

Introduction to Data-driven 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)*

June 23, 2023



Canada
Research
Chairs

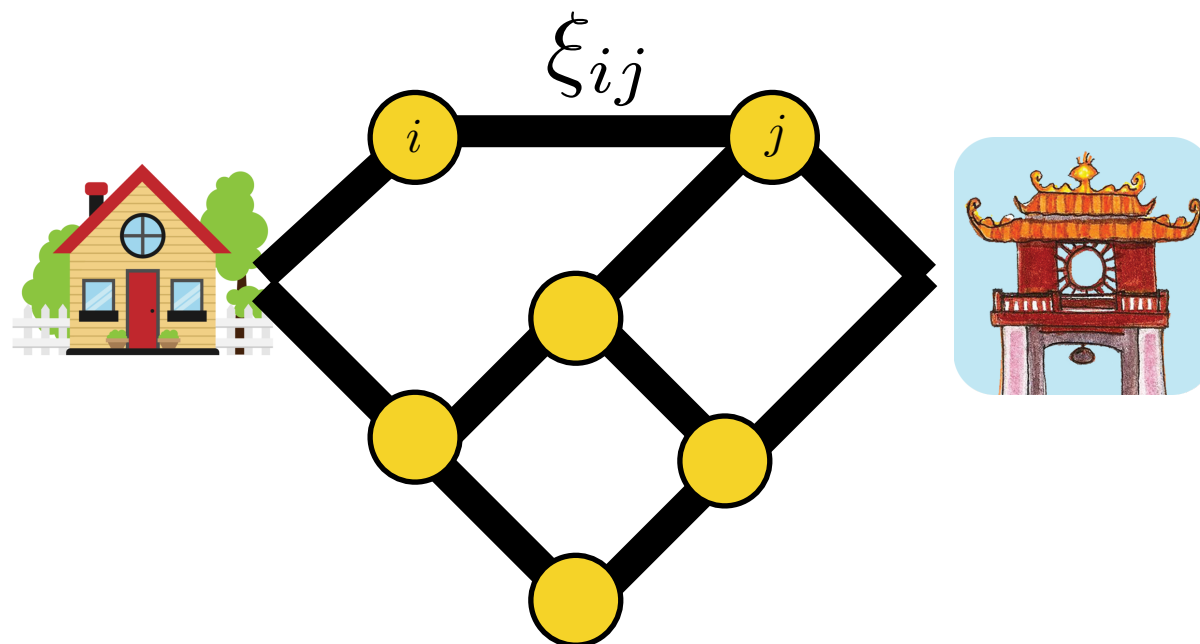
Chaires
de recherche
du Canada

Canada

**Why contextual stochastic
optimization?**

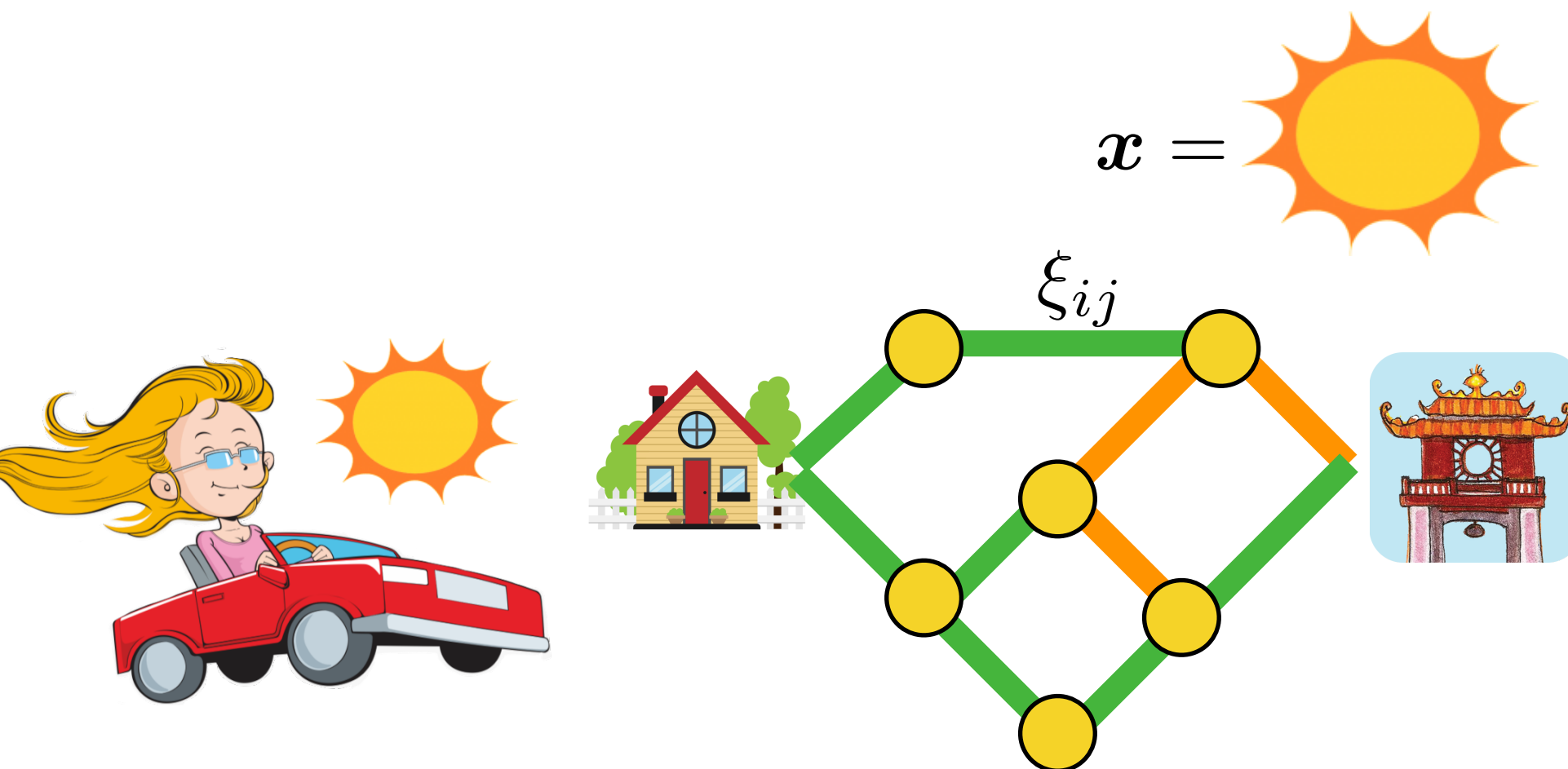
Decision Making with Contextual Information

- Revealed contextual information x
- Hidden random variables ξ



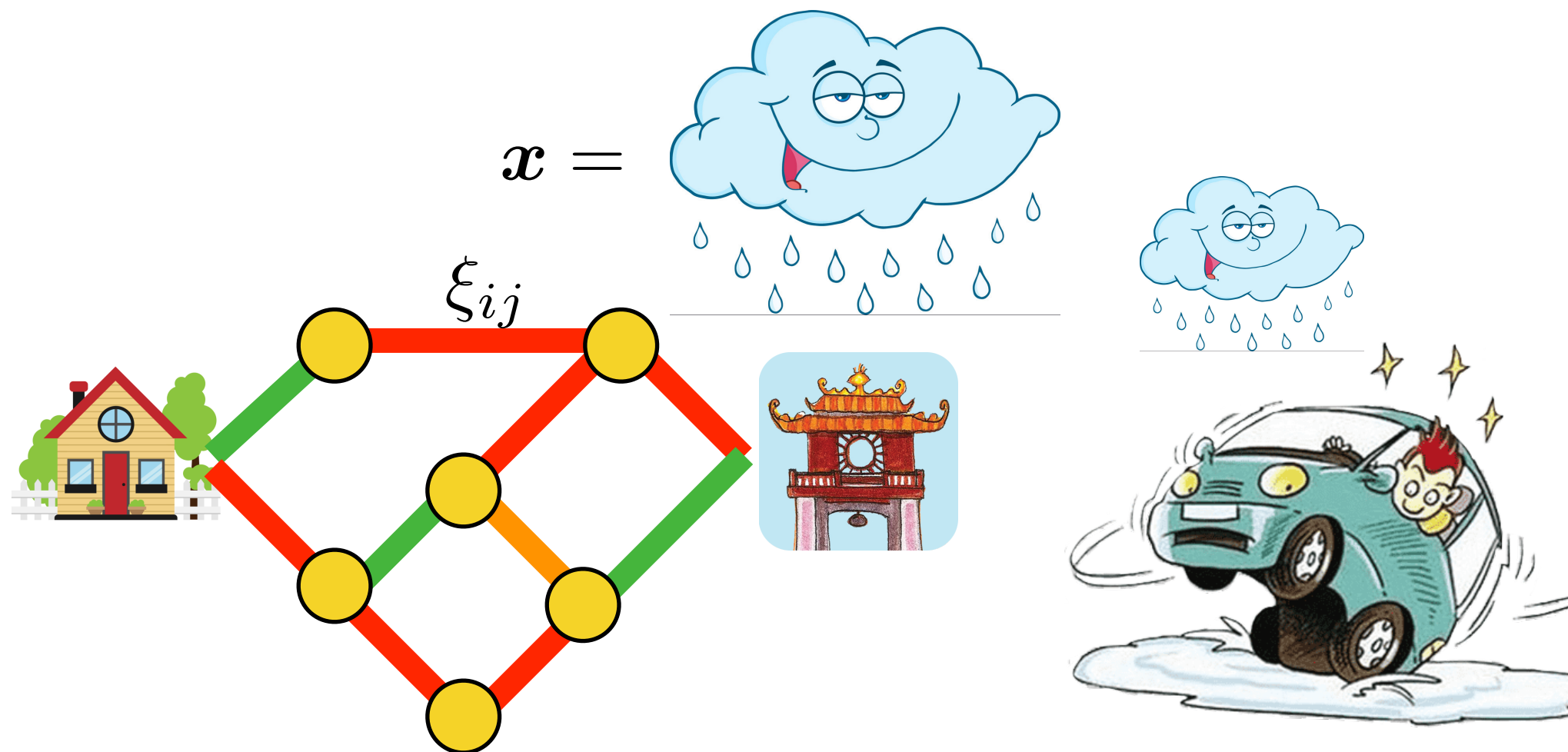
Decision Making with Contextual Information

- Revealed contextual information x
- Hidden random variables ξ



Decision Making with Contextual Information

- Revealed contextual information x
- Hidden random variables ξ

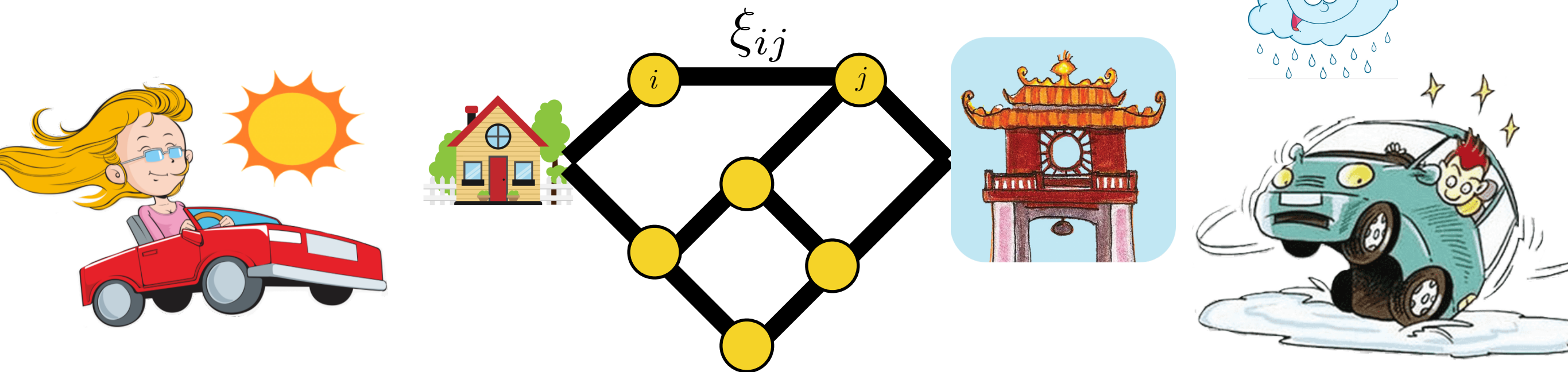


Decision Making with Contextual Information

- Revealed contextual information x
- Hidden random variables ξ

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

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

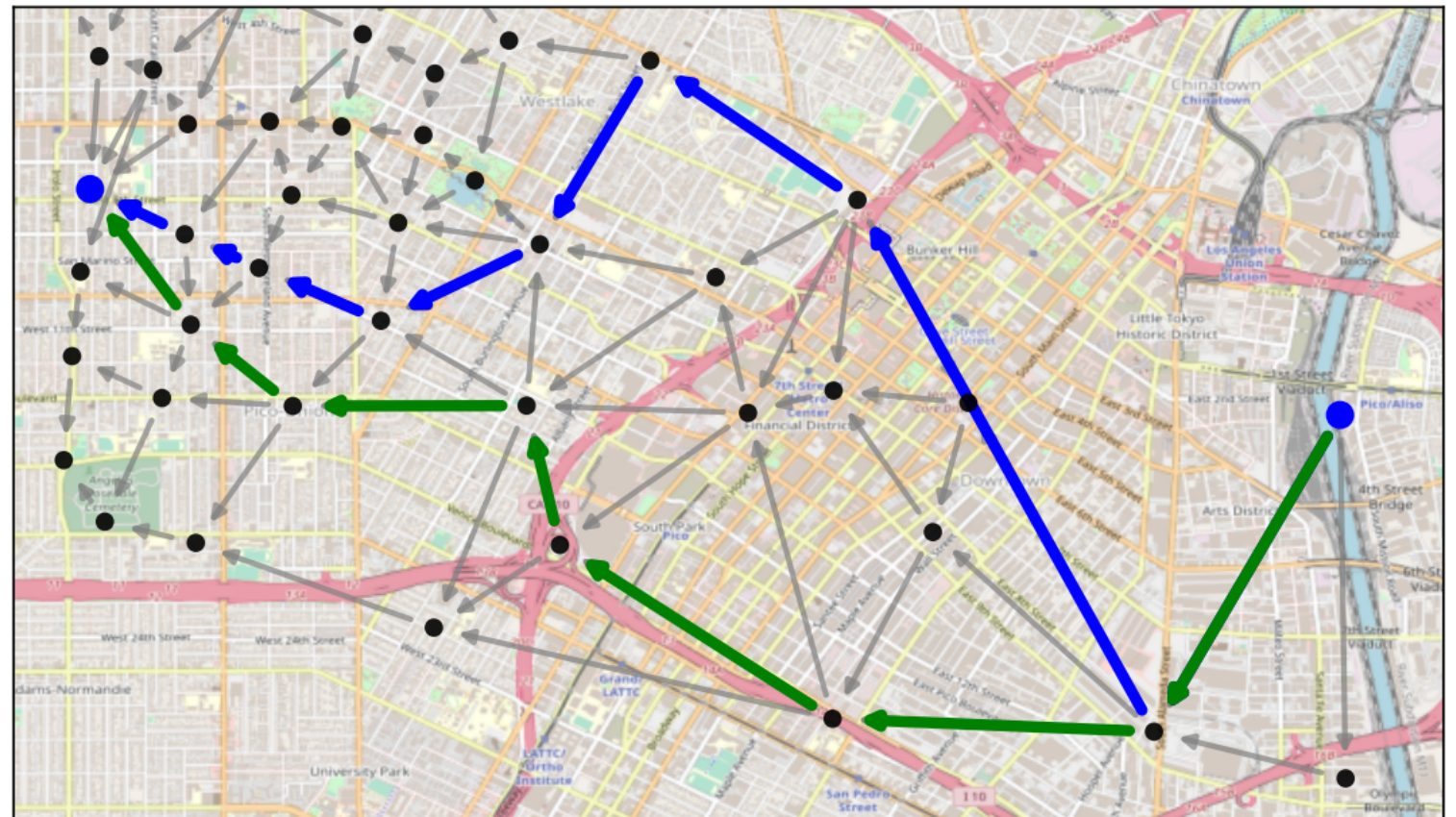


Practical motivation

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.

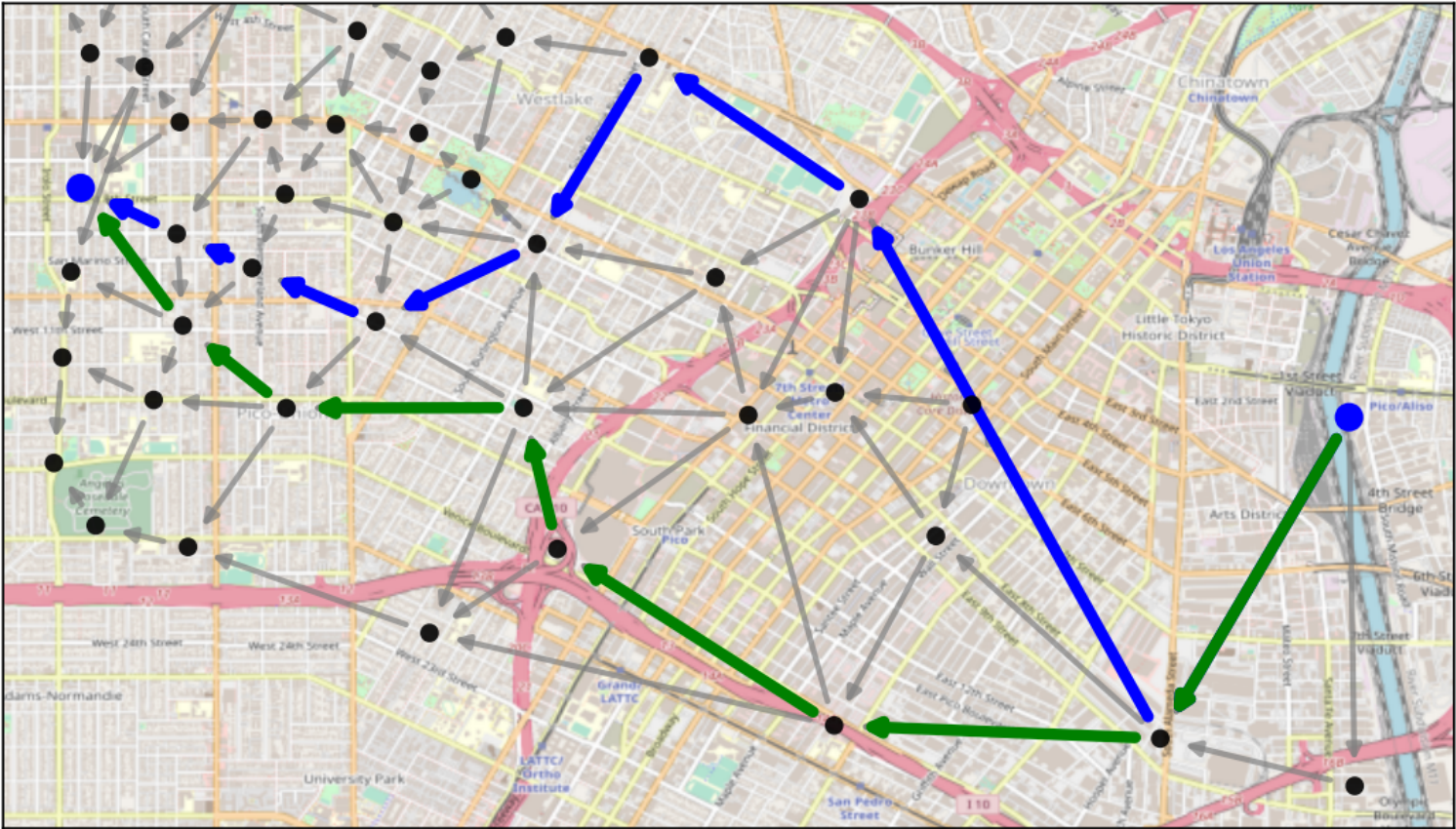


Practical motivation

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

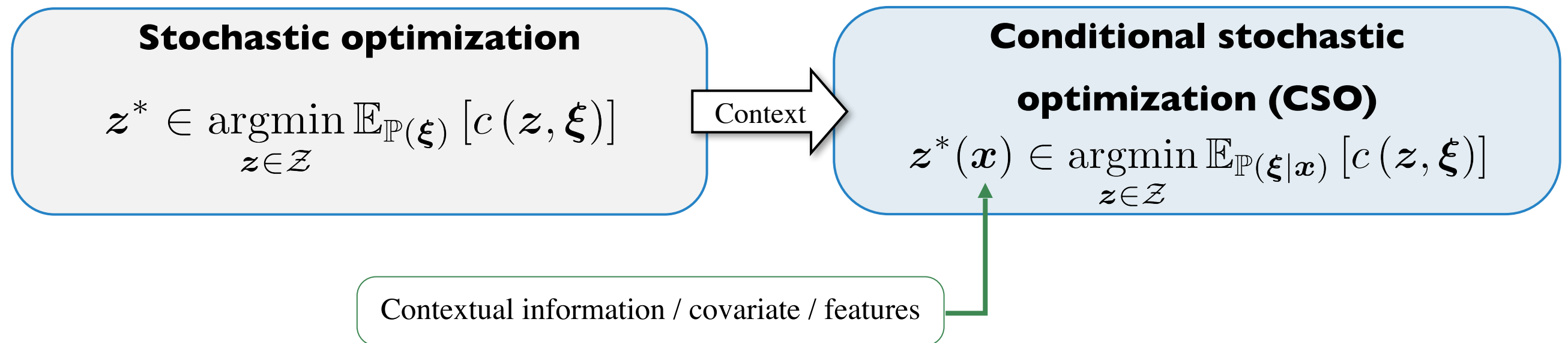
**What is contextual
optimization?**

Problem Definition

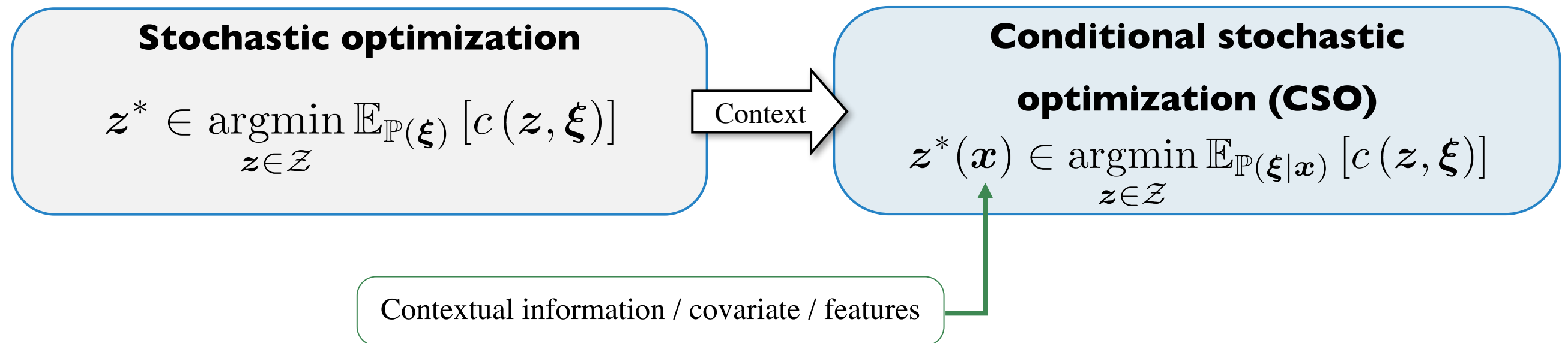
Stochastic optimization

$$z^* \in \operatorname{argmin}_{z \in \mathcal{Z}} \mathbb{E}_{\mathbb{P}(\xi)} [c(z, \xi)]$$

Problem Definition



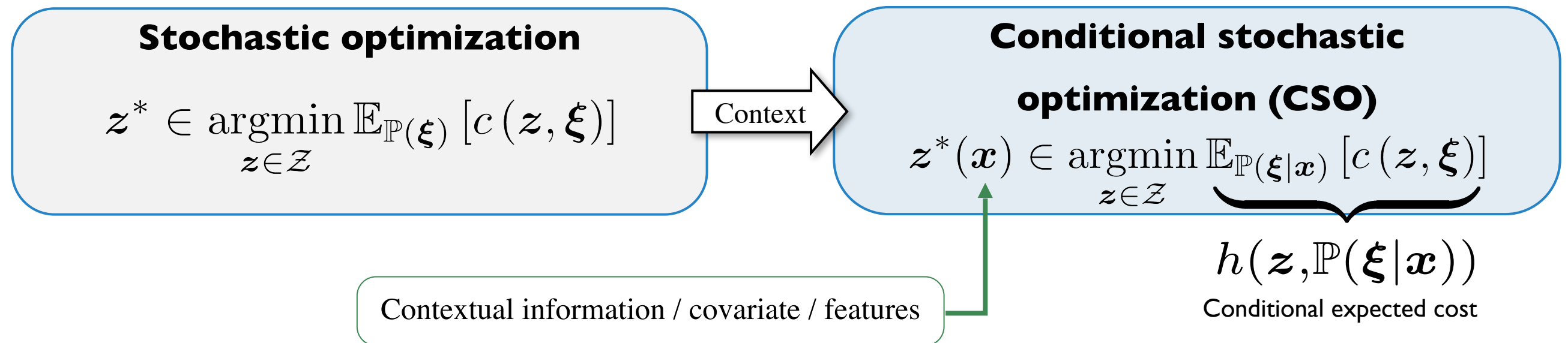
Problem Definition



Connection between CSO and policy optimization:

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

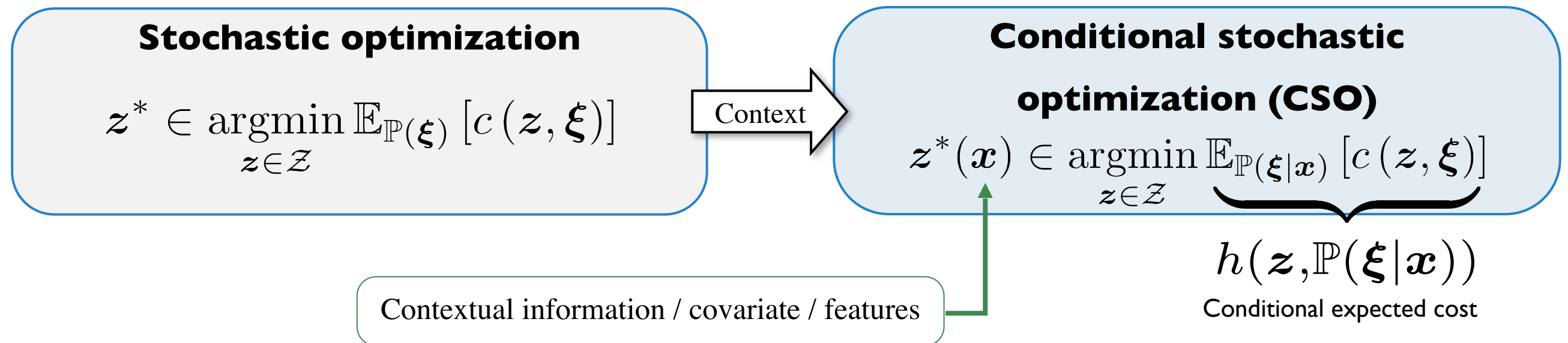
Problem Definition



Connection between CSO and policy optimization:

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

Problem Definition



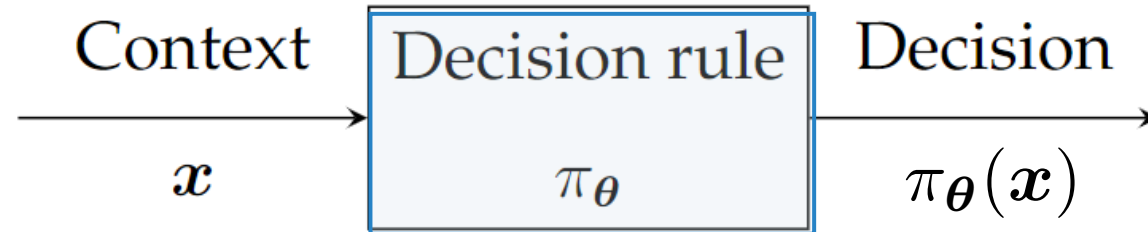
Connection between CSO and policy optimization:

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

(Unconditional) expected cost

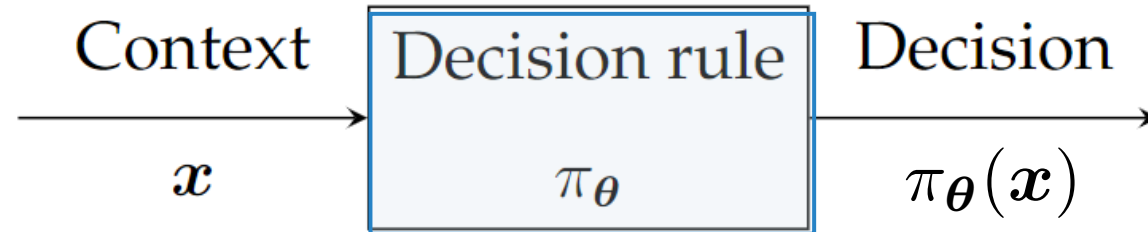
Overview of the three frameworks

Decision rule/Policy optimization

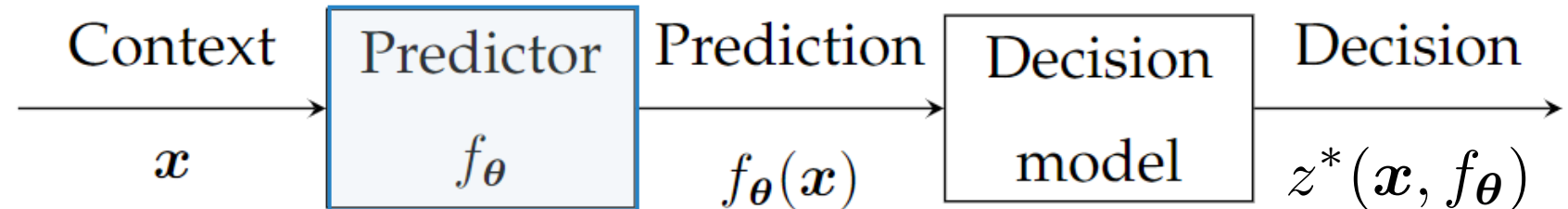


Overview of the three frameworks

Decision rule/Policy optimization

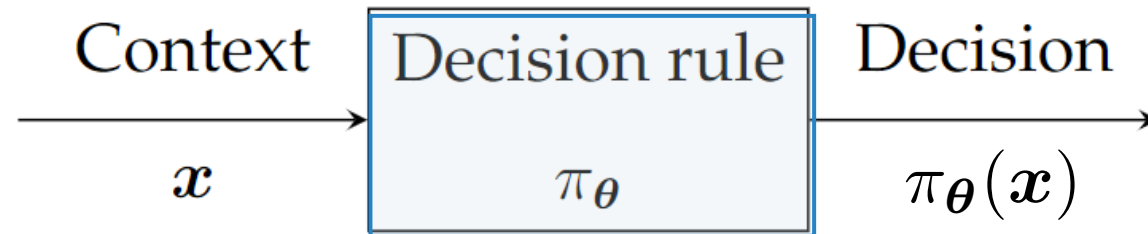


Sequential learning and optimization

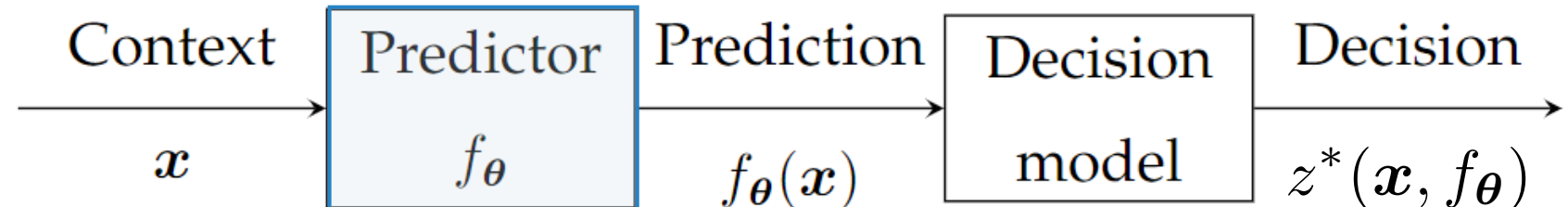


Overview of the three frameworks

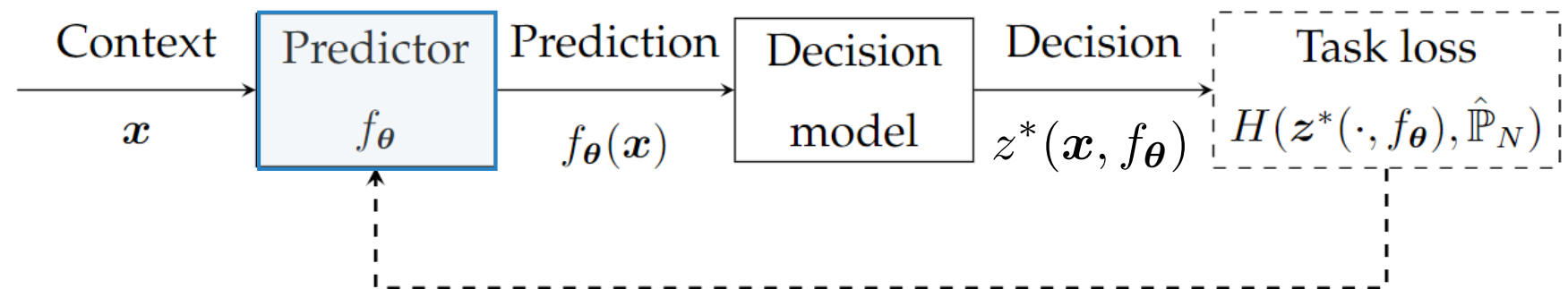
Decision rule/Policy optimization



Sequential learning and optimization



Integrated learning and optimization



Sequential learning and optimization

Learning predictors

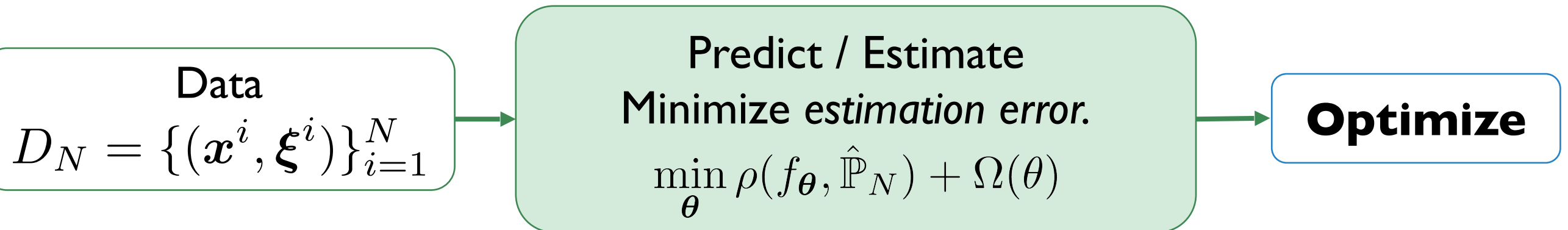
Data

$$D_N = \{(x^i, \xi^i)\}_{i=1}^N$$

Predict / Estimate
Minimize *estimation error*.
$$\min_{\theta} \rho(f_{\theta}, \hat{\mathbb{P}}_N) + \Omega(\theta)$$

Optimize

Learning predictors



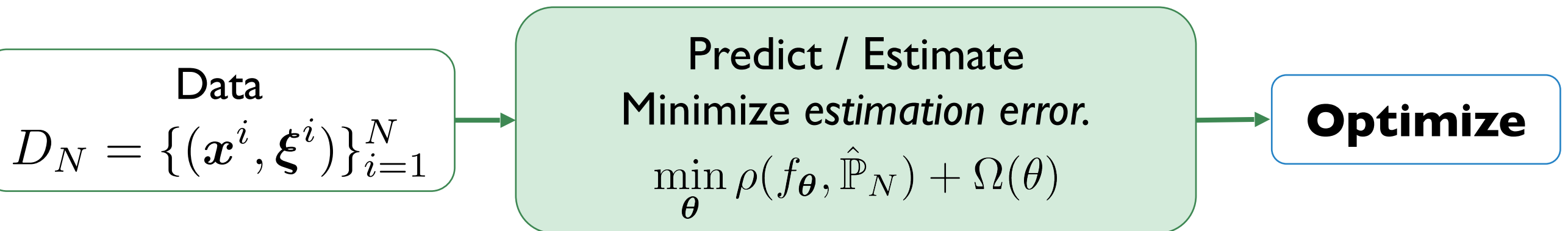
Non-linear cost function

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

Maximum Likelihood

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

Learning predictors



Non-linear cost function

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

Maximum Likelihood

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

Linear cost function

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

Mean Square Error

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

$$\mathbb{E}_{f_{\theta}(x)}[\xi^{\top} z] = \mathbb{E}_{f_{\theta}(x)}[\xi]^{\top} z = g_{\theta}(x)^{\top} z$$

Weighted SAA

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

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

Weighted SAA

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

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

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:

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

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: $w_i^{\text{KDE}}(\mathbf{x}) := \frac{\mathcal{K}((\mathbf{x} - \mathbf{x}^i)/\theta)}{\sum_{j=1}^N \mathcal{K}((\mathbf{x} - \mathbf{x}^j)/\theta)}$

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: $w_i^{\text{KDE}}(\mathbf{x}) := \frac{\mathcal{K}((\mathbf{x} - \mathbf{x}^i)/\theta)}{\sum_{j=1}^N \mathcal{K}((\mathbf{x} - \mathbf{x}^j)/\theta)}$

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

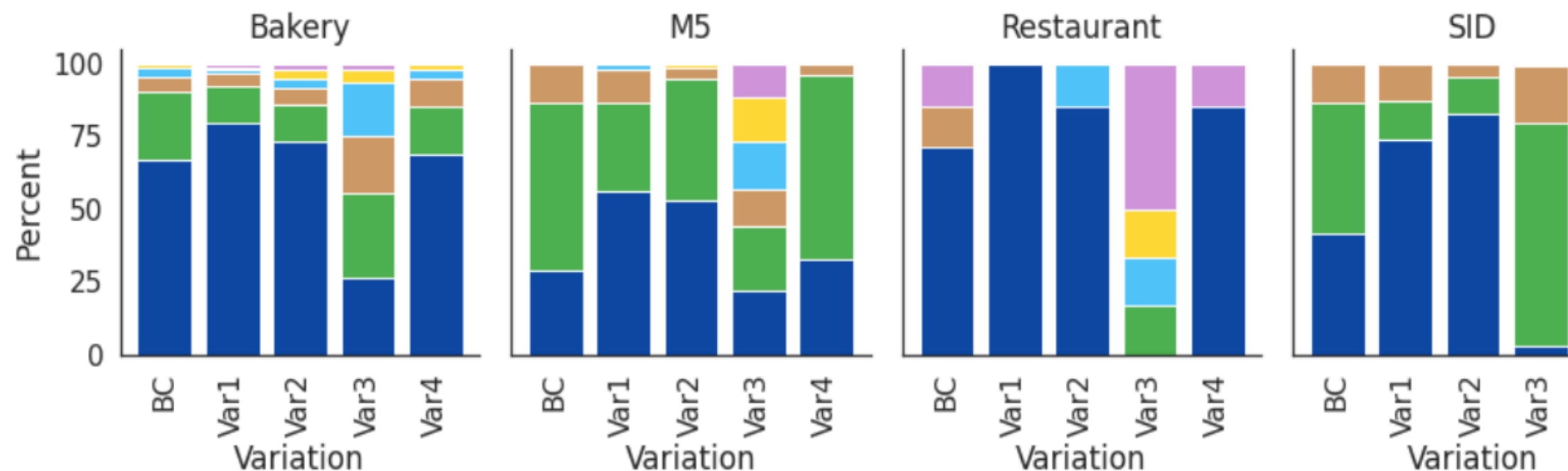
- Converges to optimal contextual policy as the size of the training set increases **when model is well specified.**

Some benchmark results (Buttler et al., 2023)

Newsvendor Problem

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

Proportion of instances where methods achieved best performance



Models:

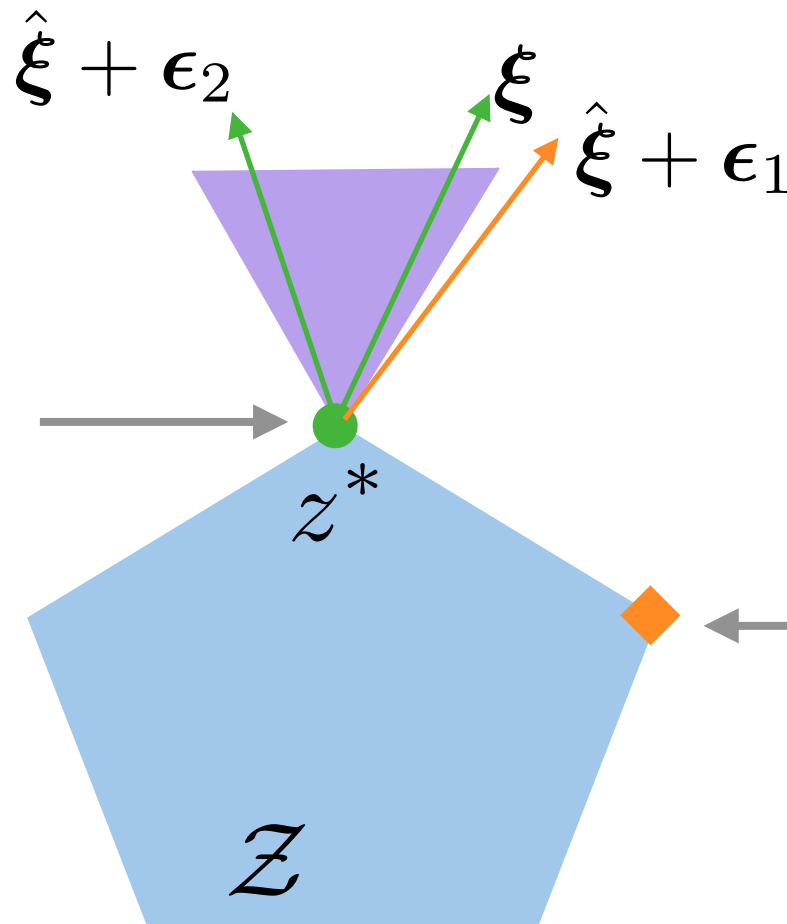
- Linear rule
- Kernel weights
- Decision tree weights
- Deep learning
- K-nearest neighbour weights
- Random forest weights

Going beyond SLO: Integrated learning and optimization

Wrong predictions lead to suboptimal decisions

$$\max_{z \in \mathcal{Z}} \xi^\top z$$

optimal
decision

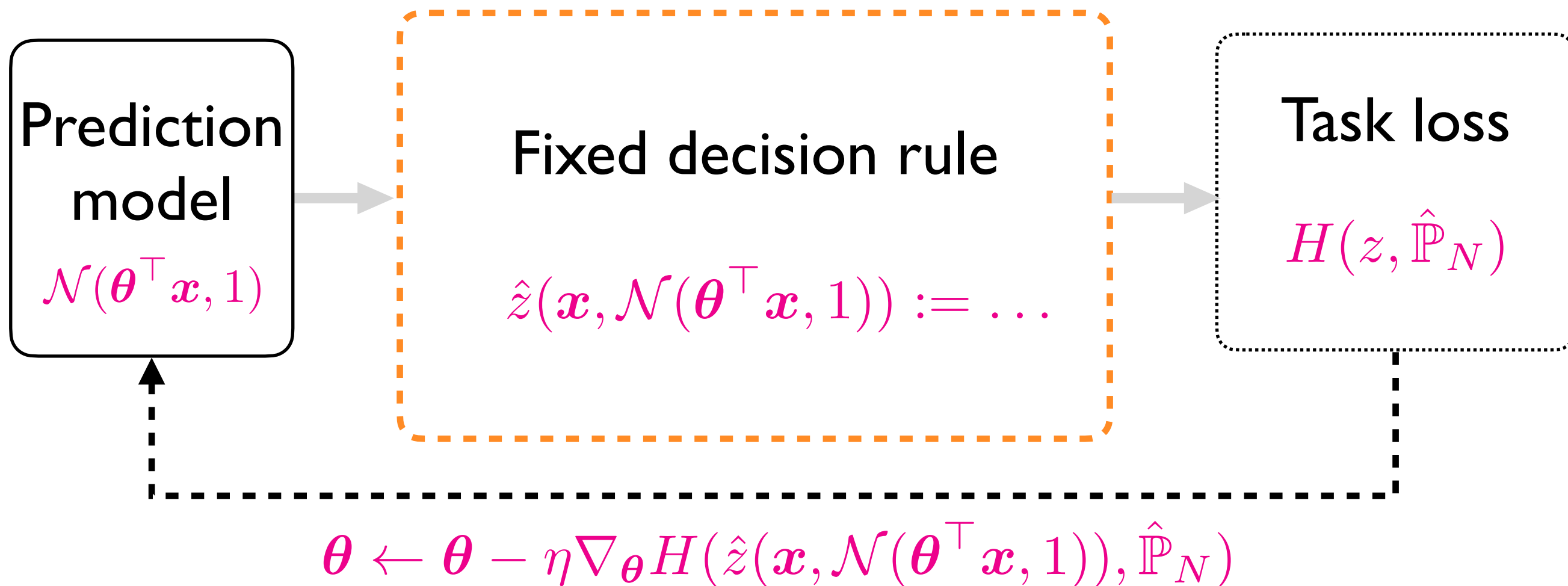


$$\|\epsilon_1\| \leq \|\epsilon_2\|$$

Suboptimal
decision

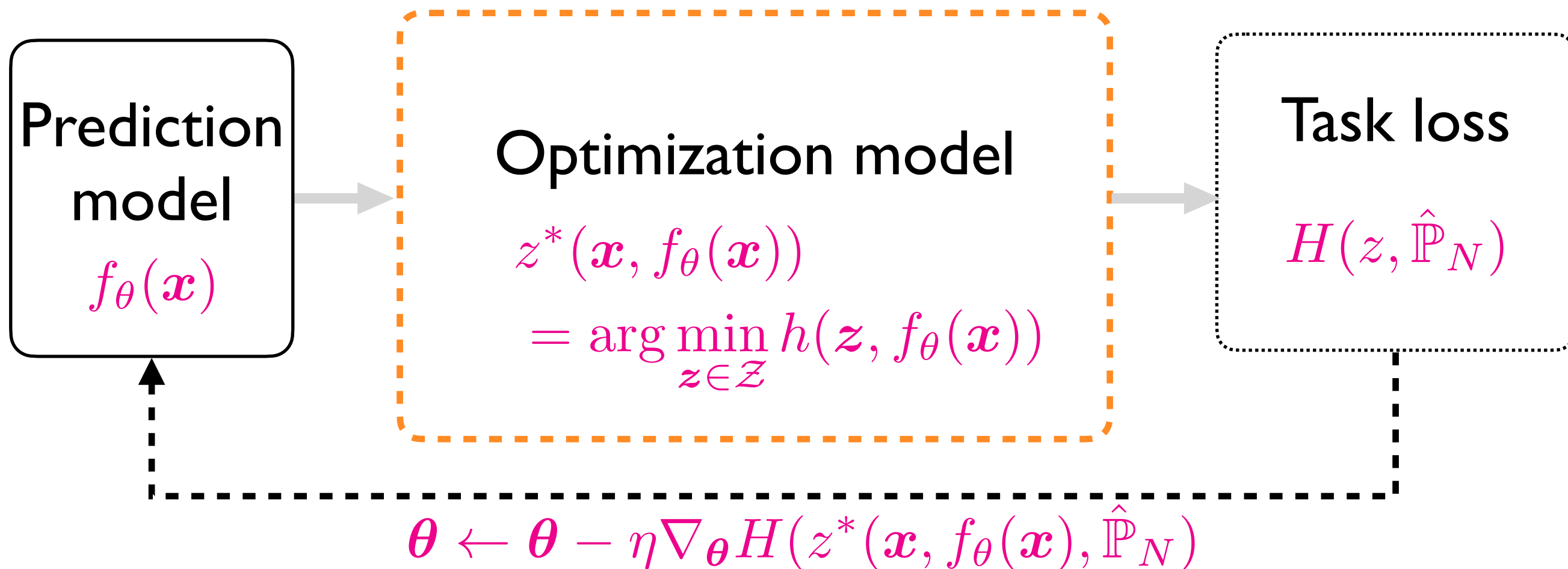
Figure adapted from [Elmachtoub and Grigas 2022]

ILO Training pipeline



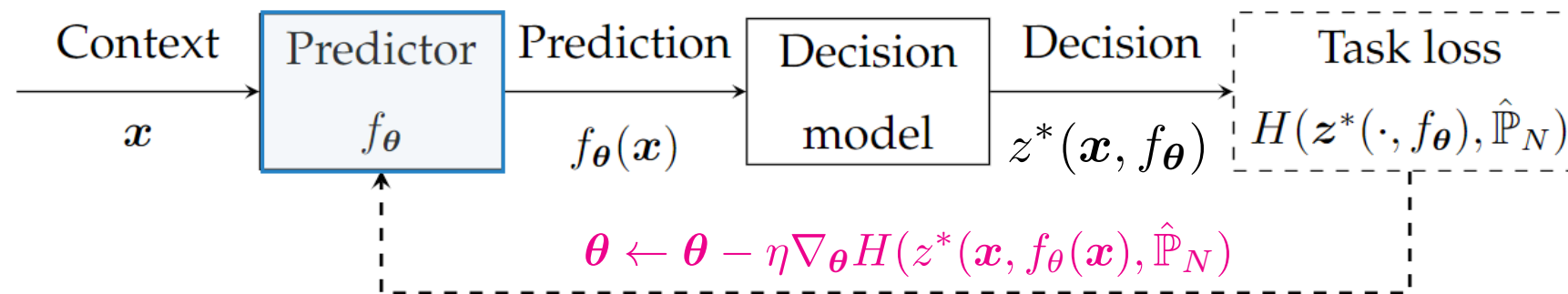
- [Bengio 1997] : Task-aware point prediction under a **fixed decision rule**

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]:

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

Smart “Predict, then optimize”

- Regret minimization [Elmachtoub & Grigas, 2022]:

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

Smart “Predict, then optimize”

- Regret minimization [Elmachtoub & Grigas, 2022]:

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

- Non-convex and discontinuous in θ

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

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

- Solve two optimization problems (MILP) at each iteration
- SPO+ has slower convergence rate when compared to sequential estimate then optimize model
- Model misspecification: SPO+ outperforms MSE

Optimal action imitation



- Imitation performance metric:

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

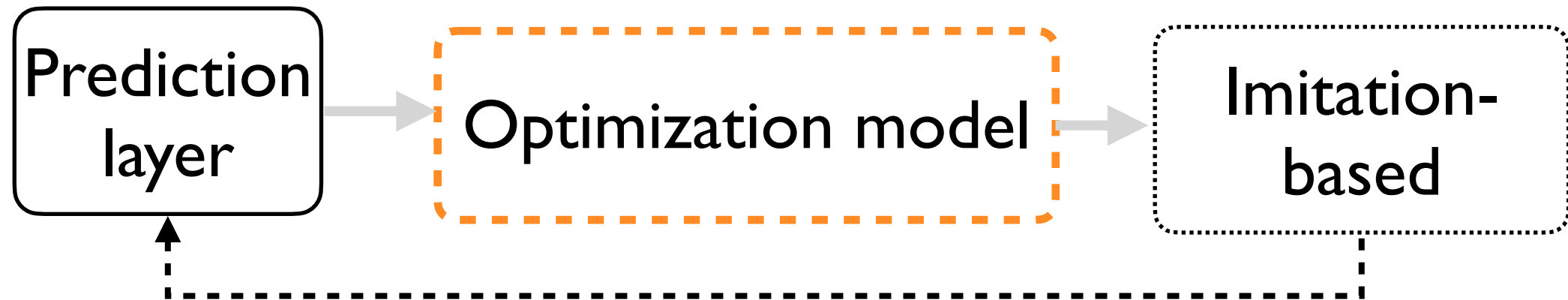
Optimal action imitation



- Imitation performance metric:

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

Optimal action imitation



- Imitation performance metric:

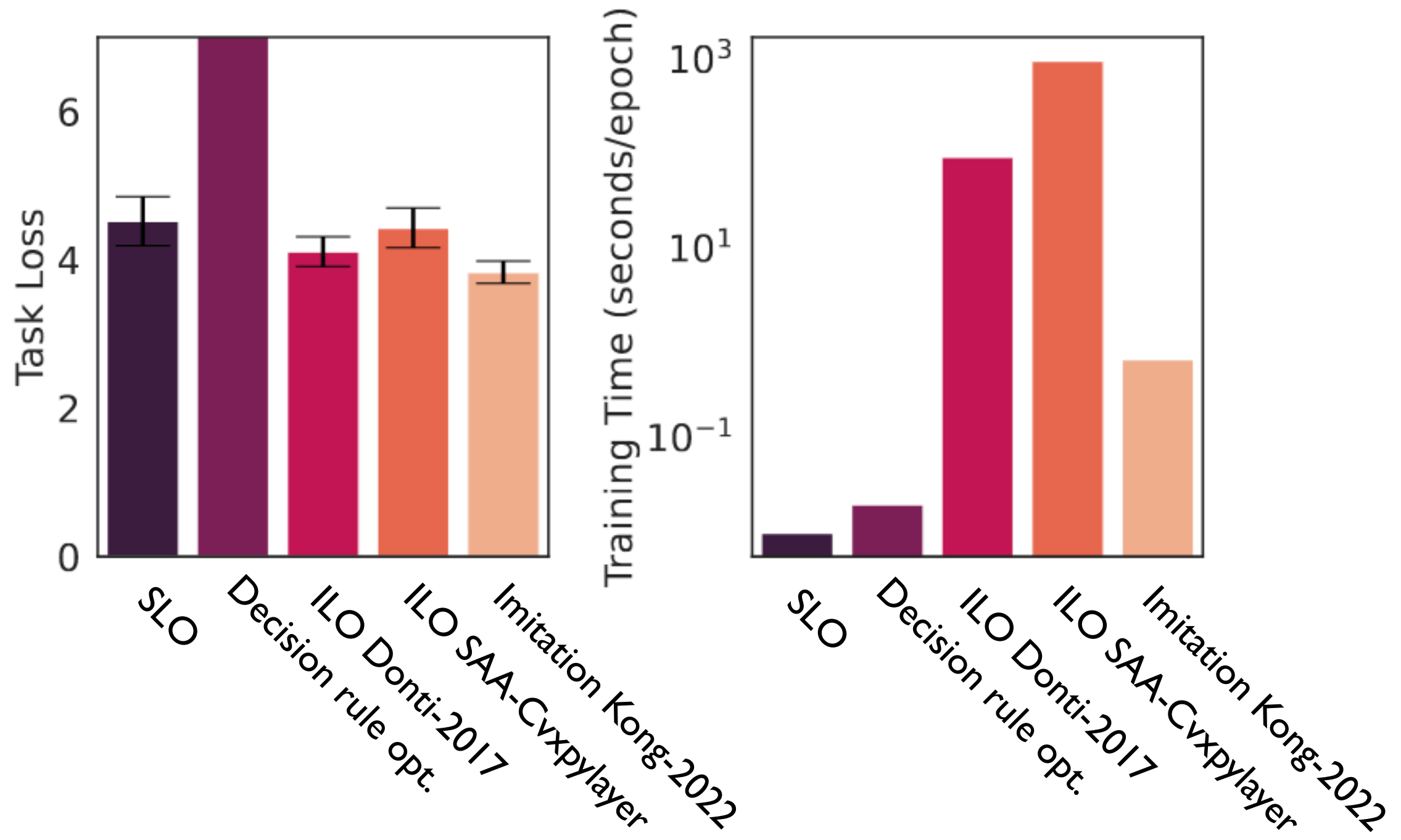
$$H(z^*(\mathbf{x}, f_\theta), \mathbb{P}) := \mathbb{E}_{\mathbb{P}}[\cancel{c(z^*(\mathbf{x}, f_\theta), \boldsymbol{\xi})}] \quad \mathbb{E}_{\hat{\mathbb{P}}_N}[d(z^*(\mathbf{x}, f_\theta), z^*(\mathbf{x}, \boldsymbol{\xi}))]$$

- 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

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

Comparison of some models

Load forecasting and generator scheduling problem
(objective similar to newsvendor 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 in mining?



(Link to survey paper)