

Bayesian Estimation, Large Deviations, and Incomplete Data

Rudolf Kulhavý

Institute of Information Theory and Automation
Academy of Sciences of the Czech Republic

Outline

- Bayesian Estimation
- Data Compression
- Approximate Posterior Distribution
- Key Properties of Approximation
- Illustrative Example
- Concluding Remarks

Bayesian Estimation

sequence of random variables

$$\mathbf{X} = (X_1, \dots, X_k), \quad X_i \in \mathcal{X}, \quad \mathcal{X} \text{ finite}$$

independent and identically distributed

$$S_\theta(\mathbf{x}), \quad \theta \in \mathcal{T}, \quad \mathcal{T} \text{ finite}, \quad S_\theta(\mathbf{x}) > 0$$

posterior distribution of Θ

$$P_{\mathbf{x}}(\theta) \propto P(\theta) S_\theta^k(\mathbf{x})$$

Some Notions

type of sequence $\mathbf{x} = (x_1, \dots, x_k)$

$$R_{\mathbf{x}}(a) = \frac{1}{k} N_{\mathbf{x}}(a) = \frac{1}{k} \sum_{i=1}^k \mathbf{1}_{x_i}(a)$$

Shannon entropy

$$H(R) = \sum_{a \in \mathcal{X}} R(a) \log \frac{1}{R(a)}$$

relative entropy, Kullback-Leibler distance, I-divergence

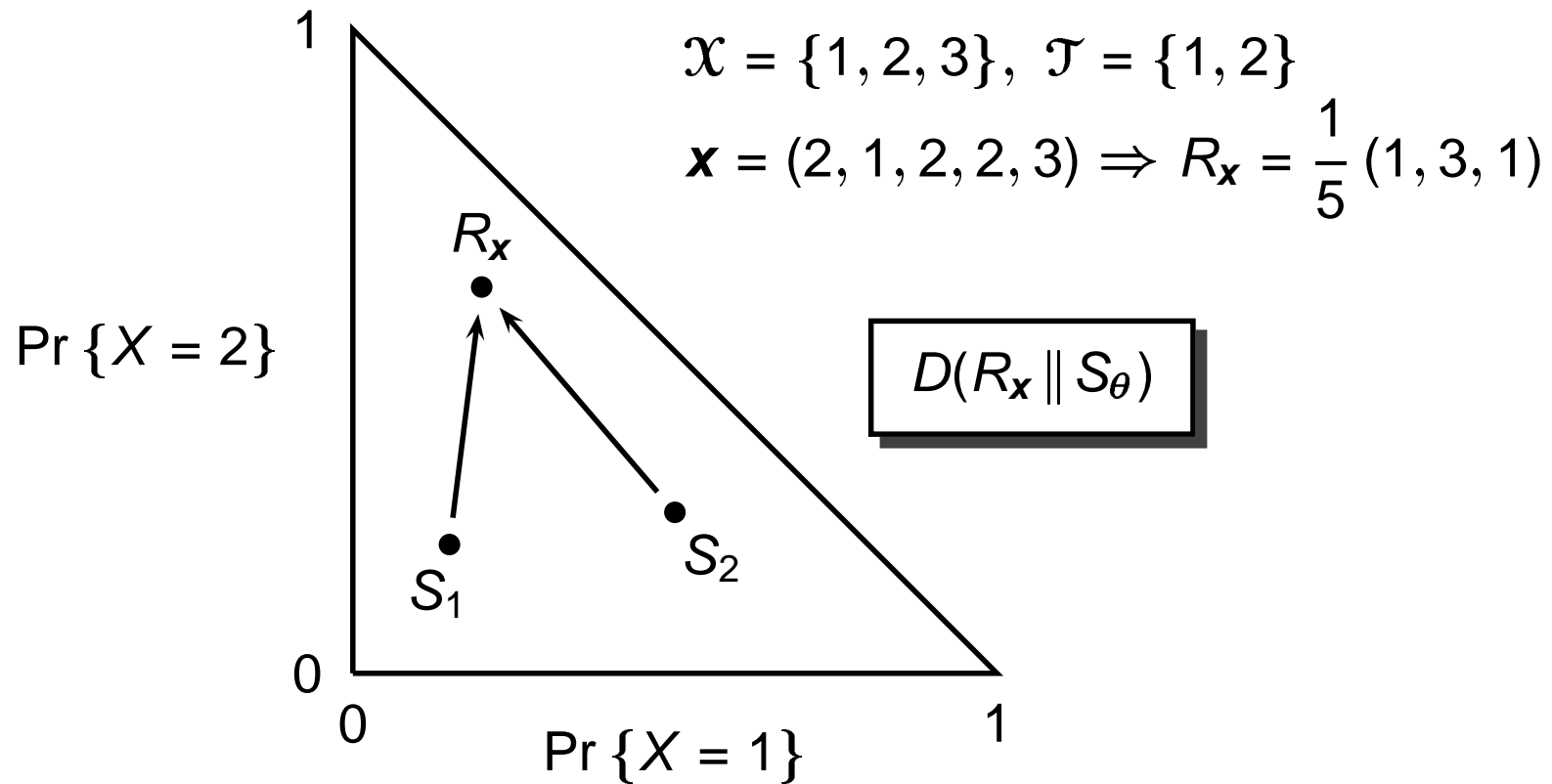
$$D(R \parallel S) = \sum_{a \in \mathcal{X}} R(a) \log \frac{R(a)}{S(a)}$$

Distribution of Sample

joint distribution of $\mathbf{X} = (X_1, \dots, X_k)$

$$\begin{aligned} S_{\theta}^k(\mathbf{x}) &= \prod_{i=1}^k S_{\theta}(x_i) \\ &= \exp \sum_{i=1}^k \log S_{\theta}(x_i) \\ &= \exp \{-k H(R_{\mathbf{x}})\} \exp \{-k D(R_{\mathbf{x}} \| S_{\theta})\} \end{aligned}$$

Probability Simplex



Posterior Distribution

$$P_{\mathbf{x}}(\theta) \propto \underbrace{P(\theta)}_{\text{prior}} \exp\left\{-\underbrace{k}_{\text{size of } \mathbf{x}} D\left(\underbrace{R_{\mathbf{x}}}_{\text{type of } \mathbf{x}} \parallel \underbrace{S_{\theta}}_{\text{model}}\right)\right\}$$

The diagram shows the equation $P_{\mathbf{x}}(\theta) \propto P(\theta) \exp\{-k D(R_{\mathbf{x}} \parallel S_{\theta})\}$. The terms $P(\theta)$, k , $R_{\mathbf{x}}$, and S_{θ} are circled. Lines connect these circles to labels below: 'prior' for $P(\theta)$, 'size of \mathbf{x} ' for k , 'type of \mathbf{x} ' for $R_{\mathbf{x}}$, and 'model' for S_{θ} .

Data Compression

data statistic

$$T(\mathbf{X}) = \frac{1}{k} \sum_{j=1}^k h(X_j) = E_{R_{\mathbf{X}}} h(X), \quad h : \mathcal{X} \mapsto \mathbb{R}^n$$

a trivial example

$$\mathcal{X} \subset \mathbb{R}, \quad h_1(X) = X, \quad h_2(X) = X^2$$

partial knowledge of \mathbf{X}

$$T(\mathbf{X}) \geq \xi \quad \Rightarrow \quad P_{\xi}(\theta) \propto P(\theta) S_{\theta}^k(\xi)$$

Distribution of Statistic

set of \mathbf{x} 's consistent with the statistic

$$\mathcal{X}_\xi = \{\mathbf{x} : T(\mathbf{x}) \geq \xi\} = \cup_{R \in \mathcal{R}_\xi} \mathcal{X}_R$$

where $\mathcal{X}_R = \{\mathbf{x} : R_{\mathbf{x}} = R\}$ and $\mathcal{R}_\xi = \{R : E_R h \geq \xi\}$

probability of all such \mathbf{x} 's

$$\begin{aligned} S_\theta^k(\xi) &= S_\theta^k(\mathcal{X}_\xi) = \sum_{\mathbf{x} \in \mathcal{X}_\xi} S_\theta^k(\mathbf{x}) = \sum_{R \in \mathcal{R}_\xi} \sum_{\mathbf{x} \in \mathcal{X}_R} S_\theta^k(\mathbf{x}) = \\ &= \sum_{R \in \mathcal{R}_\xi} |\mathcal{X}_R| \exp\{-k H(R)\} \exp\{-k D(R \| S_\theta)\} \end{aligned}$$

Approximate Distribution of Statistic

equality to the 1st order in the exponent

$$S^k \doteq \hat{S}^k \quad \text{if} \quad \lim_{k \rightarrow \infty} \frac{1}{k} \log \frac{S^k}{\hat{S}^k} = 0$$

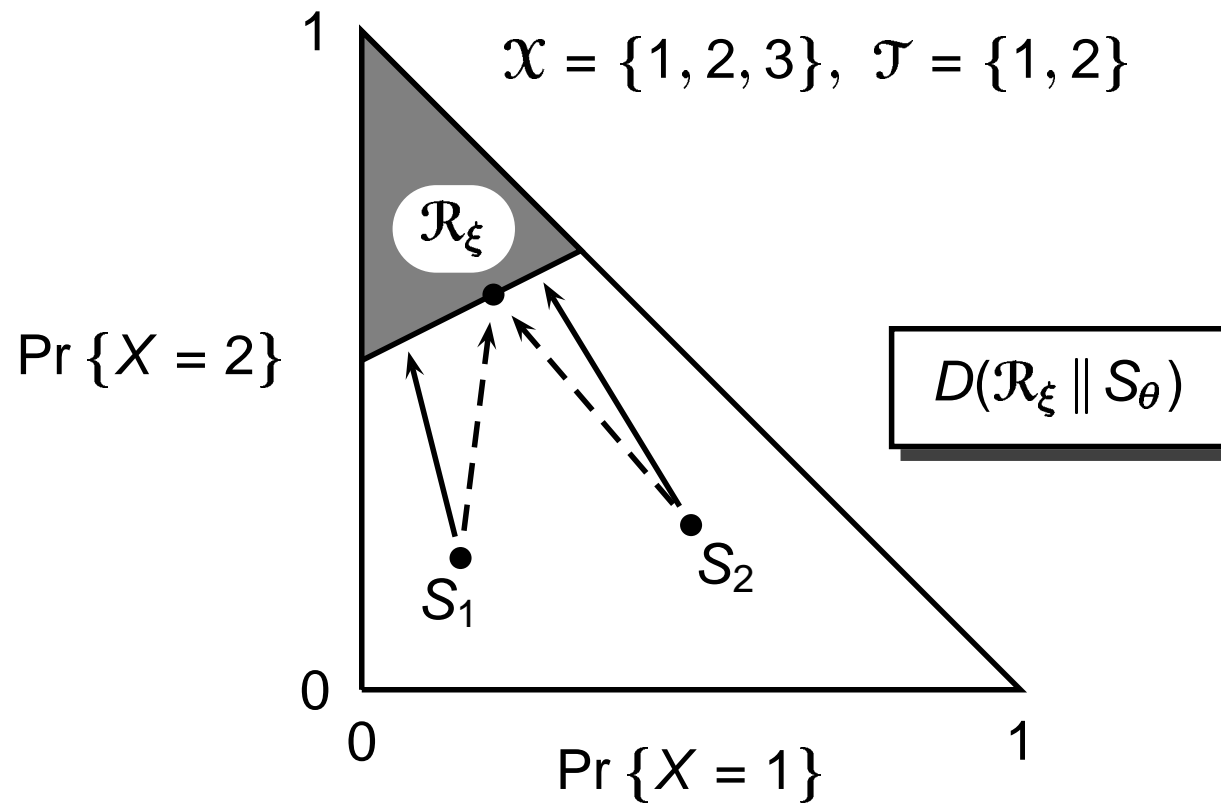
approximate expression for $S_\theta^k(\xi)$

$$S_\theta^k(\xi) \doteq \exp \{ -k D(\mathcal{R}_\xi \parallel S_\theta) \},$$

$$D(\mathcal{R}_\xi \parallel S_\theta) = \min_{R \in \mathcal{R}_\xi} D(R \parallel S_\theta)$$

large deviation theorem (Sanov, 1957)

Probability Simplex



Main Result:

Approximate Posterior Distribution

$$\hat{P}_{\mathbf{x}}(\theta) \propto P(\theta) \exp\left\{-k D(\mathcal{R}_{\xi} \parallel S_{\theta})\right\}$$

prior

size of \mathbf{x}

types of \mathbf{x} 's
consistent
with $T(\mathbf{x}) \geq \xi$

model

Key Properties of Approximation

meaningful bounds

$$0 \leq D(\mathcal{R}_\xi \parallel S_\theta) \leq D(R_x \parallel S_\theta)$$

monotonicity

$$\mathcal{R}_\xi \subset \mathcal{R}'_\xi \Rightarrow D(\mathcal{R}'_\xi \parallel S_\theta) \leq D(\mathcal{R}_\xi \parallel S_\theta)$$

extreme cases

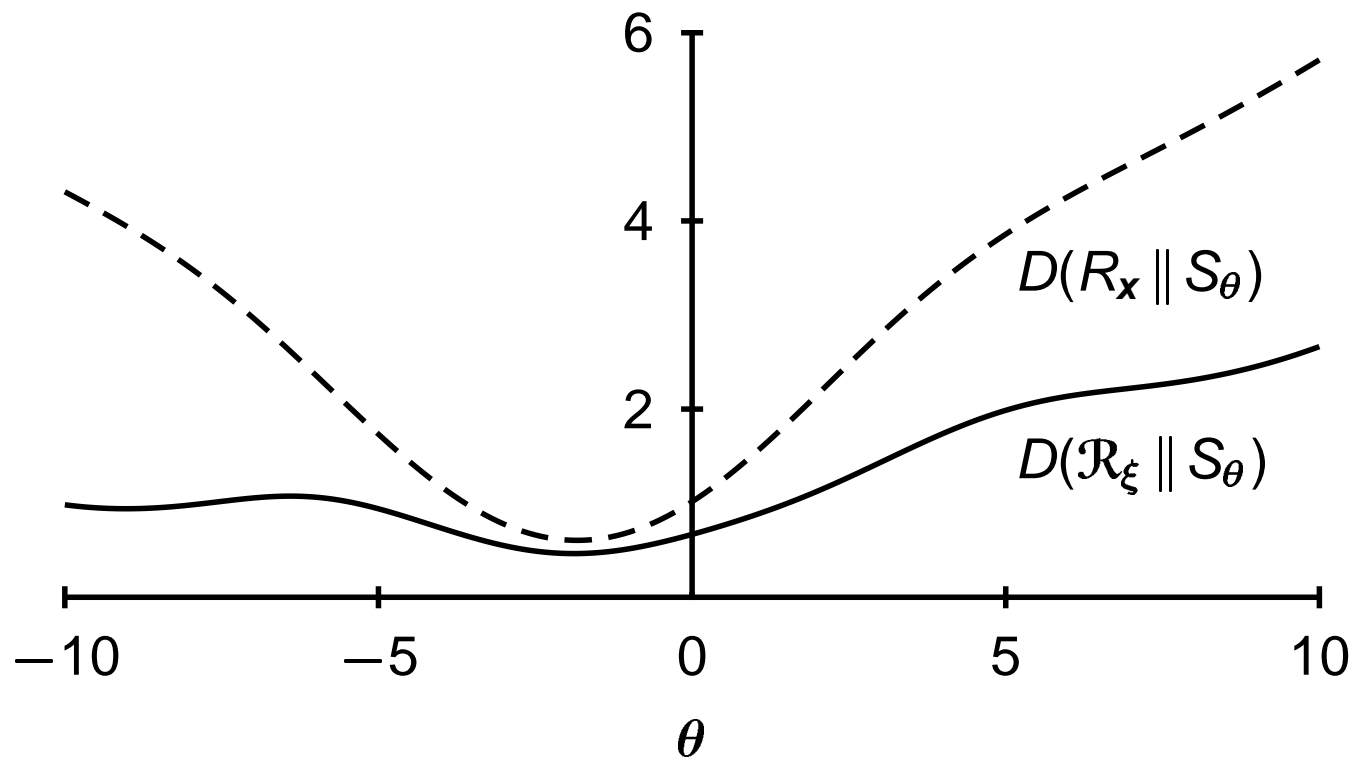
$$\mathcal{R}_\xi = \{R_x\} \Rightarrow D(\mathcal{R}_\xi \parallel S_\theta) = D(R_x \parallel S_\theta) \Rightarrow \hat{P}_\xi(\theta) = P_x(\theta)$$

$$\mathcal{R}_\xi = \mathcal{R} \Rightarrow D(\mathcal{R}_\xi \parallel S_\theta) = 0 \Rightarrow \hat{P}_\xi(\theta) = P(\theta)$$

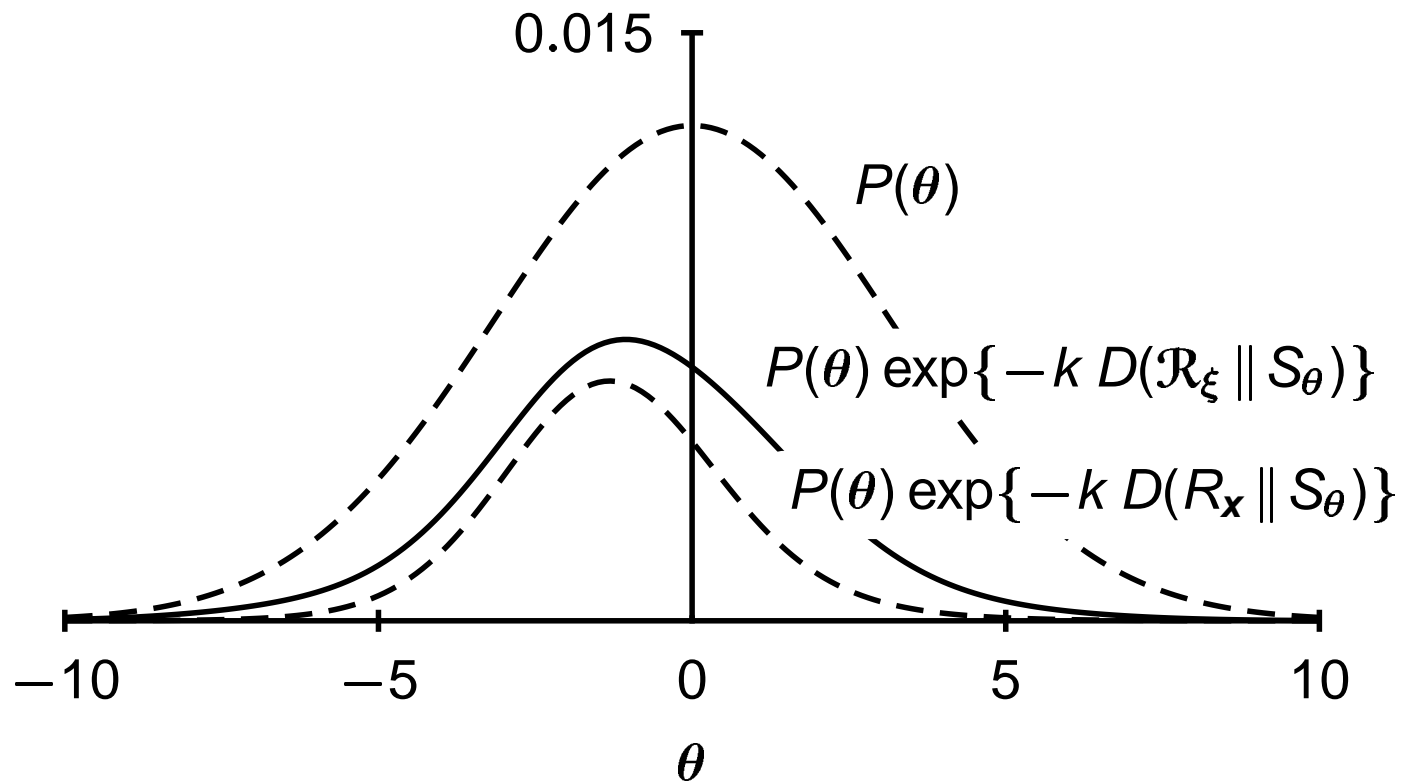
asymptotic consistency

$$R_x \rightarrow S_{\theta_0} \Rightarrow D(\mathcal{R}_\xi \parallel S_{\theta_0}) \rightarrow 0$$

Example: Relative Entropy



Example: Unnormalized Posterior



Concluding Remarks

↳ a way of coping with incomplete data

$$D(\mathcal{R}_\xi \parallel S_\theta) = \min_{R \in \mathcal{R}_\xi} D(R \parallel S_\theta)$$

↳ general approach

i.i.d. data, Markov chains, contin. θ , contin. x

↳ curse of dimensionality

computational complexity of $D(\mathcal{R}_\xi \parallel S_\theta)$

↳ shift of paradigm

from $P_x(\theta)$ to $D(\mathcal{R}_\xi \parallel S_\theta)$