# **Predict-then-Optimize:** a tour of the state-of-the-art using PyEP

### Elias B. Khalil

Department of Mechanical & Industrial Engineering SCALE AI Research Chair in Data-Driven Algorithms for Modern Supply Chains



For complete references to related work, please check our arXiv manuscript linked in the GitHub.

UNIVERSITY OF TORONTO



**Bo** Tang



## Optimization with a linear objective



Solution: use appropriate algorithm depending on type of feasible set (MILP, MIQCP, CP, custom algorithms...)

## $\min\{c_1^\mathsf{T} w : w \in \mathcal{S}\}$ W $\min\{c_{\gamma}^{\mathsf{T}}w:w\in\mathcal{S}\}$ $\mathcal{W}$ Same **Same** Decision variables Feasible set $\min\{c_n^\mathsf{T} w : w \in \mathcal{S}\}$ $\mathcal{W}$

The *n* cost vectors could be completely unrelated



## $\min\{c_1^\mathsf{T} w : w \in \mathcal{S}\}$ $\mathcal{W}$ $\min\{c_{\gamma}^{\mathsf{T}}w:w\in\mathcal{S}\}$ W Same **Same** Decision variables Feasible set $\min\{c_n^\mathsf{T} w : w \in \mathcal{S}\}$ $\mathcal{W}$

#### **Instance feature vector** (Observed)

 $x_1 \in \mathbb{R}^p$  $X_{\gamma}$ 

Function g is parametrized by a vector  $\theta$ .

# $c_i \equiv g(x_i; \theta) \in \mathbb{R}^d$

#### $\min\{c_1^\mathsf{T} w : w \in \mathcal{S}\}$ $\mathcal{W}$ $\min\{c_{\gamma}^{\mathsf{T}}w:w\in\mathcal{S}\}$ W Same

**Same** Decision variables

Feasible set

## $\min\{c_n^{\mathsf{I}}w:w\in\mathcal{S}\}$ $\mathcal{W}$

#### **Instance feature vector** (Observed)

 $x_1 \in \mathbb{R}^p$ X

 $c_i \equiv g(x_i; \theta) \in \mathbb{R}^d$ 

### $\min\{c_1^\mathsf{T} w : w \in \mathcal{S}\}$ W $\min\{c_{\gamma}^{\mathsf{T}}w:w\in\mathcal{S}\}$ $\mathcal{W}$ Same **Same** Decision variables Feasible set

 $\min\{c_n^\mathsf{T} w : w \in \mathcal{S}\}$  $\mathcal{W}$ 

#### If you know $\theta$ then you can evaluate $g(x_i; \theta)$ to obtain $c_i$ and optimize with your favorite method

**Instance feature vector** (Observed)



If you know  $\theta$  then you can evaluate  $g(x_i; \theta)$  to obtain  $c_i$  and optimize with your favorite method

 $\hat{c}_i \equiv g(x_i; \theta) \in \mathbb{R}^d$  $\min\{\hat{c}_1 | w : w \in \mathcal{S}\}$ W  $\min\{\hat{c}_2^\mathsf{T}w:w\in\mathcal{S}\}$ W Same **Same** Decision variables Feasible set

 $\min\{\hat{c}_n^\mathsf{T} w : w \in \mathcal{S}\}$  $\mathcal{W}$ 

What if you don't know  $\theta$  and you only observe  $x_i$ ?

#### **Instance feature vector** (Observed)







7



#### **Assume these** *n* **instances are from** the past, i.e., "training" instances

## $\hat{c}_i \equiv g(x_i; \theta) \in \mathbb{R}^d$ $\min\{\hat{c}_1^\mathsf{T} w : w \in \mathcal{S}\}$ $\mathcal{W}$ $\min\{\hat{c}_{\gamma}^{\mathsf{T}}w:w\in\mathcal{S}\}$ $\mathcal{W}$ Same **Same** Decision variables Feasible set

 $\min\{\hat{c}_n^\mathsf{T} w : w \in \mathcal{S}\}$ W

**True cost vector Instance feature (Observed** <u>for training</u> vector instances only) (Observed)  $C_1 \in \mathbb{R}^d$  $x_1 \in \mathbb{R}^p$  $C_{\mathcal{I}}$  $X_{\mathcal{I}}$ 



#### Assume these *n* instances are from the past, i.e., "training" instances

#### $\hat{c}_i \equiv g(x_i; \theta) \in \mathbb{R}^d$ $\min\{\hat{c}_1^\mathsf{T} w : w \in \mathcal{S}\}$ W $\min\{\hat{c}_{\gamma}^{\mathsf{T}}w:w\in\mathcal{S}\}$ W Same **Same** Decision variables Feasible set

### $\min\{\hat{c}_n^\mathsf{T} w : w \in \mathcal{S}\}$ W









## Predict-then-Optimize Training

#### Dataset $\mathcal{D}$ with data points (x, c)



X

Prediction model  $g(\mathbf{x}; \boldsymbol{\theta})$ with parameters  $\boldsymbol{\theta}$ 

## Loss function $l(\cdot)$ to measure prediction error



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





## **Predict-then-Optimize** Training and Test-time inference

#### Dataset $\mathcal{D}$ with data points ( $\mathbf{x}, \mathbf{c}$ )



X

Prediction model  $g(\mathbf{x}; \boldsymbol{\theta})$ with parameters  $\boldsymbol{\theta}$ 

## Loss function $l(\cdot)$ to measure prediction error







Google Maps, Montréal, Quebéc, Canada

## **Predict-then-Optimize s-t** shortest path: training data



## **Predict-then-Optimize** s-t shortest path: full pipeline



Loss function  $l(\cdot)$ to measure prediction error

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



**Predicted travel times** <=> s-t shortest path cost vector



## **Predict-then-Optimize** s-t shortest path: full pipeline



Loss function  $l(\cdot)$ to measure prediction error

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



**Predicted travel times** <=> s-t shortest path cost vector



## **Predict-then-Optimize** s-t shortest path: full pipeline



Loss function  $l(\cdot)$ to measure prediction error

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



**Predicted travel times** <=> s-t

shortest path cost vector

Squared error between  $C_i$ and  $\hat{c}_i$ 











#### Least-squares regression on dataset of (x, c) pairs



S

t

Edge 2





#### Least-squares regression on dataset of (x, c) pairs



t

S

Edge 2







#### Least-squares regression on dataset of (x, c) pairs

S

Edge 2



#### "Smart Predict then Optimize"





## **End-to-End Predict-then-Optimize** Training



dataset  $\mathcal{D}$  with data points  $(\mathbf{x}, \mathbf{c})$  or  $(\mathbf{x}, \mathbf{w}_{\mathbf{c}}^*)$  machine learning model  $g(\mathbf{x}, \boldsymbol{\theta})$ with parameters  $\boldsymbol{\theta}$ 

#### **Decision error**



optimization solver  $\boldsymbol{w}_{\boldsymbol{c}}^{*} = \operatorname{argmin} \boldsymbol{\hat{c}}^{T} \boldsymbol{w}$ wes



loss function  $l(\cdot)$ to measure decision error

PARTIE STARTING BI SAL RIG SA ST. BALG & S.D.





dataset  $\mathcal{D}$  with data points  $(\mathbf{x}, \mathbf{c})$  or  $(\mathbf{x}, \mathbf{w}_{\mathbf{c}}^*)$ 

#### Algorithm 1 End-to-end Gradient Descent

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

8:



dataset  $\mathcal{D}$  with data points  $(\mathbf{x}, \mathbf{c})$  or  $(\mathbf{x}, \mathbf{w}_{\mathbf{c}}^*)$ 

#### Algorithm 1 End-to-end Gradient Descent

- **Require:** coefficient matrix A, right-hand side b, data  $\mathcal{D}$ 1: Initialize predictor parameters  $\boldsymbol{\theta}$  for predictor  $g(\boldsymbol{x};\boldsymbol{\theta})$
- 2: for epochs do
- 3: for each batch of training data (x, c) do
- 4:
- Predict cost using predictor  $\hat{c} := g(x; \theta)$ 5:
- 6:
- Forward pass to compute decision loss  $l(\hat{c}, \cdot)$ 7:
- end for 9:
- 10: end for

8:

Sample batch of the cost vectors c with the corresponding features x

Forward pass to compute optimal solution  $w^*(\hat{c}) := \operatorname{argmin}_{w \in S} \hat{c}^\intercal w$ 

Backward pass from loss  $l(\hat{c}, \cdot)$  to update parameters  $\theta$  with gradient



# $rac{\partial l(\hat{c},\cdot)}{\partial oldsymbol{ heta}}$ .

#### Algorithm 1 End-to-end Gradient Descent

- **Require:** coefficient matrix A, right-hand side b, data  $\mathcal{D}$ 1: Initialize predictor parameters  $\boldsymbol{\theta}$  for predictor  $g(\boldsymbol{x};\boldsymbol{\theta})$ 
  - 2: for epochs do
  - for each batch of training data (x, c) do 3:
- 4:
- Predict cost using predictor  $\hat{c} := g(x; \theta)$ 5:
- Forward pass to compute optimal solution  $w^*(\hat{c}) := \operatorname{argmin}_{w \in S} \hat{c}^\intercal w$ 6: Forward pass to compute decision loss  $l(\hat{c}, \cdot)$ 7:

  - Backward pass from loss  $l(\hat{c}, \cdot)$  to update parameters  $\theta$  with gradient
- end for 9:
- 10: end for

8:

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

Sample batch of the cost vectors c with the corresponding features x



# $\partial l(\hat{oldsymbol{c}},\cdot)$

#### Algorithm 1 End-to-end Gradient Descent

- **Require:** coefficient matrix A, right-hand side b, data  $\mathcal{D}$ 1: Initialize predictor parameters  $\boldsymbol{\theta}$  for predictor  $g(\boldsymbol{x};\boldsymbol{\theta})$
- 2: for epochs do
- 3: for each batch of training data (x, c) do
- Sample batch of the cost vectors c with the corresponding features x4: Predict cost using predictor  $\hat{c} := g(x; \theta)$ 5:
- Forward pass to compute optimal solution  $w^*(\hat{c}) := \operatorname{argmin}_{w \in S} \hat{c}^\intercal w$ 6: Forward pass to compute decision loss  $l(\hat{c}, \cdot)$ 7:
- - Backward pass from loss  $l(\hat{c}, \cdot)$  to update parameters  $\theta$  with gradient
- end for 9:
- 10: end for

8:





## Trickier: how does the deci vary with the predicted c

# $rac{\partial l(\hat{c}, \cdot)}{\partial \theta} =$

#### Algorithm 1 End-to-end Gradient Descent

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

ision loss	Easy: gradient of predic
costs?	costs to model parameter
$\partial l(\hat{oldsymbol{c}},\cdot)$	$\partial \hat{c}$
$\partial \hat{c}$	$\overline{\partial \theta}$





# **Trickier:** how does the decision loss vary with the predicted costs? $\frac{\partial l(\hat{\boldsymbol{c}},\cdot)}{\partial \boldsymbol{\theta}} = \frac{\partial l(\hat{\boldsymbol{c}},\cdot)}{\partial \hat{\boldsymbol{c}}} \frac{\partial \hat{\boldsymbol{c}}}{\partial \boldsymbol{\theta}}$

#### Algorithm 1 End-to-end Gradient Descent

- **Require:** coefficient matrix A, right-hand side b, data  $\mathcal{D}$ 1: Initialize predictor parameters  $\boldsymbol{\theta}$  for predictor  $g(\boldsymbol{x};\boldsymbol{\theta})$ 
  - 2: for epochs do
  - for each batch of training data (x, c) do 3:
  - 4:
  - Predict cost using predictor  $\hat{c} := g(x; \theta)$ 5:
  - 6:
  - Forward pass to compute decision loss  $l(\hat{c}, \cdot)$ 7:
- end for 9:
- 10: end for

8:

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

Sample batch of the cost vectors c with the corresponding features x

Forward pass to compute optimal solution  $w^*(\hat{c}) := \operatorname{argmin}_{w \in S} \hat{c}^\intercal w$ 

Backward pass from loss  $l(\hat{c}, \cdot)$  to update parameters  $\theta$  with gradient





# Trickier: how does the decision vary with the predicted costs $rac{\partial l(\hat{m{c}},\cdot)}{\partial m{ heta}} = rac{\partial l(\hat{m{c}},\cdot)}{\partial \hat{m{c}}} rac{\partial \hat{m{c}}}{\partial m{ heta}}$

#### Algorithm 1 End-to-end Gradient Descent

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

8:

In loss **Crux:** how does the optimus  
change with the predicted co
$$l_{\text{Regret}}(\hat{\boldsymbol{c}}, \boldsymbol{c}) = \boldsymbol{c}^{\mathsf{T}} \boldsymbol{w}^*(\hat{\boldsymbol{c}}) - z^*$$



## SPO+: a principled and effective method

Elmachtoub, Adam N., and Paul Grigas. "Smart "predict, then optimize"." Management Science 68.1 (2022): 9-26.

#### Algorithm 1 End-to-end Gradient Descent

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

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



## **SPO+: a principled and effective method**

Elmachtoub, Adam N., and Paul Grigas. "Smart "predict, then optimize"." Management Science 68.1 (2022): 9-26.

$$l_{ ext{Regret}}(\hat{m{c}},m{c}) = m{c}^{\mathsf{T}}m{w}^*(\hat{m{c}}) - z^*(m{c})$$
 is not differentiable w.r.t.  $\hat{m{c}}$ 

$$l_{SPO+}(\hat{\boldsymbol{c}},\boldsymbol{c}) = -\min_{\boldsymbol{w}\in S} \{ (2\hat{\boldsymbol{c}} - \boldsymbol{c})^{\mathsf{T}}\boldsymbol{w} \} + 2\hat{\boldsymbol{c}}^{\mathsf{T}}\boldsymbol{w}^{*}(\boldsymbol{c}) - z^{*}(\boldsymbol{c})$$

#### Algorithm 1 End-to-end Gradient Descent

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

- 1: Initialize predictor parameters  $\boldsymbol{\theta}$  for predictor  $g(\boldsymbol{x}; \boldsymbol{\theta})$
- 2: for epochs do

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

9: end for

10: **end for** 

is a convex upper **bound on regret and** has subgradients



## SPO+: a principled and effective method

Elmachtoub, Adam N., and Paul Grigas. "Smart "predict, then optimize"." Management Science 68.1 (2022): 9-26.

$$l_{ ext{Regret}}(\hat{m{c}},m{c}) = m{c}^{\mathsf{T}}m{w}^*(\hat{m{c}}) - z^*(m{c})$$
 is not differentiable w.r.t.  $\hat{m{c}}$ 

$$l_{SPO+}(\hat{\boldsymbol{c}},\boldsymbol{c}) = -\lim_{\boldsymbol{w}\in S} \{ (2\hat{\boldsymbol{c}} - \boldsymbol{c})^{\mathsf{T}}\boldsymbol{w} \} + 2\hat{\boldsymbol{c}}^{\mathsf{T}}\boldsymbol{w}^{*}(\boldsymbol{c}) - z^{*}(\boldsymbol{c})$$

#### Algorithm 1 End-to-end Gradient Descent

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

- 1: Initialize predictor parameters  $\boldsymbol{\theta}$  for predictor  $g(\boldsymbol{x}; \boldsymbol{\theta})$
- 2: for epochs do

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

9: end for

10: **end for** 

is a convex upper **bound on regret and** has subgradients

Computational overhead: To compute SPO+ loss, we need to solve an optimization problem in the forward pass; this is shared with other methods.





#### Least-squares regression on dataset of (x, c) pairs



S

Edge 2

#### SPO+



Edge 2

optimal







Google Maps, Montréal, Quebéc, Canada

## **Predict-then-Optimize s-t** shortest path: training data



# **Predict-then-Optimize** s-t shortest path: training data 8 am 10 am londay, 7:35 Al RRE Saw View topography and elevation View topography and elevation



#### Google Maps, Montréal, Quebéc, Canada

## **PFYL: Perturbed Fenchel-Young Loss**

Berthet, Quentin, et al. "Learning with differentiable perturbed optimizers." NeurIPS 2020.

# **Recall the required gradient:** $\frac{\partial l(\hat{c}, \cdot)}{\partial A} =$

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

A random perturbation  $(\sigma \xi_k)$  of

#### predicted costs $\hat{c}$



Average optimal solution over K randomly perturbed costs

## **PFYL: Perturbed Fenchel-Young Loss**

Berthet, Quentin, et al. "Learning with differentiable perturbed optimizers." NeurIPS 2020.

# **Recall the required gradient:** $\frac{\partial l(\hat{c}, \cdot)}{\partial A} =$

#### Notice the true costs c do not appear here.

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

A random perturbation  $(\sigma \xi_k)$  of

#### predicted costs $\hat{c}$



Average optimal solution over K randomly perturbed costs

## **PFYL: Perturbed Fenchel-Young Loss**

Berthet, Quentin, et al. "Learning with differentiable perturbed optimizers." NeurIPS 2020.

# **Recall the required gradient:** $\frac{\partial l(\hat{c}, \cdot)}{\partial A} =$

Notice the true costs c do not appear here.

**Computational overhead:** To compute the PFYL gradient, we need to solve Koptimization problems.

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

A random perturbation  $(\sigma \xi_k)$  of

#### predicted costs $\hat{c}$



Average optimal solution over K randomly perturbed costs

## Avoiding solver calls by Learning to Rank

Mandi, Jayanta, et al. "Decision-focused learning: through the lens of learning to rank." ICML, 2022. Mulamba, Maxime, et al. "Contrastive Losses and Solution Caching for Predict-and-Optimize." IJCAI, 2021.

# Recall that both SPO+ and PFYL made one or more calls to the solver in each forward pass!

	w1	w2	w3	Objective using prediction A	Objective using prediction B	True Objective
Sol 1	0	1	0	1	3	1
Sol 2	1	0	Ο	3	2	2
Sol 3	1	1	1	2	1	3

## Av

Mar Mulam

oidin ndi, Jayanta, e ba, Maxime, e	g solve et al. "Decision et al. "Contra	<b>Verc</b> a on-focused le stive Losses	alls b earning: throu and Solution Somew better	y Lea ugh the lens Caching for hat	of learning to rank." ICML, 2022. Predict-and-Optimize." IJCAI, 2021 <b>Terrible solution</b> ranking!				
	w1	w2	<b>w3</b>	Objective using prediction A	Objective using prediction B	True Objective			
Sol 1	0	1	0	1	3	1			
Sol 2	1	0	0	3	2	2			
Sol 3	1	1	1	2	1	3			

## Av

Mar Mulam

ndi, Jayanta, e	<b>G</b> SO et al. "Decision et al. "Contra	<b>Ver C</b> a on-focused le astive Losses	alls b earning: throu and Solution	<b>y Lea</b> Ugh the lens of Caching for	<b>rning</b> of learning to Predict-and-	<b>J to R</b> rank." ICML, Optimize." IJ	<b>ank</b> 2022. CAI, 2021.
· ·			Somew better	hat	Terri	ble solu anking!	ition
	w1	w2	w3	Objective using prediction	Objective using prediction	True Objective	
Sol 1	C	Key Idea oefficien	a: Learn Its that le	to predice ead to go	ot ood	1	
Sol 2		ankings	solution	s!	.eu	2	
Sol 3	1	1	1	2	1	3	

# Predict-then-Optimize: a tour of the state-of-the-art using PyEP

Manuscript describing PyEPO, the literature, and extensive experiments: https://arxiv.org/abs/2206.14234





#### Bo Tang



#### PyTorch-based End-to-End Predict-then-Optimize Tool

Search docs

#### CONTENTS:

Introduction

Installation

**Tutorial** 

Module

Reference

**API** Reference



#### Welcome to PyEPO's documentation!

This is the documentation of PyEP0 (PyTorch-based End-to-End Predict-then-Optimize Library Tool), which aims to provide end-to-end methods for predict-then-optimize tasks.

#### Sample Code

#### import random

import gurobipy as gp from gurobipy import GRB









#### ling ML modeling

# **O** PyTorch



Training algorithms

SPO+ SPO+ with relaxations

PFYL/DPO

**DBB** (Pogancic et al. [2019])



## Benchmark generation **Based on Elmachtoub & Grigas**

# $c_{ij} = \alpha \left( (\mathscr{B}x_i)_j + 3 \right)^{deg} \cdot \epsilon_{i,j}$



 $x_i \in \mathbb{R}^p$ 

i-th instance's feature vector (Observed)

deg

 $e_{i,j}$ 

Integer in [1,6] regulating how nonlinear the mapping is

 $\mathscr{B} \in \{0,1\}^{d \times n}$  $\mathscr{B}_{p,q} \sim \operatorname{Bernoulli}(0.5)$ 

Uniform random noise









#### Image inputs (RGB features)

	0.8	0.8	0.8	0.8	0.8	0.8	0.8	0.8	1.2	1.2	1.2	0.8	
	0.8	0.8	0.8	0.8	0.8	1.2	1.2	1.2	1.2	9.2	1.2	1.2	
	0.8	0.8	0.8	1.2	1.2	1.2	1.2	1.2	1.2	1.2	1.2	1.2	
	1.2	0.8	0.8	1.2	1.2	0.8	1.2	1.2	7.7	7.7	1.2	1.2	
	1.2	0.8	1.2	1.2	0.8	0.8	0.8	1.2	1.2	7.7	1.2	1.2	
	1.2	0.8	1.2	1.2	0.8	0.8	0.8	0.8	1.2	7.7	1.2	1.2	
	1.2	0.8	1.2	1.2	1.2	0.8	0.8	1.2	1.2	1.2	7.7	7.7	
	0.8	0.8	1.2	1.2	0.8	0.8	1.2	0.8	0.8	1.2	7.7	7.7	
	0.8	0.8	1.2	1.2	1.2	1.2	1.2	1.2	1.2	1.2	7.7	7.7	
Sec. 1997	0.8	0.8	1.2	9.2				1.2	1.2	7.7	7.7	7.7	
	0.8	0.8	1.2	9.2				1.2	1.2	7.7	7.7	7.7	
	0.8	0.8	1.2	9.2				1.2	1.2	1.2	7.7	7.7	

Based on Pogančić, Marin Vlastelica, et al. "Differentiation of blackbox combinatorial solvers." ICLR 2020.

#### True costs

#### Optimal solution (NW -> SE)

## **Optimization model (in Gurobi)**

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

class myModel(optGrbModel): def \_getModel(self): # ceate a model m = gp.Model()# varibles x = m.addVars(5, name="x", vtype=GRB.BINARY) # sense (must be minimize) 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 = myModel()

$$\max_{x} \sum_{i=0}^{4} c_{i}x_{i}$$
s.t.  $3x_{0} + 4x_{1} + 3x_{2} + 6x_{3} + 4x_{4} \le 12$   
 $4x_{0} + 5x_{1} + 2x_{2} + 3x_{3} + 5x_{4} \le 10$   
 $5x_{0} + 4x_{1} + 6x_{2} + 2x_{3} + 3x_{4} \le 15$   
 $\forall x_{i} \in \{0, 1\}$ 

## Creating a dataset based on features and true costs

impo	ort pyepo	
from	n torch.utils.data import	Da
# bu data	ild dataset aset = pyepo.data.dataset.	ol
# ge data	et data loader aloader = DataLoader(datas	et
for	x, c, w, z in dataloader:	
	<pre># a batch of features print(x)</pre>	
	<pre># a batch of true costs print(c)</pre>	
	<pre># a batch of true optimal</pre>	
	<pre># a batch of true optimal</pre>	
	print(z)	

ataLoader

ptDataset(optmodel, feats, costs)

t, batch\_size=32, shuffle=True)

solutions

objective values



## Creating an ML model with PyTorch

#### from torch import nn

# construct linear model
class LinearRegression(nn.Module):
 def \_\_init\_\_(self):
 super(LinearRegression, self).\_\_init\_\_()
 # size of input and output is the size of feature and cost
 self.linear = nn.Linear(num\_feat, len\_cost)
 def forward(self, x):
 out = self.linear(x)
 return out

# init model
predmodel = LinearRegression()

## End-to-end training!

import torch

# set SGD optimizer optimizer = torch.optim.SGD(predmodel.parameters(), lr=1e-3)

# training

for epoch in range(num\_epochs):

for x, c, w, z in dataloader: # forward pass

> cp = predmodel(x) # predict costs # backward pass

optimizer.zero\_grad() # reset gradients to 0 loss.backward() # compute gradients optimizer.step() # update model parameters



```
# iterare features, costs, solutions, and objective values
    loss = spop(cp, c, w, z).mean() # calculate SPO+ loss
```

## **Experiments with TSP20**

#### Vertical axis: average regret w.r.t. true OPT on unseen test instances



## **Experiments with TSP20**

#### Vertical axis: average regret w.r.t. true OPT on unseen test instances

2-stage methods: regress on true costs (no end-to-end training)



### **SPO+** and **PFYL**, **both using a linear model** perform the best. With more data, 2-stage Random Forest is competitive!



## **Regret-accuracy tradeoff**

Vertical axis: average regret w.r.t. true OPT on unseen test instances

### Horizontal axis: Mean-Squared Error on true costs

$$l_{\text{MSE}}(\hat{c}, c) = \frac{1}{2n} \sum_{i=1}^{n} \|\hat{c}_{i} - c_{i}\|_{2}^{2}$$



L	_	0	.5

#### Image inputs (RGB features)

	0.8	0.8	0.8	0.8	0.8	0.8	0.8	0.8	1.2	1.2	1.2	0.8	
	0.8	0.8	0.8	0.8	0.8	1.2	1.2	1.2	1.2	9.2	1.2	1.2	
	0.8	0.8	0.8	1.2	1.2	1.2	1.2	1.2	1.2	1.2	1.2	1.2	
	1.2	0.8	0.8	1.2	1.2	0.8	1.2	1.2	7.7	7.7	1.2	1.2	
	1.2	0.8	1.2	1.2	0.8	0.8	0.8	1.2	1.2	7.7	1.2	1.2	
	1.2	0.8	1.2	1.2	0.8	0.8	0.8	0.8	1.2	7.7	1.2	1.2	
	1.2	0.8	1.2	1.2	1.2	0.8	0.8	1.2	1.2	1.2	7.7	7.7	
	0.8	0.8	1.2	1.2	0.8	0.8	1.2	0.8	0.8	1.2	7.7	7.7	
and the second second	0.8	0.8	1.2	1.2	1.2	1.2	1.2	1.2	1.2	1.2	7.7	7.7	
	0.8	0.8	1.2	9.2				1.2	1.2	7.7	7.7	7.7	
	0.8	0.8	1.2	9.2				1.2	1.2	7.7	7.7	7.7	
	0.8	0.8	1.2	9.2				1.2	1.2	1.2	7.7	7.7	

Based on Pogančić, Marin Vlastelica, et al. "Differentiation of blackbox combinatorial solvers." ICLR 2020.

#### True costs

#### Optimal solution (NW -> SE)

#### Image inputs







Based on Pogančić, Marin Vlastelica, et al. "Differentiation of blackbox combinatorial solvers." ICLR 2020.

#### **ResNet-18**

#### Cost predictions

0.8	0.8	0.8	0.8	0.8	0.8	0.8	0.8	1.2	1.2	1.2	0.8
0.8	0.8	0.8	0.8	0.8	1.2	1.2	1.2	1.2	9.2	1.2	1.2
0.8	0.8	0.8	1.2	1.2	1.2	1.2	1.2	1.2	1.2	1.2	1.2
1.2	0.8	0.8	1.2	1.2	0.8	1.2	1.2	7.7	7.7	1.2	1.2
1.2	0.8	1.2	1.2	0.8	0.8	0.8	1.2	1.2	7.7	1.2	1.2
1.2	0.8	1.2	1.2	0.8	0.8	0.8	0.8	1.2	7.7	1.2	1.2
1.2	0.8	1.2	1.2	1.2	0.8	0.8	1.2	1.2	1.2	7.7	7.7
0.8	0.8	1.2	1.2	0.8	0.8	1.2	0.8	0.8	1.2	7.7	7.7
0.8	0.8	1.2	1.2	1.2	1.2	1.2	1.2	1.2	1.2	7.7	7.7
0.8	0.8	1.2	9.2				1.2	1.2	7.7		7.7
0.8	0.8	1.2	9.2				1.2	1.2	7.7		7.7
	~ ~	1.0					1.0	1.0	1.0		

Based on Pogančić, Marin Vlastelica, et al. "Differentiation of blackbox combinatorial solvers." ICLR-20.

**2-Stage**, **SPO+**, and **PFYL**, **all using a** truncated ResNet-18, perform the best. DBB, originally benchmarked on this dataset, is far worse.



Image inputs





l	ļ	)	re	ec		C	<b>T</b> I	0	n.	S
	0.8	0.8	0.8	0.8	0.8	1.2	1.2	1.2	0.8	
	0.8	0.8	1.2	1.2	1.2	1.2	9.2	1.2	1.2	
	1.2	1.2	1.2	1.2	1.2	1.2	1.2	1.2	1.2	
	1.2	1.2	0.8	1.2	1.2	7.7	7.7	1.2	1.2	
	1.2	0.8	0.8	0.8	1.2	1.2	7.7	1.2	1.2	
	1.2	0.8	0.8	0.8	0.8	1.2	7.7	1.2	1.2	
	1.2	1.2	0.8	0.8	1.2	1.2	1.2	7.7	7.7	
	1.2	0.8	0.8	1.2	0.8	0.8	1.2	7.7	7.7	
	1.2	1.2	1.2	1.2	1.2	1.2	1.2	7.7	7.7	
					1.2	1.2	7.7		7.7	
					1.2	1.2	7.7	7.7	7.7	
					1.2	1.2	1.2	7.7	7.7	Ĩ

Cos

## Summary

- Predict-then-Optimize is a highly practical paradigm.
- SPO+ and PFYL are very effective end-to-end learning methods.
- The "naive" 2-stage approach is sufficient training set is large.
- Open questions:
  - Predictions in the constraints; see Hu, Xinyi, Jasper CH Lee, and Jimmy HM Lee. "Branch &
  - Reducing training time; see work by Tias Guns and collaborators.
  - More applications; see work by B. Dilkina, M. Tambe, B. Wilder, H. Bastani

Manuscript: https://arxiv.org/abs/2206.14234



Learn with Post-hoc Correction for Predict+ Optimize with Unknown Parameters in Constraints." CPAIOR 2023.



Bo Tang



## Machine Learning for Integer Programming

#### Mixed-Integer Linear Programming

[TMLR-22] [AAAI-22] [NeurIPS-20]

> [IJCAI-17] [AAAI-16]

[AAAI-22] [NeurIPS-21]

SL: Supervised Learning RL: Reinforcement Learning GNN: Graph Neural Networks

SL + GNN

**SL** + Simple models

**Custom ML** 



#### ] ] ] ]

#### Branch

Schedule heuristics

Select nodes

• Warmstart solver

Detect backdoors

## **Machine Learning for Discrete Optimization**

#### Survey on GNN for CombOpt [JMLR 2023]

#### RL + GNN

#### SL + GNN

SL + Simple models

**Custom ML** 

MILP

# ColumnStochasticGraphMultiobjectiveGenerationProgramming Optimizationoptimization

[NeurIPS-22]

[TMLR-22] [AAAI-22] [NeurIPS-20]

> [IJCAI-17] [AAAI-16]

[AAAI-22] [NeurIPS-21] [TMLR-22] [NeurIPS-17]

#### [NeurIPS-22]





# Lab Colab:

## https://tinyurl.com/ACP23-PredictAndOptimize

# PyEPO Github/Docs: https://github.com/khalil-research/PyEPO