

# Statistical Inference in Distributed and Constrained Settings

## Part 2: Unified Lower Bounds

What will we see?

- More general lower bounds that apply to high-dimensional models
- Bounds that allow interaction

The model:

(2)

- $\mathcal{P}_{\mathbb{H}} = \{P_{\theta}, \theta \in \Theta\}$

a parametric family of distributions on  $\mathcal{X}$

- $X_1, \dots, X_n$

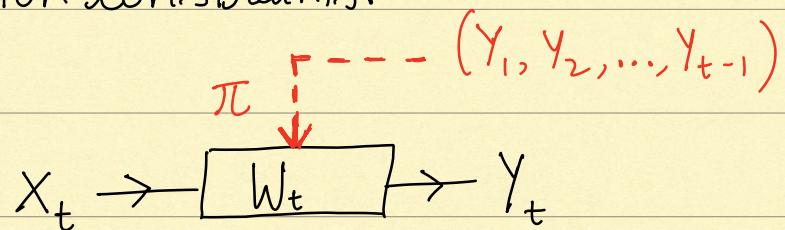
i.i.d. samples from an unknown

$$P_{\theta} \in \mathcal{P}_{\mathbb{H}}$$

- $Y_1, \dots, Y_n$

(3)

Information constraints:



*sequentially interactive*

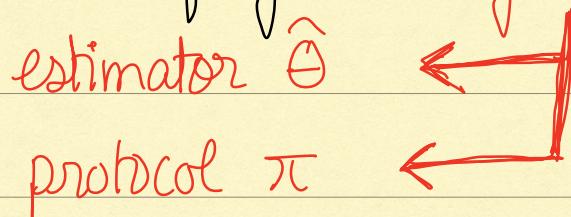
- $\hat{\theta}(Y_1, \dots, Y_n)$

Estimate for unknown  $\theta$

- $\max_{\theta \in \Theta} E_{\theta} [\|\theta - \hat{\theta}\|_p]$

worst-case loss to capture the

performance of your algorithm



(4)

How small can we make

$$\underset{\theta}{\operatorname{min}} \mathbb{E}_{\theta} [\|\theta - \hat{\theta}\|_p]$$

are allowed to use any algorithm  
in the world?

Heuristics from information theory:

(1) Uncertainty about the unknown  $\theta$

decreases as we observe  $y_1, \dots, y_n$

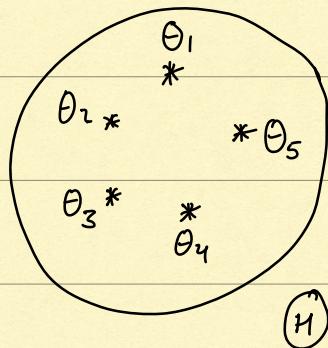
(2) This uncertainty can be quantified  
using measures of information

(3) Information is usually additive in  $n$

(4) The presence of  $w_t$ 's will reduce the  
information (data processing)

(5)

## (1) Quantifying uncertainty

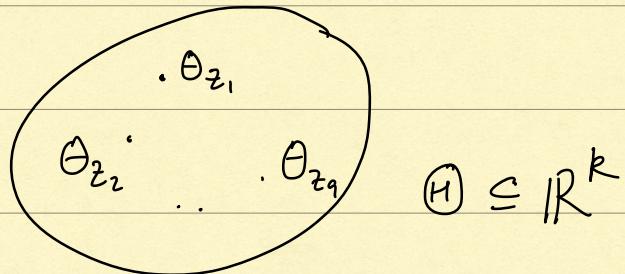


- Choose a finite set of  $M$  points in the parameter space
- Let one of them be chosen randomly and  $n$  samples be generated from it
- Uncertainty remaining upon observing  $y^n$ 
  - Two points :  $1 - d(P_{\theta_1}^{y^n}, P_{\theta_2}^{y^n}) \rightsquigarrow$  Le Cam
  - Multiple points :  $H(z|y^n) \rightsquigarrow$  Fano

(6)

→ Average uncertainty for each coordinate  
~~ Assouad's method

(Suited for high-dimensional  
+  
interactive)



Embedding:

(7)

$$|\theta_{z,i} - \theta_{z',i}| \approx r \mathbb{1}_{\{z_i \neq z'_i\}}$$

$$\|\theta_z - \theta_{z'}\|^2 \approx r^2 \sum_{i=1}^k \mathbb{1}_{\{z_i \neq z'_i\}}$$

(8)

$$\hat{z} = \underset{z \in \{-1, 1\}^k}{\operatorname{argmin}} \|\hat{\theta} - \theta_z\|$$

$$\Rightarrow E_{\theta_z} [\|\hat{\theta} - \theta_z\|^2]$$

$$\geq \frac{1}{2} E_{\theta_z} [\|\hat{\theta}_z - \theta_z\|^2]$$

$$\approx r^2 \sum_{i=1}^k P_{\theta_z} [z_i \neq \hat{z}]$$

For  $z \sim \text{unif } \{-1, 1\}^k$ , Recall:  $(z \rightarrow \theta_z \rightarrow y^n \rightarrow \hat{\theta} \rightarrow \hat{z})$

$$E_z [E_{\theta_z} [\|\hat{\theta} - \theta_z\|^2]] \quad (9)$$

$$\gtrsim r^2 \sum_{i=1}^k P[z_i \neq \hat{z}_i]$$

$P(Z_i \neq \hat{Z}_i) \rightsquigarrow$  probability of error  
for a binary  
hypothesis testing prob.

$$P_{+i}^{y^n} = p(y^n | Z_i = +)$$

$$P_{-i}^{y^n} = p(y^n | Z_i = -)$$

$$P(Z_i \neq \hat{Z}_i) \geq \frac{1}{2} (1 - d(P_{+i}^{y^n}, P_{-i}^{y^n}))$$

$$\mathbb{E}_Z [\mathbb{E}_{\theta_2} [||\hat{\theta} - \theta_2||^2]]$$

(10)

$$\gtrsim r^2 \sum_{i=1}^k (1 - d(P_{+i}^{y^n}, P_{-i}^{y^n}))$$

Thus,

loss  $\geq$  "uncertainty"

$$\approx r^2 k \left( 1 - \frac{1}{k} \sum_{i=1}^k d(P_{+i}^{y^n}, P_{-i}^{y^n}) \right)$$

## (2) Uncertainty to information

11

$$\begin{aligned} & \left( \frac{1}{k} \sum_{i=1}^k d(p_{+i}^{y^n}, p_{-i}^{y^n}) \right)^2 \\ & \leq \frac{1}{k} \sum_{i=1}^k d(p_{+i}^{y^n}, p_{-i}^{y^n})^2 \end{aligned}$$

Lemma.  $p_+, p_-$  two distributions on  $\mathcal{X}$

$$U \sim \text{unif } \{-1, +1\}$$

$$P_{X|U=+1} = p_+$$

$$P_{X|U=-1} = p_-$$

Then,

$$d(p_+, p_-)^2 \leq 2I(U \wedge X)$$

Proof.  $d(p_+, p_-)^2 \leq 2(d(p_+, q)^2 + d(p_-, q)^2)$

  $\leq D(p_+ \| q) + D(p_- \| q) = 2I(U \wedge X)$  ■

(12)

Lemma  $\Rightarrow$

$$d(p_{+i}^{Y^n}, p_{-i}^{Y^n})^2 \leq 2 I(Z_i \wedge Y^n)$$

=

$$\leq 2 I(Z_i \wedge Y^n | Z^{-i})$$

(independence of  $Z_1, \dots, Z_k$ )  $\leq 2 I(Z_i \wedge Y^n | Z^{-i})$

$$\left( \frac{1}{k} \sum_{i=1}^k d(p_{+i}^{Y^n}, p_{-i}^{Y^n}) \right)^2$$

$$\leq \frac{1}{k} \sum_{i=1}^k I(Z_i \wedge Y^n | Z^{-i}).$$

(13)

(1), (2)  $\Rightarrow$ 

$$\mathbb{E}_z \left[ \mathbb{E}_{P_{\theta_2}} [\|\theta_2 - \hat{\theta}\|^2] \right] \quad \text{"loss"} \\ \gtrsim r^2 k \left( 1 - \sqrt{ \frac{2}{k} \sum_{i=1}^k I(z_i \wedge y^n | z^{-i}) } \right) \\ \text{"information"}$$

need to derive an upper bound for this information quantity

(3) Tensorization of information

(14)

$$I(z_i \wedge y^n | z^{-i}) = \sum_{t=1}^n I(z_i \wedge y_t | z^{-i}, y^{t-1})$$

(4) How  $W_t$  s reduce information

Bounding  $I(Z_i \wedge Y_t | Z^{-i} = z^{-i}, Y^{t-1} = y^{t-1})$

Simplifying notation:

(15)

$$W_t = W^{y^{t-1}}$$

$$Z^{+i} = (z_1, \dots, z_{i-1}, 1, z_{i+1}, \dots, z_k)$$

$$Z^{-i} = (z_1, \dots, z_{i-1}, -1, z_{i+1}, \dots, z_k)$$

$$P_{Z^{+i}}^{W_t} = P(Y_t | Z_i = 1, Z^{-i}, y^{t-1})$$

$$X_t \sim P_{Z^{+i}} \rightarrow [W_t] \rightarrow Y_t$$

$$P_{Z^{-i}}^{W_t} = P(Y_t | Z_i = -1, Z^{-i}, y^{t-1})$$

$$X_t \sim P_{Z^{-i}} \rightarrow [W_t] \rightarrow Y_t$$

$I(Z_i \wedge Y_t | Z^{-i} = z^{-i}, Y^{t-1} = y^{t-1})$

(16)

$$= \frac{1}{2} D\left(P_{Z^{+i}}^{W_t} \parallel \frac{1}{2} P_{Z^{+i}}^{W_t} + \frac{1}{2} P_{Z^{-i}}^{W_t}\right)$$

$$+ \frac{1}{2} D\left(P_{Z^{-i}}^{W_t} \parallel \frac{1}{2} P_{Z^{+i}}^{W_t} + \frac{1}{2} P_{Z^{-i}}^{W_t}\right)$$

Conditioning on  $Z^{-i}$  helped - with this  
 the distribution of  $X_t$  doesn't depend on  $y^{t-1}$ :  
 $X_t - (Z_i, Z^{-i}) - y^{t-1}$

$$\Rightarrow I(Z_i \wedge Y_t | Z^{i-1}, y^{t-1}) \\ = \underset{Z}{\mathbb{E}} \left[ \underset{y^{t-1}}{\mathbb{E}} \left[ D(P_2^W || \frac{1}{2} P_2^W + \frac{1}{2} P_{Z \setminus i}^W) \right] \right]$$

- $p, q$  two distributions on  $\mathcal{X}$  (17)
- $W: \mathcal{X} \rightarrow \mathcal{Y}$

$$D(p^W || \frac{1}{2} p^W + \frac{1}{2} q^W)$$

$$\leq d_{\mathcal{X}^2} \left( p^W || \frac{1}{2} p^W + \frac{1}{2} q^W \right)$$

$$= \int \frac{\left( \frac{1}{2} p^W(y) - \frac{1}{2} q^W(y) \right)^2}{\frac{1}{2}(p^W(y) + q^W(y))} d\mu$$

$$\leq \frac{1}{2} \int \frac{(p^w(y) - q^w(y))^2}{p^w(y)} d\mu$$

$\stackrel{\text{def}}{=} \phi(x)$

$$p^w(y) - q^w(y) = E_p \left[ W(y|x) \left( \frac{q(x)}{p(x)} - 1 \right) \right]$$

$$= \frac{1}{2} \int \frac{E_p \left[ W(y|x) \phi(x) \right]^2}{E_p [W(y|x)]} d\mu$$

(18)

$$D(p^w || \frac{1}{2} p^w + \frac{1}{2} q^w) \leq$$

$$\frac{1}{2} \int \frac{E_p \left[ \phi(x) W(y|x) \right]^2}{E_p [W(y|x)]} d\mu$$

where  $\phi(x) = \frac{q(x)}{p(x)} - 1$ .

$$\Rightarrow \sum_{i=1}^k \sum_{t=1}^n I(z_i \wedge y_t) z^{-i}, y^{t-1}) \quad (19)$$

$$= \sum_{i=1}^k \sum_{t=1}^n \mathbb{E}_{z, y^{t-1}} \left[ D(P_z^{W_t} \parallel \underbrace{\frac{1}{2} P_z^{W_t} + \frac{1}{2} P_z^{W_t \oplus i}}_{}) \right]$$

$$\leq \frac{1}{2} \int \frac{\mathbb{E}_{P_z} [\phi_{z,i}(x) w_t(y|x)]^2 d\mu}{\mathbb{E}_{P_z} [w_t(y|x)]}$$

where  $\phi_{z,i}(x) = \frac{dP_z \oplus i}{dP_z}(x) - 1$ .

$$\leq \frac{1}{2} \cdot n \cdot \max_z \max_{W \in \mathcal{W}}$$

$$\sum_{i=1}^k \int \frac{\mathbb{E}_{P_z} [\phi_{z,i}(x) w_t(y|x)]^2 d\mu}{\mathbb{E}_{P_z} [w_t(y|x)]}$$

(20)

Theorem (Average information  
contraction bound)

$$|\theta_{z,i} - \theta_{z',i}| \approx r \mathbb{I}_{\{z_i \neq z'_i\}}$$

$$\frac{d p_{z^{\oplus i}}}{d p_z} = 1 + \phi_{z,i}^r$$

$$\mathbb{E}_z \mathbb{E}_{P_{\theta_2}} \left[ \|\theta_z - \hat{\theta}\|_2^2 \right]$$

$$\geq r^2 k \left( 1 - \sqrt{\frac{n}{2} \max_z \max_w \frac{1}{k} \sum_{i=1}^k \frac{\mathbb{E}_{P_z} [\phi_{z,i}^r(x) w(y|x)]^2}{\mathbb{E}_{P_z}[w(y|x)]}} \right)$$

Further bounds for

(21)

$$\frac{\sum_{i=1}^k \mathbb{E}_{P_z} [W(y|X) \phi_{z,i}(x)]^2}{\mathbb{E}_{P_z} [W(y|X)]}$$

Additional conditions:

(22)

- $\mathbb{E}_z [\phi_{z,i}^2] \leq \alpha^2$

- $\phi_{z,1}, \dots, \phi_{z,k}$  are orthogonal

Then, since  $\phi_{z,i}$  are zero-mean under  $P_z$ ,

$$\sum_{i=1}^k \mathbb{E}_{P_z} [a(x) \phi_{z,i}(x)]^2 = \sum_{i=1}^k \langle a, \phi_{z,i} \rangle^2$$

$$= \sum_{i=1}^k \langle a - \mathbb{E}[a], \phi_{z,i} \rangle^2 \stackrel{?}{\leq} \|a - \mathbb{E}[a]\|^2$$

$$\Rightarrow \frac{\sum_{i=1}^k \frac{E_{P_2}[W(y|X)\phi_{z,i}]^2}{E_{P_2}[W(y|X)]}}{\frac{\text{Var}_{P_2}(W(y|X))}{E_{P_2}[W(y|X)]}}$$

Additional conditions:

(23)

- $E_z[\phi_{z,i}^2] \leq \alpha^2$

- $\phi_{z,1}, \dots, \phi_{z,k}$  are independent and  $\sigma^2$ -subgaussian under  $P_z$

Under these assumptions, Gibbs variational formula can be used to show that

$$\begin{aligned} & \frac{\sum_{i=1}^k \frac{E_{P_2}[\phi_{z,i}(X)W(y|X)]^2}{E_{P_2}[W(y|X)]}}{\frac{\text{Var}_{P_2}(W(y|X))}{E_{P_2}[W(y|X)]}} \\ & \leq 2\alpha^2 \sigma^2 E_{P_2}\left[W(y|X) \ln \frac{W(y|X)}{E_{P_2}[W(y|X)]}\right] \end{aligned}$$

$$\mathbb{E}_{P_Z} W(y|X)$$

$$\Rightarrow \sum_{i=1}^k \frac{\int \mathbb{E}_{P_Z} [\phi_{z,i}(x) W(y|x)]^2 d\mu}{\mathbb{E}_{P_Z} [W(y|x)]} \leq 2\alpha^2 \sigma^2 I(P_Z; W)$$

Summary:

(24)

$$\mathbb{E}_Z \mathbb{E}_{P_{\theta_2}} \left[ \|\theta_2 - \hat{\theta}\|_2^2 \right]$$

$$\geq r^2 k \left( 1 - \sqrt{\frac{n}{2} \max_z \max_W \frac{1}{k} \sum_{i=1}^k \frac{\int \mathbb{E}_{P_Z} [\phi_{z,i}^r(x) W(y|x)]^2 d\mu}{\mathbb{E}_{P_Z} [W(y|x)]}} \right)$$

Assumptions:

(i)  $\mathbb{E} [\phi_{z,i}^r]^2 \leq \alpha^2$

(ii)  $\phi_{z,1}, \dots, \phi_{z,k}$  orthonormal

(iii)  $\phi_{z,1}, \dots, \phi_{z,k}$  independent and  $\sigma^2$ -subgaussian

(25)

$$(i), (ii) \Rightarrow \geq r^2 k \left( 1 - \sqrt{n \alpha^2 \max (Var_{P_Z} (W(y|x)))_{d_{11}}} \right)$$

$$| \sqrt{2R} z, w | \overline{\mathbb{E}_{P_z}[w(y|x)]}$$

$$(i), (iii) \Rightarrow \geq r^2 R \left( 1 - \sqrt{\frac{n\alpha^2}{2R}} \max_{z, w} I(P_z; W) \right)$$

So far we have made heuristics (1), (2), (3)  
concrete

#### (4) Information reduces due to $W$

- Communication constraints:

$$I(P_z; W) \leq l \text{ and } \sum_y \frac{\text{Var}(W(y|x))}{\mathbb{E}[W(y|x)]} \leq 2^l$$

$$\text{Privacy constraints: } \sum_y \frac{\text{Var}(W(y|x))}{\mathbb{E}[W(y|x)]} \leq (e^\rho - 1)^2$$