

Data-driven distributionally robust optimization: Classification of ambiguity sets

Tobias Sutter¹, Bart van Parys² & Daniel Kuhn³

Data Driven Combinatorial Optimization,
Ecole des Ponts, Paris, October 4, 2023

¹University of Konstanz

²MIT Sloan School of Management

³Risk Analytics and Optimization Chair, EPFL

Universität
Konstanz



EPFL

MIT
MANAGEMENT
SLOAN SCHOOL

arXiv:2010.06606

Data-driven decision making

Stochastic
optimization

$$\min_x c(x, \theta)$$

Family of prob.
measures

$$\{\mathbb{P}_\theta : \theta \in \Theta\}$$

Data-gen.
process

$$\{\xi_t\}_{t \in \mathbb{N}}$$

Data-driven decision making

Stochastic
optimization

$$\min_x c(x, \theta)$$

Family of prob.
measures

$$\{\mathbb{P}_\theta : \theta \in \Theta\}$$

Data-gen.
process

$$\{\xi_t\}_{t \in \mathbb{N}}$$

Examples:

- ▶ Expected loss $c(x, \theta) = \mathbb{E}_\theta[\ell(x, \xi)]$
- ▶ Risk of loss $c(x, \theta) = \rho_\theta[\ell(x, \xi)]$
- ▶ Covariate information $c(x, \theta) = \mathbb{E}_\theta[\ell(x, \xi) | C\xi \in B]$
- ▶ Long-run average loss $c(x, \theta) = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_\theta[\ell(\pi_x(s_t), s_t)]$

Data-driven decision making

Stochastic
optimization

$$\min_x c(x, \theta)$$

Family of prob.
measures

$$\{\mathbb{P}_\theta : \theta \in \Theta\}$$

Data-gen.
process

$$\{\xi_t\}_{t \in \mathbb{N}}$$

Assumptions:

- ▶ All measures defined on (Ω, \mathcal{F})
- ▶ $\Theta \subseteq \mathbb{R}^d$ open

Data-driven decision making

Stochastic
optimization

$$\min_x c(x, \theta)$$

Family of prob.
measures

$$\{\mathbb{P}_\theta : \theta \in \Theta\}$$

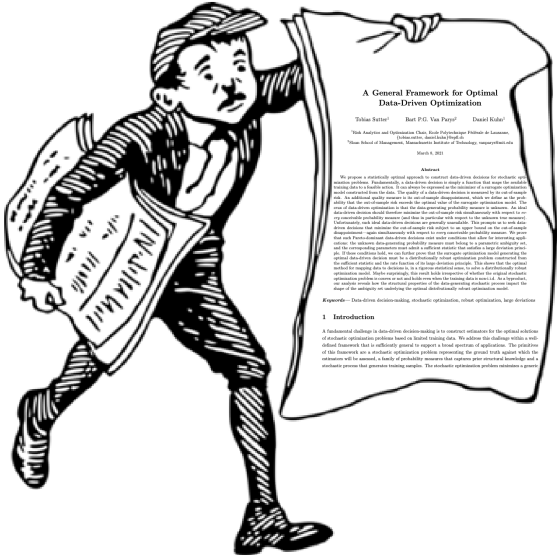
Data-gen.
process

$$\{\xi_t\}_{t \in \mathbb{N}}$$

Examples:

- ▶ Finite-state i.i.d. processes
- ▶ Finite-state Markov chains
- ▶ Vector-autoregressive processes
- ▶ I.i.d. processes with parametric distribution function

Motivating example — newsvendor problem



News vendor problem

Stochastic
optimization

$$\min_{x \in X} c(x, \theta)$$

- ▶ Order quantities $x \in X = \{1, 2, \dots, d\}$
- ▶ Demand $\xi \in \Xi = \{1, 2, \dots, d\}$
- ▶ Objective $c(x, \theta) = \mathbb{E}_{\theta}[kx - p \min\{x, \xi\}]$

Data-gen.
process

$$\{\xi_t\}_{t \in \mathbb{N}}$$

- ▶ Historical demand $\xi_t \in \Xi$

Family of prob.
measures

$$\{\mathbb{P}_{\theta} : \theta \in \Theta\}$$

- ▶ $\{\xi_t\}_{t \in \mathbb{N}}$ i.i.d. process under \mathbb{P}_{θ}
- ▶ $\mathbb{P}_{\theta}(\xi_t = i) = \theta_i$ for $i \in \Xi$

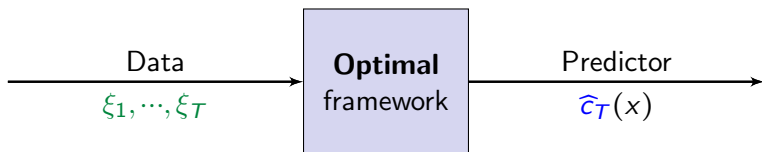
Surrogate optimization models

Original optimization problem

$$\min_{x \in X} c(x, \theta)$$

Surrogate optimization problem

$$\min_{x \in X} \widehat{c}_T(x)$$



Surrogate optimization models

Original optimization problem

$$\min_{x \in X} c(x, \theta)$$

Surrogate optimization problem

$$\min_{x \in X} \widehat{c}_T(x)$$

Construction of \widehat{c}_T

- ▶ Sample average approximation¹
- ▶ Predict-then-optimize approach²
- ▶ Neural network model³
- ▶ Distributionally robust optimization model⁴
- ▶ etc.

¹Shapiro, Annals of Statistics, 1989; ²Elmachtoub & Grigas, Management Science, 2021; ³Donti et al., NeurIPS, 2017; ⁴Delage & Ye, Operations Research, 2010

Terminology

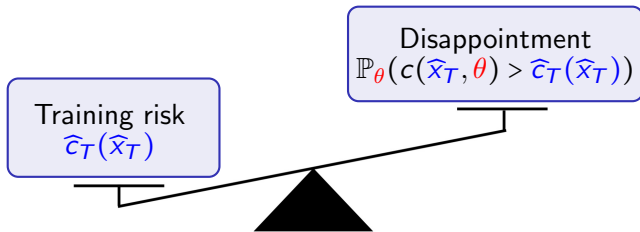
As a function of the **training data** ξ_1, \dots, ξ_T we denote

- ▶ Data-driven predictor \widehat{c}_T
- ▶ Data-driven prescriptor $\widehat{x}_T = \arg \min_{x \in X} \widehat{c}_T(x)$

Performance measures

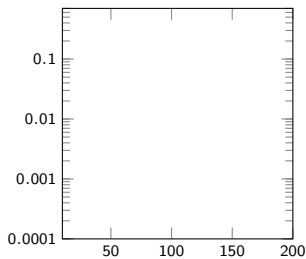
- ① In-sample (training) risk $\widehat{c}_T(\widehat{x}_T)$
- ② Out-of-sample (generalization) risk $c(\widehat{x}_T, \theta)$
- ③ Out-of-sample disappointment $\mathbb{P}_\theta(c(\widehat{x}_T, \theta) > \widehat{c}_T(\widehat{x}_T))$

A basic tradeoff

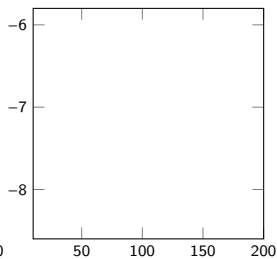


Data-driven newsvendor problem

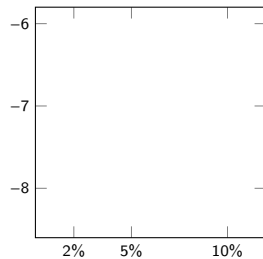
Disappointment probability
 $\mathbb{P}_{\theta^*} [c(\bar{x}_T, \theta^*) > \hat{c}_T(\bar{x}_T)]$



Conservatism
 $\mathbb{E}_{\theta^*} [\hat{c}_T(\bar{x}_T)]$

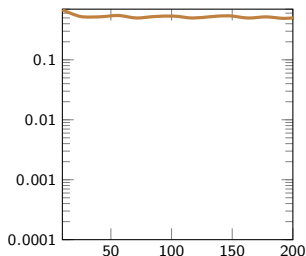


Conservatism vs
disappointment probability

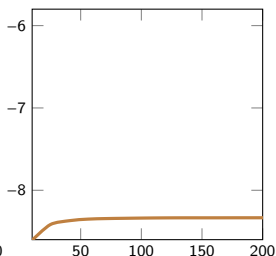


Data-driven newsvendor problem

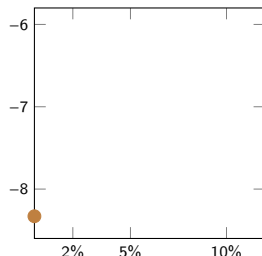
Disappointment probability
 $\mathbb{P}_{\theta^*} [c(\bar{x}_T, \theta^*) > \hat{c}_T(\bar{x}_T)]$



Conservatism
 $\mathbb{E}_{\theta^*} [\hat{c}_T(\bar{x}_T)]$



Conservatism vs
disappointment probability



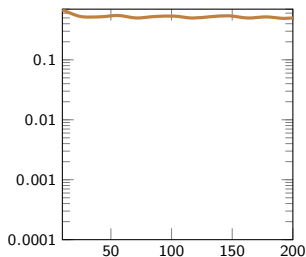
Model 1: SAA model²

$$\hat{c}_T(x) = c(x, \hat{\theta}_T)$$

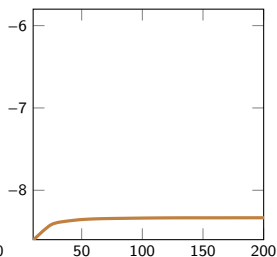
²Shapiro, Annals of Statistics, 1989

Data-driven newsvendor problem

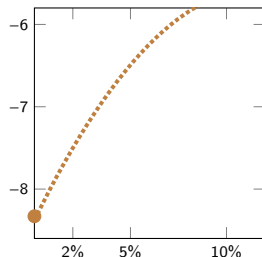
Disappointment probability
 $\mathbb{P}_{\theta^*} [c(\hat{x}_T, \theta^*) > \hat{c}_T(\hat{x}_T)]$



Conservatism
 $\mathbb{E}_{\theta^*} [\hat{c}_T(\hat{x}_T)]$



Conservatism vs
disappointment probability



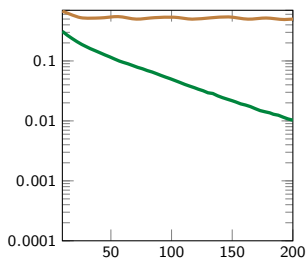
Model 1: SAA model²

$$\hat{c}_T(x) = c(x, \hat{\theta}_T) + r$$

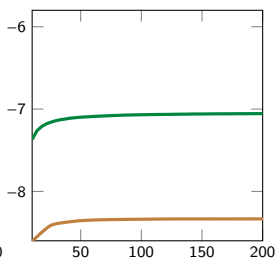
²Shapiro, Annals of Statistics, 1989

Data-driven newsvendor problem

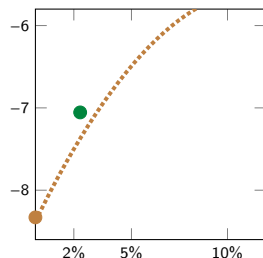
Disappointment probability
 $\mathbb{P}_{\theta^*} [c(\bar{x}_T, \theta^*) > \hat{c}_T(\bar{x}_T)]$



Conservatism
 $\mathbb{E}_{\theta^*} [\hat{c}_T(\bar{x}_T)]$



Conservatism vs
disappointment probability



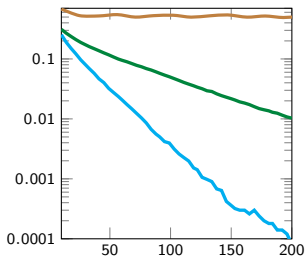
Model 2: DRO model with moment ambiguity set²

$$\hat{c}_T(x) = \sup_{\theta \in \Theta} \left\{ c(x, \theta) : |\mathbb{E}_{\hat{\theta}_T}[\xi^j] - \mathbb{E}_{\theta}[\xi^j]| \leq r \quad \forall j = 1, \dots, 4 \right\}$$

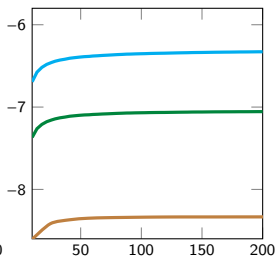
²Delage & Ye, Operations Research, 2010

Data-driven newsvendor problem

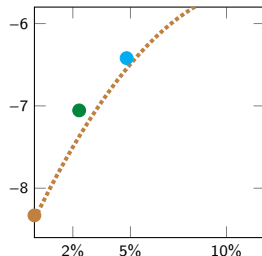
Disappointment probability
 $\mathbb{P}_{\theta^*} [c(\bar{x}_T, \theta^*) > \hat{c}_T(\bar{x}_T)]$



Conservatism
 $\mathbb{E}_{\theta^*} [\hat{c}_T(\bar{x}_T)]$



Conservatism vs
disappointment probability



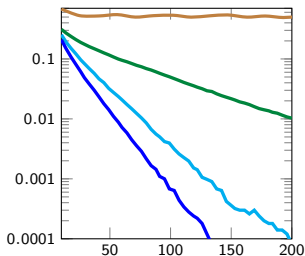
Model 2: DRO model with moment ambiguity set²

$$\hat{c}_T(x) = \sup_{\theta \in \Theta} \left\{ c(x, \theta) : |\mathbb{E}_{\hat{\theta}_T}[\xi^j] - \mathbb{E}_{\theta}[\xi^j]| \leq r \quad \forall j = 1, \dots, 4 \right\}$$

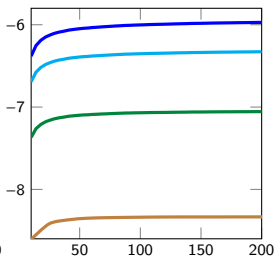
²Delage & Ye, Operations Research, 2010

Data-driven newsvendor problem

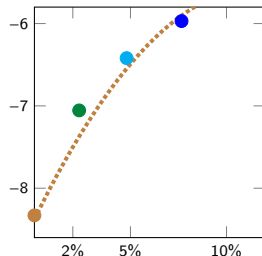
Disappointment probability
 $\mathbb{P}_{\theta^*} [c(\bar{x}_T, \theta^*) > \widehat{c}_T(\bar{x}_T)]$



Conservatism
 $\mathbb{E}_{\theta^*} [\widehat{c}_T(\bar{x}_T)]$



Conservatism vs
disappointment probability



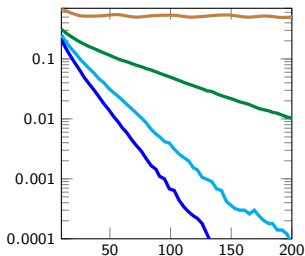
Model 2: DRO model with moment ambiguity set²

$$\widehat{c}_T(x) = \sup_{\theta \in \Theta} \left\{ c(x, \theta) : |\mathbb{E}_{\widehat{\theta}_T}[\xi^j] - \mathbb{E}_{\theta}[\xi^j]| \leq r \quad \forall j = 1, \dots, 4 \right\}$$

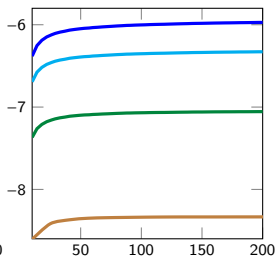
²Delage & Ye, Operations Research, 2010

Data-driven newsvendor problem

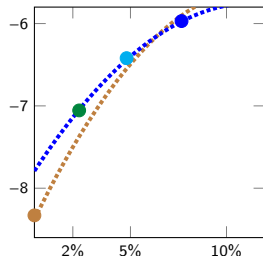
Disappointment probability
 $\mathbb{P}_{\theta^*} [c(\bar{x}_T, \theta^*) > \widehat{c}_T(\bar{x}_T)]$



Conservatism
 $\mathbb{E}_{\theta^*} [\widehat{c}_T(\bar{x}_T)]$



Conservatism vs
 disappointment probability



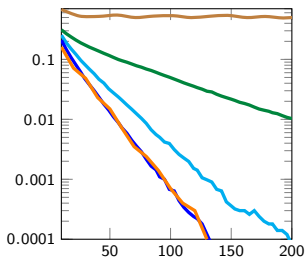
Model 2: DRO model with moment ambiguity set²

$$\widehat{c}_T(x) = \sup_{\theta \in \Theta} \left\{ c(x, \theta) : |\mathbb{E}_{\widehat{\theta}_T}[\xi^j] - \mathbb{E}_{\theta}[\xi^j]| \leq r \quad \forall j = 1, \dots, 4 \right\}$$

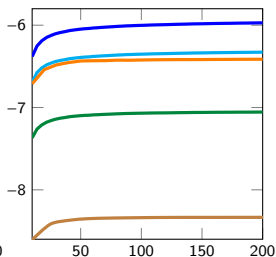
²Delage & Ye, Operations Research, 2010

Data-driven newsvendor problem

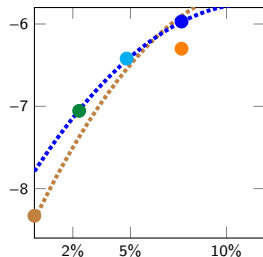
Disappointment probability
 $\mathbb{P}_{\theta^*} [c(\bar{x}_T, \theta^*) > \widehat{c}_T(\bar{x}_T)]$



Conservatism
 $\mathbb{E}_{\theta^*} [\widehat{c}_T(\bar{x}_T)]$



Conservatism vs
disappointment probability



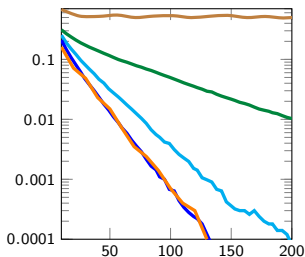
Model 3: DRO model with Wasserstein ambiguity set²

$$\widehat{c}_T(x) = \sup_{\theta \in \Theta} \{c(x, \theta) : d_W(\widehat{\theta}_T, \theta) \leq r\}$$

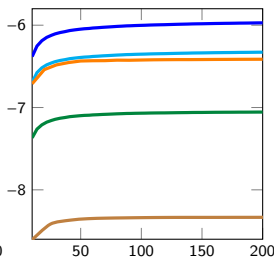
²Mohajerin Esfahani & Kuhn, Mathematical Programming, 2018

Data-driven newsvendor problem

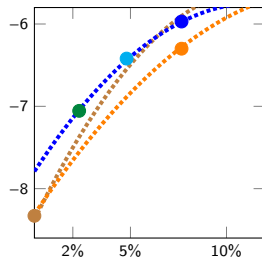
Disappointment probability
 $\mathbb{P}_{\theta^*} [c(\hat{x}_T, \theta^*) > \hat{c}_T(\hat{x}_T)]$



Conservatism
 $\mathbb{E}_{\theta^*} [\hat{c}_T(\hat{x}_T)]$



Conservatism vs
disappointment probability



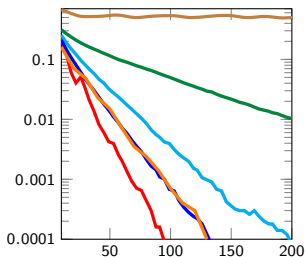
Model 3: DRO model with Wasserstein ambiguity set²

$$\hat{c}_T(x) = \sup_{\theta \in \Theta} \{c(x, \theta) : d_W(\hat{\theta}_T, \theta) \leq r\}$$

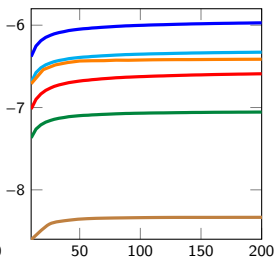
²Mohajerin Esfahani & Kuhn, Mathematical Programming, 2018

Data-driven newsvendor problem

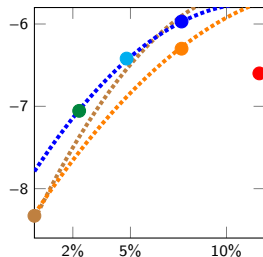
Disappointment probability
 $\mathbb{P}_{\theta^*} [c(\bar{x}_T, \theta^*) > \hat{c}_T(\bar{x}_T)]$



Conservatism
 $\mathbb{E}_{\theta^*} [\hat{c}_T(\bar{x}_T)]$



Conservatism vs
disappointment probability



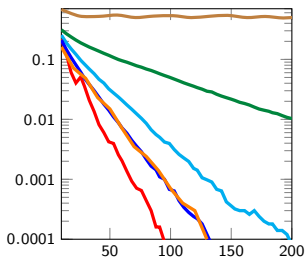
Model 4: DRO model with KL-ambiguity set²

$$\hat{c}_T(x) = \sup_{\theta \in \Theta} \{c(x, \theta) : D(\hat{\theta}_T \| \theta) \leq r\}$$

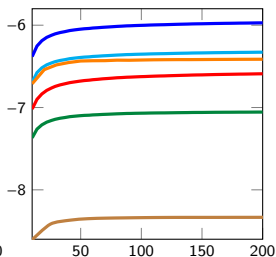
²Ben-Tal et al., Management Science, 2013

Data-driven newsvendor problem

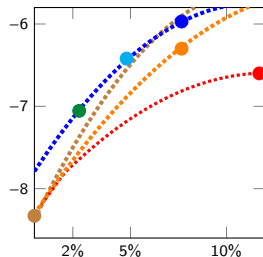
Disappointment probability
 $\mathbb{P}_{\theta^*} [c(\bar{x}_T, \theta^*) > \hat{c}_T(\bar{x}_T)]$



Conservatism
 $\mathbb{E}_{\theta^*} [\hat{c}_T(\bar{x}_T)]$



Conservatism vs
disappointment probability



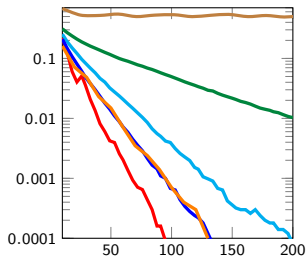
Model 4: DRO model with KL-ambiguity set²

$$\hat{c}_T(x) = \sup_{\theta \in \Theta} \{c(x, \theta) : D(\hat{\theta}_T \| \theta) \leq r\}$$

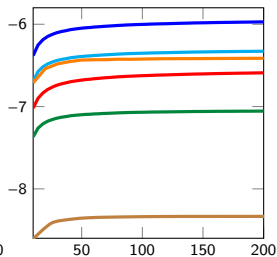
²Ben-Tal et al., Management Science, 2013

Data-driven newsvendor problem

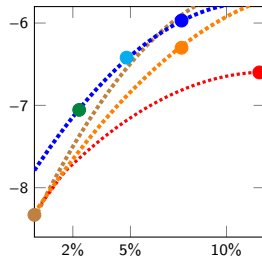
Disappointment probability
 $\mathbb{P}_{\theta^*} [c(\hat{x}_T, \theta^*) > \hat{c}_T(\hat{x}_T)]$



Conservatism
 $\mathbb{E}_{\theta^*} [\hat{c}_T(\hat{x}_T)]$

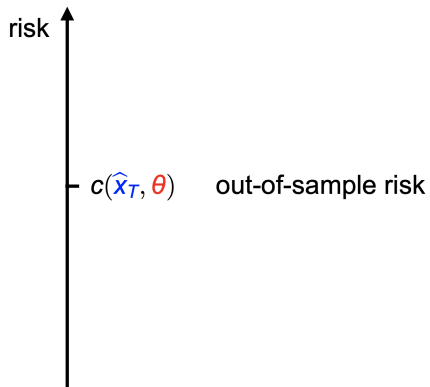


Conservatism vs
disappointment probability

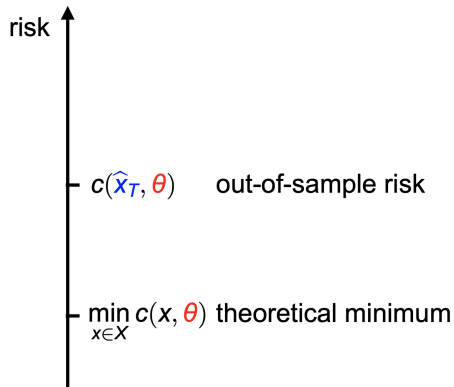


Which method is optimal?

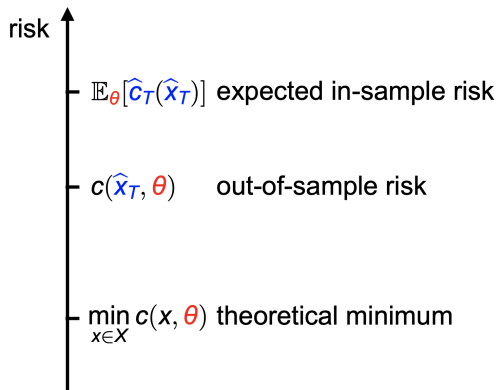
Optimal data-driven decision making



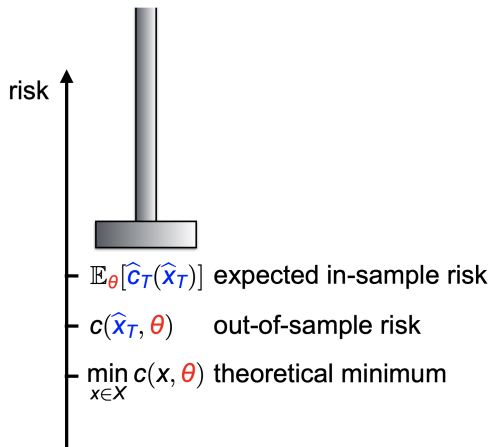
Optimal data-driven decision making



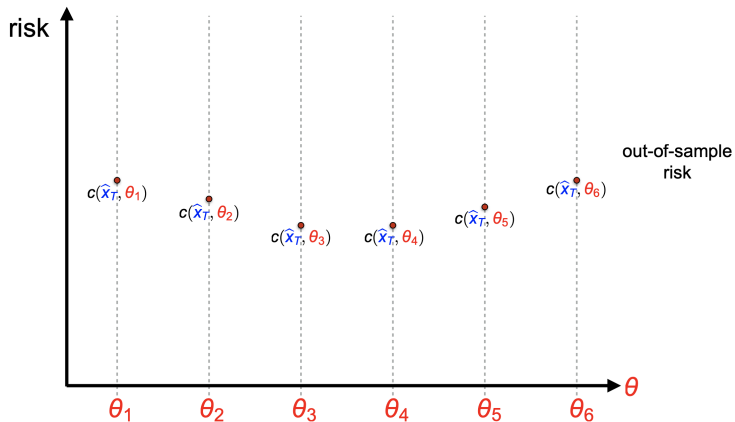
Optimal data-driven decision making



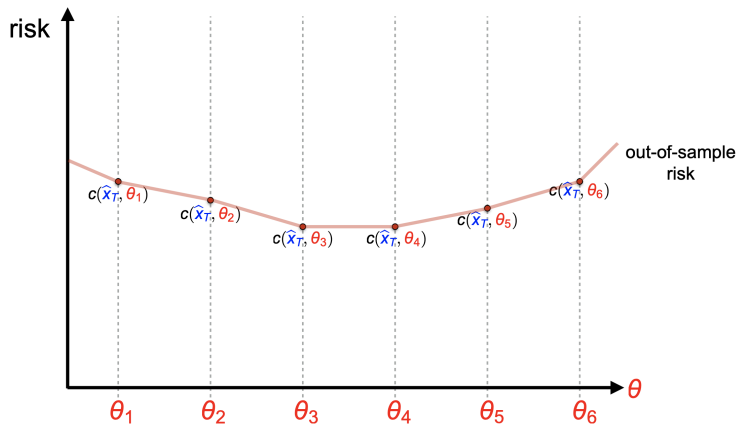
Optimal data-driven decision making



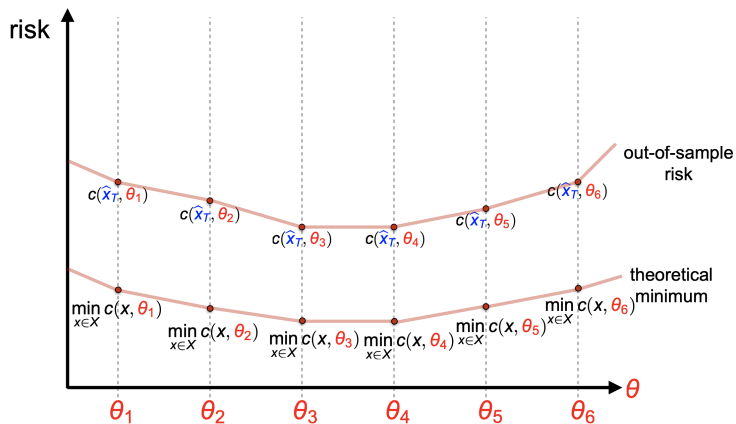
Optimal data-driven decision making



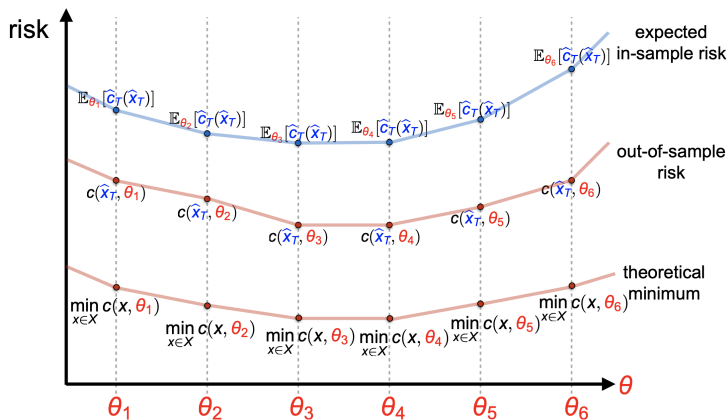
Optimal data-driven decision making



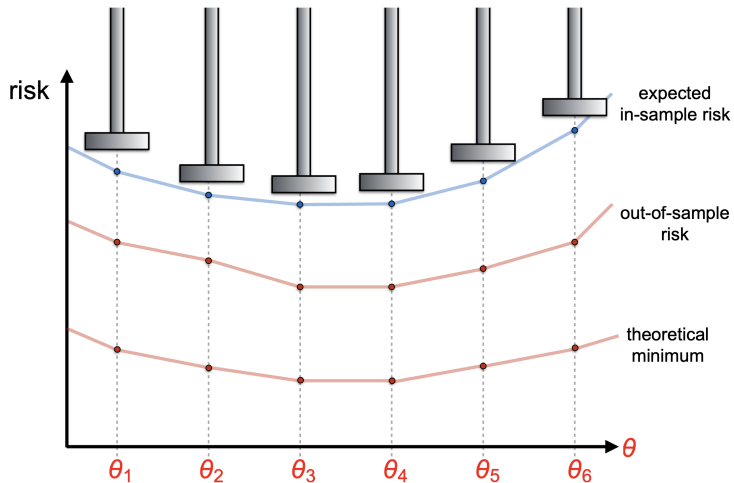
Optimal data-driven decision making



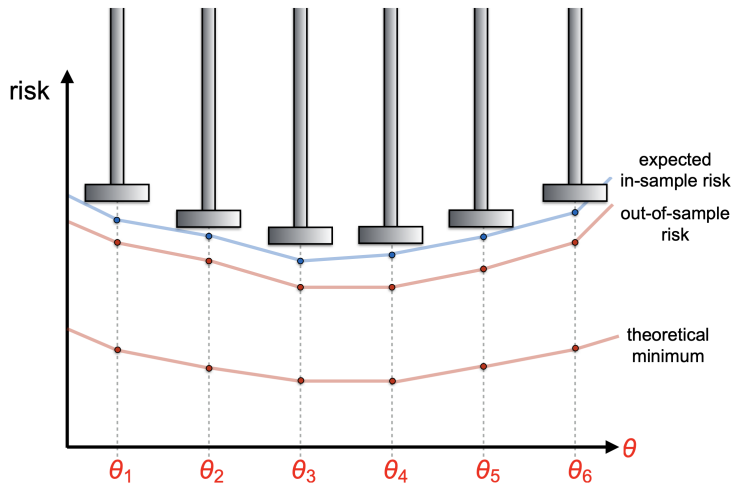
Optimal data-driven decision making



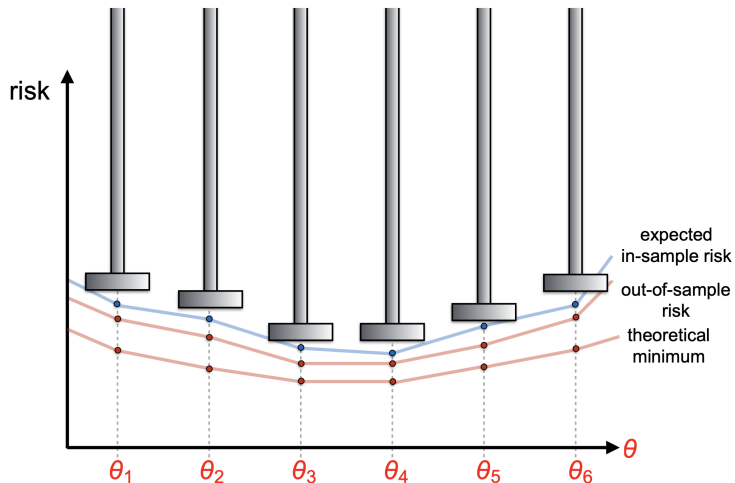
Optimal data-driven decision making



Optimal data-driven decision making



Optimal data-driven decision making



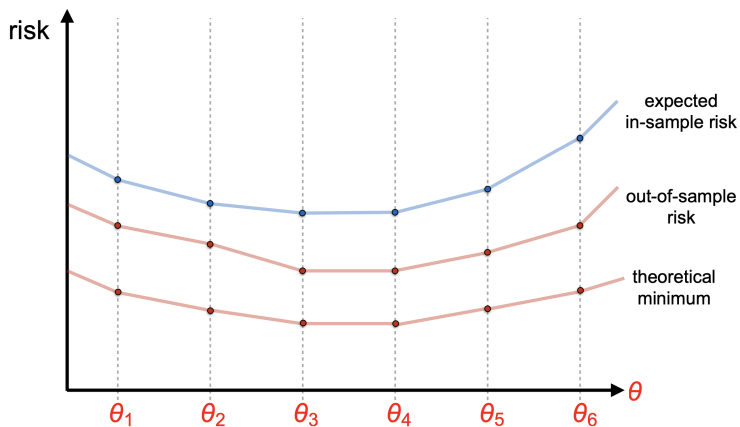
Optimal data-driven decision making (cont'd)

$$(\star) \left\{ \begin{array}{l} \min_{\hat{c}, \hat{x}} \left\{ \lim_{T \rightarrow \infty} \mathbb{E}_{\theta} [\hat{c}_T(\hat{x}_T)] \right\}_{\theta \in \Theta} \\ \text{s. t.} \quad \lim_{T \rightarrow \infty} \frac{1}{T} \log \mathbb{P}_{\theta} (c(\hat{x}_T, \theta) > \hat{c}_T(\hat{x}_T)) \leq -r \quad \forall \theta \in \Theta \end{array} \right.$$

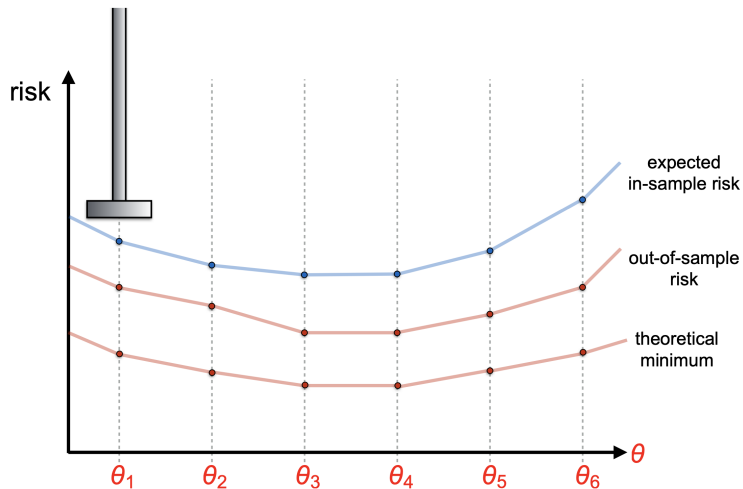
Interpretation: Among all predictors and prescriptors with **“small” disappointment** find the **least conservative** one

$$\mathbb{P}_{\theta} (c(\hat{x}_T, \theta) > \hat{c}_T(\hat{x}_T)) \leq e^{-rT+o(T)}$$

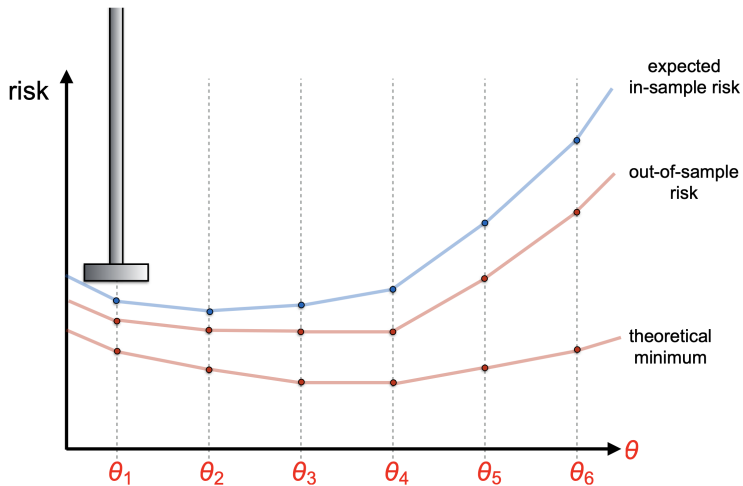
Optimal data-driven decision making (cont'd)



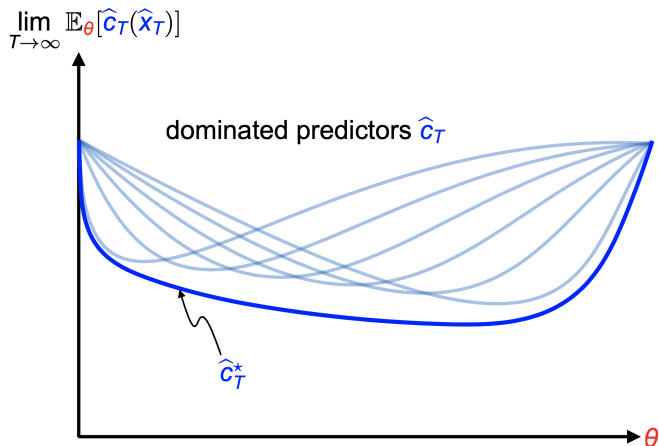
Optimal data-driven decision making (cont'd)



Optimal data-driven decision making (cont'd)



Pareto dominant solutions



\hat{C}_T^* minimizes the in-sample risk simultaneously for every θ

Optimizing over ALL surrogate models

$$(\star) \left\{ \begin{array}{l} \min_{\hat{c}, \hat{x}} \left\{ \lim_{T \rightarrow \infty} \mathbb{E}_{\theta} [\hat{c}_T(\hat{x}_T)] \right\}_{\theta \in \Theta} \\ \text{s. t.} \quad \lim_{T \rightarrow \infty} \frac{1}{T} \log \mathbb{P}_{\theta} (c(\hat{x}_T, \theta) > \hat{c}_T(\hat{x}_T)) \leq -r \quad \forall \theta \in \Theta \end{array} \right.$$

Optimizing over ALL surrogate models

$$(\star) \left\{ \begin{array}{l} \min_{\hat{c}, \hat{x}} \left\{ \lim_{T \rightarrow \infty} \mathbb{E}_{\theta} [\hat{c}_T(\hat{x}_T)] \right\}_{\theta \in \Theta} \\ \text{s. t.} \quad \lim_{T \rightarrow \infty} \frac{1}{T} \log \mathbb{P}_{\theta} (c(\hat{x}_T, \theta) > \hat{c}_T(\hat{x}_T)) \leq -r \quad \forall \theta \in \Theta \end{array} \right.$$

Strengths

- ▶ proxy for optimizing the out-of-sample risk
- ▶ admits a Pareto dominant solution in closed form
- ▶ errs on the side of caution
- ▶ facilitates separation of estimation and optimization

Optimizing over ALL surrogate models

$$(\star) \left\{ \begin{array}{l} \min_{\hat{c}, \hat{x}} \left\{ \lim_{T \rightarrow \infty} \mathbb{E}_{\theta} [\hat{c}_T(\hat{x}_T)] \right\}_{\theta \in \Theta} \\ \text{s. t.} \quad \lim_{T \rightarrow \infty} \frac{1}{T} \log \mathbb{P}_{\theta} (c(\hat{x}_T, \theta) > \hat{c}_T(\hat{x}_T)) \leq -r \quad \forall \theta \in \Theta \end{array} \right.$$

Weaknesses

- ▶ performance criteria are asymptotic
- ▶ choice of r is subjective
- ▶ feasible/optimal models are biased

Optimizing over ALL surrogate models

$$(\star) \left\{ \begin{array}{l} \min_{\hat{c}, \hat{x}} \left\{ \lim_{T \rightarrow \infty} \mathbb{E}_{\theta} [\hat{c}_T(\hat{x}_T)] \right\}_{\theta \in \Theta} \\ \text{s. t.} \quad \lim_{T \rightarrow \infty} \frac{1}{T} \log \mathbb{P}_{\theta} (c(\hat{x}_T, \theta) > \hat{c}_T(\hat{x}_T)) \leq -r \quad \forall \theta \in \Theta \end{array} \right.$$

Space of all possible predictors and prescriptors is large

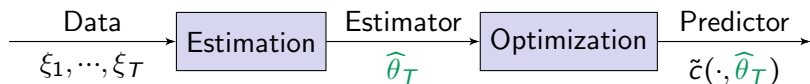
- ▶ \hat{c}_T, \hat{x}_T can be any² function depending on the available training data ξ_1, \dots, ξ_T
- ▶ Can we restrict ourselves to smaller class of functions without losing optimality?

²Some technical details, see arXiv:2010:06606

Optimizing over ALL surrogate models

$$(\star) \begin{cases} \min_{\hat{c}, \hat{x}} & \left\{ \lim_{T \rightarrow \infty} \mathbb{E}_{\theta} [\hat{c}_T(\hat{x}_T)] \right\}_{\theta \in \Theta} \\ \text{s. t.} & \lim_{T \rightarrow \infty} \frac{1}{T} \log \mathbb{P}_{\theta} (c(\hat{x}_T, \theta) > \hat{c}_T(\hat{x}_T)) \leq -r \quad \forall \theta \in \Theta \end{cases}$$

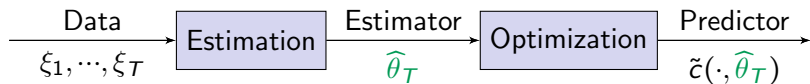
Key idea: Separation of estimation and optimization



- ① Which estimator should one pick?
- ② Is this separation without loss of optimality? Can we represent the (strong) solution to (\star) as

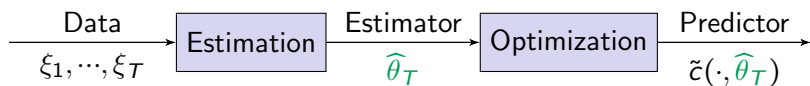
$$\hat{c}_T(x) = \tilde{c}(x, \hat{\theta}_T)$$

Separation principle - intuition



When can this separation be without loss of optimality?

Separation principle - intuition

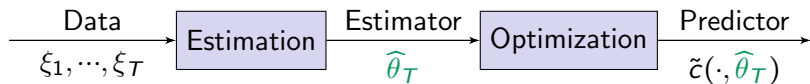


When can this separation be without loss of optimality?

- (i) **No statistical information** about θ **is lost** when considering an estimator, i.e.,

$$\theta \longrightarrow \hat{\theta}_T \longrightarrow \xi_1, \dots, \xi_T \quad \text{forms a Markov chain}$$

Separation principle - intuition



When can this separation be without loss of optimality?

- (i) **No statistical information** about θ **is lost** when considering an estimator, i.e.,

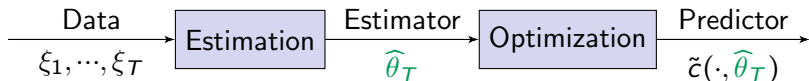
$$\theta \longrightarrow \hat{\theta}_T \longrightarrow \xi_1, \dots, \xi_T \quad \text{forms a Markov chain}$$

- (ii) **Estimator concentrates fast enough** around true model θ
 $\Rightarrow \hat{\theta}_T$ satisfies a large deviation principle

Restricted optimization problem

Original problem (★)

$$\begin{cases} \min_{\hat{c}, \hat{x}} & \left\{ \lim_{T \rightarrow \infty} \mathbb{E}_{\theta} [\hat{c}_T(\hat{x}_T)] \right\}_{\theta \in \Theta} \\ \text{s. t.} & \lim_{T \rightarrow \infty} \frac{1}{T} \log \mathbb{P}_{\theta} (c(\hat{x}_T, \theta) > \hat{c}_T(\hat{x}_T)) \leq -r \quad \forall \theta \end{cases}$$



Restricted problem (★★★)

$$\begin{cases} \min_{\tilde{c}, \tilde{x}} & \left\{ \tilde{c}(\tilde{x}(\theta), \theta) \right\}_{\theta \in \Theta} \\ \text{s. t.} & \lim_{T \rightarrow \infty} \frac{1}{T} \log \mathbb{P}_{\theta} (c(\tilde{x}(\hat{\theta}_T), \theta) > \tilde{c}(\tilde{x}(\hat{\theta}_T), \hat{\theta}_T)) \leq -r \quad \forall \theta \end{cases}$$

Large Deviations Theory

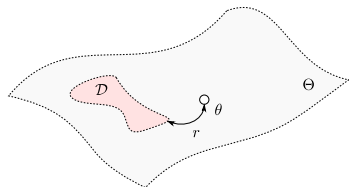
Definition: A sequence $\{\widehat{\theta}_T\}_{T \in \mathbb{N}}$ satisfies a **Large Deviation Principle (LDP)** if there is a “distance” function $I(\theta', \theta)$ such that for any Borel set $\mathcal{D} \subset \Theta'$

$$\begin{aligned} - \underbrace{\inf_{\theta' \in \text{int} \mathcal{D}} I(\theta', \theta)}_r &\leq \liminf_{T \rightarrow \infty} \frac{1}{T} \log \mathbb{P}_\theta (\widehat{\theta}_T \in \mathcal{D}) \\ &\leq \limsup_{T \rightarrow \infty} \frac{1}{T} \log \mathbb{P}_\theta (\widehat{\theta}_T \in \mathcal{D}) \leq - \underbrace{\inf_{\theta' \in \text{cl} \mathcal{D}} I(\theta', \theta)}_r \end{aligned}$$

Large Deviations Theory

Definition: A sequence $\{\widehat{\theta}_T\}_{T \in \mathbb{N}}$ satisfies a **Large Deviation Principle (LDP)** if there is a “distance” function $I(\theta', \theta)$ such that for any Borel set $\mathcal{D} \subset \Theta'$

$$\underbrace{- \inf_{\theta' \in \text{int} \mathcal{D}} I(\theta', \theta)}_r \leq \liminf_{T \rightarrow \infty} \frac{1}{T} \log \mathbb{P}_\theta (\widehat{\theta}_T \in \mathcal{D})$$
$$\leq \limsup_{T \rightarrow \infty} \frac{1}{T} \log \mathbb{P}_\theta (\widehat{\theta}_T \in \mathcal{D}) \leq - \underbrace{\inf_{\theta' \in \text{cl} \mathcal{D}} I(\theta', \theta)}_r$$



$$\mathbb{P}_\theta (\widehat{\theta}_T \in \mathcal{D}) = e^{-r \cdot T + o(T)}$$

Varadhan: *Rare events do occur every day. Someone always wins a lottery!*

Large Deviations Theory (cont'd)

Definition: $I : \Theta' \times \text{cl}\Theta \rightarrow [0, \infty]$ is called a **regular rate function** if it is

(i) **Radially monotonic** in θ , i.e.,

$$\{\theta \in \text{cl}\Theta : I(\theta', \theta) \leq r\} \subseteq \text{cl}\{\theta \in \Theta : I(\theta', \theta) < r\}$$

(ii) **Continuous** on $\Theta' \times \Theta$

(iii) **Level-compact**, i.e.,

$$\{(\theta, \theta') \in \text{cl}\Theta \times \text{cl}\Theta' : I(\theta', \theta) \leq r\} \text{ is compact } \forall r > 0$$

Large Deviations Theory (cont'd)

Definition: $I : \Theta' \times \text{cl}\Theta \rightarrow [0, \infty]$ is called a **regular rate function** if it is

(i) **Radially monotonic** in θ , i.e.,

$$\{\theta \in \text{cl}\Theta : I(\theta', \theta) \leq r\} \subseteq \text{cl}\{\theta \in \Theta : I(\theta', \theta) < r\}$$

(ii) **Continuous** on $\Theta' \times \Theta$

(iii) **Level-compact**, i.e.,

$$\{(\theta, \theta') \in \text{cl}\Theta \times \text{cl}\Theta' : I(\theta', \theta) \leq r\} \text{ is compact } \forall r > 0$$

Examples:

- ▶ Relative entropy $\Theta = \Theta' = \Delta_d$, $I(\theta', \theta) = D(\theta' || \theta)$
- ▶ Ellipsoid $\Theta = \Theta = \mathbb{R}$, $I(\theta', \theta) = (\theta - \theta')^\top \Sigma^{-1} (\theta - \theta')$
- ▶ many more ...

Pareto-dominant solution to restricted optimization problem

Restricted problem (★★)

$$\begin{cases} \min_{\tilde{c}, \tilde{x}} \{ \tilde{c}(\tilde{x}(\theta), \theta) \}_{\theta \in \Theta} \\ \text{s. t.} \quad \lim_{T \rightarrow \infty} \frac{1}{T} \log \mathbb{P}_{\theta} (c(\tilde{x}(\hat{\theta}_T), \theta) > \tilde{c}(\tilde{x}(\hat{\theta}_T), \hat{\theta}_T)) \leq -r \quad \forall \theta \end{cases}$$

Assumption: $\hat{\theta}_T$ satisfies an LDP with regular rate function I

Theorem. The Pareto-dominant solution to (★★) is given by the distributionally robust predictor

$$\tilde{c}(x, \hat{\theta}_T) = \begin{cases} \max_{\theta \in \Theta} c(x, \theta) \\ \text{s. t.} \quad I(\hat{\theta}_T, \theta) \leq r \end{cases}$$

DRO is optimal

Assumption: $\widehat{\theta}_T$ satisfies an LDP with regular rate function I

Theorem. The Pareto-dominant solution to (★★) is given by the distributionally robust predictor

$$\tilde{c}(x, \widehat{\theta}_T) = \begin{cases} \max_{\theta \in \Theta} & c(x, \theta) \\ \text{s. t.} & I(\widehat{\theta}_T, \theta) \leq r \end{cases}$$

- ▶ Shape of the ambiguity set determined by $\widehat{\theta}_T$
- ▶ Radius of ambiguity set determines the desired decay rate

Separation principle

Assumptions:

- ▶ $\widehat{\theta}_T$ satisfies an LDP with regular rate function I
- ▶ $\widehat{\theta}_T$ is a sufficient statistic for θ

Theorem. The Pareto-dominant solution to (★) is given by the distributionally robust predictor

$$\widehat{c}_T(x) = \begin{cases} \max_{\theta \in \Theta} c(x, \theta) \\ \text{s. t. } I(\widehat{\theta}_T, \theta) \leq r \end{cases}$$

- ▶ Problems (★) and (★★) have the same optimal solution
⇒ Separation of estimation and optimization is without loss of optimality
- ▶ Sufficiency restricts to exponential fam. of distributions for \mathbb{P}_θ
- ▶ Non-convex Rao-Blackwell type result

Conclusions of the Separation Theorem

① **DRO predictors are optimal** in a wide sense

$$\widehat{c}_T(x) = \begin{cases} \max_{\theta \in \Theta} c(x, \theta) \\ \text{s. t. } I(\widehat{\theta}_T, \theta) \leq r \end{cases}$$

Conclusions of the Separation Theorem

- ① **DRO predictors are optimal** in a wide sense

$$\widehat{c}_T(x) = \begin{cases} \max_{\theta \in \Theta} c(x, \theta) \\ \text{s. t. } I(\widehat{\theta}_T, \theta) \leq r \end{cases}$$

- ② **Ambiguity set is induced** by the choice of **estimator**

- ▶ $I(\cdot, \theta)$ is the rate function related to the estimator $\widehat{\theta}_T$
- ▶ size of the ambiguity set r quantifies the decay rate of the disappointment probability

Conclusions of the Separation Theorem

- ① **DRO predictors are optimal** in a wide sense

$$\widehat{c}_T(x) = \begin{cases} \max_{\theta \in \Theta} c(x, \theta) \\ \text{s. t. } I(\widehat{\theta}_T, \theta) \leq r \end{cases}$$

- ② **Ambiguity set is induced** by the choice of **estimator**

- ▶ $I(\cdot, \theta)$ is the rate function related to the estimator $\widehat{\theta}_T$
- ▶ size of the ambiguity set r quantifies the decay rate of the disappointment probability

- ③ **Invariance principle**

- ▶ $\psi : \Theta' \rightarrow \Theta'$ homeomorphism
- ▶ $\psi(\widehat{\theta}_T)$ satisfies LDP with rate function $I^\psi(\theta', \theta) = I(\psi^{-1}(\theta'), \theta)$
- ▶ DRO predictor is invariant

$$\begin{cases} \max_{\theta \in \Theta} c(x, \theta) \\ \text{s. t. } I(\widehat{\theta}_T, \theta) \leq r \end{cases} = \begin{cases} \max_{\theta \in \Theta} c(x, \theta) \\ \text{s. t. } I^\psi(\psi(\widehat{\theta}_T), \theta) \leq r \end{cases}$$

Revisit newsvendor problem

- ▶ # copies stocked: $x \in \mathbb{X} = \{0, 1, \dots, d\}$
- ▶ random daily demand: $\xi \in \Xi = \{0, 1, \dots, d\}$
- ▶ model: $\mathbb{P}_{\theta}[\xi \in i] = \theta_i$
- ▶ cost: $c(x, \theta) = \sum_{i=0}^d \theta_i (-p \min\{x, i\}) + kx$
- ▶ estimator: $(\widehat{\theta}_T)_i = \frac{1}{T} \sum_{t=1}^T 1_{\xi_t=i}, \quad i = 0, \dots, d$

Revisit newsvendor problem

- ▶ # copies stocked: $x \in \mathbb{X} = \{0, 1, \dots, d\}$
- ▶ random daily demand: $\xi \in \Xi = \{0, 1, \dots, d\}$
- ▶ model: $\mathbb{P}_\theta[\xi \in i] = \theta_i$
- ▶ cost: $c(x, \theta) = \sum_{i=0}^d \theta_i (-p \min\{x, i\}) + kx$
- ▶ estimator: $(\widehat{\theta}_T)_i = \frac{1}{T} \sum_{t=1}^T 1_{\xi_t=i}, \quad i = 0, \dots, d$

Sanov's Theorem. The estimator $\widehat{\theta}_T$ satisfies a large deviation principle with regular rate function $I(\widehat{\theta}_T, \theta) = D(\widehat{\theta}_T \parallel \theta)$

Revisit newsvendor problem

- ▶ # copies stocked: $x \in \mathbb{X} = \{0, 1, \dots, d\}$
- ▶ random daily demand: $\xi \in \Xi = \{0, 1, \dots, d\}$
- ▶ model: $\mathbb{P}_\theta[\xi \in i] = \theta_i$
- ▶ cost: $c(x, \theta) = \sum_{i=0}^d \theta_i (-p \min\{x, i\}) + kx$
- ▶ estimator: $(\widehat{\theta}_T)_i = \frac{1}{T} \sum_{t=1}^T 1_{\xi_t=i}, \quad i = 0, \dots, d$

Sanov's Theorem. The estimator $\widehat{\theta}_T$ satisfies a large deviation principle with regular rate function $I(\widehat{\theta}_T, \theta) = D(\widehat{\theta}_T \parallel \theta)$

- ▶ $\widehat{\theta}_T$ is a sufficient statistic (Fisher-Neyman)
- ▶ DRO predictor with **relative entropy** ambiguity set is optimal

$$\widehat{c}_T(x) = \tilde{c}(x, \widehat{\theta}_T) = \begin{cases} \max_{\theta \in \Theta} & c(x, \theta) \\ \text{s. t.} & D(\widehat{\theta}_T \parallel \theta) \leq r \end{cases}$$

Correlated setting: finite state Markov chain

- ▶ $\{\xi_t\}_{t \in \mathbb{N}}$ stationary Markov chain, initial state σ

Correlated setting: finite state Markov chain

- ▶ $\{\xi_t\}_{t \in \mathbb{N}}$ stationary Markov chain, initial state σ
- ▶ pairwise description $\theta_{ij} = \mathbb{P}[\xi_t = i, \xi_{t+1} = j]$

Correlated setting: finite state Markov chain

- ▶ $\{\xi_t\}_{t \in \mathbb{N}}$ stationary Markov chain, initial state σ
- ▶ pairwise description $\theta_{ij} = \mathbb{P}[\xi_t = i, \xi_{t+1} = j]$
- ▶ Models considered (irreducible)

$$\Theta = \left\{ \theta \in \mathbb{R}_{++}^{d \times d} : \sum_{i, j \in \Xi} \theta_{ij} = 1, \sum_{j \in \Xi} \theta_{ij} = \sum_{j \in \Xi} \theta_{ji} \quad \forall i \in \Xi \right\}$$

Correlated setting: finite state Markov chain

- ▶ $\{\xi_t\}_{t \in \mathbb{N}}$ stationary Markov chain, initial state σ
- ▶ pairwise description $\theta_{ij} = \mathbb{P}[\xi_t = i, \xi_{t+1} = j]$
- ▶ Models considered (irreducible)

$$\Theta = \left\{ \theta \in \mathbb{R}_{++}^{d \times d} : \sum_{i,j \in \Xi} \theta_{ij} = 1, \sum_{j \in \Xi} \theta_{ij} = \sum_{j \in \Xi} \theta_{ji} \quad \forall i \in \Xi \right\}$$

- ▶ Estimator: $(\widehat{\theta}_T)_{ij} = \frac{1}{T} (1_{\sigma=i} 1_{\xi_1=j} + \sum_{t=1}^{T-1} 1_{\xi_t=i} 1_{\xi_{t+1}=j})$
sufficient statistic for θ (Fisher-Neyman)

Correlated setting: finite state Markov chain

- ▶ $\{\xi_t\}_{t \in \mathbb{N}}$ stationary Markov chain, initial state σ
- ▶ pairwise description $\theta_{ij} = \mathbb{P}[\xi_t = i, \xi_{t+1} = j]$
- ▶ Models considered (irreducible)

$$\Theta = \left\{ \theta \in \mathbb{R}_{++}^{d \times d} : \sum_{i,j \in \Xi} \theta_{ij} = 1, \sum_{j \in \Xi} \theta_{ij} = \sum_{j \in \Xi} \theta_{ji} \quad \forall i \in \Xi \right\}$$

- ▶ Estimator: $(\widehat{\theta}_T)_{ij} = \frac{1}{T} (1_{\sigma=i} 1_{\xi_1=j} + \sum_{t=1}^{T-1} 1_{\xi_t=i} 1_{\xi_{t+1}=j})$
sufficient statistic for θ (Fisher-Neyman)
- ▶ Estimator state space

$$\Theta' = \left\{ \theta \in \mathbb{R}_+^{d \times d} : \sum_{i,j \in \Xi} \theta_{ij} = 1 \right\}$$

Correlated setting: finite state Markov chain

- ▶ $\{\xi_t\}_{t \in \mathbb{N}}$ stationary Markov chain, initial state σ
- ▶ pairwise description $\theta_{ij} = \mathbb{P}[\xi_t = i, \xi_{t+1} = j]$
- ▶ Models considered (irreducible)

$$\Theta = \left\{ \theta \in \mathbb{R}_{++}^{d \times d} : \sum_{i,j \in \Xi} \theta_{ij} = 1, \sum_{j \in \Xi} \theta_{ij} = \sum_{j \in \Xi} \theta_{ji} \quad \forall i \in \Xi \right\}$$

- ▶ Estimator: $(\widehat{\theta}_T)_{ij} = \frac{1}{T} (1_{\sigma=i} 1_{\xi_1=j} + \sum_{t=1}^{T-1} 1_{\xi_t=i} 1_{\xi_{t+1}=j})$
sufficient statistic for θ (Fisher-Neyman)
- ▶ Estimator state space

$$\Theta' = \left\{ \theta \in \mathbb{R}_+^{d \times d} : \sum_{i,j \in \Xi} \theta_{ij} = 1 \right\}$$

- ▶ “Distance measure” between estimator and underlying model
→ conditional relative entropy

Correlated setting: finite state Markov chain (cont'd)

Conditional relative entropy. For any $\theta \in \Theta, \theta' \in \Theta'$

$$D_c(\theta' \parallel \theta) = \sum_{i,j \in \Xi} \theta'_{ij} \left(\log \left(\frac{\theta'_{ij}}{\sum_{k \in \Xi} \theta'_{ik}} \right) - \log \left(\frac{\theta_{ij}}{\sum_{k \in \Xi} \theta_{ik}} \right) \right)$$

Correlated setting: finite state Markov chain (cont'd)

Conditional relative entropy. For any $\theta \in \Theta, \theta' \in \Theta'$

$$D_c(\theta' \parallel \theta) = \sum_{i,j \in \Xi} \theta'_{ij} \left(\log \left(\frac{\theta'_{ij}}{\sum_{k \in \Xi} \theta'_{ik}} \right) - \log \left(\frac{\theta_{ij}}{\sum_{k \in \Xi} \theta_{ik}} \right) \right)$$

- ▶ Similar properties as the relative entropy
 - ▶ non-negative
 - ▶ $D_c(\theta' \parallel \theta) = 0 \Leftrightarrow \theta' = \theta$
 - ▶ non-convex in $\theta \Rightarrow$ solving DRO problem [Li, S., Kuhn, ICML 2021]

Correlated setting: finite state Markov chain (cont'd)

Conditional relative entropy. For any $\theta \in \Theta, \theta' \in \Theta'$

$$D_c(\theta' \parallel \theta) = \sum_{i,j \in \Xi} \theta'_{ij} \left(\log \left(\frac{\theta'_{ij}}{\sum_{k \in \Xi} \theta'_{ik}} \right) - \log \left(\frac{\theta_{ij}}{\sum_{k \in \Xi} \theta_{ik}} \right) \right)$$

- ▶ Similar properties as the relative entropy
 - ▶ non-negative
 - ▶ $D_c(\theta' \parallel \theta) = 0 \Leftrightarrow \theta' = \theta$
 - ▶ non-convex in $\theta \Rightarrow$ solving DRO problem [Li, S., Kuhn, ICML 2021]

Lemma. The estimator $\widehat{\theta}_T$ satisfies an LDP with (regular) rate function $I(\widehat{\theta}_T, \theta) = D_c(\widehat{\theta}_T \parallel \theta)$

Correlated setting: finite state Markov chain (cont'd)

Conditional relative entropy. For any $\theta \in \Theta, \theta' \in \Theta'$

$$D_c(\theta' \parallel \theta) = \sum_{i,j \in \Xi} \theta'_{ij} \left(\log \left(\frac{\theta'_{ij}}{\sum_{k \in \Xi} \theta'_{ik}} \right) - \log \left(\frac{\theta_{ij}}{\sum_{k \in \Xi} \theta_{ik}} \right) \right)$$

- ▶ Similar properties as the relative entropy
 - ▶ non-negative
 - ▶ $D_c(\theta' \parallel \theta) = 0 \Leftrightarrow \theta' = \theta$
 - ▶ non-convex in $\theta \Rightarrow$ solving DRO problem [Li, S., Kuhn, ICML 2021]

Lemma. The estimator $\widehat{\theta}_T$ satisfies an LDP with (regular) rate function $I(\widehat{\theta}_T, \theta) = D_c(\widehat{\theta}_T \parallel \theta)$

- ▶ [Dembo & Zeitouni, Chapter 3]
- ▶ Separation Theorem holds

IID process with unknown mean

- ▶ i.i.d. process $\{\xi_t\}_{t \in \mathbb{N}} \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}_\theta$, unknown mean $\mathbb{E}_{\mathbb{P}_\theta}[\xi_1] = \theta$

IID process with unknown mean

- ▶ i.i.d. process $\{\xi_t\}_{t \in \mathbb{N}} \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}_\theta$, unknown mean $\mathbb{E}_{\mathbb{P}_\theta}[\xi_1] = \theta$
- ▶ F_θ distribution of ξ_1 under \mathbb{P}_θ

IID process with unknown mean

- ▶ i.i.d. process $\{\xi_t\}_{t \in \mathbb{N}} \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}_\theta$, unknown mean $\mathbb{E}_{\mathbb{P}_\theta}[\xi_1] = \theta$
- ▶ F_θ distribution of ξ_1 under \mathbb{P}_θ
- ▶ log-moment generating function $\Lambda(\lambda, \theta) = \log \left(\int_{\Xi} e^{\lambda^\top \xi} dF_\theta(\xi) \right)$

IID process with unknown mean

- ▶ i.i.d. process $\{\xi_t\}_{t \in \mathbb{N}} \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}_\theta$, unknown mean $\mathbb{E}_{\mathbb{P}_\theta}[\xi_1] = \theta$
- ▶ F_θ distribution of ξ_1 under \mathbb{P}_θ
- ▶ log-moment generating function $\Lambda(\lambda, \theta) = \log \left(\int_{\Xi} e^{\lambda^\top \xi} dF_\theta(\xi) \right)$
- ▶ Estimator $\widehat{\theta}_T = \frac{1}{T} \sum_{t=1}^T \xi_t$
 - ▶ consistent by Law of Large Numbers
 - ▶ for many distributions is a sufficient statistic

IID process with unknown mean

- ▶ i.i.d. process $\{\xi_t\}_{t \in \mathbb{N}} \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}_\theta$, unknown mean $\mathbb{E}_{\mathbb{P}_\theta}[\xi_1] = \theta$
- ▶ F_θ distribution of ξ_1 under \mathbb{P}_θ
- ▶ log-moment generating function $\Lambda(\lambda, \theta) = \log \left(\int_{\Xi} e^{\lambda^\top \xi} dF_\theta(\xi) \right)$
- ▶ Estimator $\widehat{\theta}_T = \frac{1}{T} \sum_{t=1}^T \xi_t$
 - ▶ consistent by Law of Large Numbers
 - ▶ for many distributions is a sufficient statistic
- ▶ Cramér function $\Lambda^*(s, \theta) = \sup_{\lambda \in \mathbb{R}^d} \{ \langle \lambda, s \rangle - \Lambda(\lambda, \theta) \}$

IID process with unknown mean

- ▶ i.i.d. process $\{\xi_t\}_{t \in \mathbb{N}} \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}_\theta$, unknown mean $\mathbb{E}_{\mathbb{P}_\theta}[\xi_1] = \theta$
- ▶ F_θ distribution of ξ_1 under \mathbb{P}_θ
- ▶ log-moment generating function $\Lambda(\lambda, \theta) = \log \left(\int_{\Xi} e^{\lambda^\top \xi} dF_\theta(\xi) \right)$
- ▶ Estimator $\widehat{\theta}_T = \frac{1}{T} \sum_{t=1}^T \xi_t$
 - ▶ consistent by Law of Large Numbers
 - ▶ for many distributions is a sufficient statistic
- ▶ Cramér function $\Lambda^*(s, \theta) = \sup_{\lambda \in \mathbb{R}^d} \{ \langle \lambda, s \rangle - \Lambda(\lambda, \theta) \}$

Lemma. The estimator $\widehat{\theta}_T$ satisfies an LDP with (regular) rate function $I(\theta', \theta) = \Lambda^*(\theta', \theta)$

IID process with unknown mean

- ▶ i.i.d. process $\{\xi_t\}_{t \in \mathbb{N}} \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}_\theta$, unknown mean $\mathbb{E}_{\mathbb{P}_\theta}[\xi_1] = \theta$
- ▶ F_θ distribution of ξ_1 under \mathbb{P}_θ
- ▶ log-moment generating function $\Lambda(\lambda, \theta) = \log \left(\int_{\Xi} e^{\lambda^\top \xi} dF_\theta(\xi) \right)$
- ▶ Estimator $\widehat{\theta}_T = \frac{1}{T} \sum_{t=1}^T \xi_t$
 - ▶ consistent by Law of Large Numbers
 - ▶ for many distributions is a sufficient statistic
- ▶ Cramér function $\Lambda^*(s, \theta) = \sup_{\lambda \in \mathbb{R}^d} \{ \langle \lambda, s \rangle - \Lambda(\lambda, \theta) \}$

Lemma. The estimator $\widehat{\theta}_T$ satisfies an LDP with (regular) rate function $I(\theta', \theta) = \Lambda^*(\theta', \theta)$

- ▶ Consequence of Cramér's Theorem
- ▶ Separation Theorem holds for many distributions

IID process with unknown mean (cont'd)

Each distribution induces an ambiguity set

$$\{\theta \in \Theta : \Lambda^*(\theta', \theta) \leq r\}$$

F_θ	$\Lambda^*(\theta', \theta)$	$\text{dom}(\Lambda^*(\cdot, \theta))$
(a) Normal	$\frac{1}{2}(\theta' - \theta)^\top \Sigma^{-1}(\theta' - \theta)$	\mathbb{R}^d
(b) Exponential	$\frac{\theta' - \theta}{\theta} + \log(\theta/\theta')$	\mathbb{R}_{++}
(c) Poisson	$\theta' \log(\theta'/\theta) - \theta' + \theta$	\mathbb{R}_{++}
(d) Bernoulli	$\theta' \log\left(\frac{\theta'(1-\theta)}{\theta(1-\theta')}\right) - \log\left(\frac{1-\theta}{1-\theta'}\right)$	$(0, 1)$

Many more possible, e.g., Gamma, Geometric, Binomial, ...

Summary

Meta-optimization problem

- ▶ optimizes over surrogate optimization models
- ▶ balances in-sample risk vs. out-of-sample disappointment
- ▶ pushes down the out-of-sample risk

Separation of estimation and optimization

- ▶ holds if $\widehat{\theta}_T$ is a sufficient statistic that obeys an LDP
- ▶ reminiscent of Rao-Blackwell theorem

Pareto-dominant solution is a DRO model

- ▶ ambiguity set is a rate-ball around $\widehat{\theta}_T$
- ▶ radius = decay rate of the out-of-sample disappointment
- ▶ invariant under homeomorphic transformations

Conclusion

- ① **DRO is optimal:** Optimal data-driven predictor is given as a DRO problem centred around an estimator
- ② **Ambiguity set**
 - ▶ **Structure** induced by underlying stochastic process via LDP
 - ▶ **Size** has operational meaning as the decay rate of the disappointment probability
- ③ Data-driven DRO framework for **non-i.i.d. data**

Outlook

- ▶ The proposed prescriptor is **not consistent**

$$\widehat{c}_T(x) \not\rightarrow c(x, \theta) \quad \text{as } T \rightarrow \infty$$

- ▶ **Idea:** Can we trade speed in the decay of

$$\mathbb{P}_\theta (c(\widehat{x}_T, \theta) > \widehat{c}_T(\widehat{x}_T))$$

to achieve consistency?

- ▶ Are there other statistical criteria for optimality?

⇒ A. Ganguly and T. Sutter, *Optimal learning via Moderate Deviations Theory*, *arXiv:2305.14496*, 2023

Reference

This talk

- ▶ T. Sutter, B.P. Van Parys, and D. Kuhn, *A General Framework for Optimal Data-Driven Optimization*, arXiv:2010.06606, 2020
- ▶ A. Ganguly and T. Sutter, *Optimal learning via Moderate Deviations Theory* arXiv:2305.14496, 2023
- ▶ M. Li, T. Sutter, and D. Kuhn, *Distributionally Robust Optimization with Markovian Data*, ICML, 2021

Appendix

Interval estimation

Goal: Estimate cost $c(\theta)$ via a confidence interval $\widehat{\mathcal{I}}_T^* = \mathcal{I}_T^*(\widehat{\theta}_T)$ where $\mathcal{I}_T^*(\theta) = [\underline{c}_{T,r}^*(\theta), \bar{c}_{T,r}^*(\theta)]$ with the following properties

① **Exponential accuracy:**

$$\mathbb{P}_\theta(c(\theta) \notin \widehat{\mathcal{I}}_T^*) \leq e^{-rb_T}, \quad 1 \ll b_T \ll T$$

Interval estimation

Goal: Estimate cost $c(\theta)$ via a confidence interval $\widehat{\mathcal{I}}_T^* = \mathcal{I}_T^*(\widehat{\theta}_T)$ where $\mathcal{I}_T^*(\theta) = [\underline{c}_{T,r}^*(\theta), \bar{c}_{T,r}^*(\theta)]$ with the following properties

① **Exponential accuracy:**

$$\mathbb{P}_\theta(c(\theta) \notin \widehat{\mathcal{I}}_T^*) \leq e^{-rb_T}, \quad 1 \ll b_T \ll T$$

② **Minimality:** Any interval $\mathcal{I}_T(\theta) = [\underline{c}_{T,r}(\theta), \bar{c}_{T,r}(\theta)]$ satisfying ① is eventually larger than $\mathcal{I}_T^*(\theta)$

Interval estimation

Goal: Estimate cost $c(\theta)$ via a confidence interval $\widehat{\mathcal{I}}_T^* = \mathcal{I}_T^*(\widehat{\theta}_T)$ where $\mathcal{I}_T^*(\theta) = [\underline{c}_{T,r}^*(\theta), \bar{c}_{T,r}^*(\theta)]$ with the following properties

① **Exponential accuracy:**

$$\mathbb{P}_{\theta}(c(\theta) \notin \widehat{\mathcal{I}}_T^*) \leq e^{-rb_T}, \quad 1 \ll b_T \ll T$$

② **Minimality:** Any interval $\mathcal{I}_T(\theta) = [\underline{c}_{T,r}(\theta), \bar{c}_{T,r}(\theta)]$ satisfying ① is eventually larger than $\mathcal{I}_T^*(\theta)$

③ **Consistency:** $\widehat{\mathcal{I}}_T^* \rightarrow \{c(\theta)\}$ as $T \rightarrow \infty$

Interval estimation

Goal: Estimate cost $c(\theta)$ via a confidence interval $\widehat{\mathcal{I}}_T^* = \mathcal{I}_T^*(\widehat{\theta}_T)$ where $\mathcal{I}_T^*(\theta) = [\underline{c}_{T,r}^*(\theta), \bar{c}_{T,r}^*(\theta)]$ with the following properties

① **Exponential accuracy:**

$$\mathbb{P}_\theta(c(\theta) \notin \widehat{\mathcal{I}}_T^*) \leq e^{-rb_T}, \quad 1 \ll b_T \ll T$$

② **Minimality:** Any interval $\mathcal{I}_T(\theta) = [\underline{c}_{T,r}(\theta), \bar{c}_{T,r}(\theta)]$ satisfying ① is eventually larger than $\mathcal{I}_T^*(\theta)$

③ **Consistency:** $\widehat{\mathcal{I}}_T^* \rightarrow \{c(\theta)\}$ as $T \rightarrow \infty$

④ **Mischaracterization probability:**

$$\mathbb{P}_\theta(c(\theta') \notin \widehat{\mathcal{I}}_T^*) > e^{-rb_T}, \quad \forall \theta' : c(\theta') \neq c(\theta)$$

Interval estimation

Goal: Estimate cost $c(\theta)$ via a confidence interval $\widehat{\mathcal{I}}_T^* = \mathcal{I}_T^*(\widehat{\theta}_T)$ where $\mathcal{I}_T^*(\theta) = [\underline{c}_{T,r}^*(\theta), \bar{c}_{T,r}^*(\theta)]$ with the following properties

① **Exponential accuracy:**

$$\mathbb{P}_\theta(c(\theta) \notin \widehat{\mathcal{I}}_T^*) \leq e^{-rb_T}, \quad 1 \ll b_T \ll T$$

② **Minimality:** Any interval $\mathcal{I}_T(\theta) = [\underline{c}_{T,r}(\theta), \bar{c}_{T,r}(\theta)]$ satisfying ① is eventually larger than $\mathcal{I}_T^*(\theta)$

③ **Consistency:** $\widehat{\mathcal{I}}_T^* \rightarrow \{c(\theta)\}$ as $T \rightarrow \infty$

④ **Mischaracterization probability:**

$$\mathbb{P}_\theta(c(\theta') \notin \widehat{\mathcal{I}}_T^*) > e^{-rb_T}, \quad \forall \theta' : c(\theta') \neq c(\theta)$$

⑤ **Uniformly most accurate (UMA):** Any interval $\widehat{\mathcal{I}}_T$ satisfying ① is such that

$$\mathbb{P}_\theta(c(\theta') \in \widehat{\mathcal{I}}_T^*) \leq \mathbb{P}_\theta(c(\theta') \in \widehat{\mathcal{I}}_T), \quad \forall \theta' : c(\theta') \neq c(\theta)$$

Heuristic CLT based confidence intervals

- ▶ Given a **fixed** α , CLT guarantees

$$\lim_{T \rightarrow \infty} \mathbb{P}_{\theta}(c(\theta) \notin \mathcal{I}_{T,\alpha}^{\text{CLT}}(\widehat{\theta}_T)) \leq \alpha$$

- ▶ for the CLT-based interval

$$\mathcal{I}_{T,\alpha}^{\text{CLT}}(\widehat{\theta}_T) = \left[c(\widehat{\theta}_T) - \kappa_T^{\text{CLT}}(\alpha), c(\widehat{\theta}_T) + \kappa_T^{\text{CLT}}(\alpha) \right]$$

$$\kappa_T^{\text{CLT}}(\alpha) = \Phi^{-1}(1 - \alpha/2) \sqrt{\nabla c(\widehat{\theta}_T)^{\top} S(\widehat{\theta}_T) \nabla c(\widehat{\theta}_T)} / \sqrt{T}$$

- ▶ Heuristic choice $\alpha = e^{-rbT}$

Question: Does the CI $\mathcal{I}_{T,\alpha}^{\text{CLT}}(\widehat{\theta}_T)$ for $\alpha = e^{-rT}$ satisfy any of the properties ① – ⑤?

Optimal confidence interval

- ▶ Interval $\widehat{\mathcal{I}}_T^* = \mathcal{I}_T^*(\widehat{\theta}_T)$ for $\mathcal{I}_T^*(\theta) = [\underline{c}_{T,r}^*(\theta), \bar{c}_{T,r}^*(\theta)]$
 - ▶ $\underline{c}_{T,r}^*(\theta') = \inf_{\theta \in \Theta} \{c(\theta) : I^M(a_T(\theta' - \theta), \theta) \leq r\}$
 - ▶ $\bar{c}_{T,r}^*(\theta') = \sup_{\theta \in \Theta} \{c(\theta) : I^M(a_T(\theta' - \theta), \theta) \leq r\}$
- ▶ $I^M(\cdot, \theta)$: Moderate deviation rate function of $\widehat{\theta}_T$
- ▶ $a_T = \sqrt{T/b_T}$, $1 \ll b_T \ll T$

The confidence interval $\widehat{\mathcal{I}}_T^*$ satisfies the properties ① – ⑤

- ▶ Θ can be infinite dimensional
- ▶ mild assumptions on $I^M(\cdot, \theta)$

Example: Asymptotic variance of OU process

- ▶ Ornstein-Uhlenbeck process

$$dX_t = -\theta X_t dt + dW_t, \quad X_0 = 0$$

- ▶ Asymptotic variance

$$c(\theta) = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} (\mathbb{E}_\theta[X_t^2] - \mathbb{E}_\theta[X_t]^2) = \frac{1}{2\theta}$$

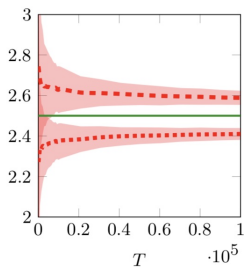
- ▶ Maximum likelihood estimator

$$\widehat{\theta}_T = -\frac{X_T^2 - T}{2 \int_0^T X_t^2 dt}, \quad I^M(\vartheta, \theta) = \frac{\vartheta^2}{2\theta}$$

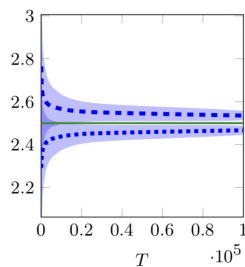
- ▶ Optimal CI $\mathcal{I}_T^*(\widehat{\theta}_T) = [c(\widehat{\theta}_T) + \kappa_T^-, c(\widehat{\theta}_T) + \kappa_T^+]$

$$\kappa_T^\pm = \frac{1}{2\widehat{\theta}_T^2} \left(r_T \pm \sqrt{r_T^2 + 2\widehat{\theta}_T r_T} \right), \quad r_T = r b_T / T$$

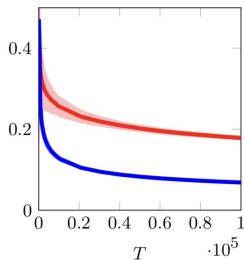
Example: Asymptotic variance of OU process



(a) Upper and lower confidence bound



(b) Upper and lower confidence bound



(c) Interval length

- ▶ Optimal value $c(\theta)$
- ▶ Optimal interval $\mathcal{I}_T^*(\hat{\theta}_T)$
- ▶ CLT interval $\mathcal{I}_{T,\alpha}^{\text{CLT}}(\hat{\theta}_T)$, $\alpha = e^{-rb_T}$