

Lecture 14: Communication/privacy constrained estimation

Lecturer: Yanjun Han

May 12, 2021

Today's plan

constrained estimation:

- communication and privacy constraints
- distributed estimation and types of communication protocols
- tool I: strong data processing inequality
- tool II: van Trees inequality + quantized Fisher information
- tool III: direct modeling + Assouad's lemma

Constrained estimation

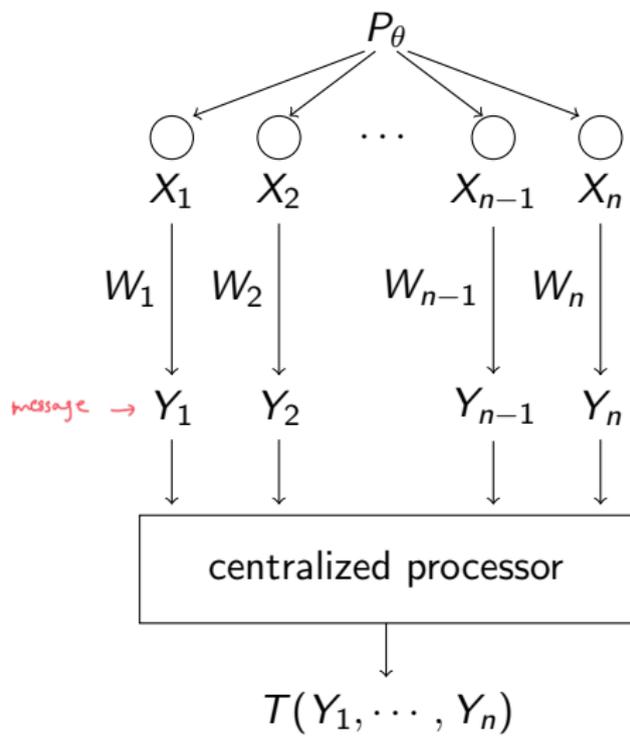
- usual statistical model: $X \sim P_\theta$
- however, the learner only has access to Y , which is the output of sending X via a channel $W \in \mathcal{W}$ up to learner's choice
- target: given loss $L(\theta, T)$, **jointly** design W and $T(Y)$
- communication-constrained channel \mathcal{W}_k :

$$W \in \mathcal{W}_k : \mathcal{X} \rightarrow \{0, 1\}^k$$

- privacy-constrained channel \mathcal{W}_ϵ :

$$W \in \mathcal{W}_\epsilon : \mathcal{X} \rightarrow \mathcal{Y} \quad \sup_{x, x' \in \mathcal{X}} \sup_{A \subseteq \mathcal{Y}} \frac{W(A | x)}{W(A | x')} \leq e^\epsilon.$$

Distributed estimation



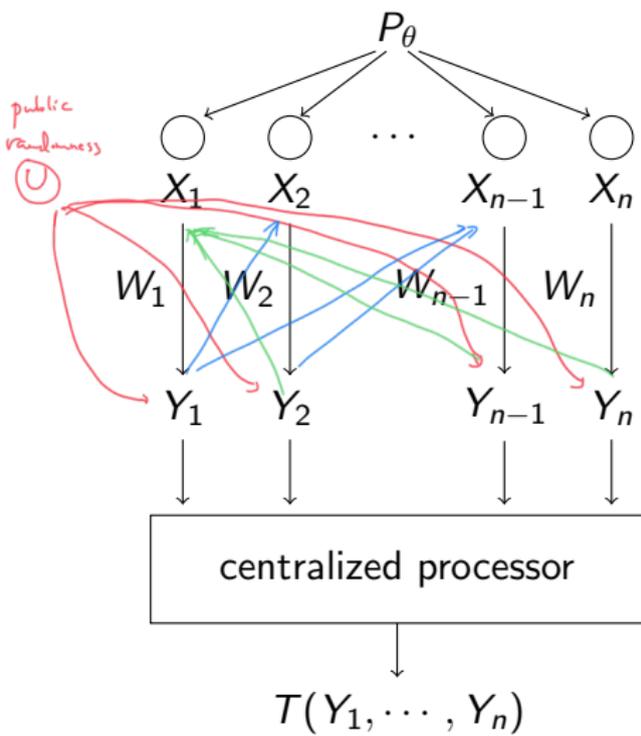
minimax risk:

$$\inf_T \inf_{W_1, \dots, W_n} \sup_{\theta} \mathbb{E}_{\theta} [L(\theta, T(Y^n))]$$

examples:

- distributed Gaussian mean estimation $P_\theta = \mathcal{N}(\theta, I_d)$
- distributed discrete distribution estimation $P_\theta = (\theta_1, \dots, \theta_d)$
- distributed uniformity testing

Communication protocols

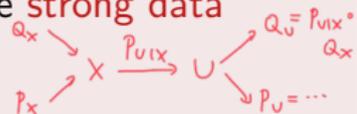


- simultaneous message passing (SMP) with private randomness
- SMP with public randomness
- sequential communication protocol
- blackboard (interactive) communication protocol

Tool I: strong data processing inequality

Definition

Given a f -divergence and $P_{U|X}$ (or in addition P_X), the **strong data processing coefficient** is defined as



$$\eta_f(P_{U|X}) \triangleq \sup_{P_X, Q_X} \frac{D_f(P_{U|X} \circ Q_X, P_{U|X} \circ P_X)}{D_f(Q_X, P_X)}$$

$$\eta_f(P_{U|X}, P_X) \triangleq \sup_{Q_X} \frac{D_f(P_{U|X} \circ Q_X, P_{U|X} \circ P_X)}{D_f(Q_X, P_X)}$$

- data-processing inequality: $\eta_f \leq 1$
- Polyanskiy and Wu (2017):

$$\eta_{\chi^2} = \eta_{\text{KL}} \leq \eta_f \leq \eta_{\text{TV}} = \sup_{x, x'} \|P_{U|X=x} - P_{U|X=x'}\|_{\text{TV}}$$

- Ordentlich and Polyanskiy (2021): suffice to consider binary (P_X, Q_X) for any $\eta_f(P_{U|X})$

The KL case

Lemma

$$\eta_{\text{KL}}(P_{U|X}) = \sup_{P_{XY}: U-X-Y} \frac{I(U; Y)}{I(X; Y)} = \frac{\mathbb{E}_Y[D_{\text{KL}}(P_{U|Y} \| P_U)]}{\mathbb{E}_Y[D_{\text{KL}}(P_{X|Y} \| P_X)]}$$
$$\eta_{\text{KL}}(P_{U|X}, P_X) = \sup_{P_{Y|X}: U-X-Y} \frac{I(U; Y)}{I(X; Y)} \leq \max_{Q_X} \frac{D_{\text{KL}}(P_{U|Y} \| P_U)}{D_{\text{KL}}(P_{X|Y} \| P_X)}$$

examples:

- $\mathcal{X} = \mathcal{U} = \{0, 1\}$ and $\mathbb{P}(U \neq x | X = x) = \varepsilon$:

$$\eta_{\text{KL}}(P_{U|X}) = (1 - 2\varepsilon)^2$$

- (X, U) joint Gaussian with correlation $\rho \in [-1, 1]$:

$$\eta_{\text{KL}}(P_{U|X}, P_X) = \rho^2$$

Why mutual information?

- upper bounds of $I(X; Y)$ available for both communication and privacy constrained estimation
- communication-constrained channel \mathcal{W}_k :

$$I(X; Y) \leq k$$

$$I(X; Y) \leq \sup_x D_{KL}(P_{Y|X=x} \parallel \text{Unif}(Y)) \leq k$$

- privacy-constrained channel \mathcal{W}_ϵ :

$$I(X; Y) \leq \min\{\epsilon, 2\epsilon^2\}$$

$$\leq \max_{x, x'} D_{KL}(P_{Y|X=x} \parallel P_{Y|X=x'})$$

$$e^{-\epsilon} \leq \frac{p(y)}{q(y)} \leq e^\epsilon$$

$$1) D_{KL}(p \parallel q) = \sum_x p(x) \log \frac{p(x)}{q(x)} \leq \epsilon$$

$$2) D_{KL}(p \parallel q) = \sum_x \left[p(x) \log \frac{p(x)}{q(x)} - p(x) + q(x) \right] \\ = \sum_x q(x) \left[\underbrace{\log \frac{p(x)}{q(x)} - \frac{p(x)}{q(x)} + 1}_{\leq O(\epsilon^2)} \right] \leq 2\epsilon^2$$

A simple example

- setting: $X_1, \dots, X_n \sim \text{Bern}(p)$ with $p \in [0, 1]$, and $W_i(Y_i | X_i) \in \mathcal{W}_\varepsilon$
- claim:

$$\inf_T \inf_{W_1, \dots, W_n} \sup_p \mathbb{E}_p[(T(Y^n) - p)^2] \gtrsim \frac{1}{n} + \frac{1}{n\varepsilon^2}$$

- analysis: let $U \sim \text{Unif}(\{0, 1\})$, and $P_{X_i|U=0} = \text{Bern}(1/2 - \delta)$,
 $P_{X_i|U=1} = \text{Bern}(1/2 + \delta)$
- strong data-processing inequality: $U - X^n - Y^n$

$$I(U; Y^n) \leq \underbrace{\eta_{\text{KL}}(P_{U|X^n}, P_{X^n})}_{\text{red underline}} \cdot I(X^n; Y^n)$$

- however, computing $\eta_{\text{KL}}(P_{U|X^n}, P_{X^n})$ is very complicated...

Tensorization under SMP

- key observation: under SMP, $(Y_1, \dots, Y_n) | U$ are independent
- upper bound of mutual information:



$$\begin{aligned} I(U; Y^n) &\leq \sum_{i=1}^n I(U; Y_i) \\ &\leq \sum_{i=1}^n \underbrace{\eta_{\text{KL}}(P_{U|X_i}, P_{X_i})}_{\text{KL divergence}} \cdot I(X_i; Y_i) \\ &\leq n \cdot \underbrace{(2\delta)^2}_{\text{KL bound}} \cdot 2\epsilon^2 \\ &= 8n\delta^2\epsilon^2 \end{aligned}$$

- two-point method: choosing $\delta^2 \asymp (n\epsilon^2)^{-1}$

Solution for the interactive protocols

Theorem (Braverman, Garg, Ma, Nguyen, Woodruff, 2016)

Suppose that $U \sim \text{Unif}(\{0, 1\})$ and $P_{X|U=0}(x)/P_{X|U=1}(x) \in [1/c, c]$ for some $c > 1$ and any $x \in \mathcal{X}$. Then if Y^n is obtained from X^n via any interactive communication protocol, we have

$$H^2(P_{Y^n|U=0}, P_{Y^n|U=1}) \lesssim_c \eta_{\text{KL}}(P_{U|X}, P_X) \cdot I(X^n; Y^n)$$

upper bound on $I(X^n; Y^n)$ under interactive protocols:

- communication constraints: $I(X^n; Y^n) \leq \log |\mathcal{Y}^n| = nk$
- privacy constraints:

$$I(X^n; Y^n) = \sum_t I(X^n; Y_t | Y^{t-1}) \leq n \min\{\varepsilon, 2\varepsilon^2\}$$

Product Bernoulli example

- setting: $X_1, \dots, X_n \sim \prod_{i=1}^d \text{Bern}(p_i)$, with $(p_1, \dots, p_d) \in [0, 1]^d$
- constraints: communication or privacy
- claim:

$$\inf_T \inf_{W_1, \dots, W_n \in \mathcal{W}} \sup_p \mathbb{E}_p [\|T(Y^n) - p\|_2^2] \gtrsim \begin{cases} \frac{d}{n} \cdot \frac{d}{\min\{k, d\}} & \text{if } \mathcal{W} = \mathcal{W}_k \\ \frac{d}{n} \cdot \frac{d}{\min\{\varepsilon^2, \varepsilon, d\}} & \text{if } \mathcal{W} = \mathcal{W}_\varepsilon \end{cases}$$

- analysis under SMP: apply Fano's method to $U \sim \text{Unif}(\{1/2 \pm \delta\}^d)$, but need to compute $\eta_{\text{KL}}(P_{U|X}, P_X)$
- analysis under interactive protocol: $\stackrel{d}{=} \sum_{i=1}^d P_{U_i|X_i}$ direct sum argument + previous data processing inequality

Tensorization property of SDPI

Theorem

$$\eta_{\text{KL}} \left(\prod_i P_{U_i|X_i}, \prod_i P_{X_i} \right) = \max_i \eta_{\text{KL}}(P_{U_i|X_i}, P_{X_i})$$

proof:

$$\begin{aligned}
 \frac{I(Y; U^d)}{I(Y; X^d)} &= \frac{\mathbb{E}_Y [D_{\text{KL}}(P_{U^d|Y} \| P_{U^d})]}{\mathbb{E}_Y [D_{\text{KL}}(P_{X^d|Y} \| P_{X^d})]} \\
 &= \frac{\sum_i \mathbb{E}_{(U^{i-1}, Y)} [D_{\text{KL}}(P_{U_i|U^{i-1}, Y} \| P_{U_i})]}{\sum_i \mathbb{E}_{(X^{i-1}, Y)} [D_{\text{KL}}(P_{X_i|X^{i-1}, Y} \| P_{X_i})]} \\
 &\leq \frac{\sum_i \mathbb{E}_{(U^{i-1}, Y)} \mathbb{E}_{X^{i-1}|U^{i-1}, Y} [D_{\text{KL}}(P_{U_i|X^{i-1}, Y} \| P_{U_i})]}{\sum_i \mathbb{E}_{(X^{i-1}, Y)} [D_{\text{KL}}(P_{X_i|X^{i-1}, Y} \| P_{X_i})]} \\
 &\leq \max_{i, y, x^{i-1}} \frac{D_{\text{KL}}(P_{U_i|X^{i-1}=x^{i-1}, Y=y} \| P_{U_i})}{D_{\text{KL}}(P_{X_i|X^{i-1}=x^{i-1}, Y=y} \| P_{X_i})} \\
 &\leq \max_i \eta_{\text{KL}}(P_{U_i|X_i}, P_{X_i})
 \end{aligned}$$

$P_{U_i|U^{i-1}} = P_{U_i}$
 $P_{U_i|U^{i-1}, Y} = \sum_{x^{i-1}} P_{U_i|X^{i-1}, Y} \times P_{X^{i-1}|U^{i-1}, Y}$
 $U_1 \rightarrow X_1 \rightarrow Y$
 $U_2 \rightarrow X_2 \rightarrow Y$
 $U_d \rightarrow X_d \rightarrow Y$

Analysis under interactive protocol

- Assouad's lemma under interactive protocol:

$$\sum_{j=1}^d \mathbb{P}(\hat{v}_j(Y^n) \neq v_j) \geq \frac{1}{2} \sum_{j=1}^d (1 - \|P_{+j} - P_{-j}\|_{\text{TV}})$$

where P_{+j}, P_{-j} are the distributions of the final message Y^n under $v_j = 1$ and $v_j = -1$, respectively

- idea of upper bounding $\|P_{+j} - P_{-j}\|_{\text{TV}}$: if P_θ has a product structure, $= P_{\theta_1} \times P_{\theta_2} \times \dots \times P_{\theta_d}$ then every node could discard all but the j -th coordinate of X
- consequence: reduce to 1-dimensional problem
- so-called “direct-sum” result in literature

Gaussian location model

- setting: $X_1, \dots, X_n \sim \mathcal{N}(\theta, I_d)$, with $\theta \in \mathbb{R}^d$
- constraints: communication or privacy
- result:

$$\inf_T \inf_{W_1, \dots, W_n \in \mathcal{W}} \sup_{\theta} \mathbb{E}_{\theta} [\|T(Y^n) - \theta\|_2^2] \gtrsim \begin{cases} \frac{d}{n} \cdot \frac{d}{\min\{k, d\}} & \text{if } \mathcal{W} = \mathcal{W}_k \\ \frac{d}{n} \cdot \frac{d}{\min\{\varepsilon^2, \varepsilon, d\}} & \text{if } \mathcal{W} = \mathcal{W}_{\varepsilon} \end{cases}$$

- analysis exactly the same, with SDPI specialized to Gaussian model

Tool II: quantized Fisher information

- recall the definition of Fisher information matrix:

$$I_X(\theta) = \mathbb{E}_{X \sim P_\theta} \left[\left(\frac{\partial \log p(X | \theta)}{\partial \theta} \right) \left(\frac{\partial \log p(X | \theta)}{\partial \theta} \right)^\top \right]$$

- quantized Fisher information for Y :

$$I_Y(\theta) = \mathbb{E}_{Y \sim W \circ P_\theta} \left[\left(\frac{\partial \log p(Y | \theta)}{\partial \theta} \right) \left(\frac{\partial \log p(Y | \theta)}{\partial \theta} \right)^\top \right]$$

- relationship between $p(Y | \theta)$ and $p(X | \theta)$:

$$p(y | \theta) = \sum_{x \in \mathcal{X}} W(y | x) p(x | \theta)$$

Lemma

$$I_Y(\theta) = \mathbb{E}_Y \left[\left(\mathbb{E}_\theta \left[\overset{\frac{\partial}{\partial \theta} \log p(x|\theta)}{\parallel} \dot{\ell}_\theta(X) \mid Y \right] \right) \left(\mathbb{E}_\theta \left[\dot{\ell}_\theta(X) \mid Y \right] \right)^\top \right]$$

van Trees inequality

Theorem (van Trees inequality, a.k.a. Bayesian Cramér–Rao bound)

For a differentiable prior density π on $[a, b]$ with $\pi(a) = \pi(b) = 0$, it holds that

$$\mathbb{E}_{\theta \sim \pi} \mathbb{E}_{X \sim P_\theta} [(T(X) - \theta)^2] \geq \frac{1}{\mathbb{E}_{\theta \sim \pi} [I_X(\theta)] + I(\pi)},$$

where $I(\pi) = \int_a^b [\pi'(t)]^2 / \pi(t) dt$ is the Fisher information of π .

Proof: consider two expressions of the quantity

$$S = \int_a^b \int_{\mathcal{X}} (T(x) - \theta) \frac{\partial [p(x | \theta) \pi(\theta)]}{\partial \theta} dx d\theta = 1$$
$$S^2 \leq \left(\int_a^b \int_{\mathcal{X}} (T(x) - \theta)^2 \pi(\theta) p(x|\theta) dx d\theta \right) \underbrace{\left(\int_a^b \int_{\mathcal{X}} \left(\frac{\partial [p(x|\theta)\pi(\theta)]}{\partial \theta} \right)^2 / \pi(\theta)p(x|\theta) dx d\theta \right)}_{= \mathbb{E}_{\theta \sim \pi} [I(\theta)] + I(\pi)}$$

Corollaries and generalizations

- Corollary 1: for 1-D model,

$$\sup_{\theta \in [a,b]} \mathbb{E}_{X \sim P_\theta} [(T(X) - \theta)^2] \geq \frac{1}{\sup_{\theta \in [a,b]} I_X(\theta) + \frac{4\pi^2}{(b-a)^2}}$$

$\pi(x) \propto \sin^2\left(\frac{x-a}{b-a}\pi\right)$

- Corollary 2: for high-dimensional model,

$$\sup_{\theta \in [a,b]^d} \mathbb{E}_{X \sim P_\theta} [\|T(X) - \theta\|_2^2] \geq \frac{d^2}{\sup_{\theta \in [a,b]^d} \text{Tr}(I_X(\theta)) + 4d\pi^2/(b-a)^2}$$

- application in constrained estimation: upper bound

$$\sup_{\theta \in [a,b]^d} \sup_{W \in \mathcal{W}} \text{Tr}(I_{Y^n}(\theta))$$

Tensorization property

Even under an interactive protocol, we have

$$\text{Tr}(I_{Y^n}(\theta)) \leq n \cdot \sup_{W \in \mathcal{W}} \text{Tr}(I_Y(\theta))$$

Upper bounds of individual Fisher information

Upper bound I (bounded variance)

If $\dot{\ell}_\theta(x)$ has a bounded second moment σ^2 along any direction under $x \sim P_\theta$, then

$$\sup_{W \in \mathcal{W}} \text{Tr}(I_Y(\theta)) \leq \begin{cases} 2^k \sigma^2 & \text{if } \mathcal{W} = \mathcal{W}_k, \\ \min\{e^\varepsilon, \varepsilon^2\} \sigma^2 & \text{if } \mathcal{W} = \mathcal{W}_\varepsilon. \end{cases}$$

Upper bound II (sub-Gaussian)

If $\dot{\ell}_\theta(x)$ is subGaussian with parameter σ^2 under $x \sim P_\theta$, then

$$\sup_{W \in \mathcal{W}} \text{Tr}(I_Y(\theta)) \lesssim \begin{cases} k \sigma^2 & \text{if } \mathcal{W} = \mathcal{W}_k, \\ \min\{\varepsilon, \varepsilon^2\} \sigma^2 & \text{if } \mathcal{W} = \mathcal{W}_\varepsilon. \end{cases}$$

Example I: Gaussian location model

- van Trees inequality:

$$R_{n,d}^* \geq \frac{d^2}{n \cdot \sup_{\theta \in [a,b]^d} \sup_{W \in \mathcal{W}} \text{Tr}(I_Y(\theta)) + 4d\pi^2/(b-a)^2}$$

- score function $\dot{\ell}_\theta(x) = \theta - x$ is 1-sub-Gaussian
- communication constraint: $\text{Tr}(I_Y(\theta)) \lesssim k$, giving

$$R_{n,d,k}^* \gtrsim \frac{d}{n} \cdot \frac{d}{\min\{d, k\}}$$

- privacy constraint: $\text{Tr}(I_Y(\theta)) \lesssim \min\{\varepsilon, \varepsilon^2\}$, giving

$$R_{n,d,\varepsilon}^* \gtrsim \frac{d}{n} \cdot \frac{d}{\min\{d, \varepsilon, \varepsilon^2\}}$$

Example II: discrete distribution model

- setting: $X_1, \dots, X_n \sim \theta = (\theta_1, \dots, \theta_d)$, which is a probability vector
- van Trees inequality:

$$R_{n,d}^* \geq \frac{d^2}{n \cdot \sup_{\theta \in [a,b]^d} \sup_{W \in \mathcal{W}} \text{Tr}(I_Y(\theta)) + 4d\pi^2/(b-a)^2}$$

- score function only finite second moment d
- communication constraint: $\text{Tr}(I_Y(\theta)) \lesssim d2^k$, giving

$$R_{n,d,k}^* \gtrsim \frac{1}{n} \cdot \frac{d}{\min\{d, \underline{2^k}\}}$$

- privacy constraint: $\text{Tr}(I_Y(\theta)) \lesssim \min\{e^\epsilon, \epsilon^2\}d$, giving

$$R_{n,d,k}^* \gtrsim \frac{1}{n} \cdot \frac{d}{\min\{d, \underline{\epsilon^2}, \underline{e^\epsilon}\}}$$

Tool III: direct modeling + Assouad's lemma

- setting: $X_1, \dots, X_n \sim p = (p_1, \dots, p_d)$
- target: uniformity testing $p = \text{unif}_d$ vs. $\|p - \text{unif}_d\|_{\text{TV}} \geq \delta$
- result: table of sample complexities

	\mathcal{W}_k	\mathcal{W}_ϵ
SMP w/ <u>public coin</u> (this lecture)	$\frac{\sqrt{d}}{\delta^2} \cdot \sqrt{\frac{d}{\min\{d, 2^k\}}}$	$\frac{\sqrt{d}}{\delta^2} \cdot \frac{\sqrt{d}}{\min\{\sqrt{d}, \epsilon^2\}}$
SMP w/ <u>private coin</u> (future lecture)	$\frac{\sqrt{d}}{\delta^2} \cdot \frac{d}{\min\{d, 2^k\}}$	$\frac{\sqrt{d}}{\delta^2} \cdot \frac{d}{\min\{d, \epsilon^2\}}$

Direct modeling the channel

- assume $\mathcal{W} = \mathcal{W}_k$
- Paninski's construction: for $v \in \{\pm 1\}^{d/2}$, let

$$X \sim P_v = \left(\frac{1 + v_1 \delta}{d}, \frac{1 - v_1 \delta}{d}, \dots, \frac{1 + v_{d/2} \delta}{d}, \frac{1 - v_{d/2} \delta}{d} \right)$$

- distribution of $Y \sim Q_v$: for $y \in \{0, 1\}^k$,

$$Q_v(y) = \sum_{i \in [d/2]} \left(\frac{1 + v_i \delta}{d} W(y | 2i - 1) + \frac{1 - v_i \delta}{d} W(y | 2i) \right)$$

References: strong data processing inequality

- Maxim Raginsky. “Strong data processing inequalities and Φ -Sobolev inequalities for discrete channels.” *IEEE Transactions on Information Theory* 62.6 (2016): 3355-3389.
- Yury Polyanskiy and Yihong Wu. “Strong data-processing inequalities for channels and Bayesian networks.” *Convexity and Concentration*. Springer, New York, NY, 2017. 211-249.
- John C. Duchi, Michael I. Jordan, and Martin J. Wainwright. “Local privacy and statistical minimax rates.” 2013 IEEE 54th Annual Symposium on Foundations of Computer Science. IEEE, 2013.
- Mark Braverman, Ankit Garg, Tengyu Ma, Huy L. Nguyen, and David P. Woodruff. “Communication lower bounds for statistical estimation problems via a distributed data processing inequality.” In *Proceedings of the forty-eighth annual ACM symposium on Theory of Computing*, pp. 1011-1020. 2016.

References: more recent approaches

- Yanjun Han, Ayfer Özgür, and Tsachy Weissman. “Geometric lower bounds for distributed parameter estimation under communication constraints.” *Conference On Learning Theory*. PMLR, 2018.
- Leighton Pate Barnes, Yanjun Han, and Ayfer Özgür. “Lower bounds for learning distributions under communication constraints via Fisher information.” *Journal of Machine Learning Research* 21.236 (2020): 1-30.
- Jayadev Acharya, Clément L. Canonne, and Himanshu Tyagi. “Inference under information constraints I: Lower bounds from chi-square contraction.” *IEEE Transactions on Information Theory* 66.12 (2020): 7835-7855.
- FOCS'20 tutorial:
<http://www.cs.columbia.edu/~ccanonne/tutorial-focs2020/>

Next lecture: scandiction (Tsachy)