ECE 598 Computational Power Systems

Decision-focused analytics

Vladimir Dvorkin

University of Michigan

Last lecture recap

- Look around you and form teams of 2 people (1 min)
- Quickly review your notes or the slide deck (1 min)
- Share your three personal highlights with your partner (3 min)
- Get iClicker app ready

Closed-loop analytics

Recap: Why forecasting?

- Forecasting is the first step in **decision-making**
- Resolves (some) uncertainty of decision-making inputs
- Brings confidence to **decision-making**

Closed-loop analytics

Recap: Why forecasting?

- Forecasting is the first step in decision-making
- Resolves (some) uncertainty of decision-making inputs
- Brings confidence to decision-making



Decision feedback of the closed-loop analytics

- Inform forecast on the quality of decisions induced on that forecast
- Requires integration of forecast and decision-making optimization problems

Example: Wind power forecast & electricity pricing

Given a wind power forecast \widehat{w} , solve the DC-OPF problem:

$$\begin{array}{ll} \underset{\underline{p} \leqslant p \leqslant \overline{p}}{\text{minimize}} & p^\top Cp + c^\top p & \textit{generator dispatch cost} \\ \text{subject to} & \mathbf{1}^\top (p + \widehat{w} - d) = 0 : \lambda, & \textit{power balance condition} \\ & |F(p + \widehat{w} - d)| \leqslant \overline{f} : \lambda_{\overline{f}}, \lambda_{\underline{f}}, & \textit{power flow limits} \\ \end{array}$$

Example: Wind power forecast & electricity pricing

Given a wind power forecast \widehat{w} , solve the DC-OPF problem:

$$\begin{array}{ll} \underset{\underline{p} \leq p \leq \overline{p}}{\text{minimize}} & p^{\top} Cp + c^{\top} p & \text{generator dispatch cost} \\ \\ \text{subject to} & \mathbf{1}^{\top} (p + \widehat{w} - d) = 0 : \lambda, & \text{power balance condition} \\ & |F(p + \widehat{w} - d)| \leqslant \overline{f} : \lambda_{\overline{f}}, \lambda_{\underline{f}}, & \text{power flow limits} \\ \end{array}$$

Location marginal prices (LMPs) are derived from the dual solution:

$$\pi(\widehat{\mathbf{w}}) = \underbrace{\lambda_b \cdot \mathbf{1}}_{\text{uniform price}} - \underbrace{\mathbf{F}^{\top}(\lambda_{\overline{\mathbf{f}}} - \lambda_{\underline{\mathbf{f}}})}_{\text{adjustment due to congestion}}$$

which are unique w.r.t forecast \hat{w} under reasonable assumptions!

Example: Wind power forecast & electricity pricing

Given a wind power forecast \hat{w} , solve the DC-OPF problem:

$$\begin{array}{ll} \underset{\underline{p} \leq p \leq \overline{p}}{\text{minimize}} & p^{\top} Cp + c^{\top} p & \text{generator dispatch cost} \\ \\ \text{subject to} & \mathbf{1}^{\top} (p + \widehat{w} - d) = 0 : \lambda, & \text{power balance condition} \\ & |F(p + \widehat{w} - d)| \leqslant \overline{f} : \lambda_{\overline{f}}, \lambda_{\underline{f}}, & \text{power flow limits} \\ \end{array}$$

Location marginal prices (LMPs) are derived from the dual solution:

$$\pi(\widehat{\mathbf{w}}) = \underbrace{\lambda_b \cdot \mathbf{1}}_{\text{uniform price}} - \underbrace{\mathsf{F}^{\top}(\lambda_{\overline{\mathsf{f}}} - \lambda_{\underline{\mathsf{f}}})}_{\text{adjustment due to congestion}}$$

which are unique w.r.t forecast \hat{w} under reasonable assumptions!

The LMP error is then defined as:

$$\delta \pi = \pi(\widehat{\mathbf{w}}) - \pi(\mathbf{w})$$

gap between LMPs induced on forecast (\widehat{w}) and actual realization (w) of wind power.

Example: Wind power forecast & electricity pricing (cont'd)



- IEEE 118-Bus system with one installed wind farm (bus 39)
- We forecast wind power output using weather forecast (×2 forecast error)
- LMP errors varying across the grid due to congestion (Q1. Is it fair?)
- The slack bus has the smallest LMP error (Q2. Why?)

Integrating forecasting and decision-making

Forecast problem

Dataset $\mathcal{D} = \{(\varphi_1, \mathbf{w}_1), \dots, (\varphi_k, \mathbf{w}_k)\}$ with features φ and observations \mathbf{w}

 $\blacksquare Machine learning model <math>\mathbb{W}_{\theta} : \mathcal{F} \mapsto \mathcal{W} \text{ maps features to predictions, i.e.,}$

$$\mathbb{W}_{\boldsymbol{ heta}}(\widehat{\boldsymbol{arphi}}) = \widehat{\mathbf{w}}_{i}$$

where heta is a model parameter to be optimized

Prediction-focused learning:

$$\underset{\boldsymbol{\theta}}{\mathsf{minimize}} \quad \frac{1}{2k} \sum_{i=1}^{k} \| \mathbb{W}_{\boldsymbol{\theta}}(\boldsymbol{\varphi}_i) - \mathbf{w}_i \|_2^2 \, .$$

Decision problem

w-parametrized optimization problem:

 $\begin{array}{ll} \underset{\mathbf{x}}{\text{minimize}} & c(\mathbf{x}) \\ \text{subject to} & g(\mathbf{x}, \mathbf{w}) \leqslant \mathbf{0} \end{array}$

• $\mathbf{x} = (x_1, \dots, x_n) \in \mathbb{R}^n$ is the vector of decision variable

- $\mathbf{w} \in \mathbb{R}^m$ is the input parameter (e.g., forecast from $\mathbb{W}_{\boldsymbol{\theta}}$)
- Objective function c and constraint function g

Decision is the optimal solution $\mathbf{x}^{\star}(\mathbf{w})$, parametrized by \mathbf{w} .

Decision problem

w-parametrized optimization problem:

 $\begin{array}{ll} \underset{\mathbf{x}}{\text{minimize}} & c(\mathbf{x}) \\ \text{subject to} & g(\mathbf{x}, \mathbf{w}) \leqslant \mathbf{0} \end{array}$

x = $(x_1, \ldots, x_n) \in \mathbb{R}^n$ is the vector of decision variable

- $\mathbf{w} \in \mathbb{R}^m$ is the input parameter (e.g., forecast from $\mathbb{W}_{\boldsymbol{\theta}}$)
- Objective function c and constraint function g

Decision is the optimal solution $\mathbf{x}^{\star}(\mathbf{w})$, parametrized by \mathbf{w} .

Decision regret is the difference of two objective function values

 $c(\mathbf{x}^{\star}(\widehat{\mathbf{w}})) - c(\mathbf{x}^{\star}(\mathbf{w}))$

with prediction \widehat{w} and the actual realization w

Decision-focused learning

Minimize the average decision regret across k historical records

This calls for the following **bilevel** optimization structure

$$\begin{array}{ll} \underset{\boldsymbol{\theta}}{\text{minimize}} & \frac{1}{2k} \sum_{i=1}^{k} \| c(\mathbf{x}^{\star}(\widehat{\mathbf{w}}_{i})) - c(\mathbf{x}^{\star}(\mathbf{w}_{i})) \|_{2}^{2} \\ \text{subject to} & \widehat{\mathbf{w}}_{i} = \mathbb{W}_{\boldsymbol{\theta}}(\boldsymbol{\varphi}_{i}) & \text{upper level} \\ & \mathbf{x}^{\star}(\widehat{\mathbf{w}}_{i}) = \operatorname*{argmin}_{\mathbf{x}} c(\mathbf{x}) & \text{lower level} \\ & \text{subject to} & g(\mathbf{x}, \widehat{\mathbf{w}}_{i}) \leq \mathbf{0} & \forall i = 1, \dots, k \end{array}$$

- **Upper-level:** optimize the parameter of the ML model θ
- Lower-level: respond with decision \mathbf{x} to parameter $\boldsymbol{\theta}$ for each record *i*
- **Q1.** How is it different from the prediction-focused optimization?
- **Q2.** How to solve this bilevel optimization?

KKT-based reformulation

The Karush-Kuhn-Tucker conditions (KKTs) of the lower-level problem:

- **I** Primal feasibility: $g(\mathbf{x}^{\star}, \widehat{\mathbf{w}}_i) \leq \mathbf{0}$
- 2 Dual feasibility: $\lambda^{\star} \succeq 0$
- Lagrangian optimality:

$$abla_{\mathbf{x}} c(\mathbf{x}^{\star}) + {\boldsymbol{\lambda}^{\star}}^{ op}
abla_{\mathbf{x}} g(\mathbf{x}^{\star}, \widehat{\mathbf{w}}_{i}) = \mathbf{0}$$

4 Complementary slackness:
$$\lambda^{\star \top} g(\mathbf{x}^{\star}, \widehat{\mathbf{w}}_i) = 0$$

For convex optimization, the KKTs are necessary and sufficient for optimality

Thus, we can replace the lower-level problem with its KKTs!

KKT-based reformulation (cont'd)

Assume linear regression, i.e., $\mathbb{W}_{\theta}(\varphi_i) = \theta^{\top} \varphi_i$

KKT-based reformulation of the bilevel problem

$$\begin{array}{ll} \underset{\boldsymbol{\theta}, \mathbf{x}_{i}, \boldsymbol{\lambda}_{i}, \widehat{\mathbf{w}}_{i}}{\text{minimize}} & \frac{1}{2k} \sum_{i=1}^{k} \| c(\mathbf{x}_{i}) - c(\mathbf{x}^{\star}(\mathbf{w}_{i})) \|_{2}^{2} \\ \text{subject to} & \widehat{\mathbf{w}}_{i} = \boldsymbol{\theta}^{\top} \boldsymbol{\varphi}_{i}, \\ & \underset{\boldsymbol{g}(\mathbf{x}_{i}, \widehat{\mathbf{w}}_{i}) \leq \mathbf{0}, \quad \boldsymbol{\lambda}_{i} \geq \mathbf{0}, \\ \nabla_{\mathbf{x}} c(\mathbf{x}_{i}) + \boldsymbol{\lambda}^{\top} \nabla_{\mathbf{x}_{i}} g(\mathbf{x}, \widehat{\mathbf{w}}_{i}) = \mathbf{0}, \\ & \boldsymbol{\lambda}_{i}^{\top} g(\mathbf{x}_{i}, \widehat{\mathbf{w}}_{i}) = \mathbf{0}, \quad \forall i = 1, \dots, k \end{array}$$

KKT-based reformulation (cont'd)

Assume linear regression, i.e., $\mathbb{W}_{\theta}(\varphi_i) = \theta^{\top} \varphi_i$

KKT-based reformulation of the bilevel problem

$$\begin{array}{ll} \underset{\boldsymbol{\theta}, \mathbf{x}_{i}, \boldsymbol{\lambda}_{i}, \widehat{\mathbf{w}}_{i}}{\text{minimize}} & \frac{1}{2k} \sum_{i=1}^{k} \| c(\mathbf{x}_{i}) - c(\mathbf{x}^{\star}(\mathbf{w}_{i})) \|_{2}^{2} \\ \text{subject to} & \widehat{\mathbf{w}}_{i} = \boldsymbol{\theta}^{\top} \boldsymbol{\varphi}_{i}, \\ & \underset{\boldsymbol{g}(\mathbf{x}_{i}, \widehat{\mathbf{w}}_{i}) \leq \mathbf{0}, \quad \boldsymbol{\lambda}_{i} \geq \mathbf{0}, \\ & \nabla_{\mathbf{x}} c(\mathbf{x}_{i}) + \boldsymbol{\lambda}^{\top} \nabla_{\mathbf{x}_{i}} g(\mathbf{x}, \widehat{\mathbf{w}}_{i}) = \mathbf{0}, \\ & \boldsymbol{\lambda}_{i}^{\top} g(\mathbf{x}_{i}, \widehat{\mathbf{w}}_{i}) = \mathbf{0}, \quad \forall i = 1, \dots, k \end{array}$$

Q1. Where does $c(\mathbf{x}^*(\mathbf{w}_i))$ come from?

KKT-based reformulation (cont'd)

Assume linear regression, i.e., $\mathbb{W}_{\theta}(\varphi_i) = \theta^{\top} \varphi_i$

KKT-based reformulation of the bilevel problem

$$\begin{array}{ll} \underset{\boldsymbol{\theta}, \mathbf{x}_{i}, \boldsymbol{\lambda}_{i}, \widehat{\mathbf{w}}_{i}}{\text{minimize}} & \frac{1}{2k} \sum_{i=1}^{k} \| c(\mathbf{x}_{i}) - c(\mathbf{x}^{\star}(\mathbf{w}_{i})) \|_{2}^{2} \\ \text{subject to} & \widehat{\mathbf{w}}_{i} = \boldsymbol{\theta}^{\top} \boldsymbol{\varphi}_{i}, \\ & \\ & g(\mathbf{x}_{i}, \widehat{\mathbf{w}}_{i}) \leq \mathbf{0}, \quad \boldsymbol{\lambda}_{i} \geq \mathbf{0}, \\ & \nabla_{\mathbf{x}} c(\mathbf{x}_{i}) + \boldsymbol{\lambda}^{\top} \nabla_{\mathbf{x}_{i}} g(\mathbf{x}, \widehat{\mathbf{w}}_{i}) = \mathbf{0}, \\ & \boldsymbol{\lambda}_{i}^{\top} g(\mathbf{x}_{i}, \widehat{\mathbf{w}}_{i}) = \mathbf{0}, \quad \forall i = 1, \dots, k \end{array}$$

- **Q1.** Where does $c(\mathbf{x}^*(\mathbf{w}_i))$ come from?
- Non-linear optimization program with hard complementarity condition.
- **Q2.** How can we make it more tractable? (demonstrate in-class)
- **Q3.** How to introduce regularization?

Gradient descent-based solution

Recall the optimization objective

$$\begin{array}{ll} \underset{\boldsymbol{\theta}}{\mathsf{minimize}} \quad \mathcal{R} \triangleq \frac{1}{2k} \sum_{i=1}^{k} \| \boldsymbol{c}(\mathbf{x}^{\star}(\mathbb{W}_{\boldsymbol{\theta}}(\boldsymbol{\varphi}_{i})) - \boldsymbol{c}_{i}^{\star} \|_{2}^{2} \end{array}$$

Gradient descent-based solution

Recall the optimization objective

$$\begin{array}{ll} \underset{\boldsymbol{\theta}}{\text{minimize}} \quad \mathcal{R} \triangleq \frac{1}{2k} \sum_{i=1}^{k} \| \boldsymbol{c}(\mathbf{x}^{\star}(\mathbb{W}_{\boldsymbol{\theta}}(\boldsymbol{\varphi}_{i})) - \boldsymbol{c}_{i}^{\star} \|_{2}^{2} \end{array}$$

Derivatives of the regret w.r.t. model parameter θ (linear regression)

$$\begin{split} \frac{\partial \mathcal{R}}{\partial \theta} &= \frac{\partial \mathcal{R}}{\partial c} \cdot \frac{\partial c}{\partial \mathbf{x}^{\star}} \cdot \frac{\partial \mathbf{x}^{\star}}{\partial \mathbb{W}_{\theta}(\varphi_{i})} \cdot \frac{\partial \mathbb{W}_{\theta}(\varphi_{i})}{\partial \theta} \\ &= \frac{1}{k} \sum_{i=1}^{k} \left(c(\mathbf{x}^{\star}(\mathbb{W}_{\theta}(\varphi_{i}))) - c_{i}^{\star} \right) \cdot \frac{\partial c}{\partial \mathbf{x}^{\star}} \cdot \frac{\partial \mathbf{x}^{\star}}{\partial \mathbb{W}_{\theta}(\varphi_{i})} \cdot \frac{\partial \mathbb{W}_{\theta}(\varphi_{i})}{\partial \theta} \\ &= \frac{1}{k} \sum_{i=1}^{k} \left(c(\mathbf{x}^{\star}(\mathbb{W}_{\theta}(\varphi_{i}))) - c_{i}^{\star} \right) \cdot \nabla_{\mathbf{x}} c(\mathbf{x}) \cdot \frac{\partial \mathbf{x}^{\star}}{\partial \mathbb{W}_{\theta}(\varphi_{i})} \cdot \frac{\partial \mathbb{W}_{\theta}(\varphi_{i})}{\partial \theta} \\ &= \frac{1}{k} \sum_{i=1}^{k} \left(c(\mathbf{x}^{\star}(\mathbb{W}_{\theta}(\varphi_{i}))) - c_{i}^{\star} \right) \cdot \nabla_{\mathbf{x}} c(\mathbf{x}_{i}) \cdot \frac{\partial \mathbf{x}^{\star}}{\partial \mathbb{W}_{\theta}(\varphi_{i})} \cdot \varphi_{i} \end{split}$$

Gradient descent-based solution

Recall the optimization objective

$$\begin{array}{ll} \underset{\boldsymbol{\theta}}{\text{minimize}} \quad \mathcal{R} \triangleq \frac{1}{2k} \sum_{i=1}^{k} \| \boldsymbol{c}(\mathbf{x}^{\star}(\mathbb{W}_{\boldsymbol{\theta}}(\boldsymbol{\varphi}_{i})) - \boldsymbol{c}_{i}^{\star} \|_{2}^{2} \end{array}$$

Derivatives of the regret w.r.t. model parameter θ (linear regression)

$$\begin{split} \frac{\partial \mathcal{R}}{\partial \theta} &= \frac{\partial \mathcal{R}}{\partial c} \cdot \frac{\partial c}{\partial \mathbf{x}^{\star}} \cdot \frac{\partial \mathbf{x}^{\star}}{\partial \mathbb{W}_{\theta}(\varphi_{i})} \cdot \frac{\partial \mathbb{W}_{\theta}(\varphi_{i})}{\partial \theta} \\ &= \frac{1}{k} \sum_{i=1}^{k} \left(c(\mathbf{x}^{\star}(\mathbb{W}_{\theta}(\varphi_{i}))) - c_{i}^{\star} \right) \cdot \frac{\partial c}{\partial \mathbf{x}^{\star}} \cdot \frac{\partial \mathbf{x}^{\star}}{\partial \mathbb{W}_{\theta}(\varphi_{i})} \cdot \frac{\partial \mathbb{W}_{\theta}(\varphi_{i})}{\partial \theta} \\ &= \frac{1}{k} \sum_{i=1}^{k} \left(c(\mathbf{x}^{\star}(\mathbb{W}_{\theta}(\varphi_{i}))) - c_{i}^{\star} \right) \cdot \nabla_{\mathbf{x}} c(\mathbf{x}) \cdot \frac{\partial \mathbf{x}^{\star}}{\partial \mathbb{W}_{\theta}(\varphi_{i})} \cdot \frac{\partial \mathbb{W}_{\theta}(\varphi_{i})}{\partial \theta} \\ &= \frac{1}{k} \sum_{i=1}^{k} \left(c(\mathbf{x}^{\star}(\mathbb{W}_{\theta}(\varphi_{i}))) - c_{i}^{\star} \right) \cdot \nabla_{\mathbf{x}} c(\mathbf{x}_{i}) \cdot \frac{\partial \mathbf{x}^{\star}}{\partial \mathbb{W}_{\theta}(\varphi_{i})} \cdot \varphi_{i} \end{split}$$

Regret is then minimized over iterations \$\theta \leftarrow \theta - \alpha \frac{\partial \mathcal{R}}{\partial \mathcal{V}_{\mathcal{L}}}\$ How to compute the last unknown gradient \$\frac{\partial \partial \partia

Computing $\frac{\partial \mathbf{x}^{\star}}{\partial \mathbb{W}_{\theta}(\varphi_i)}$ for convex quadratic programs

$$\begin{array}{ll} \underset{\mathbf{x}}{\text{minimize}} & \frac{1}{2} \mathbf{x}^{\top} \mathbf{C} \mathbf{x} + \mathbf{c}^{\top} \mathbf{x} \\ \\ \text{subject to} & \mathbf{A} \mathbf{x} = \mathbf{b}, \quad \mathbf{G} \mathbf{x} \leqslant \mathbf{h} \end{array}$$

The Lagrangian function $\mathcal{L}(\mathbf{x}, \nu, \lambda) = \frac{1}{2} \mathbf{x}^\top \mathbf{C} \mathbf{x} + \mathbf{c}^\top \mathbf{x} + \nu^\top (\mathbf{A} \mathbf{x} - \mathbf{b}) + \lambda^\top (\mathbf{G} \mathbf{x} - \mathbf{h})$

The stationarity, primal feasibility, and complementarity slackness:

$$egin{aligned} \mathsf{C} \mathsf{x}^{\star} + \mathsf{c} + \mathsf{A}^{ op} \, oldsymbol{
u}^{\star} + \mathsf{G}^{ op} \, \lambda^{\star} = 0 \ & \mathsf{A} \mathsf{x}^{\star} - \mathsf{b} = 0 \ & \mathsf{diag}[\lambda^{\star}](\mathsf{G} \mathsf{x}^{\star} - \mathsf{h}) = 0 \end{aligned}$$

Implicit differentiation (using a product rule, i.e., $\partial Cx^* = \partial C \cdot x^* + C \cdot \partial x^*)^1$

$$\begin{bmatrix} \mathbf{C} & \mathbf{G}^{\top} & \mathbf{A}^{\top} \\ \operatorname{diag}[\lambda^{*}]\mathbf{G} & \operatorname{diag}[\mathbf{Gx}^{*} - \mathbf{h}] & \mathbf{0} \\ \mathbf{A} & \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \partial \mathbf{x} \\ \partial \lambda \\ \partial \nu \end{bmatrix} = \begin{bmatrix} -\partial \mathbf{Cx}^{*} - \partial \mathbf{c} - \partial \mathbf{G}^{\top} \lambda^{*} - \partial \mathbf{A}^{\top} \nu^{*} \\ -\operatorname{diag}[\lambda^{*}]\partial \mathbf{Gx}^{*} + \operatorname{diag}[\lambda^{*}]\partial \mathbf{h} \\ -\partial \mathbf{Ax}^{*} + \partial \mathbf{b} \end{bmatrix}$$

¹Amos B., & Kolter J. Z. Optnet: Differentiable optimization as a layer in neural networks. 2017

Computing $\frac{\partial \mathbf{x}^{\star}}{\partial \mathbb{W}_{\theta}(\varphi_i)}$ for convex quadratic programs (cont'd)

$$\begin{bmatrix} \partial \mathbf{x} \\ \partial \lambda \\ \partial \nu \end{bmatrix} = \begin{bmatrix} \mathbf{C} & \mathbf{G}^\top & \mathbf{A}^\top \\ \mathsf{diag}[\lambda^*]\mathbf{G} & \mathsf{diag}[\mathbf{G}\mathbf{x}^* - \mathbf{h}] & \mathbf{0} \\ \mathbf{A} & \mathbf{0} & \mathbf{0} \end{bmatrix}^{-1} \begin{bmatrix} -\partial \mathbf{C}\mathbf{x}^* - \partial \mathbf{c} - \partial \mathbf{G}^\top \lambda^* - \partial \mathbf{A}^\top \nu^* \\ -\mathsf{diag}[\lambda^*]\partial \mathbf{G}\mathbf{x}^* + \mathsf{diag}[\lambda^*]\partial \mathbf{h} \\ -\partial \mathbf{A}\mathbf{x}^* + \partial \mathbf{b} \end{bmatrix}$$

Suppose that vector **b** is given by the prediction of our model, i.e., $\mathbf{b} = \mathbb{W}_{\theta}(\varphi_i)$ Focus only on $\partial \mathbf{b}$, and set other differentials to zero

$$\begin{bmatrix} \partial \mathbf{x} \\ \partial \boldsymbol{\lambda} \\ \partial \boldsymbol{\nu} \end{bmatrix} = \begin{bmatrix} \mathbf{C} & \mathbf{G}^{\top} & \mathbf{A}^{\top} \\ \operatorname{diag}[\boldsymbol{\lambda}^{\star}]\mathbf{G} & \operatorname{diag}[\mathbf{G}\mathbf{x}^{\star} - \mathbf{h}] & \mathbf{0} \\ \mathbf{A} & \mathbf{0} & \mathbf{0} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \\ \partial \mathbf{b} \end{bmatrix} = \begin{bmatrix} (\mathbf{A}^{\top})^{-1}\partial \mathbf{b} \\ \mathbf{0} \\ \mathbf{0} \end{bmatrix}$$
Hence, the $\frac{\partial \mathbf{x}}{\partial \mathbf{b}} = \frac{\partial \mathbf{x}}{\partial W_{\boldsymbol{\theta}}(\varphi_i)} = (\mathbf{A}^{\top})^{-1}$

By analogy, compute the derivative of optimizer x to any other QP problem data

Gradient descent for decision-focused learning

$$\begin{array}{ll} \underset{\boldsymbol{\theta}}{\text{minimize}} \quad \mathcal{R} \triangleq \frac{1}{2k} \sum_{i=1}^{k} \| \boldsymbol{c}(\mathbf{x}^{\star}(\mathbb{W}_{\boldsymbol{\theta}}(\boldsymbol{\varphi}_{i})) - \boldsymbol{c}_{i}^{\star} \|_{2}^{2} \end{array}$$

The gradient descent-like algorithms use the update

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \alpha \frac{\partial \mathcal{R}}{\partial \boldsymbol{\theta}}$$

for some step size $\alpha > 0$

- Requires solving an optimization problem at every iteration (Q. Why?)
- **Vet**, more versatile solution (no structural assumption on model \mathbb{W}_{θ})
- Also, scales up to much larger problems (in terms of dimensions n, m and k)

Illustrative example: Setting

Dataset
$$\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_i, y_i), \dots, (\mathbf{x}_n, y_n)\}$$
 of *n* records

Prediction-focused learning: find the linear model $\hat{\mathbf{y}} = \mathbf{w}^{\top} \hat{\mathbf{x}}$ by solving

$$\underset{\mathbf{w}}{\text{minimize}} \quad \frac{1}{2n} \sum_{i=1}^{n} (\mathbf{w}^{\top} \mathbf{x}_{i} - y_{i})^{2}$$

The downstream \hat{y} -parametrized decision problem:

minimize
$$\frac{1}{2}(z-\hat{y})^2 + \lambda z^2$$

subject to $z \ge 0$

Decision-focused learning: seeks the linear model by solving

$$\begin{array}{ll} \underset{\mathbf{w}, \hat{\mathbf{y}}}{\text{minimize}} & \frac{1}{2n} \sum_{i=1}^{n} (z^{\star}(\hat{y}_{i}) - z^{\star}(y_{i}))^{2} \\ \text{subject to} & \hat{y}_{i} = \mathbf{w}^{\top} \mathbf{x}_{i} \\ & z^{\star}(\hat{y}_{i}) = \underset{z}{\operatorname{argmin}} & \frac{1}{2} (z - \hat{y}_{i})^{2} + \lambda z^{2} \\ & \text{subject to} & z \ge 0, \end{array} \qquad \forall i = 1, \dots, n \end{array}$$

Illustrative example: Solution via bilevel optimization

■ The KKTs of the lower-level problem are

$z - \widehat{y} + \lambda - \mu = 0$	(Lagrange optimality)
$z \ge 0$	(primal feasibility)
$\mu \geqslant 0$	(dual feasibility)
$z\mu = 0$	(complementarity slackness)

Illustrative example: Solution via bilevel optimization

The KKTs of the lower-level problem are

$z - \widehat{y} + \lambda - \mu = 0$	(Lagrange optimality)
$z \ge 0$	(primal feasibility)
$\mu \geqslant 0$	(dual feasibility)
$z\mu=0$	(complementarity slackness)

Big-M reformulation of complementarity slackness:

$$\begin{cases} \mu \leqslant Mu \\ z \leqslant M(1-u) \end{cases}$$

where $u \in \{0,1\}$ is an auxiliary binary variable, and M is a large enough constant

Illustrative example: Solution via bilevel optimization

The KKTs of the lower-level problem are

 $\begin{array}{ll} z-\widehat{y}+\lambda-\mu=0 & (\mbox{Lagrange optimality})\\ z\geqslant 0 & (\mbox{primal feasibility})\\ \mu\geqslant 0 & (\mbox{dual feasibility})\\ z\mu=0 & (\mbox{complementarity slackness}) \end{array}$

Big-M reformulation of complementarity slackness:

$$\begin{cases} \mu \leqslant Mu \\ z \leqslant M(1-u) \end{cases}$$

where $u \in \{0, 1\}$ is an auxiliary binary variable, and M is a large enough constant I The bilevel reformulation of the decision-focused learning problem:

$$\begin{array}{ll} \underset{\mathbf{w},\widehat{\mathbf{y}},\mathbf{z},\boldsymbol{\mu},\mathbf{u}}{\text{minimize}} & \frac{1}{2n} \sum_{i=1}^{n} (z^{\star}(\widehat{y}_{i}) - z^{\star}(y_{i}))^{2} \\ \text{subject to} & \widehat{y}_{i} = \mathbf{w}^{\top} \widehat{\mathbf{x}}_{i}, \quad \forall i = 1, \dots, n \\ & \mathbf{z} - \widehat{\mathbf{y}} + 2\lambda \mathbf{z} - \boldsymbol{\mu} = \mathbf{0} \\ & \mathbf{z} \ge \mathbf{0}, \quad \boldsymbol{\mu} \ge \mathbf{0}, \quad \mathbf{u} \in \{0, 1\}^{n} \\ & \boldsymbol{\mu} \le M \mathbf{u}, \quad \mathbf{z} \le M(\mathbf{1} - \mathbf{u}) \end{array}$$

Illustrative example: Solution via stochastic gradient descent

$$\begin{array}{ll} \underset{\mathbf{w},\widehat{\mathbf{y}}}{\text{minimize}} & \frac{1}{2n} \sum_{i=1}^{n} (z^{\star}(\widehat{y}_{i}) - z^{\star}(y_{i}))^{2} \\ \text{subject to} & \widehat{y}_{i} = \mathbf{w}^{\top} \mathbf{x}_{i} \\ & z^{\star}(\widehat{y}_{i}) = \underset{z}{\operatorname{argmin}} & \frac{1}{2} (z - \widehat{y}_{i})^{2} + \lambda z^{2} \\ & \text{subject to} & z \ge 0, \end{array}$$

Closed-form solution for the lower-level problem

$$z^{\star}(\widehat{y_i}) = \max\left\{0, rac{\widehat{y_i}}{1+2\lambda}
ight\}$$

The gradient $\partial z^*(\hat{y}_i)/\partial \hat{y}_i$ is obtained using LogSumExp approximation:

$$z^{\star}(\widehat{y_{i}}) = \max\left\{0, \frac{\widehat{y_{i}}}{1+2\lambda}\right\} \approx \log\left(1 + \exp\left(\frac{\widehat{y_{i}}}{1+2\lambda}\right)\right) \qquad (\text{EECS 559})$$
$$z^{\star}(\widehat{y_{i}})/\partial\widehat{y_{i}} \approx \frac{\exp\left(\frac{\widehat{y_{i}}}{1+2\lambda}\right)}{1 + \exp\left(\frac{\widehat{y_{i}}}{1+2\lambda}\right)} \frac{1}{1+2\lambda}$$

Illustrative example: Solution via stochastic gradient descent

$$\begin{array}{ll} \underset{\mathbf{w}, \widehat{\mathbf{y}}}{\text{minimize}} & \frac{1}{2n} \sum_{i=1}^{n} (z^{\star}(\widehat{y}_{i}) - z^{\star}(y_{i}))^{2} \\ \text{subject to} & \widehat{y}_{i} = \mathbf{w}^{\top} \mathbf{x}_{i} \\ & z^{\star}(\widehat{y}_{i}) = \underset{z}{\operatorname{argmin}} & \frac{1}{2} (z - \widehat{y}_{i})^{2} + \lambda z^{2} \\ & \text{subject to} & z \ge 0, \end{array}$$

\blacksquare For SGD, we draw one sample *i* from a dataset and minimize

$$\mathcal{C}_i = (z^*(\widehat{y}_i) - z^*(y_i))^2$$

By the chain rule, the gradient is

$$\frac{\partial \mathcal{C}_{i}}{\partial \mathbf{w}} = (z^{\star}(\hat{y}_{i}) - z^{\star}(y_{i})) \cdot \frac{\partial z^{\star}}{\partial y_{i}} \cdot \frac{\partial y_{i}}{\partial \mathbf{w}} \\ = \left(\underbrace{\max\left\{0, \frac{\hat{y}_{i}}{1+2\lambda}\right\}}_{\text{closed form for } z^{\star}} - \underbrace{z^{\star}(y_{i})}_{\text{precomputed offline}}\right) \cdot \underbrace{\frac{\exp\left(\frac{\hat{y}_{i}}{1+2\lambda}\right)}{1+\exp\left(\frac{\hat{y}_{i}}{1+2\lambda}\right)} \frac{1}{1+2\lambda}}{\partial z^{\star}(\hat{y}_{i})/\partial \hat{y}_{i}} \cdot \mathbf{x}_{i}$$

The SGD perform the following update:

$$\mathbf{w} \leftarrow \mathbf{w} - \alpha \frac{\partial \mathcal{C}_i}{\partial \mathbf{w}}$$

with some step size $\alpha > 0$.

Illustrative example: Results

- Experiments on synthetic data (100 samples)
- Prediction space: displays how the model explains the underlying data
- Decision space: displays how the model prescribes decisions



- Prediction-focused model is better in predictions, decision-focused in decisions
- Bilevel optimization provides the **global** optimal solution
- Gradient descent-based solution does not have global guarantees

Tutorial time

- Look around you and form teams of 2 people (1 min)
- Pick one person to code; the other one guides
- Work in pairs for the whole tutorial session

Goal: replicate the results of the illustrative example

Multiparametric programming to ease computational burden

Challenge: Large-scale downstream optimization

■ Wind power forecast \hat{w} -parametrized DC-OPF problem

$$\begin{array}{ll} \mbox{minimize} & \frac{1}{2} \mathbf{p}^\top \mathbf{C} \mathbf{p} + \mathbf{c}^\top \mathbf{p} & generator \ dispatch \ cost \\ \mbox{subject to} & \underline{\mathbf{p}} \leqslant \mathbf{p} \leqslant \overline{\mathbf{p}} & generation \ limits \\ & \mathbf{1}^\top (\mathbf{p} + \widehat{\mathbf{w}} - \mathbf{d}) = 0 & power \ balance \ condition \\ & |\mathbf{F} (\mathbf{p} + \widehat{\mathbf{w}} - \mathbf{d})| \leqslant \overline{\mathbf{f}} & power \ flow \ limits \\ \end{array}$$

For practical-sized systems, there are thousands of variables and constraintsSolving DC-OPF at every SGD update significantly slows convergence

Multiparametric programming to ease computational burden

Challenge: Large-scale downstream optimization

■ Wind power forecast ŵ-parametrized DC-OPF problem

$$\begin{array}{ll} \mbox{minimize} & \frac{1}{2} \mathbf{p}^\top \mathbf{C} \mathbf{p} + \mathbf{c}^\top \mathbf{p} & generator \ dispatch \ cost \\ \mbox{subject to} & \underline{\mathbf{p}} \leqslant \mathbf{p} \leqslant \overline{\mathbf{p}} & generation \ limits \\ & \mathbf{1}^\top (\mathbf{p} + \widehat{\mathbf{w}} - \mathbf{d}) = 0 & power \ balance \ condition \\ & |\mathbf{F} (\mathbf{p} + \widehat{\mathbf{w}} - \mathbf{d})| \leqslant \overline{\mathbf{f}} & power \ flow \ limits \\ \end{array}$$

For practical-sized systems, there are thousands of variables and constraints

Solving DC-OPF at every SGD update significantly slows convergence

Solution: Multiparametric programming (MPP)

- Studies the solutions of an optimization problem as a function of its parameters
- Let $\widehat{\mathbf{w}}$ belong to the union $\mathcal{W} \triangleq \bigcup_{i=1}^{k} [0, 1]_i$ of normalized output ranges of k wind farms. A sampled wind power output $\widehat{\mathbf{w}}$ is seen as a random draw from \mathcal{W}
- For convex quadratic programs, set *W* can be portioned into polyhedral regions, each having a closed-form solution for primal and dual variables ⇒ speed up!

Compact DC-OPF formulation for MPP

$$\begin{array}{ll} \underset{p}{\text{minimize}} & \frac{1}{2} \mathbf{p}^{\top} \mathbf{C} \mathbf{p} + \mathbf{c}^{\top} \mathbf{p} \\\\ \text{subject to} & \mathbf{A} \mathbf{p} \leqslant \mathbf{B} \widehat{\mathbf{w}} + \mathbf{b} & : \boldsymbol{\lambda} \end{array}$$

where problem data is

$$\mathbf{A} = \begin{bmatrix} \mathbf{I} \\ -\mathbf{I} \\ -\mathbf{I}^{\top} \\ \mathbf{F} \\ -\mathbf{F} \end{bmatrix}, \quad \mathbf{B} = \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \\ \mathbf{1}^{\top} \\ \mathbf{F} \\ -\mathbf{F} \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} \mathbf{\overline{p}} \\ -\mathbf{\underline{p}} \\ -\mathbf{1}^{\top} \mathbf{d} \\ \mathbf{\overline{f}} \\ \mathbf{\overline{f}} \end{bmatrix}$$

Q. Why can we replace the power balance constraint with inequality?

Compact DC-OPF formulation for MPP

 $\begin{array}{ll} \mbox{minimize} & \frac{1}{2} \mathbf{p}^\top \mathbf{C} \mathbf{p} + \mathbf{c}^\top \mathbf{p} \\ \mbox{subject to} & \mathbf{A} \mathbf{p} \leqslant \mathbf{B} \widehat{\mathbf{w}} + \mathbf{b} & : \lambda \end{array}$

where problem data is

$$\mathbf{A} = \begin{bmatrix} \mathbf{I} \\ -\mathbf{I} \\ -\mathbf{1}^{\top} \\ \mathbf{F} \\ -\mathbf{F} \end{bmatrix}, \quad \mathbf{B} = \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \\ \mathbf{1}^{\top} \\ \mathbf{F} \\ -\mathbf{F} \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} \mathbf{\overline{p}} \\ -\mathbf{\underline{p}} \\ -\mathbf{1}^{\top} \mathbf{\overline{d}} \\ \mathbf{\overline{f}} \\ \mathbf{\overline{f}} \end{bmatrix}$$

Q. Why can we replace the power balance constraint with inequality?

- A multiparametric programming problem (MPP) dependent on ŵ
- For some realization \widehat{w} , MPP expresses the primal **p** and dual λ solutions as affine functions of some \widehat{w} in the *neighborhood* of \widehat{w}

Closed-form DC-OPF solutions from MPP

Lagrangian function of the compact DC-OPF:

$$\mathcal{L}(\mathbf{p}, \boldsymbol{\lambda}) = \frac{1}{2} \mathbf{p}^{\top} \mathbf{C} \mathbf{p} + \mathbf{c}^{\top} \mathbf{p} + \boldsymbol{\lambda}^{\top} \left(\mathbf{A} \mathbf{p} - \mathbf{B} \widehat{\mathbf{w}} - \mathbf{b} \right)$$

Let $\tilde{A}, \tilde{B}, \tilde{b}$ be sub-matrices/vectors containing only those entries of A, B, b that correspond to the *active* (binding) constraints ($\lambda > 0$)

A, $\overline{\mathbf{B}}$, $\overline{\mathbf{b}}$ collect the remaining entries corresponding to *inactive* constraints ($\lambda = 0$).

Closed-form DC-OPF solutions from MPP

Lagrangian function of the compact DC-OPF:

$$\mathcal{L}(\mathbf{p}, \boldsymbol{\lambda}) = \frac{1}{2} \mathbf{p}^{\top} \mathbf{C} \mathbf{p} + \mathbf{c}^{\top} \mathbf{p} + \boldsymbol{\lambda}^{\top} \left(\mathbf{A} \mathbf{p} - \mathbf{B} \widehat{\mathbf{w}} - \mathbf{b} \right)$$

Let $\tilde{A}, \tilde{B}, \tilde{b}$ be sub-matrices/vectors containing only those entries of A, B, b that correspond to the *active* (binding) constraints ($\lambda > 0$)

A, $\overline{\mathbf{B}}$, $\overline{\mathbf{b}}$ collect the remaining entries corresponding to *inactive* constraints ($\lambda = 0$).

From the Karush–Kuhn–Tucker (KKT) conditions we have a system of equations:

Q. From which KKT conditions do these equation come from?

Manipulating these equations, we get the closed-form solutions

$$\mathbf{C}\mathbf{p} + \mathbf{c} + \mathbf{\tilde{A}}^{\top} \mathbf{\tilde{\lambda}} = \mathbf{0}$$
 (1a)

$$\tilde{\mathbf{A}}\mathbf{p} - \tilde{\mathbf{B}}\widehat{\mathbf{w}} - \widetilde{\mathbf{b}} = \mathbf{0} \tag{1b}$$

From (1a) we express **p** (assuming **C** is invertible)

$$\mathbf{p} = -\mathbf{C}^{-1} \left(\mathbf{c} + \tilde{\mathbf{A}}^{\top} \tilde{\boldsymbol{\lambda}} \right)$$
(2)

and substitute it into (1b):

$$\begin{split} &-\tilde{A}C^{-1}c - \tilde{A}C^{-1}\tilde{A}^{\top}\tilde{\lambda} - \tilde{B}\hat{w} - \tilde{b} = 0 \\ \Leftrightarrow \quad \tilde{\lambda} = \underbrace{-\left(\tilde{A}C^{-1}\tilde{A}^{\top}\right)^{-1}\tilde{B}}_{D}\hat{w} \underbrace{-\left(\tilde{A}C^{-1}\tilde{A}^{\top}\right)^{-1}\left(\tilde{b} + \tilde{A}C^{-1}c\right)}_{d} \\ \Leftrightarrow \quad \tilde{\lambda} = D\hat{w} + d \quad (\text{closed-form solution for duals}) \end{split}$$
(3)

Substituting (3) into (2), we have the solution for p :

$$\mathbf{p} = -\mathbf{C}^{-1} \left(\mathbf{c} + \tilde{\mathbf{A}}^{\top} \left(\mathbf{D} \widehat{\mathbf{w}} + \mathbf{d} \right) \right) = \underbrace{-\mathbf{C}^{-1} \tilde{\mathbf{A}}^{\top} \mathbf{D}}_{\Pi} \widehat{\mathbf{w}} \underbrace{-\mathbf{C}^{-1} \mathbf{c} - \mathbf{C}^{-1} \tilde{\mathbf{A}}^{\top} \mathbf{d}}_{\pi}$$

$$\Leftrightarrow \quad \mathbf{p} = \Pi \widehat{\mathbf{w}} + \pi \quad \text{(closed-form solution for primals)} \tag{4}$$

The primal-dual solution depends linearly on $\widehat{\mathbf{w}}$:

$$\begin{bmatrix} \mathbf{p} \\ \tilde{\boldsymbol{\lambda}} \end{bmatrix} = \begin{bmatrix} \mathbf{\Pi} \\ \mathbf{D} \end{bmatrix} \widehat{\mathbf{w}} + \begin{bmatrix} \boldsymbol{\pi} \\ \mathbf{d} \end{bmatrix}$$
 (5)

Q. How do we obtain matrices and vectors Π , **D**, π and **d**?

The primal-dual solution depends linearly on $\widehat{\mathbf{w}}$:

$$\begin{bmatrix} \mathbf{p} \\ \tilde{\boldsymbol{\lambda}} \end{bmatrix} = \begin{bmatrix} \mathbf{\Pi} \\ \mathbf{D} \end{bmatrix} \widehat{\mathbf{w}} + \begin{bmatrix} \boldsymbol{\pi} \\ \mathbf{d} \end{bmatrix}$$
 (5)

Q. How do we obtain matrices and vectors Π , **D**, π and **d**?

For a given \hat{w} , we obtain them by solving an optimization problem (**bad** news)

Eq. (5) valid not only for $\hat{\mathbf{w}}$ but for all $\hat{\mathbf{w}}$ in the neighborhood of $\hat{\mathbf{w}}$ (good news)

The primal-dual solution depends linearly on $\widehat{\mathbf{w}}$:

$$\begin{bmatrix} \mathbf{p} \\ \tilde{\boldsymbol{\lambda}} \end{bmatrix} = \begin{bmatrix} \mathbf{\Pi} \\ \mathbf{D} \end{bmatrix} \widehat{\mathbf{w}} + \begin{bmatrix} \boldsymbol{\pi} \\ \mathbf{d} \end{bmatrix}$$
 (5)

Q. How do we obtain matrices and vectors Π , **D**, π and **d**?

For a given \hat{w} , we obtain them by solving an optimization problem (bad news)

Eq. (5) valid not only for $\hat{\mathbf{w}}$ but for all $\hat{\mathbf{w}}$ in the neighborhood of $\hat{\mathbf{w}}$ (good news)

The neighborhood originates from two KKT conditions:

$$(\overline{A}\Pi - \overline{B})\widehat{w} < \overline{b} - \overline{A}\pi$$
 (primal feasibility of inactive constraints) (6a)
 $D\widehat{w} + d \ge 0$ (dual feasibility) (6b)

We term the neighborhood of \hat{w} a *critical region (CR)*, formally defined as:

$$\mathcal{R} = \left\{ \widehat{\mathbf{w}} \in \mathcal{W} \mid (\overline{\mathbf{A}} \Pi - \overline{\mathbf{B}}) \widehat{\mathbf{w}} < \overline{\mathbf{b}} - \overline{\mathbf{A}} \pi, \ \mathbf{D} \widehat{\mathbf{w}} + \mathbf{d} \ge \mathbf{0} \right\}$$
(7)

In fact, there are many CRs. The task is to identify the CR corresponding to a particular realization of wind power, and then compute solution from (5).

Identifying critical regions

- 5-Bus PJM testbed with 3 variable net loads
- Initialize the set of critical regions $C\mathcal{R} = \{\emptyset\}$

■ To identify critical regions, perform the following iterations

```
for i = 1, ..., 1000 do
Sample \hat{w} from W
if \hat{w} does not belong to any \mathcal{R} in \mathcal{CR} then
Solve DC-OPF on \hat{w}, compute \mathcal{R} and add to \mathcal{CR}
end if
end for
```

Q1. What is the complexity of this algorithm?

Identifying critical regions

- 5-Bus PJM testbed with 3 variable net loads
- Initialize the set of critical regions $C\mathcal{R} = \{\emptyset\}$

To identify critical regions, perform the following iterations

```
for i = 1, ..., 1000 do
Sample \hat{w} from W
if \hat{w} does not belong to any \mathcal{R} in \mathcal{CR} then
Solve DC-OPF on \hat{w}, compute \mathcal{R} and add to \mathcal{CR}
end if
end for
```

Q1. What is the complexity of this algorithm?

■ 7 distinct critical regions, each having a fixed set of active constraints



Q2. How many DC-OPF problems did we solve in this case?

MPP-Powered SGD for decision-focused learning

Machine learning model $\widehat{\mathbf{w}} = \mathbb{W}_{\boldsymbol{\theta}}(\widehat{\boldsymbol{\varphi}})$

For a given $\widehat{\mathbf{w}}$, the primal solution and the gradient are

$$\mathbf{p} = \mathbf{\Pi}\widehat{\mathbf{w}} + m{\pi} \qquad \quad rac{\partial \mathbf{p}^{\star}}{\partial \widehat{\mathbf{w}}} = \mathbf{\Pi}^{ op}$$

Decision-focused problem minimizing the primal decision error

$$\underset{\boldsymbol{\theta}}{\text{minimize}} \quad \mathcal{R} \triangleq \frac{1}{2k} \sum_{i=1}^{k} \|\mathbf{p}^{\star}(\mathbb{W}_{\boldsymbol{\theta}}(\widehat{\varphi}_{i})) - \mathbf{p}_{i}^{\star}\|_{2}^{2}$$

■ The gradient update in SGD for a particular sample *i* takes the form:

$$egin{aligned} &rac{\partial \mathcal{R}}{\partial oldsymbol{ heta}} = \left(\mathbf{p}^{\star}(\mathbb{W}_{oldsymbol{ heta}}(\widehat{arphi}_{i})) - \mathbf{p}^{\star}_{i}
ight) \cdot rac{\partial \mathbf{p}^{\star}}{\partial \mathbb{W}_{oldsymbol{ heta}}(arphi_{i})} \cdot rac{\partial \mathbb{W}_{oldsymbol{ heta}}(arphi_{i})}{\partial oldsymbol{ heta}} \ &= \left(\mathbf{\Pi} \mathbb{W}_{oldsymbol{ heta}}(\widehat{arphi}_{i}) + \pi - \mathbf{p}^{\star}_{i}
ight) \cdot \mathbf{\Pi}^{ op} \cdot rac{\partial \mathbb{W}_{oldsymbol{ heta}}(arphi_{i})}{\partial oldsymbol{ heta}} \end{aligned}$$

- **I** No need to solve optimization if we know the critical region of $\mathbb{W}_{\theta}(\widehat{\varphi}_i)$!
- Significant time saving after we explore a sufficient number of critical regions
- By analogy, we can minimize the dual decision error

Example (cont'd): Wind power forecast & electricity pricing



Example (cont'd): Price-awareness for wind power forecast

- Dataset $\{(\varphi_1, \mathbf{w}_1), \dots, (\varphi_k, \mathbf{w}_k)\}$ of wind power records
- Two deep learning architectures **DeepWP** and **DeepWP+** for forecasting²

loss function:

 $\|\widehat{\mathbf{w}} - \mathbf{w}\|$



DeepWP+ informs wind power predictions about the downstream pricing error

²Dvorkin V., & Fioretto F. Price-aware deep learning for electricity markets. 2023

Example (cont'd): Price-awareness for wind power forecast

- Dataset $\{(\varphi_1, \mathbf{w}_1), \dots, (\varphi_k, \mathbf{w}_k)\}$ of wind power records
- Two deep learning architectures DeepWP and DeepWP+ for forecasting²



DeepWP+ informs wind power predictions about the downstream pricing errors

²Dvorkin V., & Fioretto F. Price-aware deep learning for electricity markets. 2023

Example (cont'd): Numerical experiments

Standard PowerModels.jl test cases

- 1,000 wind power records
 Active power output
 Wind speed and direction
 - vind speed and directic
 - Blade pitch angle
- DeepWP has 4 hidden layers with 30 neurons each. DeepWP+ additionally includes an opt. layer
- ADAM optimizer with varying learning rate



Example (cont'd): IEEE 118-bus system



DeepWP: Forecast error minimization yields $\delta \pi \in [-4, 1]$ \$/MWh DeepWP+: Price error minimization yields $\delta \pi \in [-1, 1]$ \$/MWh

- Q1. Which prediction is fairer?
- Q2. What prevents us from completely eliminating the price error?

Example (cont'd): Wind power forecasts



DeepWP: Minimizes the average forecast deviation

DeepWP+: Intentionally over-predicts in certain range of wind speeds

Example (cont'd): Bias of DeepWP+ model



DeepWP+ training starts at iteration 500 using a pre-trained DeepWP model
 RMSE(ŵ) and RMSE(â) are conflicting objectives which are kept in balance

Example (cont'd): Forecast and price error trade-off

	DeepWP				DeepWP+							
case	$RMSE(\hat{w})$	$RMSE(\hat{\pi})$	$CVaR(\hat{\pi})$	$\alpha-value$	RMSE(ŵ)		$RMSE(\hat{\pi})$		$CVaR(\hat{\pi})$		$\alpha-value$	
	MWh	\$/MWh	\$/MWh	\$/MWh	MWh	gain	\$/MWh	gain	\$/MWh	gain	\$/MWh	gain
14_ieee	0.35	0.62	1.52	0	0.35	+0.6%	0.61	-0.6%	1.50	-0.8%	0	_
57_ieee	2.31	11.03	34.64	32.08	2.60	+11.2%	10.72	-2.9%	33.59	-3.1%	30.92	-3.8%
24_ieee	4.08	8.62	37.70	27.48	4.51	+9.6%	8.33	-3.5%	36.35	-3.7%	26.26	-4.6%
39_epri	5.94	11.15	31.21	17.53	6.43	+7.6%	10.19	-9.4%	28.02	-11.4%	15.84	-10.7%
73_ieee	4.02	5.12	16.21	32.83	5.51	+26.9%	4.24	-20.8%	13.41	-20.9%	26.63	-23.3%
118_ieee	2.29	3.59	11.32	17.91	2.60	+12.1%	2.88	-24.7%	9.06	-25.0%	14.09	-27.2%

- Worst-case improvement exceeds that of the average case
- Price error reduction and fairness improves with the size of the network

Resources

- Mandi, J. et al. (2024). Decision-focused learning: Foundations, state of the art, benchmark and future opportunities. Journal of Artificial Intelligence Research, 80, 1623-1701.
- 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).
- Tøndel, P., Johansen, T. A., & Bemporad, A. (2003). An algorithm for multi-parametric quadratic programming and explicit MPC solutions. Automatica, 39(3), 489-497.
- Dvorkin, V., & Fioretto, F. (2023). Price-aware deep learning for electricity markets. arXiv preprint arXiv:2308.01436.