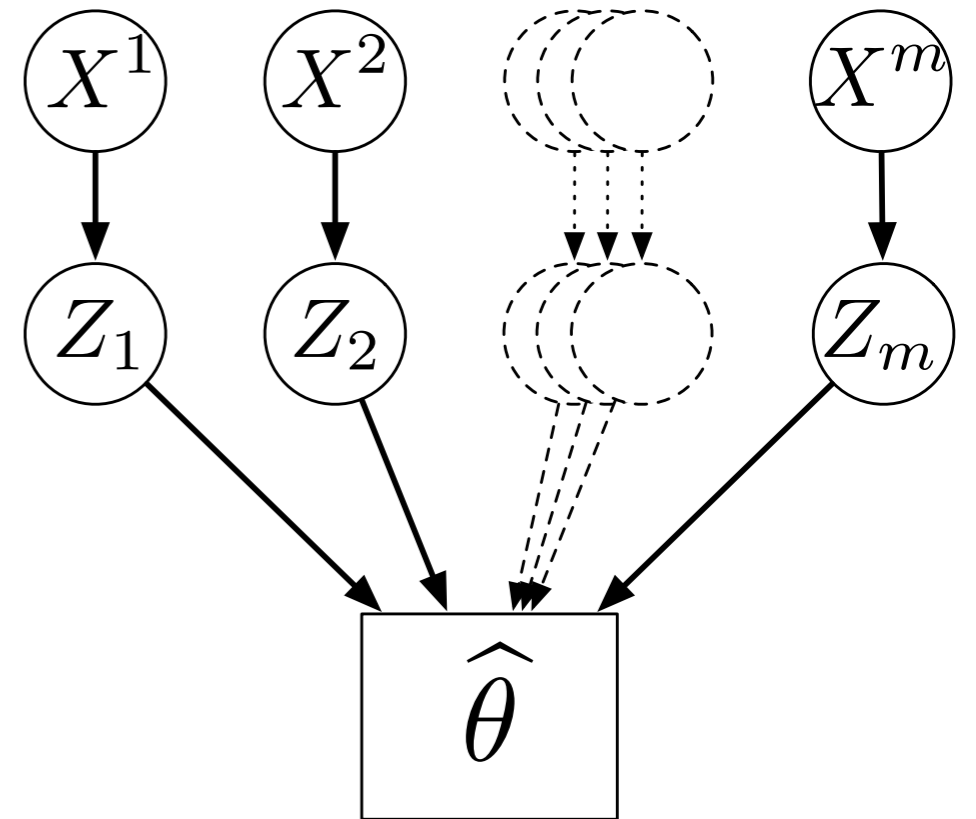
$$\mathfrak{M}_n(\theta(\mathcal{P}), B) := \inf_{\pi \in \Pi_B} \inf_{\hat{\theta}} \sup_{P \in \mathcal{P}} \mathbb{E}_P \left[\|\hat{\theta}(Z_1^m) - \theta(P)\|_2^2 \right]$$

Best protocol $Z_i = \pi(X^i)$ with Z_i smaller than B bits

Vignette: mean estimation

Consider estimation in normal location family, $\theta \in [-1, 1]^d$

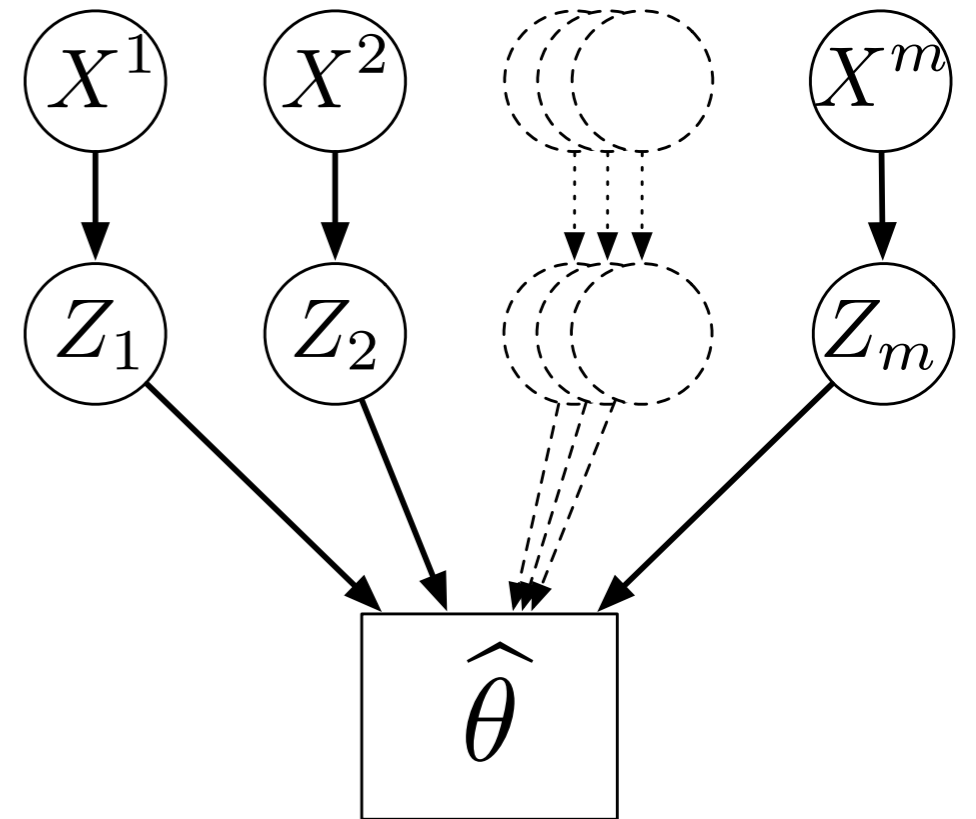
$$X_j^i \stackrel{\text{iid}}{\sim} \mathbf{N}(\theta, \sigma^2 I_{d \times d})$$



Vignette: mean estimation

Consider estimation in normal location family, $\theta \in [-1, 1]^d$

$$X_j^i \stackrel{\text{iid}}{\sim} \mathbf{N}(\theta, \sigma^2 I_{d \times d})$$



Theorem: when each agent has sample of size n

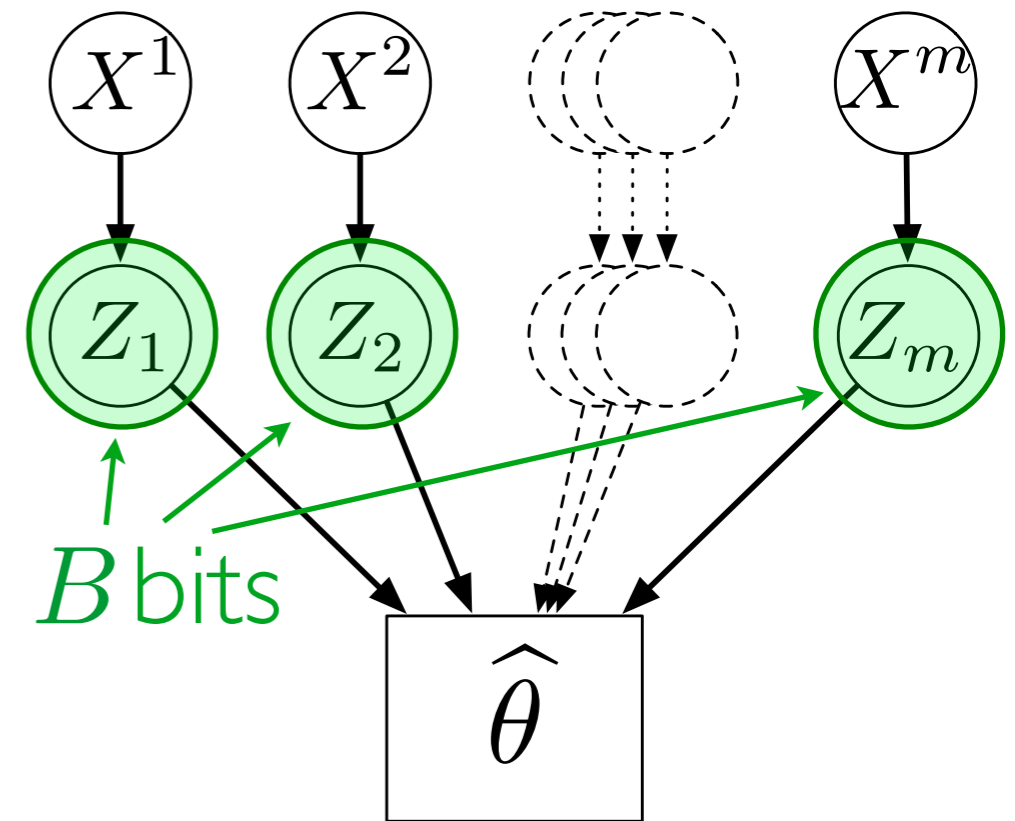
Minimax rate

$$\mathbb{E}[\|\hat{\theta}(X^1, \dots, X^m) - \theta\|_2^2] \asymp \frac{\sigma^2 d}{nm}$$

Vignette: mean estimation

Consider estimation in normal location family, $\theta \in [-1, 1]^d$

$$X_j^i \stackrel{\text{iid}}{\sim} \mathbf{N}(\theta, \sigma^2 I_{d \times d})$$



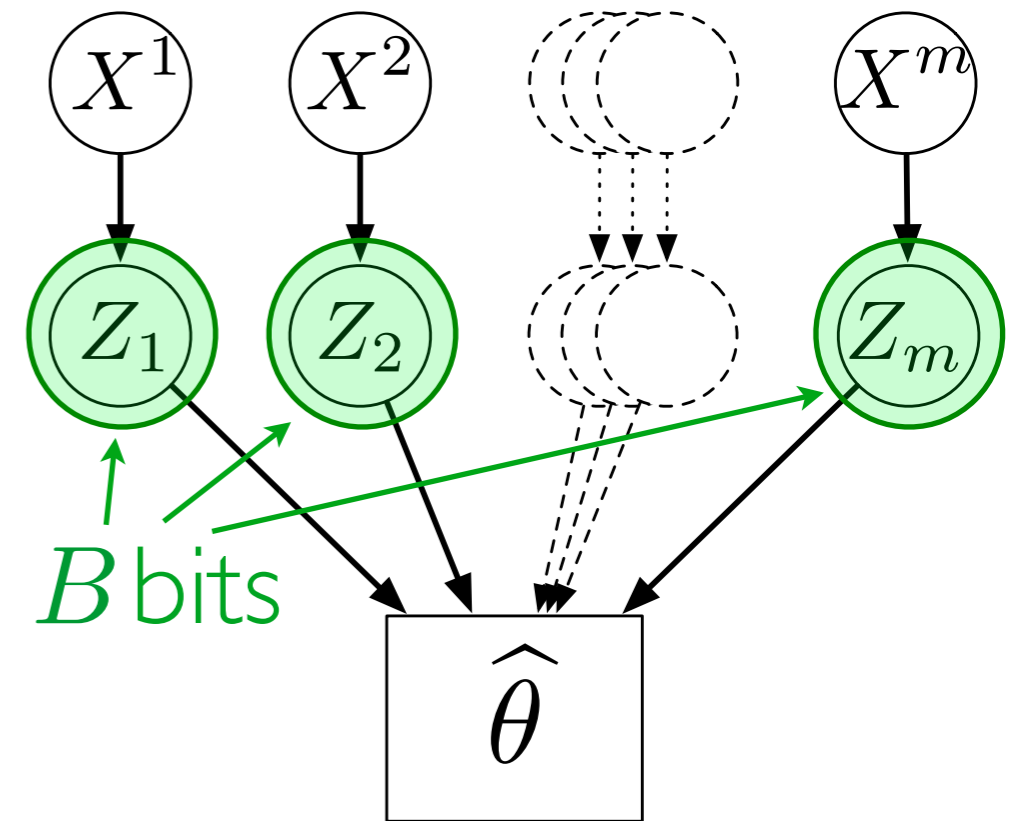
Theorem: when each agent has sample of size n
Minimax rate with B -bounded communication

$$\frac{d}{B \wedge d} \frac{1}{\log m} \frac{\sigma^2 d}{nm} \lesssim \mathfrak{M}_n(\mathcal{N}_d, B) \lesssim \frac{d \log m}{B \wedge d} \frac{\sigma^2 d}{nm}$$

Vignette: mean estimation

Consider estimation in normal location family, $\theta \in [-1, 1]^d$

$$X_j^i \stackrel{\text{iid}}{\sim} \mathbf{N}(\theta, \sigma^2 I_{d \times d})$$



Theorem: when each agent has sample of size n
Minimax rate with B -bounded communication

$$\frac{d}{B \wedge d} \frac{1}{\log m} \frac{\sigma^2 d}{nm} \lesssim \mathfrak{M}_n(\mathcal{N}_d, B) \lesssim \frac{d \log m}{B \wedge d} \frac{\sigma^2 d}{nm}$$

Consequence: each sends $\approx d$ bits for optimal estimation