

Mean Estimation from One-bit Measurements

Alon Kipnis (Stanford)
John Duchi (Stanford)

Allerton
October 2017

Table of Contents

Introduction

Motivation

Preliminary

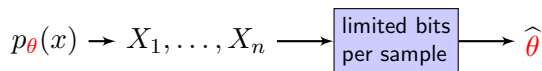
Adaptive Encoding

Distributed Encoding

Summary

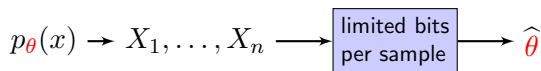
Motivation

Point estimation under communication constraints:



Motivation

Point estimation under communication constraints:

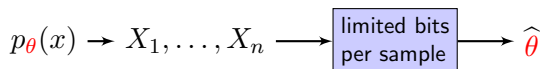


Estimation error is due to:

- (i) limited data
- (ii) limited bits

Motivation

Point estimation under communication constraints:



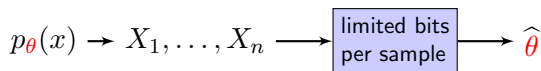
Estimation error is due to:

- (i) limited data
- (ii) limited bits

Relevant scenarios:

Motivation

Point estimation under communication constraints:



Estimation error is due to:

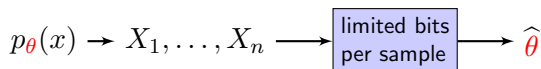
- (i) limited data
- (ii) limited bits

Relevant scenarios:

- ▶ big data

Motivation

Point estimation under communication constraints:



Estimation error is due to:

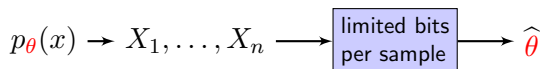
- (i) limited data
- (ii) limited bits

Relevant scenarios:

- ▶ big data
- ▶ low-power sensors

Motivation

Point estimation under communication constraints:



Estimation error is due to:

- (i) limited data
- (ii) limited bits

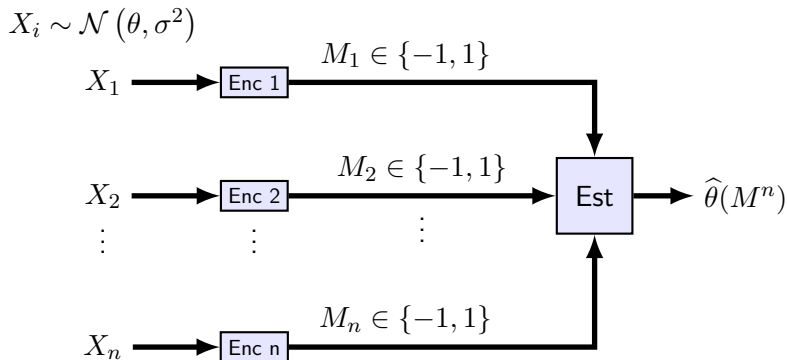
Relevant scenarios:

- ▶ big data
- ▶ low-power sensors
- ▶ distributed computing / optimization

This talk:

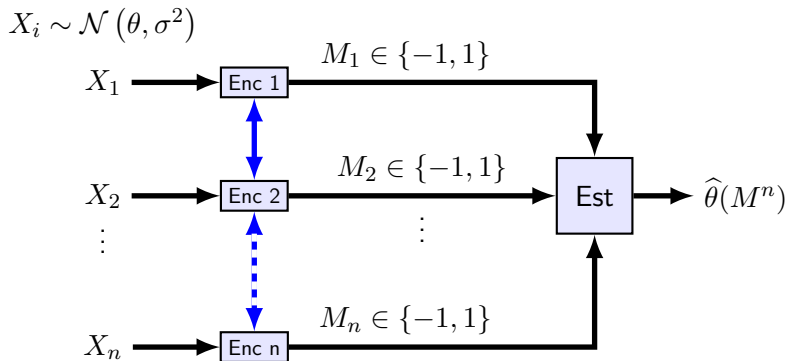
Estimating the mean θ of a normal distribution $\mathcal{N}(\theta, \sigma^2)$ from one-bit per sample (σ is known)

Three Encoding Scenarios



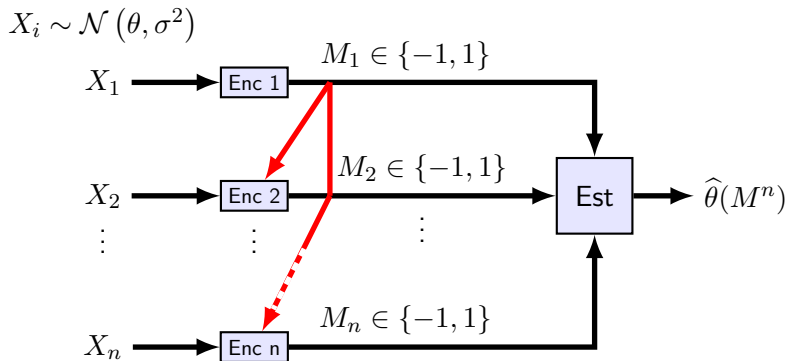
- Distributed: $M_i = f_i(X_i)$

Three Encoding Scenarios



- Distributed: $M_i = f_i(X_i)$
- Centralized: $M^n = (M_1, \dots, M_n) = f(X_1, \dots, X_n)$

Three Encoding Scenarios



- ▶ Distributed: $M_i = f_i(X_i)$
- ▶ Centralized: $M^n = (M_1, \dots, M_n) = f(X_1, \dots, X_n)$
- ▶ Adaptive / Sequential: $M_i = f_i(X_i, M^{i-1})$

Related Work

Related Work

- ▶ Estimation via compressed information [Han '87], [Zhang & Berger '88] (centralized)

Related Work

- ▶ Estimation via compressed information [Han '87], [Zhang & Berger '88] (centralized)
- ▶ Estimation from multiple machines subject to a bit constraint [Zhang, Duchi, Jordan, Wainwright '13] (distributed / adaptive)

Related Work

- ▶ Estimation via compressed information [Han '87], [Zhang & Berger '88] (centralized)
- ▶ Estimation from multiple machines subject to a bit constraint [Zhang, Duchi, Jordan, Wainwright '13] (distributed / adaptive)
- ▶ Distributed hypothesis testing under quantization [Tsitsiklis '88] (distributed)

Related Work

- ▶ Estimation via compressed information [Han '87], [Zhang & Berger '88] (centralized)
- ▶ Estimation from multiple machines subject to a bit constraint [Zhang, Duchi, Jordan, Wainwright '13] (distributed / adaptive)
- ▶ Distributed hypothesis testing under quantization [Tsitsiklis '88] (distributed)
- ▶ Remote multiterminal source coding (CEO) [Berger, Zhang, Viswanathan '96], [Oohama '97] (distributed)

Consistency

Q: in what setting consistent estimation is possible?

Consistency

Q: in what setting consistent estimation is possible?

A: all !

Consistency

Q: in what setting consistent estimation is possible?

A: all ! **Proof:**

$$M_i = \mathbf{1}(X_i > 0), \quad i = 1, \dots, n$$

(as in the distributed setting)

$$\frac{1}{n} \sum_{i=1}^n M_i \rightarrow \mathbb{P}(X > 0) = \Phi(\theta/\sigma)$$

Efficiency

Definition: *asymptotic relative efficiency (ARE)* of an estimator:

$$\text{ARE}(\hat{\theta}) \triangleq \lim_{n \rightarrow \infty} \mathbb{E} \left[\left(\hat{\theta} - \theta \right)^2 \right] / (\sigma^2/n)$$

(relative to minimax risk without bit constraint)

ARE under Centralized Encoding

Proposition

If the parameter space Θ is bounded, then the ARE under centralized encoding is 1

ARE under Centralized Encoding

Proposition

If the parameter space Θ is bounded, then the ARE under centralized encoding is 1

Proof:

$$\mathbb{E} \left(\theta - \hat{\theta} \right)^2 = \overbrace{\mathbb{E} \left(\theta - \bar{\theta} \right)^2}^{\sigma^2/n} + \underbrace{\mathbb{E} \left(\bar{\theta} - \hat{\theta} \right)^2}_{O(2^{-2n})}$$

ARE under Centralized Encoding

Proposition

If the parameter space Θ is bounded, then the ARE under centralized encoding is 1

Proof:

$$\mathbb{E} \left(\theta - \hat{\theta} \right)^2 = \overbrace{\mathbb{E} \left(\theta - \bar{\theta} \right)^2}^{\sigma^2/n} + \underbrace{\mathbb{E} \left(\bar{\theta} - \hat{\theta} \right)^2}_{O(2^{-2n})}$$

- ▶ Encoder is required to describe $\bar{\theta}$ using n bits
 - ▶ divide parameter space Θ into 2^n regions of equal size
 - ▶ send region index where $\bar{\theta}$ falls

ARE under Centralized Encoding

Proposition

If the parameter space Θ is bounded, then the ARE under centralized encoding is 1

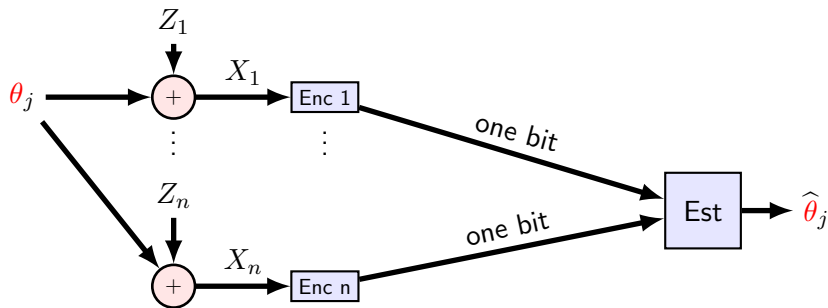
Proof:

$$\mathbb{E} \left(\theta - \hat{\theta} \right)^2 = \overbrace{\mathbb{E} \left(\theta - \bar{\theta} \right)^2}^{\sigma^2/n} + \underbrace{\mathbb{E} \left(\bar{\theta} - \hat{\theta} \right)^2}_{O(2^{-2n})}$$

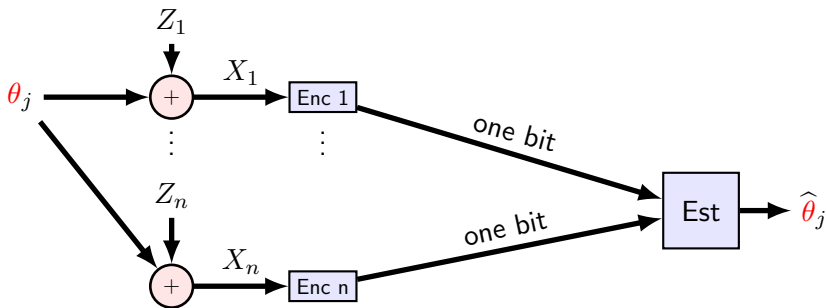
- ▶ Encoder is required to describe $\bar{\theta}$ using n bits
 - ▶ divide parameter space Θ into 2^n regions of equal size
 - ▶ send region index where $\bar{\theta}$ falls

Note: mean of $\bar{\theta}$ is unknown – hard to derive a globally optimal strategy

Relation to CEO



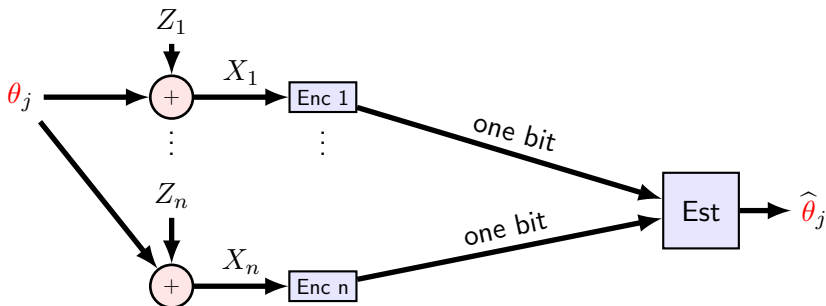
Relation to CEO



Assume:

- ▶ $Z_1, \dots, Z_n \stackrel{\text{i.i.d}}{\sim} \mathcal{N}(0, \sigma^2)$

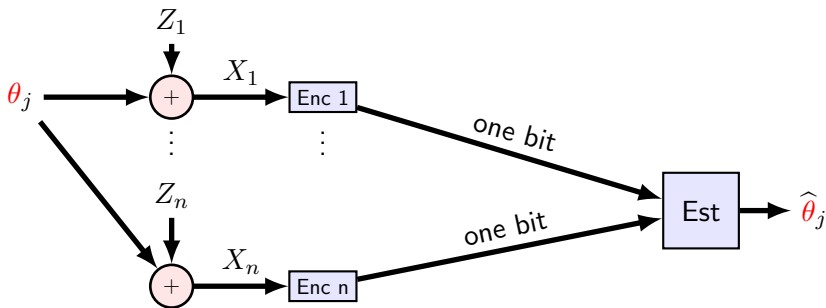
Relation to CEO



Assume:

- ▶ $Z_1, \dots, Z_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2)$
- ▶ Replicate k times: $\theta_1, \dots, \theta_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_0^2)$

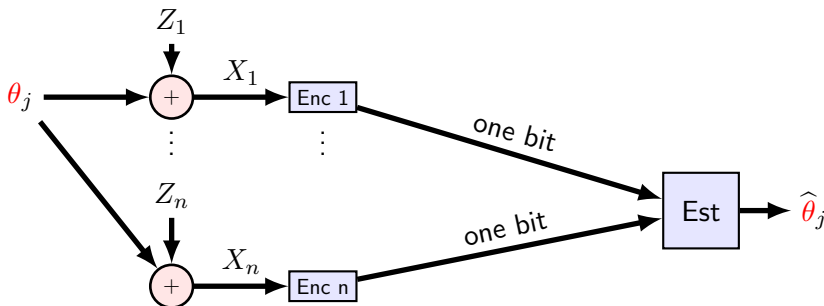
Relation to CEO



Assume:

- ▶ $Z_1, \dots, Z_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2)$
- ▶ Replicate k times: $\theta_1, \dots, \theta_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_0^2)$
- ▶ Block encode $X_{j,1}, \dots, X_{j,n}$ using k bits, $j = 1, \dots, k$

Relation to CEO



Assume:

- ▶ $Z_1, \dots, Z_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2)$
- ▶ Replicate k times: $\theta_1, \dots, \theta_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_0^2)$
- ▶ Block encode $X_{j,1}, \dots, X_{j,n}$ using k bits, $j = 1, \dots, k$
- ▶ $D_{CEO} = \inf_{k, \hat{\theta}} \frac{1}{k} \sum_{j=1}^k \mathbb{E} \left(\theta_j - \hat{\theta}_j \right)^2$

Quadratic Gaussian CEO under optimal rate allocation [Chen, Zhang, Berger, Wicker '04] :

$$D_{CEO} \geq \frac{4}{3} \frac{\sigma^2}{n} + o(1/n)$$

Quadratic Gaussian CEO under optimal rate allocation [Chen, Zhang, Berger, Wicker '04] :

$$D_{CEO} \geq \frac{4}{3} \frac{\sigma^2}{n} + o(1/n)$$

Conclusion

Distributed encoding will hurt you (even if you can repeat experiment and encode over blocks)

Quadratic Gaussian CEO under optimal rate allocation [Chen, Zhang, Berger, Wicker '04] :

$$D_{CEO} \geq \frac{4}{3} \frac{\sigma^2}{n} + o(1/n)$$

Conclusion

Distributed encoding will hurt you (even if you can repeat experiment and encode over blocks)

- In fact [K., Rini, Goldsmith '17]:

$$D_{CEO} \leq \frac{4}{3} \frac{\sigma^2}{n} + \frac{\sigma_0^2}{3n} + o(1/n)$$

(ARE of 4/3 is achievable with coding over blocks)

Table of Contents

Introduction

Motivation

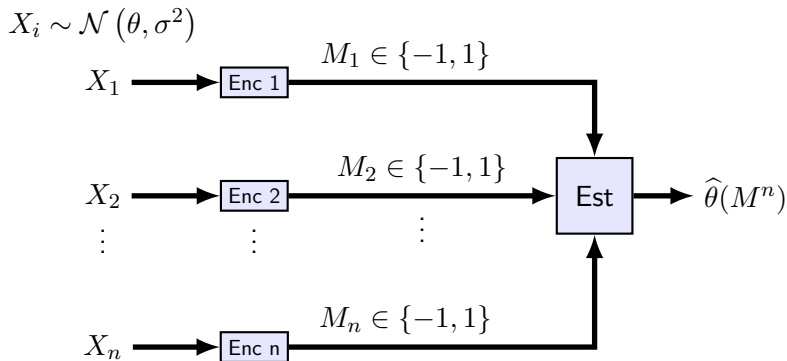
Preliminary

Adaptive Encoding

Distributed Encoding

Summary

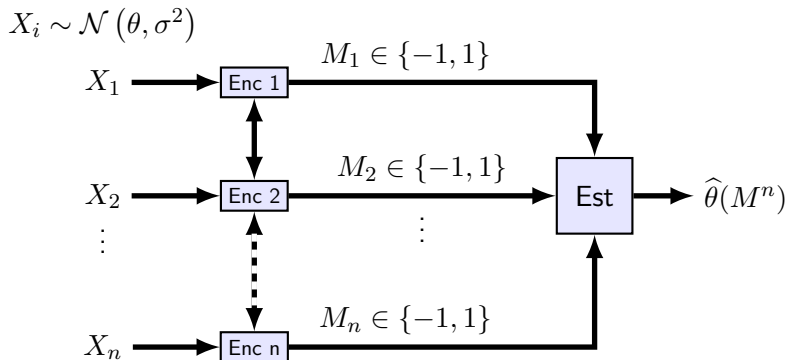
Three Encoding Scenarios



► Distributed: $M_i = f_i(X_i)$

$\text{ARE} \geq 4/3$

Three Encoding Scenarios

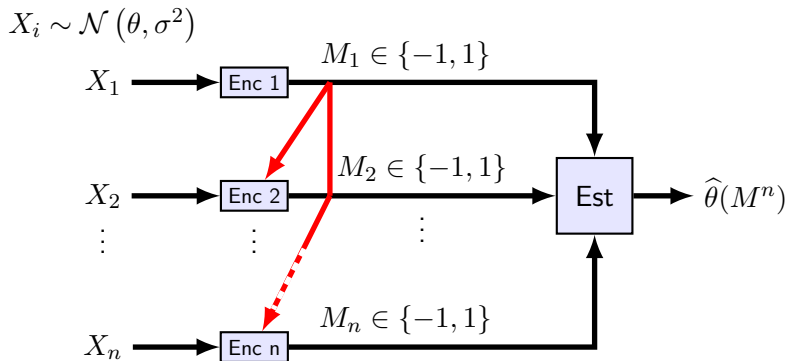


- Distributed: $M_i = f_i(X_i)$
- Centralized: $M^n = f(X_1, \dots, X_n)$

$$\text{ARE} \geq 4/3$$

$$\text{ARE} = 1$$

Three Encoding Scenarios



- ▶ Distributed: $M_i = f_i(X_i)$ ARE $\geq 4/3$
- ▶ Centralized: $M^n = f(X_1, \dots, X_n)$ ARE = 1
- ▶ Adaptive / Sequential: $M_i = f_i(X_i, M^{i-1})$ ARE = $\pi/2$

Main Results (adaptive encoding)

Theorem (achievability)

There exists an estimator with ARE $\pi/2$

Theorem (converse)

No estimator have ARE lower than $\pi/2$

Achievability

existence of an estimator with $\text{ARE} = \pi/2$

(i) For $X \sim \mathcal{N}(\theta, \sigma^2)$, $\text{med}(X) = \theta$

Achievability

existence of an estimator with $\text{ARE} = \pi/2$

(i) For $X \sim \mathcal{N}(\theta, \sigma^2)$, $\text{med}(X) = \theta$

(ii) $\text{med}(X) = \underset{m}{\operatorname{argmin}} \mathbb{E} |X - m|$

Achievability

existence of an estimator with $\text{ARE} = \pi/2$

- (i) For $X \sim \mathcal{N}(\theta, \sigma^2)$, $\text{med}(X) = \theta$
- (ii) $\text{med}(X) = \underset{m}{\operatorname{argmin}} \mathbb{E} |X - m|$
- (iii) Stochastic gradient descent on $\mathbb{E} |X - \theta|$:

$$\theta_n = \theta_{n-1} + \gamma_n \text{sign}(X_n - \theta_n)$$

$$\hat{\theta}_n = \frac{1}{n} \sum_{i=1}^n \theta_i$$

Achievability

existence of an estimator with $\text{ARE} = \pi/2$

- (i) For $X \sim \mathcal{N}(\theta, \sigma^2)$, $\text{med}(X) = \theta$
- (ii) $\text{med}(X) = \underset{m}{\operatorname{argmin}} \mathbb{E} |X - m|$
- (iii) Stochastic gradient descent on $\mathbb{E} |X - \theta|$:

$$\theta_n = \theta_{n-1} + \gamma_n \text{sign}(X_n - \theta_n)$$

$$\hat{\theta}_n = \frac{1}{n} \sum_{i=1}^n \theta_i$$

From [Polyak & Juditsky '92] (under conditions on (γ_n)):

$$\sqrt{n}(\theta - \hat{\theta}_n) \rightarrow \mathcal{N}(0, \sigma^2 \pi/2)$$

Converse

$$\text{ARE} \geq \pi/2$$

The van-Trees inequality (e.g. [Tsybakov '08]) implies

$$\mathbb{E} \left(\theta - \hat{\theta} \right)^2 \geq \frac{1}{\mathbb{E} I_{\theta}(M^n) + I_{\pi}} \geq \frac{1}{\sum_{i=1}^n I_{\theta}(M_i | M^{i-1}) + I_{\pi}}$$

I_{π} is the location Fisher information with w.r.t. some prior $\pi(d\theta)$ on Θ

Lemma (K. & Duchi '17)

$$I_{\theta}(M_i | M^{i-1}) \leq \frac{2}{\pi \sigma^2}$$

Proof:

Stein identity implies that portion maximizing the information is a threshold: $M_i^{-1}(1) = (\theta, \infty)$

The rest follows by induction over number of portions

Table of Contents

Introduction

Motivation

Preliminary

Adaptive Encoding

Distributed Encoding

Summary

Distributed Encoding

Threshold Detection

We consider only messages of the form

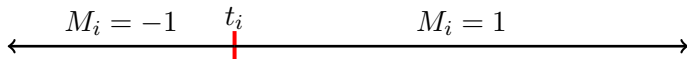
$$M_i = \text{sign}(X_i - t_i), \quad i = 1, \dots, n$$

Distributed Encoding

Threshold Detection

We consider only messages of the form

$$M_i = \text{sign}(X_i - t_i), \quad i = 1, \dots, n$$



Assume:

$$\lambda_n([a, b]) = \frac{1}{n} \text{card}([a, b] \cap \{t_i\})$$

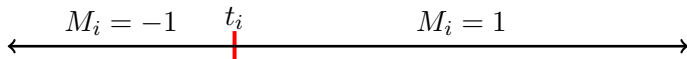
converges weakly to a probability distribution λ

Distributed Encoding

Threshold Detection

We consider only messages of the form

$$M_i = \text{sign}(X_i - t_i), \quad i = 1, \dots, n$$



Assume:

$$\lambda_n([a, b]) = \frac{1}{n} \text{card}([a, b] \cap \{t_i\})$$

converges weakly to a probability distribution λ

Example: t_1, \dots, t_n drawn i.i.d. from a distribution λ on \mathbb{R}

Main Results (distributed encoding)

Theorem

(i) *The Maximum likelihood estimator $\hat{\theta}_{ML}$ satisfies*

$$\sqrt{n}(\theta - \hat{\theta}_{ML}) \rightarrow \mathcal{N}(0, \sigma^2 / K_{\lambda}(\theta))$$

where:

$$K_{\lambda}(\theta) = \int_{\mathbb{R}} \eta\left(\frac{t - \theta}{\sigma}\right) \lambda(dt)$$

$$\eta(x) = \frac{\phi^2(x)}{\Phi(x)\Phi(-x)}$$

Main Results (distributed encoding)

Theorem

(i) *The Maximum likelihood estimator $\hat{\theta}_{ML}$ satisfies*

$$\sqrt{n}(\theta - \hat{\theta}_{ML}) \rightarrow \mathcal{N}(0, \sigma^2 / K_{\lambda}(\theta))$$

where:

$$K_{\lambda}(\theta) = \int_{\mathbb{R}} \eta\left(\frac{t - \theta}{\sigma}\right) \lambda(dt)$$
$$\eta(x) = \frac{\phi^2(x)}{\Phi(x)\Phi(-x)}$$

(ii) *For any estimator $\hat{\theta}(M_1, \dots, M_n)$:*

$$\liminf_{c \rightarrow \infty} \liminf_{n \rightarrow \infty} \sup_{\tau: |\tau - \theta| \leq \frac{c}{\sqrt{n}}} n \mathbb{E} \left(\hat{\theta} - \tau \right)^2 \geq \sigma^2 / K_{\lambda}(\theta),$$

Interpretations

Interpretations

- ▶ ML estimator is local asymptotically minimax

Interpretations

- ▶ ML estimator is local asymptotically minimax
- ▶ ARE of ML is $1/K_\lambda(\theta)$ – only depends on asymptotic threshold density λ

Interpretations

- ▶ ML estimator is local asymptotically minimax
- ▶ ARE of ML is $1/K_\lambda(\theta)$ – only depends on asymptotic threshold density λ



$$\text{ARE} = \frac{1}{K_\lambda(\theta)} = \frac{\sigma^2}{\int \eta\left(\frac{t-\theta}{\sigma}\right) \lambda(dt)} > \pi/2$$

(equality iff $\lambda = \delta_\theta$)

Interpretations

- ▶ ML estimator is local asymptotically minimax
- ▶ ARE of ML is $1/K_\lambda(\theta)$ – only depends on asymptotic threshold density λ



$$\text{ARE} = \frac{1}{K_\lambda(\theta)} = \frac{\sigma^2}{\int \eta\left(\frac{t-\theta}{\sigma}\right) \lambda(dt)} > \pi/2$$

(equality iff $\lambda = \delta_\theta$)

- ▶ Minimax λ and ARE can be found using a convex program – depends on radius of Θ

Table of Contents

Introduction

Motivation

Preliminary

Adaptive Encoding

Distributed Encoding

Summary

Summary

Summary

- ▶ ARE in adaptive setting is $\pi/2$ regardless of size of parameter space

Summary

- ▶ ARE in adaptive setting is $\pi/2$ regardless of size of parameter space
- ▶ ~ 1.57 more samples are required due to 1-bit constraints

Summary

- ▶ ARE in adaptive setting is $\pi/2$ regardless of size of parameter space
- ▶ ~ 1.57 more samples are required due to 1-bit constraints
- ▶ MLE is local asymptotically optimal for threshold detection

Summary

- ▶ ARE in adaptive setting is $\pi/2$ regardless of size of parameter space
- ▶ ~ 1.57 more samples are required due to 1-bit constraints
- ▶ MLE is local asymptotically optimal for threshold detection
- ▶ ARE of MLE characterized by density of threshold values

Summary

- ▶ ARE in adaptive setting is $\pi/2$ regardless of size of parameter space
- ▶ ~ 1.57 more samples are required due to 1-bit constraints
- ▶ MLE is local asymptotically optimal for threshold detection
- ▶ ARE of MLE characterized by density of threshold values

Open question

Is there a distributed encoding scheme with ARE that is both finite and independent of size of parameter space ?

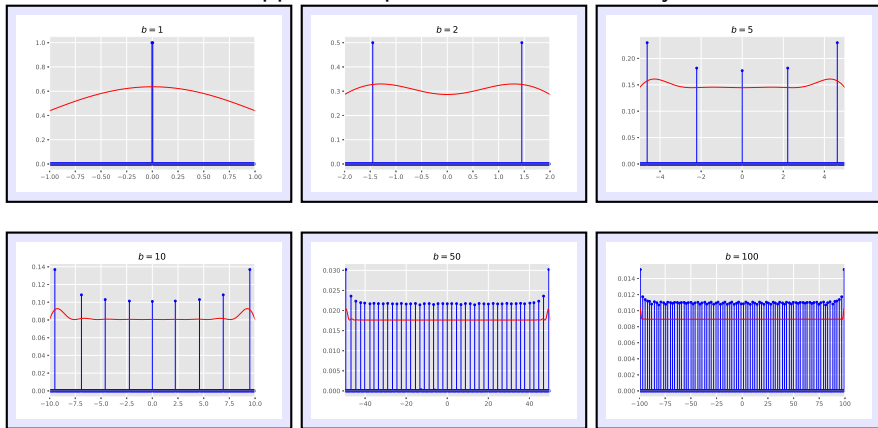
Minimax threshold density

Minimax λ for $\theta \in (-b\sigma, b\sigma)$:

$$\begin{aligned} &\text{maximize} && \inf_{\tau \in (-b, b)} \int \eta(t - \tau) \lambda(dt) \\ &\text{subject to} && \lambda(dt) \geq 0, \quad \int \lambda(dt) \leq 1. \end{aligned}$$

Minimax λ

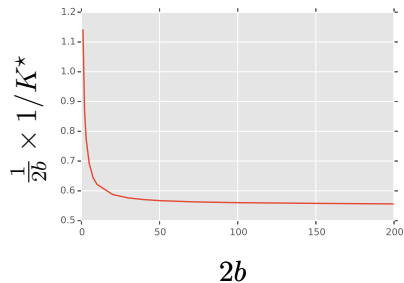
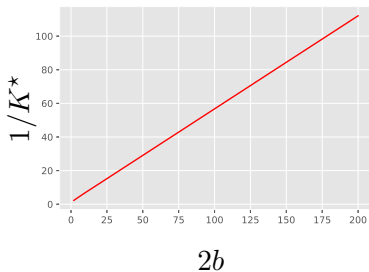
support of optimal threshold density λ^*



$$K^* = \inf_{\theta} K^*(\theta) = \inf_{\theta} \int \eta(t - \theta) \lambda^*(dt)$$

Minimax λ

Minimax ARE vs size of parameter space



- ▶ ARE increases with size of parameter space