

# A General Framework For Optimal Data-Driven Optimization

Tobias Sutter,<sup>1)</sup> Bart Van Parys,<sup>2)</sup> Daniel Kuhn <sup>1)</sup>

<sup>1)</sup> Risk Analytics and Optimization Chair, EPFL  
[www.epfl.ch/labs/rao/](http://www.epfl.ch/labs/rao/)

<sup>2)</sup> MIT Sloan School of Management  
[web.mit.edu/vanparys/www/](http://web.mit.edu/vanparys/www/)

# Data-Driven Decision-Making

Stochastic optimization problem

$$\underset{x \in X}{\text{minimize}} \quad c(x, \theta)$$

# Data-Driven Decision-Making

Stochastic optimization problem

$$\underset{x \in X}{\text{minimize}} \quad c(x, \theta)$$

Family of probability measures

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

# Data-Driven Decision-Making

Stochastic optimization problem

$$\underset{x \in X}{\text{minimize}} \quad c(x, \theta)$$

Family of probability measures

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

Data-generating process

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



# Data-Driven Decision-Making

Stochastic optimization problem

$$\underset{x \in X}{\text{minimize}} \quad c(x, \theta)$$

Family of probability measures

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

Data-generating process

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

**Examples:**

# Data-Driven Decision-Making

Stochastic optimization problem

$$\underset{x \in X}{\text{minimize}} \quad c(x, \theta)$$

Family of probability measures

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

Data-generating process

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

## Examples:

- Expected loss

$$c(x, \theta) = \mathbb{E}_\theta[\ell(x, \xi)]$$

# Data-Driven Decision-Making

Stochastic optimization problem

$$\underset{x \in X}{\text{minimize}} \quad c(x, \theta)$$

Family of probability measures

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

Data-generating process

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

## Examples:

- Expected loss
- Risk of loss

$$c(x, \theta) = \mathbb{E}_\theta[\ell(x, \xi)]$$

$$c(x, \theta) = \rho_\theta[\ell(x, \xi)]$$

# Data-Driven Decision-Making

Stochastic optimization problem

$$\underset{x \in X}{\text{minimize}} \quad c(x, \theta)$$

Family of probability measures

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

Data-generating 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

$$c(x, \theta) = \mathbb{E}_\theta[\ell(x, \xi) | C\xi \in B]$$

# Data-Driven Decision-Making

Stochastic optimization problem

$$\underset{x \in X}{\text{minimize}} \quad c(x, \theta)$$

Family of probability measures

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

Data-generating 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

$$c(x, \theta) = \mathbb{E}_\theta[\ell(x, \xi) | \mathbf{C}\xi \in B]$$

► Long-run average loss

$$c(x, \theta) = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \mathbb{E}_\theta[\ell(\pi_x(\mathbf{s}_t), \mathbf{s}_t)]$$

# Data-Driven Decision-Making

Stochastic optimization problem

$$\underset{x \in X}{\text{minimize}} \quad c(x, \theta)$$

Family of probability measures

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

Data-generating process

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

**Assumptions:**

# Data-Driven Decision-Making

Stochastic optimization problem

$$\underset{x \in X}{\text{minimize}} \quad c(x, \theta)$$

Family of probability measures

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

Data-generating process

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

## Assumptions:

- All measures defined on  $(\Omega, \mathcal{F})$

# Data-Driven Decision-Making

Stochastic optimization problem

$$\underset{x \in X}{\text{minimize}} \quad c(x, \theta)$$

Family of probability measures

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

Data-generating process

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

## Assumptions:

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



# Data-Driven Decision-Making

Stochastic optimization problem

$$\underset{x \in X}{\text{minimize}} \quad c(x, \theta)$$

Family of probability measures

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

Data-generating process

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

**Examples:**

# Data-Driven Decision-Making

Stochastic optimization problem

$$\underset{x \in X}{\text{minimize}} \quad c(x, \theta)$$

Family of probability measures

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

Data-generating process

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

## Examples:

- Finite-state i.i.d. processes

# Data-Driven Decision-Making

Stochastic optimization problem

$$\underset{x \in X}{\text{minimize}} \quad c(x, \theta)$$

Family of probability measures

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

Data-generating process

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

## Examples:

- ▶ Finite-state i.i.d. processes
- ▶ Finite-state Markov chains

# Data-Driven Decision-Making

Stochastic optimization problem

$$\underset{x \in X}{\text{minimize}} \quad c(x, \theta)$$

Family of probability measures

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

Data-generating process

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

## Examples:

- ▶ Finite-state i.i.d. processes
- ▶ Finite-state Markov chains
- ▶ Vector-autoregressive processes

# Data-Driven Decision-Making

Stochastic optimization problem

$$\underset{x \in X}{\text{minimize}} \quad c(x, \theta)$$

Family of probability measures

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

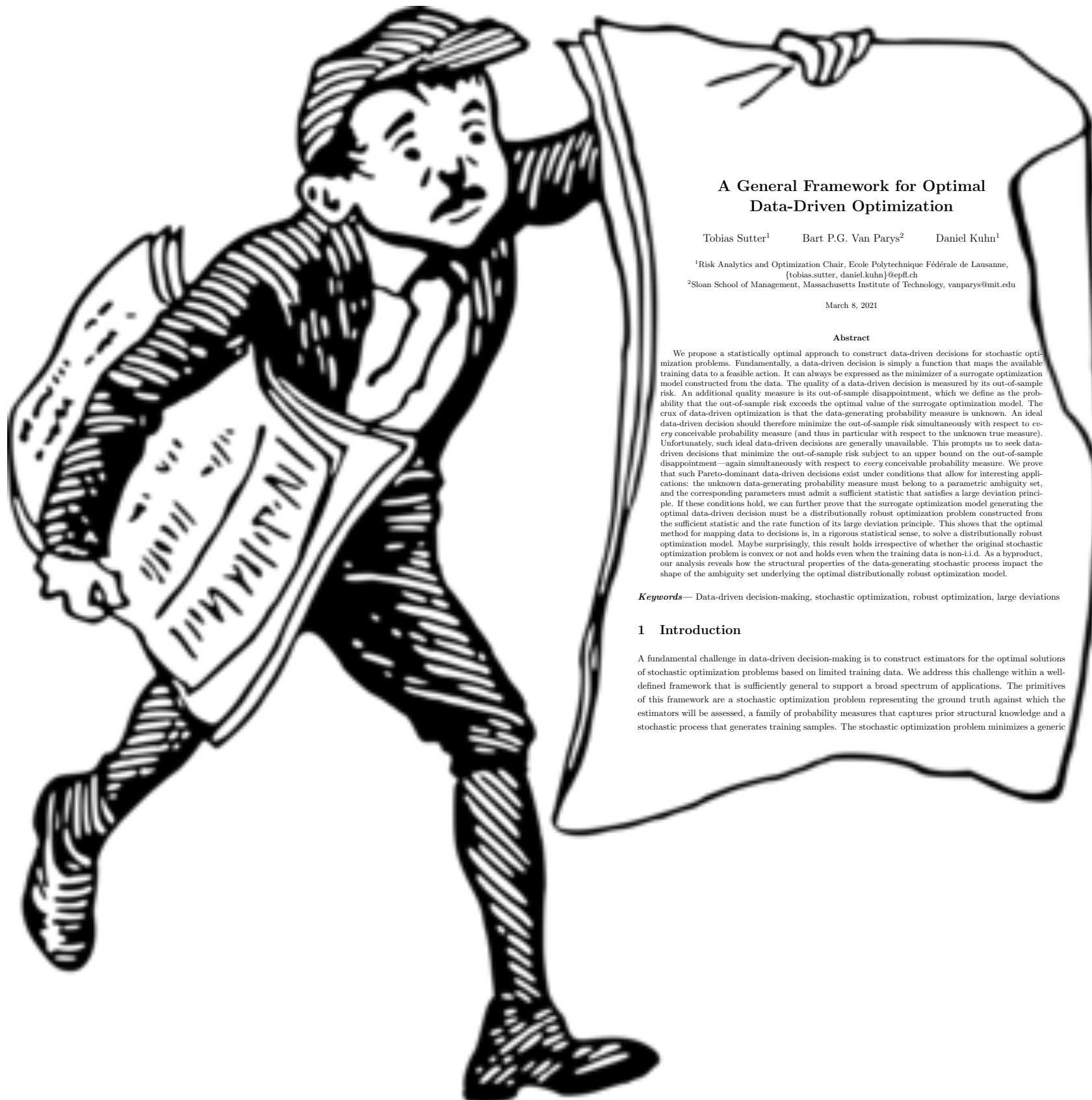
Data-generating process

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

## Examples:

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

# Example: Newsvendor Problem



## A General Framework for Optimal Data-Driven Optimization

Tobias Sutter<sup>1</sup>   Bart P.G. Van Parys<sup>2</sup>   Daniel Kuhn<sup>1</sup>

<sup>1</sup>Risk Analytics and Optimization Chair, Ecole Polytechnique Fédérale de Lausanne,  
{tobias.sutter, daniel.kuhn}@epfl.ch

<sup>2</sup>Sloan School of Management, Massachusetts Institute of Technology, vanparys@mit.edu

March 8, 2021

### Abstract

We propose a statistically optimal approach to construct data-driven decisions for stochastic optimization problems. Fundamentally, a data-driven decision is simply a function that maps the available training data to a feasible action. It can always be expressed as the minimizer of a surrogate optimization model constructed from the data. The quality of a data-driven decision is measured by its out-of-sample risk. An additional quality measure is its out-of-sample disappointment, which we define as the probability that the out-of-sample risk exceeds the optimal value of the surrogate optimization model. The crux of data-driven optimization is that the data-generating probability measure is unknown. An ideal data-driven decision should therefore minimize the out-of-sample risk simultaneously with respect to *every* conceivable probability measure (and thus in particular with respect to the unknown true measure). Unfortunately, such ideal data-driven decisions are generally unavailable. This prompts us to seek data-driven decisions that minimize the out-of-sample risk subject to an upper bound on the out-of-sample disappointment—again simultaneously with respect to *every* conceivable probability measure. We prove that such Pareto-dominant data-driven decisions exist under conditions that allow for interesting applications: the unknown data-generating probability measure must belong to a parametric ambiguity set, and the corresponding parameters must admit a sufficient statistic that satisfies a large deviation principle. If these conditions hold, we can further prove that the surrogate optimization model generating the optimal data-driven decision must be a distributionally robust optimization problem constructed from the sufficient statistic and the rate function of its large deviation principle. This shows that the optimal method for mapping data to decisions is, in a rigorous statistical sense, to solve a distributionally robust optimization model. Maybe surprisingly, this result holds irrespective of whether the original stochastic optimization problem is convex or not and holds even when the training data is non-i.i.d. As a byproduct, our analysis reveals how the structural properties of the data-generating stochastic process impact the shape of the ambiguity set underlying the optimal distributionally robust optimization model.

**Keywords**— Data-driven decision-making, stochastic optimization, robust optimization, large deviations

## 1 Introduction

A fundamental challenge in data-driven decision-making is to construct estimators for the optimal solutions of stochastic optimization problems based on limited training data. We address this challenge within a well-defined framework that is sufficiently general to support a broad spectrum of applications. The primitives of this framework are a stochastic optimization problem representing the ground truth against which the estimators will be assessed, a family of probability measures that captures prior structural knowledge and a stochastic process that generates training samples. The stochastic optimization problem minimizes a generic

# Example: Newsvendor Problem

Stochastic optimization problem

$$\underset{x \in X}{\text{minimize}} \quad c(x, \theta)$$

# Example: Newsvendor Problem

Stochastic optimization problem

$$\underset{x \in X}{\text{minimize}} \quad c(x, \theta)$$

- Order quantities  $x \in X = \{1, \dots, d\}$



# Example: Newsvendor Problem

Stochastic optimization problem

minimize  $c(x, \theta)$   
 $x \in X$

- ▶ Order quantities  $x \in X = \{1, \dots, d\}$
- ▶ Demand  $\xi \in \Xi = \{1, \dots, d\}$

# Example: Newsvendor Problem

Stochastic optimization problem

$$\underset{x \in X}{\text{minimize}} \quad c(x, \theta)$$

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

# Example: Newsvendor Problem

Stochastic optimization problem

$$\underset{x \in X}{\text{minimize}} \quad c(x, \theta)$$

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

wholesale price      retail price



# Example: Newsvendor Problem

Stochastic optimization problem

$$\underset{x \in X}{\text{minimize}} \quad c(x, \theta)$$

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

order  
quantity

sales

# Example: Newsvendor Problem

Stochastic optimization problem

$$\underset{x \in X}{\text{minimize}} \quad c(x, \theta)$$

Data-generating process

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

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

# Example: Newsvendor Problem

Stochastic optimization problem

$$\underset{x \in X}{\text{minimize}} \quad c(x, \theta)$$

Data-generating process

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

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

# Example: Newsvendor Problem

Stochastic optimization problem

$$\underset{x \in X}{\text{minimize}} \quad c(x, \theta)$$

Data-generating process

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

Family of probability measures

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

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

# Example: Newsvendor Problem

Stochastic optimization problem

$$\underset{x \in X}{\text{minimize}} \quad c(x, \theta)$$

Data-generating process

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

Family of probability measures

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

- ▶ Order quantities  $x \in X = \{1, \dots, d\}$
- ▶ Demand  $\xi \in \Xi = \{1, \dots, d\}$
- ▶ Expected cost  $c(x, \theta) = \mathbb{E}_\theta [kx - p \min\{x, \xi\}]$
- ▶ Historical demands  $\xi_t \in \Xi$
- ▶  $\{\xi_t\}_{t \in \mathbb{N}}$  i.i.d. process under  $\mathbb{P}_\theta$



# Example: Newsvendor Problem

Stochastic optimization problem

$$\underset{x \in X}{\text{minimize}} \quad c(x, \theta)$$

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

Data-generating process

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

- ▶ Historical demands  $\xi_t \in \Xi$

Family of probability measures

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

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

# Surrogate Optimization Models

# Surrogate Optimization Models

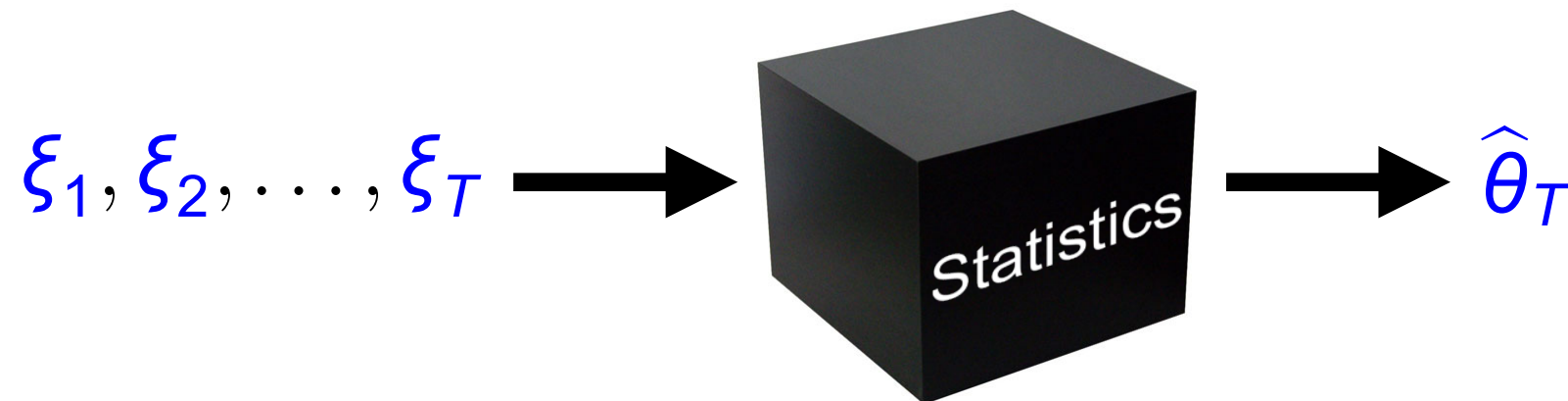
**Original** optimization problem:

$$\underset{x \in X}{\text{minimize}} \quad c(x, \theta)$$

# Surrogate Optimization Models

**Surrogate** optimization problem:

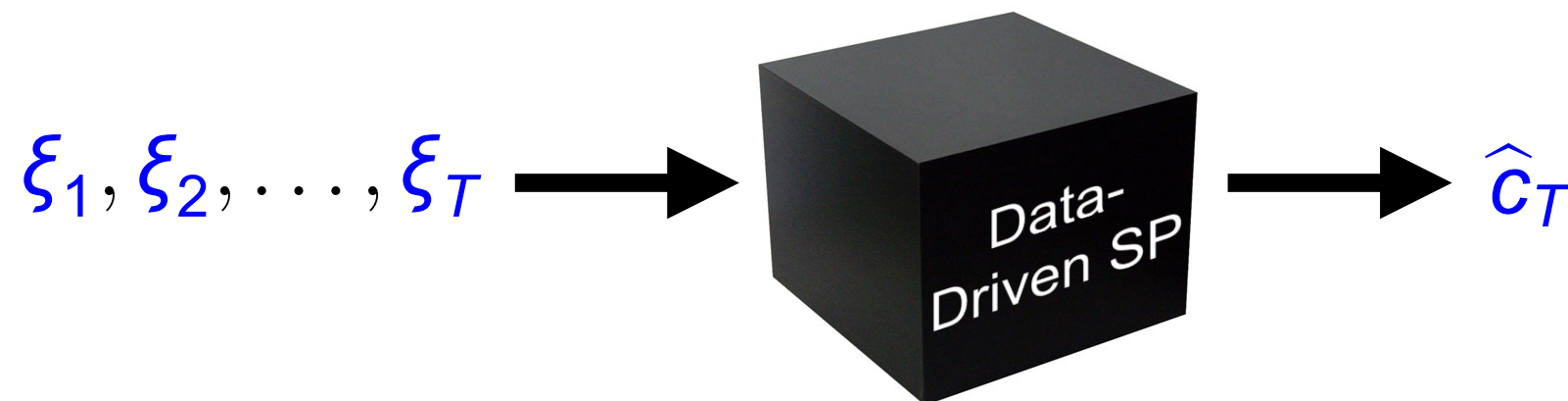
$$\underset{x \in X}{\text{minimize}} \quad c(x, \hat{\theta}_T)$$



# Surrogate Optimization Models

**Surrogate optimization problem:**

$$\underset{x \in X}{\text{minimize}} \quad \hat{c}_T(x)$$



# Surrogate Optimization Models

## Surrogate optimization problem:

$$\underset{x \in X}{\text{minimize}} \quad \hat{c}_T(x)$$

## Construction of $\hat{c}_T$ :

- ▶ Sample average approximation<sup>1)</sup>
- ▶ Regularized nominal model<sup>2)</sup>
- ▶ Predict-then-optimize approach<sup>3)</sup>
- ▶ Neural network model<sup>4)</sup>
- ▶ Distributionally robust optimization model<sup>5)</sup>
- ▶ etc.

---

<sup>1)</sup> Shapiro, *Annals of Statistics*, 1989; <sup>2)</sup> Hoerl & Kennard, *Technometrics*, 1970; <sup>3)</sup> Elmach-  
toub & Grigas, *Management Science*, 2021; <sup>4)</sup> Donti et al., *NIPS*, 2017; <sup>5)</sup> Delage & Ye, *Op-  
erations Research*, 2010; Mohajerin Esfahani & Kuhn, *Mathematical Programming*, 2018.

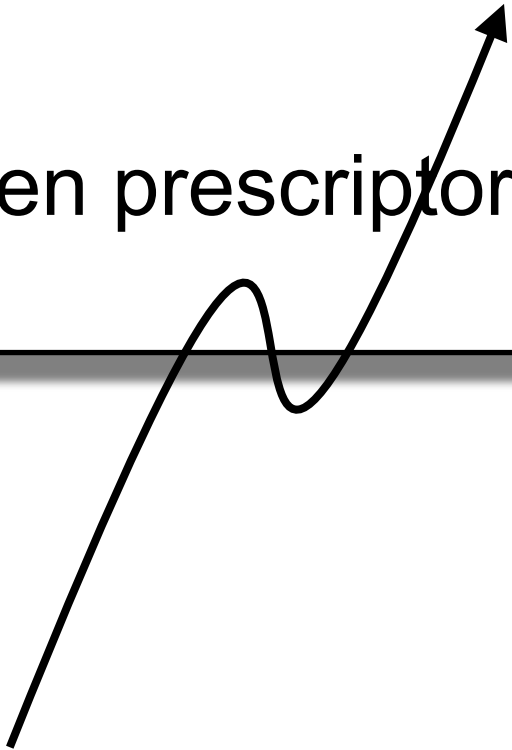
## Definitions:

- ▶ Data-driven predictor  $\hat{c}_T$
- ▶ Data-driven prescriptor  $\hat{x}_T \in \operatorname{argmin}_{x \in X} \hat{c}_T(x)$

## Definitions:

- ▶ Data-driven predictor  $\hat{c}_T$
- ▶ Data-driven prescriptor  $\hat{x}_T \in \operatorname{argmin}_{x \in X} \hat{c}_T(x)$

determines the surrogate  
optimization model

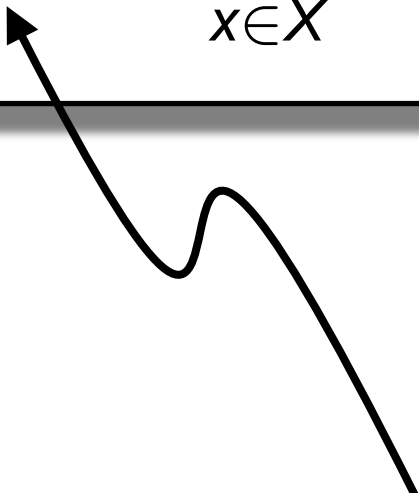




## Definitions:

- ▶ Data-driven predictor  $\hat{c}_T$
- ▶ Data-driven prescriptor  $\hat{x}_T \in \operatorname{argmin}_{x \in X} \hat{c}_T(x)$

*any function that maps*  
 $\xi_1, \xi_2, \dots, \xi_T$  to  $X$



## Definitions:

- ▶ Data-driven predictor  $\hat{c}_T$
- ▶ Data-driven prescriptor  $\hat{x}_T \in \operatorname{argmin}_{x \in X} \hat{c}_T(x)$

## Performance measures:

## Definitions:

- ▶ Data-driven predictor  $\hat{c}_T$
- ▶ Data-driven prescriptor  $\hat{x}_T \in \operatorname{argmin}_{x \in X} \hat{c}_T(x)$

## Performance measures:

In-sample risk  $\hat{c}_T(\hat{x}_T)$

## Definitions:

- ▶ Data-driven predictor  $\hat{c}_T$
- ▶ Data-driven prescriptor  $\hat{x}_T \in \operatorname{argmin}_{x \in X} \hat{c}_T(x)$

## Performance measures:

In-sample risk  $\hat{c}_T(\hat{x}_T)$

Out-of-sample risk  $c(\hat{x}_T, \theta)$

## Definitions:

- ▶ Data-driven predictor  $\hat{c}_T$
- ▶ Data-driven prescriptor  $\hat{x}_T \in \operatorname{argmin}_{x \in X} \hat{c}_T(x)$

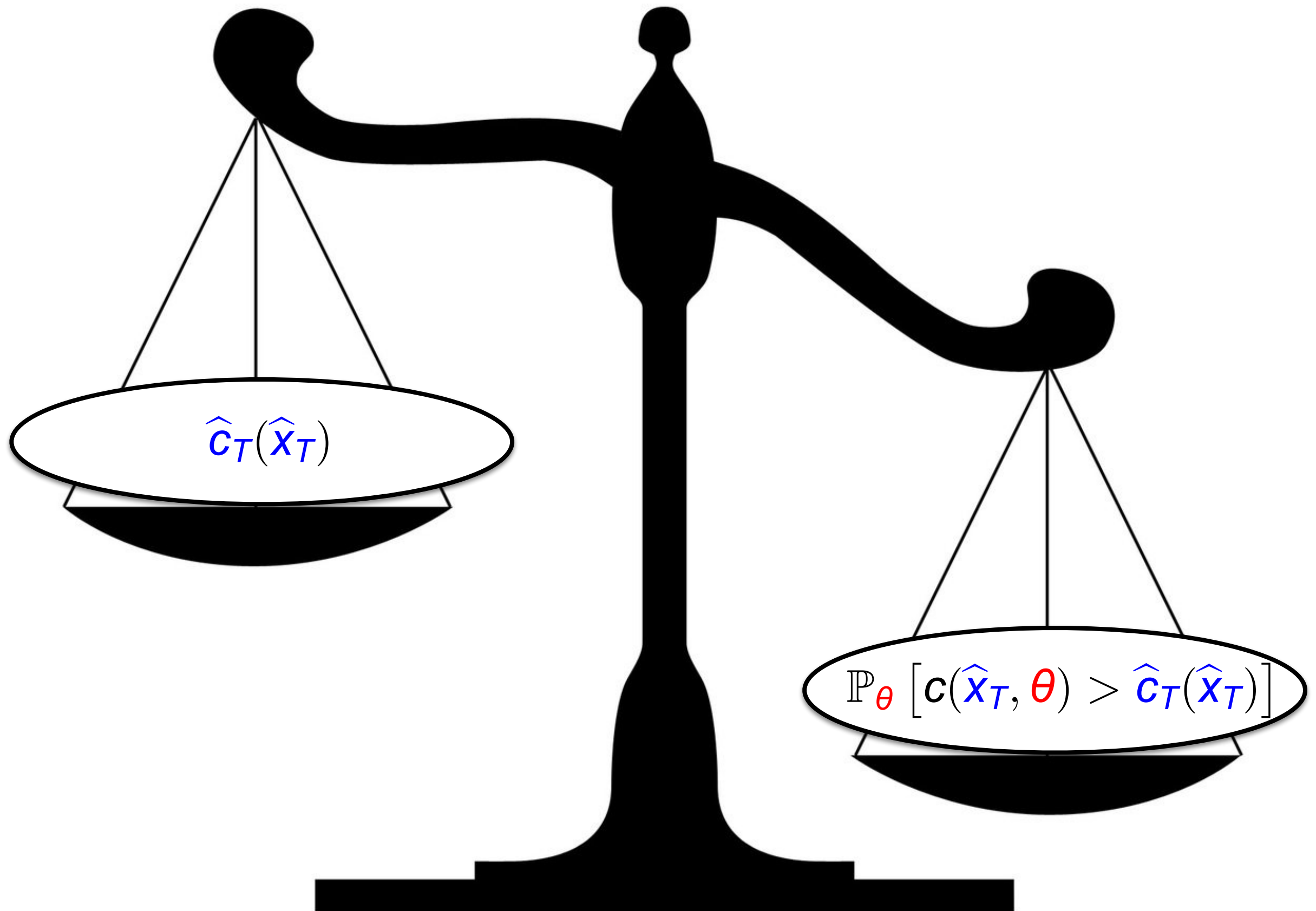
## Performance measures:

In-sample risk  $\hat{c}_T(\hat{x}_T)$

Out-of-sample risk  $c(\hat{x}_T, \theta)$

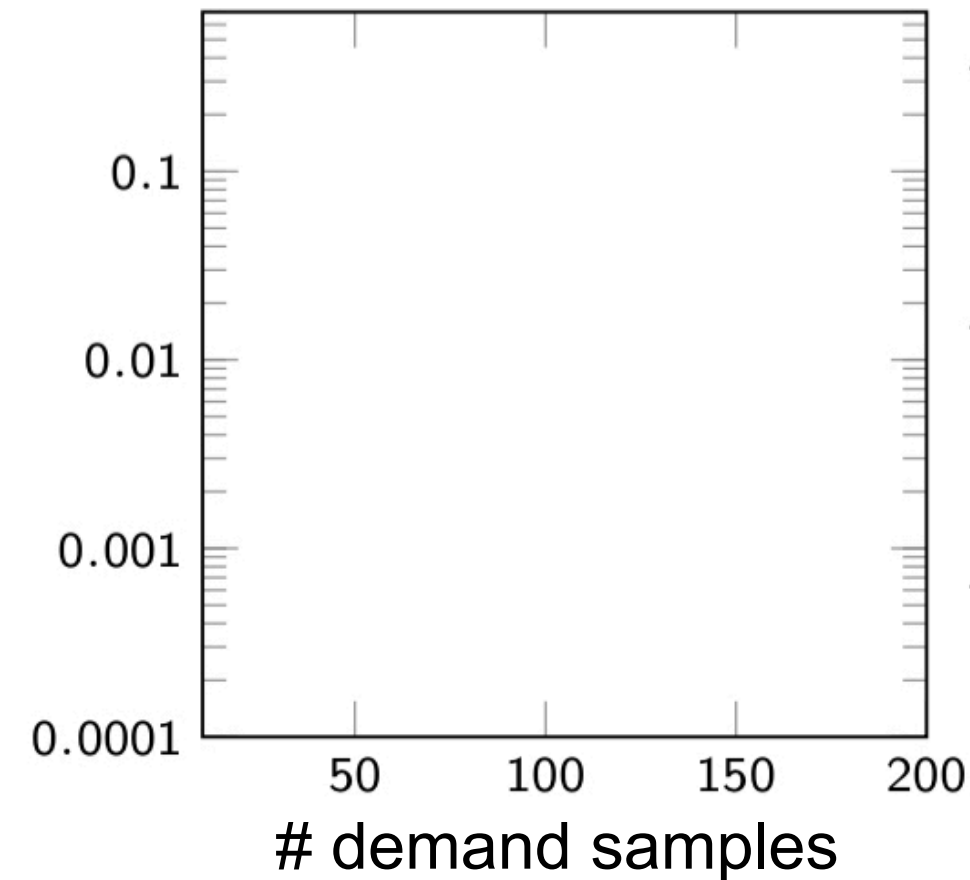
Out-of-sample disappointment  $\mathbb{P}_\theta [c(\hat{x}_T, \theta) > \hat{c}_T(\hat{x}_T)]$

# A Basic Trade-Off

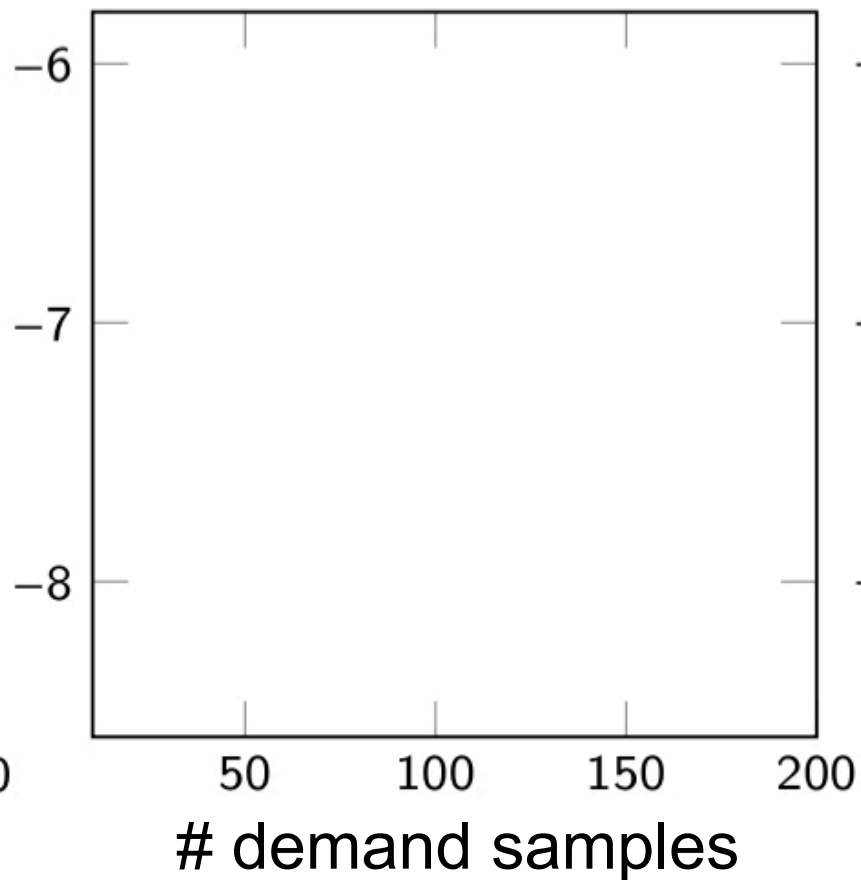


# Data-Driven Newsvendor Problem

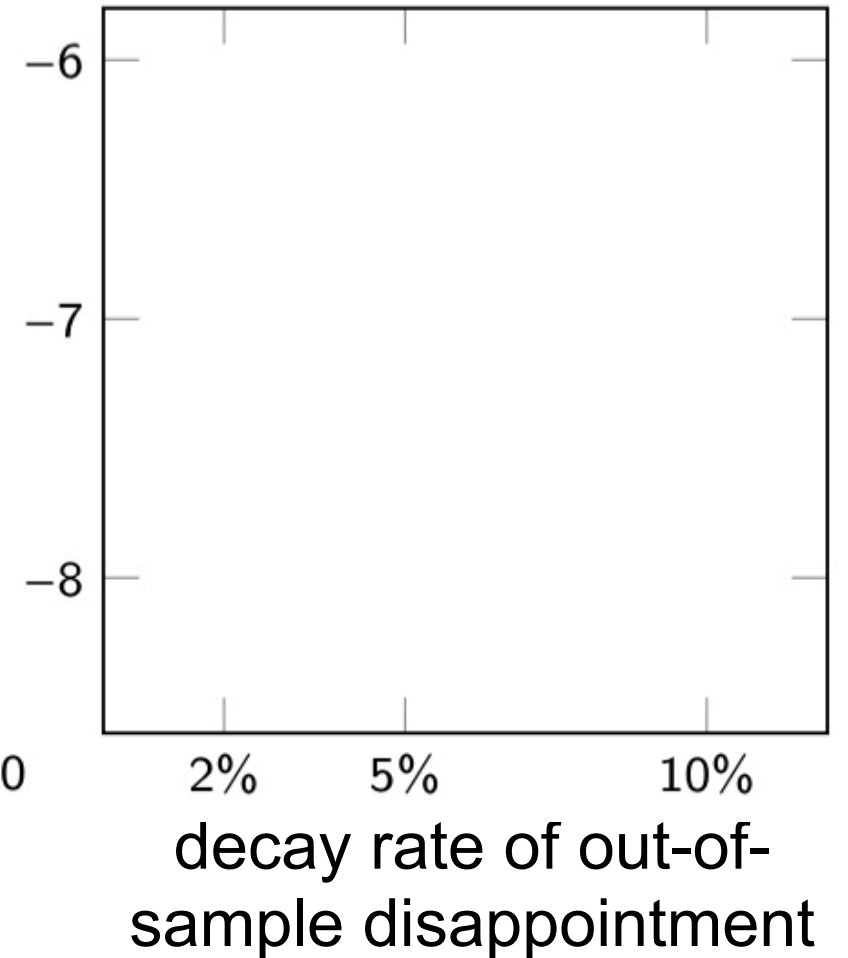
out-of-sample  
disappointment



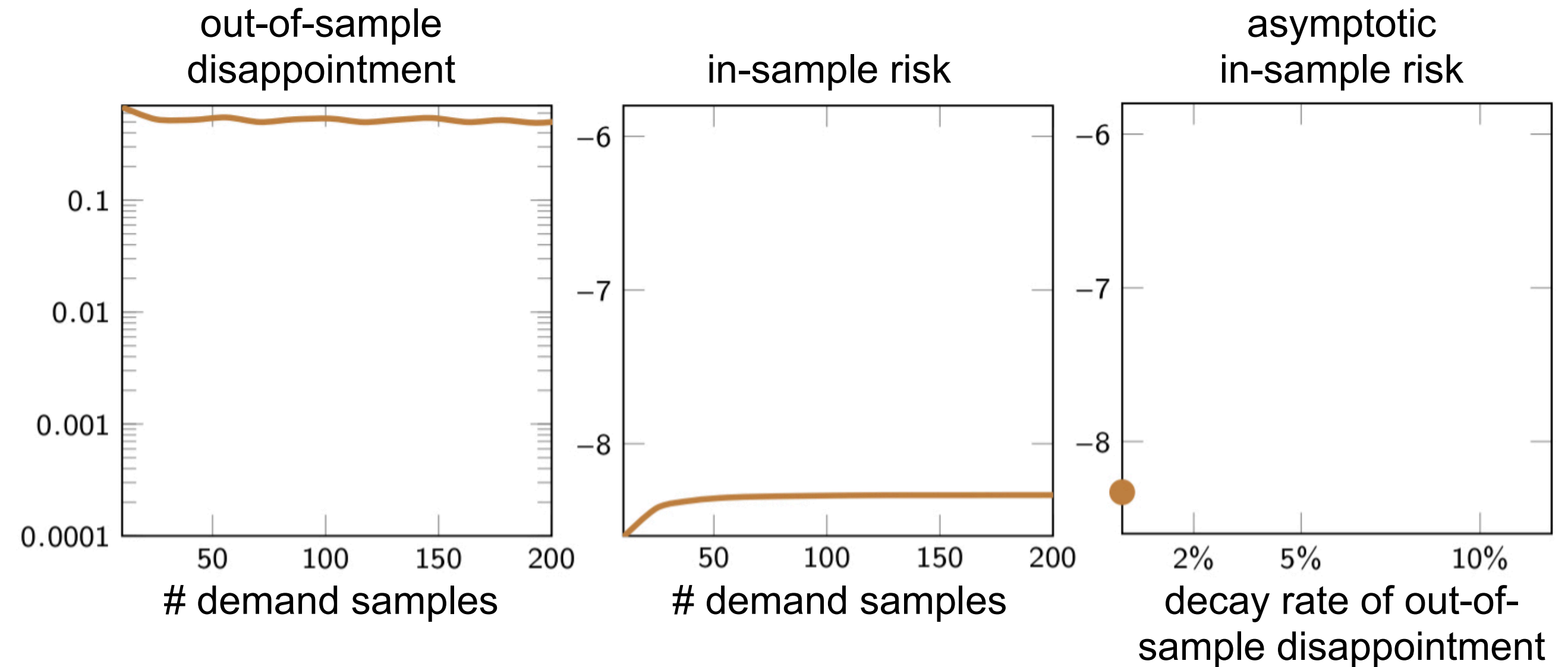
in-sample risk



asymptotic  
in-sample risk



# Data-Driven Newsvendor Problem



**Model 1: SAA model<sup>1)</sup>**

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

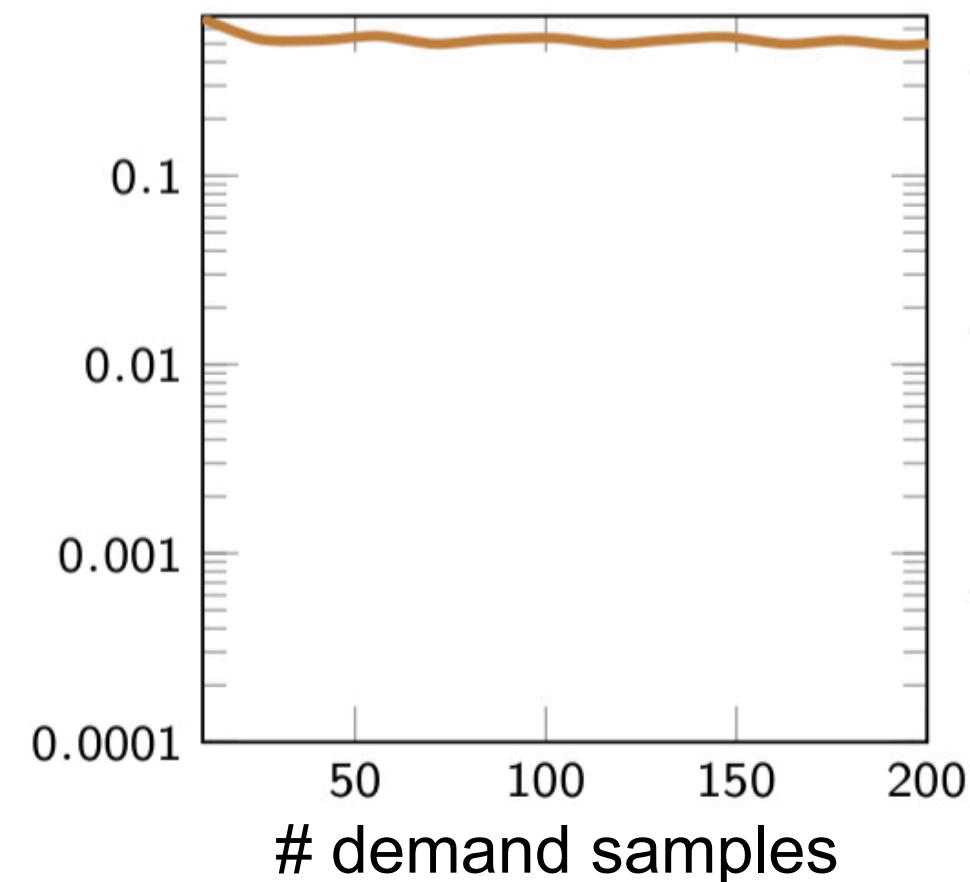
---

<sup>1)</sup> Shapiro, *Annals of Statistics*, 1989.

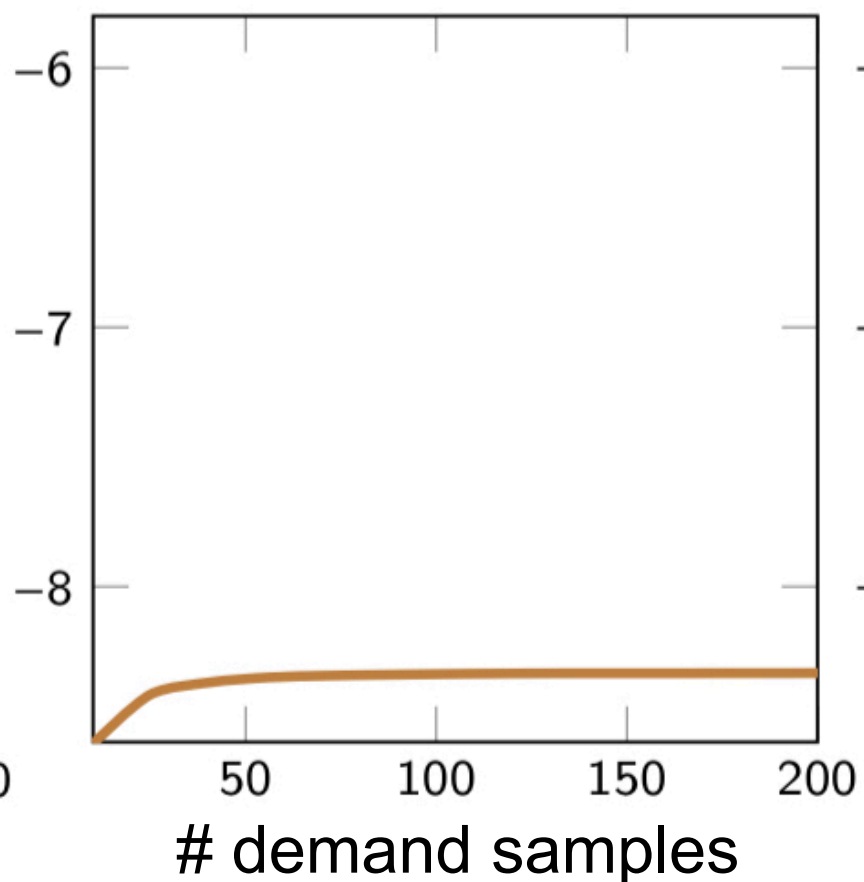


# Data-Driven Newsvendor Problem

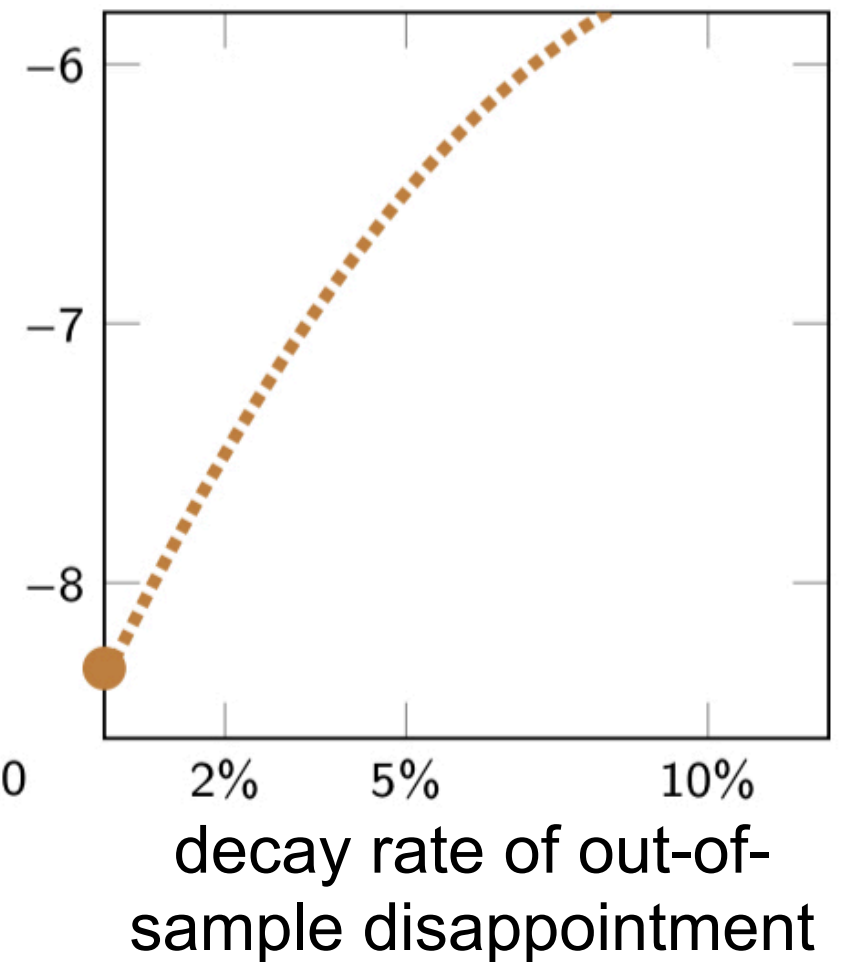
out-of-sample  
disappointment



in-sample risk



asymptotic  
in-sample risk

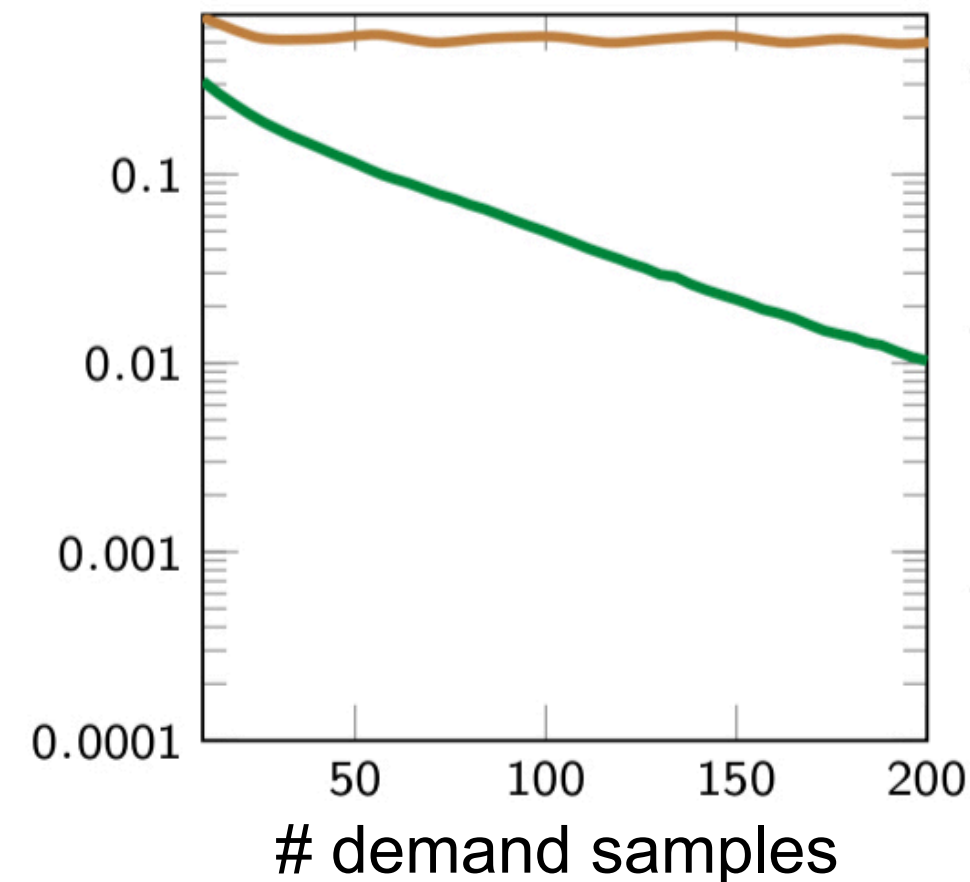


**Model 2:** SAA model with offset

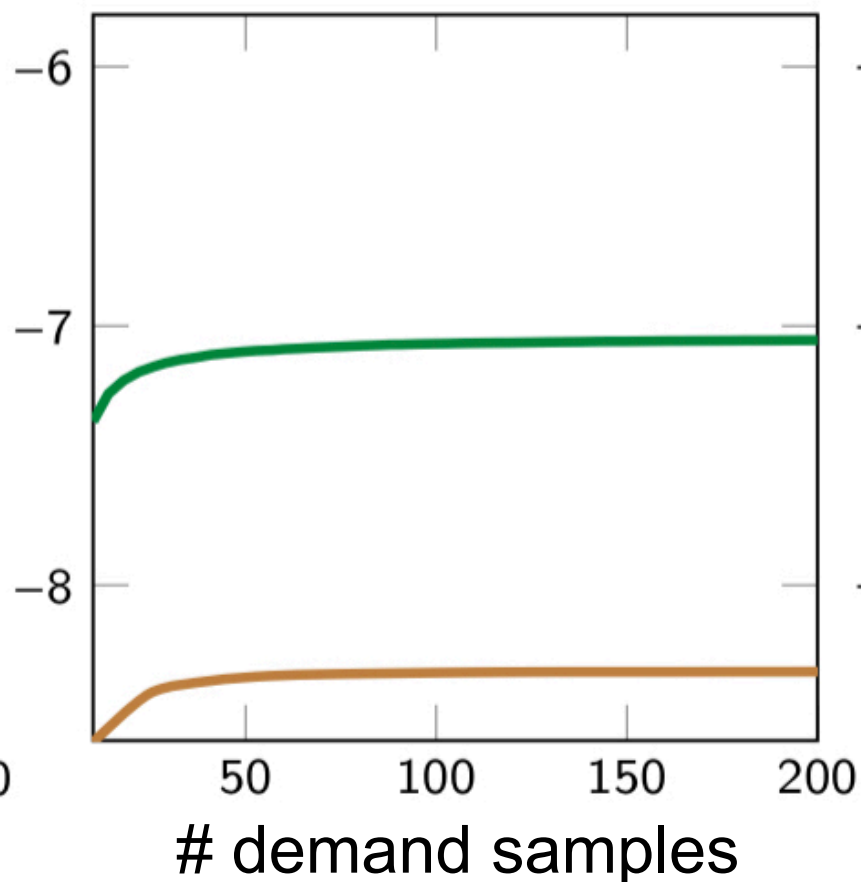
$$\hat{c}_T(x) = c(x, \hat{\theta}_T) + r$$

# Data-Driven Newsvendor Problem

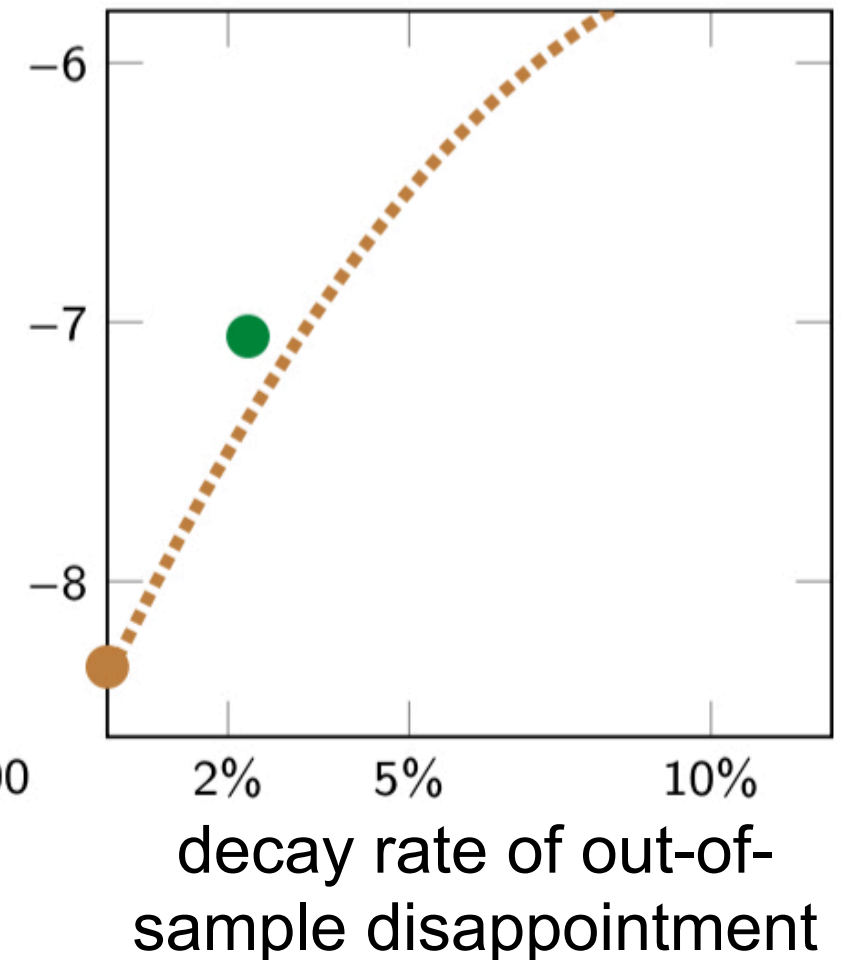
out-of-sample  
disappointment



in-sample risk



asymptotic  
in-sample risk

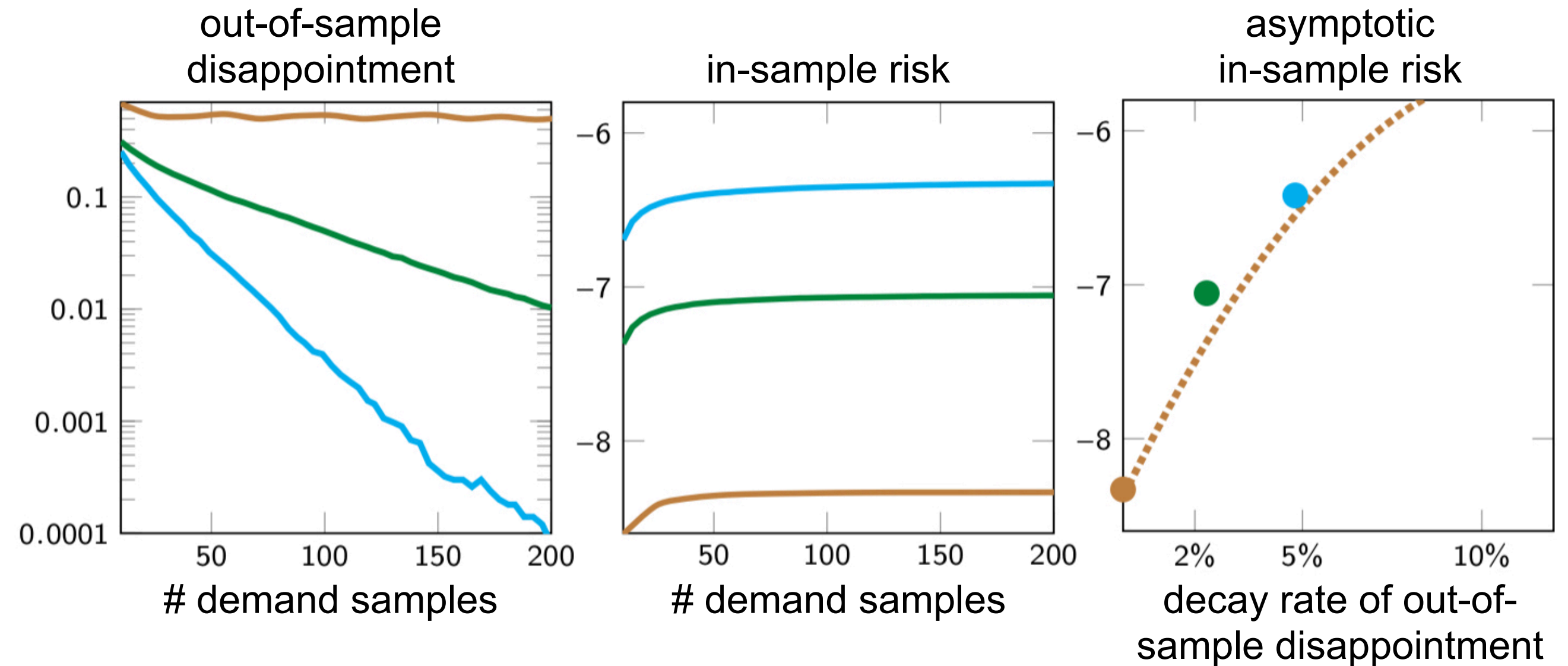


**Model 3:** DRO model with moment ambiguity set<sup>1)</sup>

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

<sup>1)</sup> Delage & Ye, *Operations Research*, 2010.

# Data-Driven Newsvendor Problem



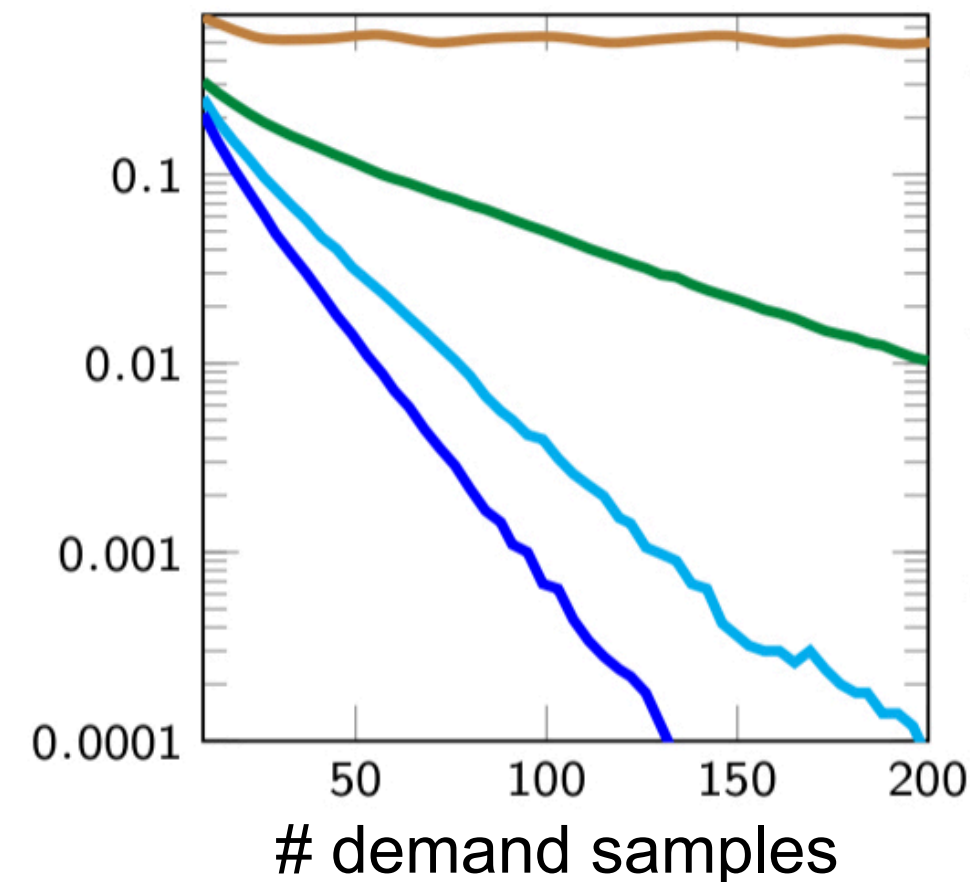
**Model 3:** DRO model with moment ambiguity set<sup>1)</sup>

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

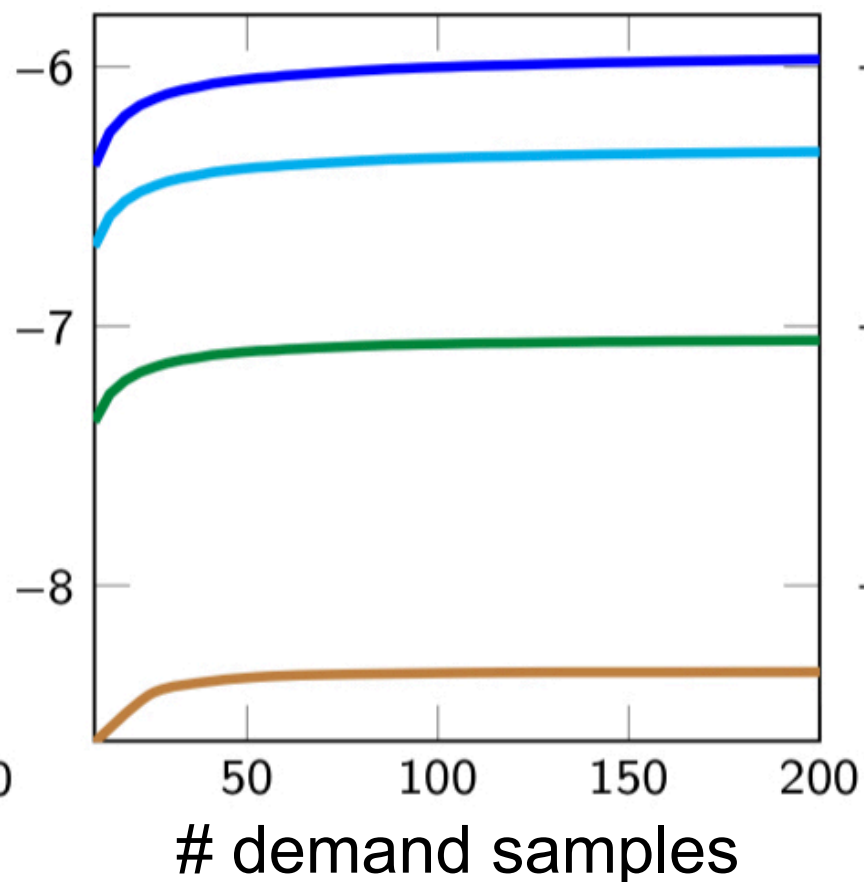
<sup>1)</sup> Delage & Ye, *Operations Research*, 2010.

# Data-Driven Newsvendor Problem

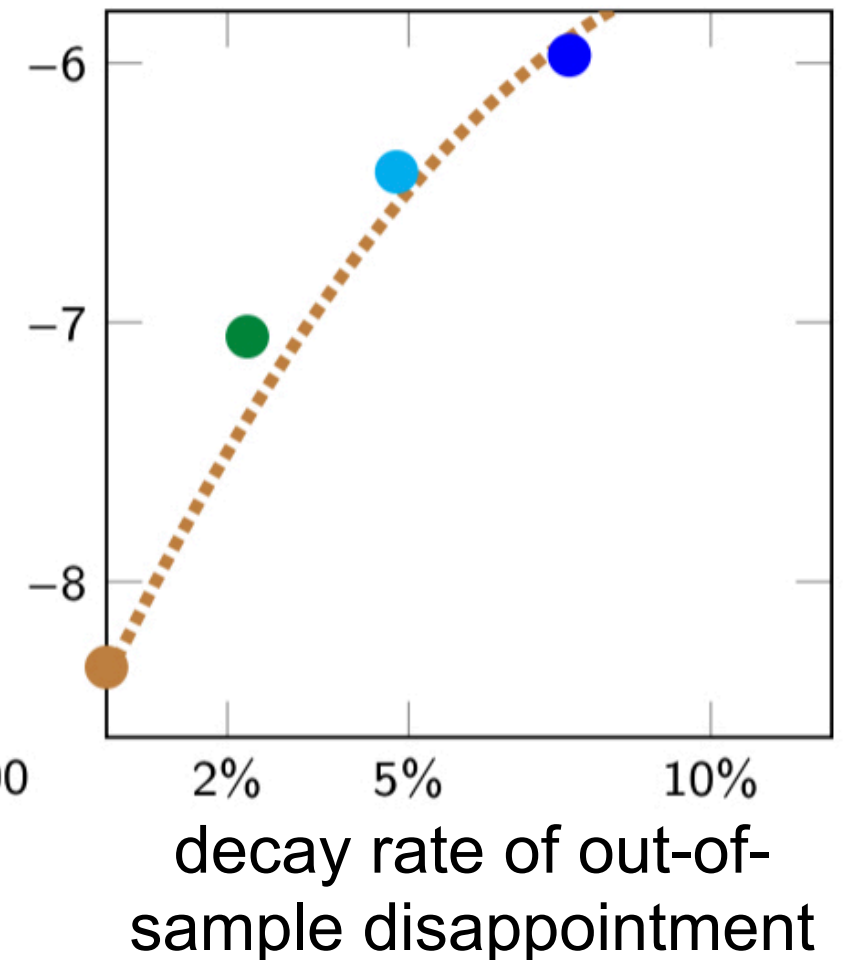
out-of-sample  
disappointment



in-sample risk



asymptotic  
in-sample risk



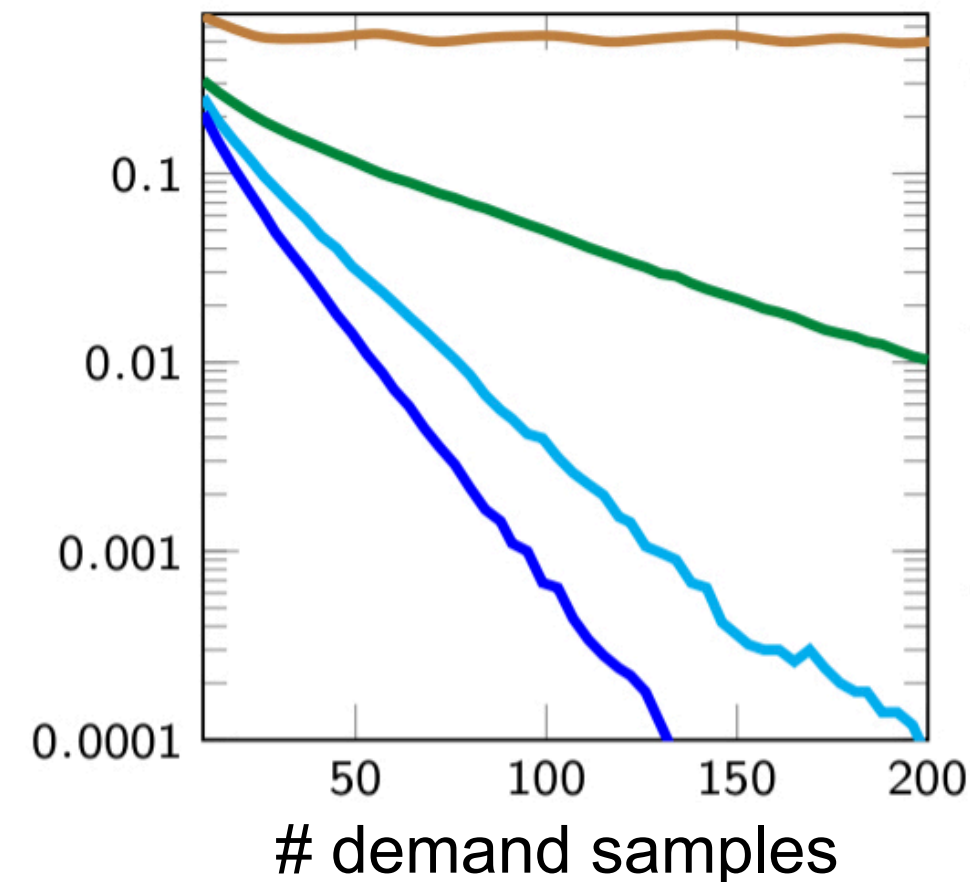
**Model 3:** DRO model with moment ambiguity set<sup>1)</sup>

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

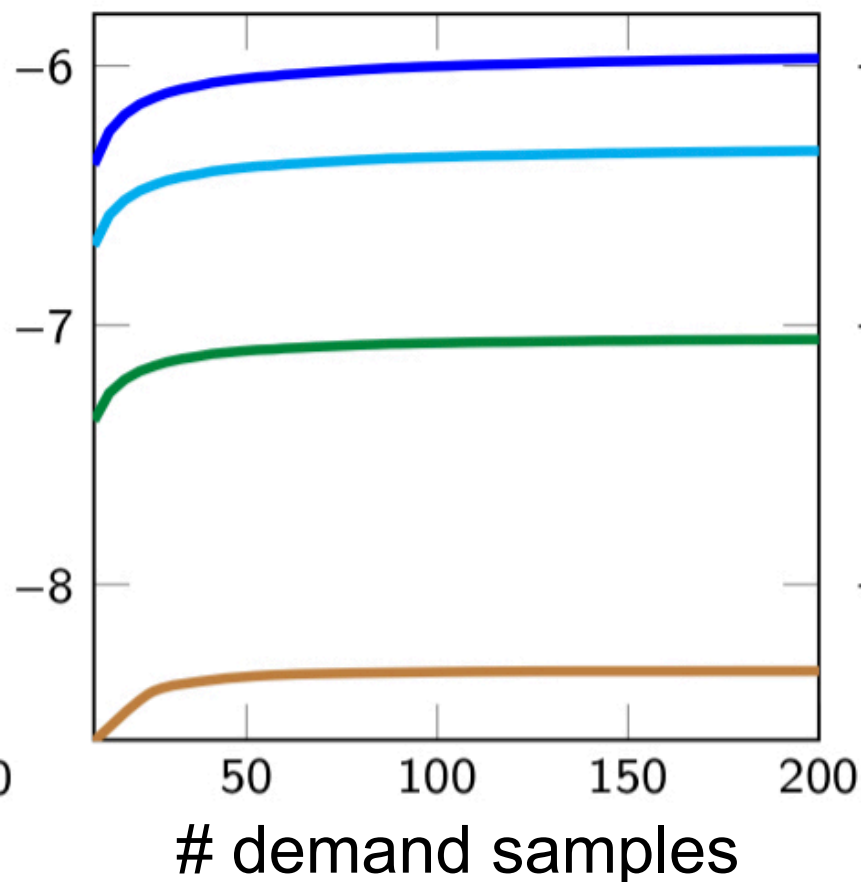
<sup>1)</sup> Delage & Ye, *Operations Research*, 2010.

# Data-Driven Newsvendor Problem

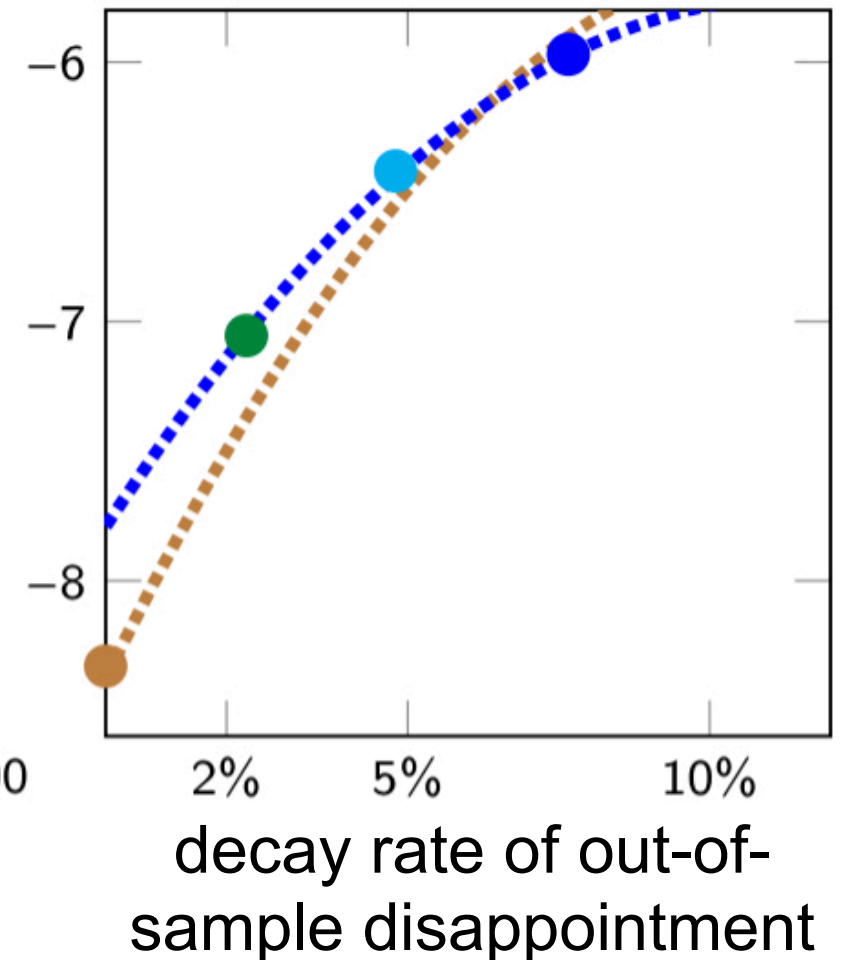
out-of-sample  
disappointment



in-sample risk



asymptotic  
in-sample risk



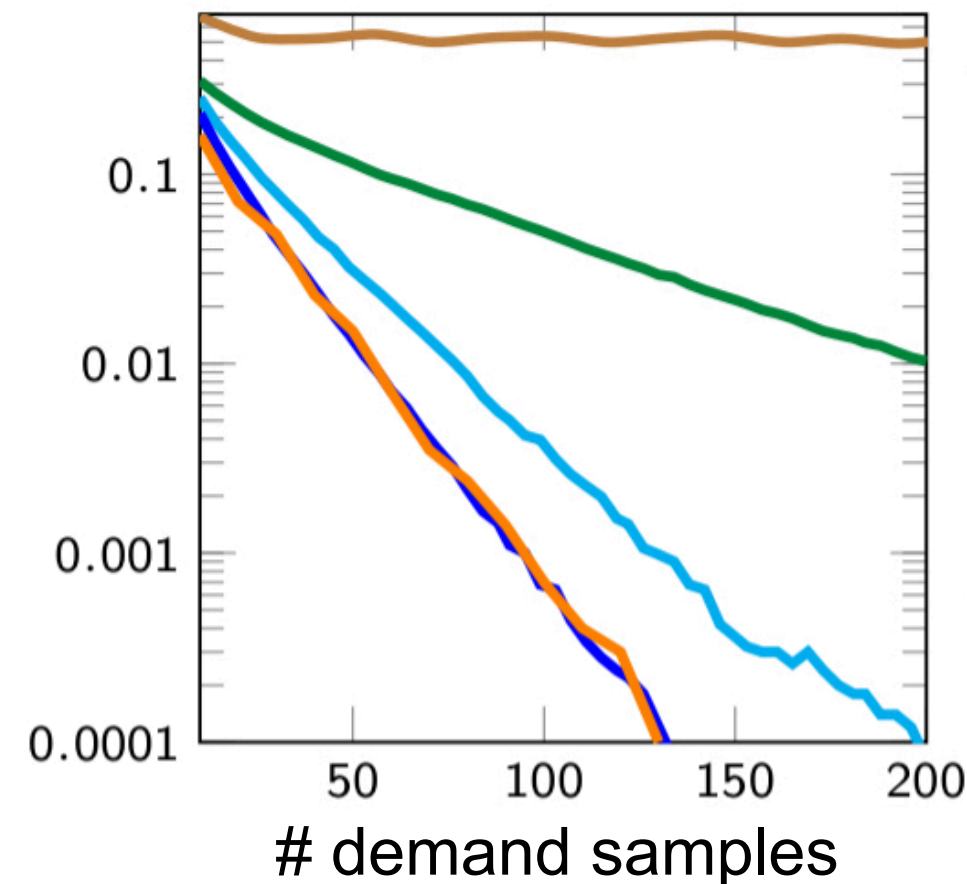
**Model 3:** DRO model with moment ambiguity set<sup>1)</sup>

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

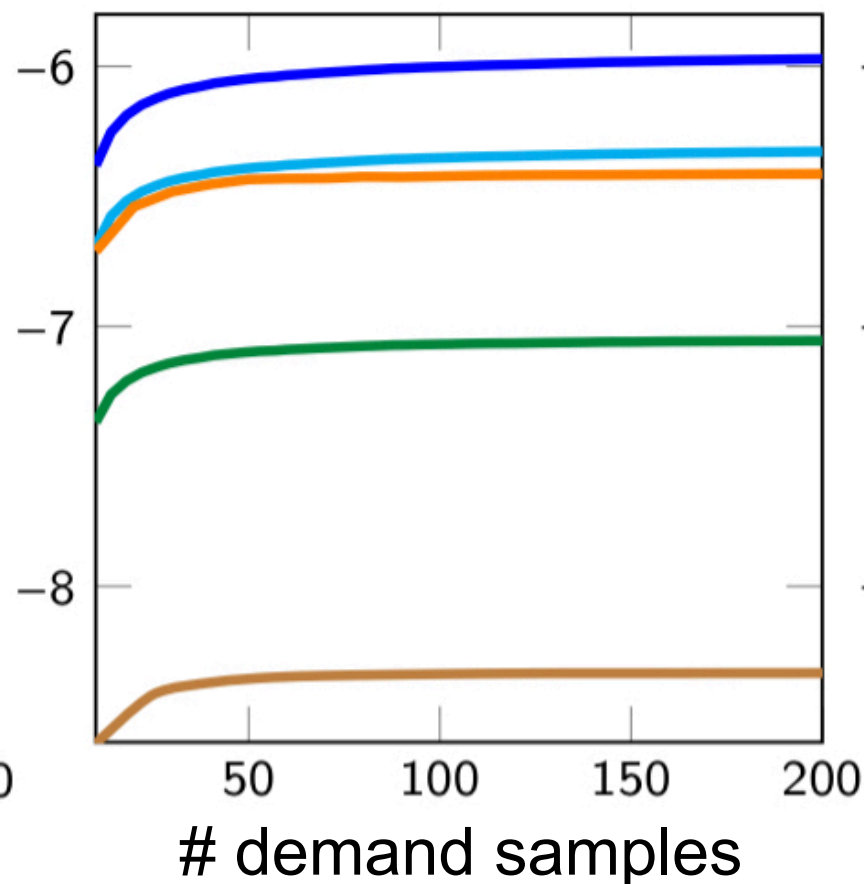
<sup>1)</sup> Delage & Ye, *Operations Research*, 2010.

# Data-Driven Newsvendor Problem

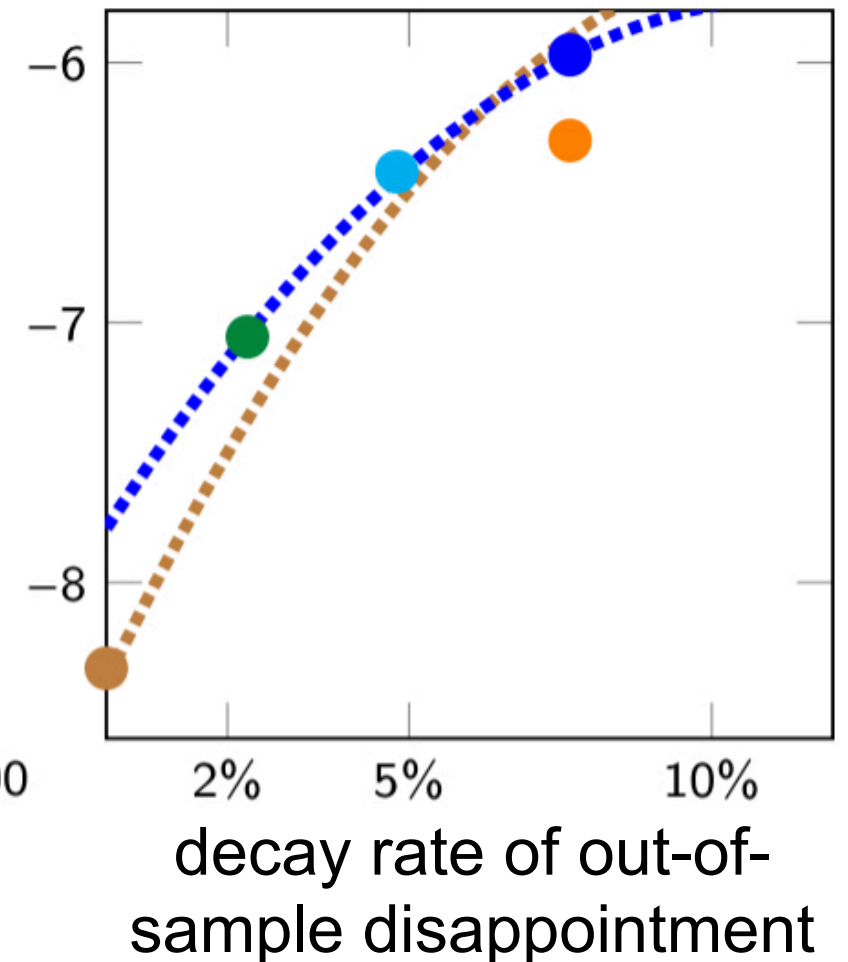
out-of-sample  
disappointment



in-sample risk



asymptotic  
in-sample risk



**Model 4:** DRO model with Wasserstein ambiguity set<sup>1)</sup>

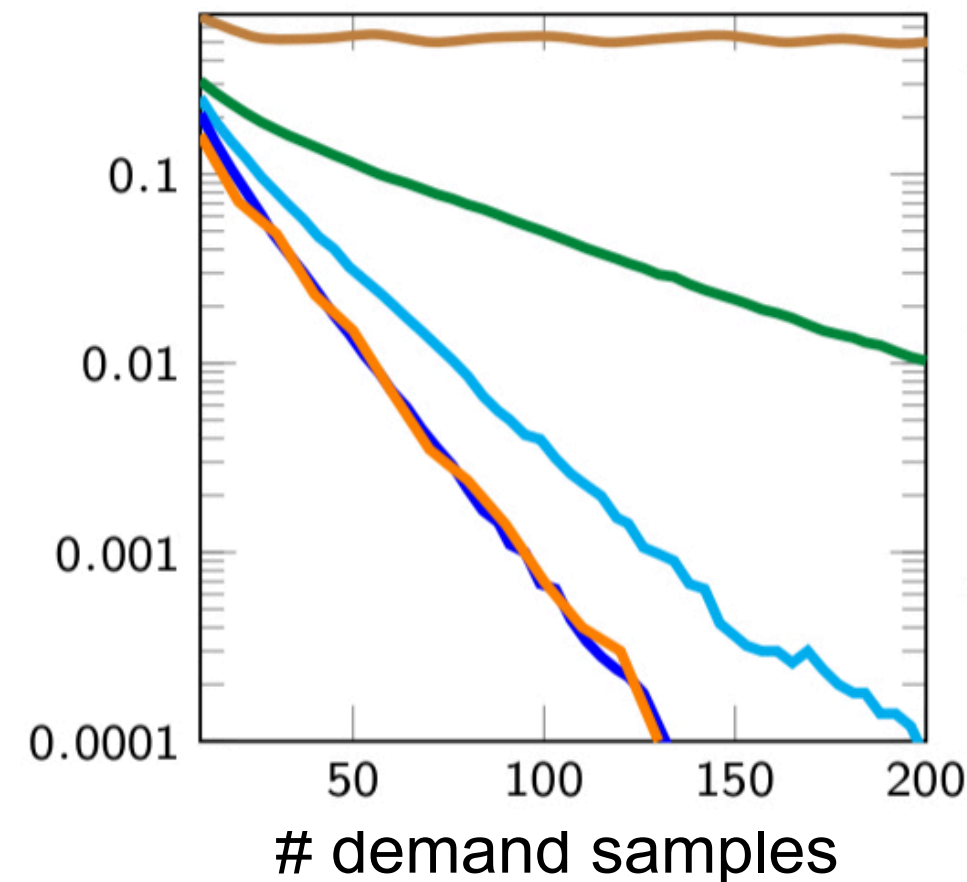
$$\hat{c}_T(x) = \sup_{\theta \in \Theta} \left\{ c(x, \theta) : d_W(\hat{\theta}_T || \theta) \leq r \right\}$$

<sup>1)</sup> Mohajerin Esfahani & Kuhn, *Mathematical Programming*, 2018.

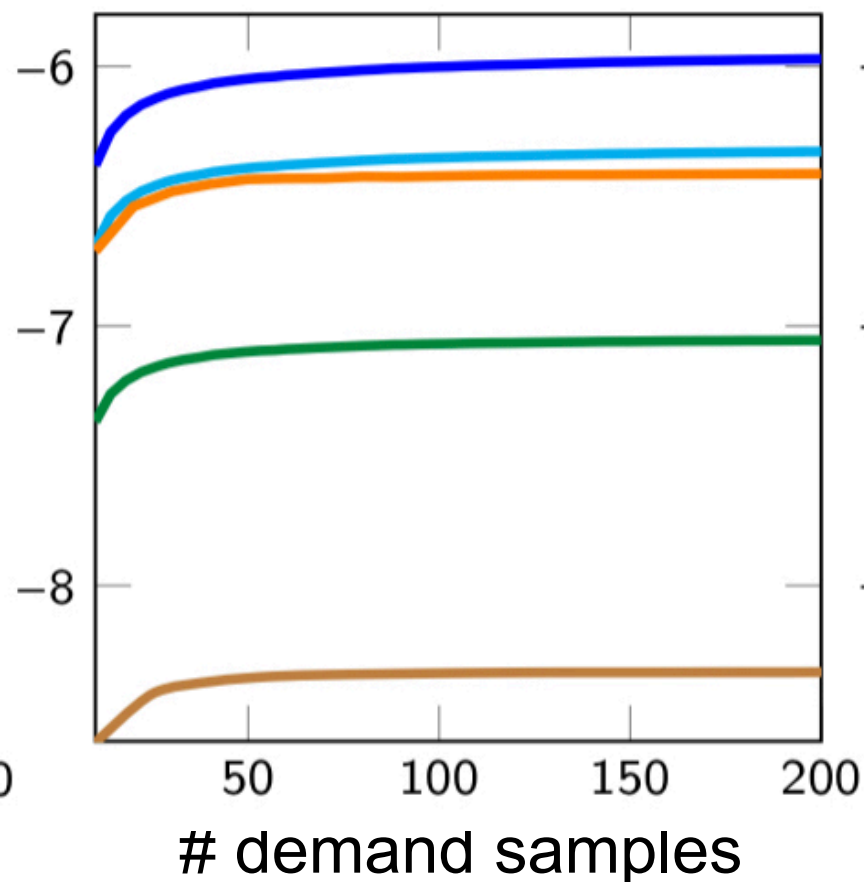


# Data-Driven Newsvendor Problem

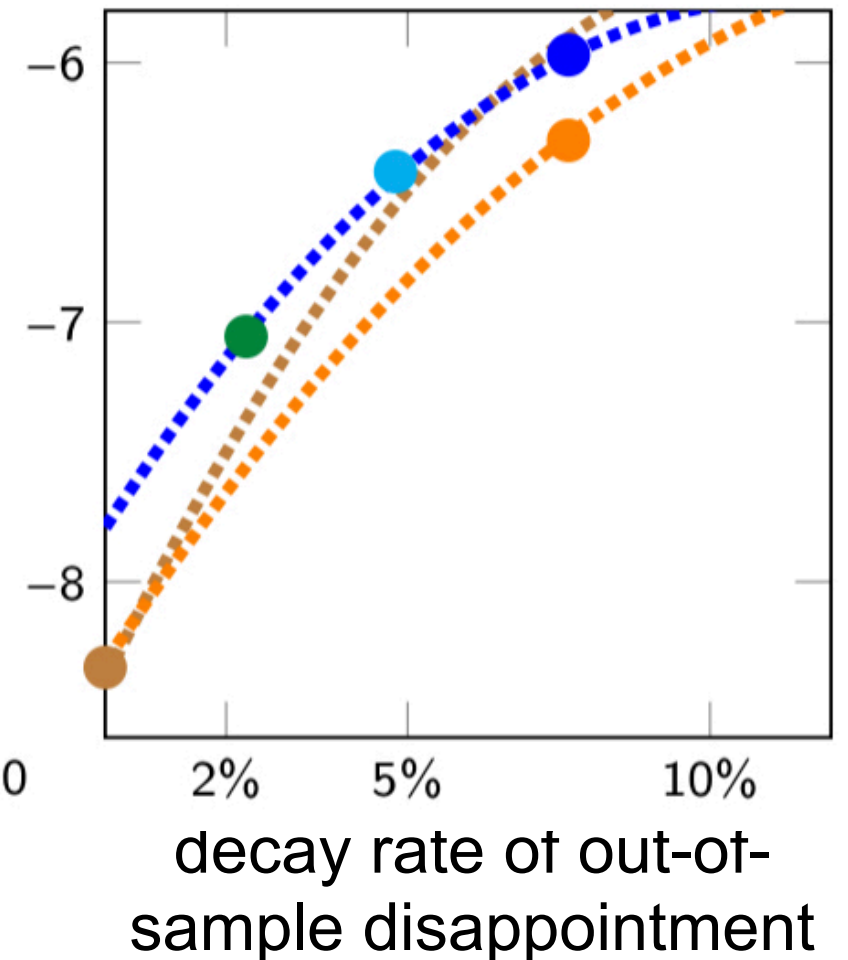
out-of-sample  
disappointment



in-sample risk



asymptotic  
in-sample risk

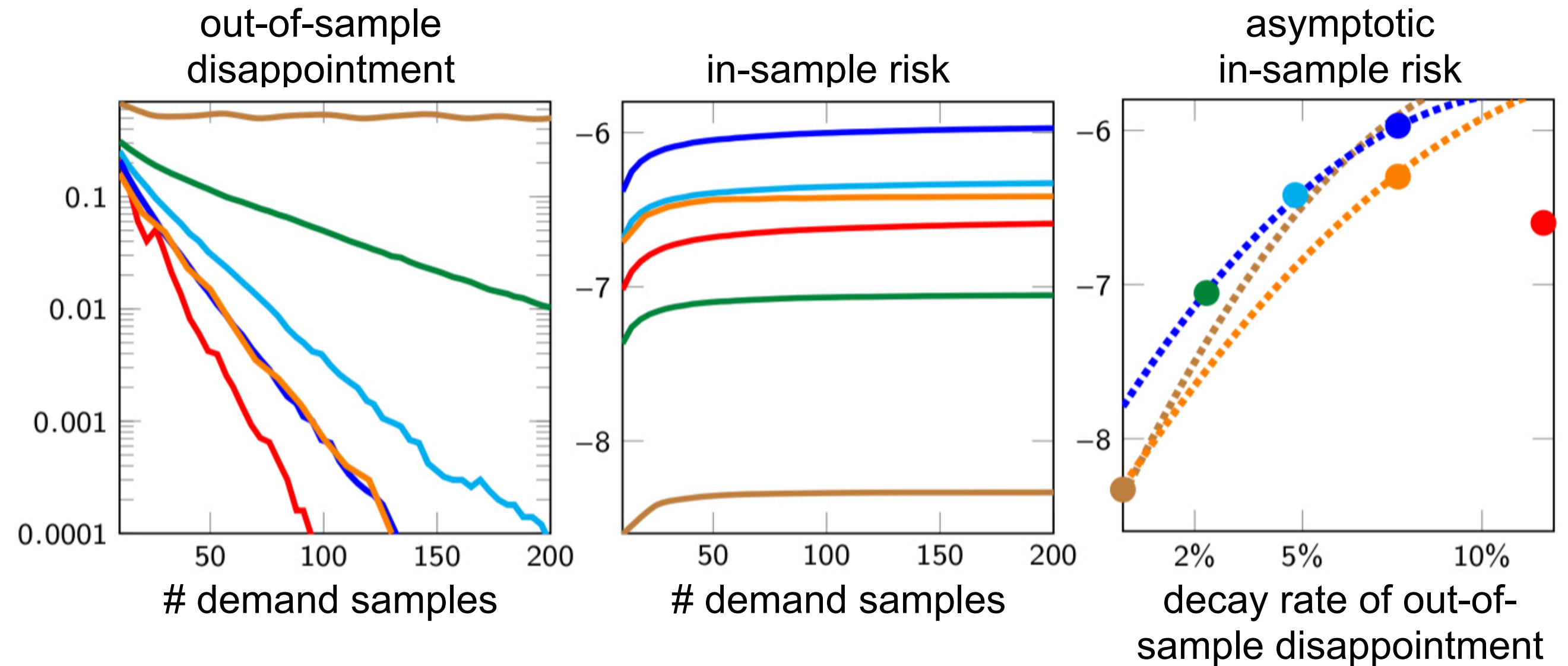


**Model 4:** DRO model with Wasserstein ambiguity set<sup>1)</sup>

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

<sup>1)</sup> Mohajerin Esfahani & Kuhn, *Mathematical Programming*, 2018.

# Data-Driven Newsvendor Problem



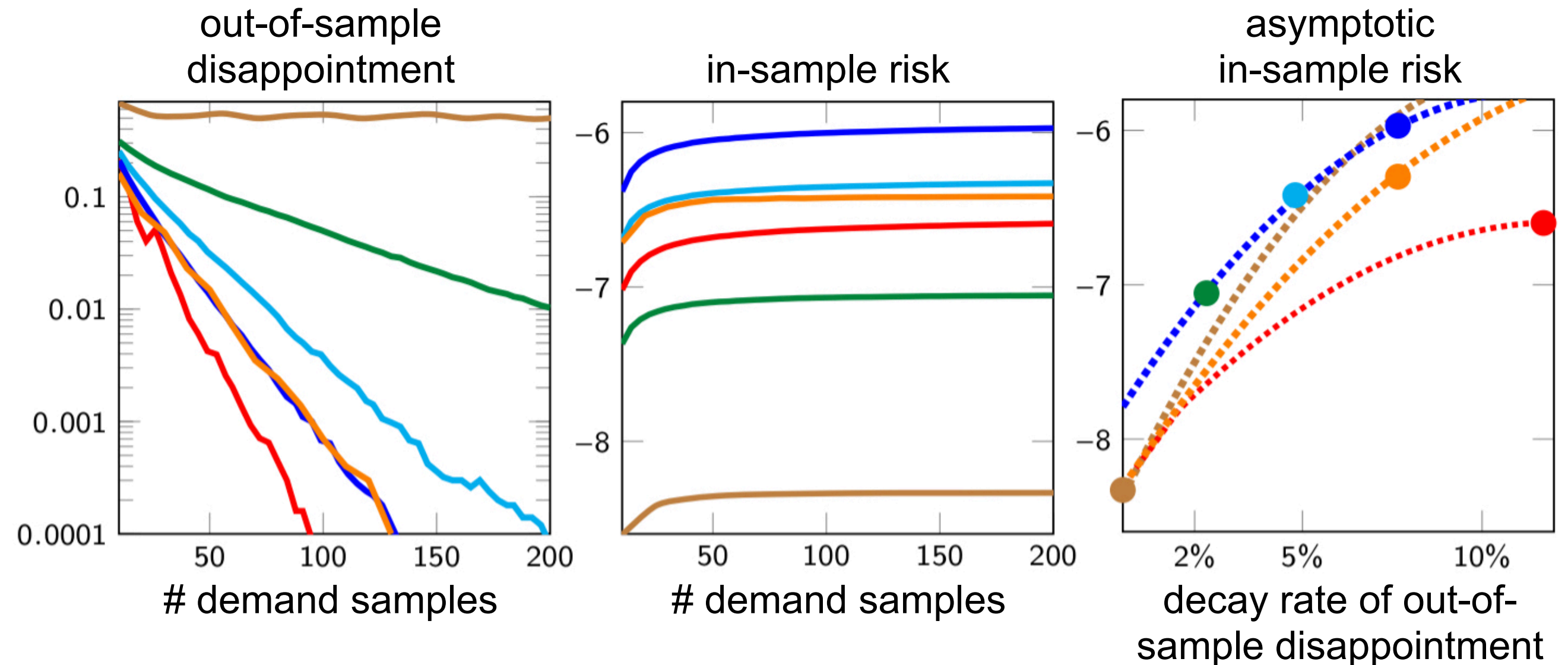
**Model 5:** DRO model with KL ambiguity set<sup>1)</sup>

$$\hat{c}_T(x) = \sup_{\theta \in \Theta} \left\{ c(x, \theta) : D_{\text{KL}}(\hat{\theta}_T \| \theta) \leq r \right\}$$

<sup>1)</sup> Ben-Tal et al., *Management Science*, 2013.



# Data-Driven Newsvendor Problem



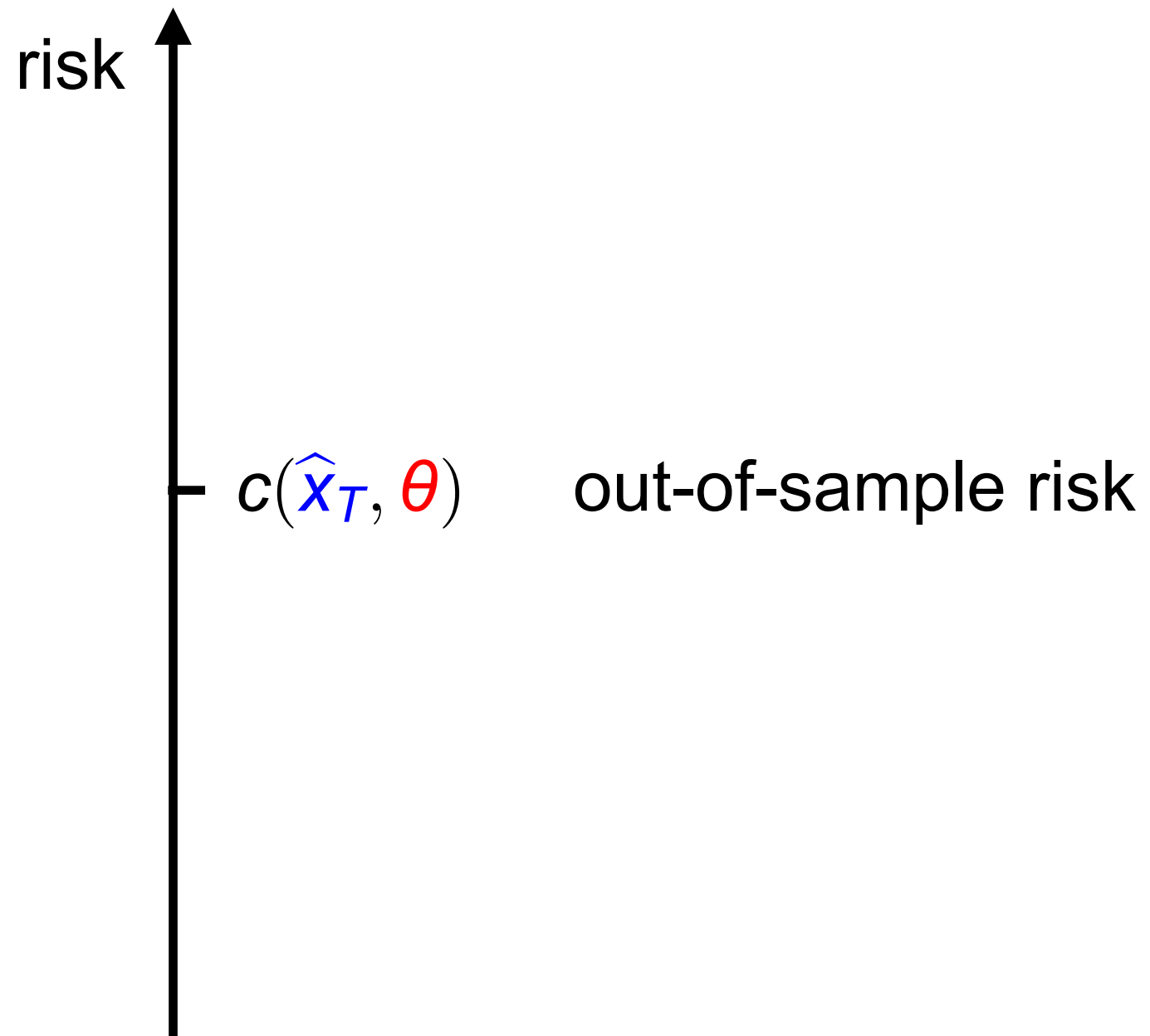
**Model 5:** DRO model with KL ambiguity set<sup>1)</sup>

$$\hat{c}_T(x) = \sup_{\theta \in \Theta} \left\{ c(x, \theta) : D_{\text{KL}}(\hat{\theta}_T \| \theta) \leq r \right\}$$

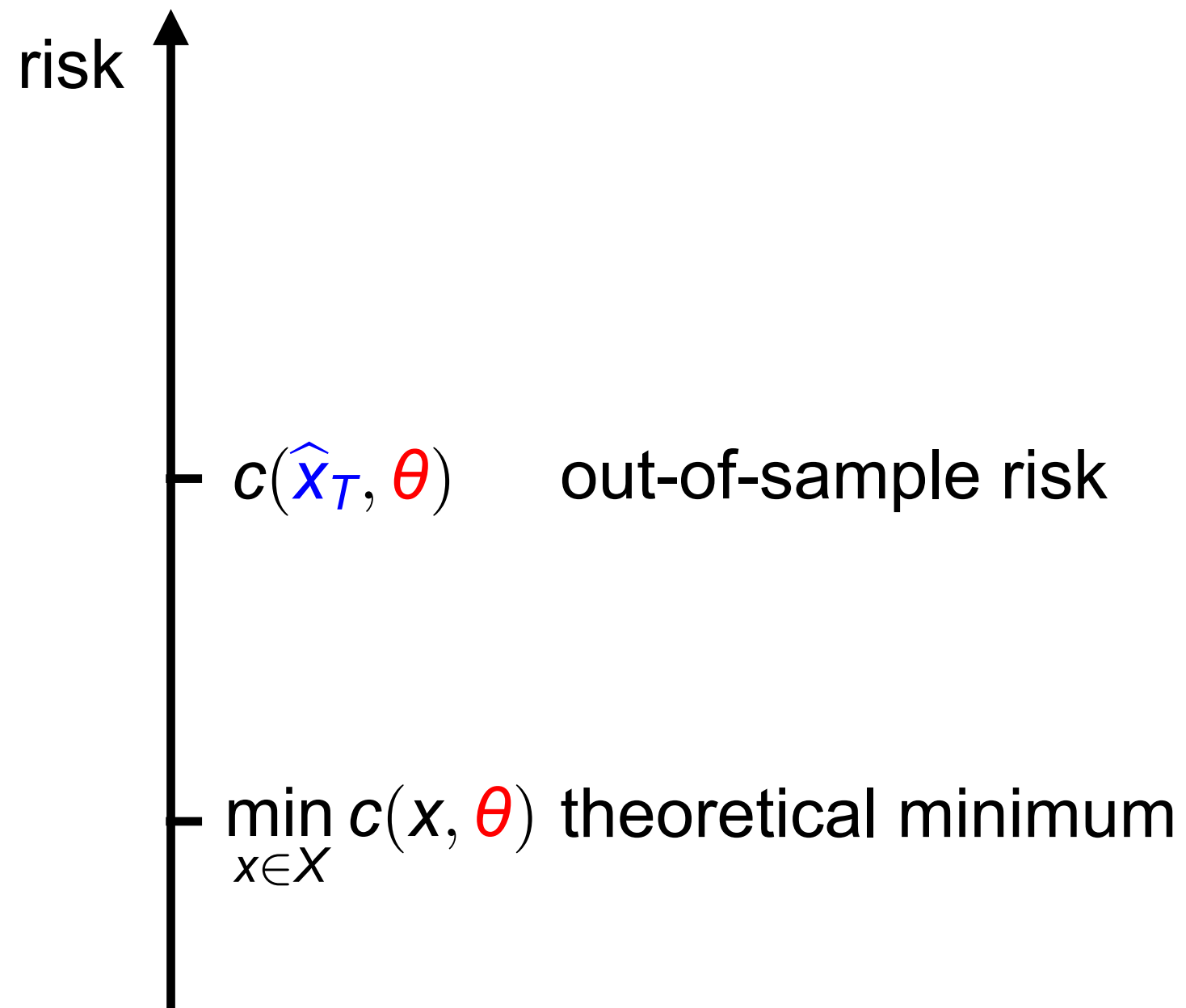
<sup>1)</sup> Ben-Tal et al., *Management Science*, 2013.

# Constructing “Optimal” Surrogate Optimization Models

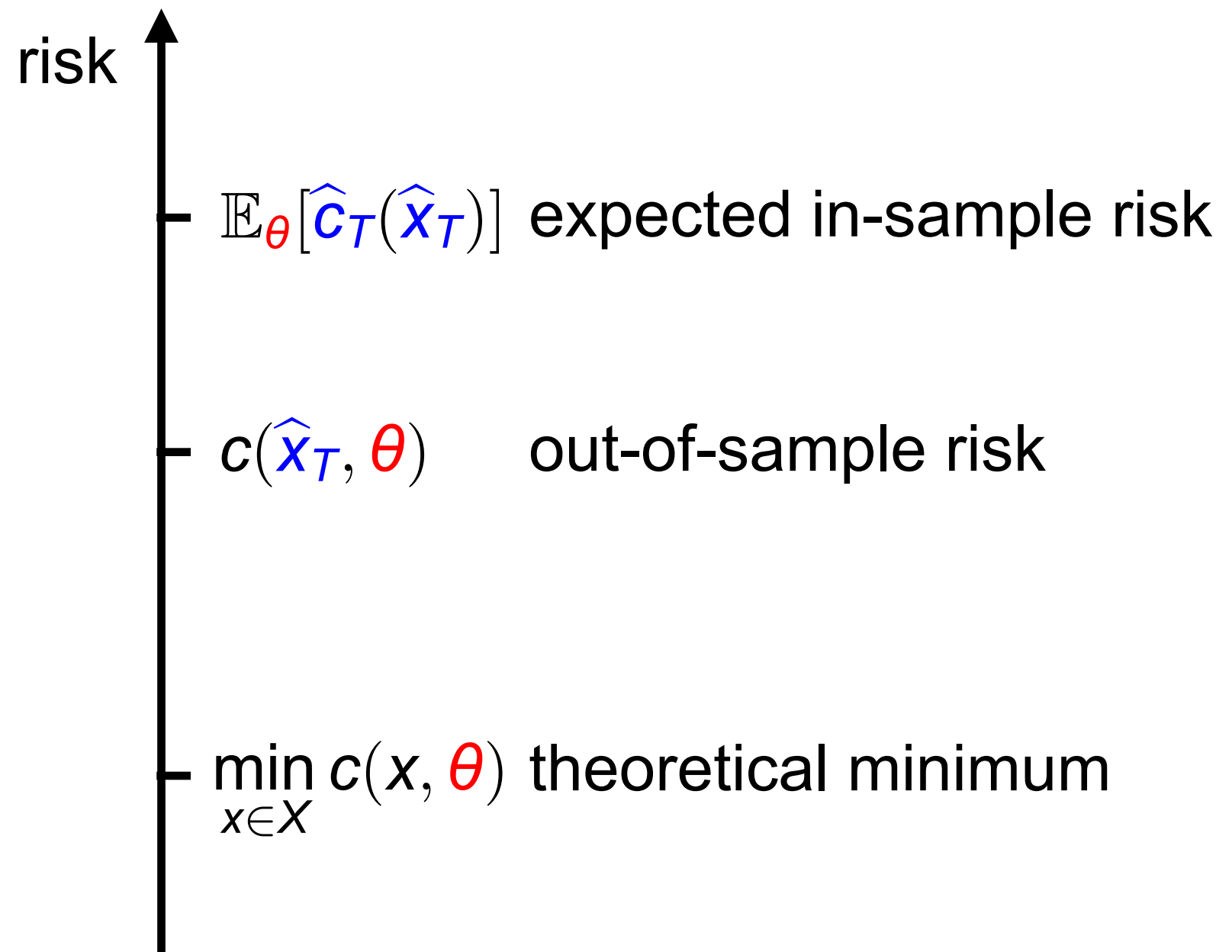
# Optimal Data-Driven Optimization



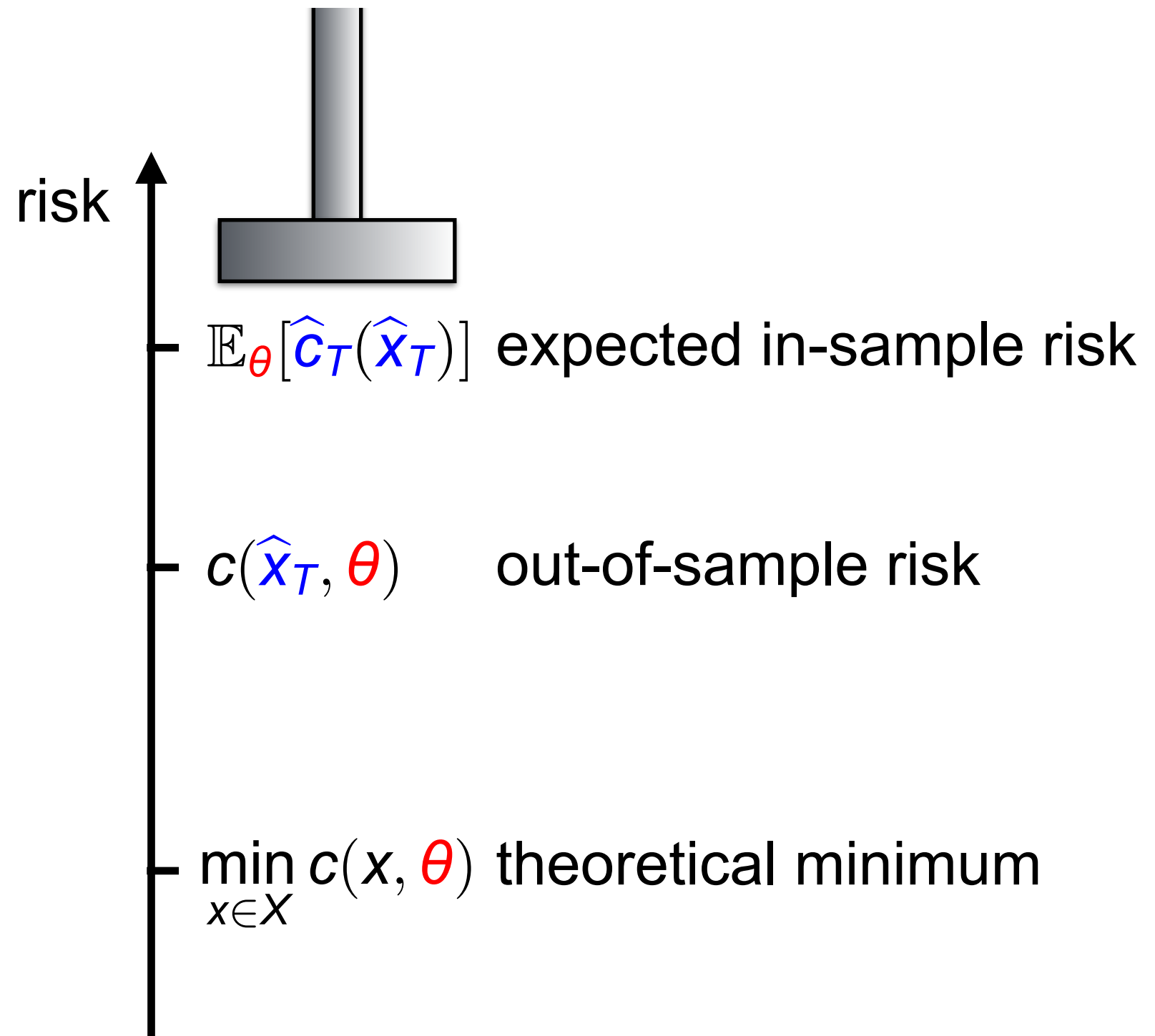
# Optimal Data-Driven Optimization



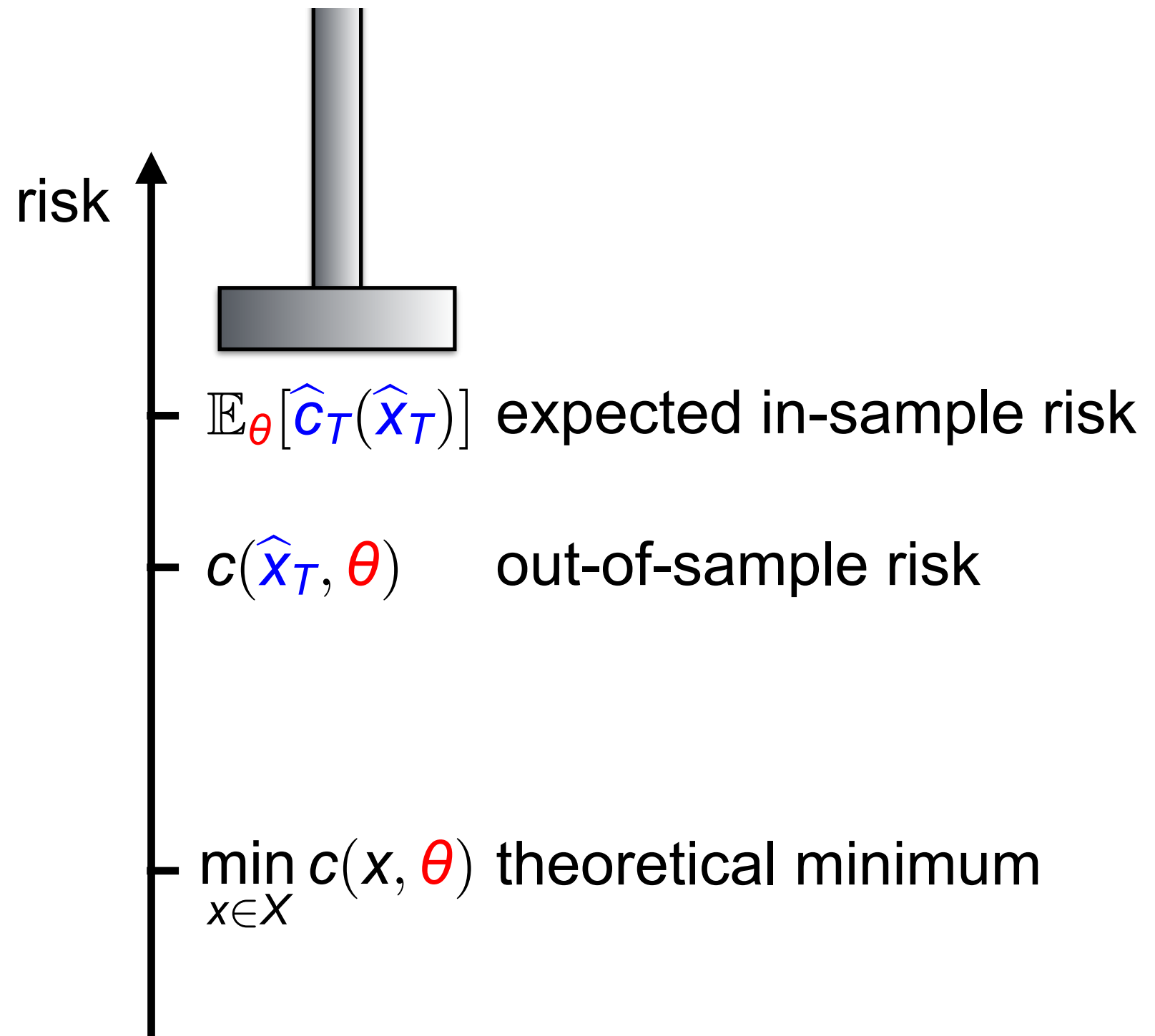
# Optimal Data-Driven Optimization



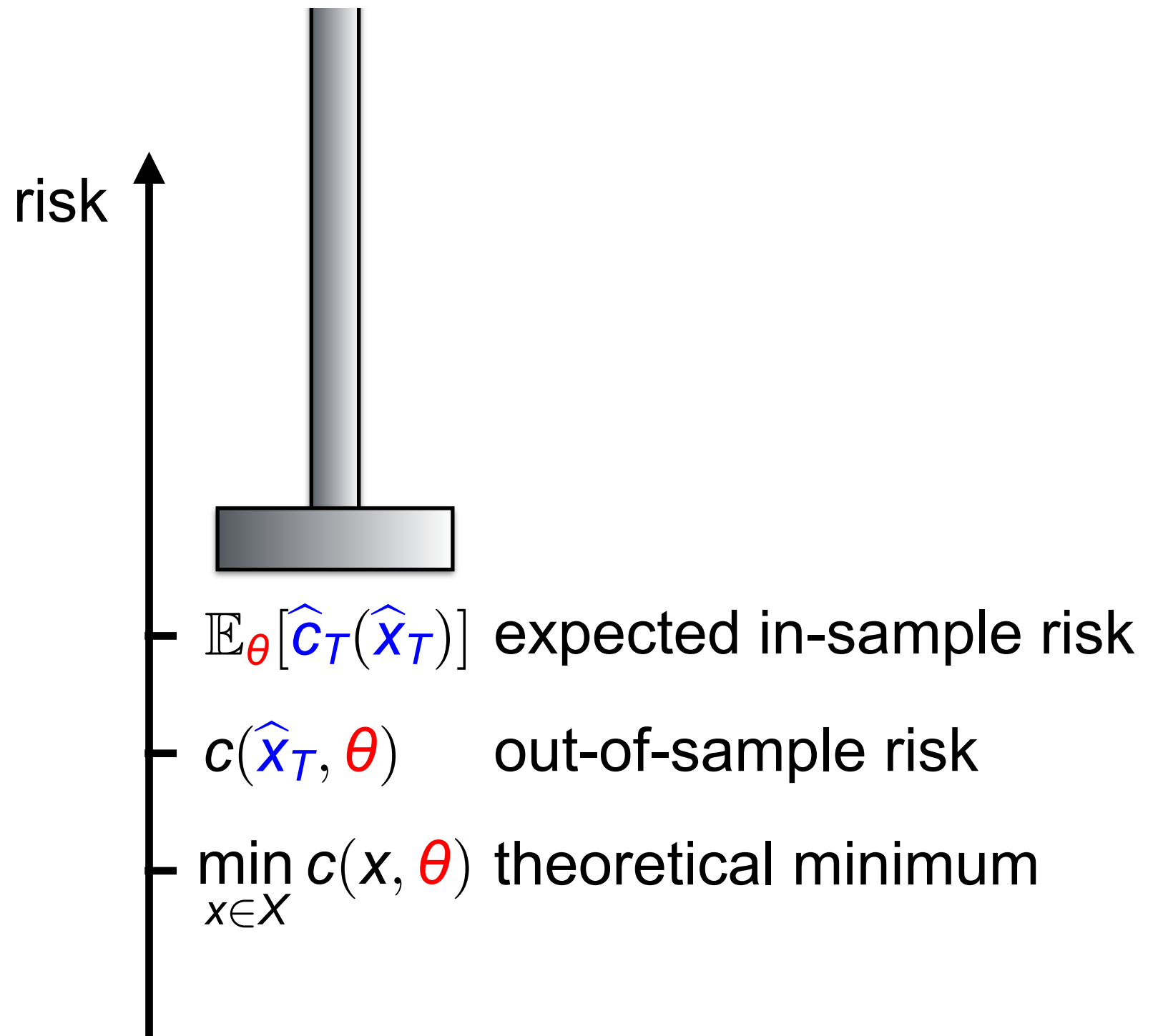
# Optimal Data-Driven Optimization



# Optimal Data-Driven Optimization

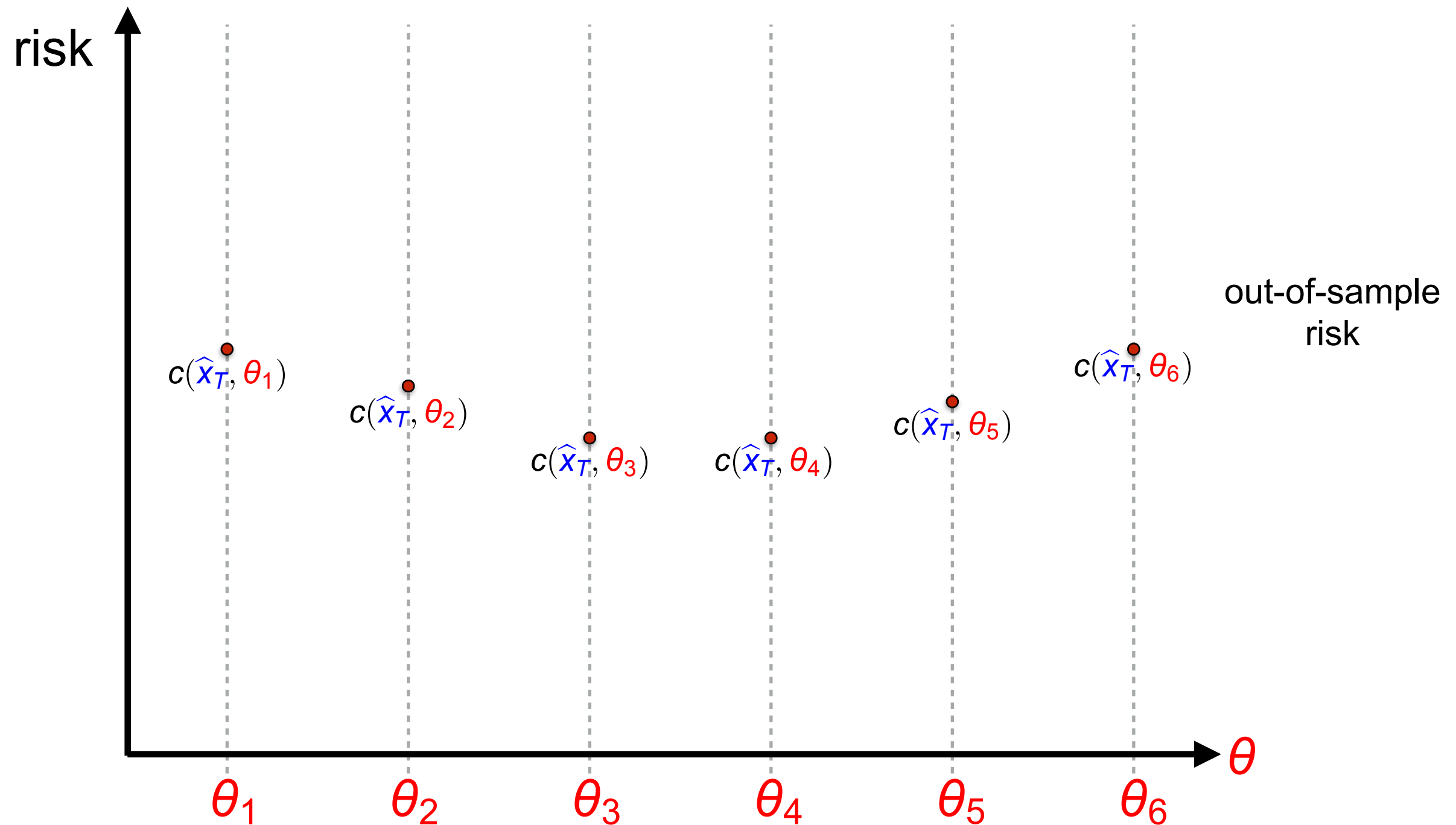


# Optimal Data-Driven Optimization

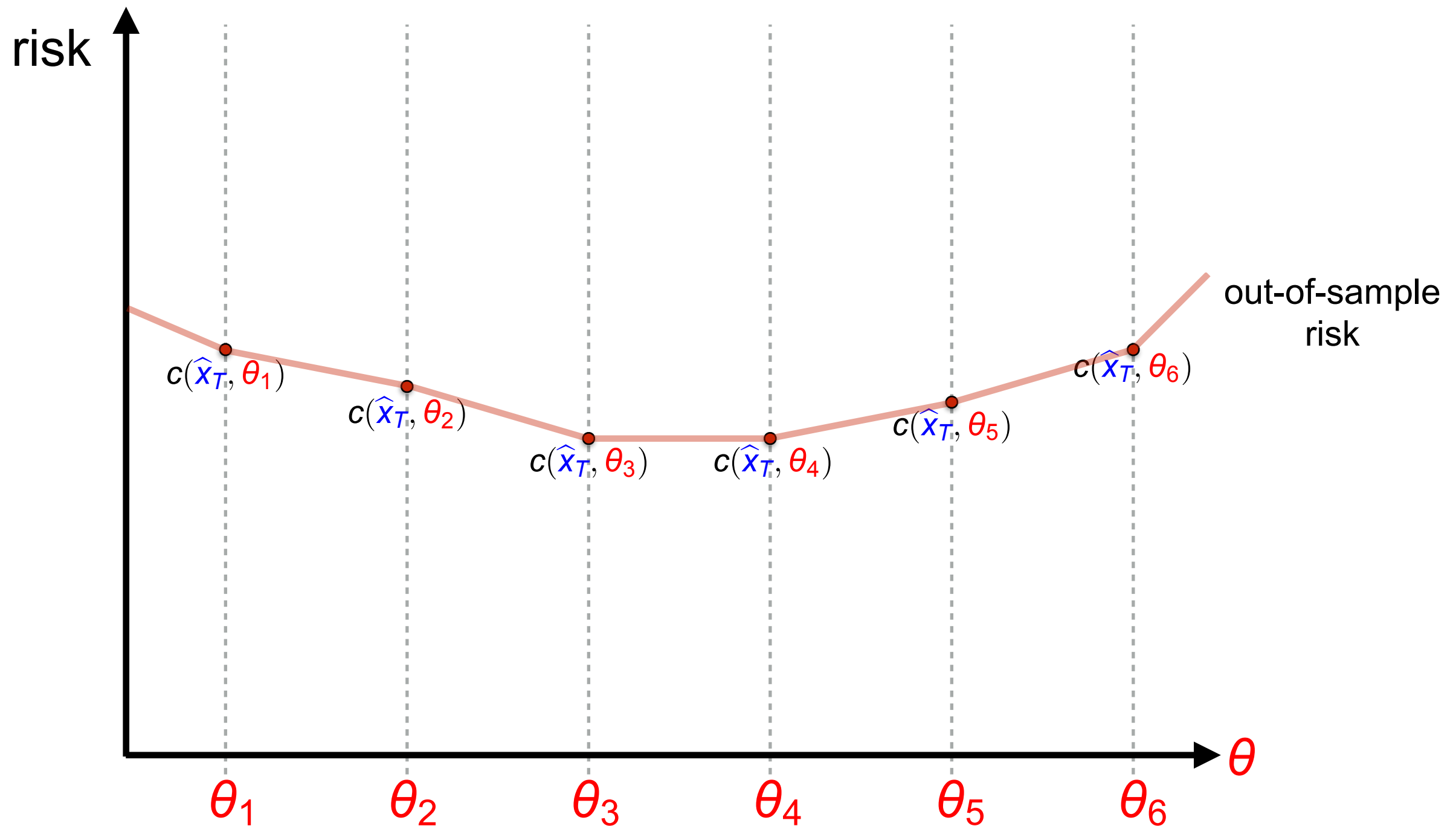




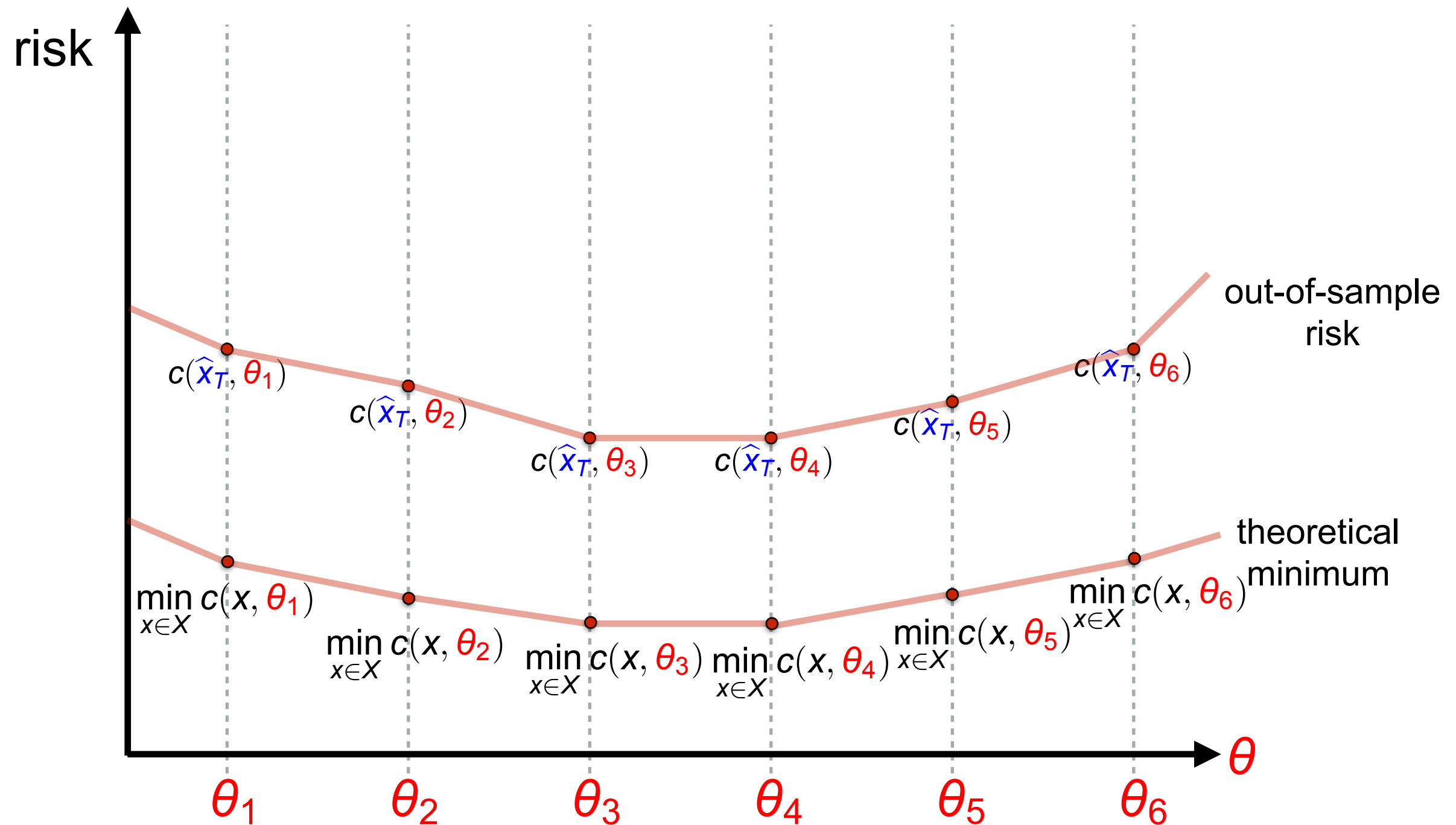
# Optimal Data-Driven Optimization



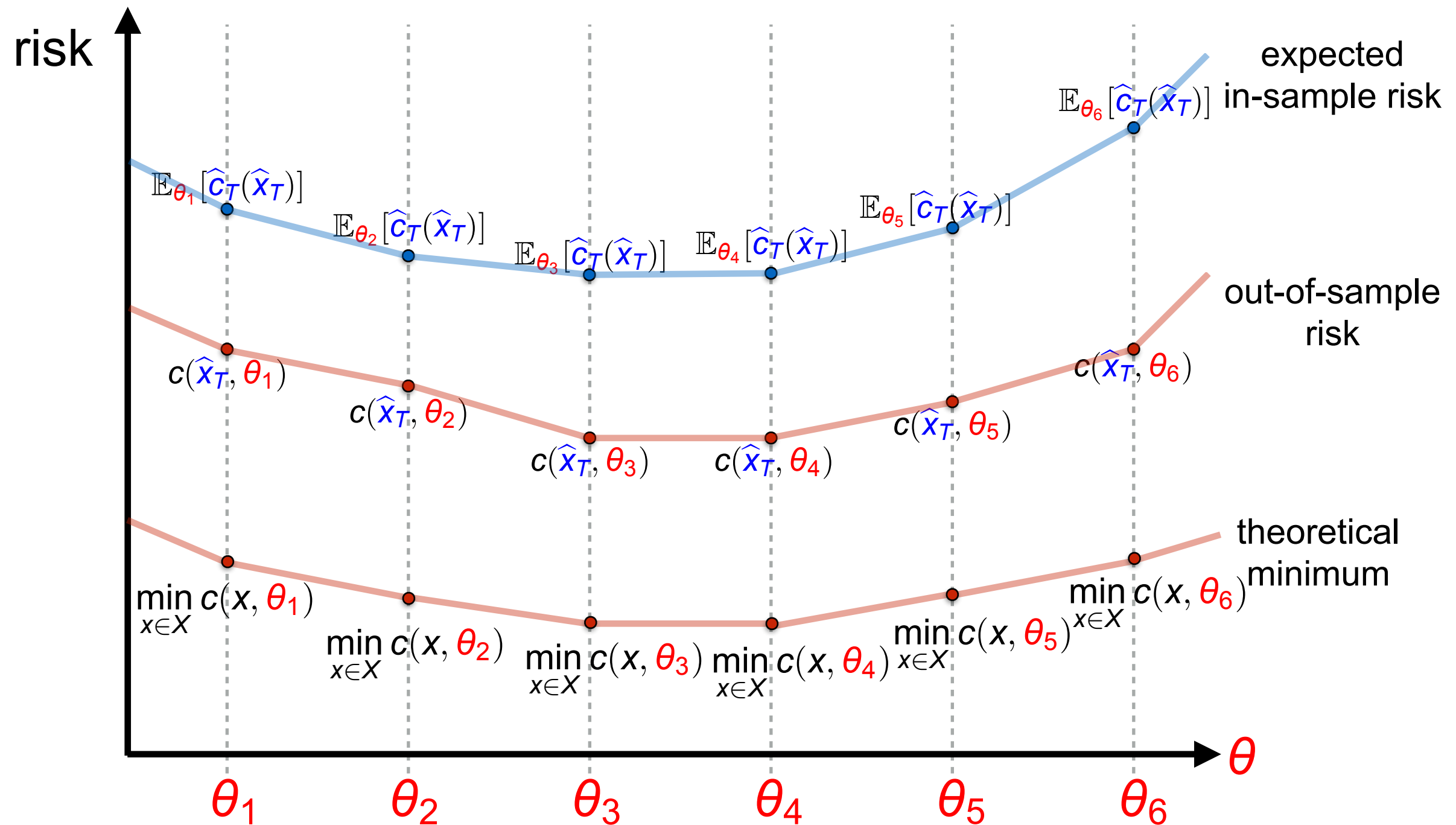
# Optimal Data-Driven Optimization



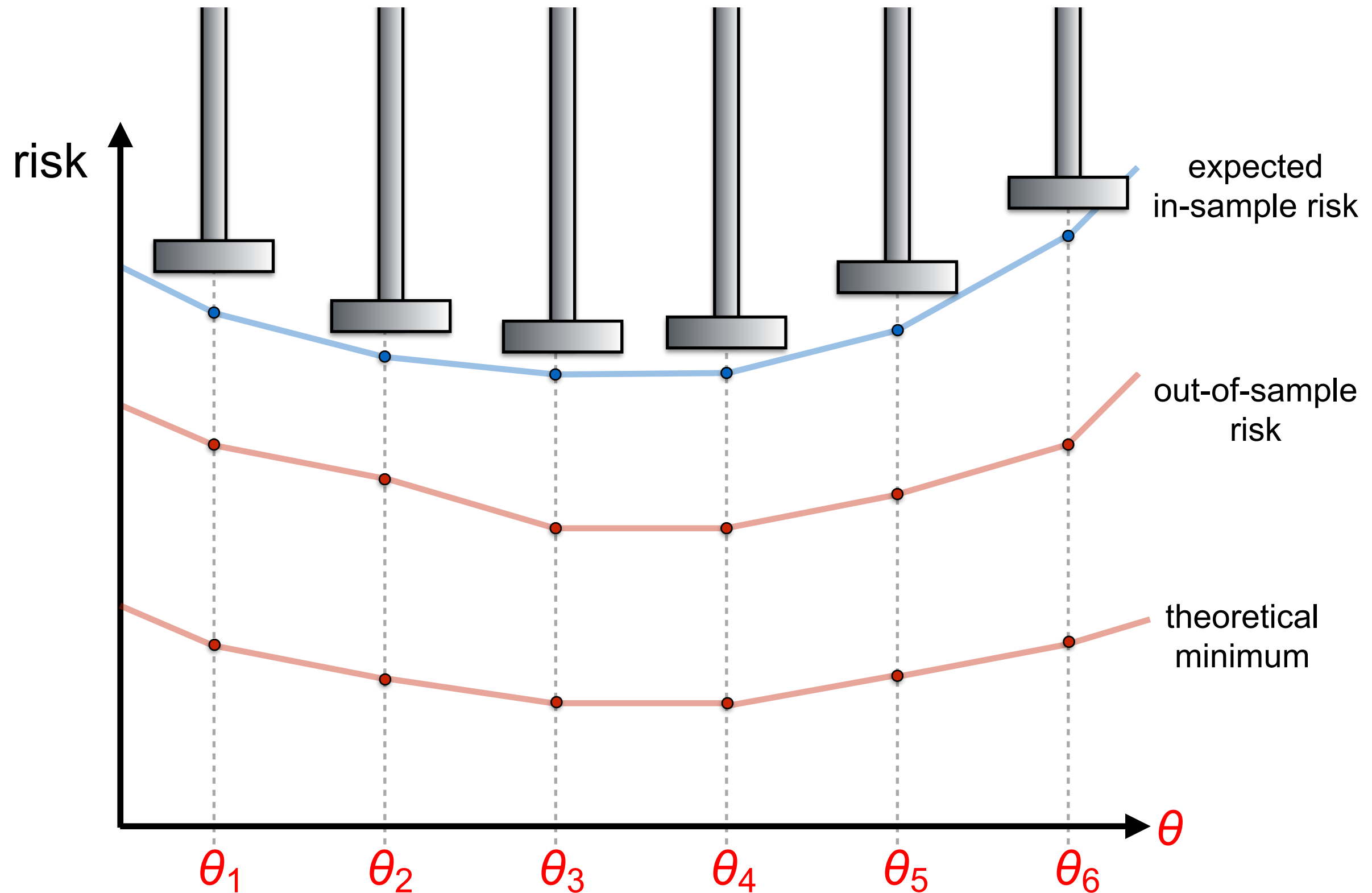
# Optimal Data-Driven Optimization



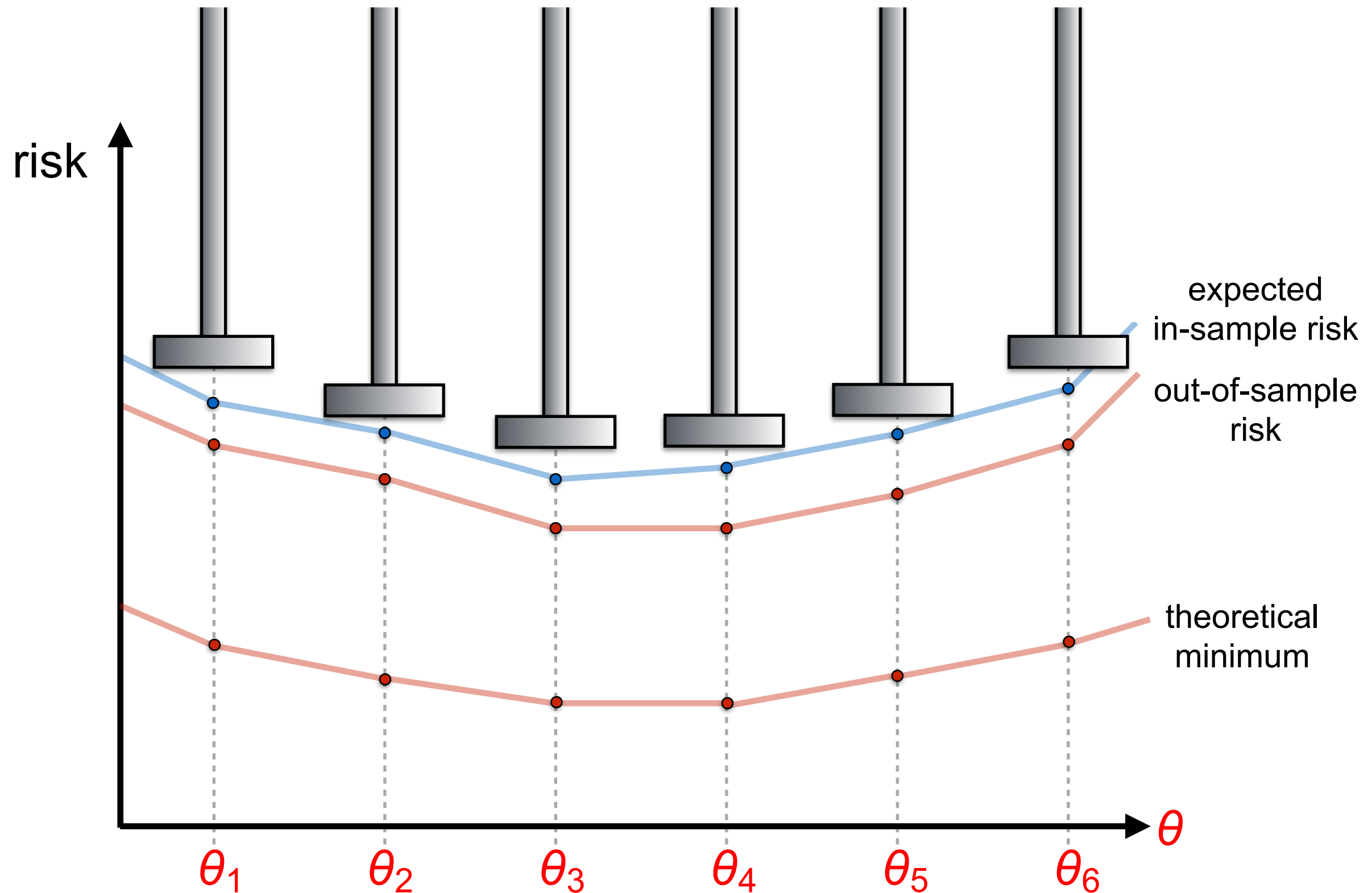
# Optimal Data-Driven Optimization



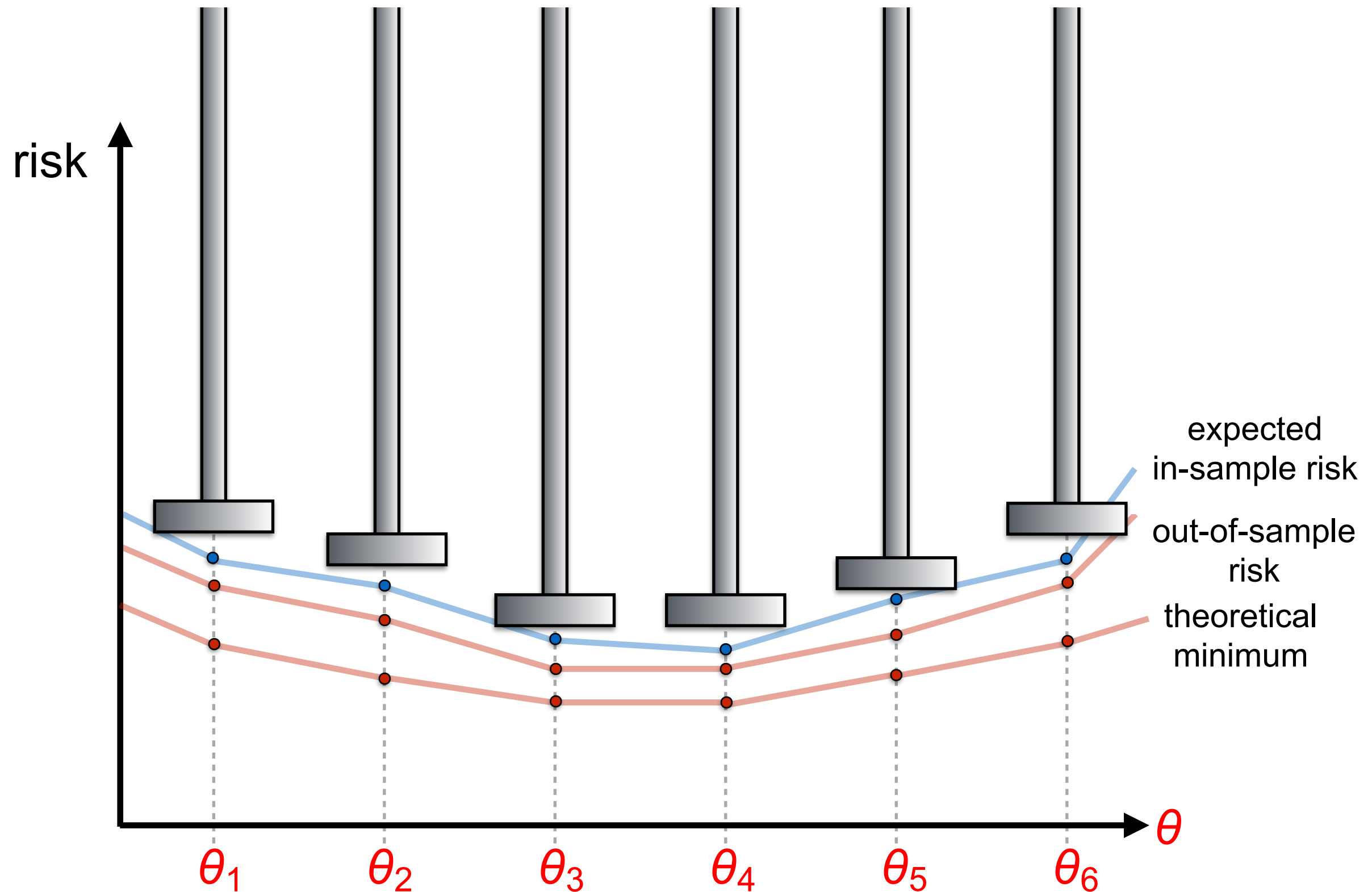
# Optimal Data-Driven Optimization



# Optimal Data-Driven Optimization



# Optimal Data-Driven Optimization



# Optimal Data-Driven Optimization

Minimize the in-sample risk and require that the out-of-sample disappointment decays exponentially at rate  $r$

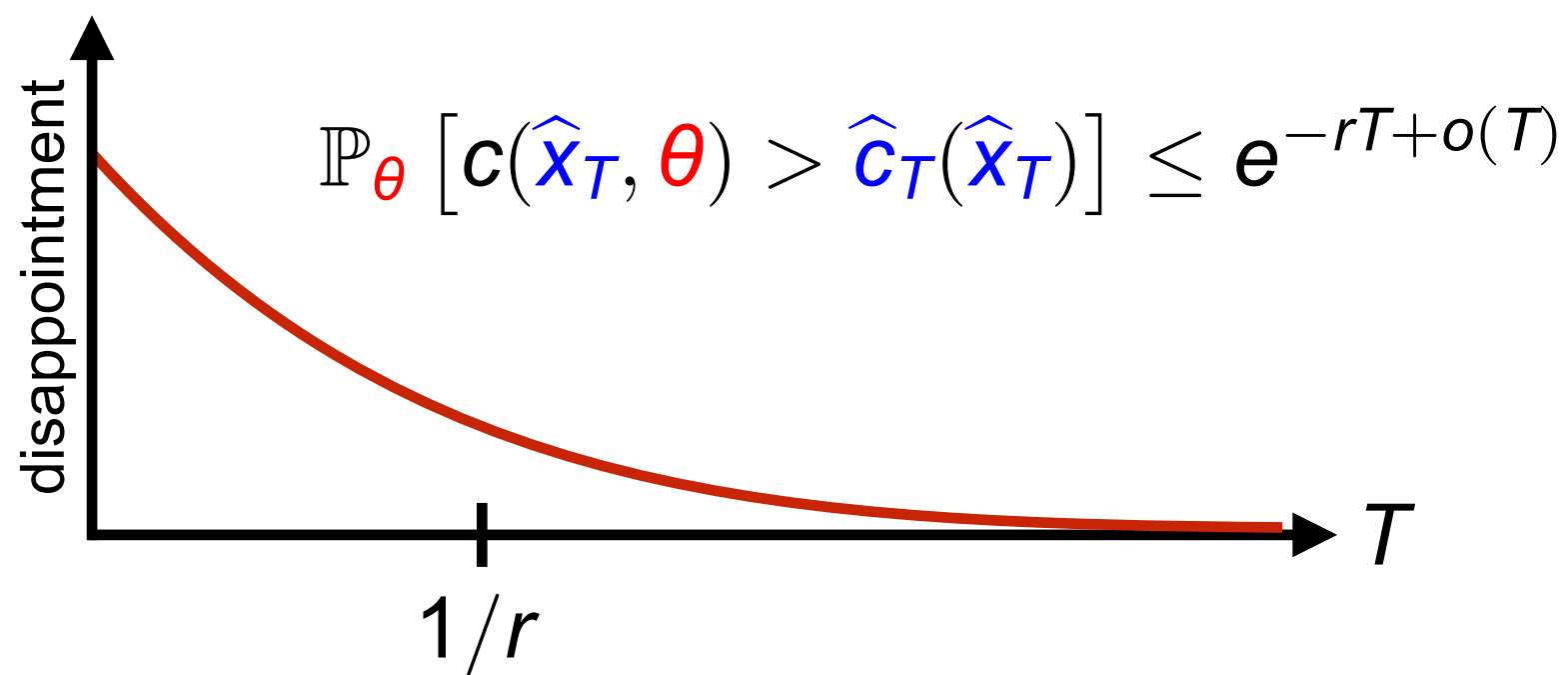
$$\begin{aligned} & \underset{\hat{c}_T, \hat{x}_T}{\text{minimize}} \quad \left\{ \lim_{T \rightarrow \infty} \mathbb{E}_{\theta} [\hat{c}_T(\hat{x}_T)] \right\}_{\theta \in \Theta} \\ & \text{subject to} \quad \limsup_{T \rightarrow \infty} \frac{1}{T} \log \mathbb{P}_{\theta} [c(\hat{x}_T, \theta) > \hat{c}_T(\hat{x}_T)] \leq -r \quad \forall \theta \in \Theta \end{aligned}$$



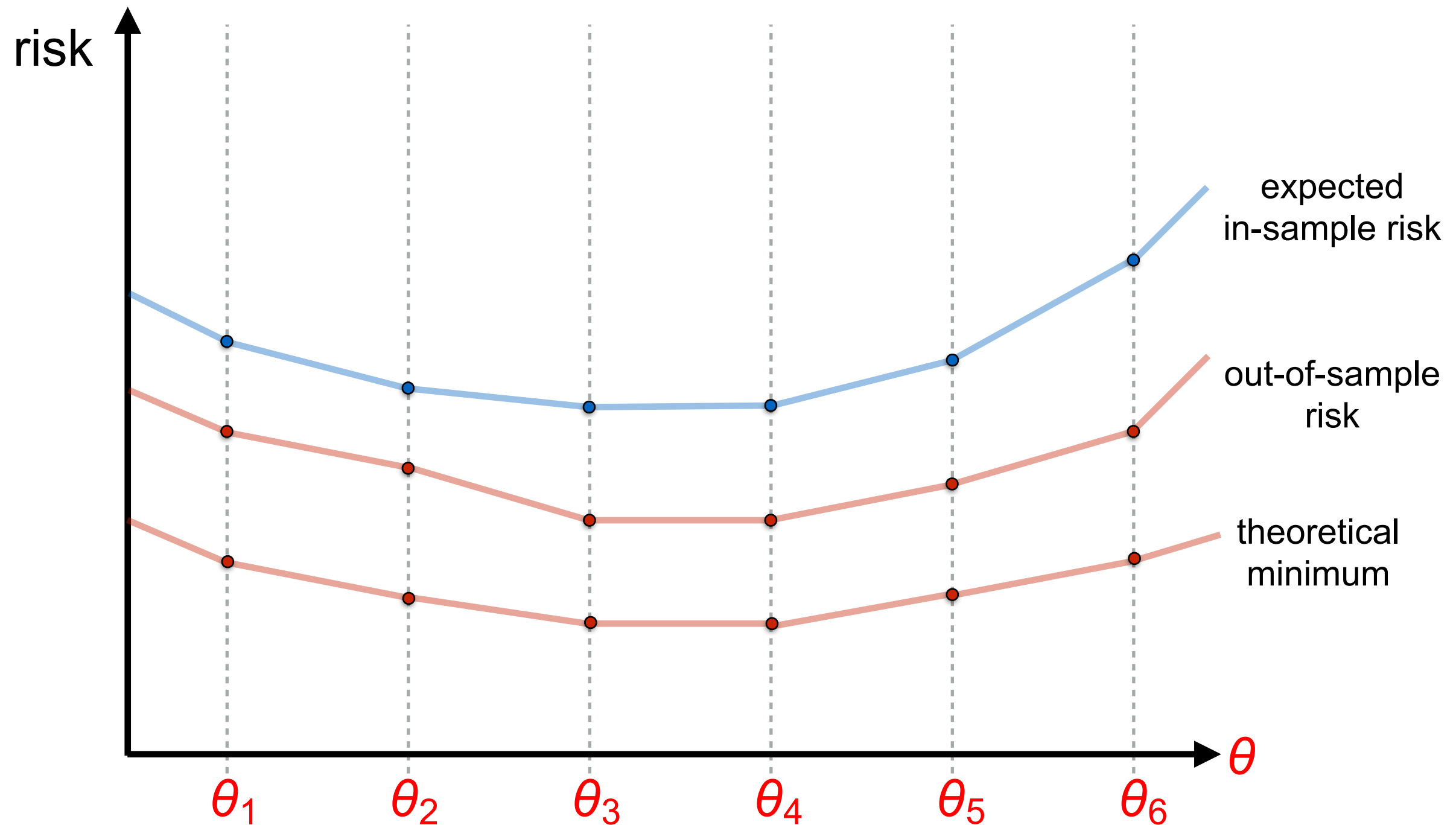
# Optimal Data-Driven Optimization

Minimize the in-sample risk and require that the out-of-sample disappointment decays exponentially at rate  $r$

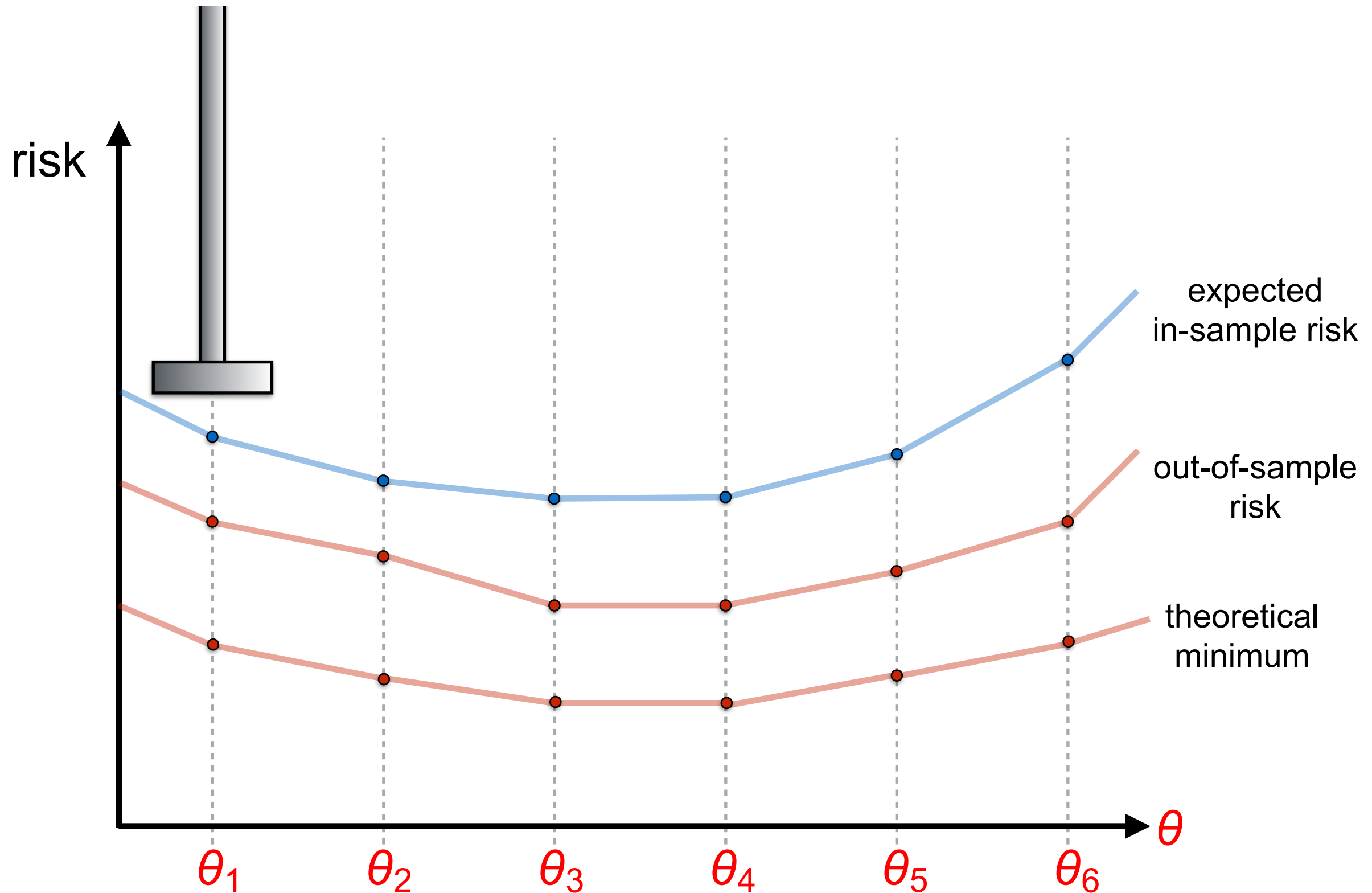
$$\begin{aligned} & \underset{\hat{c}_T, \hat{x}_T}{\text{minimize}} \quad \left\{ \lim_{T \rightarrow \infty} \mathbb{E}_{\theta} [\hat{c}_T(\hat{x}_T)] \right\}_{\theta \in \Theta} \\ & \text{subject to} \quad \limsup_{T \rightarrow \infty} \frac{1}{T} \log \mathbb{P}_{\theta} [c(\hat{x}_T, \theta) > \hat{c}_T(\hat{x}_T)] \leq -r \quad \forall \theta \in \Theta \end{aligned}$$



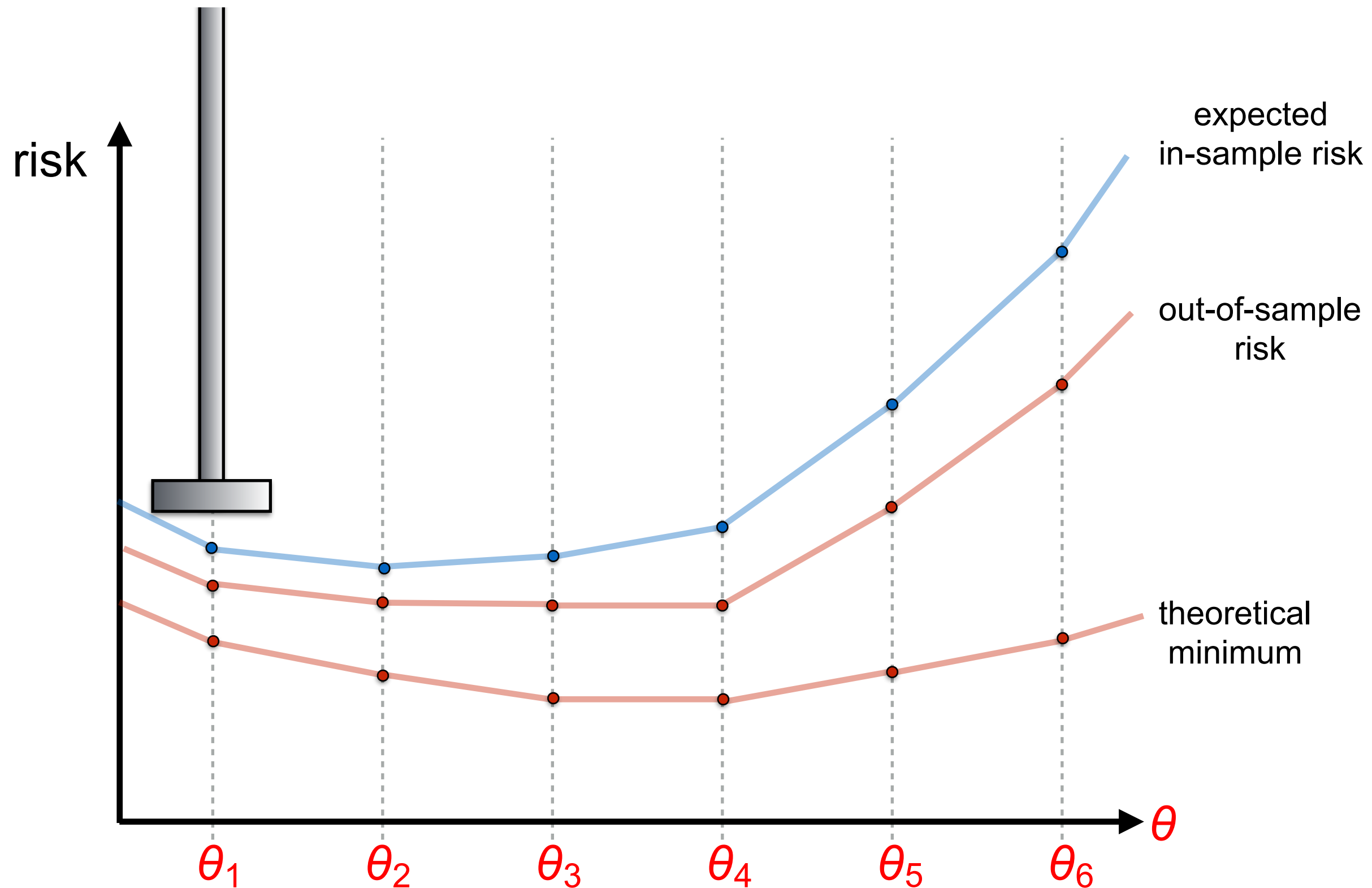
# Optimal Data-Driven Optimization



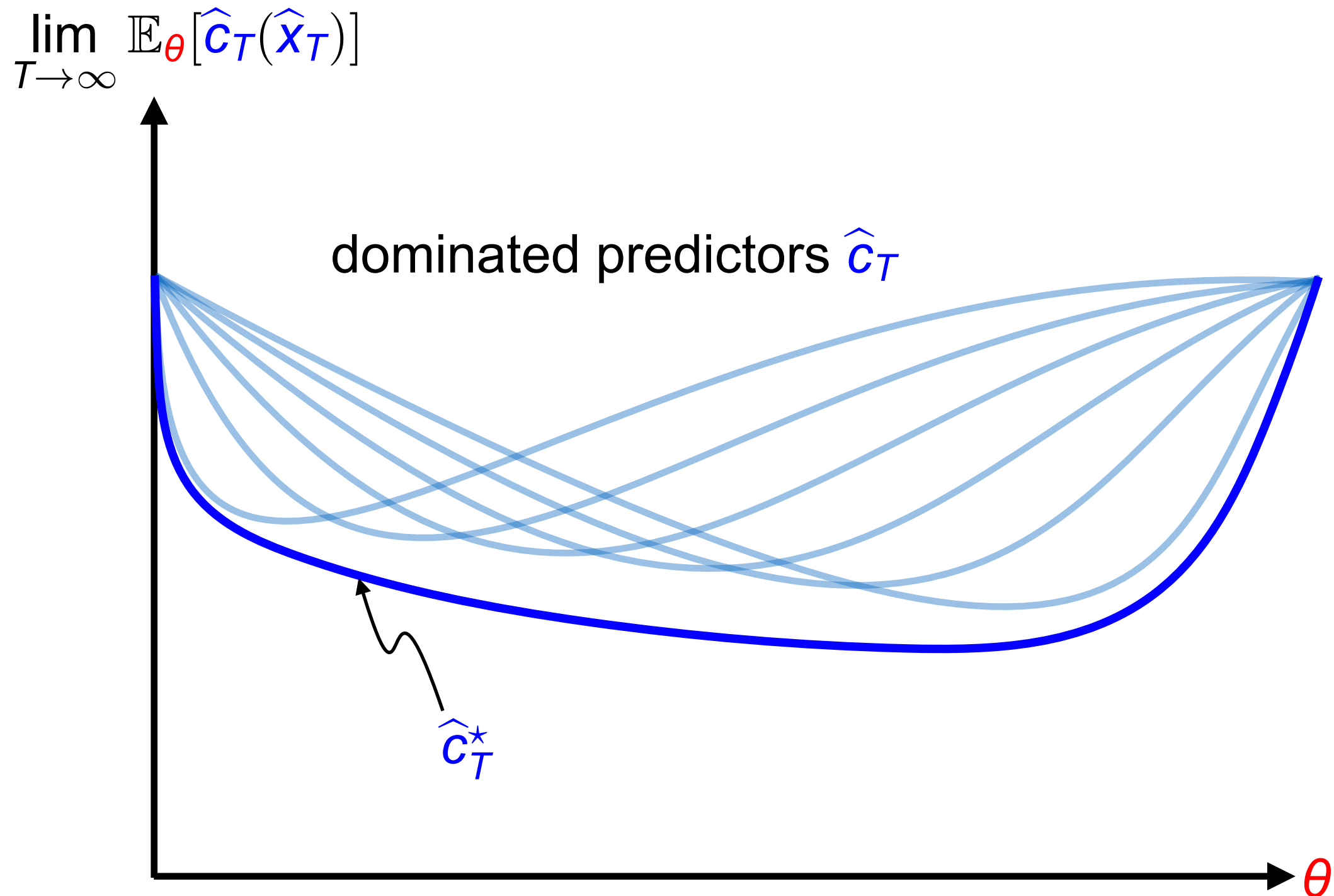
# Optimal Data-Driven Optimization



# Optimal Data-Driven Optimization



# Pareto-Dominant Solutions



$\hat{c}_T^*$  minimizes the in-sample risk simultaneously for every  $\theta$

# Meta-Optimization Problem (MOP)

MOP optimizes over all surrogate optimization models

$$\begin{aligned} & \underset{\hat{c}_T, \hat{x}_T}{\text{minimize}} \quad \left\{ \lim_{T \rightarrow \infty} \mathbb{E}_{\theta} [\hat{c}_T(\hat{x}_T)] \right\}_{\theta \in \Theta} \\ & \text{subject to} \quad \limsup_{T \rightarrow \infty} \frac{1}{T} \log \mathbb{P}_{\theta} [c(\hat{x}_T, \theta) > \hat{c}_T(\hat{x}_T)] \leq -r \quad \forall \theta \in \Theta \end{aligned}$$

## Strengths:

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

# Meta-Optimization Problem (MOP)

MOP optimizes over all surrogate optimization models

$$\begin{aligned} & \underset{\hat{c}_T, \hat{x}_T}{\text{minimize}} \quad \left\{ \lim_{T \rightarrow \infty} \mathbb{E}_{\theta} [\hat{c}_T(\hat{x}_T)] \right\}_{\theta \in \Theta} \\ & \text{subject to} \quad \limsup_{T \rightarrow \infty} \frac{1}{T} \log \mathbb{P}_{\theta} [c(\hat{x}_T, \theta) > \hat{c}_T(\hat{x}_T)] \leq -r \quad \forall \theta \in \Theta \end{aligned}$$

## Weaknesses:

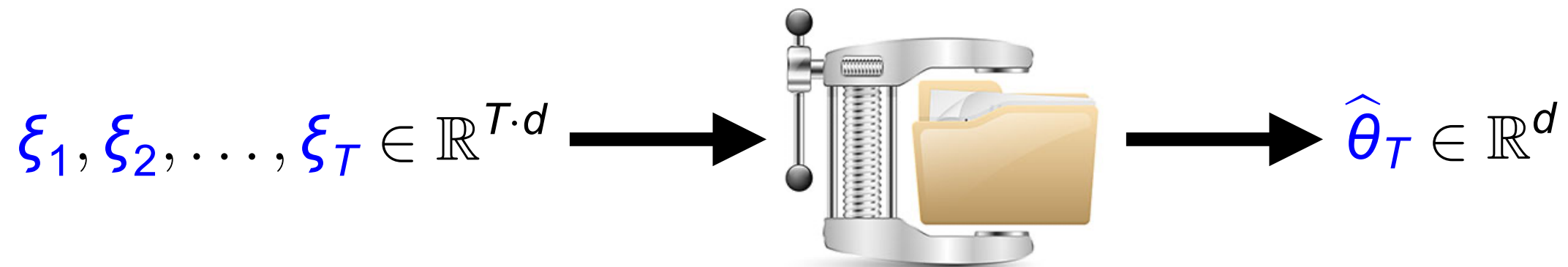
- ▶ performance criteria are asymptotic
- ▶ choice of  $r$  is subjective
- ▶ why insist on exponential decay?
- ▶ feasible/optimal models are biased

# Restricted Meta-Optimization Problems



# Data Compression

Compress the raw data to an estimator of  $\theta$



$\implies$  compressed predictors depend on  $\xi_1, \xi_2, \dots, \xi_T$  and on  $T$  only indirectly through the summary statistic  $\hat{\theta}_T$

# Compressed Predictors and Prescriptors

- ▶ Set  $\hat{c}_T(x) = \tilde{c}(x, \hat{\theta}_T)$  for some continuous function  $\tilde{c}$
- ▶ Set  $\hat{x}_T = \tilde{x}(\hat{\theta}_T)$  for some quasi-continuous function  $\tilde{x}$  with

$$\tilde{x}(\hat{\theta}_T) \in \operatorname{argmin}_{x \in X} \tilde{c}(x, \hat{\theta}_T)$$

# Restricted MOP

Restricted MOP over compressed predictors/prescriptors:

$$\begin{aligned} & \underset{\tilde{c}, \tilde{x}}{\text{minimize}} && \{ \tilde{c}(\tilde{x}(\theta), \theta) \}_{\theta \in \Theta} \\ & \text{subject to} && \limsup_{T \rightarrow \infty} \frac{1}{T} \log \mathbb{P}_{\theta} \left[ c(\tilde{x}(\hat{\theta}_T), \theta) > \tilde{c}(\tilde{x}(\hat{\theta}_T), \hat{\theta}_T) \right] \leq -r \quad \forall \theta \in \Theta \end{aligned}$$

# Large Deviation Principle (LDP)

**Definition:**<sup>1)</sup> The estimators  $\hat{\theta}_T$ ,  $T \in \mathbb{N}$ , satisfy an LDP if there exists a rate function  $I(\theta', \theta)$  such that for all Borel sets  $\mathcal{D} \subseteq \Theta$

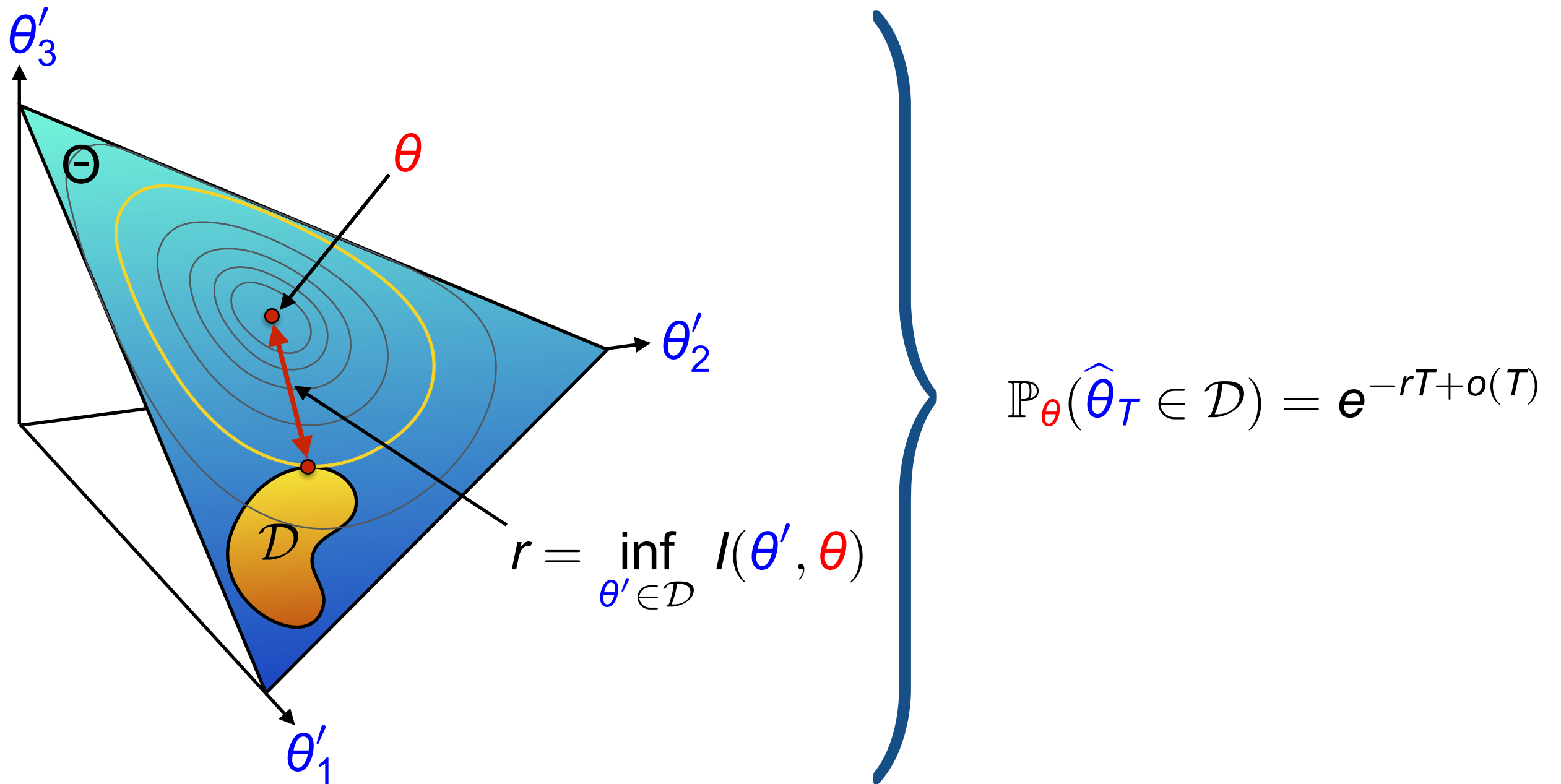
$$\limsup_{T \rightarrow \infty} \frac{1}{T} \log \mathbb{P}_{\theta}(\hat{\theta}_T \in \mathcal{D}) \leq - \inf_{\theta' \in \text{cl } \mathcal{D}} I(\theta', \theta)$$

$$\liminf_{T \rightarrow \infty} \frac{1}{T} \log \mathbb{P}_{\theta}(\hat{\theta}_T \in \mathcal{D}) \geq - \inf_{\theta' \in \text{int } \mathcal{D}} I(\theta', \theta)$$

---

<sup>1)</sup> den Hollander, *American Mathematical Society*, 2008; Dembo & Zeitouni, *Springer*, 2009.

# Large Deviation Principle (LDP)



# DRO is Optimal

## Assumption:

- ▶  $\hat{\theta}_T$  satisfies an LDP with a “regular” rate function

**Theorem 1** (DRO is optimal): The following distributionally robust compressed predictor is a Pareto-dominant solution for the restricted MOP.

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

# DRO is Optimal

## Assumption:

- ▶  $\hat{\theta}_T$  satisfies an LDP with a “regular” rate function

**Theorem 1** (DRO is optimal): The following distributionally robust compressed predictor is a Pareto-dominant solution for the restricted MOP.

$$\tilde{c}^*(x, \hat{\theta}_T) = \begin{cases} \sup_{\theta \in \Theta} c(x, \theta) \\ \text{s.t. } l(\hat{\theta}_T, \theta) \leq r \end{cases}$$

## Note:

- ▶ The shape of the ambiguity set is determined by  $\hat{\theta}_T$
- ▶ The “radius” of the ambiguity set is given by the decay rate  $r$

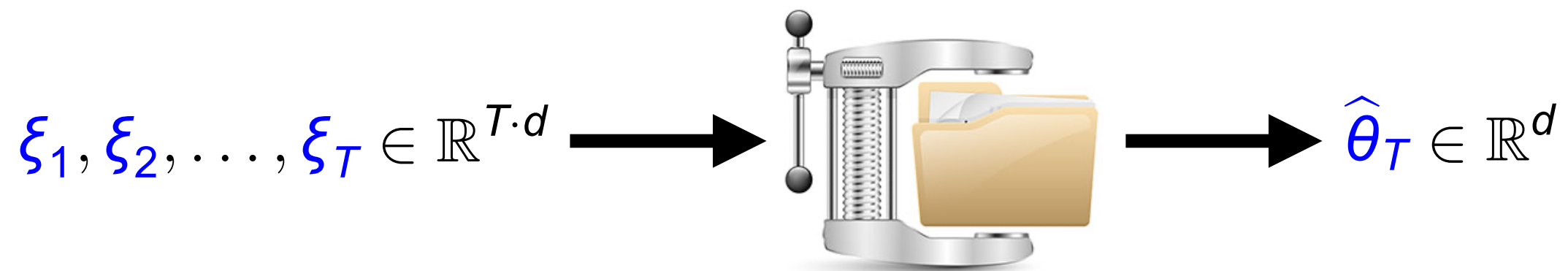
# Separation of Estimation and Optimization



# Sufficient Statistic

**Definition:**  $\hat{\theta}_T$  is a sufficient statistic for  $\theta$  if the distribution of  $\xi_1, \xi_2, \dots, \xi_T$  conditional on  $\hat{\theta}_T = \theta'$  is independent of  $\theta \in \Theta$ .

$\Rightarrow$  Lossless compression



# DRO is Optimal

## Assumptions:

- ▶  $\hat{\theta}_T$  is a sufficient statistic for  $\theta$
- ▶  $\hat{\theta}_T$  satisfies an LDP with a “regular” rate function

**Theorem 2** (DRO is optimal): The following distributionally robust surrogate optimization model is a Pareto-dominant solution for the original MOP.

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

# DRO is Optimal

## Assumptions:

- ▶  $\hat{\theta}_T$  is a sufficient statistic for  $\theta$
- ▶  $\hat{\theta}_T$  satisfies an LDP

Separation of estimation and optimization:

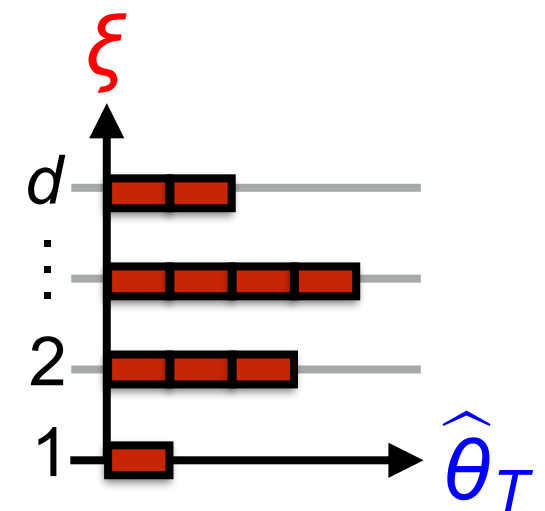
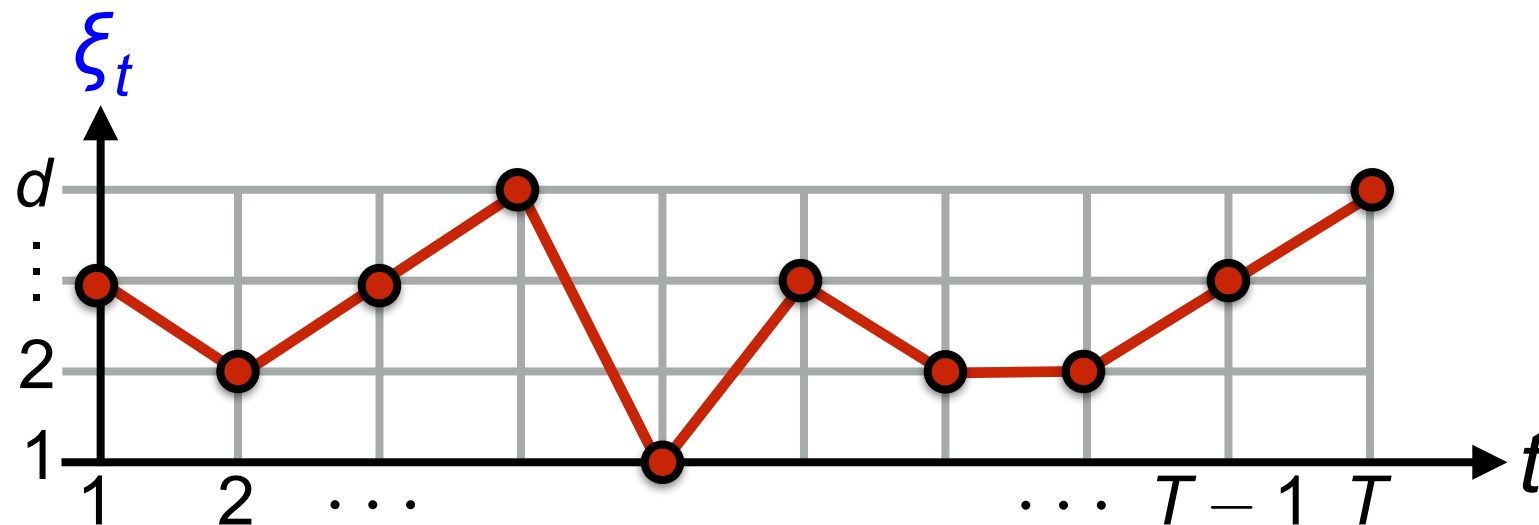
- 1) Evaluate the estimator
- 2) Solve the DRO problem

# Data-Generating Processes

# Newsvendor Problem Revisited

“Compressing” the raw data:

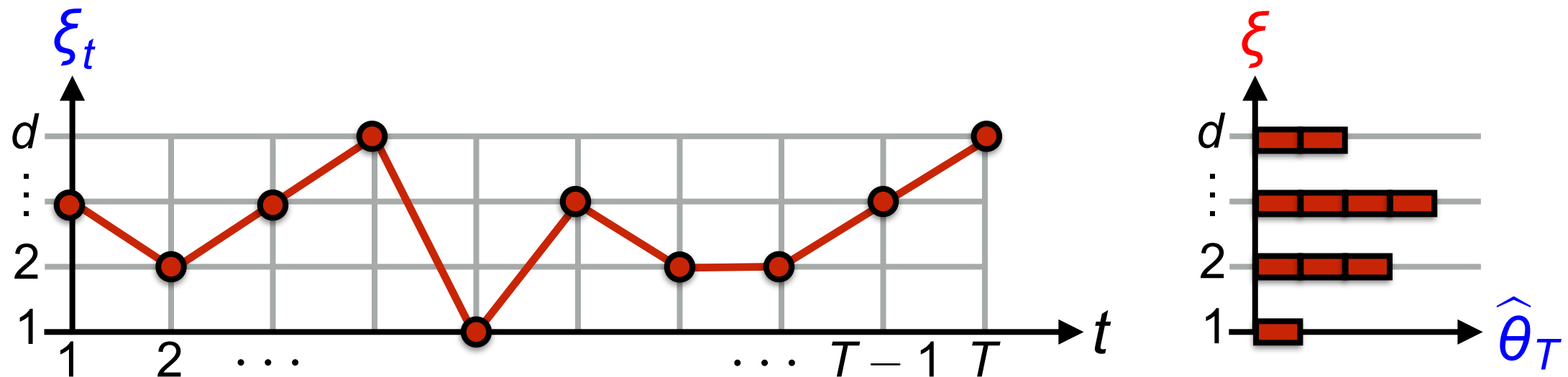
$$\left. \begin{array}{l} \text{demand observations} \\ (\xi_1, \xi_2, \dots, \xi_T) \end{array} \right\} \longleftrightarrow \left\{ \begin{array}{l} \text{empirical distribution} \\ (\hat{\theta}_T)_i = \frac{1}{T} \sum_{t=1}^T \mathbf{1}_{\xi_t=i} \end{array} \right.$$



# Newsvendor Problem Revisited

“Compressing” the raw data:

$$\left. \begin{array}{l} \text{demand observations} \\ (\xi_1, \xi_2, \dots, \xi_T) \end{array} \right\} \longleftrightarrow \left\{ \begin{array}{l} \text{empirical distribution} \\ (\hat{\theta}_T)_i = \frac{1}{T} \sum_{t=1}^T \mathbf{1}_{\xi_t=i} \end{array} \right.$$



- Fisher-Neyman:<sup>1)</sup>  $\hat{\theta}_T$  is a sufficient statistic for  $\theta$
- Sanov:<sup>2)</sup>  $\hat{\theta}_T$  satisfies an LDP with  $I(\hat{\theta}_T, \theta) = D_{\text{KL}}(\hat{\theta}_T \parallel \theta)$

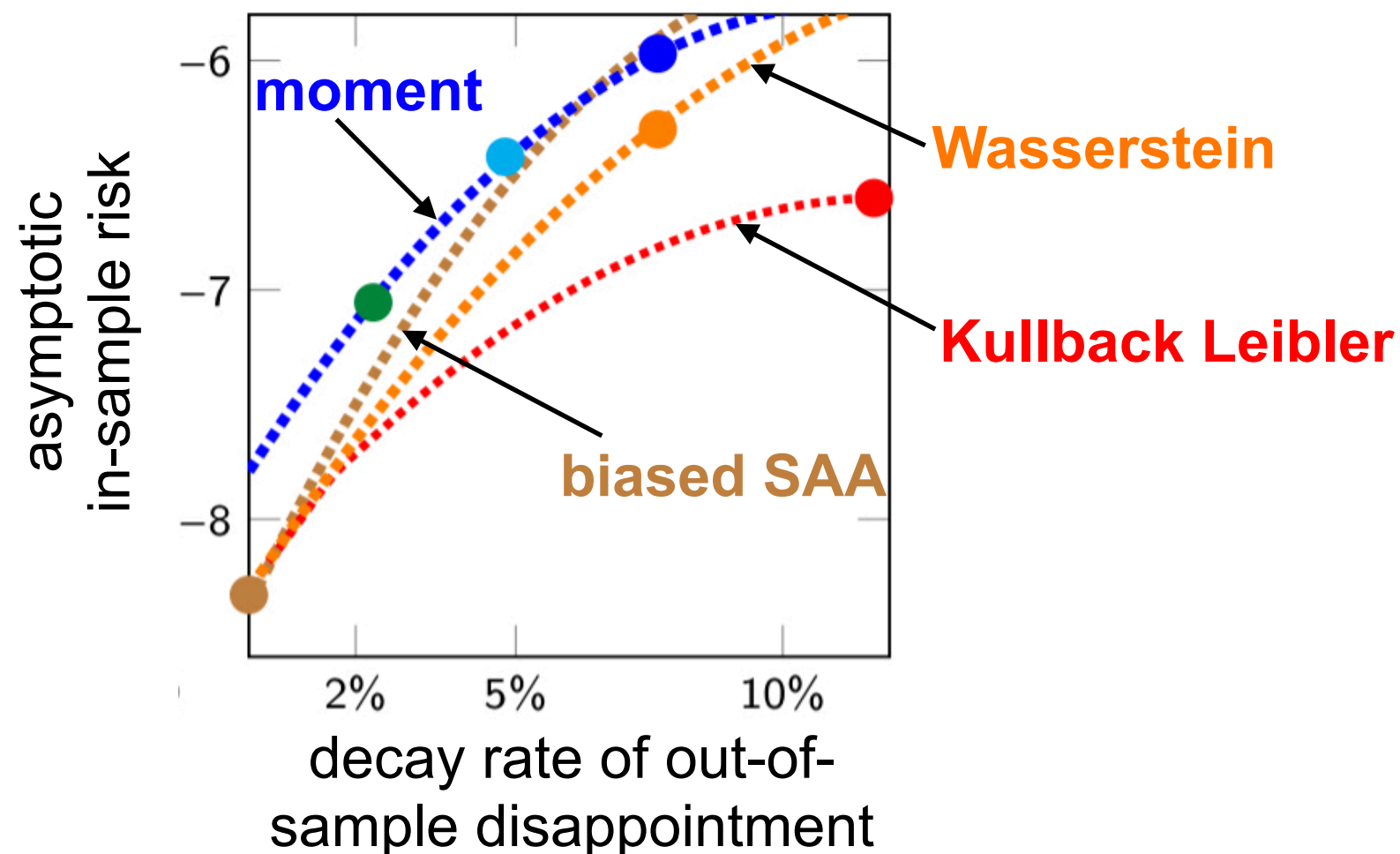
<sup>1)</sup> Lehmann & Casella, *Springer*, 1998;

<sup>2)</sup> Sanov, *Matematicheskii Sbornik*, 1957.

# Newsvendor Problem Revisited

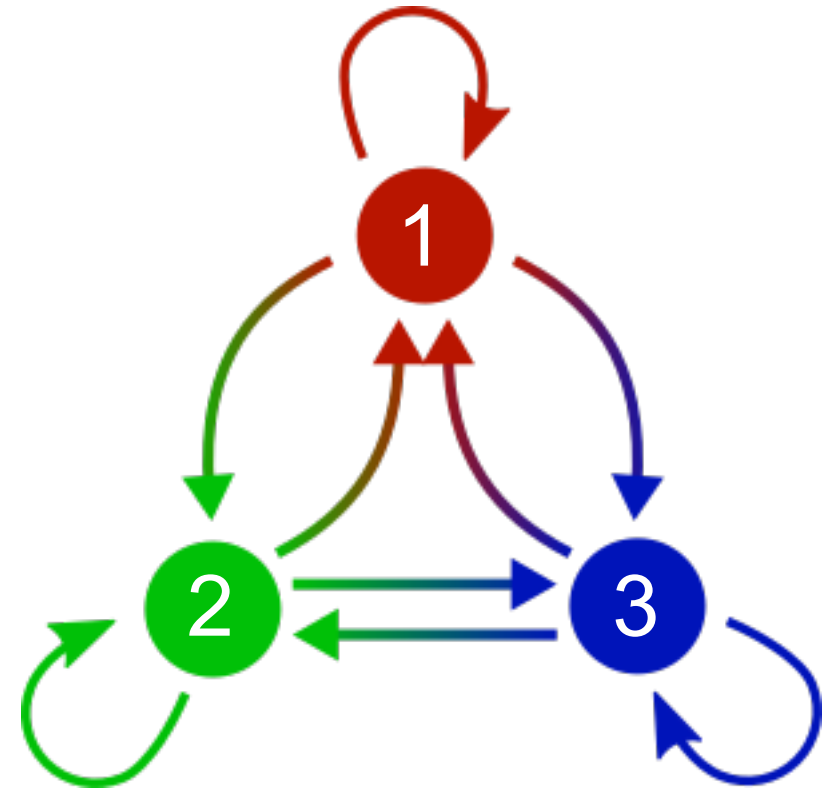
- ▶ The separation principle holds
- ▶ The optimal data-driven predictor is

$$\hat{c}_T(x) = \begin{cases} \sup_{\theta \in \Theta} c(x, \theta) \\ \text{s.t. } D_{\text{KL}}(\hat{\theta}_T \| \theta) \leq r \end{cases}$$



# Finite-State Markov Chains

Assume that  $\{\xi_t\}_{t \in \mathbb{N}}$  is a Markov chain on  $\{1, 2, \dots, d\}$

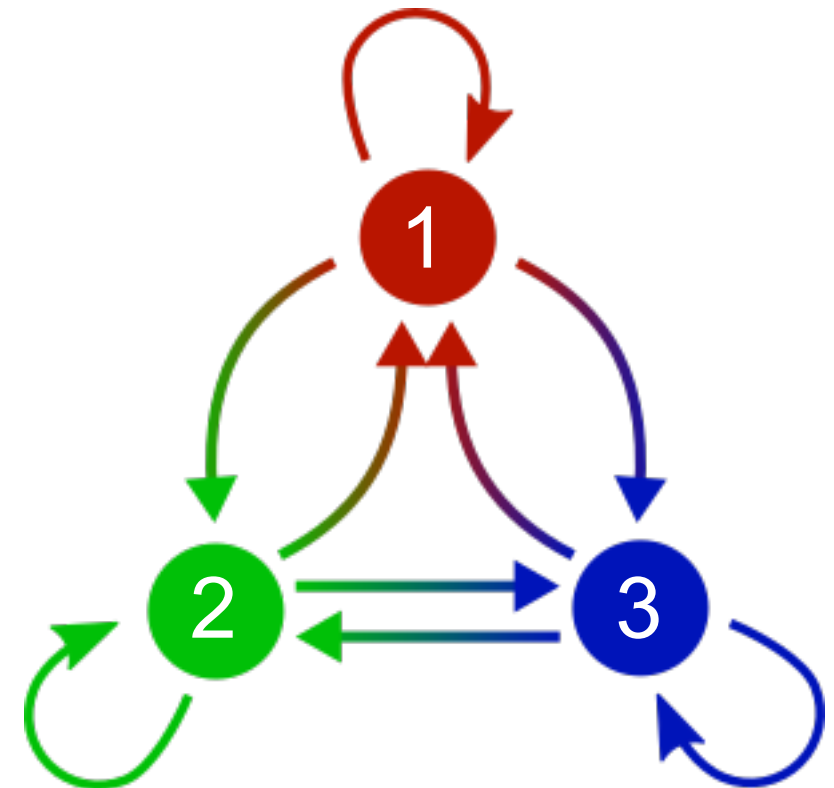




# Finite-State Markov Chains

Assume that  $\{\xi_t\}_{t \in \mathbb{N}}$  is a Markov chain on  $\{1, 2, \dots, d\}$

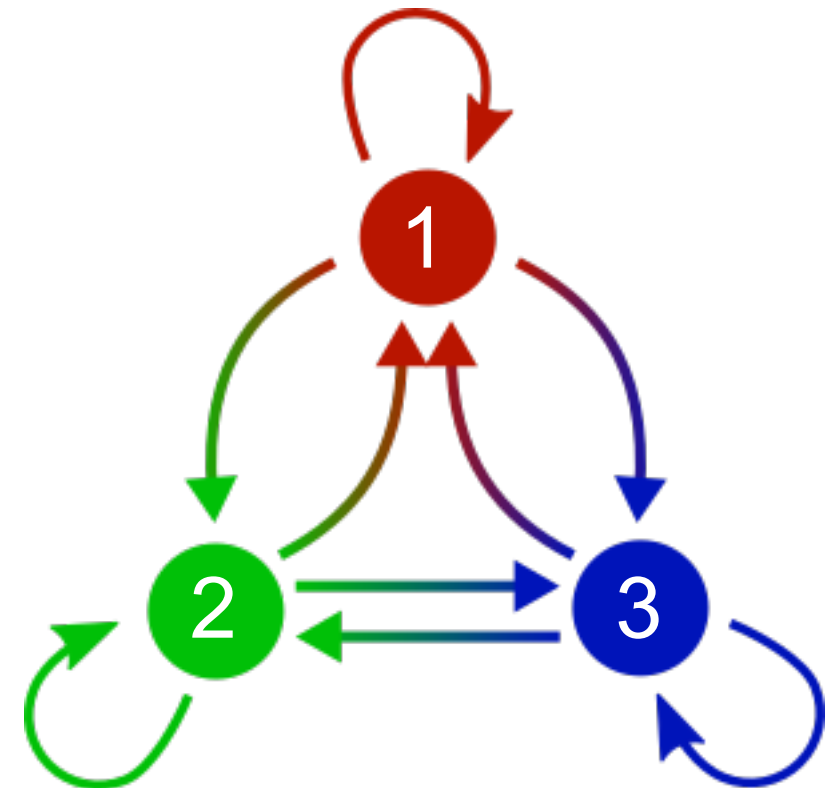
- ▶  $\theta_{ij} = \lim_{T \rightarrow \infty} \mathbb{P}_{\theta}(\xi_t = i, \xi_{t+1} = j)$
- ▶ all one-step transitions possible



# Finite-State Markov Chains

Assume that  $\{\xi_t\}_{t \in \mathbb{N}}$  is a Markov chain on  $\{1, 2, \dots, d\}$

- ▶  $\theta_{ij} = \lim_{T \rightarrow \infty} \mathbb{P}_{\theta}(\xi_t = i, \xi_{t+1} = j)$
- ▶ all one-step transitions possible

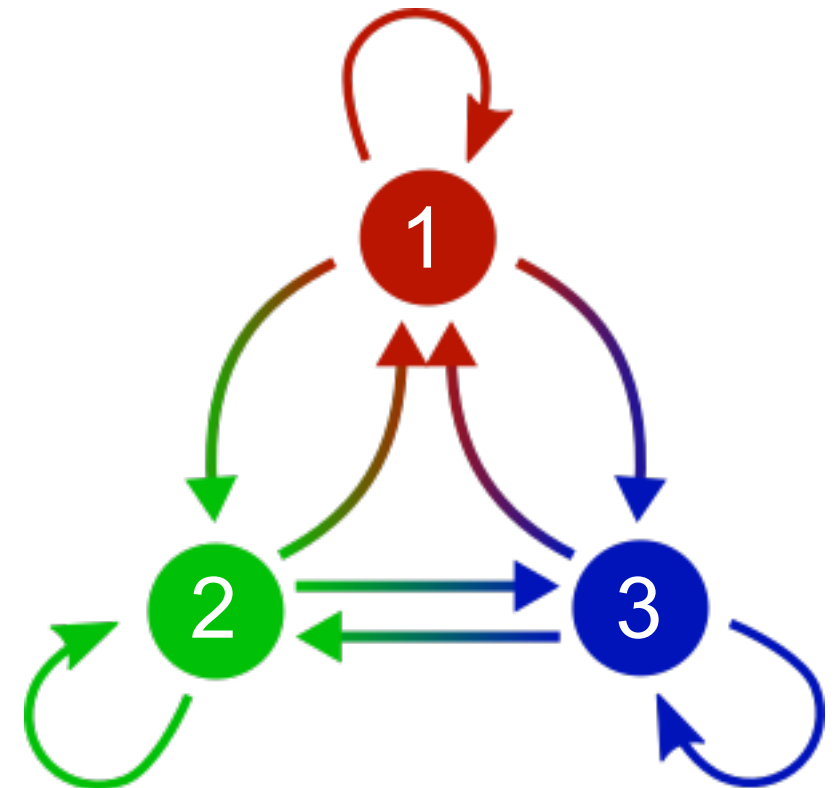


$$\Theta = \left\{ \theta \in \mathbb{R}_{++}^{d \times d} : \sum_{i,j} \theta_{ij} = 1, \sum_j \theta_{ij} = \sum_j \theta_{ji} \forall i \right\}$$

# Finite-State Markov Chains

Assume that  $\{\xi_t\}_{t \in \mathbb{N}}$  is a Markov chain on  $\{1, 2, \dots, d\}$

- ▶  $\theta_{ij} = \lim_{T \rightarrow \infty} \mathbb{P}_{\theta}(\xi_t = i, \xi_{t+1} = j)$
- ▶ all one-step transitions possible



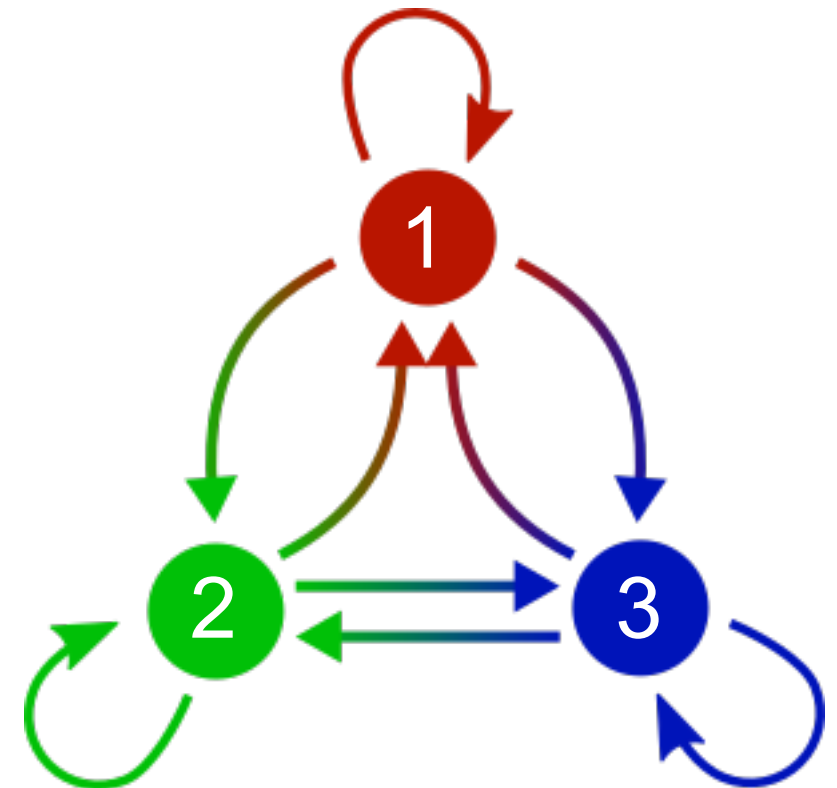
$$\Theta = \left\{ \theta \in \mathbb{R}_{++}^{d \times d} : \sum_{i,j} \theta_{ij} = 1, \sum_j \theta_{ij} = \sum_j \theta_{ji} \forall i \right\}$$

all transitions have probability  $> 0$

# Finite-State Markov Chains

Assume that  $\{\xi_t\}_{t \in \mathbb{N}}$  is a Markov chain on  $\{1, 2, \dots, d\}$

- ▶  $\theta_{ij} = \lim_{T \rightarrow \infty} \mathbb{P}_{\theta}(\xi_t = i, \xi_{t+1} = j)$
- ▶ all one-step transitions possible  
sum of all entries = 1,



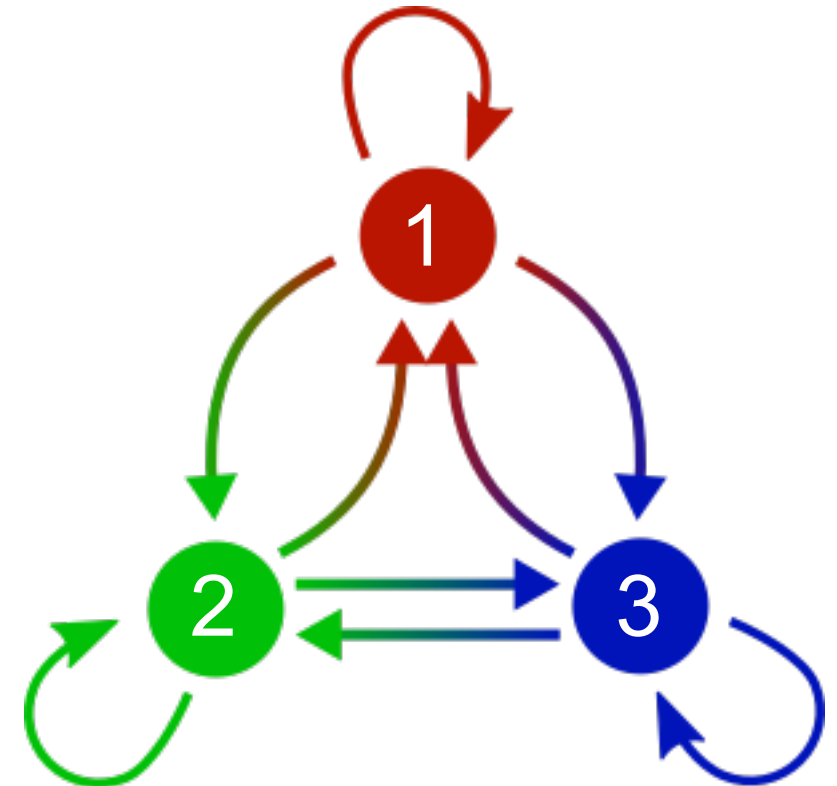
$$\Theta = \left\{ \theta \in \mathbb{R}_{++}^{d \times d} : \underbrace{\sum_{i,j} \theta_{ij}} = 1, \sum_j \theta_{ij} = \sum_j \theta_{ji} \forall i \right\}$$

normalization

# Finite-State Markov Chains

Assume that  $\{\xi_t\}_{t \in \mathbb{N}}$  is a Markov chain on  $\{1, 2, \dots, d\}$

- ▶  $\theta_{ij} = \lim_{T \rightarrow \infty} \mathbb{P}_{\theta}(\xi_t = i, \xi_{t+1} = j)$
- ▶ all one-step transitions possible



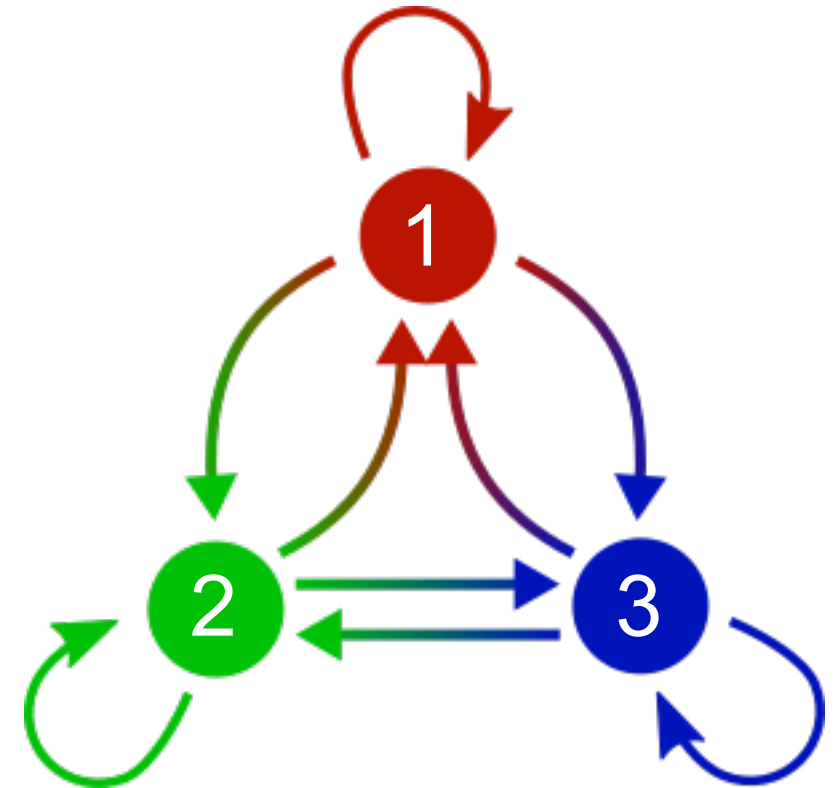
$$\Theta = \left\{ \theta \in \mathbb{R}_{++}^{d \times d} : \sum_{i,j} \theta_{ij} = 1, \underbrace{\sum_j \theta_{ij}} = \sum_j \theta_{ji} \forall i \right\}$$

invariant probability of state  $i$

# Finite-State Markov Chains

Assume that  $\{\xi_t\}_{t \in \mathbb{N}}$  is a Markov chain on  $\{1, 2, \dots, d\}$

- ▶  $\theta_{ij} = \lim_{T \rightarrow \infty} \mathbb{P}_{\theta}(\xi_t = i, \xi_{t+1} = j)$
- ▶ all one-step transitions possible



$$\Theta = \left\{ \theta \in \mathbb{R}_{++}^{d \times d} : \sum_{i,j} \theta_{ij} = 1, \sum_j \theta_{ij} = \underbrace{\sum_j \theta_{ji}}_{\text{invariant probability of state } i} \forall i \right\}$$

invariant probability of state  $i$

# Finite-State Markov Chains

**“Compressing” the raw data:**

$$\left. \begin{array}{l} \text{available observations} \\ (\xi_1, \xi_2, \dots, \xi_T) \end{array} \right\} \longleftrightarrow \left\{ \begin{array}{l} \text{empirical doublet distribution} \\ (\hat{\theta}_T)_{ij} = \frac{1}{T} \sum_{t=1}^T \mathbf{1}_{(\xi_{t-1}, \xi_t) = (i, j)} \end{array} \right.$$

# Finite-State Markov Chains

“Compressing” the raw data:

$$\left. \begin{array}{l} \text{available observations} \\ (\xi_1, \xi_2, \dots, \xi_T) \end{array} \right\} \longleftrightarrow \left\{ \begin{array}{l} \text{empirical doublet distribution} \\ (\hat{\theta}_T)_{ij} = \frac{1}{T} \sum_{t=1}^T \mathbf{1}_{(\xi_{t-1}, \xi_t) = (i,j)} \end{array} \right.$$

- Fisher-Neyman:<sup>1)</sup>  $\hat{\theta}_T$  is a sufficient statistic for  $\theta$
- Dembo & Zeitouni:<sup>2)</sup>  $\hat{\theta}_T$  satisfies an LDP with  $I(\hat{\theta}_T, \theta) = D_c(\hat{\theta}_T \parallel \theta)$

**Definition:** Conditional relative entropy

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

---

<sup>1)</sup> Lehmann & Casella, *Springer*, 1998;

<sup>2)</sup> Dembo & Zeitouni, *Springer*, 1998.



# Autoregressive Gaussian Processes

## Vector autoregressive processes with unknown drift:

- ▶  $\xi_{t+1} = \theta + A\xi_t + \varepsilon_{t+1}$  stationary, driven by Gaussian noise
- ▶  $\hat{\theta}_T = (\mathbb{I}_d - A) \frac{1}{T} \sum_{t=1}^T \xi_t$  satisfies LDP but is *not* sufficient<sup>1)</sup>

## Scalar autoregressive processes with unknown coefficient:

- ▶  $\xi_{t+1} = \theta\xi_t + \varepsilon_{t+1}$  stationary, driven by Gaussian noise
- ▶ Least squares and Yule-Walker estimators satisfy LDPs but are *not* sufficient<sup>2)</sup>

---

<sup>1)</sup> Dembo & Zeitouni, *Springer*, 1998;

<sup>2)</sup> Bercu et al., *Stochastic Processes and their Applications*, 1997.

# I.I.D. Processes with Parametric CDFs

Assume that the  $\{\xi_t\}_{t \in \mathbb{N}}$  are i.i.d. with any of the following CDFs:

- ▶ normal distribution with mean  $\theta$
- ▶ exponential distribution with rate parameter  $\theta$
- ▶ gamma distribution with scale parameter  $\theta$
- ▶ Poisson distribution with rate parameter  $\theta$
- ▶ Bernoulli distribution with success probability  $\theta$
- ▶ geometric distribution with success probability  $\theta$
- ▶ binomial distribution with success probability  $\theta$

# I.I.D. Processes with Parametric CDFs

Then,  $\hat{\theta}_T = \frac{1}{T} \sum_{t=1}^T \xi_t$  is sufficient<sup>1)</sup> and satisfies an LDP,<sup>2)</sup> where

- ▶  $\Lambda(\lambda, \theta) = \log \mathbb{E}_{\theta} \left[ \exp(\lambda^{\top} \xi_t) \right]$  is the log-MGF, and
- ▶  $I(\theta', \theta) = \sup_{\lambda} \theta'^{\top} \lambda - \Lambda(\lambda, \theta)$  is a “regular” rate function.

---

<sup>1)</sup> Lehmann & Casella, *Springer*, 1998;

<sup>2)</sup> Cramér, *Actualités scientifiques et industrielles*, 1938.

# I.I.D. Processes with Parametric CDFs

Then,  $\hat{\theta}_T = \frac{1}{T} \sum_{t=1}^T \xi_t$  is sufficient<sup>1)</sup> and satisfies an LDP,<sup>2)</sup> where

►  $\Lambda(\lambda, \theta) = \log \mathbb{E}_{\theta} [\exp(\lambda^T \xi_t)]$  is the log MGF

►  $I(\hat{\theta})$

Many convex uncertainty sets can be constructed in this way and are thus optimal for some i.i.d. process!

---

<sup>1)</sup> Lehmann & Casella, *Springer*, 1998;

<sup>2)</sup> Cramér, *Actualités scientifiques et industrielles*, 1938.

# Summary & Conclusions

## ► **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  $\hat{\theta}_T$  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  $\hat{\theta}_T$
- radius = decay rate of the out-of-sample disappointment
- invariant under homeomorphic transformations

## ► **Data efficiency**

- Pareto dominance reminiscent of Bahadur efficiency

## ► **Generality of results**

- hold even for non-convex decision problems
- hold even for non-i.i.d. data processes

## ► **Theoretical justification of DRO**

- shape of ambiguity set depends on the data process
- radius of ambiguity set has physical interpretation

## ► **Computation**

- customized algorithms for new DRO models<sup>1)</sup>

---

<sup>1)</sup> Li et al., *ICML*, 2021.

# This Talk is Based on...

- [1] M. Li, T. Sutter and D. Kuhn. **Distributionally Robust Optimization Based on Markovian Data**. ICML, 2021.
- [2] T. Sutter, W. Jongeneel, S. Shafieezadeh Abadeh and D. Kuhn. **From Moderate Deviations Theory to Distributionally Robust Optimization: Learning from Correlated Data**. *Working paper*, 2021.
- [3] T. Sutter, B. Van Parys and D. Kuhn. **A General Framework for Optimal Data-Driven Optimization**. *arXiv:2010.06606*, 2020
- [4] B. Van Parys, P. Mohajerin Esfahani and D. Kuhn. **From Data to Decisions: Distributionally Robust Optimization is Optimal**. *Management Science*, 2020.





# Appendix: Proof Ideas

# Optimizing over Optimization Problems

**Restricted MOP for a fixed decision:**

$$\begin{array}{ll} \underset{\tilde{c}}{\text{minimize}} & \{\tilde{c}(\theta)\}_{\theta \in \Theta} \\ \text{subject to} & \limsup_{T \rightarrow \infty} \frac{1}{T} \log \mathbb{P}_{\theta} \left[ c(\theta) > \tilde{c}(\hat{\theta}_T) \right] \leq -r \quad \forall \theta \in \Theta \end{array}$$

# Optimizing over Optimization Problems

**Restricted MOP for a fixed decision:**

$$\begin{aligned} & \underset{\tilde{c}}{\text{minimize}} && \{ \tilde{c}(\theta) \}_{\theta \in \Theta} \\ & \text{subject to} && \limsup_{T \rightarrow \infty} \frac{1}{T} \log \mathbb{P}_{\theta} \left[ c(\theta) > \tilde{c}(\hat{\theta}_T) \right] \leq -r \quad \forall \theta \in \Theta \end{aligned}$$

**Pareto-dominant solution:**

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

# Feasibility

**Theorem:**  $\tilde{c}^*$  is feasible in the MOP.

# Feasibility

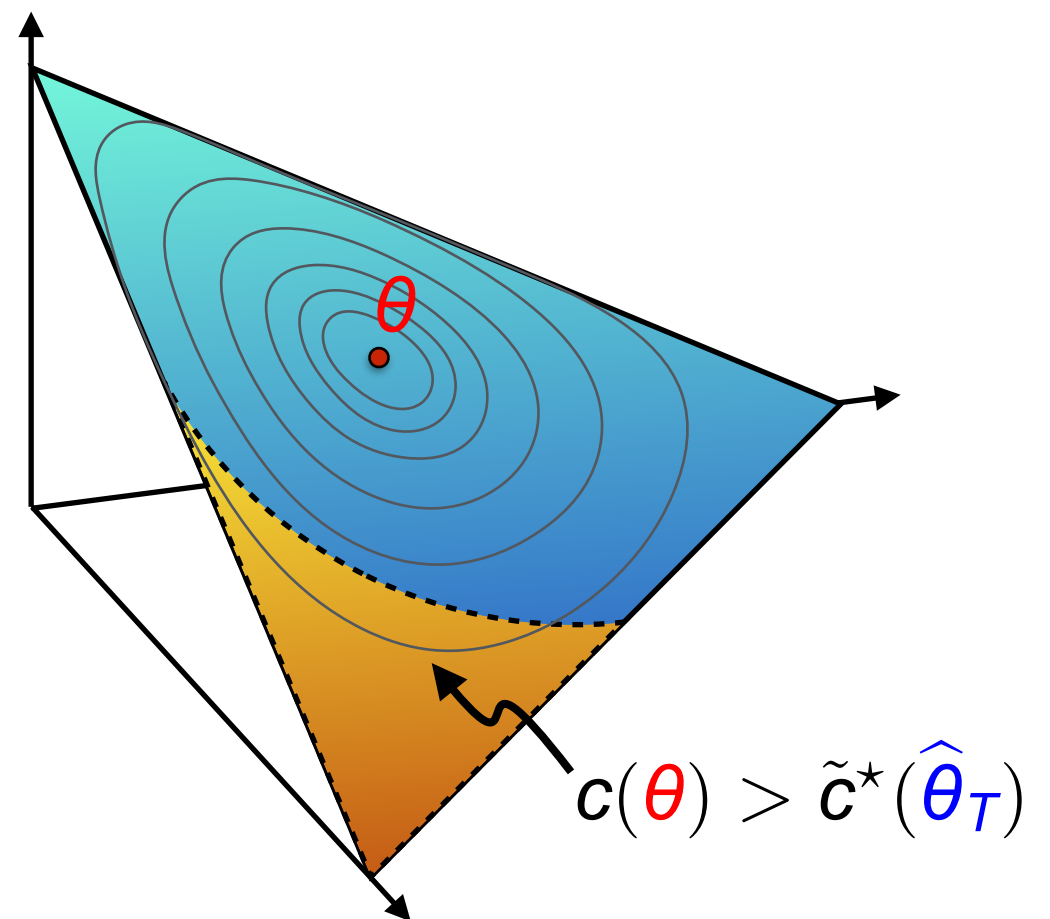
**Theorem:**  $\tilde{c}^*$  is feasible in the MOP.

$$\begin{aligned} c(\theta) > \tilde{c}^*(\hat{\theta}_T) &\implies c(\theta) > \sup_{\theta' \in \Theta} \left\{ c(\theta') : l(\hat{\theta}_T, \theta') \leq r \right\} \\ &\implies l(\hat{\theta}_T, \theta) > r \end{aligned}$$

# Feasibility

**Theorem:**  $\tilde{c}^*$  is feasible in the MOP.

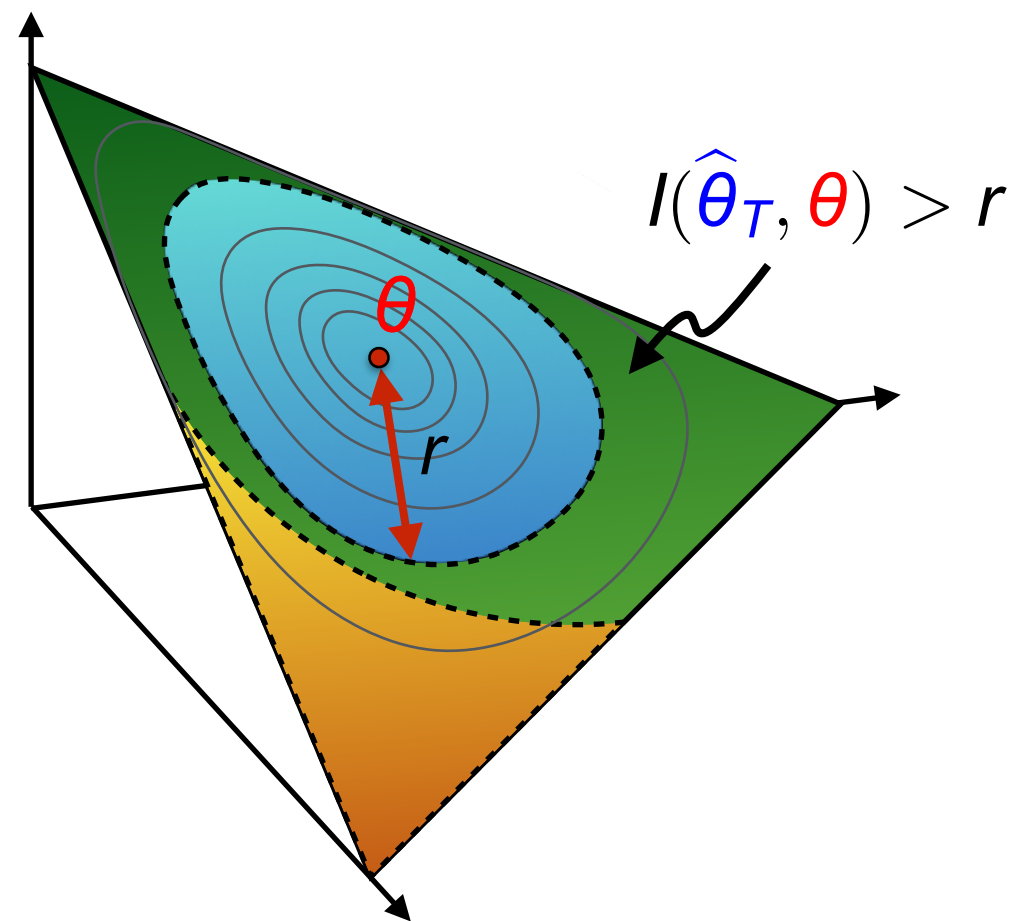
$$\begin{aligned} c(\theta) > \tilde{c}^*(\hat{\theta}_T) &\implies c(\theta) > \sup_{\theta' \in \Theta} \left\{ c(\theta') : I(\hat{\theta}_T, \theta') \leq r \right\} \\ &\implies I(\hat{\theta}_T, \theta) > r \end{aligned}$$



# Feasibility

**Theorem:**  $\tilde{c}^*$  is feasible in the MOP.

$$\begin{aligned} c(\theta) > \tilde{c}^*(\hat{\theta}_T) &\implies c(\theta) > \sup_{\theta' \in \Theta} \left\{ c(\theta') : l(\hat{\theta}_T, \theta') \leq r \right\} \\ &\implies l(\hat{\theta}_T, \theta) > r \end{aligned}$$

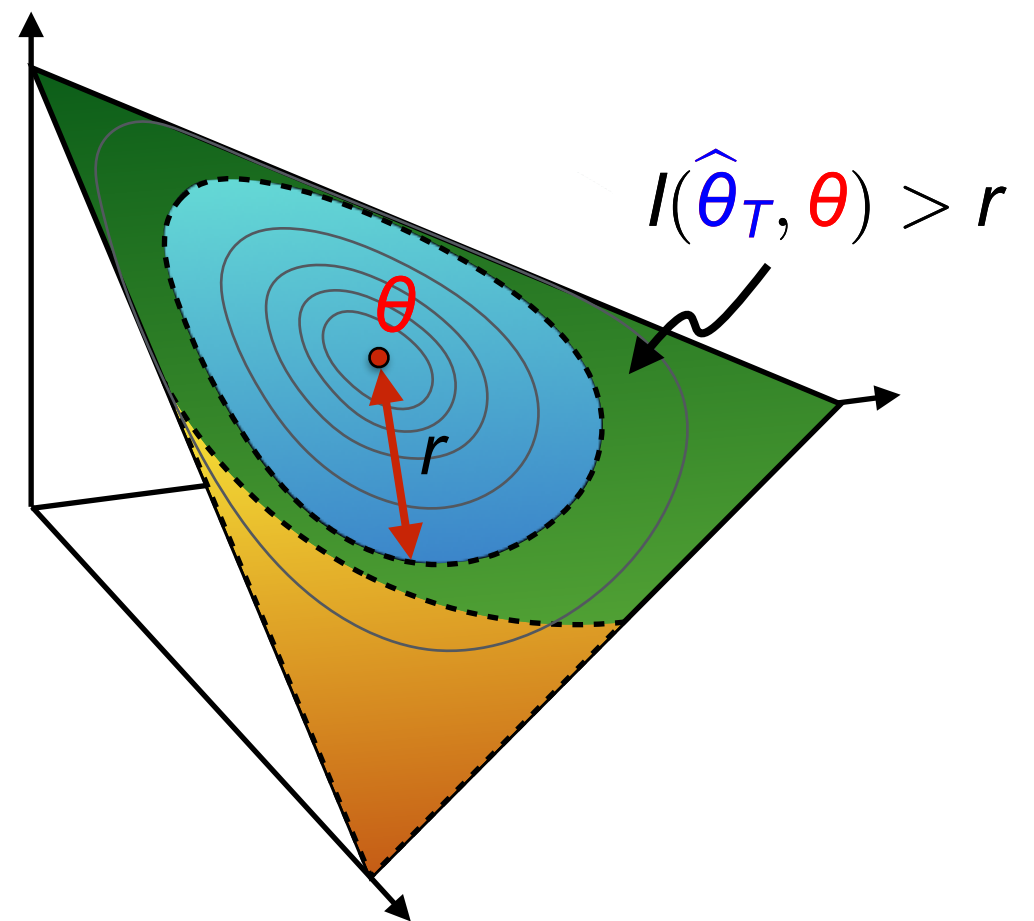


# Feasibility

**Theorem:**  $\tilde{c}^*$  is feasible in the MOP.

$$\begin{aligned} c(\theta) > \tilde{c}^*(\hat{\theta}_T) &\implies c(\theta) > \sup_{\theta' \in \Theta} \left\{ c(\theta') : I(\hat{\theta}_T, \theta') \leq r \right\} \\ &\implies I(\hat{\theta}_T, \theta) > r \end{aligned}$$

$$\begin{aligned} &\implies \mathbb{P}_{\theta} \left[ c(\theta) > \tilde{c}^*(\hat{\theta}_T) \right] \\ &\leq \underbrace{\mathbb{P}_{\theta} \left[ I(\hat{\theta}_T, \theta) > r \right]} \\ &\leq e^{-rT+o(T)} \quad \forall \theta \in \Theta \end{aligned}$$





# Optimality

**Theorem:** If  $r > 0$ , then  $\tilde{c}^*$  is Pareto-dominant in the MOP.

# Optimality

**Theorem:** If  $r > 0$ , then  $\tilde{c}^*$  is Pareto-dominant in the MOP.

Assume that there is  $\tilde{c}$  feasible in MOP and  $\theta_1$  with  $\tilde{c}(\theta_1) < \tilde{c}^*(\theta_1)$

# Optimality

**Theorem:** If  $r > 0$ , then  $\tilde{c}^*$  is Pareto-dominant in the MOP.

Assume that there is  $\tilde{c}$  feasible in MOP and  $\theta_1$  with  $\tilde{c}(\theta_1) < \tilde{c}^*(\theta_1)$

$$\implies \tilde{c}(\theta_1) < \sup_{\theta \in \Theta} \{c(\theta) : l(\theta_1, \theta) \leq r\}$$

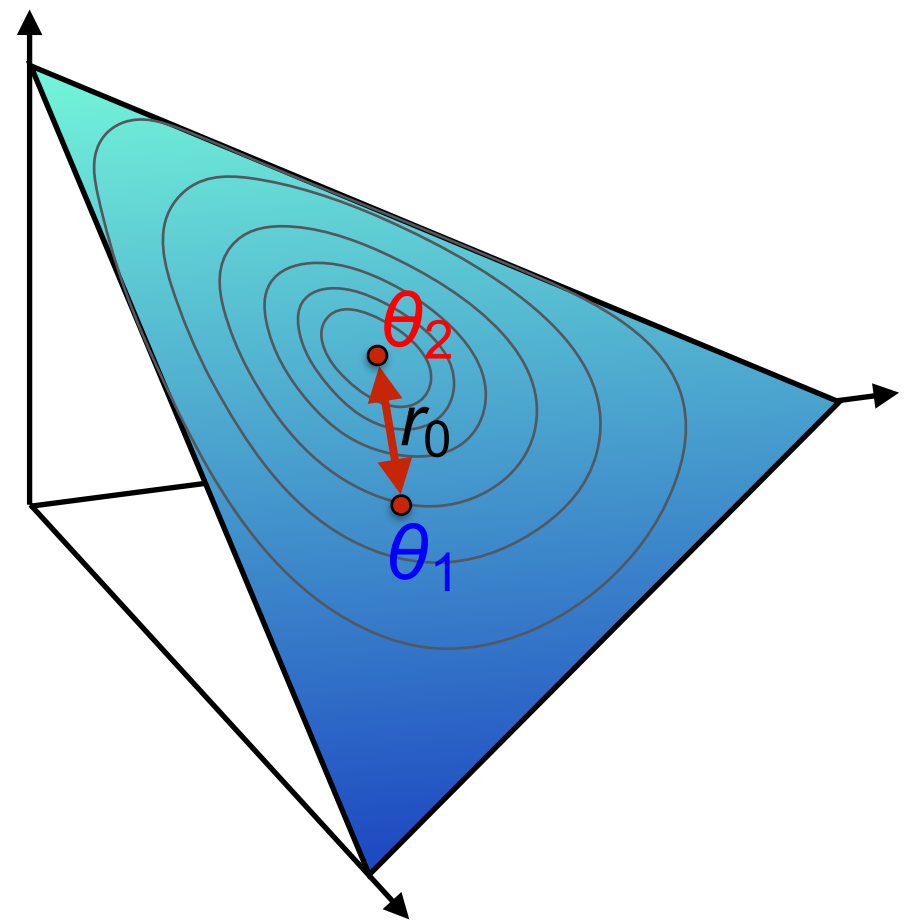
# Optimality

**Theorem:** If  $r > 0$ , then  $\tilde{c}^*$  is Pareto-dominant in the MOP.

Assume that there is  $\tilde{c}$  feasible in MOP and  $\theta_1$  with  $\tilde{c}(\theta_1) < \tilde{c}^*(\theta_1)$

$$\implies \tilde{c}(\theta_1) < \sup_{\theta \in \Theta} \{c(\theta) : l(\theta_1, \theta) \leq r\}$$

$$\implies \exists \theta_2 : \tilde{c}(\theta_1) < c(\theta_2) \text{ and } l(\theta_1, \theta_2) = r_0 < r$$



# Optimality

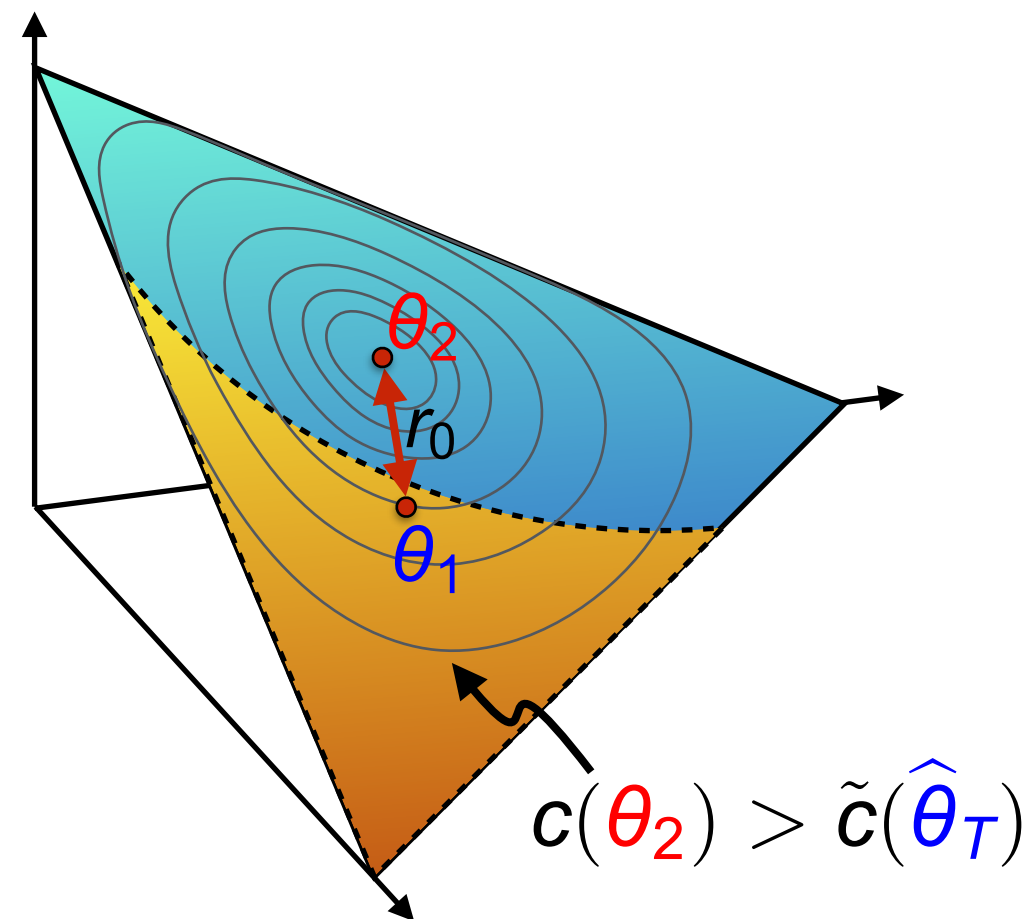
**Theorem:** If  $r > 0$ , then  $\tilde{c}^*$  is Pareto-dominant in the MOP.

Assume that there is  $\tilde{c}$  feasible in MOP and  $\theta_1$  with  $\tilde{c}(\theta_1) < \tilde{c}^*(\theta_1)$

$$\implies \tilde{c}(\theta_1) < \sup_{\theta \in \Theta} \{c(\theta) : l(\theta_1, \theta) \leq r\}$$

$$\implies \exists \theta_2 : \tilde{c}(\theta_1) < c(\theta_2) \text{ and } l(\theta_1, \theta_2) = r_0 < r$$

$$\implies \mathbb{P}_{\theta_2} \left[ c(\theta_2) > \tilde{c}(\hat{\theta}_T) \right] \geq e^{-r_1 T + o(T)}$$



# Optimality

**Theorem:** If  $r > 0$ , then  $\tilde{c}^*$  is Pareto-dominant in the MOP.

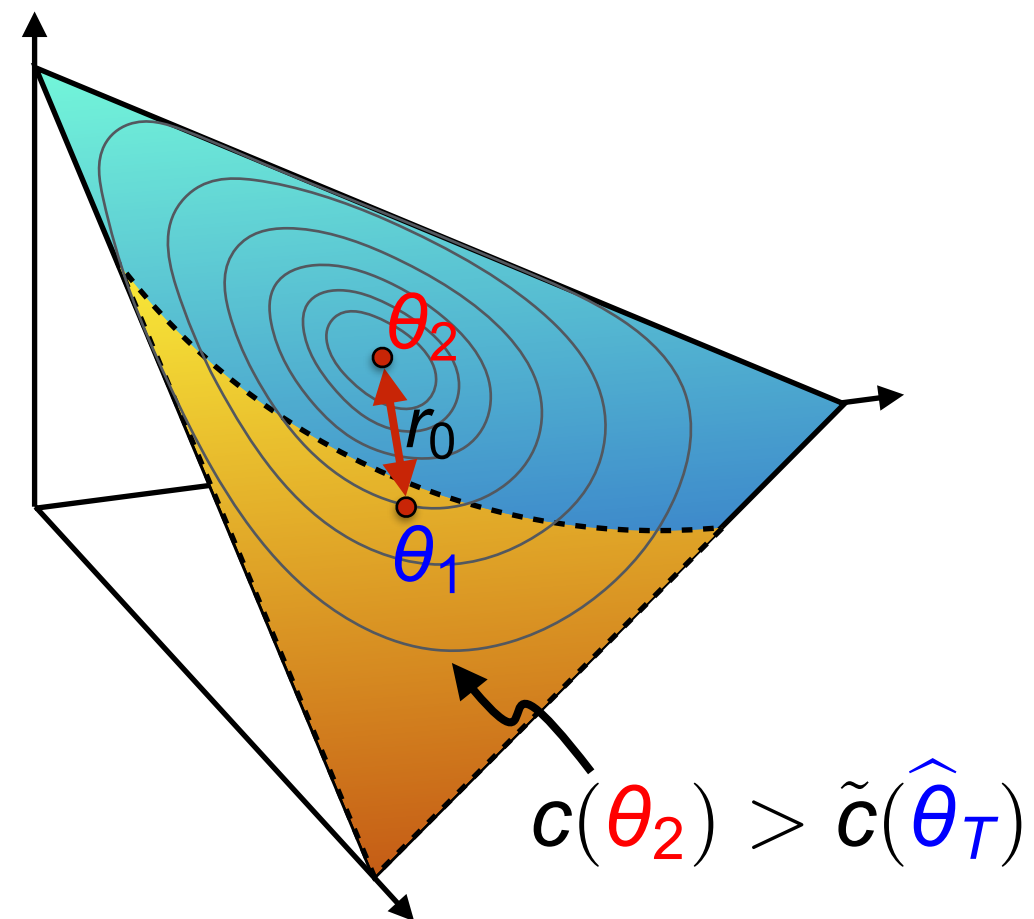
Assume that there is  $\tilde{c}$  feasible in MOP and  $\theta_1$  with  $\tilde{c}(\theta_1) < \tilde{c}^*(\theta_1)$

$$\implies \tilde{c}(\theta_1) < \sup_{\theta \in \Theta} \{c(\theta) : l(\theta_1, \theta) \leq r\}$$

$$\implies \exists \theta_2 : \tilde{c}(\theta_1) < c(\theta_2) \text{ and } l(\theta_1, \theta_2) = r_0 < r$$

$$\implies \mathbb{P}_{\theta_2} \left[ c(\theta_2) > \tilde{c}(\hat{\theta}_T) \right] \geq e^{-r_1 T + o(T)}$$

$$\implies \tilde{c} \text{ infeasible in MOP} \quad \text{⚡}$$



# Appendix: Data-Driven Control

# From Data to Controllers?

Closed-loop LTI system:  $\mathbf{x}_{t+1} = \boldsymbol{\theta} \mathbf{x}_t + \mathbf{w}_t$

Least squares estimator:  $\hat{\boldsymbol{\theta}}_T = \left( \sum_{t=1}^T \hat{\mathbf{x}}_t \hat{\mathbf{x}}_{t-1}^\top \right) \left( \sum_{t=1}^T \hat{\mathbf{x}}_{t-1} \hat{\mathbf{x}}_{t-1}^\top \right)^{-1}$

**Theorem:** The modified least squares estimator  $\boldsymbol{\theta} + \sqrt[4]{T}(\hat{\boldsymbol{\theta}}_T - \boldsymbol{\theta})$  satisfies a moderate deviations principle with rate function

$$I(\boldsymbol{\theta}', \boldsymbol{\theta}) = \frac{1}{2} \text{tr} \left( \mathbf{S}_{\mathbf{w}}^{-1} (\boldsymbol{\theta}' - \boldsymbol{\theta}) \mathbf{S}_{\boldsymbol{\theta}} (\boldsymbol{\theta}' - \boldsymbol{\theta})^\top \right),$$

where  $\mathbf{S}_{\boldsymbol{\theta}}$  solves the Lyapunov equation  $\mathbf{S}_{\boldsymbol{\theta}} = \boldsymbol{\theta} \mathbf{S}_{\boldsymbol{\theta}} \boldsymbol{\theta}^\top + \mathbf{S}_{\mathbf{w}}$ .

$$\implies \text{DRO bounds on } J(\boldsymbol{\theta}) = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_{\boldsymbol{\theta}} \left[ \mathbf{x}_t^\top (\mathbf{Q} + \mathbf{K}^\top \mathbf{R} \mathbf{K}) \mathbf{x}_t \right]$$