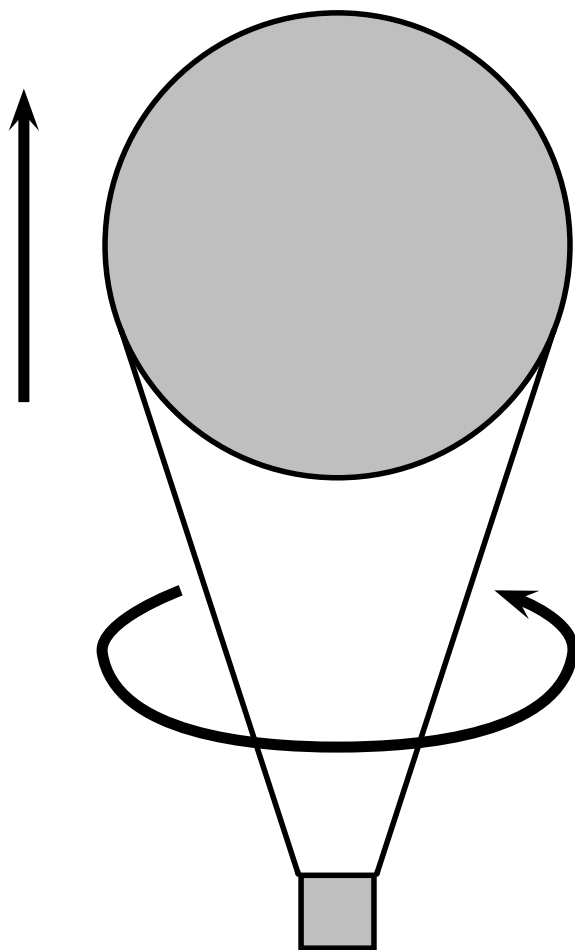


A Kullback-Leibler Distance Approach to System Identification

Rudolf Kulhavý

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

Example 1: BALLOON



Sun radiation measurements

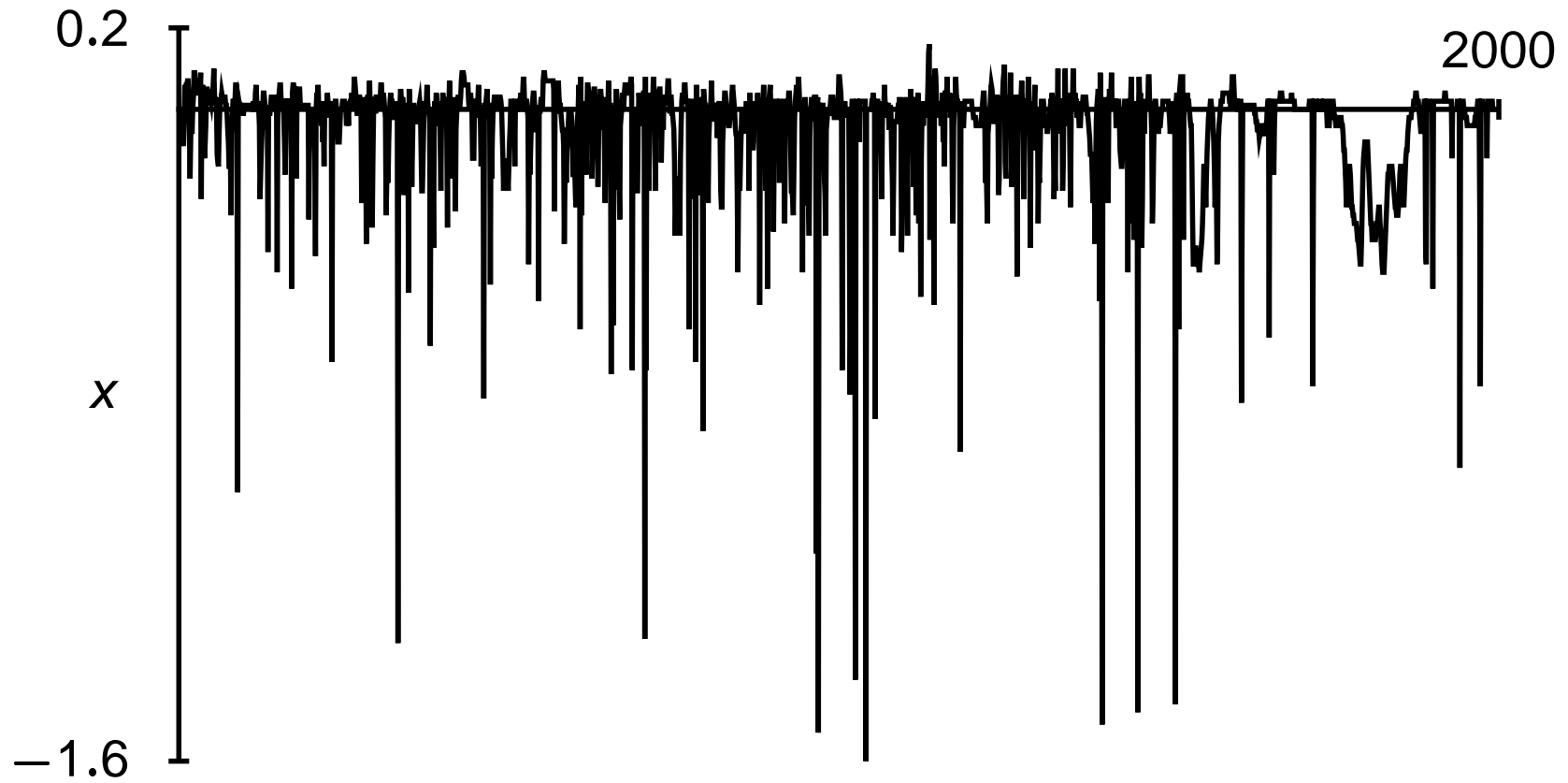
Source:

<http://lib.stat.cmu.edu>

`/Datasets/balloon`

L. Davis and U. Gather, JASA, 1993

BALLOON: Residuals

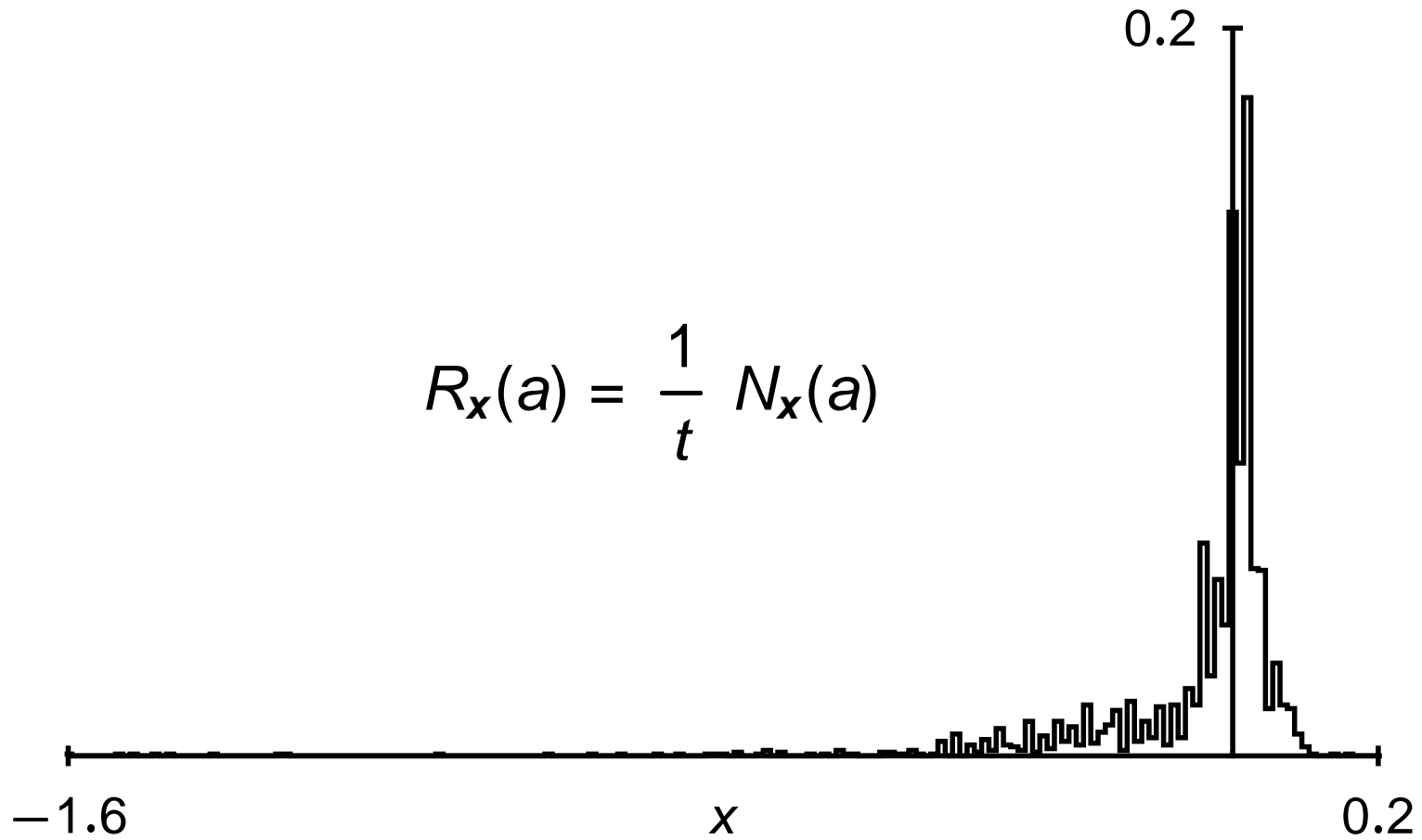




***What's the actual
distribution of Data?***

BALLOON: Empirical Distribution

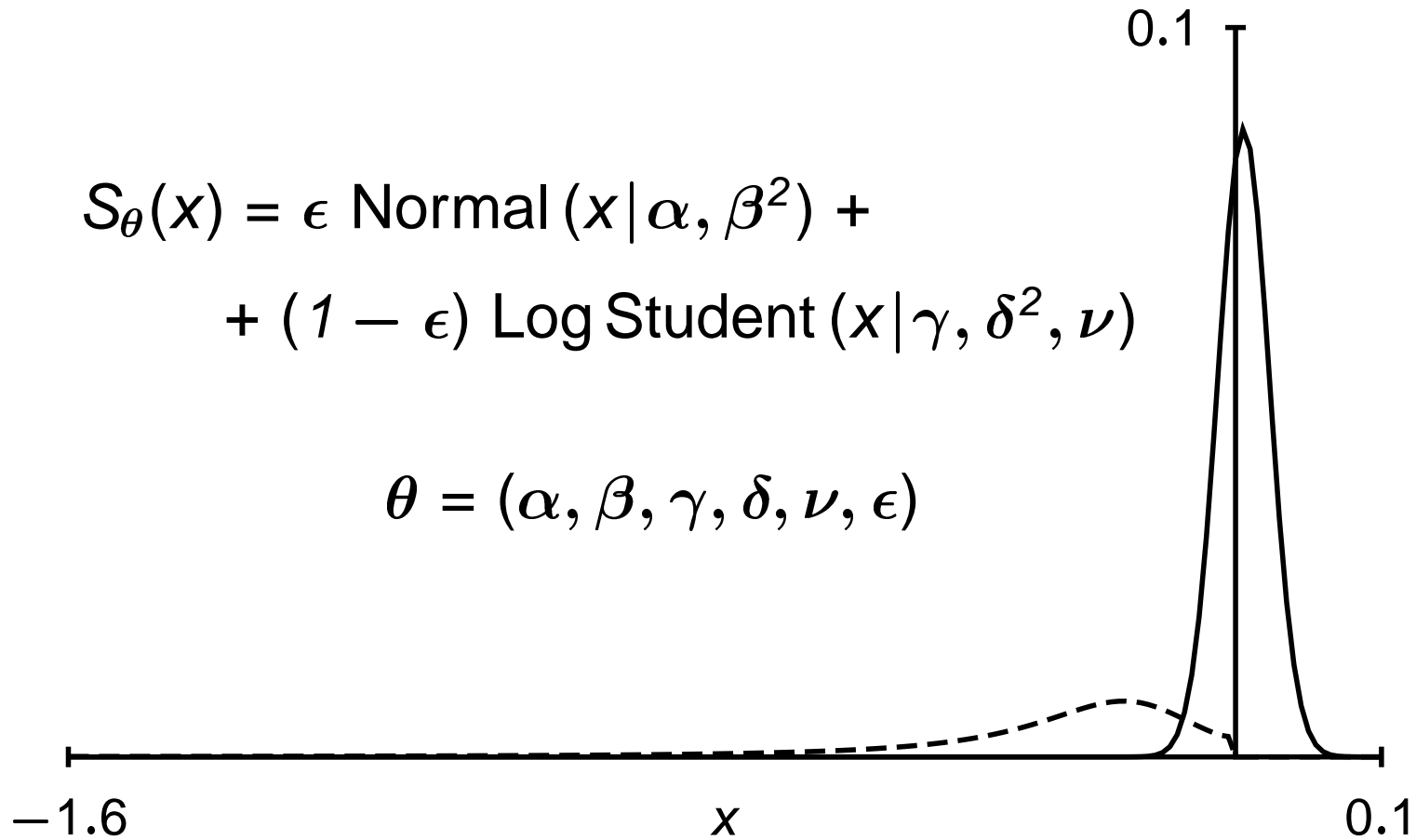
$$R_x(a) = \frac{1}{t} N_x(a)$$



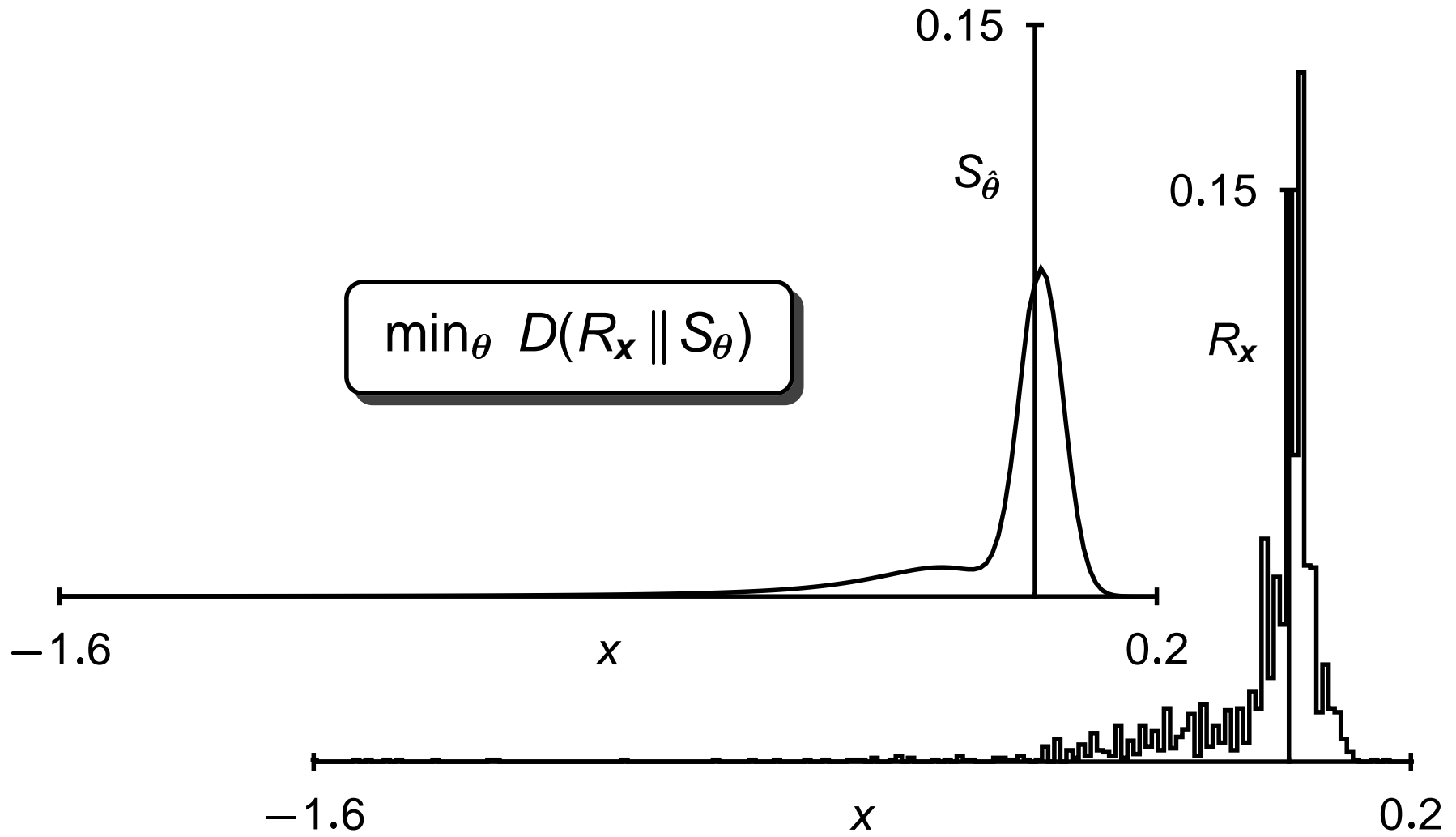
BALLOON: Model Distribution

$$S_{\theta}(x) = \epsilon \text{ Normal}(x|\alpha, \beta^2) + \\ + (1 - \epsilon) \text{ Log Student}(x|\gamma, \delta^2, \nu)$$

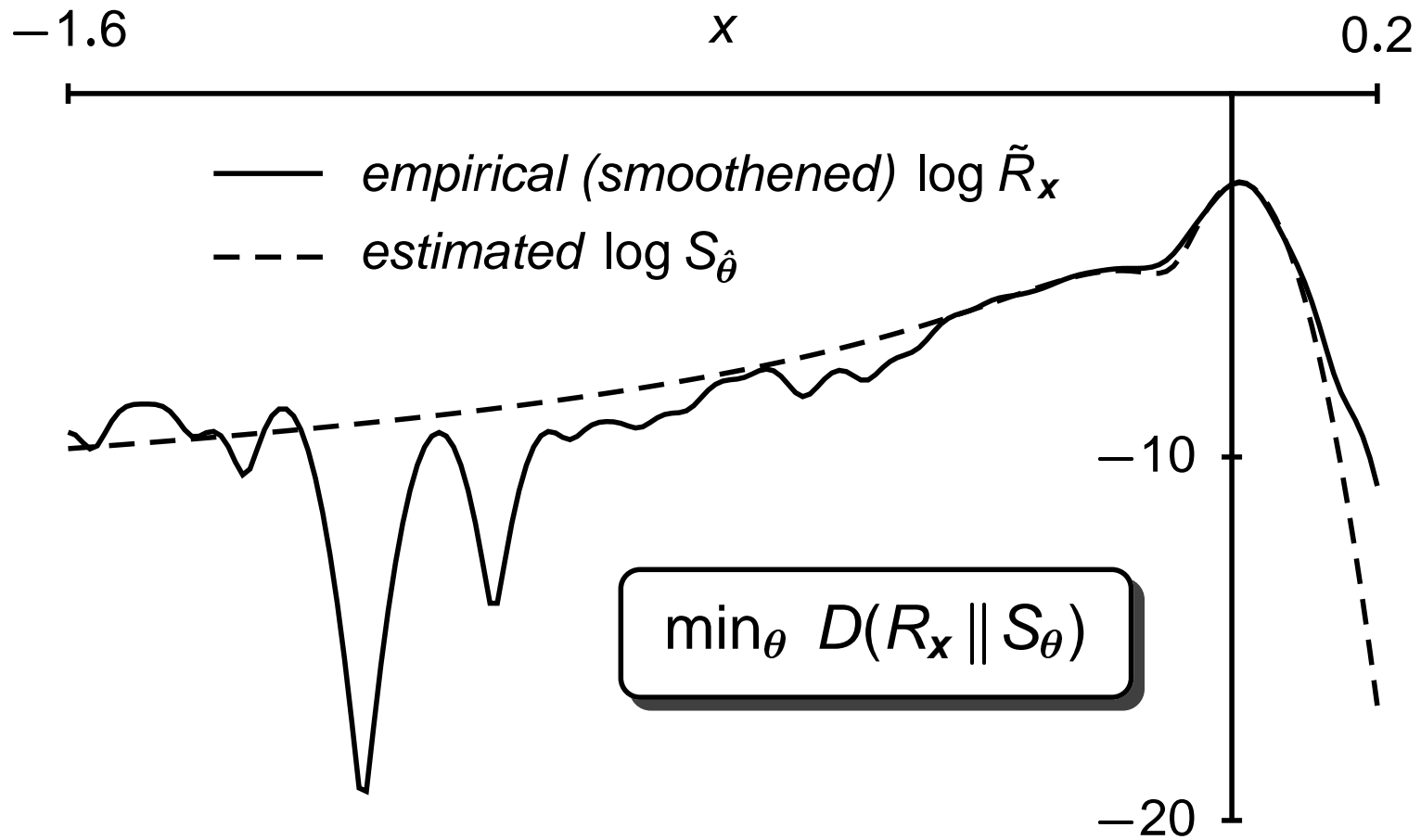
$$\theta = (\alpha, \beta, \gamma, \delta, \nu, \epsilon)$$



BALLOON: Parameter Estimation



BALLOON: Parameter Estimation



***What's a proper distance
of empirical and model
distributions?***

Ann. Math. Statist. **22** (1951) 79–86

ON INFORMATION AND SUFFICIENCY

BY S. KULLBACK AND R. A. LEIBLER

The George Washington University and Washington, D. C.

Measuring a “Distance” of Distributions

$$D(R \parallel S) = \sum_{x \in \mathcal{X}} R(x) \log \frac{R(x)}{S(x)} \quad \left| \begin{array}{l} \textit{Kullback-Leibler distance} \\ \textit{relative entropy} \end{array} \right.$$
$$D(R \parallel S) = \int_{\mathcal{X}} r(x) \log \frac{r(x)}{s(x)} dx \quad \left| \begin{array}{l} \textit{informational divergence} \\ \textit{discrimination information} \end{array} \right.$$

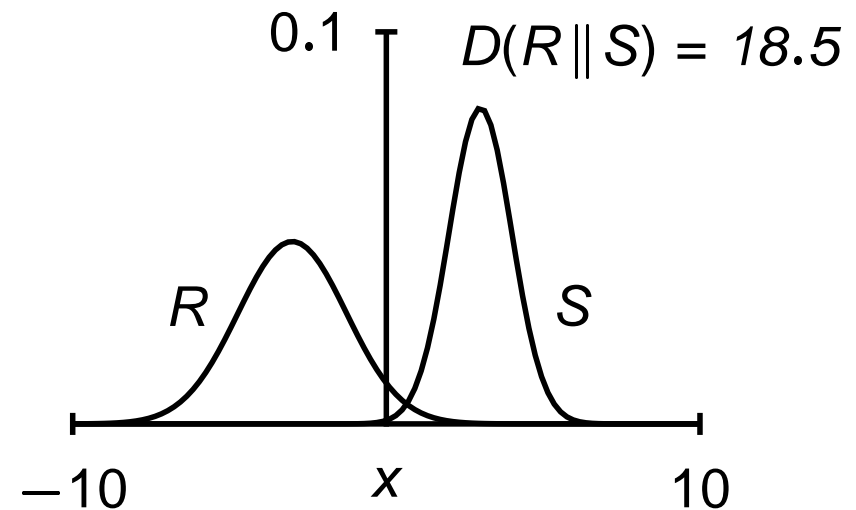
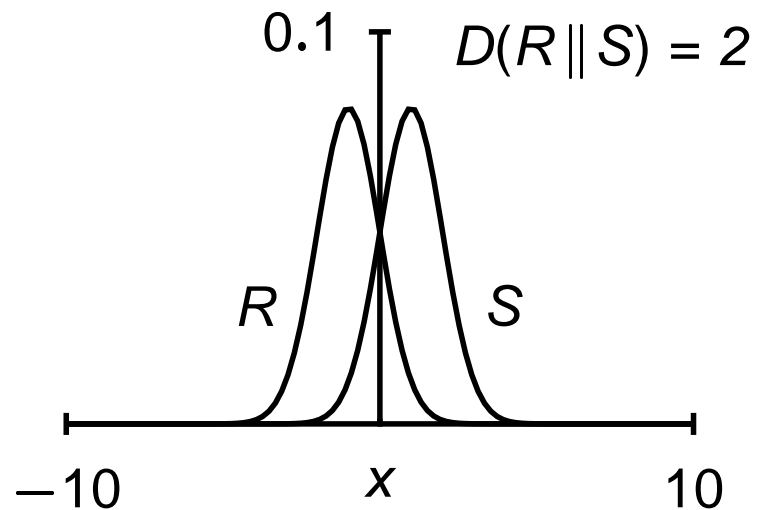
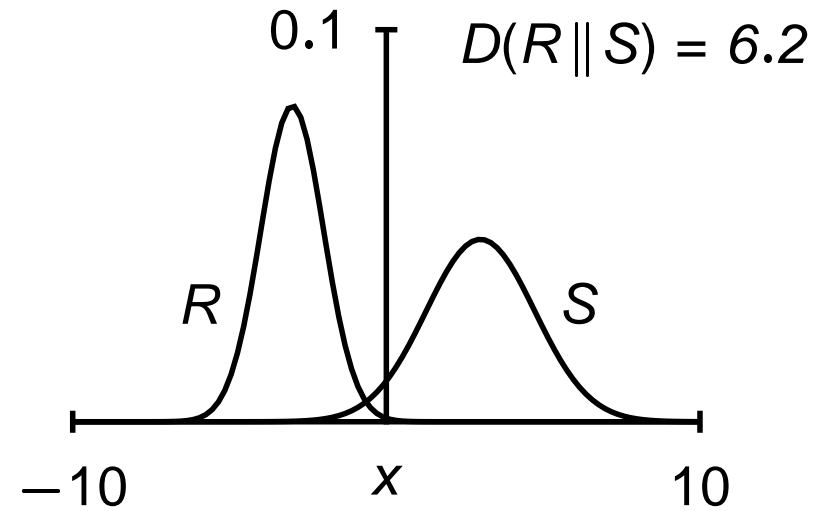
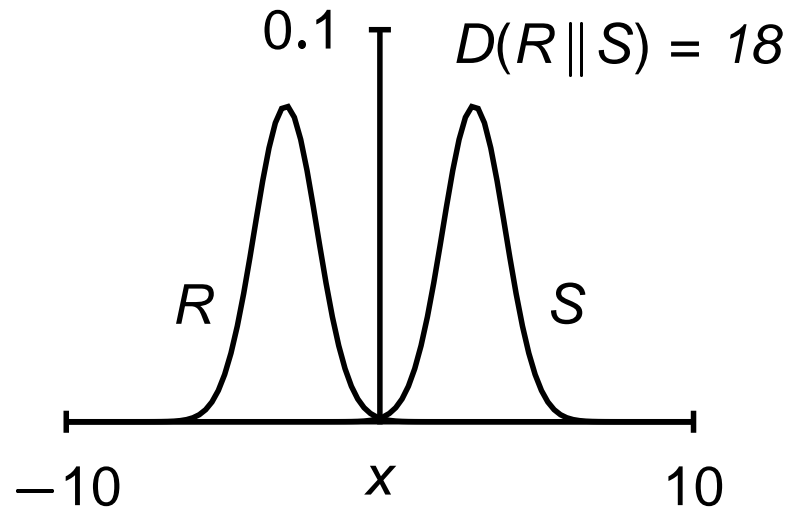
K.-L. Distance of Normal Distributions

$$R \sim \text{Normal}(\mu_R, \sigma_R^2)$$

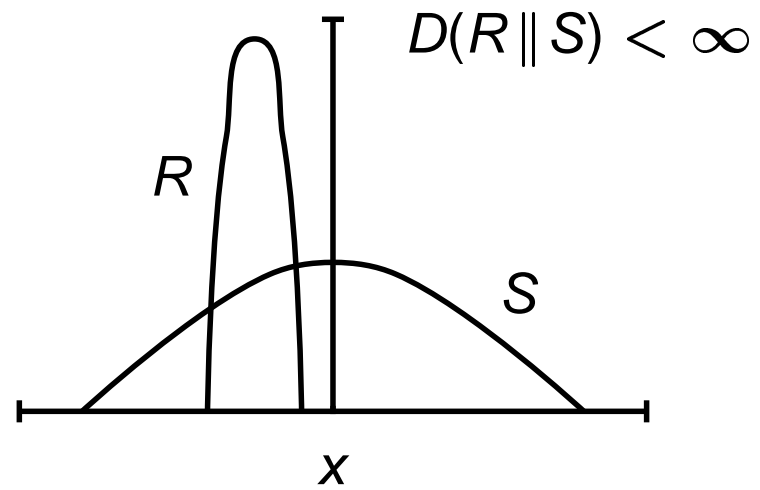
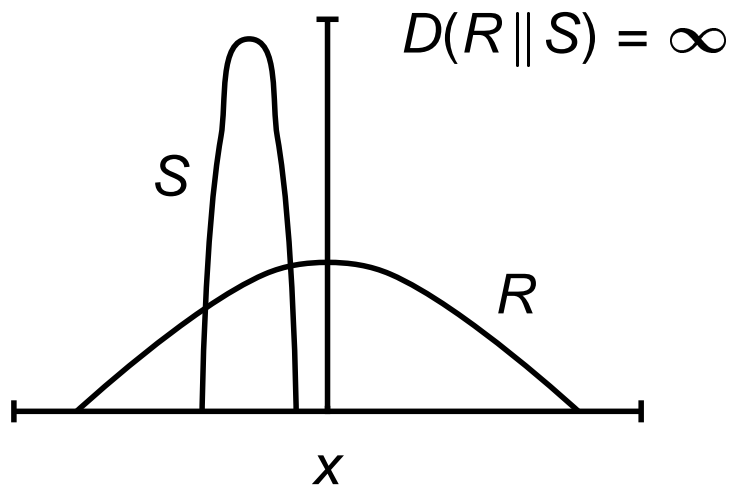
$$S \sim \text{Normal}(\mu_S, \sigma_S^2)$$

$$D(R \parallel S) = \frac{1}{2} \left(\frac{\sigma_R^2}{\sigma_S^2} - \log \frac{\sigma_R^2}{\sigma_S^2} - 1 \right) + \frac{1}{2} \frac{(\mu_R - \mu_S)^2}{\sigma_S^2}$$

K.-L. Distance of Normal Distributions



Distributions with Different Support



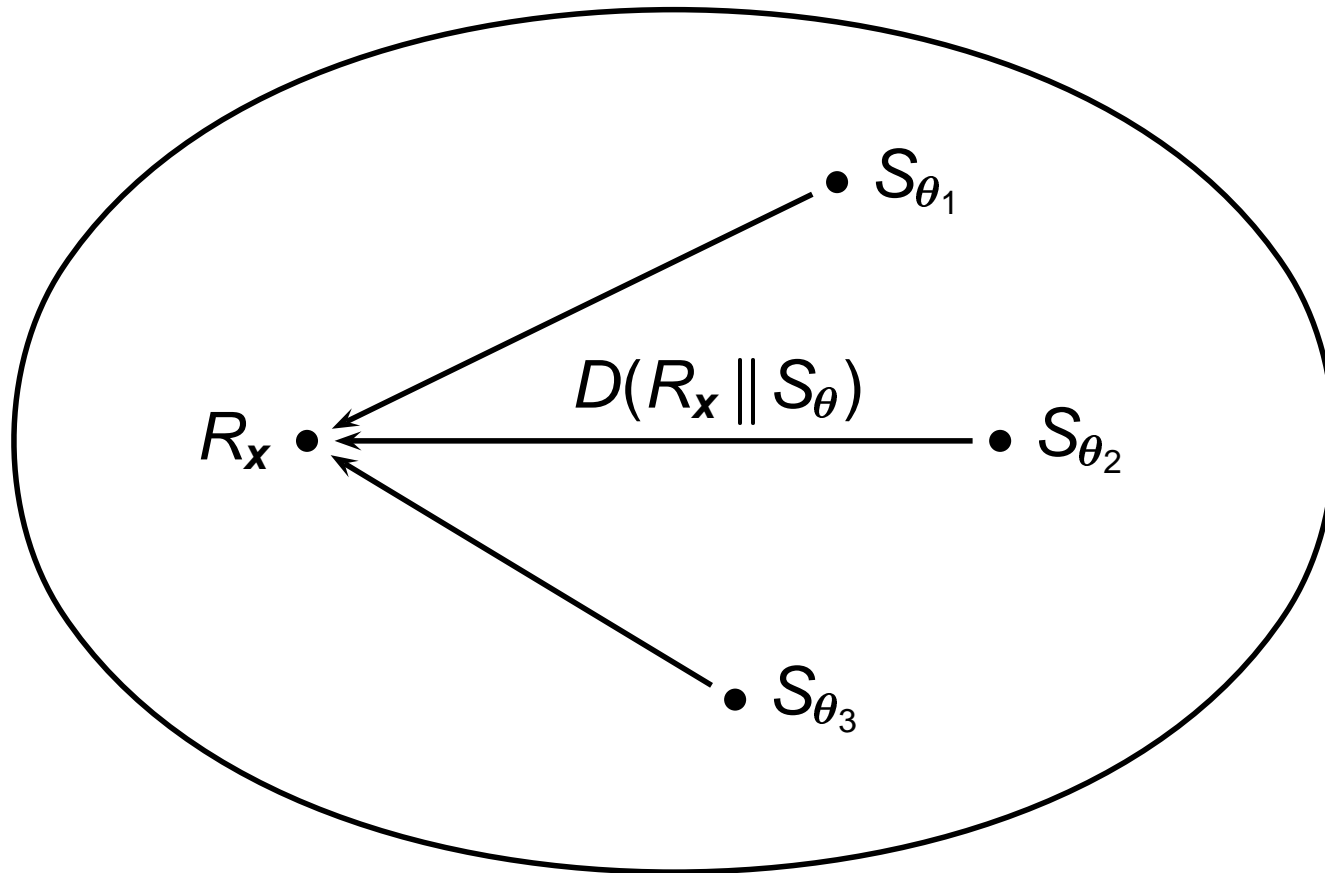
Joint Distribution of Sample

$$S_{\theta}^t(\mathbf{x}) = \underbrace{C(\mathbf{x})}_{\substack{\theta\text{-indep.} \\ \text{factor}}} \exp\left(-\underbrace{t}_{\substack{\text{amount} \\ \text{of Data}}} D\left(\underbrace{R_{\mathbf{x}}}_{\substack{\text{empirical} \\ \text{distribution}}} \parallel \underbrace{S_{\theta}}_{\substack{\text{model} \\ \text{distribution}}}\right)\right)$$

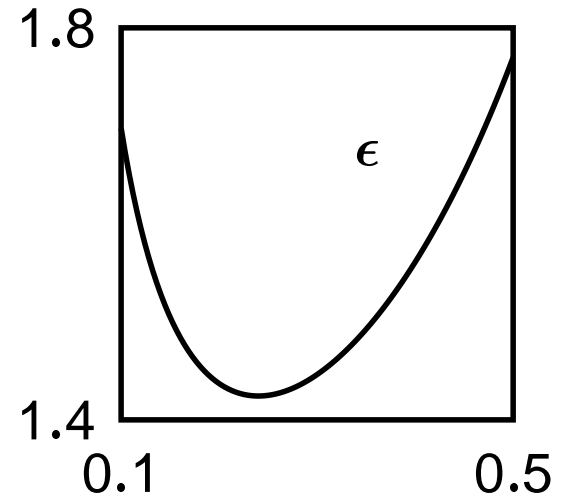
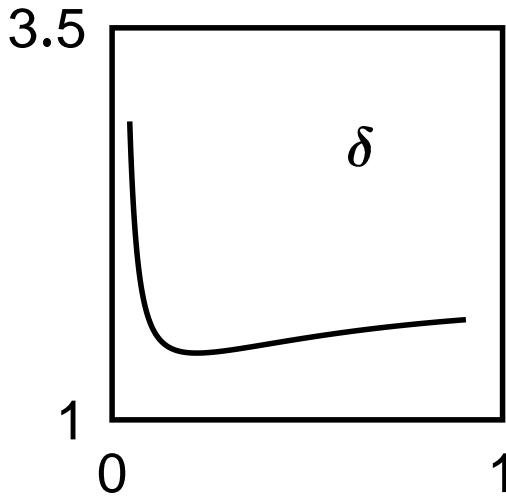
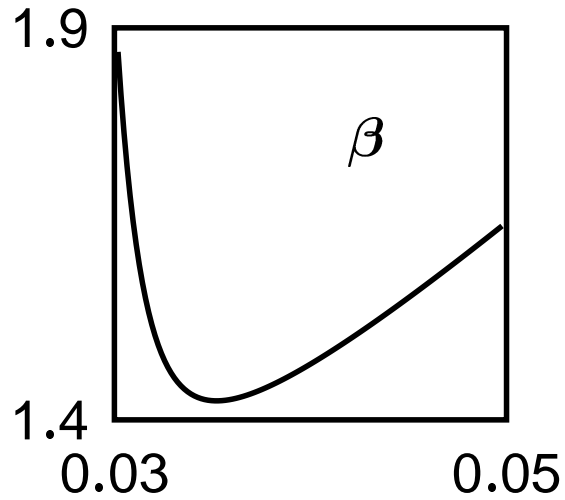
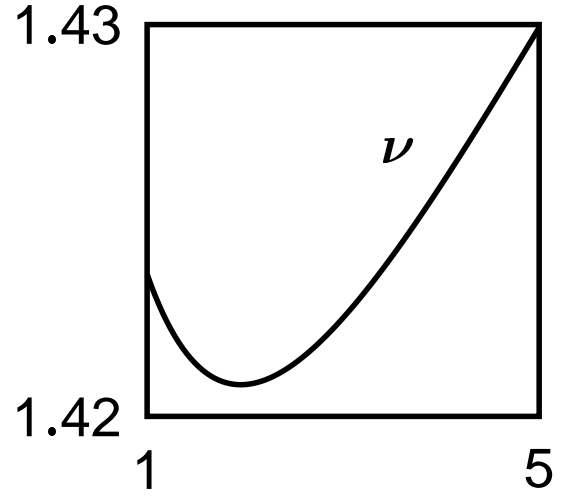
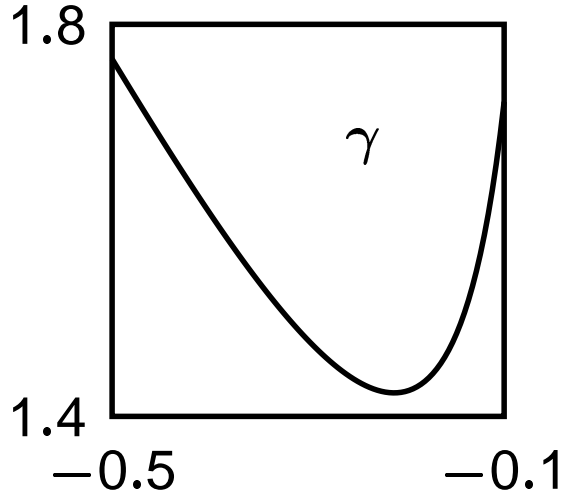
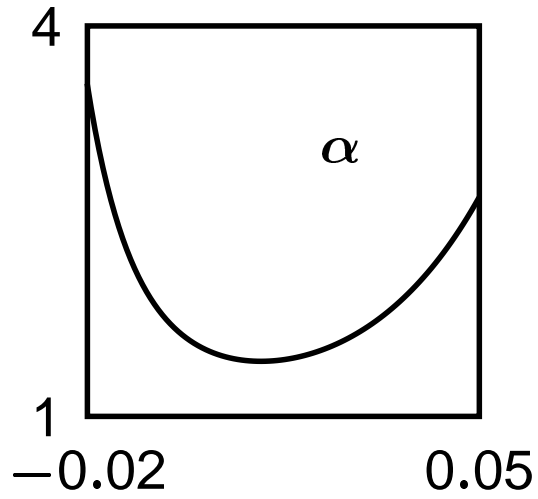
Posterior Distribution of Parameter

$$P_x(\theta) \propto \underbrace{P(\theta)}_{\substack{\text{prior} \\ \text{distribution}}} \exp\left(-\underbrace{t}_{\substack{\text{amount} \\ \text{of Data}}} D\left(\underbrace{R_x}_{\substack{\text{empirical} \\ \text{distribution}}} \parallel \underbrace{S_\theta}_{\substack{\text{model} \\ \text{distribution}}}\right)\right)$$

Minimum Distance Estimation



$$D(R_x \parallel \epsilon \text{ Normal}(\alpha, \beta^2) + (1 - \epsilon) \text{ Log Student}(\gamma, \delta^2, \nu))$$



K.-L. Distance and Likelihood

$$D(R_{\mathbf{x}} \parallel S_{\theta}) = C(\mathbf{x}) - \frac{1}{t} \log S_{\theta}^t(\mathbf{x})$$



$$\min_{\theta} D(R_{\mathbf{x}} \parallel S_{\theta}) \Leftrightarrow \max_{\theta} S_{\theta}^t(\mathbf{x})$$

minimum
distance

maximum
likelihood

Minimum Distance Estimation

"true" parametric distribution	empirical distribution
$\min_{\theta} D(S_{\theta_0} \ S_{\theta})$	$D(R_x \ S_{\theta})$
\Updownarrow	↓
$\max_{\theta} E_{S_{\theta_0}}(\log S_{\theta})$	$R_x \rightarrow S_{\theta_0}$
$E_{S_{\theta_0}}(\cdot) \approx E_{R_x}(\cdot)$	$D(S_{\theta_0} \ S_{\theta})$

***What to do when
the empirical distribution is
only partially known?***

Data Statistic

$$E_{R_x}(h) = \frac{1}{t} \sum_{k=1}^t h(x_k) \triangleq \xi \quad \text{given } h : \mathcal{X} \rightarrow \mathbf{R}^n$$

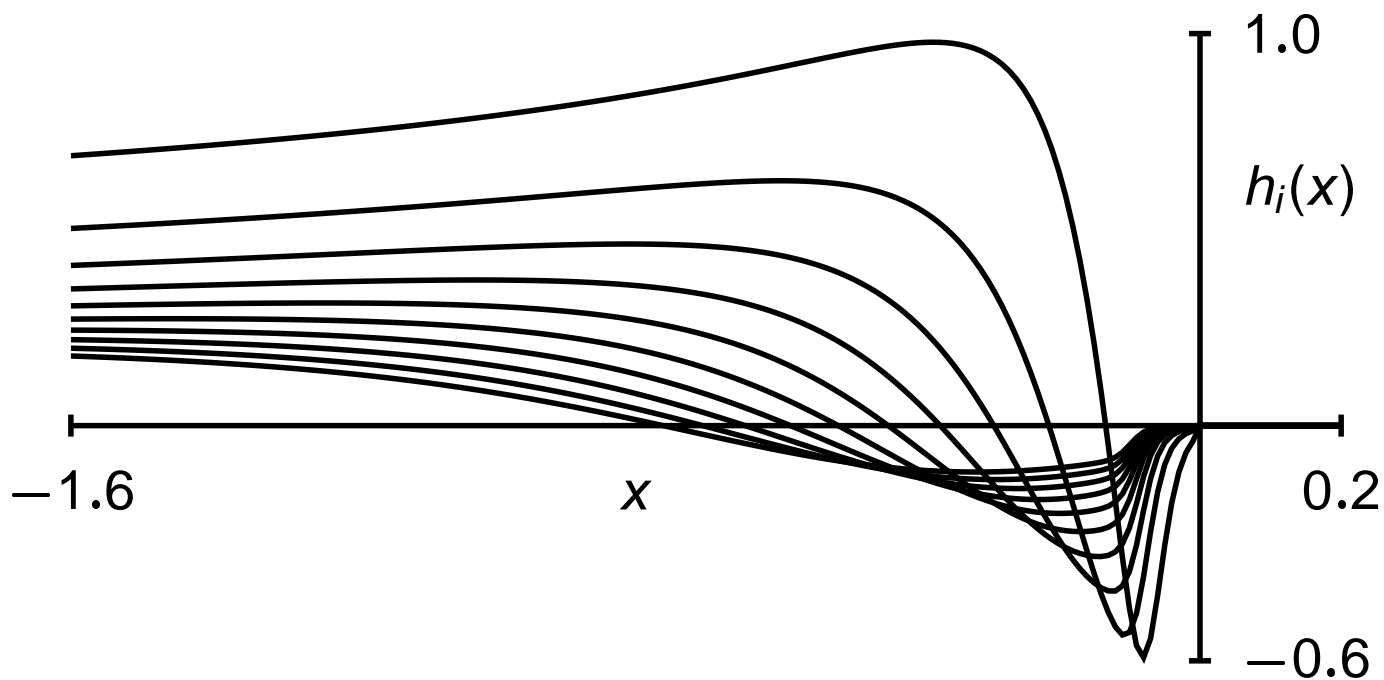
implies ambiguity

$$R_x \in \mathcal{R}_\xi \triangleq \{R : E_R(h) = \xi\}$$

BALLOON: Choice of Statistic

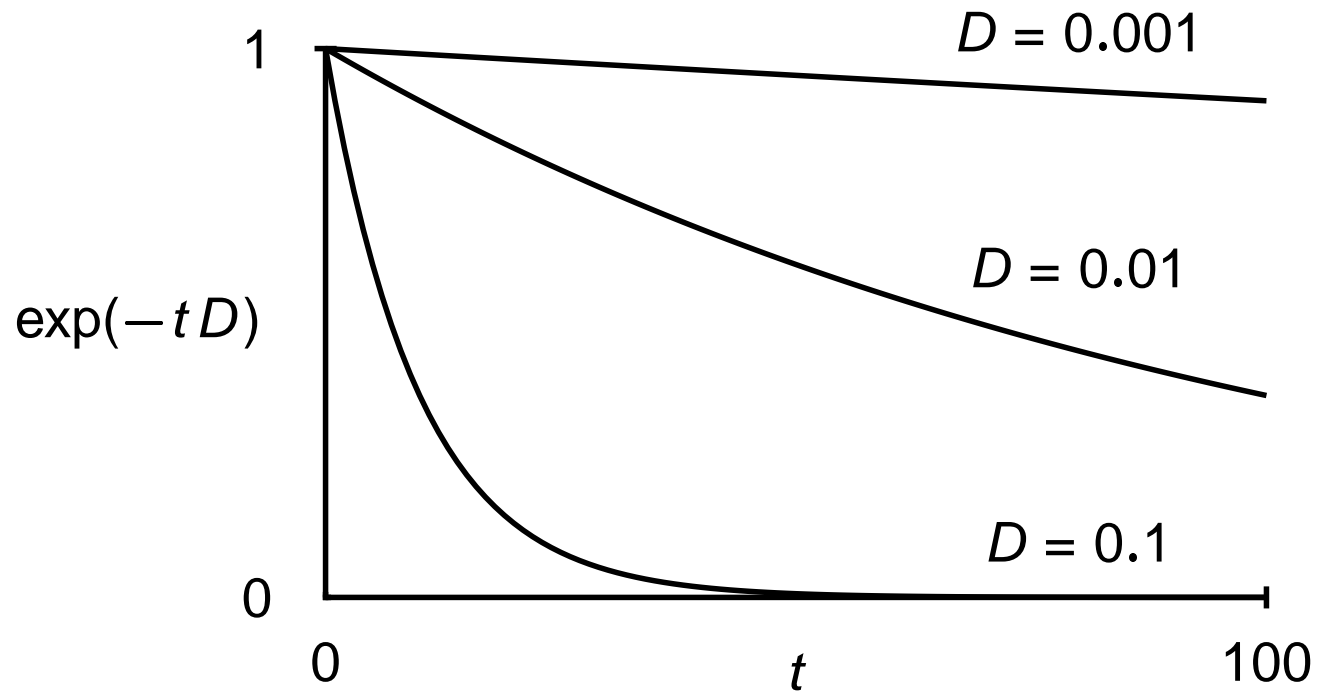
$$h_i(x) = \log S_{\gamma_i}(x) - \log S_{\gamma_{i+1}}(x), \quad i = 1, \dots, n$$

$$E_{R_x}(h_i) = \log S_{\gamma_i}^t(\mathbf{x}) - \log S_{\gamma_{i+1}}^t(\mathbf{x}), \quad i = 1, \dots, n$$



Probability of Large Deviations

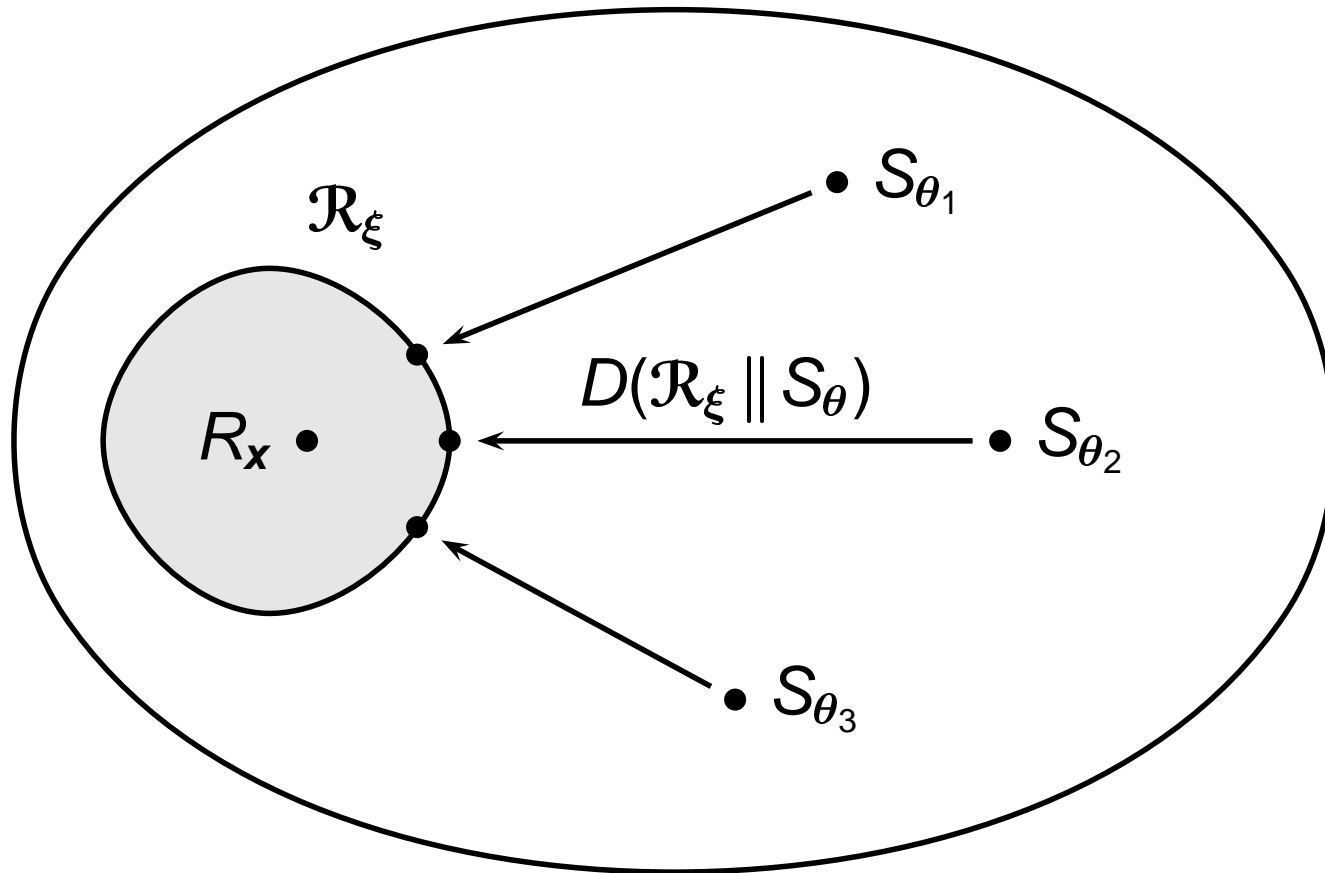
$$S_{\theta}^t(R_x \in \mathcal{R}_{\xi}) = \exp\left(-t \min_{R \in \mathcal{R}_{\xi}} D(R \| S_{\theta})\right) \exp(-t o(1))$$



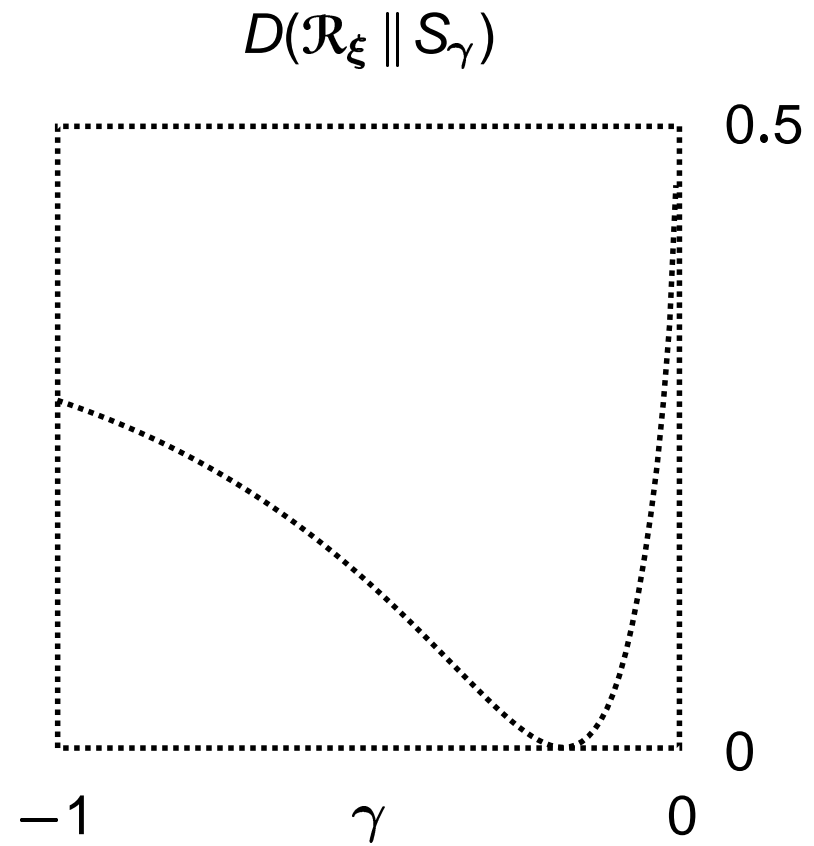
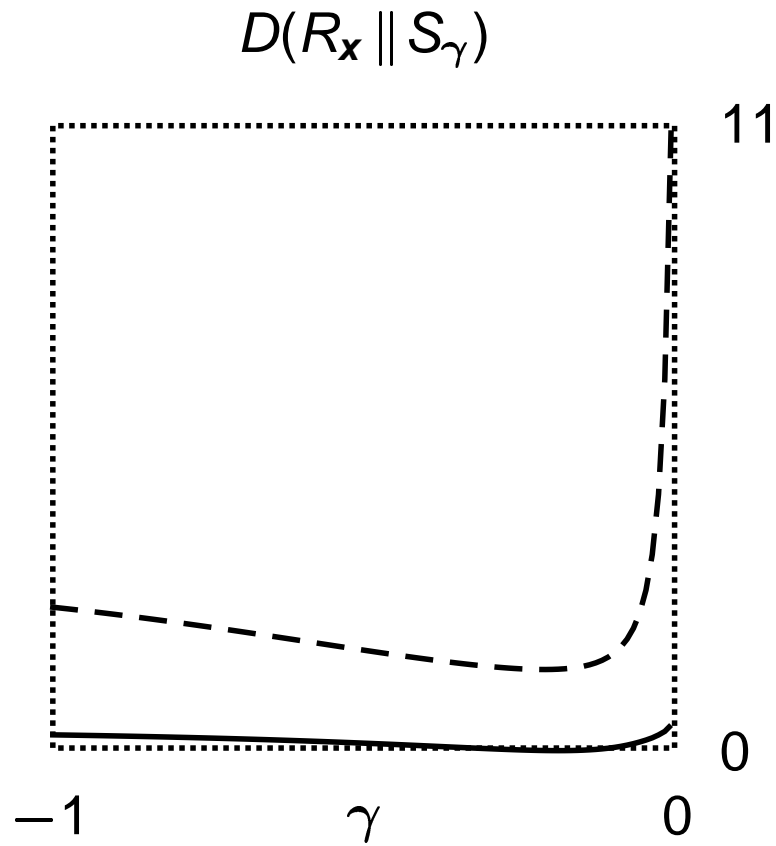
Joint Distribution of a Sample Set

$$S_{\theta}^t(R_x \in \mathcal{R}_{\xi}) \doteq \exp\left(- \underbrace{t}_{\substack{\text{amount} \\ \text{of Data}}} D\left(\underbrace{\mathcal{R}_{\xi}}_{\substack{\text{set of} \\ \text{empirical} \\ \text{distributions}}} \parallel \underbrace{S_{\theta}}_{\substack{\text{model} \\ \text{distribution}}}\right)\right)$$

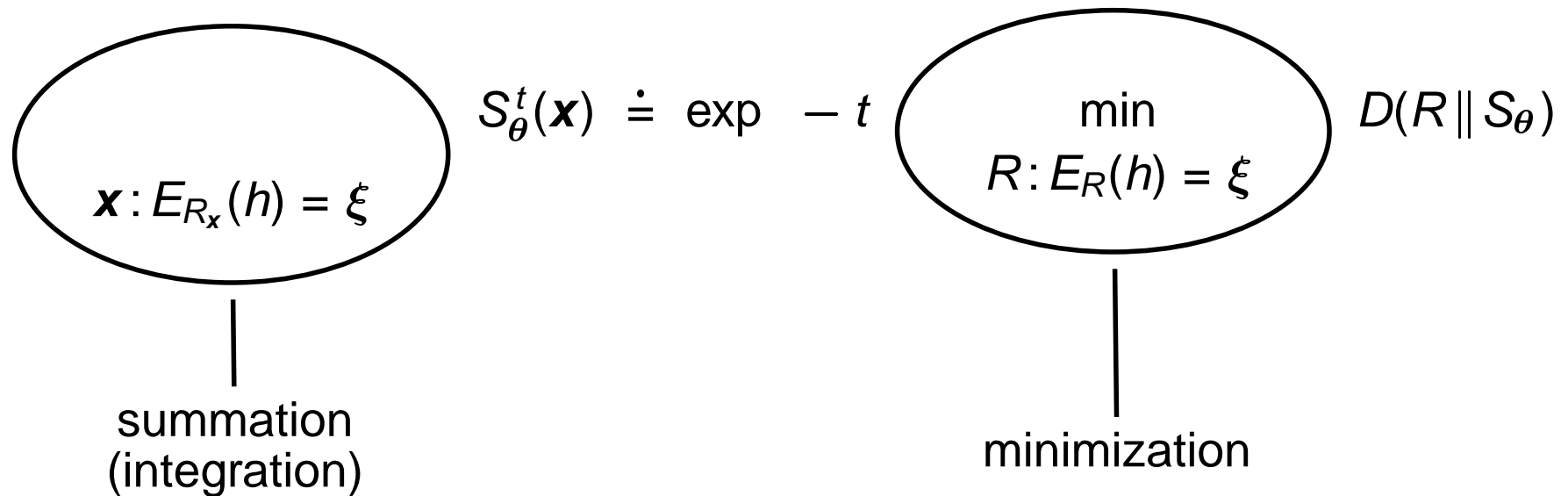
Minimum Distance Estimation



BALLOON: Minimum K.-L. Distance



Why the Large Deviation Approximation?



a significant drop in complexity!

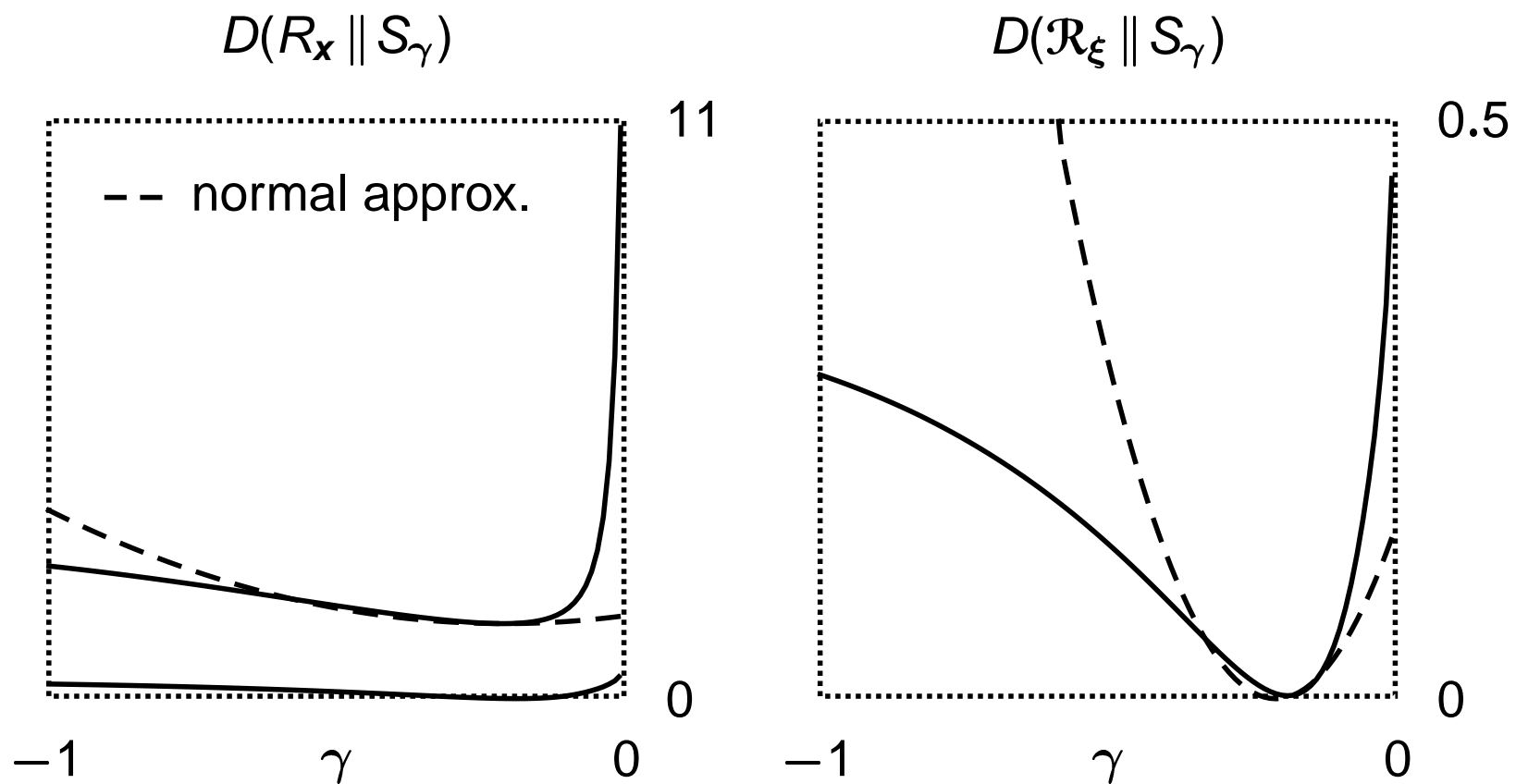
Large Deviation Approximation

$$S_{\theta}^t(R_x \in \mathcal{R}_{\xi}) = \exp\left(-t \min_{R \in \mathcal{R}_{\xi}} D(R \| S_{\theta})\right) \exp\left(-t o(1)\right)$$

Normal Approximation

$$S_{\theta}^t(\mathbf{x}) = \exp\left(-t D(R_x \| S_{\hat{\theta}})\right) \exp\left(-t \frac{1}{2} I(\hat{\theta}) (\theta - \hat{\theta})^2\right) \\ \cdot \exp\left(-t O(|\theta - \hat{\theta}|^3)\right)$$

BALLOON: Minimum K.-L. Distance





***How to deal with
dependent Data?***

Example 2: SUNSPOTS

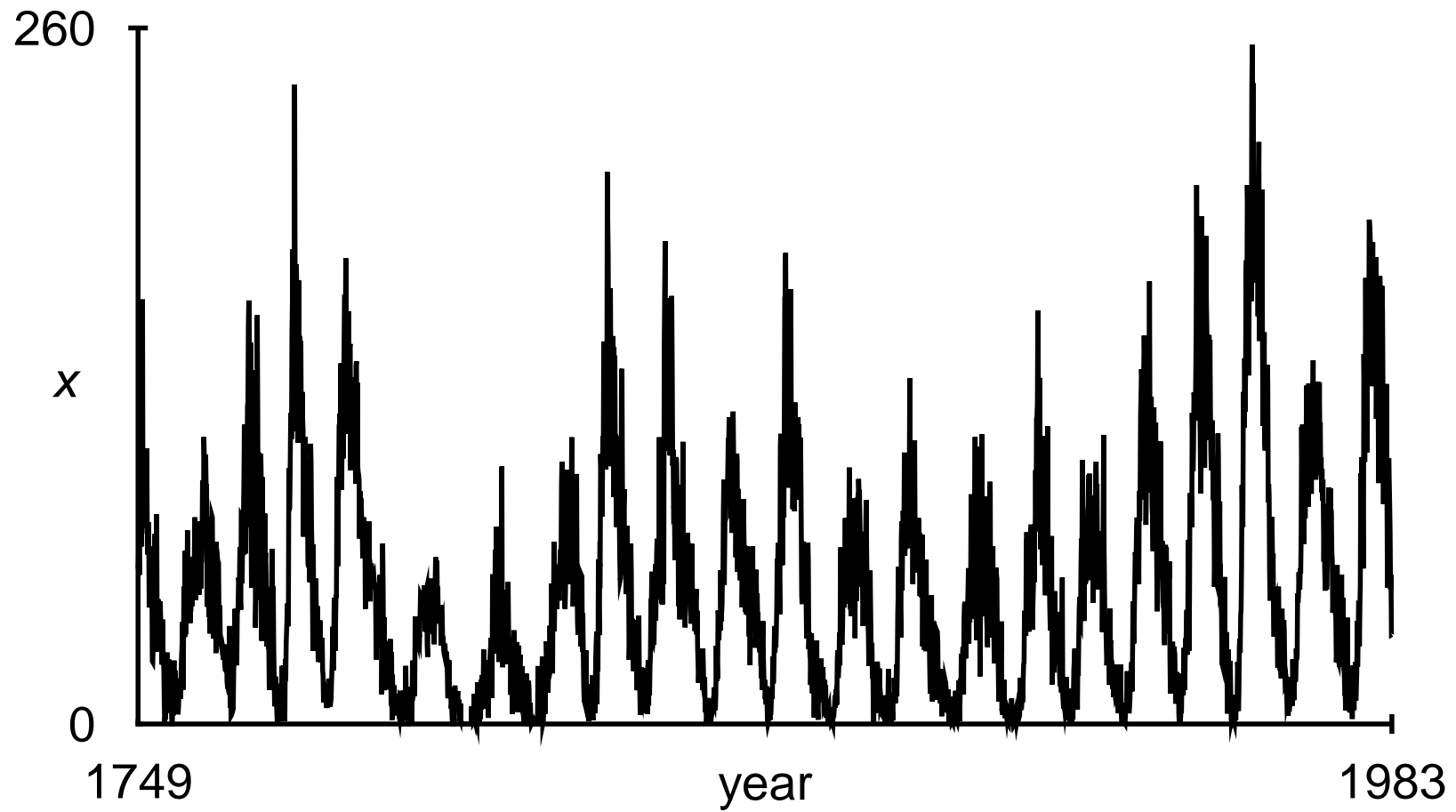
Zurich monthly sunspot numbers in 1749–1983

Source:

<http://lib.stat.cmu.edu/Datasets/Andrews/T11.1>

Andrews and Herzberg, DATA

SUNSPOTS: Time Series

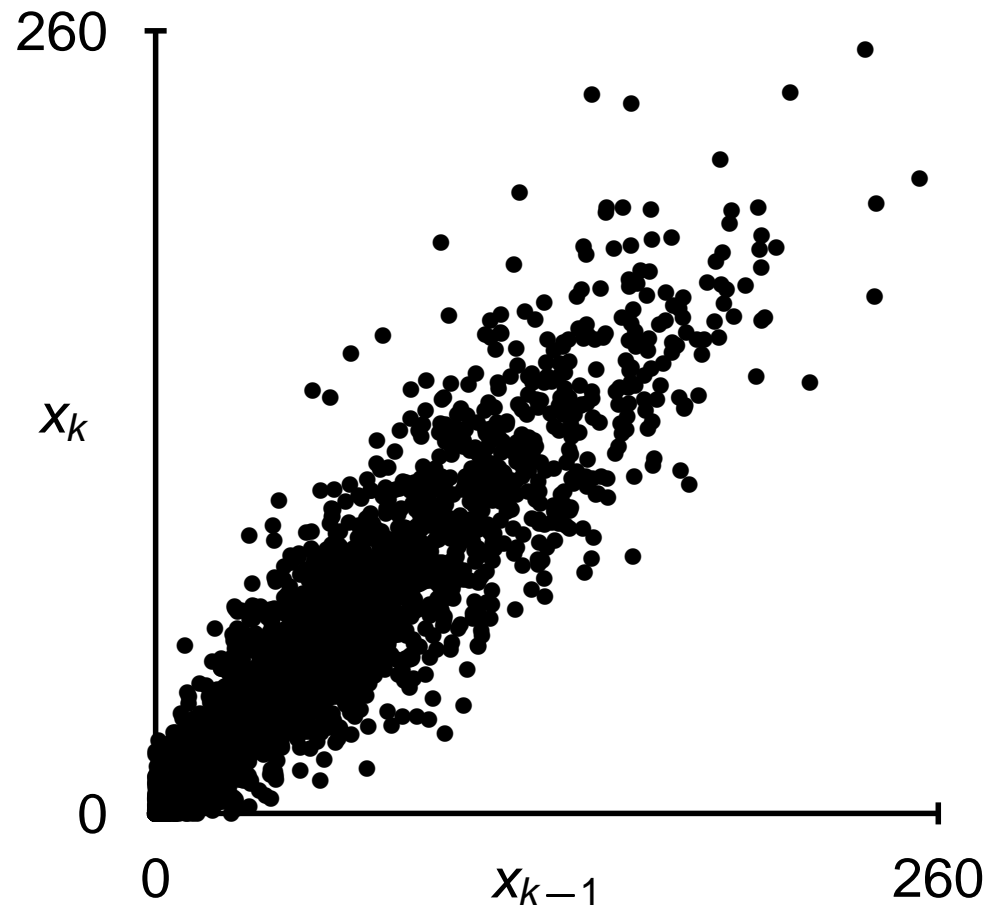


Joint Distribution of Sample

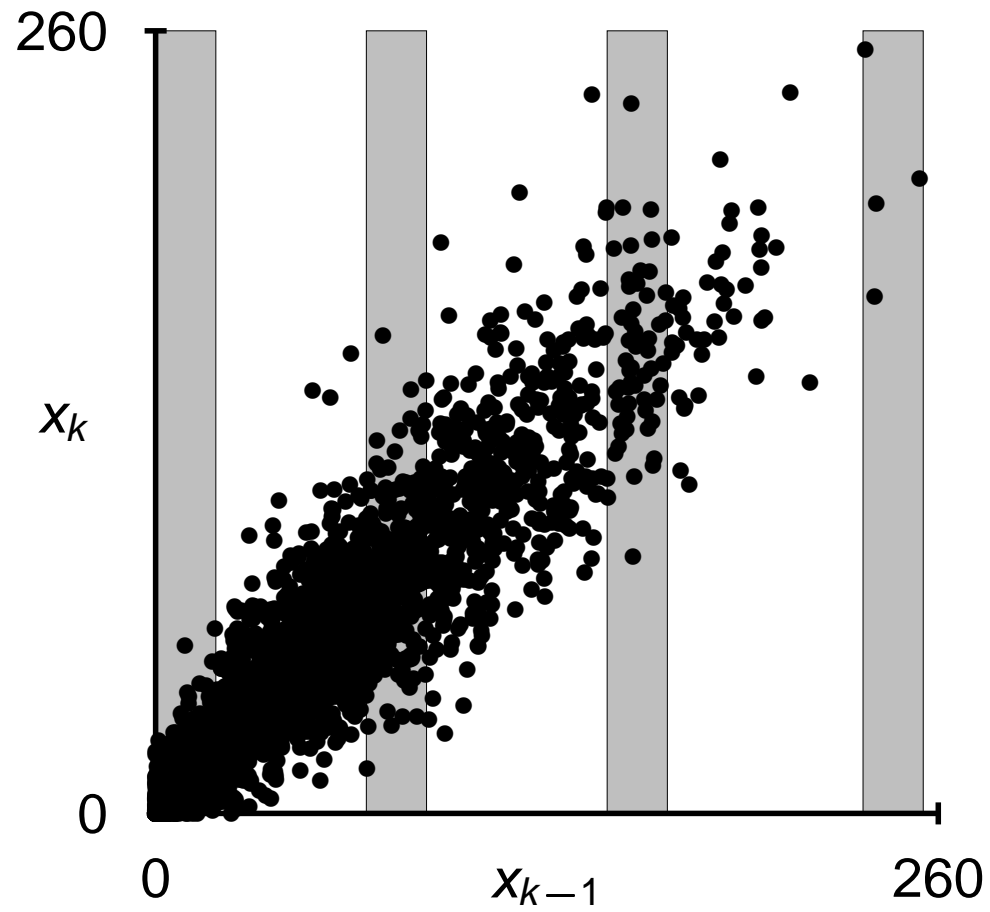
model distribution	$S_{\theta}(x_k x_{k-1}, \dots, x_{k-m})$
empirical distribution	$R_x(x_k, x_{k-1}, \dots, x_{k-m})$
“weighting” distribution	$W(x_{k-1}, \dots, x_{k-m})$

$$\begin{aligned} S_{\theta}^{t-m}(x_t, \dots, x_{m+1} | x_m, \dots, x_1) &= \\ &= C(\mathbf{x}) \exp\left(- (t-m) \min_W D(R_x \| S_{\theta} W)\right) \end{aligned}$$

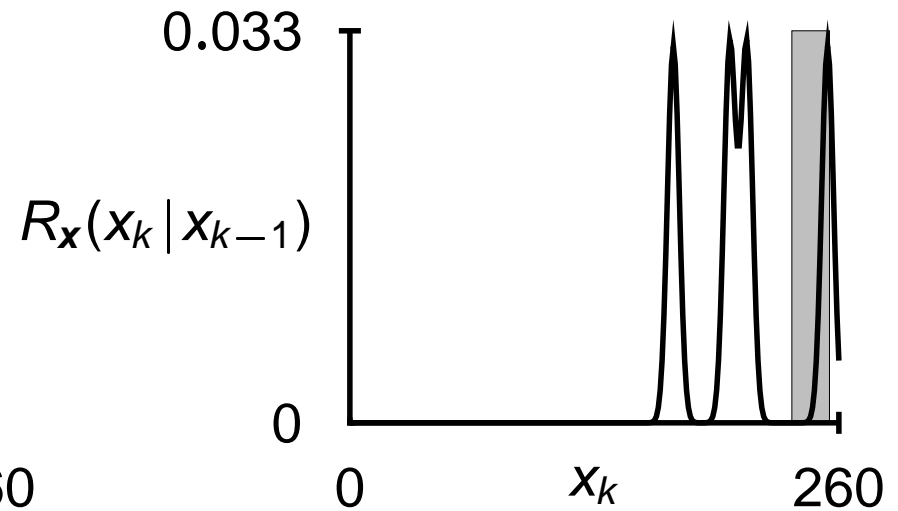
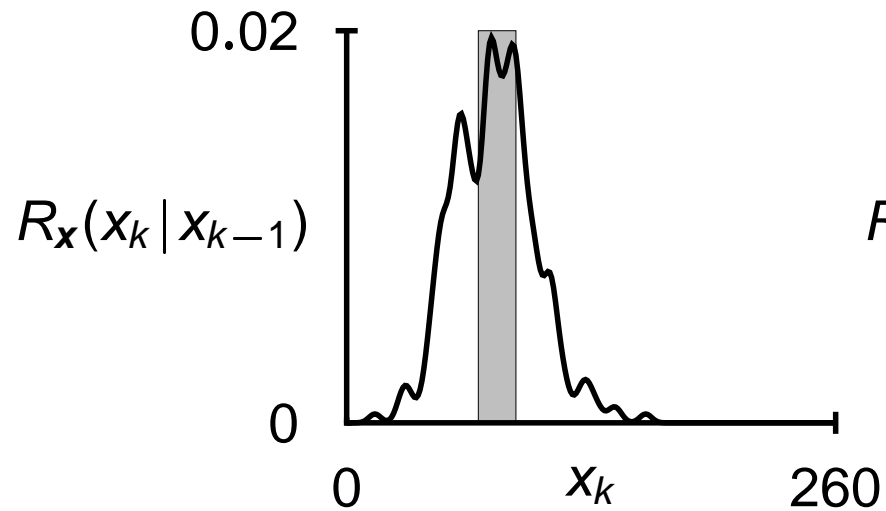
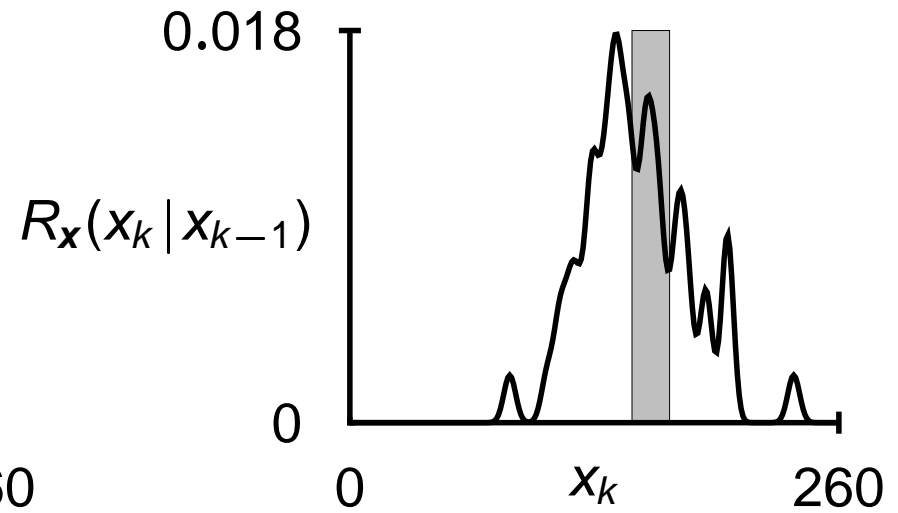
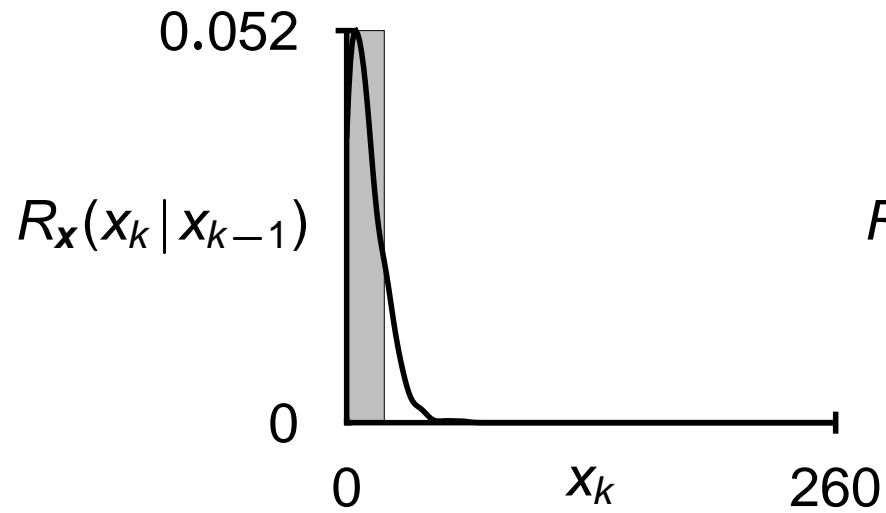
SUNSPOTS: Empirical Distribution



SUNSPOTS: Empirical Distribution



SUNSPOTS: Conditional Empirical Distrib's



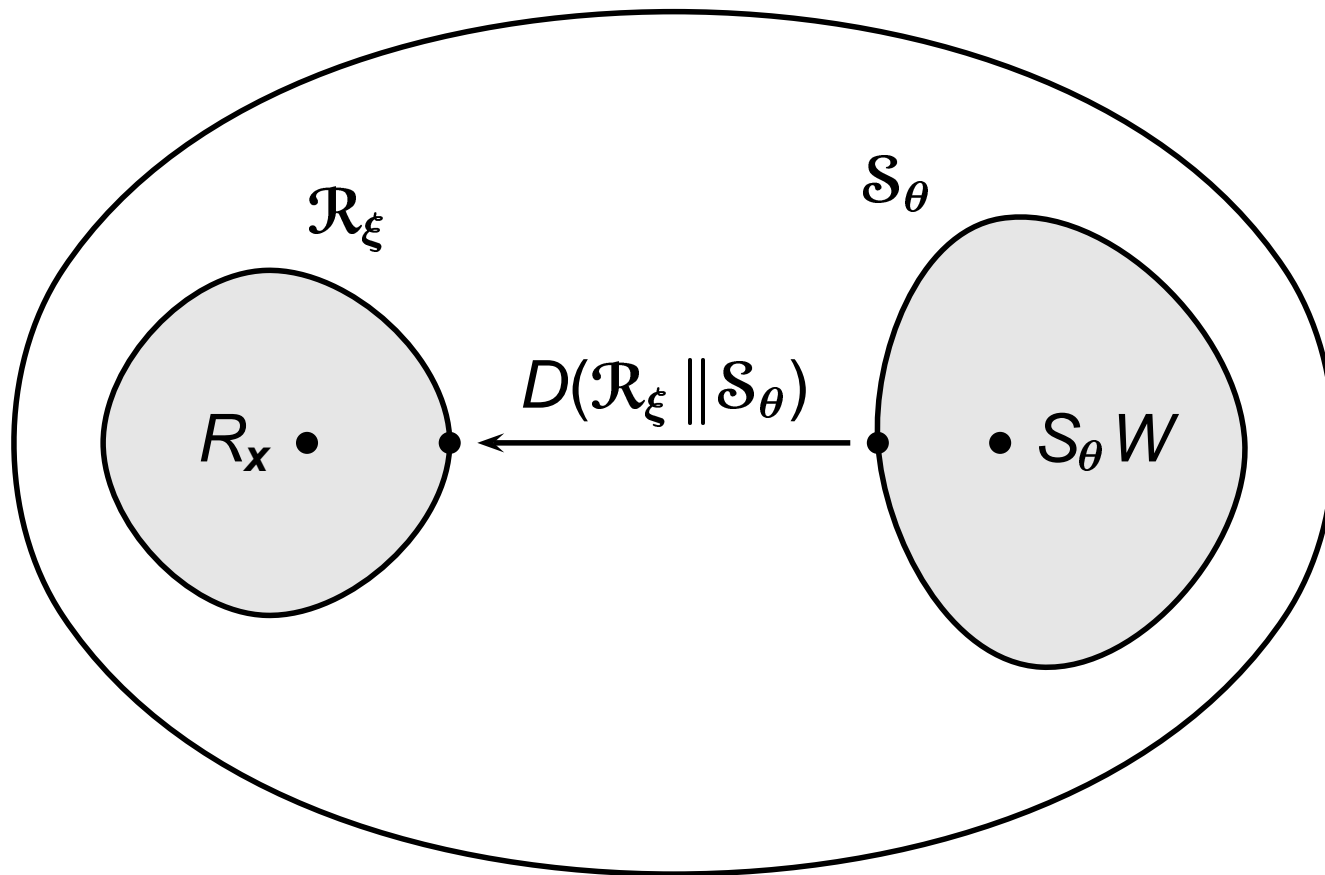
Probability of Large Deviations

$$\text{given } E_{R_x}(h) = \frac{1}{t} \sum_{k=m+1}^t h(x_k, x_{k-1}, \dots, x_{k-m}) \triangleq \xi$$

where $h : \mathcal{X}^{m+1} \rightarrow \mathbf{R}^n$

$$S_{\theta}^t(R_x \in \mathcal{R}_{\xi}) \doteq \exp\left(-t \min_{R \in \mathcal{R}_{\xi}} \min_W D(R \| S_{\theta} W)\right)$$

Minimum Distance Estimation





***What are the
pros and cons?***

K.-L. Distance

- minimum distance
 - ⇒ optimization

☞ *computationally easier*

Probability

- marginal probability
 - ⇒ summation
(integration)

K.-L. Distance

□ approximation theory

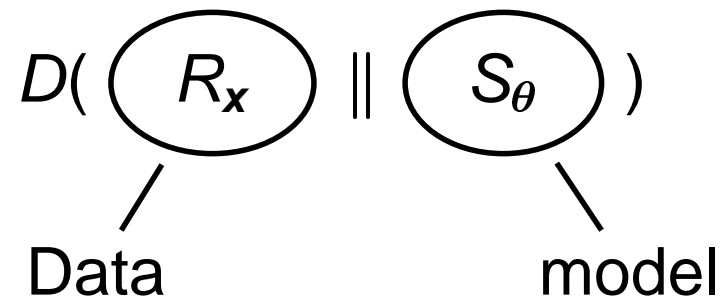
☞ *conceptually simpler*

☞ *ready for further approximation*

Probability

□ probability theory & statistics

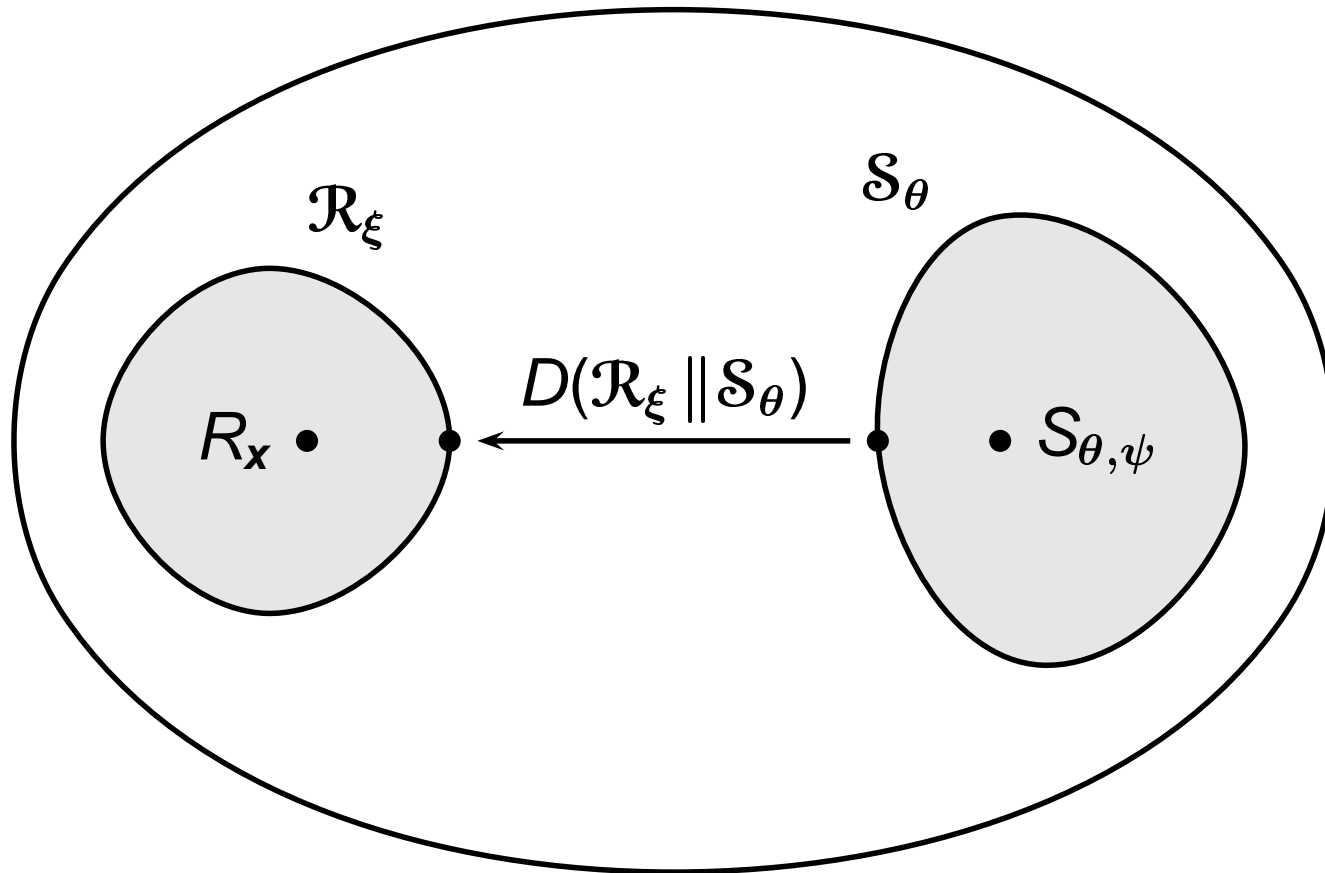
Symmetric Role of Data & Model



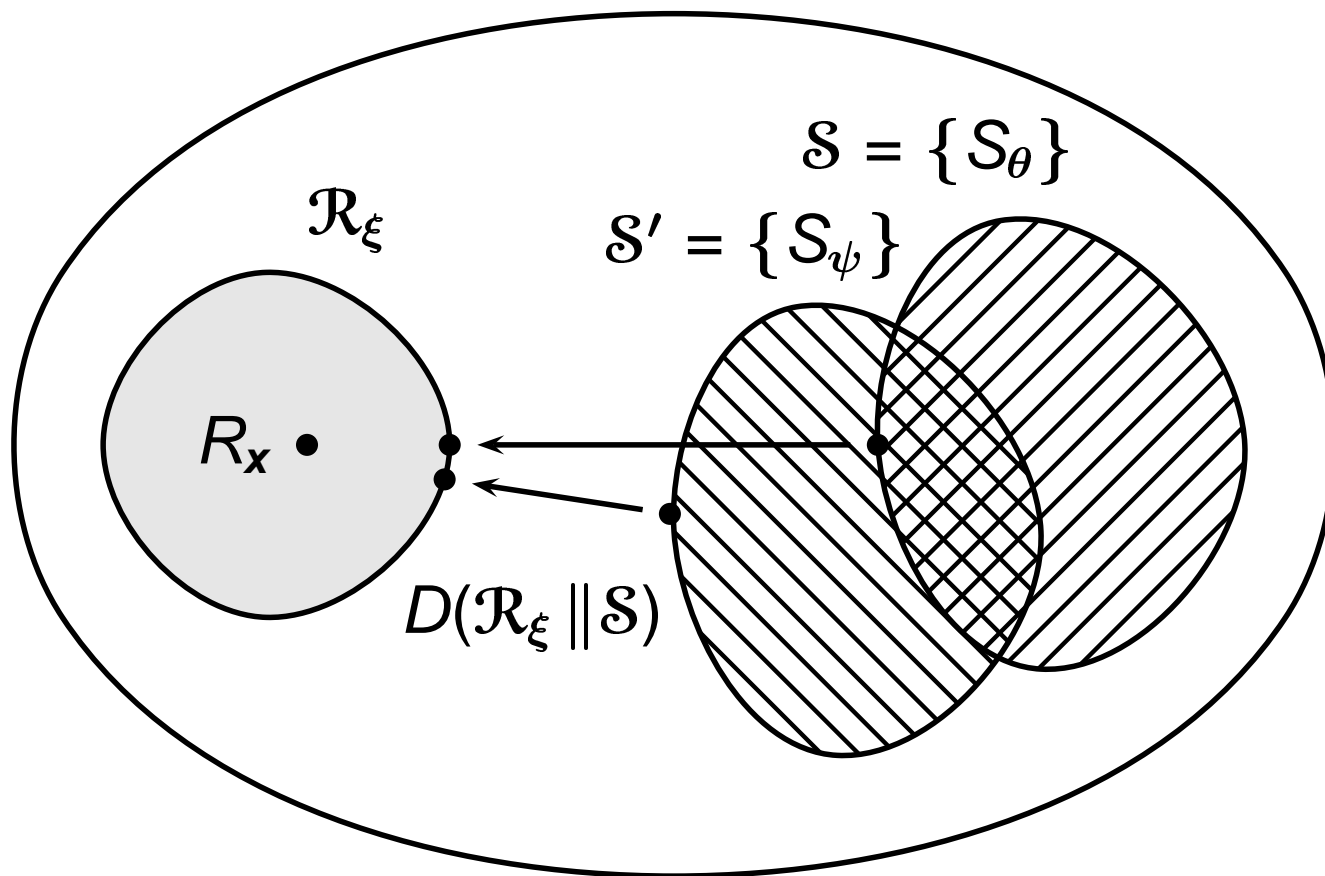
admits straightforward extensions

$$D(\{ R : E_R(h) = \xi \} \parallel \{ S_{\theta, \psi} : \psi \in \mathcal{N} \})$$

Robust Identification



Adaptive Identification



K.-L. Distance

- *absolute* measure of goodness-of-fit

$$D(R_{\mathbf{x}} \parallel S_{\theta}) \Leftrightarrow S_{\theta}^t(\mathbf{x})$$

$$D(\mathcal{R}_{\xi} \parallel S_{\theta}) \Leftrightarrow S_{\theta}^t(\mathcal{R}_{\xi})$$

Posterior Prob.

- *relative* measure of goodness-of-fit

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

$$P_{\xi}(S_{\theta}) \propto P(\theta) S_{\theta}^t(\mathcal{R}_{\xi})$$

A Big Open Issue

approximation of Kullback-Leibler distance

$$D(\mathcal{R}_\xi \parallel S_\theta) \doteq ?$$

- $D(\mathcal{R}_\xi \parallel S_\theta) \doteq D(\mathcal{R}_\xi \parallel \hat{S}_\theta)$, $\{\hat{S}_\theta\}$ exponential
- $D(\mathcal{R}_\xi \parallel S_\theta) \doteq \|E_{R_x}(h) - E_{S_\theta}(h)\|_W^2$