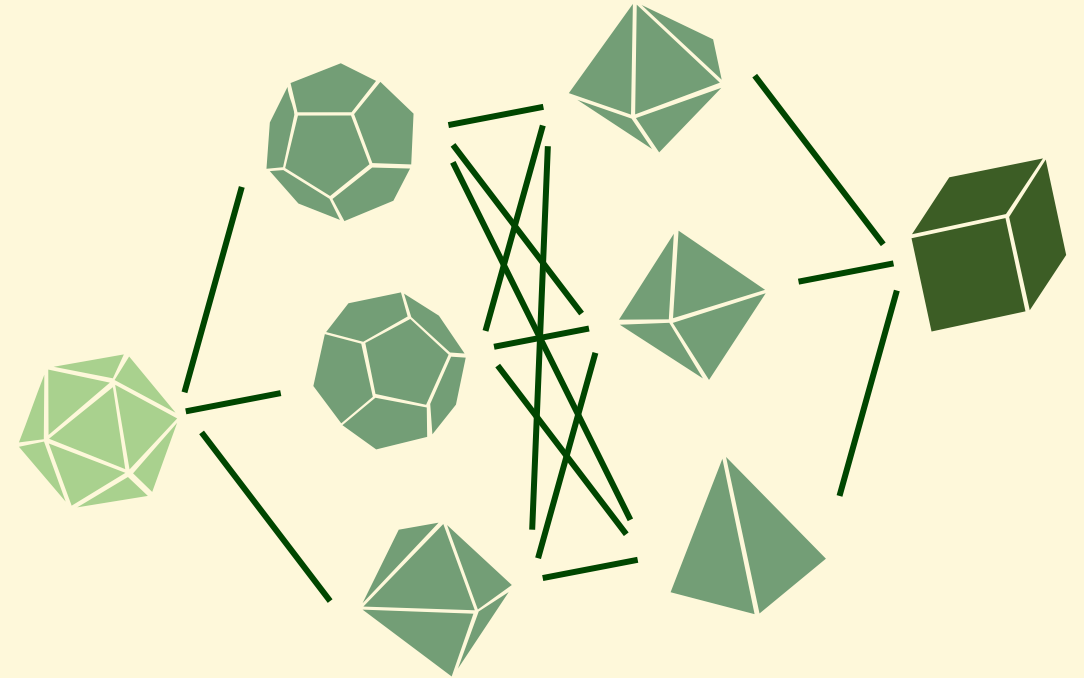




A PyTorch-based End-to-End Predict-then-Optimize Library



Presented by Bo Tang
Feb 26, 2024



Authors



Bo Tang

PhD Candidate
Department of Mechanical &
Industrial Engineering,
University of Toronto



Elias B. Khalil

Assistant Professor
Department of Mechanical &
Industrial Engineering, University
of Toronto

SCALE AI Research Chair
Data-Driven Algorithms for
Modern Supply Chains

Linear Objective Function

Linear Objective
Function

$$\min_{\underbrace{\mathbf{w}}_{\text{Decision Variables}}} \{ \underbrace{\mathbf{c}^T \mathbf{w}}_{\text{Linear Objective Function}} : \underbrace{\mathbf{w} \in \mathcal{S}}_{\text{Feasible Region}} \}$$

Decision Variables

Feasible Region

Use appropriate algorithm (MILP, MIQCP, CP, custom algorithms...) to obtain optimal solution $\mathbf{w}^*(\mathbf{c})$

Unknown Cost Coefficients

$$\min_w \{ \mathbf{c}_1^\top \mathbf{w} : \mathbf{w} \in S \}$$

$$\min_w \{ \mathbf{c}_2^\top \mathbf{w} : \mathbf{w} \in S \}$$

$$\min_w \{ \mathbf{c}_3^\top \mathbf{w} : \mathbf{w} \in S \}$$

$$\vdots$$

Unknown Cost Coefficients

Unknown Cost

$$\min_w \{ \mathbf{c}_1^\top \mathbf{w} : \mathbf{w} \in S \}$$
$$\min_w \{ \mathbf{c}_2^\top \mathbf{w} : \mathbf{w} \in S \}$$
$$\min_w \{ \mathbf{c}_3^\top \mathbf{w} : \mathbf{w} \in S \}$$

\therefore Identical
Constraints

Unknown Cost Coefficients

Unknown Cost

$$\min_w \{ \mathbf{c}_1^\top \mathbf{w} : \mathbf{w} \in S \}$$

$$\min_w \{ \mathbf{c}_2^\top \mathbf{w} : \mathbf{w} \in S \}$$

$$\min_w \{ \mathbf{c}_3^\top \mathbf{w} : \mathbf{w} \in S \}$$

\therefore Identical Constraints

Observed Feature Vector

\mathbf{x}_1

\mathbf{x}_2

\mathbf{x}_3

\vdots

Unknown Cost Coefficients

Unknown Cost

$$\min_w \{ \mathbf{c}_1^\top \mathbf{w} : \mathbf{w} \in S \}$$

$$\min_w \{ \mathbf{c}_2^\top \mathbf{w} : \mathbf{w} \in S \}$$

$$\min_w \{ \mathbf{c}_3^\top \mathbf{w} : \mathbf{w} \in S \}$$

\therefore Identical Constraints

Observed Feature Vector

\mathbf{x}_1

\mathbf{x}_2

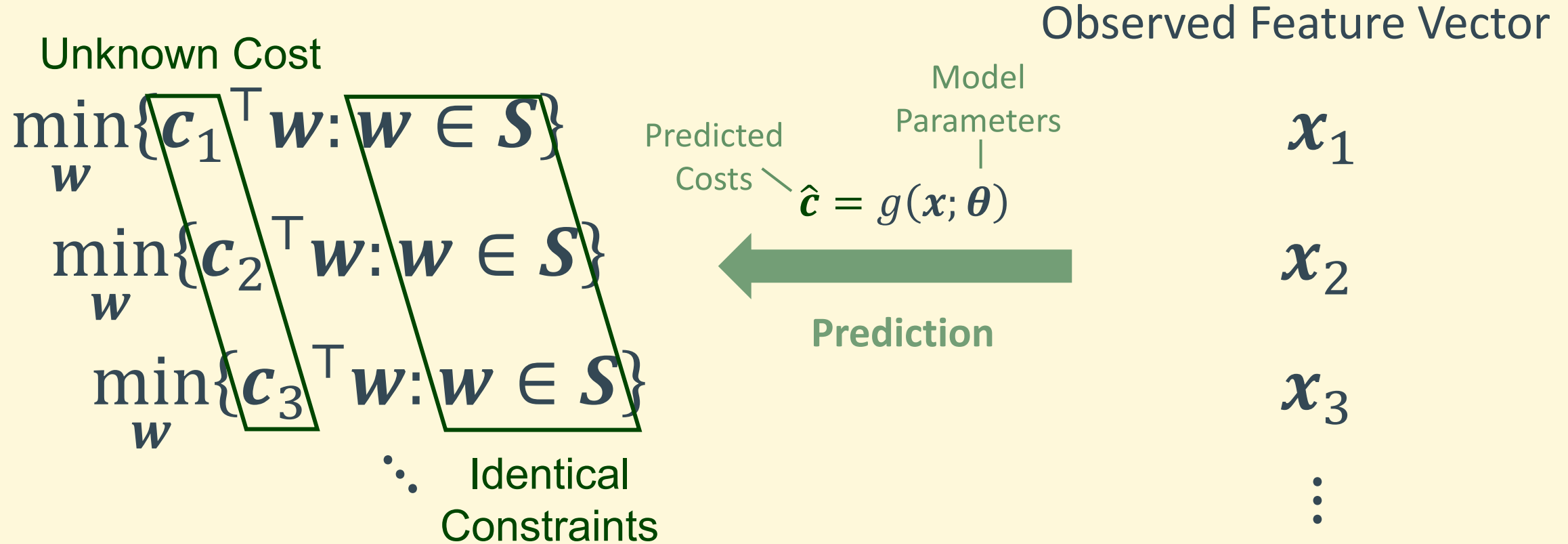
\mathbf{x}_3

\vdots

$$\hat{\mathbf{c}} = g(\mathbf{x}; \theta)$$

Prediction

Unknown Cost Coefficients



Examples



❖ Vehicle Routing



❖ Energy Scheduling



❖ Portfolio Optimization

Examples



❖ Vehicle Routing



❖ Energy Scheduling



❖ Portfolio Optimization

? Unknown Costs: Travel Time, Electricity Prices, Asset Returns

Examples



❖ Vehicle Routing



❖ Energy Scheduling



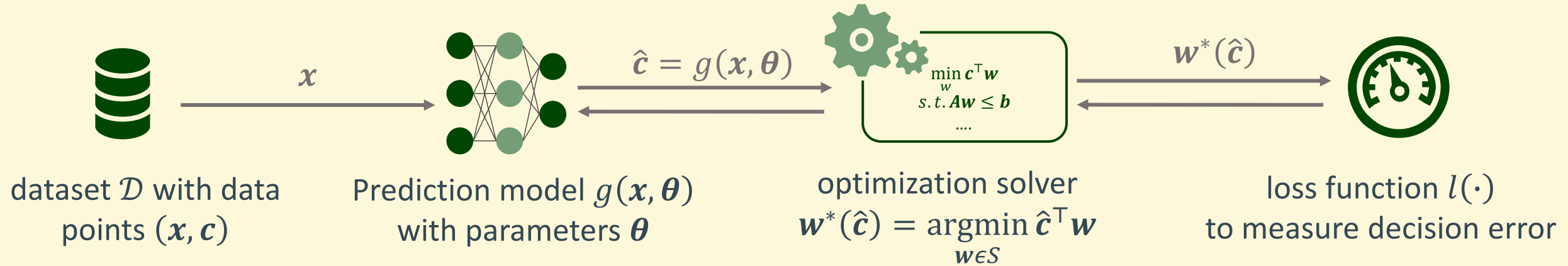
❖ Portfolio Optimization

? Unknown Costs: Travel Time, Electricity Prices, Asset Returns

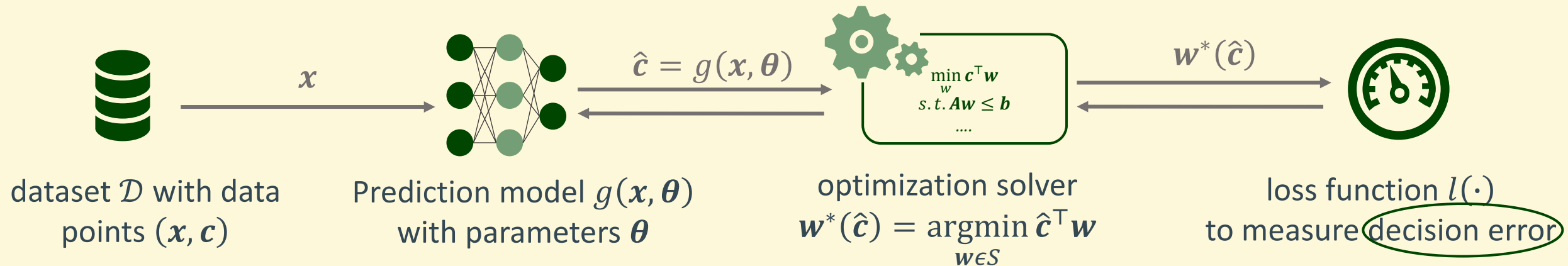


👁 Observed Features: Distance, Time, Weather, Financial Factors...

End-to-End Predict-then-Optimize



End-to-End Predict-then-Optimize



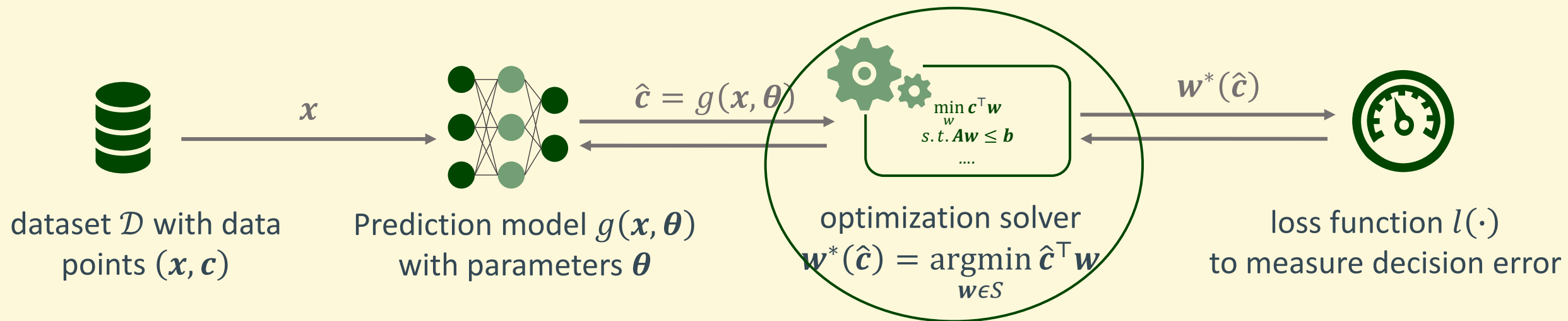
Optimal solution if you optimize with \hat{c} **Optimal value with true cost c**

e.g.

$$l_{\text{Regret}}(\hat{c}, c) = \overbrace{c^T w^*(\hat{c})}^{\text{Optimal value with true cost } c} - \overbrace{c^T w^*(c)}^{\text{Optimal solution if you optimize with } \hat{c}}$$

$$l_{\text{Square}}(\hat{c}, c) = \frac{1}{2} \|w^*(c) - w^*(\hat{c})\|_2^2$$

End-to-End Predict-then-Optimize

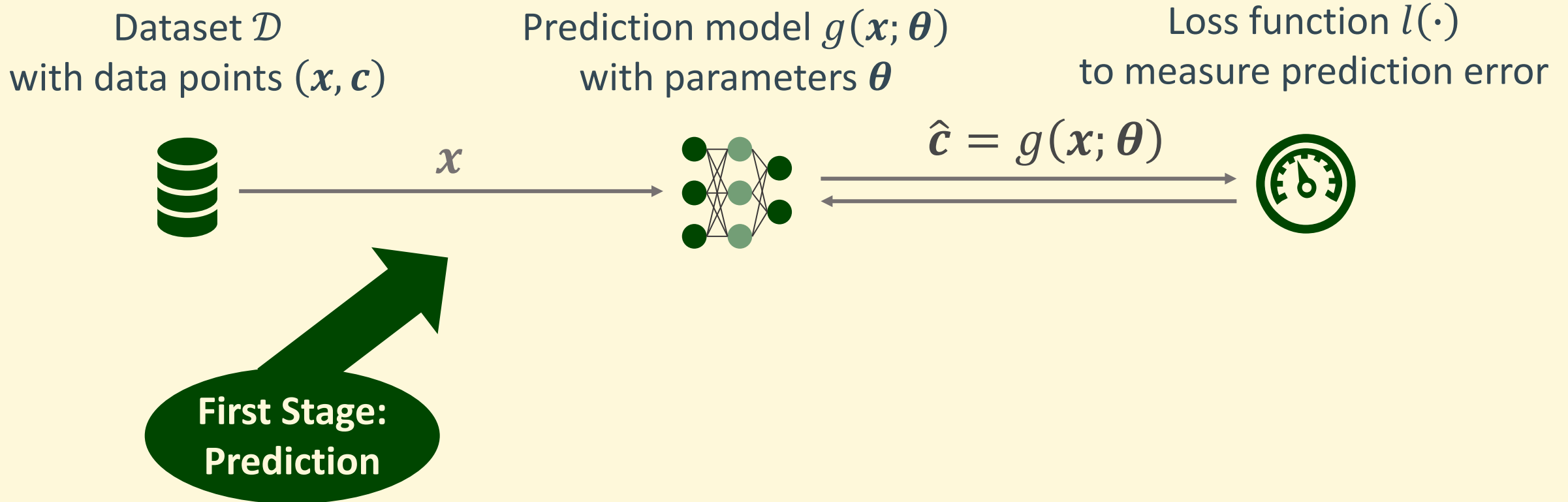


Algorithm 1 End-to-end Learning

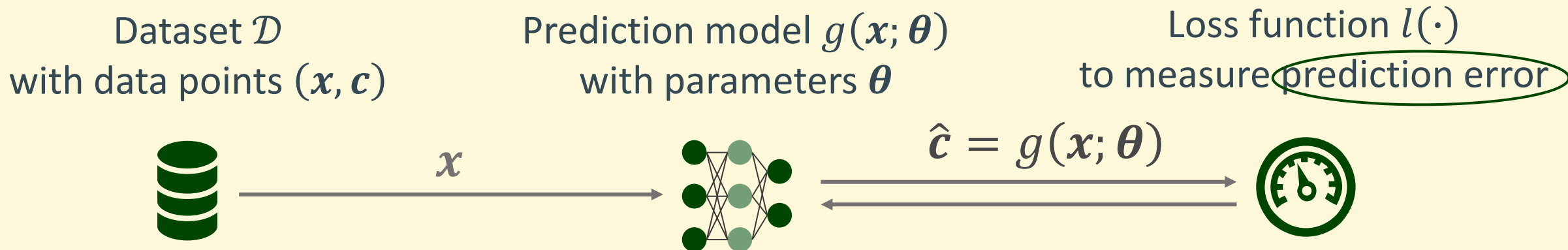
Require: coefficient matrix \mathbf{A} , right-hand side \mathbf{b} , data \mathcal{D}

- 1: Initialize predictor parameters $\boldsymbol{\theta}$ for predictor $g(\mathbf{x}; \boldsymbol{\theta})$
- 2: **for** epochs **do**
- 3: **for** each batch of training data (\mathbf{x}, \mathbf{c}) **do**
- 4: Sample batch of the cost vectors \mathbf{c} with the corresponding features \mathbf{x}
- 5: Predict cost using predictor $\hat{\mathbf{c}} := g(\mathbf{x}; \boldsymbol{\theta})$
- 6: Forward pass to compute optimal solution $\mathbf{w}_{\hat{\mathbf{c}}}^* := \operatorname{argmin}_{\mathbf{w} \in S} \hat{\mathbf{c}}^T \mathbf{w}$
- 7: Forward pass to compute decision loss $l(\cdot)$
- 8: Backward pass from loss $l(\cdot)$ to update parameters $\boldsymbol{\theta}$ with gradient
- 9: **end for**
- 10: **end for**

Two-Stage Predict-then-Optimize



Two-Stage Predict-then-Optimize



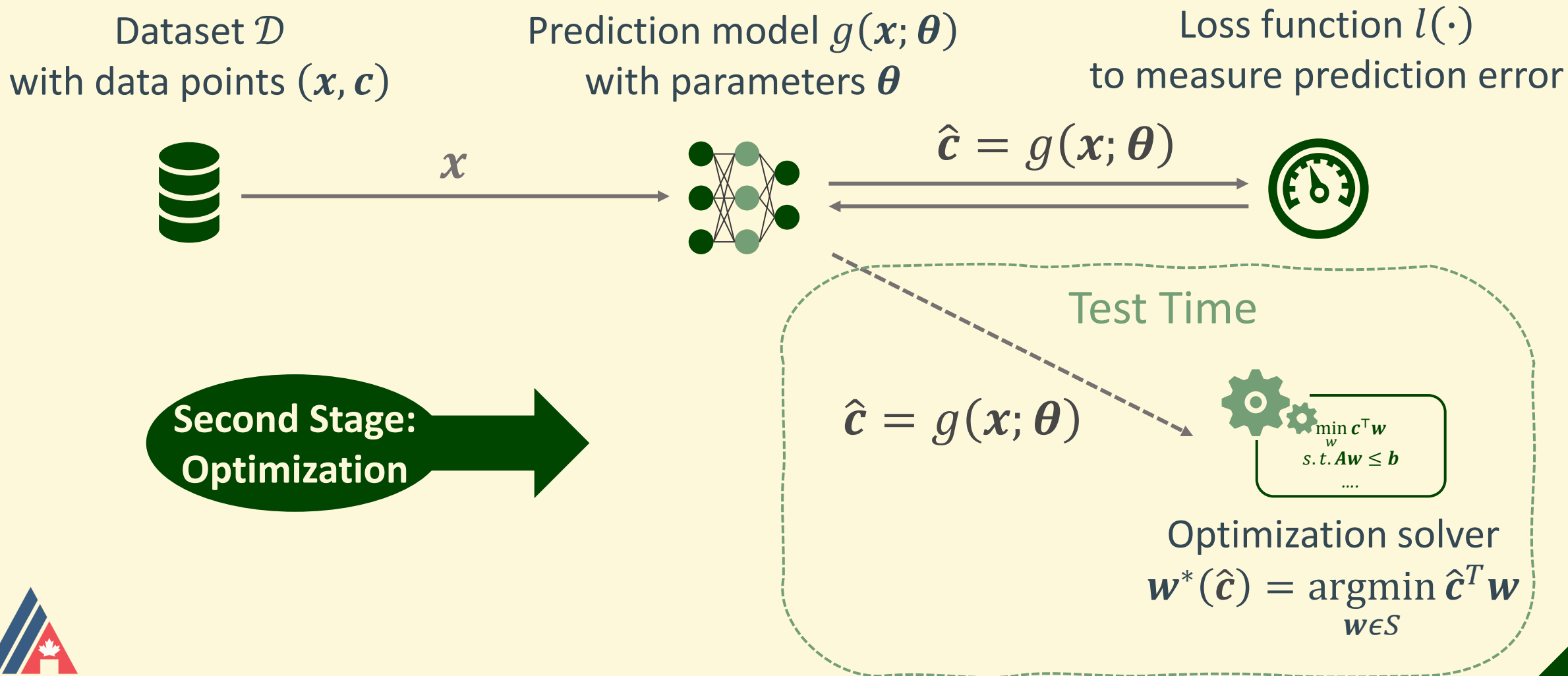
**First Stage:
Prediction**

e.g.

$$l_{\text{MSE}}(\hat{\mathbf{c}}, \mathbf{c}) = \frac{1}{2} \|\mathbf{c} - \hat{\mathbf{c}}\|_2^2$$

$$l_{\text{MAE}}(\hat{\mathbf{c}}, \mathbf{c}) = \|\mathbf{c} - \hat{\mathbf{c}}\|_1$$

Two-Stage Predict-then-Optimize



Mismatch of Prediction and Decision Error

For example:

$$\begin{aligned} \max_{w_1, w_2} \quad & c_1 w_1 + c_2 w_2 \\ \text{s.t.} \quad & w_1 + w_2 \leq 1 \\ & w_1, w_2 \geq 0 \end{aligned}$$

Assume the true cost is $\mathbf{c} = (0,1)$, the optimal solution is $\mathbf{w}^*(\mathbf{c}) = (0,1)$

If the prediction $\hat{\mathbf{c}} = (1,0)$, the solution $\mathbf{w}^*(\hat{\mathbf{c}}) =$ and $l_{\text{MSE}}(\hat{\mathbf{c}}, \mathbf{c}) =$;

If the prediction $\hat{\mathbf{c}} = (0,3)$, the solution $\mathbf{w}^*(\hat{\mathbf{c}}) =$ and $l_{\text{MSE}}(\hat{\mathbf{c}}, \mathbf{c}) =$

Mismatch of Prediction and Decision Error

For example:

$$\begin{aligned} \max_{w_1, w_2} \quad & c_1 w_1 + c_2 w_2 \\ \text{s.t.} \quad & w_1 + w_2 \leq 1 \\ & w_1, w_2 \geq 0 \end{aligned}$$

Assume the true cost is $\mathbf{c} = (0,1)$, the optimal solution is $\mathbf{w}^*(\mathbf{c}) = (0,1)$

If the prediction $\hat{\mathbf{c}} = (1,0)$, the solution $\mathbf{w}^*(\hat{\mathbf{c}}) = (1,0)$ and $l_{\text{MSE}}(\hat{\mathbf{c}}, \mathbf{c}) = 1$;

If the prediction $\hat{\mathbf{c}} = (0,3)$, the solution $\mathbf{w}^*(\hat{\mathbf{c}}) = (0,1)$ and $l_{\text{MSE}}(\hat{\mathbf{c}}, \mathbf{c}) = 2$

Mismatch of Prediction and Decision Error

For example:

$$\begin{aligned} \max_{w_1, w_2} \quad & c_1 w_1 + c_2 w_2 \\ \text{s.t.} \quad & w_1 + w_2 \leq 1 \\ & w_1, w_2 \geq 0 \end{aligned}$$

Assume the true cost is $\mathbf{c} = (0,1)$, the optimal solution is $\mathbf{w}^*(\mathbf{c}) = (0,1)$

If the prediction $\hat{\mathbf{c}} = (1,0)$, the solution $\mathbf{w}^*(\hat{\mathbf{c}}) = (1,0)$ and $l_{\text{MSE}}(\hat{\mathbf{c}}, \mathbf{c}) = 1$;

If the prediction $\hat{\mathbf{c}} = (0,3)$, the solution $\mathbf{w}^*(\hat{\mathbf{c}}) = (0,1)$ and $l_{\text{MSE}}(\hat{\mathbf{c}}, \mathbf{c}) = 2$

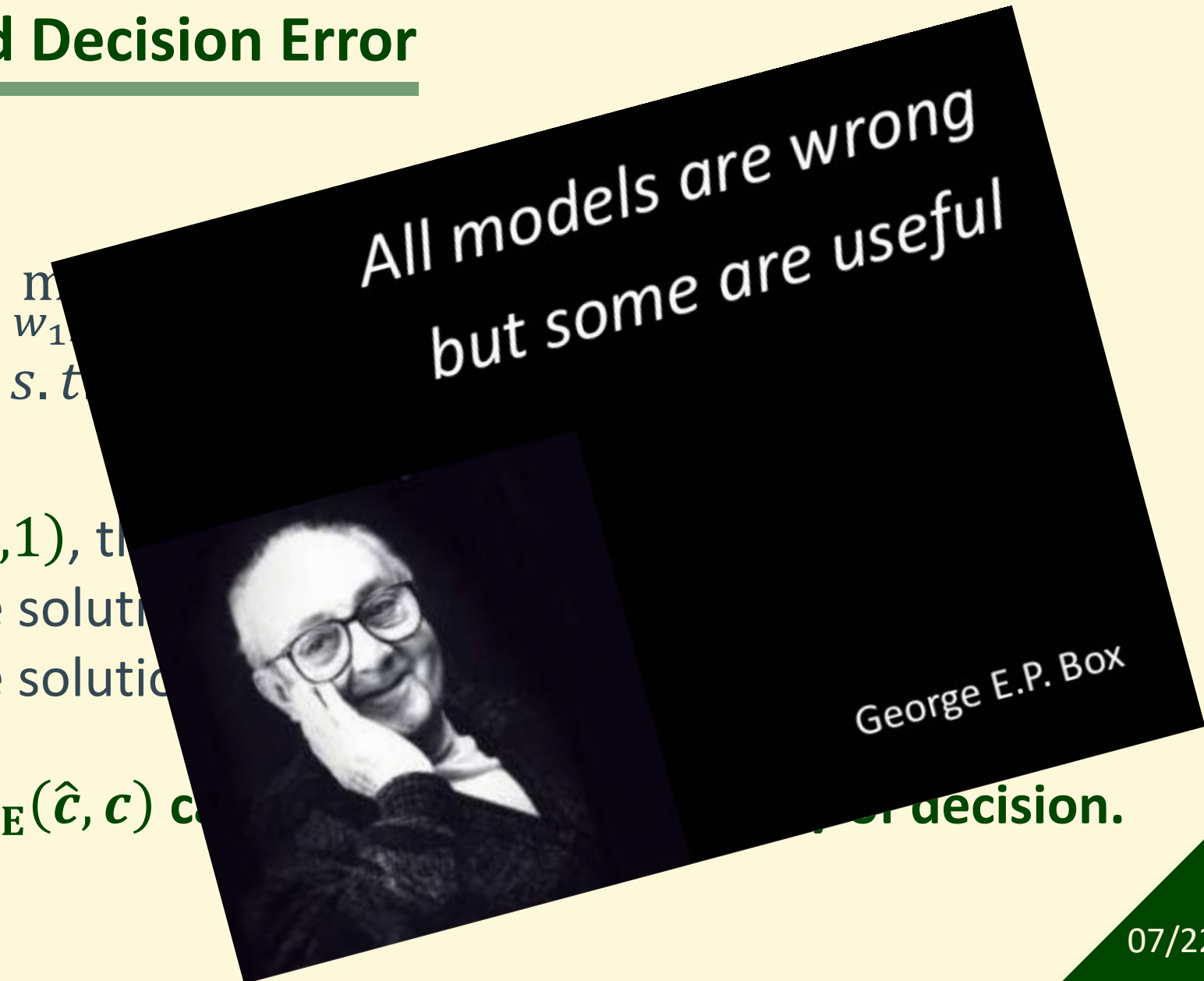
Prediction error such as $l_{\text{MSE}}(\hat{\mathbf{c}}, \mathbf{c})$ cannot measure the quality of decision.

Mismatch of Prediction and Decision Error

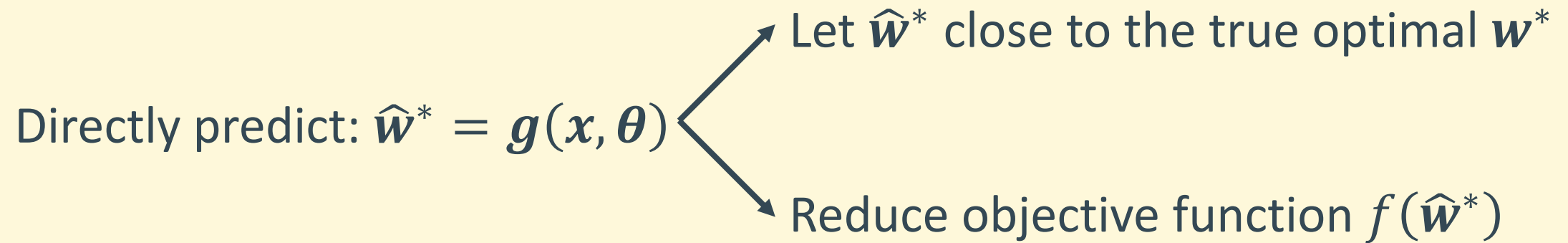
For example:

Assume the true cost is $\mathbf{c} = (0,1)$, the solution is \mathbf{w}_1
 If the prediction $\hat{\mathbf{c}} = (1,0)$, the solution is \mathbf{s}^*
 If the prediction $\hat{\mathbf{c}} = (0,3)$, the solution is \mathbf{s}^*

Prediction error such as $l_{\text{MSE}}(\hat{\mathbf{c}}, \mathbf{c})$ can be a poor proxy for decision error.



Imitation Learning and Parametric Optimization



Imitation Learning and Parametric Optimization

Directly predict: $\hat{\mathbf{w}}^* = g(\mathbf{x}, \theta)$

- Let $\hat{\mathbf{w}}^*$ close to the true optimal \mathbf{w}^*
- Reduce objective function $f(\hat{\mathbf{w}}^*)$



Efficiency

Avoid the major bottleneck
in computational efficiency:
optimizing

Imitation Learning and Parametric Optimization

Directly predict: $\hat{\mathbf{w}}^* = g(\mathbf{x}, \theta)$

- Let $\hat{\mathbf{w}}^*$ close to the true optimal \mathbf{w}^*
- Reduce objective function $f(\hat{\mathbf{w}}^*)$



Efficiency

Avoid the major bottleneck
in computational efficiency:
optimizing



Feasibility

Prediction often faces
feasibility issues

Chain-Rule of Gradient

Algorithm 1 End-to-end Learning

Require: coefficient matrix \mathbf{A} , right-hand side \mathbf{b} , data \mathcal{D}

- 1: Initialize predictor parameters θ for predictor $g(\mathbf{x}; \theta)$
 - 2: **for** epochs **do**
 - 3: **for** each batch of training data (\mathbf{x}, \mathbf{c}) **do**
 - 4: Sample batch of the cost vectors \mathbf{c} with the corresponding features \mathbf{x}
 - 5: Predict cost using predictor $\hat{\mathbf{c}} := g(\mathbf{x}; \theta)$
 - 6: Forward pass to compute optimal solution $\mathbf{w}_{\hat{\mathbf{c}}}^* := \operatorname{argmin}_{\mathbf{w} \in S} \hat{\mathbf{c}}^T \mathbf{w}$
 - 7: Forward pass to compute decision loss $l(\cdot)$
 - 8: Backward pass from loss $l(\cdot)$ to update parameters θ with gradient
 - 9: **end for**
 - 10: **end for**
-

Chain Rule:

$$\frac{\partial l(\cdot)}{\partial \theta} =$$

Chain-Rule of Gradient

Algorithm 1 End-to-end Learning

Require: coefficient matrix \mathbf{A} , right-hand side \mathbf{b} , data \mathcal{D}

- 1: Initialize predictor parameters θ for predictor $g(\mathbf{x}; \theta)$
 - 2: **for** epochs **do**
 - 3: **for** each batch of training data (\mathbf{x}, \mathbf{c}) **do**
 - 4: Sample batch of the cost vectors \mathbf{c} with the corresponding features \mathbf{x}
 - 5: Predict cost using predictor $\hat{\mathbf{c}} := g(\mathbf{x}; \theta)$
 - 6: Forward pass to compute optimal solution $\mathbf{w}_{\hat{\mathbf{c}}}^* := \operatorname{argmin}_{\mathbf{w} \in S} \hat{\mathbf{c}}^T \mathbf{w}$
 - 7: Forward pass to compute decision loss $l(\cdot)$
 - 8: Backward pass from loss $l(\cdot)$ to update parameters θ with gradient
 - 9: **end for**
 - 10: **end for**
-

Chain Rule:

$$\frac{\partial l(\cdot)}{\partial \theta} = \frac{\partial l(\cdot)}{\partial \hat{\mathbf{c}}} \frac{\partial \hat{\mathbf{c}}}{\partial \theta}$$

Chain-Rule of Gradient

Algorithm 1 End-to-end Learning

Require: coefficient matrix \mathbf{A} , right-hand side \mathbf{b} , data \mathcal{D}

- 1: Initialize predictor parameters θ for predictor $g(\mathbf{x}; \theta)$
 - 2: **for** epochs **do**
 - 3: **for** each batch of training data (\mathbf{x}, \mathbf{c}) **do**
 - 4: Sample batch of the cost vectors \mathbf{c} with the corresponding features \mathbf{x}
 - 5: Predict cost using predictor $\hat{\mathbf{c}} := g(\mathbf{x}; \theta)$
 - 6: Forward pass to compute optimal solution $\mathbf{w}_{\hat{\mathbf{c}}}^* := \operatorname{argmin}_{\mathbf{w} \in S} \hat{\mathbf{c}}^T \mathbf{w}$
 - 7: Forward pass to compute decision loss $l(\cdot)$
 - 8: Backward pass from loss $l(\cdot)$ to update parameters θ with gradient
 - 9: **end for**
 - 10: **end for**
-

Chain Rule:

$$\frac{\partial l(\cdot)}{\partial \theta} = \frac{\partial l(\cdot)}{\partial \hat{\mathbf{c}}} \left(\frac{\partial \hat{\mathbf{c}}}{\partial \theta} \right)$$

Easy:

Gradient of predicted costs to model parameters

Chain-Rule of Gradient

Algorithm 1 End-to-end Learning

Require: coefficient matrix \mathbf{A} , right-hand side \mathbf{b} , data \mathcal{D}

```

1: Initialize predictor parameters  $\theta$  for predictor  $g(\mathbf{x}; \theta)$ 
2: for epochs do
3:   for each batch of training data  $(\mathbf{x}, \mathbf{c})$  do
4:     Sample batch of the cost vectors  $\mathbf{c}$  with the corresponding features  $\mathbf{x}$ 
5:     Predict cost using predictor  $\hat{\mathbf{c}} := g(\mathbf{x}; \theta)$ 
6:     Forward pass to compute optimal solution  $\mathbf{w}_{\hat{\mathbf{c}}}^* := \operatorname{argmin}_{\mathbf{w} \in S} \hat{\mathbf{c}}^T \mathbf{w}$ 
7:     Forward pass to compute decision loss  $l(\cdot)$ 
8:     Backward pass from loss  $l(\cdot)$  to update parameters  $\theta$  with gradient
9:   end for
10: end for

```

Chain Rule:

$$\frac{\partial l(\cdot)}{\partial \theta} = \left(\frac{\partial l(\cdot)}{\partial \hat{\mathbf{c}}} \right) \frac{\partial \hat{\mathbf{c}}}{\partial \theta}$$

Hard:

Decision loss vary with the predicted costs

Chain-Rule of Gradient

Algorithm 1 End-to-end Learning

Require: coefficient matrix \mathbf{A} , right-hand side \mathbf{b} , data \mathcal{D}

- 1: Initialize predictor parameters θ for predictor $g(\mathbf{x}; \theta)$
- 2: **for** epochs **do**
- 3: **for** each batch of training data (\mathbf{x}, \mathbf{c}) **do**
- 4: Sample batch of the cost vectors \mathbf{c} with the corresponding features \mathbf{x}
- 5: Predict cost using predictor $\hat{\mathbf{c}} := g(\mathbf{x}; \theta)$
- 6: Forward pass to compute optimal solution $\mathbf{w}_{\hat{\mathbf{c}}}^* := \operatorname{argmin}_{\mathbf{w} \in S} \hat{\mathbf{c}}^T \mathbf{w}$
- 7: Forward pass to compute decision loss $l(\cdot)$
- 8: Backward pass from loss $l(\cdot)$ to update parameters θ with gradient
- 9: **end for**
- 10: **end for**

Chain Rule:

$$\frac{\partial l(\cdot)}{\partial \theta} = \left(\frac{\partial l(\cdot)}{\partial \hat{\mathbf{c}}} \right) \frac{\partial \hat{\mathbf{c}}}{\partial \theta}$$

Hard:

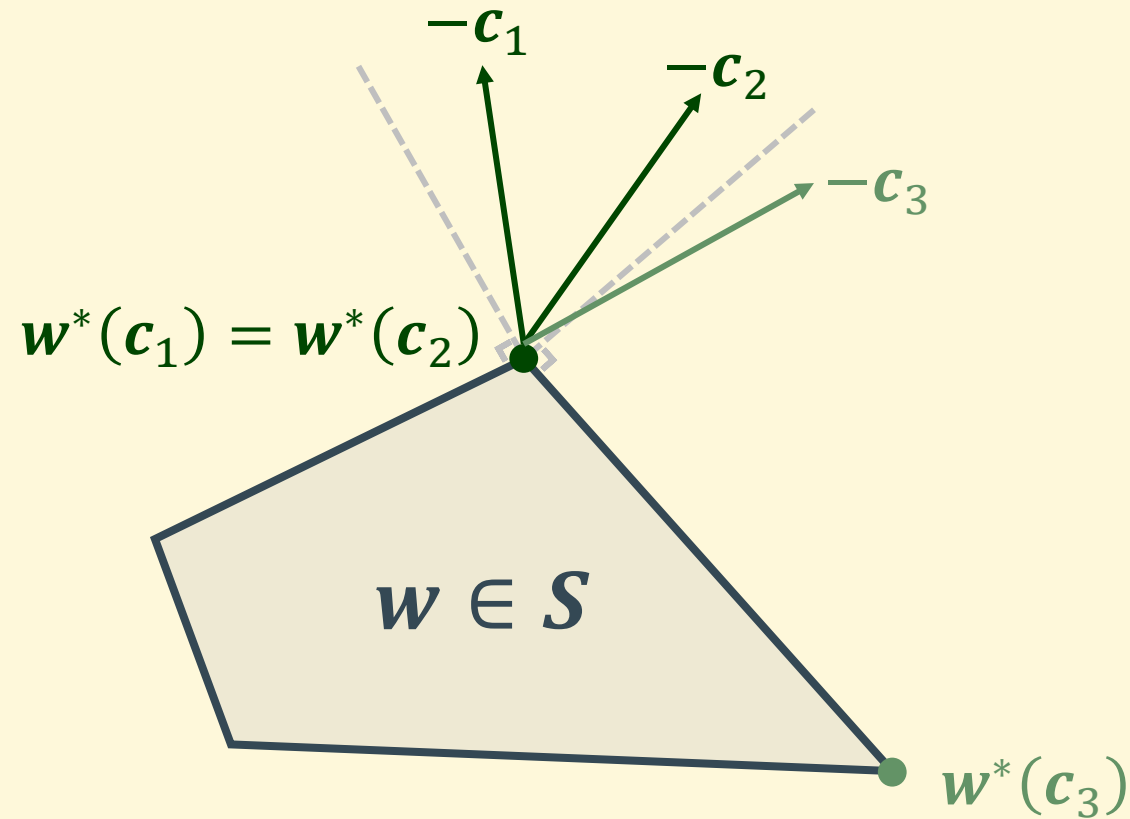
Decision loss vary with the predicted costs

$$l_{\text{Regret}}(\hat{\mathbf{c}}, \mathbf{c}) = \mathbf{c}^T \mathbf{w}^*(\hat{\mathbf{c}}) - \mathbf{c}^T \mathbf{w}^*(\mathbf{c})$$

$$l_{\text{Square}}(\hat{\mathbf{c}}, \mathbf{c}) = \frac{1}{2} \|\mathbf{w}^*(\mathbf{c}) - \mathbf{w}^*(\hat{\mathbf{c}})\|_2^2$$

with respect to $\mathbf{w}^*(\hat{\mathbf{c}})$!

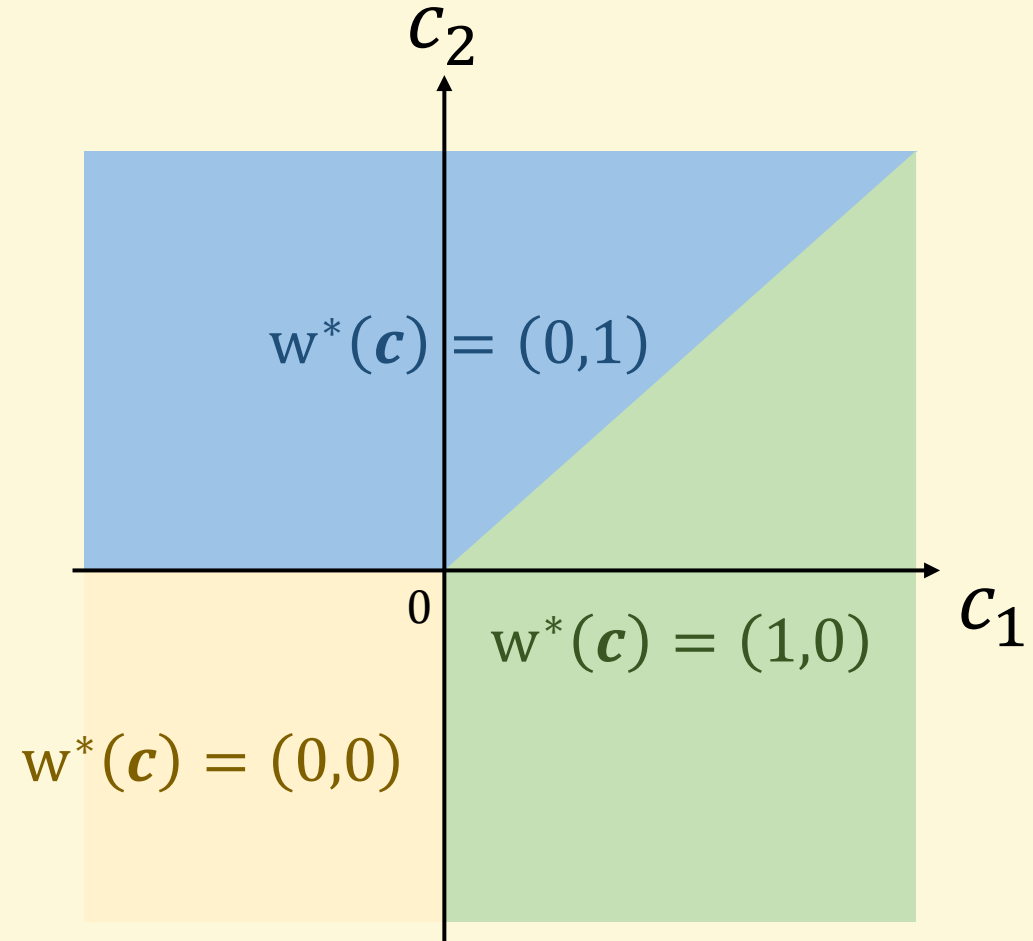
Piecewise Constant Function



$w^*(c)$ is a piecewise constant function!

Piecewise Constant Function

$$\begin{aligned} \max_{w_1, w_2} \quad & c_1 w_1 + c_2 w_2 \\ \text{s.t.} \quad & w_1 + w_2 \leq 1 \\ & w_1, w_2 \geq 0 \end{aligned}$$



Piecewise Constant Function

Is that the end of the story?

Piecewise Constant Function

Is that the end of the story?

Surrogate Losses / Gradients!

Derivative of Implicit Functions

OptNet:

- Solve the partial derivative matrix linear equations to calculate the solution and gradients for both **forward** and **backward** pass.
- Add a quadratic term to the linear objective function to obtain the **nonzero** gradient.

Karush-Kuhn-Tucker conditions

Given general problem

$$\begin{aligned} \min_{x \in \mathbb{R}^n} \quad & f(x) \\ \text{subject to} \quad & h_i(x) \leq 0, \quad i = 1, \dots, m \\ & \ell_j(x) = 0, \quad j = 1, \dots, r \end{aligned}$$

The **Karush-Kuhn-Tucker conditions** or **KKT conditions** are:

- $0 \in \partial f(x) + \sum_{i=1}^m u_i \partial h_i(x) + \sum_{j=1}^r v_j \partial \ell_j(x)$ (stationarity)
- $u_i \cdot h_i(x) = 0$ for all i (complementary slackness)
- $h_i(x) \leq 0, \ell_j(x) = 0$ for all i, j (primal feasibility)
- $u_i \geq 0$ for all i (dual feasibility)

- Amos, B., & Kolter, J. Z. (2017, July). Optnet: Differentiable optimization as a layer in neural networks. In International Conference on Machine Learning (pp. 136-145). PMLR.
- Wilder, B., Dilkina, B., & Tambe, M. (2019, July). Melding the data-decisions pipeline: Decision-focused learning for combinatorial optimization. In Proceedings of the AAAI Conference on Artificial Intelligence (Vol. 33, No. 01, pp. 1658-1665).

Smart “predict, then optimize”

A convex upper bound of $l_{\text{Regret}}(\hat{\mathbf{c}}, \mathbf{c})$:

- Elmachtoub, A. N., & Grigas, P. (2021). Smart “predict, then optimize”. Management Science.

Smart “predict, then optimize”

A convex upper bound of $l_{\text{Regret}}(\hat{\mathbf{c}}, \mathbf{c})$:

$$l_{\text{SPO}+}(\hat{\mathbf{c}}, \mathbf{c}) = -\min_{\mathbf{w} \in \mathcal{W}} (2\hat{\mathbf{c}} - \mathbf{c})^{\top} \mathbf{w} + 2\hat{\mathbf{c}}^{\top} \mathbf{w}^*(\mathbf{c}) - \mathbf{c}^{\top} \mathbf{w}^*(\mathbf{c})$$

- Elmachtoub, A. N., & Grigas, P. (2021). Smart “predict, then optimize”. Management Science.

Smart “predict, then optimize”

A convex upper bound of $l_{\text{Regret}}(\hat{\mathbf{c}}, \mathbf{c})$:

$$l_{\text{SPO}+}(\hat{\mathbf{c}}, \mathbf{c}) = - \underbrace{\min_{\mathbf{w} \in W} (2\hat{\mathbf{c}} - \mathbf{c})^\top \mathbf{w}} + 2\hat{\mathbf{c}}^\top \mathbf{w}^*(\mathbf{c}) - \mathbf{c}^\top \mathbf{w}^*(\mathbf{c})$$

Computational overhead: We need to solve an optimization problem $\min_{\mathbf{w} \in W} (2\hat{\mathbf{c}} - \mathbf{c})^\top \mathbf{w}$ per iteration.

Smart “predict, then optimize”

A convex upper bound of $l_{\text{Regret}}(\hat{\mathbf{c}}, \mathbf{c})$:

$$l_{\text{SPO}+}(\hat{\mathbf{c}}, \mathbf{c}) = -\min_{\mathbf{w} \in \mathcal{W}} (2\hat{\mathbf{c}} - \mathbf{c})^\top \mathbf{w} + 2\hat{\mathbf{c}}^\top \mathbf{w}^*(\mathbf{c}) - \mathbf{c}^\top \mathbf{w}^*(\mathbf{c})$$

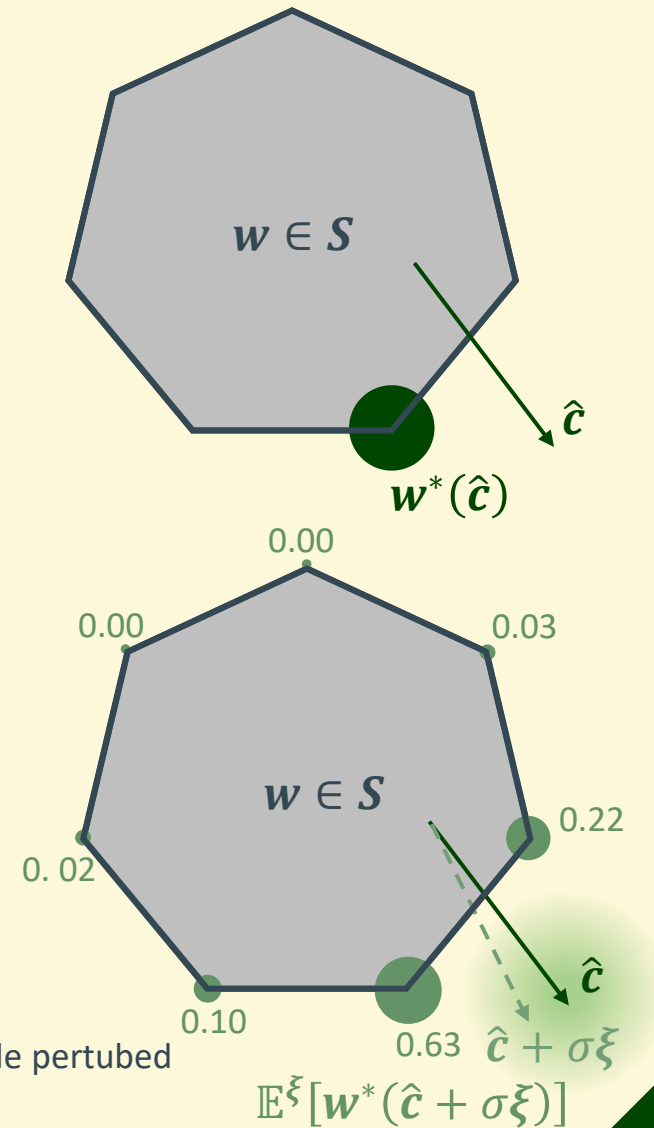
The subgradient of $l_{\text{SPO}+}(\hat{\mathbf{c}}, \mathbf{c})$:

$$2\mathbf{w}^*(\mathbf{c}) - 2\mathbf{w}^*(2\hat{\mathbf{c}} - \mathbf{c}) \in \frac{\partial l_{\text{SPO}+}(\hat{\mathbf{c}}, \mathbf{c})}{\partial \hat{\mathbf{c}}}$$

- Elmachtoub, A. N., & Grigas, P. (2021). Smart “predict, then optimize”. Management Science.

Perturbed Method

Use random perturbation to deal with the cost vector \hat{c} .

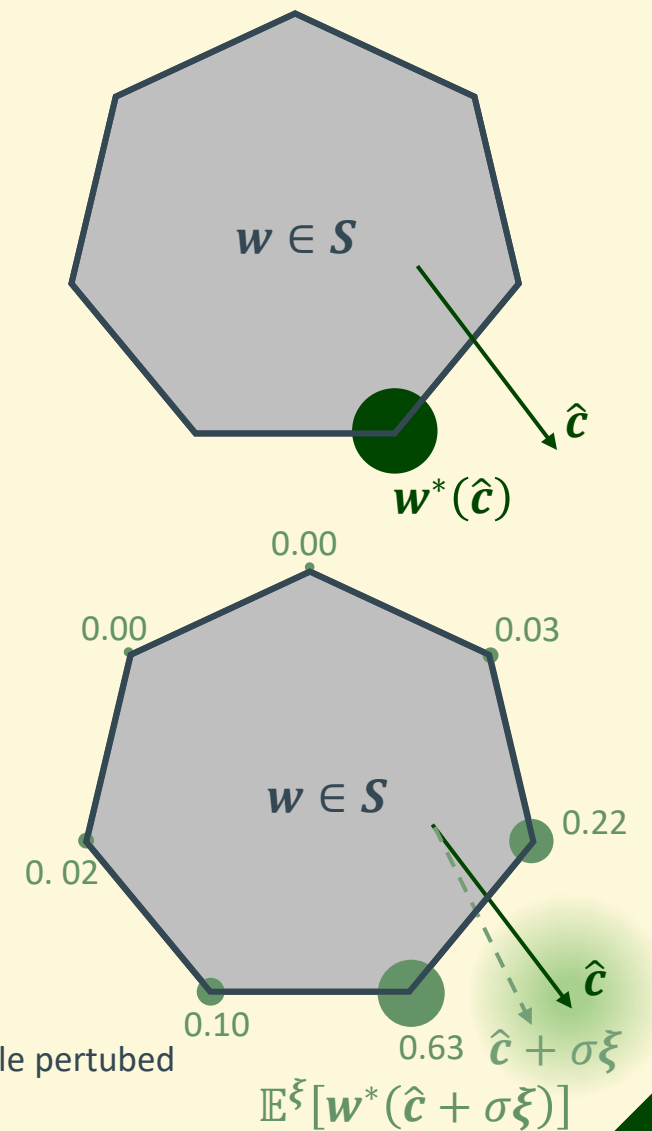


- Berthet, Q., Blondel, M., Teboul, O., Cuturi, M., Vert, J. P., & Bach, F. (2020). Learning with differentiable perturbed optimizers. Advances in neural information processing systems, 33, 9508-9519.
- Dalle, G., Baty, L., Bouvier, L., & Parmentier, A. (2022). Learning with combinatorial optimization layers: a probabilistic approach. arXiv preprint arXiv:2207.13513.c

Perturbed Method

Use random perturbation to deal with the cost vector \hat{c} .

The expected value $\mathbb{E}^{\xi}[\mathbf{w}^*(\hat{c} + \sigma\xi)]$ replaces $\mathbf{w}^*(\hat{c})$, which is the weighted average (convex combination) of the extreme points of feasible region.



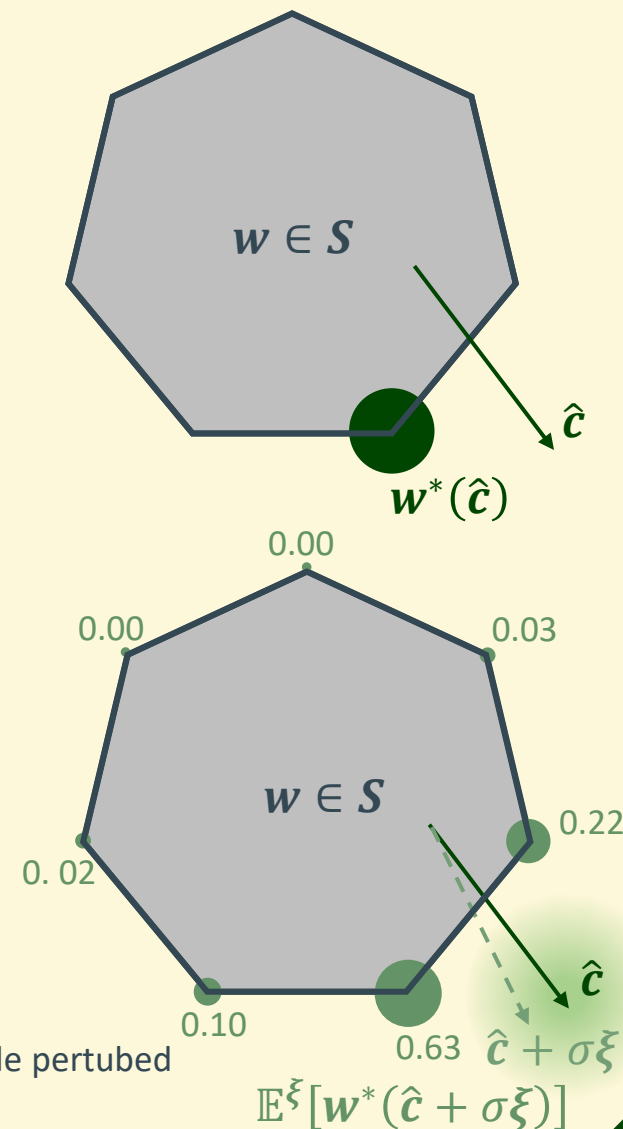
- Berthet, Q., Blondel, M., Teboul, O., Cuturi, M., Vert, J. P., & Bach, F. (2020). Learning with differentiable perturbed optimizers. Advances in neural information processing systems, 33, 9508-9519.
- Dalle, G., Baty, L., Bouvier, L., & Parmentier, A. (2022). Learning with combinatorial optimization layers: a probabilistic approach. arXiv preprint arXiv:2207.13513.c

Perturbed Method

Use random perturbation to deal with the cost vector \hat{c} .

A random perturbation $\xi \sim \mathcal{N}(\mathbf{0}, \mathbf{1})$

The expected value $\mathbb{E}^{\xi}[\mathbf{w}^*(\hat{c} + \sigma\xi)]$ replaces $\mathbf{w}^*(\hat{c})$, which is the weighted average (convex combination) of the extreme points of feasible region.



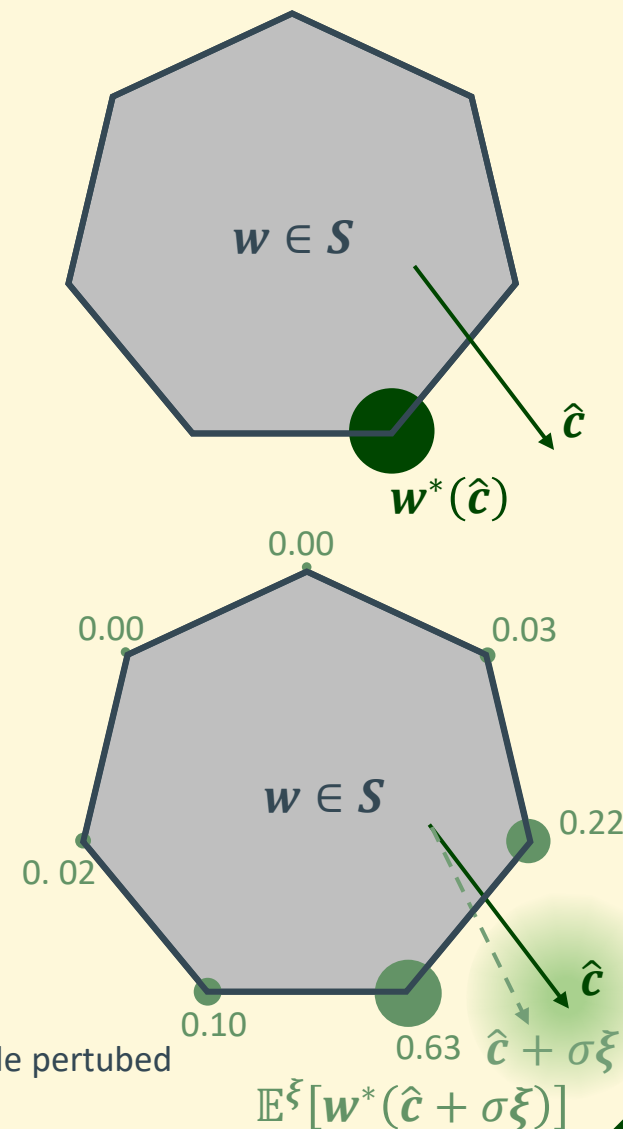
- Berthet, Q., Blondel, M., Teboul, O., Cuturi, M., Vert, J. P., & Bach, F. (2020). Learning with differentiable perturbed optimizers. Advances in neural information processing systems, 33, 9508-9519.
- Dalle, G., Baty, L., Bouvier, L., & Parmentier, A. (2022). Learning with combinatorial optimization layers: a probabilistic approach. arXiv preprint arXiv:2207.13513.c

Perturbed Method

Use random perturbation to deal with the cost vector \hat{c} .

The expected value $\mathbb{E}^{\xi}[\mathbf{w}^*(\hat{c} + \sigma\xi)]$ replaces $\mathbf{w}^*(\hat{c})$, which is the weighted average (convex combination) of the extreme points of feasible region.

How to calculate the expectation?



- Berthet, Q., Blondel, M., Teboul, O., Cuturi, M., Vert, J. P., & Bach, F. (2020). Learning with differentiable perturbed optimizers. Advances in neural information processing systems, 33, 9508-9519.
- Dalle, G., Baty, L., Bouvier, L., & Parmentier, A. (2022). Learning with combinatorial optimization layers: a probabilistic approach. arXiv preprint arXiv:2207.13513.c

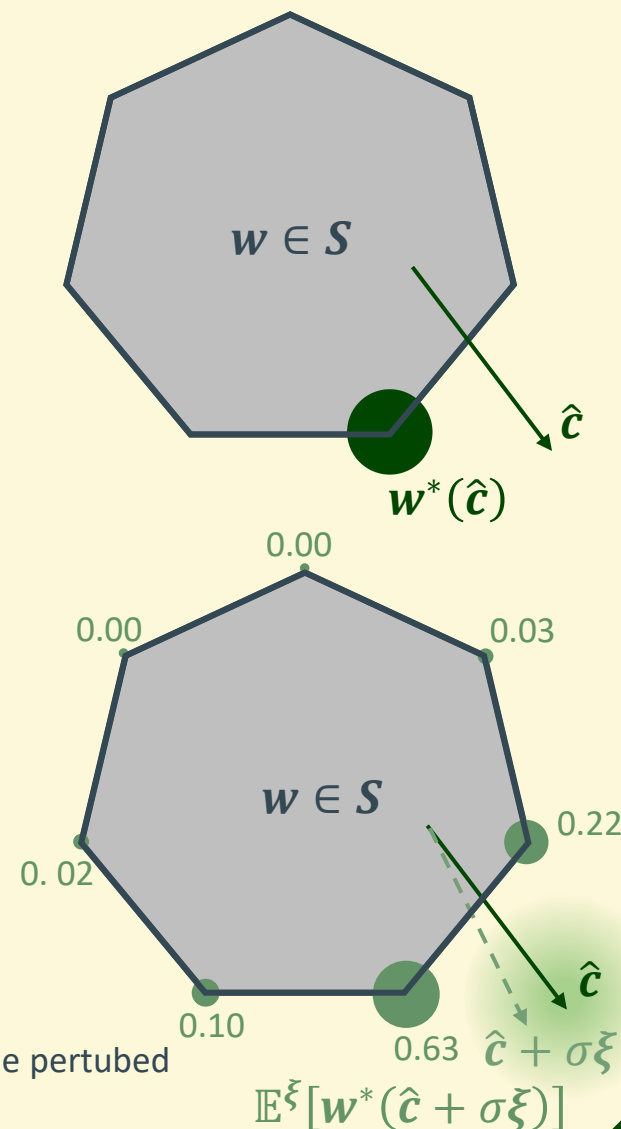
Perturbed Method

Use random perturbation to deal with the cost vector $\hat{\mathbf{c}}$.

The expected value $\mathbb{E}^{\xi}[\mathbf{w}^*(\hat{\mathbf{c}} + \sigma\xi)]$ replaces $\mathbf{w}^*(\hat{\mathbf{c}})$, which is the weighted average (convex combination) of the extreme points of feasible region.

How to calculate the expectation?

$$\mathbb{E}^{\xi}[\mathbf{w}^*(\hat{\mathbf{c}} + \sigma\xi)] \approx \frac{1}{K} \sum_{\kappa}^K \mathbf{w}^*(\hat{\mathbf{c}} + \sigma\xi_{\kappa})$$



- Berthet, Q., Blondel, M., Teboul, O., Cuturi, M., Vert, J. P., & Bach, F. (2020). Learning with differentiable perturbed optimizers. Advances in neural information processing systems, 33, 9508-9519.
- Dalle, G., Baty, L., Bouvier, L., & Parmentier, A. (2022). Learning with combinatorial optimization layers: a probabilistic approach. arXiv preprint arXiv:2207.13513.c

Perturbed Method

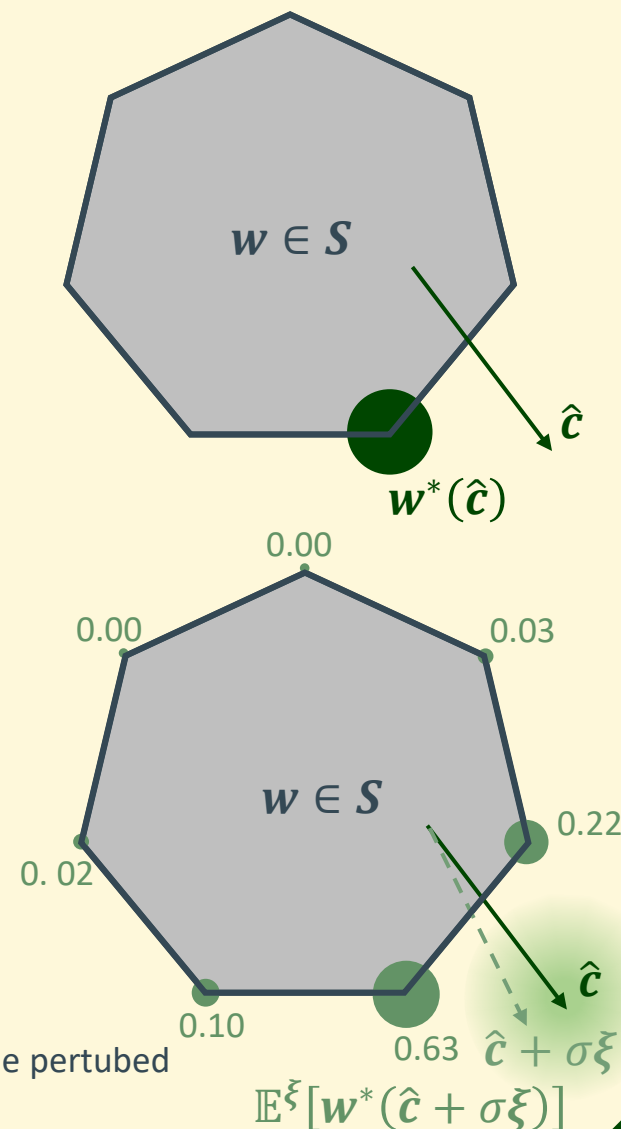
Use random perturbation to deal with the cost vector $\hat{\mathbf{c}}$.

The expected value $\mathbb{E}^{\xi}[\mathbf{w}^*(\hat{\mathbf{c}} + \sigma\xi)]$ replaces $\mathbf{w}^*(\hat{\mathbf{c}})$, which is the weighted average (convex combination) of the extreme points of feasible region.

How to calculate the expectation?

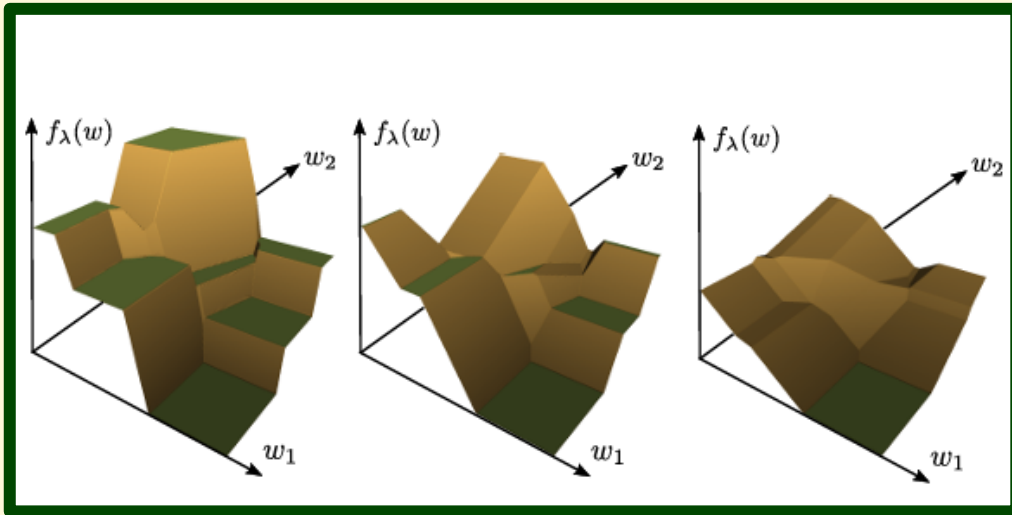
$$\mathbb{E}^{\xi}[\mathbf{w}^*(\hat{\mathbf{c}} + \sigma\xi)] \approx \frac{1}{K} \sum_{\kappa}^K \mathbf{w}^*(\hat{\mathbf{c}} + \sigma\xi_{\kappa})$$

A non-negativity requirement for $\hat{\mathbf{c}}$: multiplication perturbation



- Berthet, Q., Blondel, M., Teboul, O., Cuturi, M., Vert, J. P., & Bach, F. (2020). Learning with differentiable perturbed optimizers. Advances in neural information processing systems, 33, 9508-9519.
- Dalle, G., Baty, L., Bouvier, L., & Parmentier, A. (2022). Learning with combinatorial optimization layers: a probabilistic approach. arXiv preprint arXiv:2207.13513.c

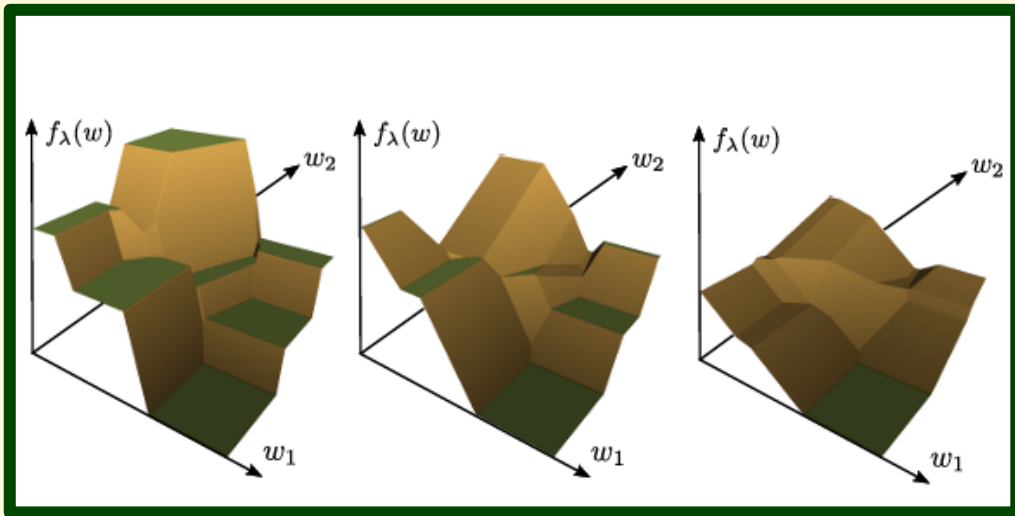
Black-box Method



- **Differentiable Black-box:** Perform continuous interpolation on the piecewise constant loss function to transform it into a piecewise linear function.

- Pogančić, M. V., Paulus, A., Musil, V., Martius, G., & Rolinek, M. (2019, September). Differentiation of blackbox combinatorial solvers. In International Conference on Learning Representations.
- Sahoo, S. S., Paulus, A., Vlastelica, M., Musil, V., Kuleshov, V., & Martius, G. (2022). Backpropagation through combinatorial algorithms: Identity with projection works. arXiv preprint arXiv:2205.15213.

Black-box Method



- **Differentiable Black-box:** Perform continuous interpolation on the piecewise constant loss function to transform it into a piecewise linear function.
- **Negative Identity:** Replace the solver gradient $\frac{\partial w^*(\hat{c})}{\partial \hat{c}}$ with the negative unit matrix $-I$.

Contrastive & Ranking Method

During the training process, we can naturally collect a large number of feasible solutions, forming a solution set Γ .

- Mulamba, M., Mandi, J., Diligenti, M., Lombardi, M., Bucarey, V., & Guns, T. (2021). Contrastive losses and solution caching for predict-and-optimize. Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence.
- Mandi, J., Bucarey, V., Mulamba, M., & Guns, T. (2022). Decision-focused learning: through the lens of learning to rank. Proceedings of the 39th International Conference on Machine Learning.

Contrastive & Ranking Method

During the training process, we can naturally collect a large number of feasible solutions, forming a solution set Γ .

- **Contrastive Method:**

Take the subset of suboptimal solutions $\Gamma \setminus \mathbf{w}^*(\mathbf{c})$ as **negative samples**, to maximize the difference between the **optimal solution** and the **suboptimal solutions**.

$$l_{NCE}(\hat{\mathbf{c}}, \mathbf{c}) = \frac{1}{|\Gamma| - 1} \sum_{\mathbf{w}^\gamma \in \Gamma \setminus \mathbf{w}^*(\mathbf{c})} (\hat{\mathbf{c}}^T \mathbf{w}^*(\mathbf{c}) - \hat{\mathbf{c}}^T \mathbf{w}^\gamma)$$

- Mulamba, M., Mandi, J., Diligenti, M., Lombardi, M., Bucarey, V., & Guns, T. (2021). Contrastive losses and solution caching for predict-and-optimize. Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence.
- Mandi, J., Bucarey, V., Mulamba, M., & Guns, T. (2022). Decision-focused learning: through the lens of learning to rank. Proceedings of the 39th International Conference on Machine Learning.

Contrastive & Ranking Method

During the training process, we can naturally collect a large number of feasible solutions, forming a solution set Γ .

- **Contrastive Method:**

Take the subset of suboptimal solutions $\Gamma \setminus \mathbf{w}^*(c)$ as **negative samples**, to maximize the difference between the **optimal solution** and the **suboptimal solutions**.

$$l_{NCE}(\hat{\mathbf{c}}, c) = \frac{1}{|\Gamma| - 1} \sum_{\mathbf{w}^\gamma \in \Gamma \setminus \mathbf{w}^*(c)} (\hat{\mathbf{c}}^T \mathbf{w}^*(c) - \hat{\mathbf{c}}^T \mathbf{w}^\gamma)$$

- **Ranking Method:**

Transform the predict-then-optimize task as Learning to Rank, with the **objective value** (such as $\hat{\mathbf{c}}^T \mathbf{w}$) as score, in order to correctly rank the subset of feasible solutions Γ .

Software



An open-sourced library to facilitate predict-then-optimize, bridging the gap between optimization and machine learning.

Software

Benchmarks



Shortest Path



Knapsack



Traveling Salesperson

Opt Modeling



Any custom
solver

ML Modeling



Algorithms

- SPO+
- DBB
- NID
- DPO
- PFYL
- NCE
- LTR

Autograd Function

Algorithms

- SPO+
- DBB
- NID
- DPO
- PFYL
- NCE
- LTR

`pyepo.func.perturbedFenchelYoung` allows us to set a Fenchel-Young loss for training, which requires parameters:

- `optmodel` : a PyEPO optimization model
- `n_samples` : number of Monte-Carlo samples
- `sigma` : the amplitude of the perturbation for costs
- `processes` : number of processors for multi-thread, 1 for single-core, 0 for all of the cores
- `seed` : random state seed for perturbations

```
import pyepo

# init SPO+ loss
spop = pyepo.func.SPOPlus(optmodel, processes=2)
# init PFY loss
pfy = pyepo.func.perturbedFenchelYoung(optmodel, n_samples=3, sigma=1.0, processes=2)
# init NCE loss
nce = pyepo.func.NCE(optmodel, processes=2, solve_ratio=0.05, dataset=dataset_train)
```

Modeling

```
import gurobipy as gp
from gurobipy import GRB
from pyepo.model.grb import optGrbModel

class myOptModel(optGrbModel):
    def _getModel(self):
        # ceate a model
        m = gp.Model()
        # variables
        x = m.addVars(5, name="x", vtype=GRB.BINARY)
        # sense
        m.modelSense = GRB.MAXIMIZE
        # constraints
        m.addConstr(3*x[0]+4*x[1]+3*x[2]+6*x[3]+4*x[4]<=12)
        m.addConstr(4*x[0]+5*x[1]+2*x[2]+3*x[3]+5*x[4]<=10)
        m.addConstr(5*x[0]+4*x[1]+6*x[2]+2*x[3]+3*x[4]<=15)
        return m, x

optmodel = myOptModel()
```

$$\max_w \sum_{i=0}^4 c_i w_i$$

$$s.t. \quad 3w_0 + 4w_1 + 3w_2 + 6w_3 + 4w_4 \leq 12$$

$$4w_0 + 5w_1 + 2w_2 + 3w_3 + 5w_4 \leq 10$$

$$5w_0 + 4w_1 + 6w_2 + 2w_3 + 3w_4 \leq 10$$

$$w_0, w_1, w_2, w_3, w_4 \in \{0,1\}$$

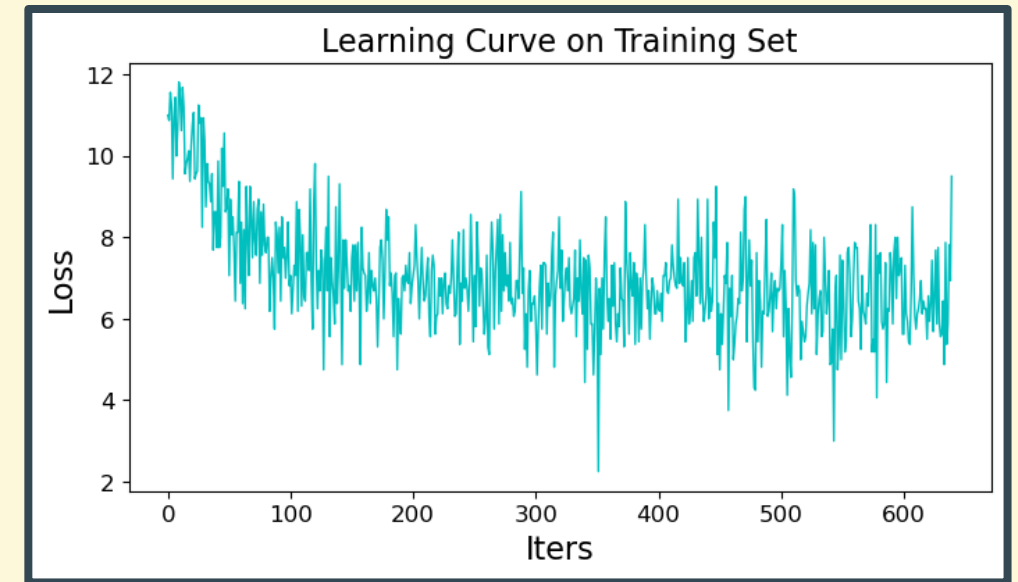



[Optimization Model](#)

End-to-End Training

```
# set adam optimizer
optimizer = torch.optim.Adam(reg.parameters(), lr=5e-3)

# train mode
reg.train()
for epoch in range(5):
    # load data
    for i, data in enumerate(loader_train):
        x, c, w, z = data # feat, cost, sol, obj
        # cuda
        if torch.cuda.is_available():
            x, c, w, z = x.cuda(), c.cuda(), w.cuda(), z.cuda()
        # forward pass
        cp = reg(x)
        loss = pfy(cp, w)
        # backward pass
        optimizer.zero_grad()
        loss.backward()
        optimizer.step()
```



 [2D knapsack Visualization](#)

Thank You

