

Introduction to Contextual (Stochastic) Optimization



Erick Delage
GERAD & Department of Decision Sciences

HEC MONTRÉAL



*(Joint work with Utsav Sadana, Abhilash Chenreddy, Alexandre Forel,
Emma Frejinger, Thibaut Vidal)*

*ICSP Tutorial on End-to-end Learning
July 23, 2023*



Canada
Research
Chairs

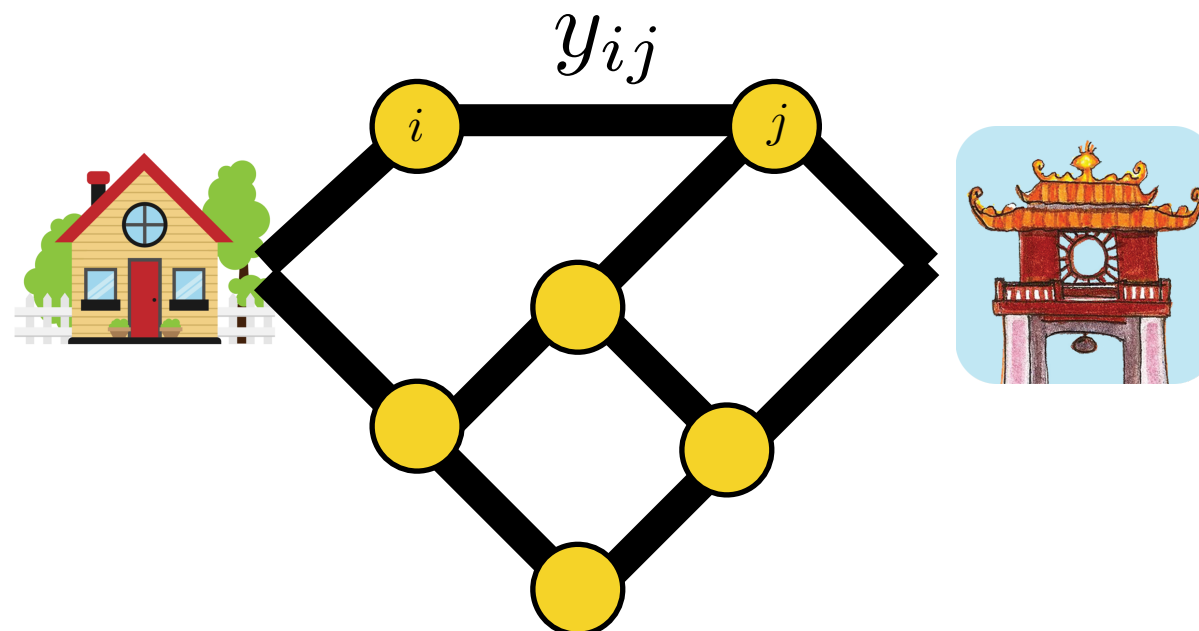
Chaires
de recherche
du Canada

Canada

**Why contextual stochastic
optimization?**

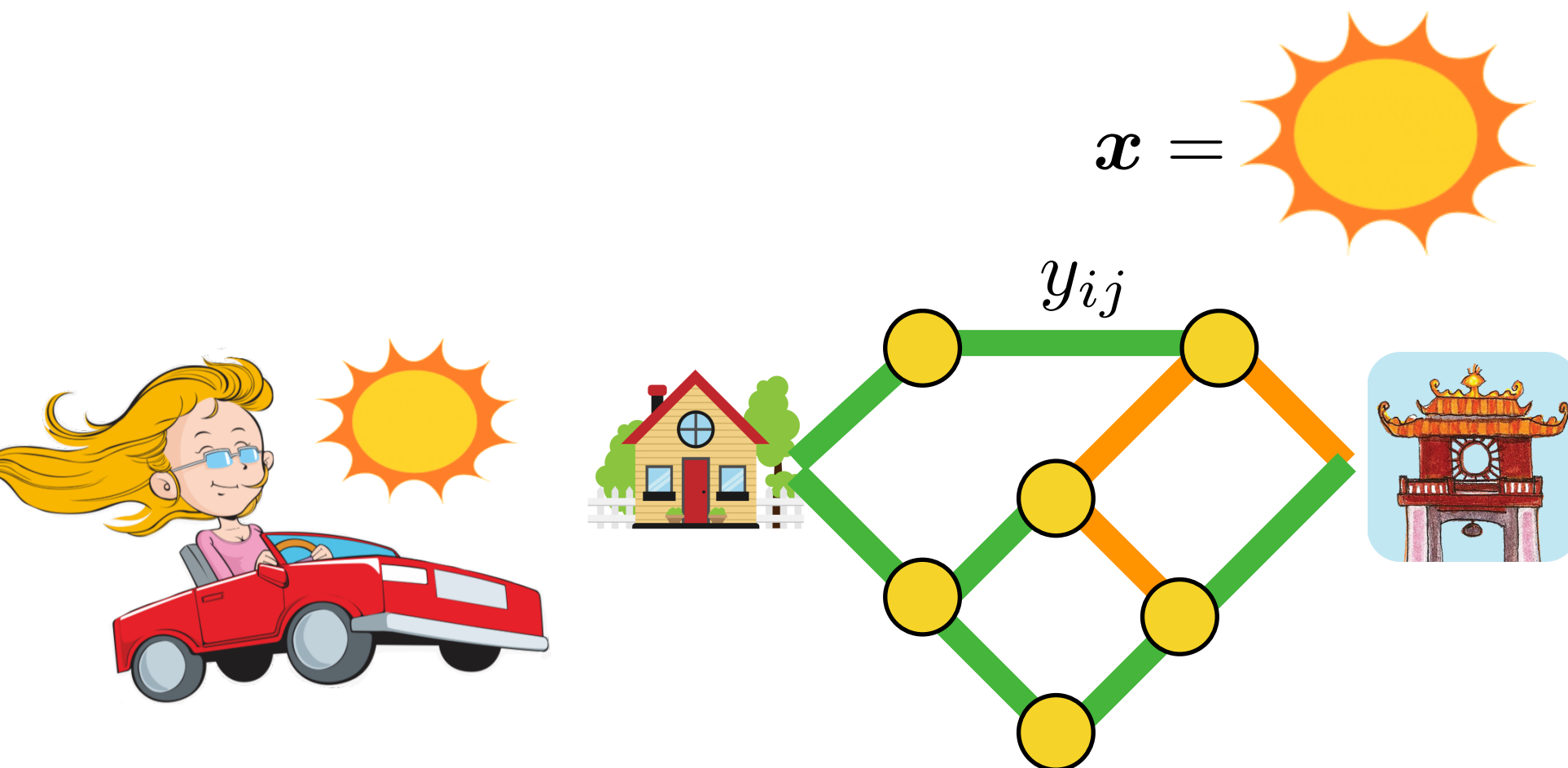
Decision Making with Contextual Information

- Revealed contextual information x
- Hidden random variables y



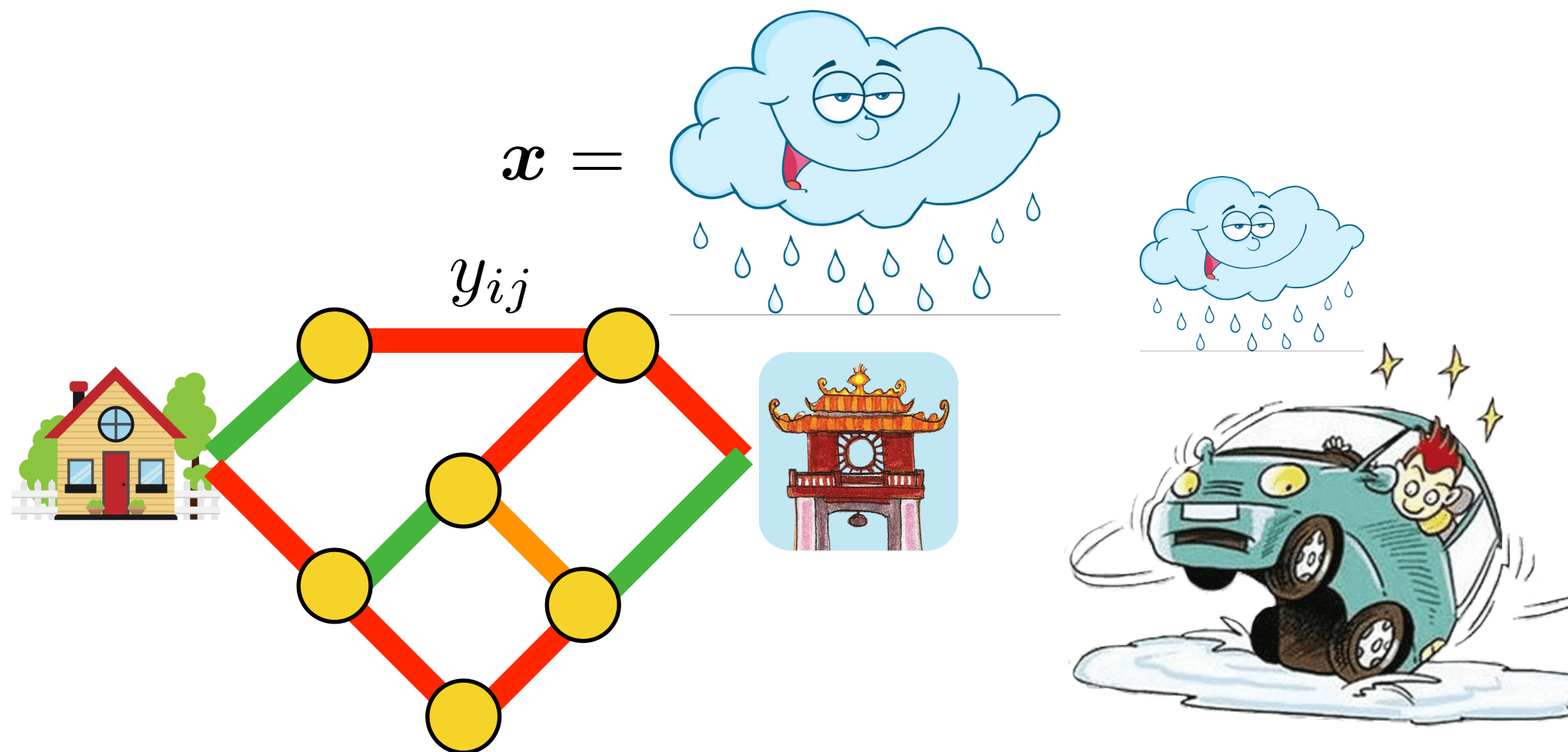
Decision Making with Contextual Information

- Revealed contextual information x
- Hidden random variables y



Decision Making with Contextual Information

- Revealed contextual information x
- Hidden random variables y

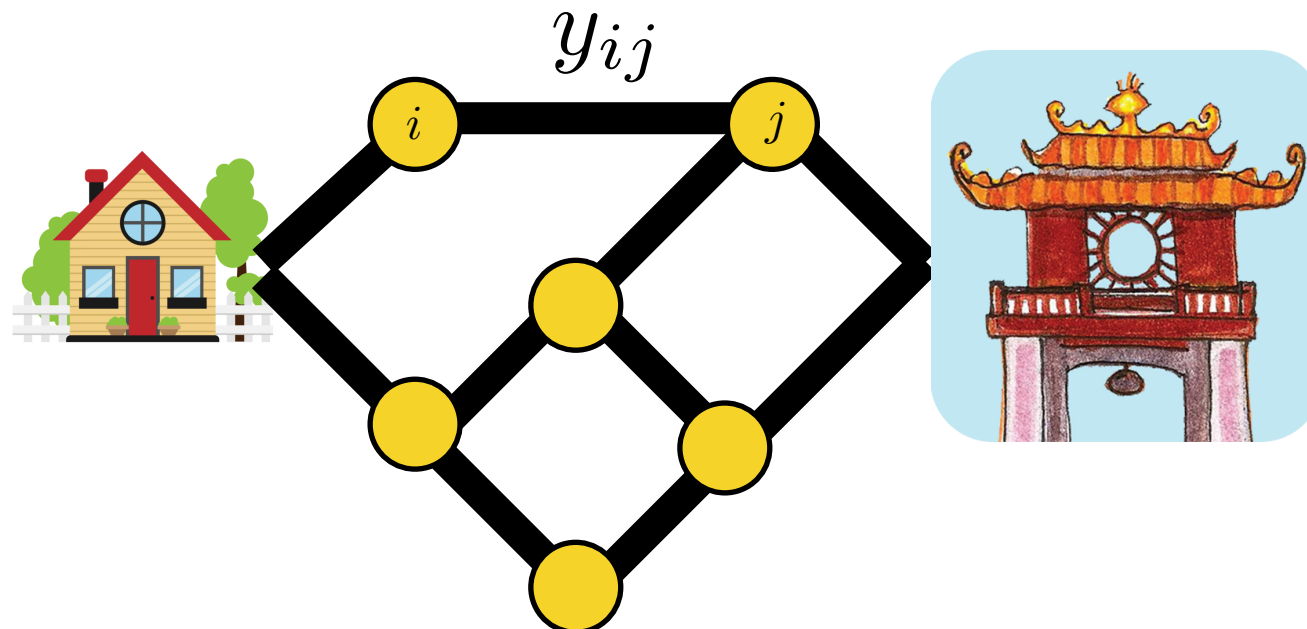
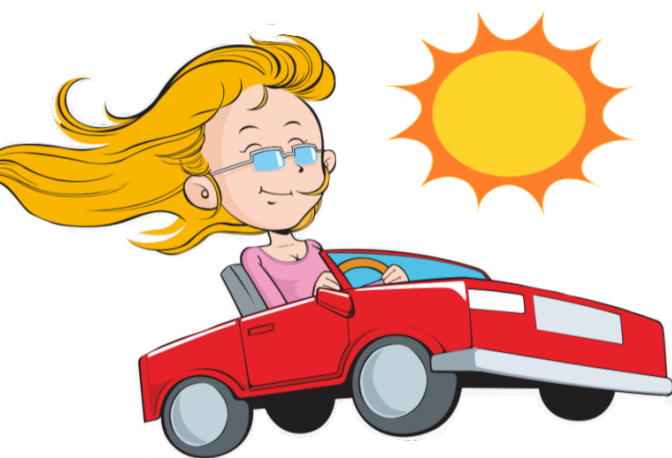


Decision Making with Contextual Information

- Revealed contextual information x
- Hidden random variables y

$$\mathbb{E}^{\mathbb{P}}[y | x = \text{Sun}]$$

$$\mathbb{E}^{\mathbb{P}}[y | x = \text{Rain}]$$

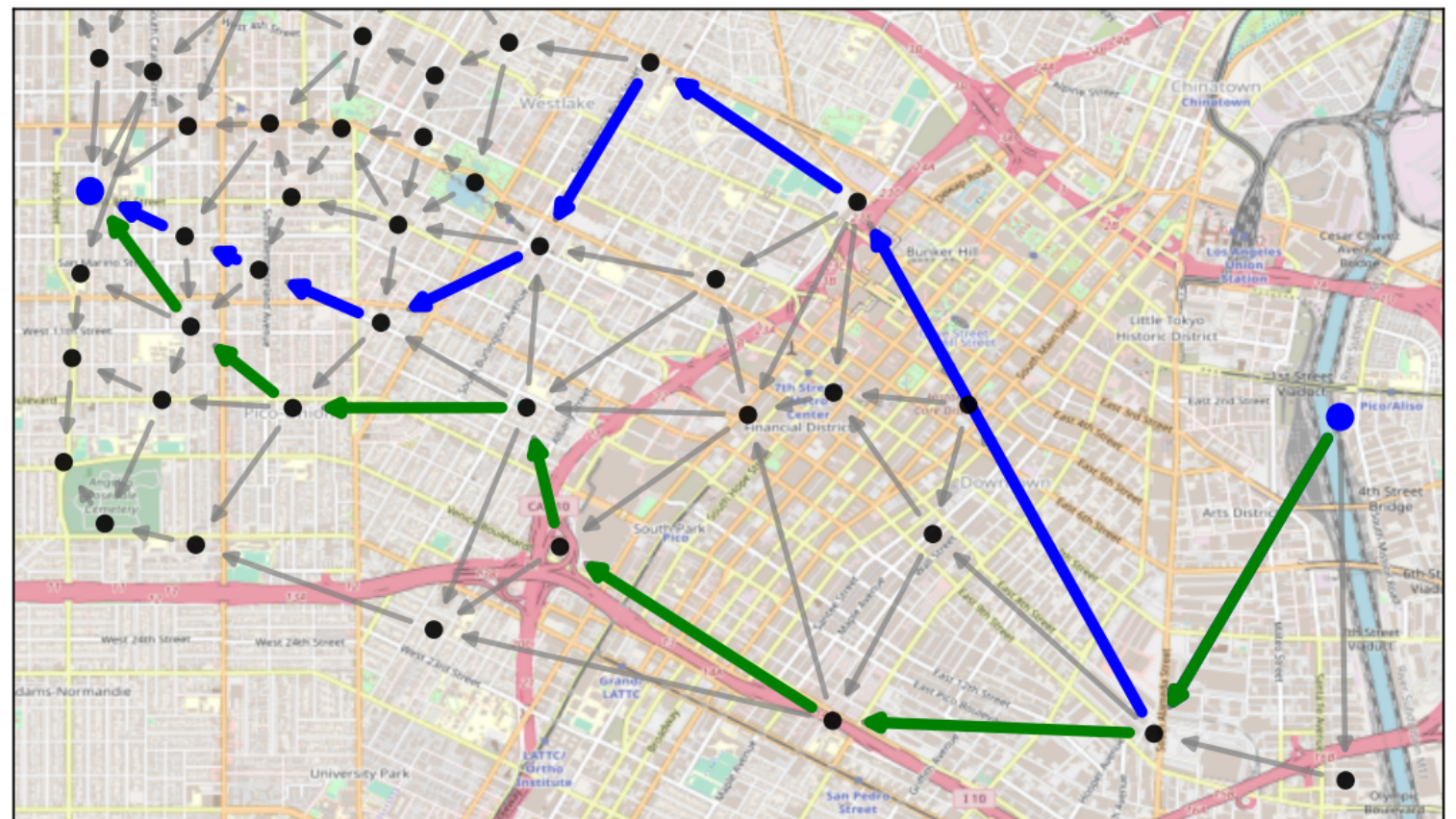


Practical motivation

Example I:
Shortest path over Los Angeles downtown (Kallus & Mao, 2022)

Problem: find shortest path
traversing Los Angeles downtown area
from East to West

Travel times over all arcs are uncertain. We
have relevant contextual information.

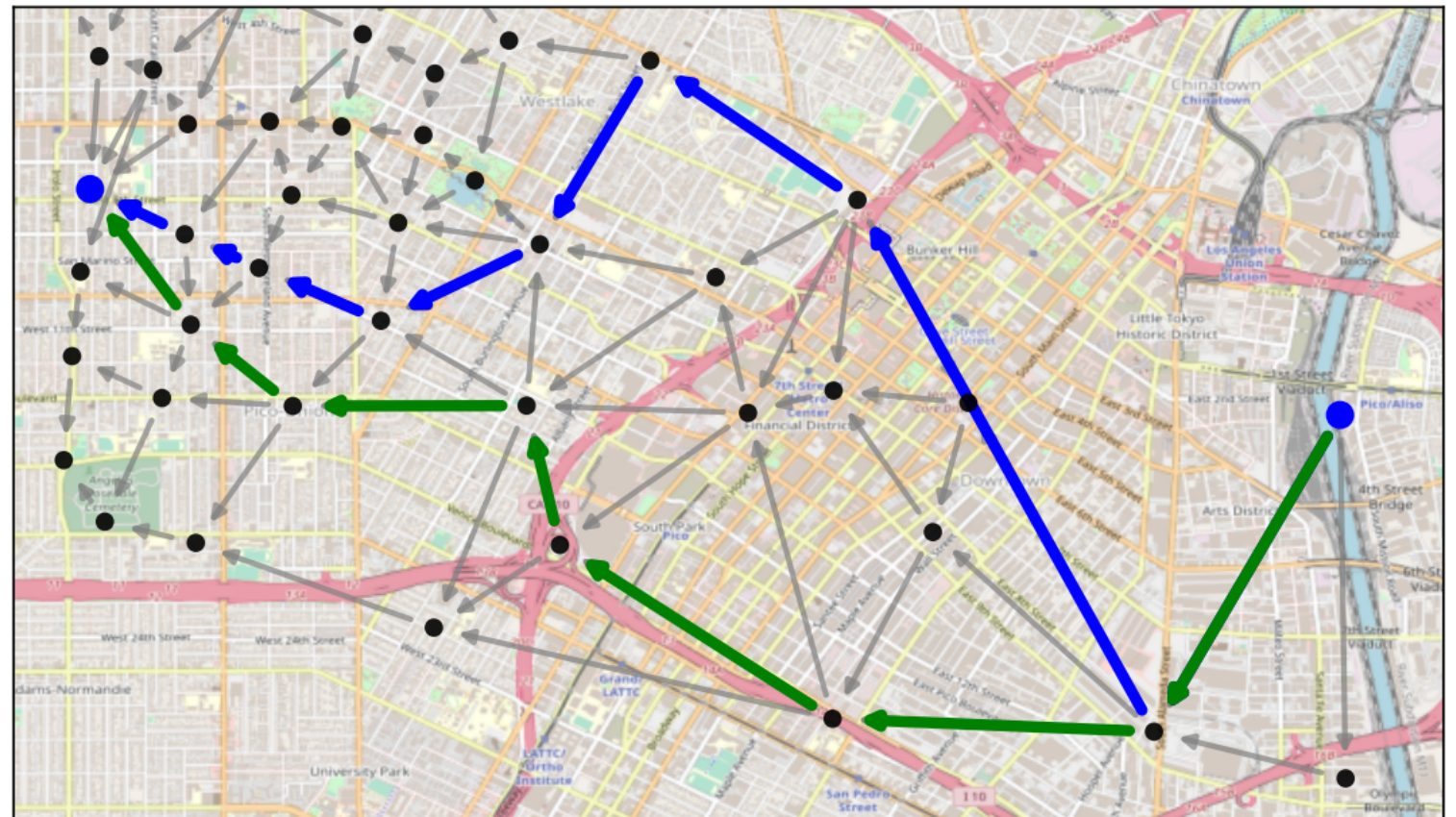


Practical motivation

Example I:
Shortest path over Los Angeles downtown (Kallus & Mao, 2022)

Problem: find shortest path
traversing Los Angeles downtown area
from East to West

Travel times over all arcs are uncertain. We
have relevant contextual information.



Green path is optimal

Blue path is optimal

Period	Temp.	Wind speed	Rain	Visibility	Day	Month
Midday	57.17	4	0	6.99	2	11
AM	57.17	4	0	6.99	2	11

Practical motivation

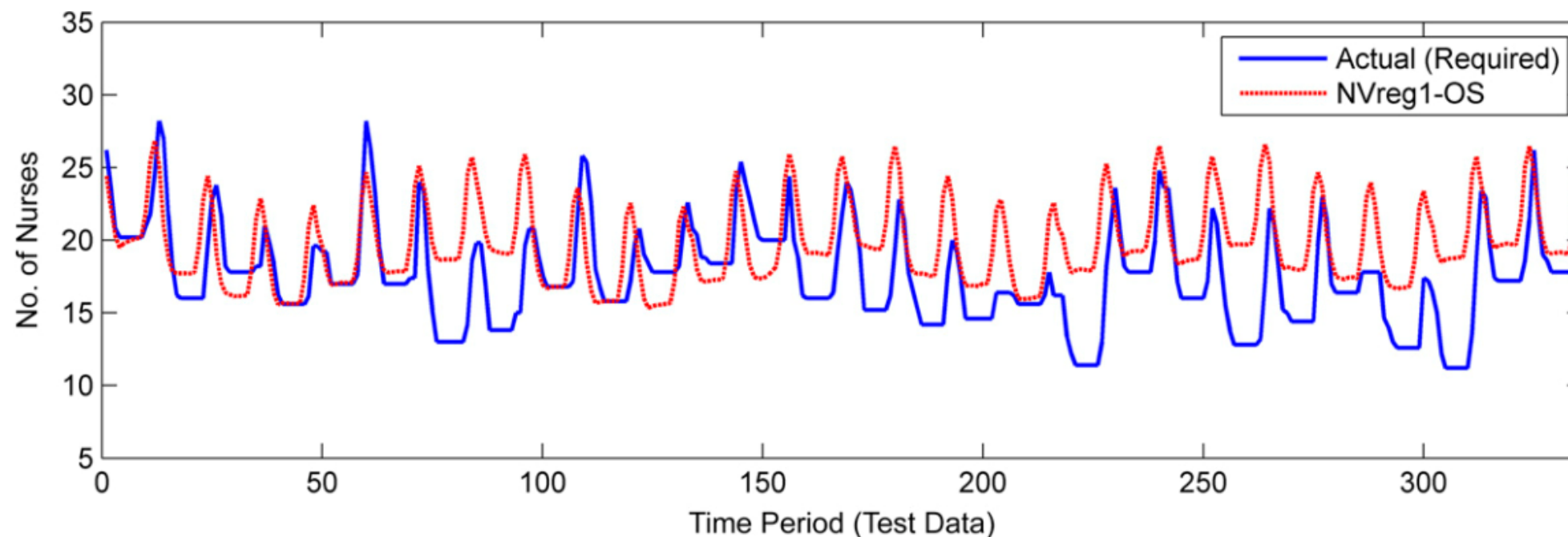
Example 2:

Nurse Staffing in a Hospital (Ban & Rudin, 2019)



Decide how many nurse to schedule on a given day:
large penalty for under-/over-staffing
➤ **A newsvendor model** with uncertain demand

Historical data:
Demand and context



Features

Day of the week

Time of the day

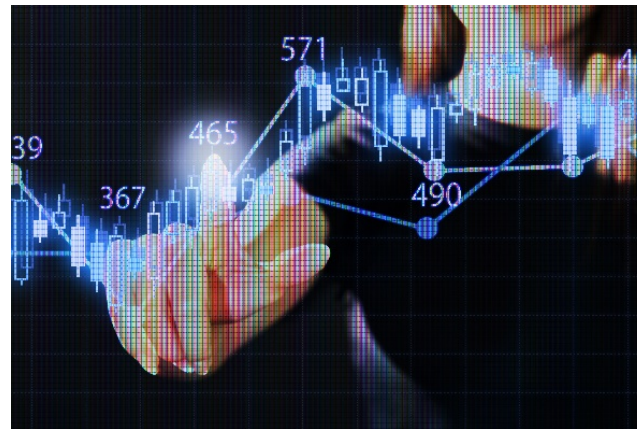
Past demand observations

Practical motivation

In uncertain environments: we should use available contextual information to improve decisions



Manage inventory



Build portfolio



Deliver packages

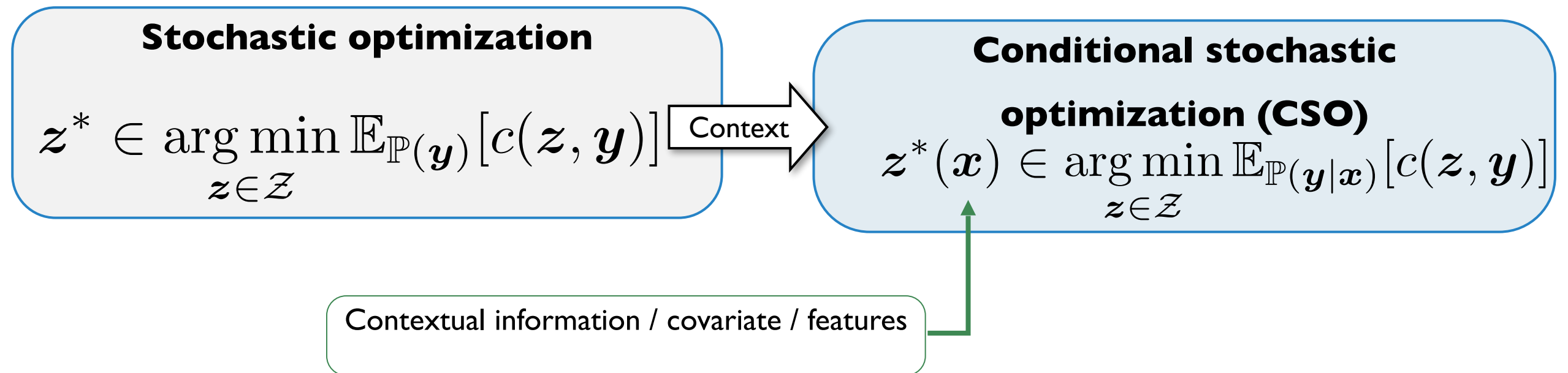
**What is contextual
optimization?**

Problem Definition

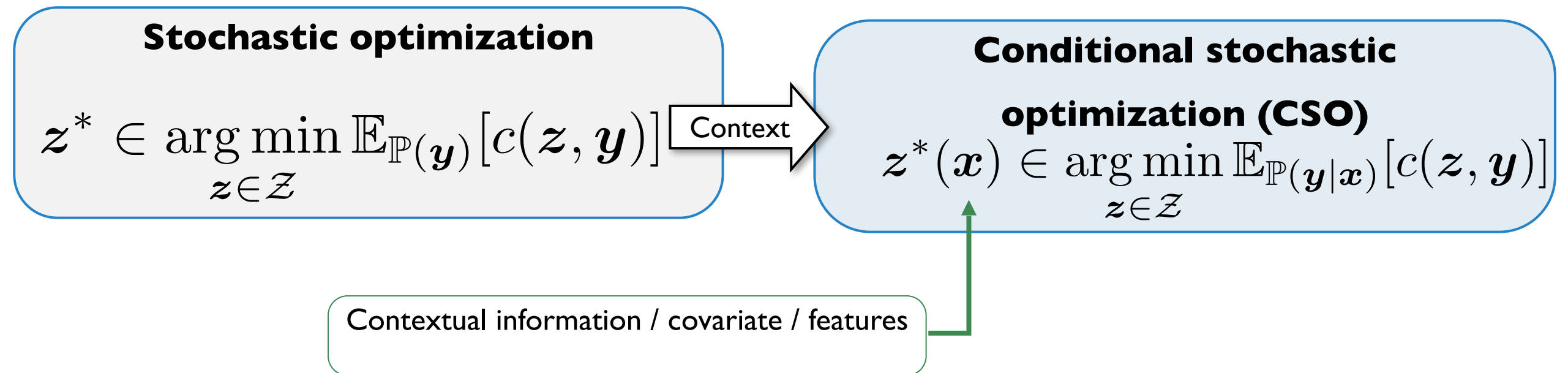
Stochastic optimization

$$\mathbf{z}^* \in \arg \min_{\mathbf{z} \in \mathcal{Z}} \mathbb{E}_{\mathbb{P}(\mathbf{y})} [c(\mathbf{z}, \mathbf{y})]$$

Problem Definition



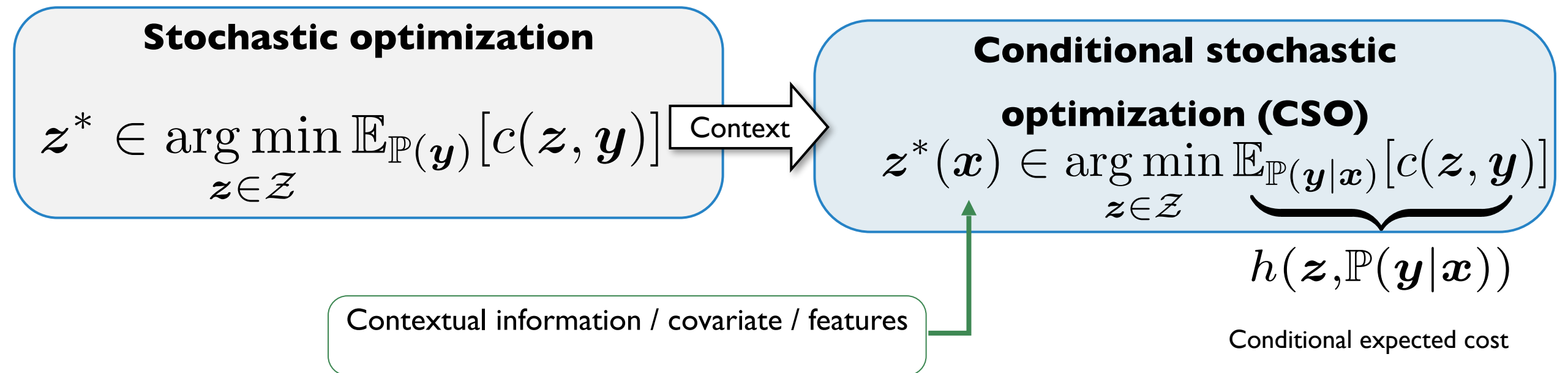
Problem Definition



Connection between CSO and policy optimization:

$$\pi^* \in \arg \min_{\pi: \mathcal{X} \rightarrow \mathcal{Z}} \mathbb{E}_{\mathbb{P}} [c(\pi(\mathbf{x}), \mathbf{y})] \Leftrightarrow \pi^*(\mathbf{x}) \in \arg \min_{\mathbf{z} \in \mathcal{Z}} \mathbb{E}_{\mathbb{P}(\mathbf{y}|\mathbf{x})} [c(\mathbf{z}, \mathbf{y})] \text{ a.s.}$$

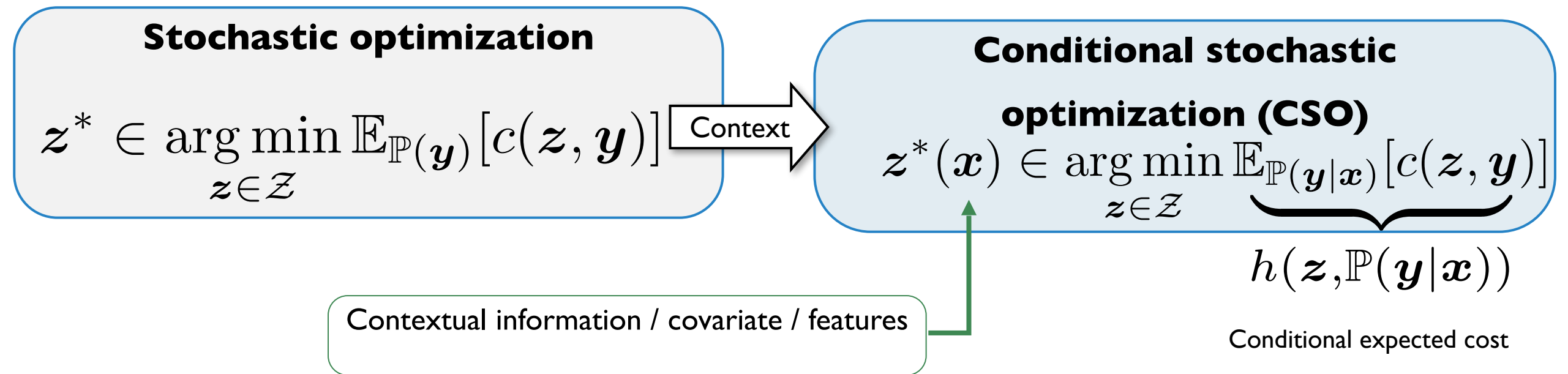
Problem Definition



Connection between CSO and policy optimization:

$$\pi^* \in \arg \min_{\pi: \mathcal{X} \rightarrow \mathcal{Z}} \mathbb{E}_{\mathbb{P}} [c(\pi(\mathbf{x}), \mathbf{y})] \Leftrightarrow \pi^*(\mathbf{x}) \in \arg \min_{z \in \mathcal{Z}} \mathbb{E}_{\mathbb{P}(\mathbf{y}|\mathbf{x})} [c(z, \mathbf{y})] \text{ a.s.}$$

Problem Definition



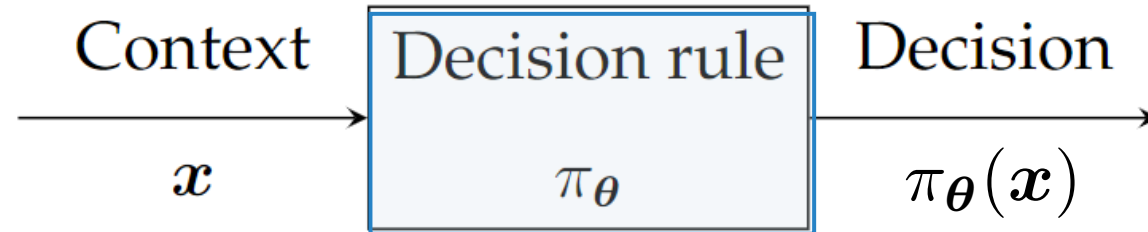
Connection between CSO and policy optimization:

$$\pi^* \in \arg \min_{\pi: \mathcal{X} \rightarrow \mathcal{Z}} \underbrace{\mathbb{E}_{\mathbb{P}} [c(\pi(\mathbf{x}), \mathbf{y})]}_{H(\pi, \mathbb{P})} \Leftrightarrow \pi^*(\mathbf{x}) \in \arg \min_{z \in \mathcal{Z}} \mathbb{E}_{\mathbb{P}(\mathbf{y}|\mathbf{x})} [c(\mathbf{z}, \mathbf{y})] \text{ a.s.}$$

(Unconditional) expected cost

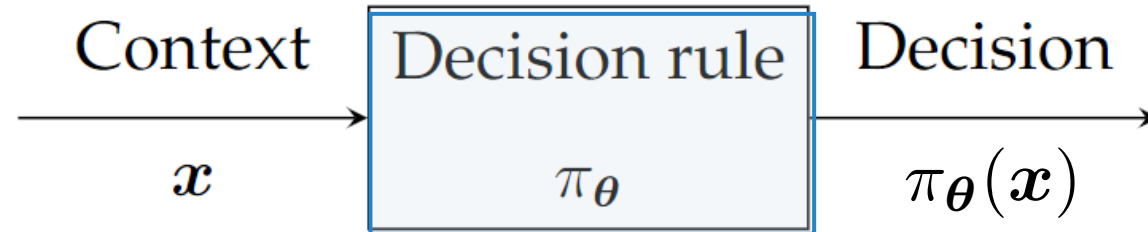
Overview of the three frameworks

Decision rule/Policy optimization

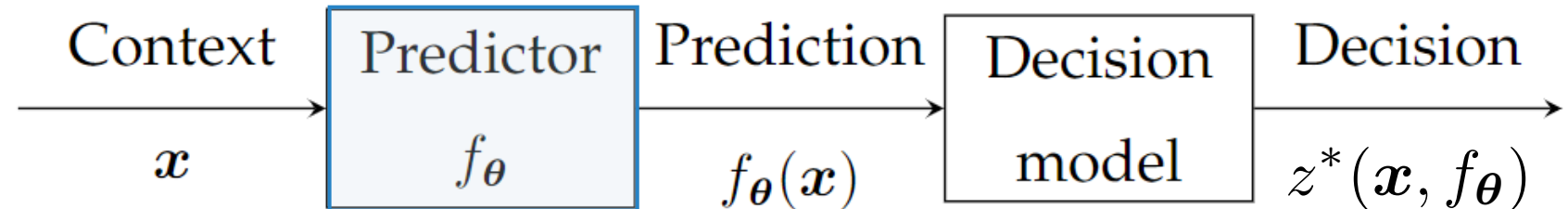


Overview of the three frameworks

Decision rule/Policy optimization

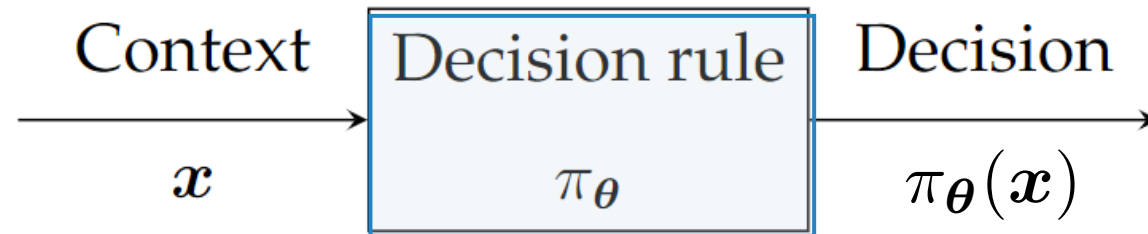


Sequential learning and optimization

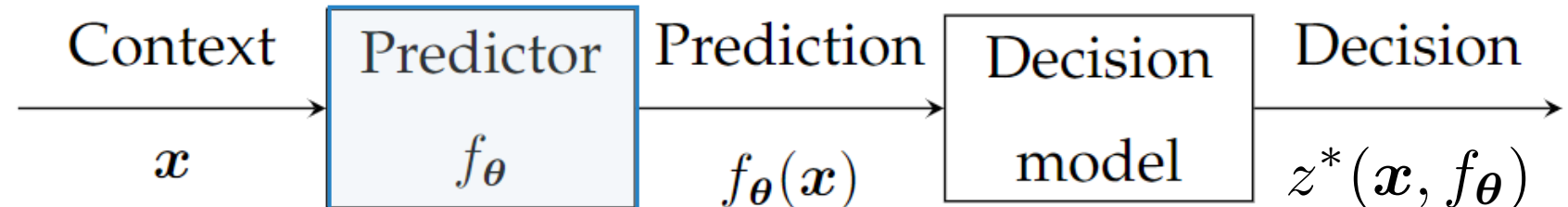


Overview of the three frameworks

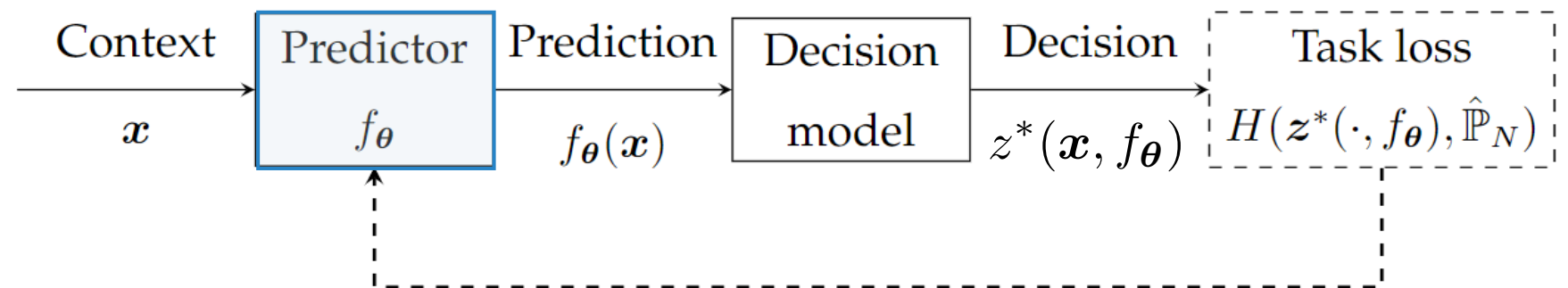
Decision rule/Policy optimization



Sequential learning and optimization



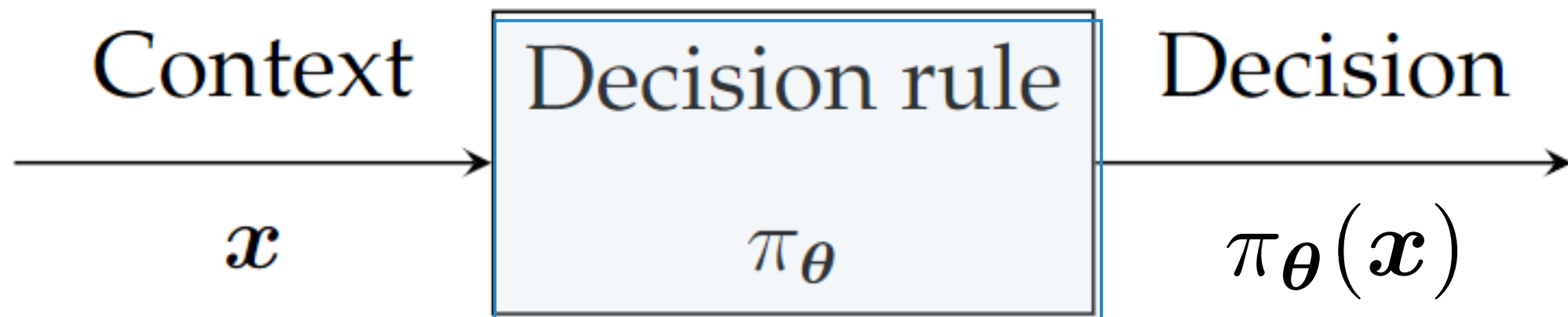
Integrated learning and optimization



Outline of the Tutorial

- Decision rule optimization
- Sequential learning and optimization
- Integrated learning and optimization
- Take-away messages

Decision rule optimization



Learning decision rules (LDRs)

- Find policy to minimize the expected cost
 - Infinite dimensional problem
- Linear DRs to solve newsvendor problem [Ban & Rudin, 2019]

$$\min_{\pi: \pi(\mathbf{x}) = \mathbf{q}^\top \mathbf{x}} H(\pi, \hat{\mathbb{P}}_N) + \lambda \Omega(\boldsymbol{\pi}) := \min_{\mathbf{q}} \frac{1}{N} \sum_{i=1}^N c(\mathbf{q}^\top \mathbf{x}^i, \mathbf{y}^i) + \lambda \|\mathbf{q}\|_k$$

- Linear DR have finite sample guarantees
- Linear DRs are asymptotically suboptimal in general

Decision rules on lifted space

- Linear in **transformation of features**: [Ban & Rudin, 2019]
- Policies in the reproducing kernel Hilbert space (**RKHS**) [Bertsimas & Koduri, 2023]
- **Piecewise affine** decision rules [Zhang et al., 2023]
 - Outperforms models with policy in the RKHS
- **Policy Net** [Oroojlooyjadid et al., 2020]
 - Lack interpretability
- Challenge: Ensure constraints are satisfied

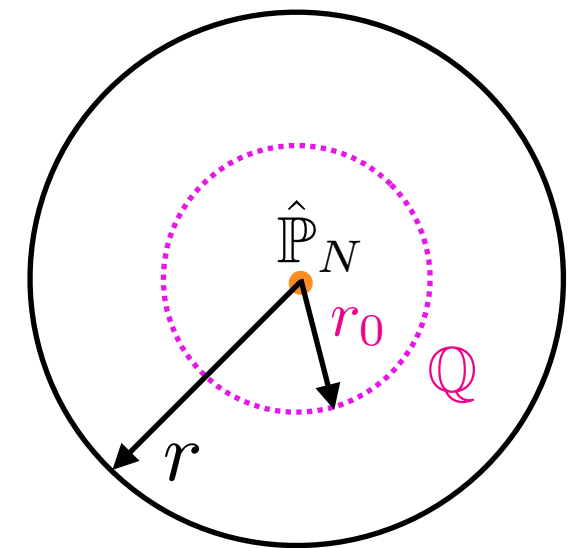
Distributionally robust optimization

- **Estimation error:** Empirical distribution biased in low data regime
- One can **robustify** against all distributions in an ambiguity set:

$$\min_{\pi \in \Pi} \sup_{Q \in \mathcal{D}} H(\pi, Q)$$

- E.g.: Wasserstein ambiguity set [Mohajerin and Kuhn 2018]

$$\mathcal{D} := \{Q \in \mathcal{P}(\mathcal{X} \times \mathcal{Y}) : \mathcal{W}(Q, \hat{\mathbb{P}}_N) \leq r\}$$



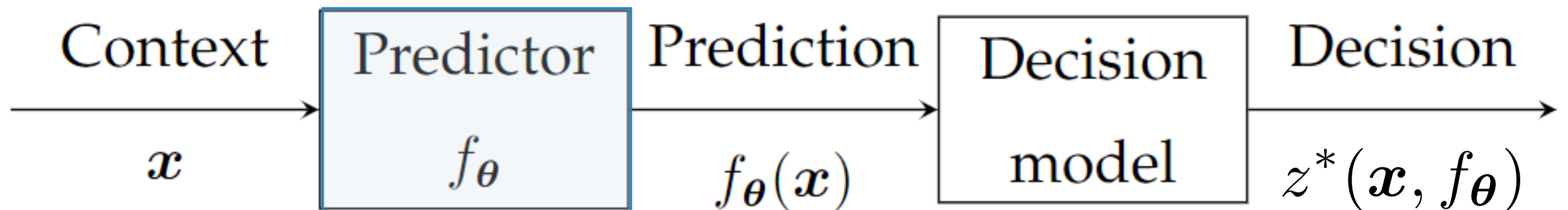
DR-Newsvendor

- Two-step procedure [Zhang et al., 2023]
 - Solve DRO problem with policy defined on historical observations of features
 - Use **Shapley extension to** interpolate to all unobserved realizations of features
- Outperforms linear decision rules, kNN, random forest, StochOptForest

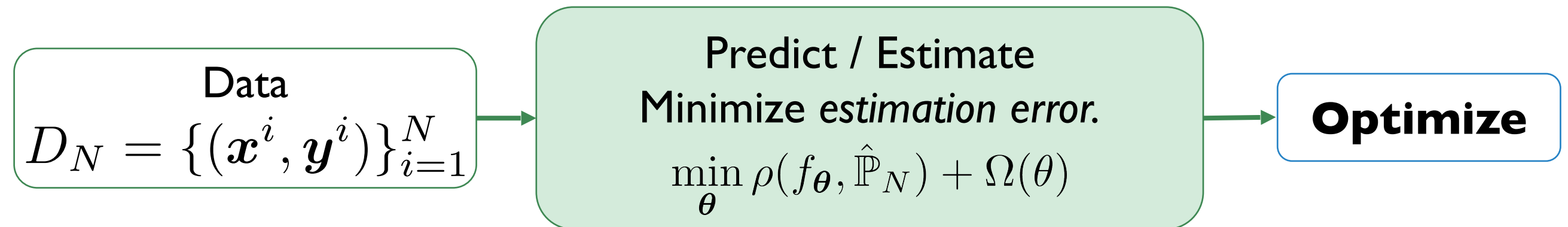
DRO with causal transport

- [Yang et al. 2023] raises issue that Wasserstein distance distorts the conditional information structure
- They suggest using a Causal transport metric, which protects causal effects found in the data
- Tractable reformulations obtained when:
 - Linear decision rules
 - Cost function is affine

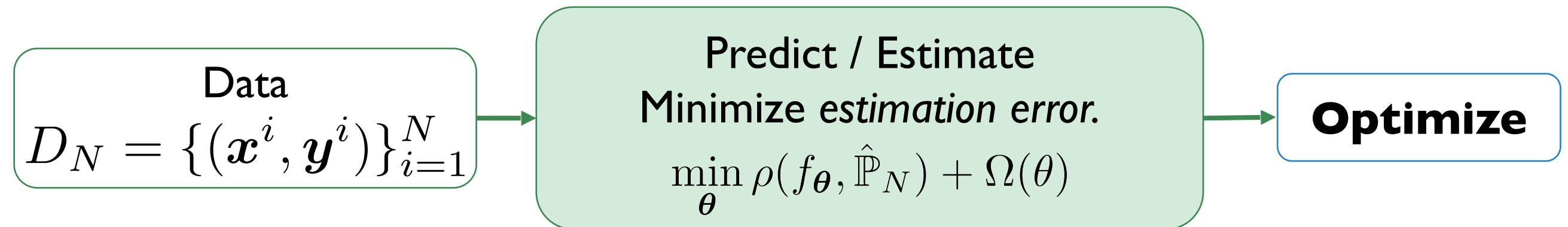
Sequential learning and optimization



Learning predictors



Learning predictors



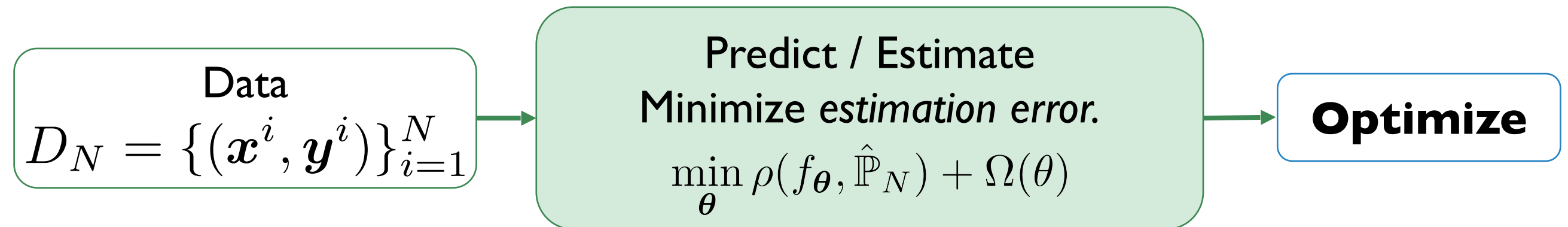
Non-linear cost function

f_{θ} is a **conditional density estimator**

Maximum Log-Likelihood

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{N} \sum_{i=1}^N -\log(\mathbb{P}_{f_{\theta}(\mathbf{x}^i)}(\mathbf{y}^i)) + \Omega(\theta)$$

Learning predictors



Non-linear cost function

f_{θ} is a **conditional density estimator**

Maximum Log-Likelihood

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{N} \sum_{i=1}^N -\log(\mathbb{P}_{f_{\theta}(x^i)}(y^i)) + \Omega(\theta)$$

Linear cost function

f_{θ} replaced with **point predictor**
(denoted g_{θ})

Mean Square Error

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{N} \sum_{i=1}^N \|g_{\theta}(x^i) - y^i\|^2 + \Omega(\theta)$$

$$h(z, f_{\theta}) = \mathbb{E}_{f_{\theta}(x)}[\mathbf{y}^{\top} \mathbf{z}] = \mathbb{E}_{f_{\theta}(x)}[\mathbf{y}]^{\top} \mathbf{z} = g_{\theta}(x)^{\top} \mathbf{z} = h(z, g_{\theta})$$

Weighted SAA

Minimizing expected costs w.r.t. a distribution is often done through SAA:

$$\min_{\mathbf{z} \in \mathcal{Z}} \mathbb{E}_{f_{\theta}(\mathbf{x})} [c(\mathbf{z}, \mathbf{y})] \text{ with } f_{\theta}(\mathbf{x}) := \frac{1}{N} \sum_{i=1}^N \delta_{\mathbf{y}^i(\mathbf{x})}$$

Weighted SAA

Minimizing expected costs w.r.t. a distribution is often done through SAA:

$$\min_{\mathbf{z} \in \mathcal{Z}} \mathbb{E}_{f_{\theta}(\mathbf{x})} [c(\mathbf{z}, \mathbf{y})] \text{ with } f_{\theta}(\mathbf{x}) := \frac{1}{N} \sum_{i=1}^N \delta_{\mathbf{y}^i}(\mathbf{x})$$

Residual based

Measure the **error of a trained regression model** on the historical data

$$f_{\theta}(\mathbf{x}) := \frac{1}{N} \sum_{i=1}^N \delta_{g_{\theta}(\mathbf{x}) + \epsilon_i}$$

Weighted SAA

Minimizing expected costs w.r.t. a distribution is often done through SAA:

$$\min_{z \in \mathcal{Z}} \mathbb{E}_{f_\theta(\mathbf{x})} [c(\mathbf{z}, \mathbf{y})] \text{ with } f_\theta(\mathbf{x}) := \sum_{i=1}^N \delta_{\mathbf{y}^i} \cdot w_i(\mathbf{x})$$

Residual based

Measure the **error of a trained regression model** on the historical data

$$f_\theta(\mathbf{x}) := \frac{1}{N} \sum_{i=1}^N \delta_{g_\theta(\mathbf{x}) + \epsilon^i}$$

Weight based

Measure **proximity in feature space** between \mathbf{x} and historical covariates \mathbf{x}^i

Weighted SAA

Proximity in feature space

- k -nearest neighbor: $w_i^{\text{kNN}}(\mathbf{x}) := (1/k) \mathbb{1}[\mathbf{x}^i \in \mathcal{N}_k(\mathbf{x})]$
- Kernel density estimation: $w_i^{\text{KDE}}(\mathbf{x}) := \frac{\mathcal{K}(\mathbf{x}, \mathbf{x}^i)}{\sum_{j=1}^N \mathcal{K}(\mathbf{x}, \mathbf{x}^j)}$

Weighted SAA

Proximity in feature space

- k -nearest neighbor: $w_i^{\text{kNN}}(\mathbf{x}) := (1/k) \mathbb{1}[\mathbf{x}^i \in \mathcal{N}_k(\mathbf{x})]$
- Kernel density estimation: $w_i^{\text{KDE}}(\mathbf{x}) := \frac{\mathcal{K}(\mathbf{x}, \mathbf{x}^i)}{\sum_{j=1}^N \mathcal{K}(\mathbf{x}, \mathbf{x}^j)}$

Supervised learning

- Decision tree: $w_i^{\text{DT}}(\mathbf{x}) := \frac{\mathbb{1}[\mathcal{R}(\mathbf{x}) = \mathcal{R}(\mathbf{x}^i)]}{\sum_{j=1}^N \mathbb{1}[\mathcal{R}(\mathbf{x}) = \mathcal{R}(\mathbf{x}^j)]}$
- Random forest: average over set of decision trees.

Why do sequential learning and optimization?

It's fast!

- Train once on historical data:
no need to solve optimization models during training

It works

- Can perform better than non-contextual approach
- Can be trained using less data when model is well specified

Theoretical guarantees

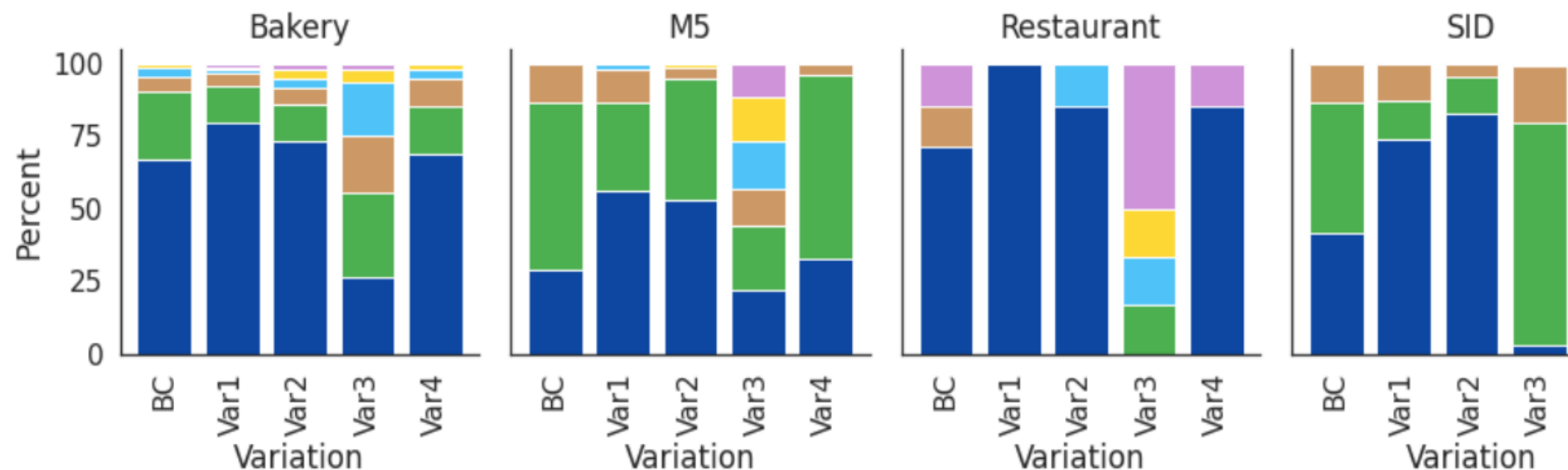
- Converges to optimal contextual policy as the size of the training set increases **when model is well specified.**

Some benchmark results (Buttler et al., 2023)

Newsvendor Problem

Compare **sequential** L&O and **decision rules** on 4 data sets.

Proportion of instances where methods achieved best performance



Models:

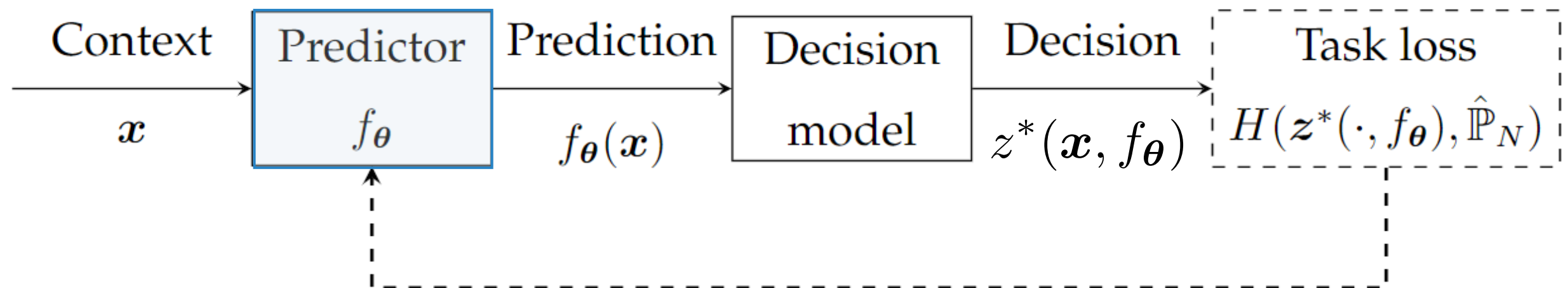
- Linear rule
- Kernel weights
- Decision tree weights
- Deep learning
- K-nearest neighbour weights
- Random forest weights

Sequential Learning & Optimization References

	Method			Regularization		Learning model					
	rCSO	wSAA	EVB	Reg. CSO	DRO	General	Linear	Kernel	kNN	DT	RF
Hannah et al. (2010)	X	✓	X	X	X	X	X	✓	X	X	X
Ferreira et al. (2016)	X	X	✓	X	X	X	X	X	X	✓	X
Ban et al. (2019)	✓	X	X	X	X	X	✓	X	X	X	X
Chen and Paschalidis (2019)	X	✓	X	X	✓	X	X	X	✓	X	X
Bertsimas and Kallus (2020)	X	✓	X	X	X	X	✓	X	✓	✓	✓
Kannan et al. (2020)	✓	X	X	X	X	✓	✓	✓	✓	✓	✓
Kannan et al. (2021)	✓	X	X	X	✓	✓	✓	✓	✓	✓	✓
Liu et al. (2021)	X	X	✓	X	X	X	✓	X	X	✓	X
Srivastava et al. (2021)	X	✓	X	✓	X	X	X	✓	X	X	X
Wang et al. (2021)	X	✓	X	X	✓	X	X	✓	X	X	X
Bertsimas and Van Parys (2022)	X	✓	X	X	✓	X	X	✓	✓	X	X
Deng and Sen (2022)	✓	X	X	X	X	✓	✓	✓	✓	✓	✓
Esteban-Pérez and Morales (2022)	X	✓	X	X	✓	X	X	✓	✓	X	X
Kannan et al. (2022)	✓	X	X	X	✓	✓	✓	✓	✓	✓	✓
Lin et al. (2022)	X	✓	X	✓	X	X	X	X	✓	✓	✓
Nguyen et al. (2021)	X	✓	X	X	✓	X	X	X	✓	X	X
Notz and Pibernik (2022)	X	✓	X	X	X	X	X	✓	X	✓	X
Zhu et al. (2022)	X	X	✓	X	✓	✓	✓	✓	✓	✓	✓
Perakis et al. (2023)	✓	X	X	X	✓	X	✓	X	X	X	X

Going beyond SLO: Integrated learning and optimization

Going beyond SLO: Integrated learning and optimization



Wrong predictions lead to suboptimal decisions

$$\max_{z \in \mathcal{Z}} \mathbf{y}^\top \mathbf{z}$$

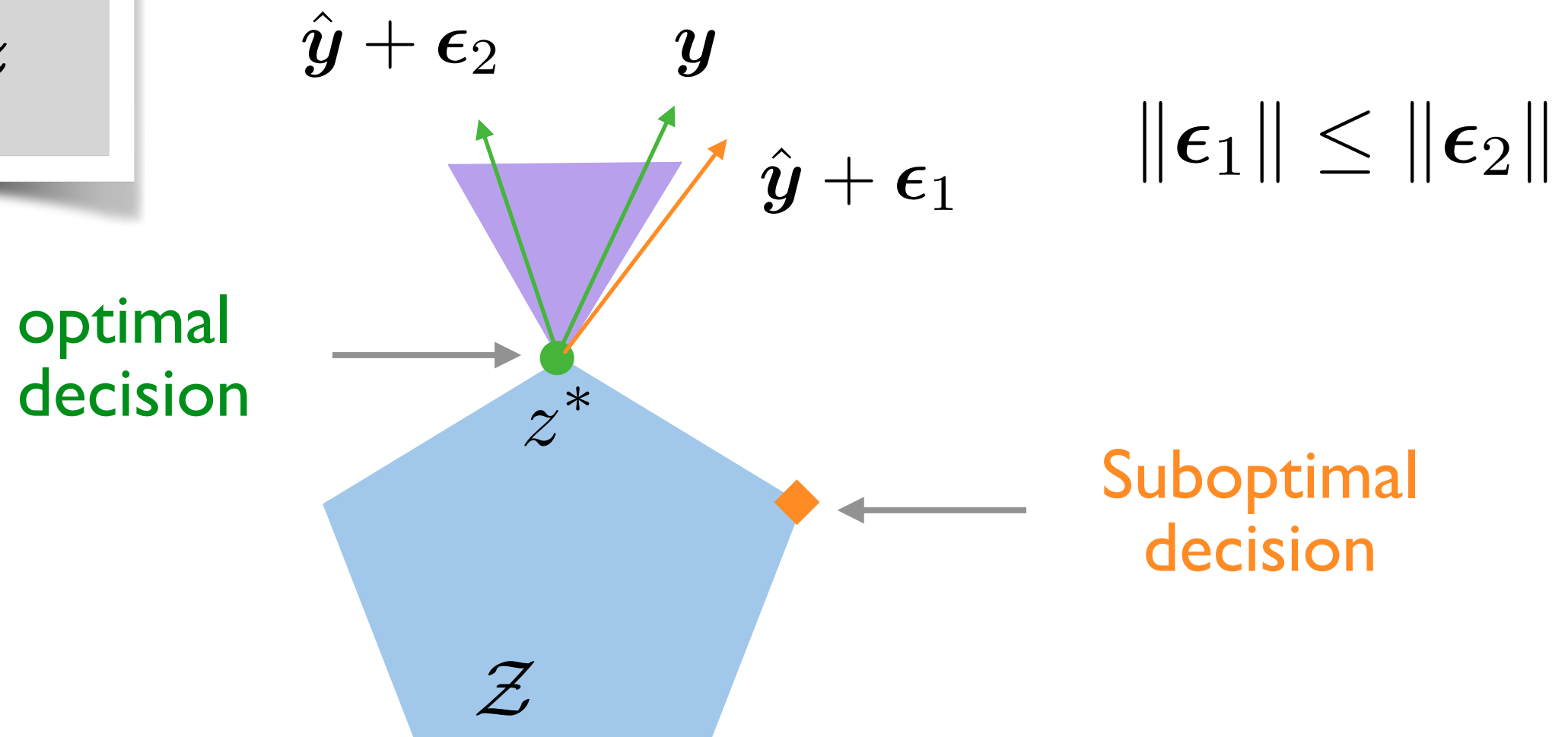
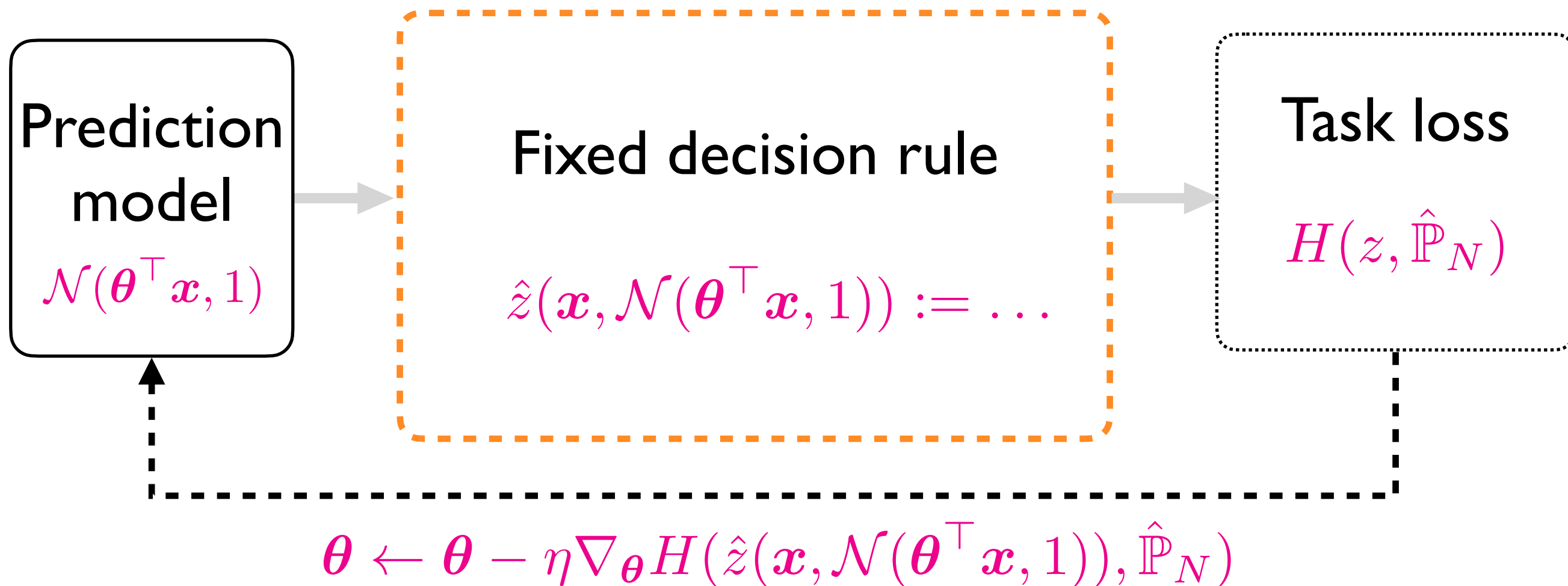


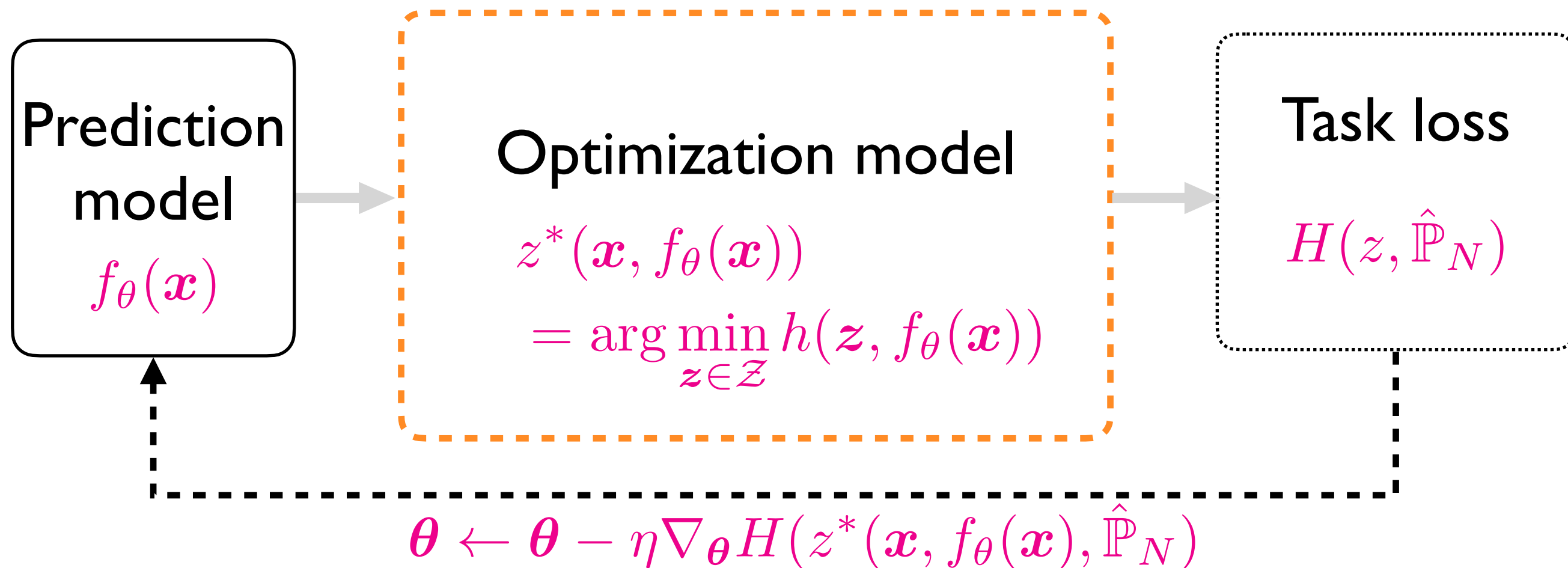
Figure adapted from [Elmachtoub and Grigas 2022]

ILO training pipeline



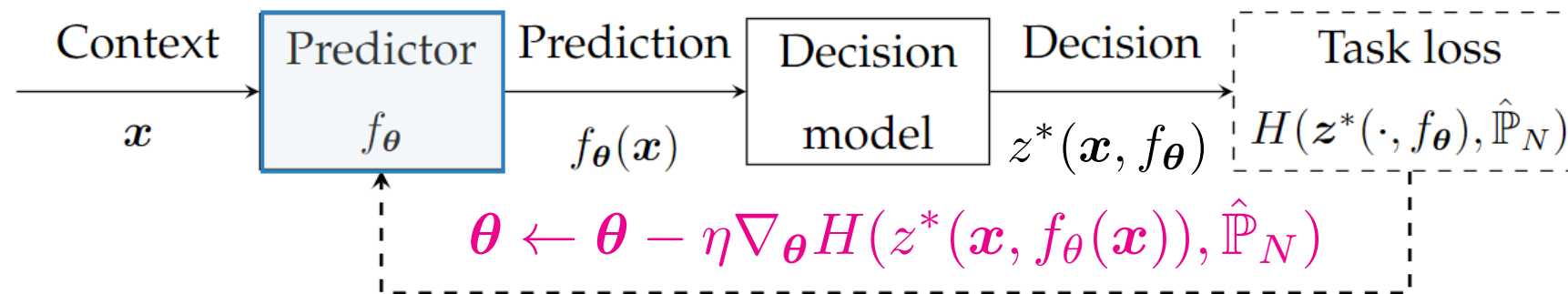
- [Bengio 1997] : Task-aware point prediction under a **fixed decision rule**

ILO training pipeline



- [Bengio 1997] : Task-aware point prediction under a **fixed decision rule**
- [Donti et al. 2017] : Task-aware conditional density prediction under **CSO model**

How to differentiate through argmin operation?



- Implicit differentiation through KKT conditions for convex problems
- Unroll the operations made by the optimization process:
 - Differentiate through its computational graph
 - Implicit differentiation of the fixed point equations at local optimum [Butler and Kwon, 2023] and [Kotary et al. 2023]
- Replace optimizer with a differentiable deep neural network [Grigas et al. 2021]
- Libraries: TorchOpt [Bilevel], CvxpyLayer [Convex], PyEPO [Linear]

Smart “Predict, then optimize”

- Regret minimization [Elmachtoub & Grigas, 2022]:

$$H(z^*(\mathbf{x}, f_\theta), \mathbb{P}) := \mathbb{E}_{\mathbb{P}}[c(z^*(\mathbf{x}, f_\theta), \mathbf{y})]$$

Smart “Predict, then optimize”

- Regret minimization [Elmachtoub & Grigas, 2022]:

$$H(z^*(\mathbf{x}, f_\theta), \mathbb{P}) := \mathbb{E}_{\mathbb{P}}[c(z^*(\mathbf{x}, f_\theta), \mathbf{y})] - \mathbb{E}_{\mathbb{P}}[c(z^*(\mathbf{x}, f_\theta), \mathbf{y}) - \min_{z \in \mathcal{Z}} c(z, \mathbf{y})]$$

Smart “Predict, then optimize”

- Regret minimization [Elmachtoub & Grigas, 2022]:

$$H(z^*(\mathbf{x}, f_\theta), \mathbb{P}) := \mathbb{E}_{\mathbb{P}}[c(z^*(\mathbf{x}, f_\theta), \mathbf{y})] - \mathbb{E}_{\mathbb{P}}[c(z^*(\mathbf{x}, f_\theta), \mathbf{y}) - \min_{\mathbf{z} \in \mathcal{Z}} c(\mathbf{z}, \mathbf{y})]$$

- Non-convex and discontinuous in θ
- Replace with SPO+:

$$\min_{\theta} \mathbb{E}_{\mathbb{P}} [\ell_{\text{SPO}+}(g_{\theta}(\mathbf{x}), \mathbf{y})]$$

where

$$\ell_{\text{SPO}+}(\hat{\mathbf{y}}, \mathbf{y}) := \sup_{\mathbf{z} \in \mathcal{Z}} (\mathbf{y} - 2\hat{\mathbf{y}})^T \mathbf{z} + 2\hat{\mathbf{y}}^T \mathbf{z}^*(\mathbf{x}, \mathbf{y}) - \mathbf{y}^T \mathbf{z}^*(\mathbf{x}, \mathbf{y}),$$

Smart “Predict, then optimize”

- Regret minimization [Elmachetoub & Grigas, 2022]:

$$H(z^*(\mathbf{x}, f_\theta), \mathbb{P}) := \mathbb{E}_{\mathbb{P}}[c(z^*(\mathbf{x}, f_\theta), \mathbf{y})] - \mathbb{E}_{\mathbb{P}}[c(z^*(\mathbf{x}, f_\theta), \mathbf{y}) - \min_{\mathbf{z} \in \mathcal{Z}} c(\mathbf{z}, \mathbf{y})]$$

- Non-convex and discontinuous in θ

- Replace with SPO+:

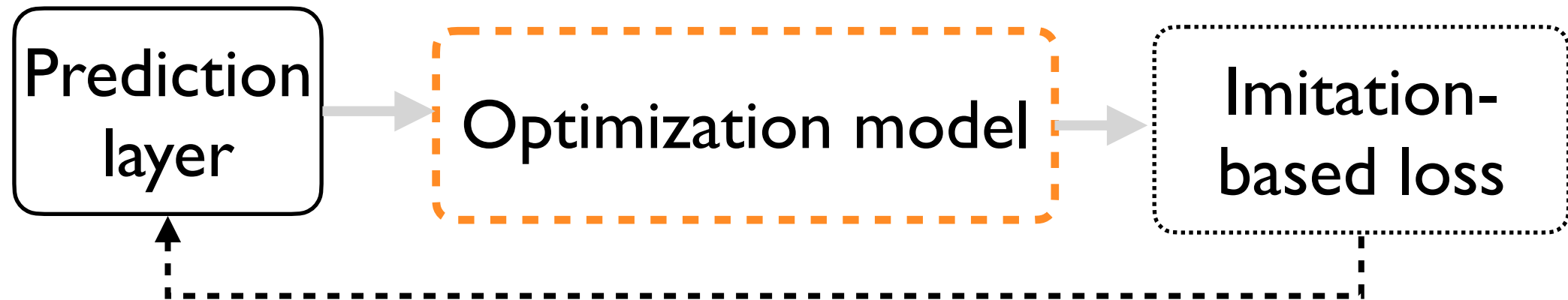
$$\min_{\theta} \mathbb{E}_{\mathbb{P}} [\ell_{\text{SPO}+}(g_{\theta}(\mathbf{x}), \mathbf{y})]$$

where

$$\ell_{\text{SPO}+}(\hat{\mathbf{y}}, \mathbf{y}) := \sup_{\mathbf{z} \in \mathcal{Z}} (\mathbf{y} - 2\hat{\mathbf{y}})^T \mathbf{z} + 2\hat{\mathbf{y}}^T \mathbf{z}^*(\mathbf{x}, \mathbf{y}) - \mathbf{y}^T \mathbf{z}^*(\mathbf{x}, \mathbf{y}),$$

- Solve an optimization problem at each iteration
- SPO+ has slower convergence rate than SLO approach
- If model misspecified, SPO+ can outperform SLO

Optimal action imitation



- Imitation performance metric:

$$H(z^*(\mathbf{x}, f_\theta), \mathbb{P}) := \mathbb{E}_{\mathbb{P}}[c(z^*(\mathbf{x}, f_\theta), \mathbf{y})]$$

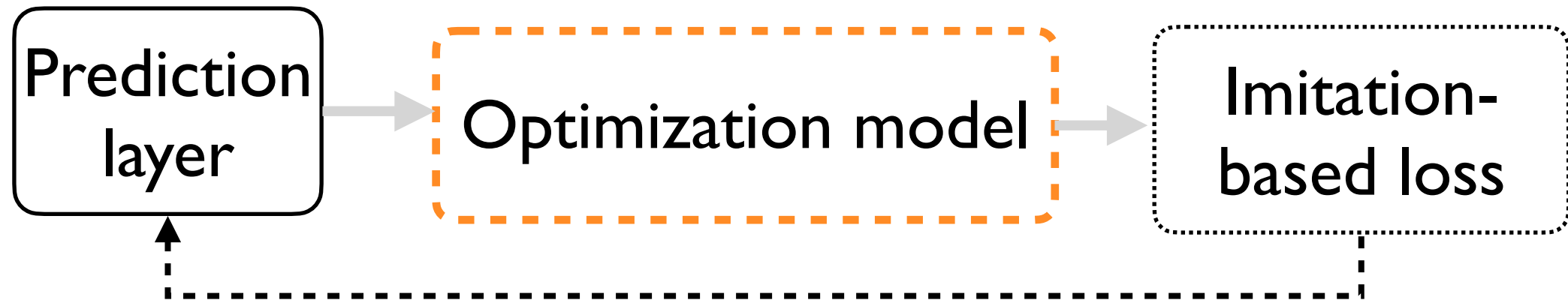
Optimal action imitation



- Imitation performance metric:

$$H(z^*(\mathbf{x}, f_\theta), \mathbb{P}) := \mathbb{E}_{\mathbb{P}}[c(z^*(\mathbf{x}, f_\theta), \mathbf{y})] \quad \mathbb{E}_{\hat{\mathbb{P}}_N}[d(z^*(\mathbf{x}, f_\theta), z^*(\mathbf{x}, \mathbf{y}))]$$

Optimal action imitation



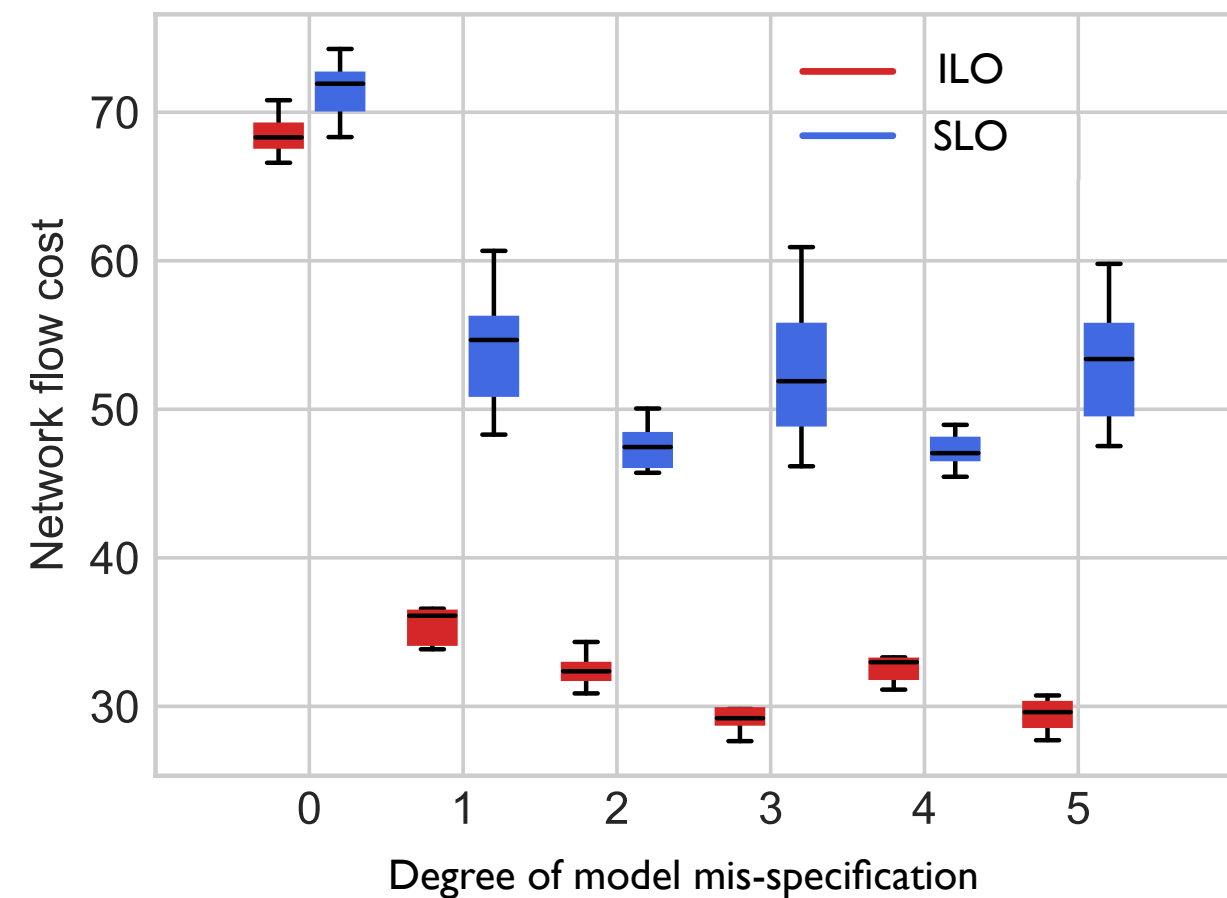
- Imitation performance metric:

$$H(z^*(\mathbf{x}, f_\theta), \mathbb{P}) := \mathbb{E}_{\mathbb{P}}[c(z^*(\mathbf{x}, f_\theta), \mathbf{y})] \quad \mathbb{E}_{\hat{\mathbb{P}}_N}[d(z^*(\mathbf{x}, f_\theta), z^*(\mathbf{x}, \mathbf{y}))]$$

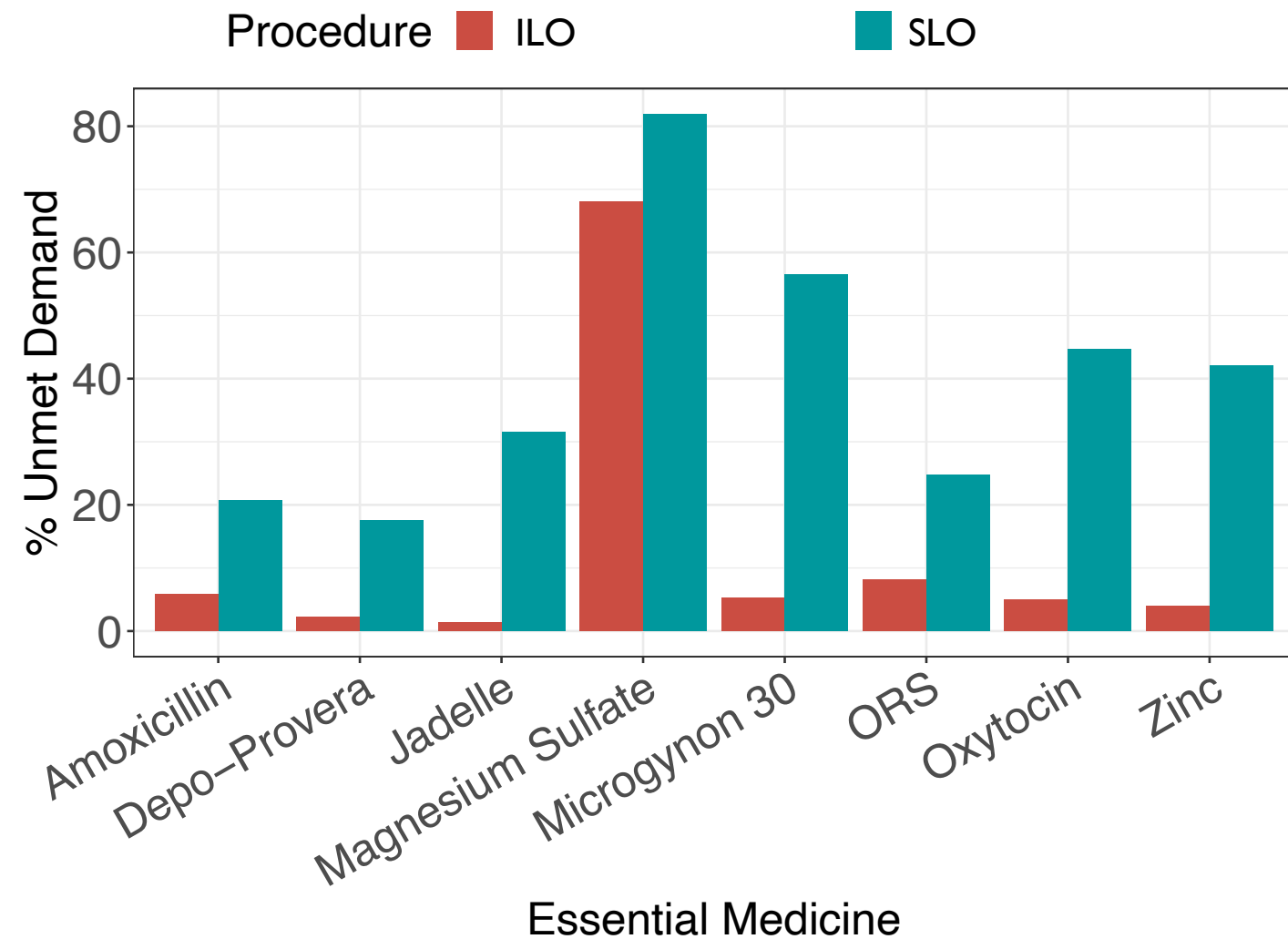
- Training based on perturbed optimizers:
 - [Berthet et al., 2020] uses additive perturbation of point prediction
 - [Dalle et al., 2022] uses multiplicative perturbations
 - [Mulamba et al., 2021] and [Kong et al., 2022] uses energy-based optimizer

$$\tilde{z}(\mathbf{x}, f_\theta) \sim \frac{\exp(-\alpha h(\mathbf{z}, f_\theta(\mathbf{x})))}{\int \exp(-\alpha h(\mathbf{z}, f_\theta(\mathbf{x}))) d\mathbf{z}}$$

ILO outperforms SLO



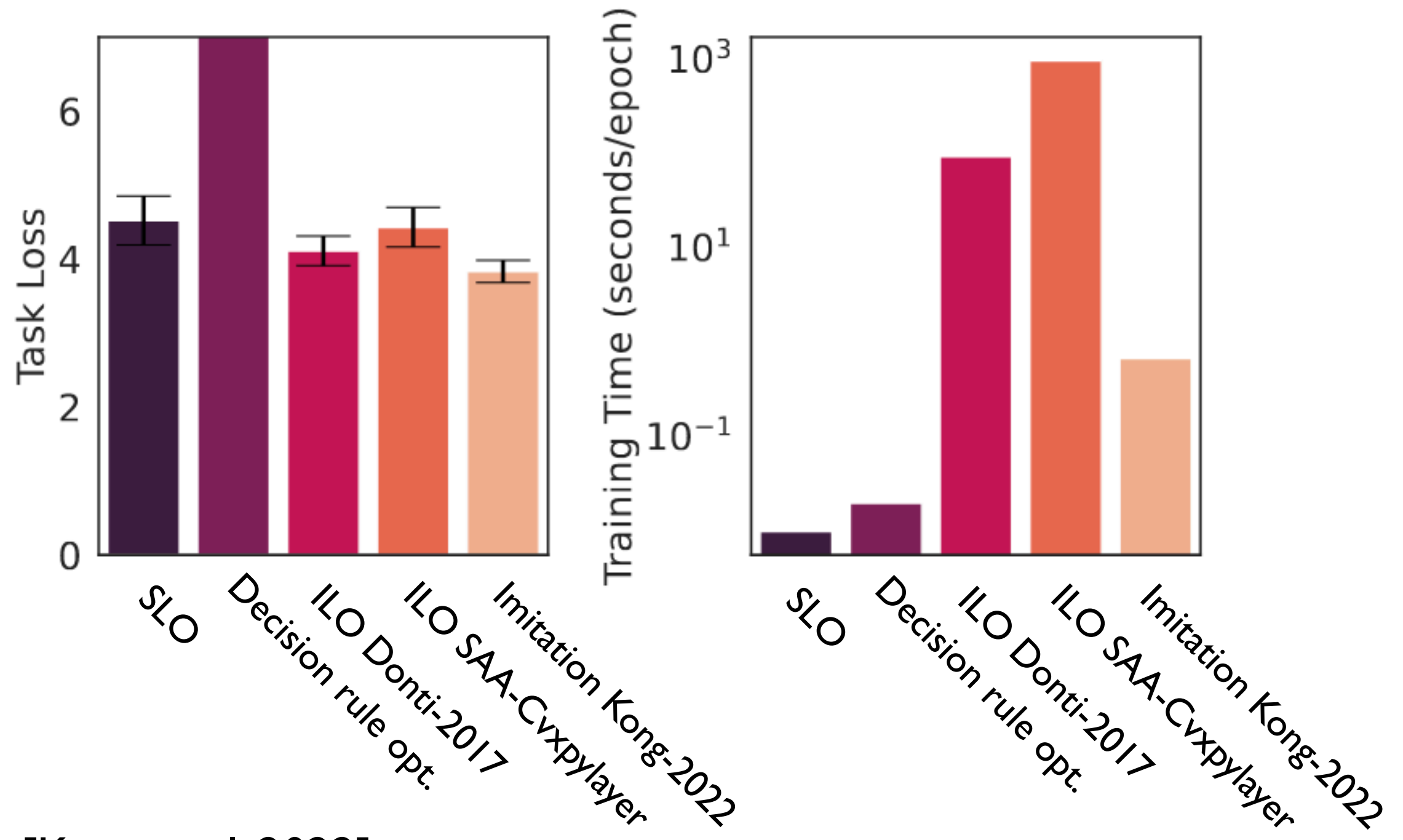
Source: [Grigas et al. 2021]



Source: [Chung et al. 2022]

Comparison of different approaches

Load forecasting and generator scheduling problem



Source: [Kong et al. 2022]

Take-away messages

- Contextual stochastic optimization is a rapidly evolving field that provides methods for identifying data-driven decision that exploit most recently available information.
- Three types of approaches:
 - Decision rule/policy optimization
 - Sequential learning and optimization
 - Integrated learning and optimization
- Four types of performance measures:
 - Statistical accuracy of prediction model
 - Task-based expected cost of induced policy
 - Task-based expected regret of induced policy
 - Quality of imitation
- Many potential applications ?



(Link to survey paper)

References

- Ban GY, Rudin C (2019) The Big Data Newsvendor: Practical Insights from Machine Learning. *Operations Research* 67(1):90–108
- Bertsimas D, Koduri N (2022) Data-driven optimization: A reproducing kernel Hilbert space approach. *Operations Research* 70(1):454–471.
- Buttler S, Philippi A, Stein N, Pibernik R (2022) A meta analysis of data-driven newsvendor approaches. *ICLR 2022 Workshop on Setting up ML Evaluation Standards to Accelerate Progress*
- Chung TH, Rostami V, Bastani H, Bastani O (2022) Decision-aware learning for optimizing health supply chains. *arXiv preprint arXiv:2211.08507*.
- Elmachtoub AN, Liang JCN, McNellis R (2020) Decision trees for decision-making under the predict-then-optimize framework. *ICML*, 2858–2867
- Grigas P, Qi M, Shen M (2021) Integrated conditional estimation-optimization. *arXiv preprint arXiv:2110.12351*
- Kallus N, Mao X (2022) Stochastic optimization forests. *Management Science* 69(4):1975–1994.
- Kong L, Cui J, Zhuang Y, Feng R, Prakash BA, Zhang C (2022) End-to-end stochastic optimization with energy-based model. *Advances in Neural Information Processing Systems*, volume 35, 11341–11354 (Curran Associates, Inc.).
- Mohajerin Esfahani P, Kuhn D (2018) Data-driven distributionally robust optimization using the wasserstein metric: Performance guarantees and tractable reformulations. *Mathematical Programming* 171(1- 2):115–166.
- Oroojlooyjadid A, Snyder LV, Takáč M (2020) Applying deep learning to the newsvendor problem. *IIE Transactions* 52(4):444–463.
- Shah S, Wang K, Wilder B, Perrault A, Tambe M (2022) Decision-focused learning without decision-making: Learning locally optimized decision losses. *NeurIPS*.
- Yang J, Zhang L, Chen N, Gao R, Hu M (2023) Decision-making with side information: A causal transport robust approach.
- Zhang L, Yang J, Gao R (2023) Optimal robust policy for feature-based newsvendor. *Management Science* (Forthcoming)
- Zhang Y, Liu J, Zhao X (2023b) Data-driven piecewise affine decision rules for stochastic programming with covariate information. *arXiv:2304.13646*.