Optimization over Trained Neural Networks: Difference-of-Convex Algorithm and Application to Data Center Scheduling

Xinwei Liu. Vladmir Dvorkin

Department of Electrical Engineering and Computer Science

University of Michigan - Ann Arbor

INFORMS 2025. Atlanta. GA

October 24, 2025

Overview



- 1. Motivation and Background
- 2. Mathematical Formulation

3. Application

Optimization over trained neural networks



- Usually, we know all elements (c, g, \mathcal{X}) of an optimization problem.
- What should we do if the constraint function g is unknown?
- In practice, we can approximate the function g from the data.

Optimization over trained neural networks



$$\begin{array}{cccc} \underset{\mathbf{x} \in \mathcal{X}}{\text{minimize}} & c(\mathbf{x}) & & \underset{\mathbf{x} \in \mathcal{X}}{\text{minimize}} & c(\mathbf{x}) \\ \text{subject to} & g(\mathbf{x}) \leqslant \mathbf{0} & & \text{subject to} & \mathbb{NN}(\mathbf{x}) \leqslant \mathbf{0} \end{array}$$

- ullet Usually, we know all elements (c,g,\mathcal{X}) of an optimization problem.
- What should we do if the constraint function g is unknown?
- ullet In practice, we can approximate the function g from the data.
- Consider an approximation by a neural network $\mathbb{NN}: \mathcal{X} \mapsto \mathbb{R}$
- The problem becomes an optimization over trained neural network.



A neural network mapping d to λ is mathematically represented as

$$\pmb{\lambda} = \pmb{\mathsf{W}}_{\mathcal{N}+1} \mathsf{ReLU}(\dots \mathsf{ReLU}(\pmb{\mathsf{W}}_1 \mathbf{d} + \mathbf{b}_1) \dots) + \mathbf{b}_{\mathcal{N}+1}$$

with fixed weights and biases **W** and **b**, and variable input **d** and output λ



A neural network mapping d to λ is mathematically represented as

$$\lambda = \mathbf{W}_{N+1} \mathsf{ReLU}(\dots \mathsf{ReLU}(\mathbf{W}_1 \mathbf{d} + \mathbf{b}_1) \dots) + \mathbf{b}_{N+1}$$

with fixed weights and biases **W** and **b**, and variable input **d** and output λ

y = a + v $v = ReLU(a) \implies$

$$y=a+v$$
 y output variable $0\leqslant y\perp v\geqslant 0$ v auxiliary variable \perp complementarity



A neural network mapping d to λ is mathematically represented as

$$oldsymbol{\lambda} = oldsymbol{\mathsf{W}}_{\mathcal{N}+1} \mathsf{ReLU}(\dots \mathsf{ReLU}(oldsymbol{\mathsf{W}}_1 \mathbf{d} + \mathbf{b}_1) \dots) + \mathbf{b}_{\mathcal{N}+1}$$

with fixed weights and biases W and b. and variable input d and output λ

We can write the ReLU activation function using complementarity logic:

$$y = \mathsf{ReLU}(a) \implies egin{array}{cccc} y = a + v & y & \mathsf{output} & \mathsf{variable} \ 0 \leqslant y \perp v \geqslant 0 & v & \mathsf{auxiliary} & \mathsf{variable} \ & \perp & \mathsf{complementarity} \end{array}$$

Extremely difficult problem to solve in practice due to large computational burden

a input variable

State of the art approaches on complementarity constraint



State of the art approaches

- Mix-integer constraint[Tjeng et al., 2017, Grimstad and Andersson, 2019]
- Big-M reformulation of ReLU-numerical issue
 [Tjeng et al., 2017, Grimstad and Andersson, 2019]
- SOS1 reformulation with decision trees[Turner et al., 2024]
- Semidefinite relaxation of ReLU [Dathathri et al., 2020, Fazlyab et al., 2020]
- Bound propagation [Wang et al., 2021]

>>> X. Liu 5 / 3

Contribution



Our contribution

- Use difference of convex approach to optimize over trained neural network.
- Relax the ReLU constraints by penalizing them in the objective function.
- Avoid the computational complexity of standard MIP formulations.
- Formulate an algorithm to compute the hyperparameter to achieve fast convergence.
- Demonstrate in the data center load allocation example.

>>> X. Liu 6 / 3



```
\label{eq:minimize} \begin{split} \underset{\mathbf{d} \in \mathcal{D}, \mathbf{y}, \boldsymbol{\lambda}}{\text{minimize}} & c(\boldsymbol{\lambda}) \quad \text{convex function} \\ \text{subject to} & \mathbf{y}_1 = \text{ReLU}(\mathbf{W}_1\mathbf{d} + \mathbf{b}_1), \\ & \mathbf{y}_i = \text{ReLU}(\mathbf{W}_i\mathbf{y}_{i-1} + \mathbf{b}_i), \\ & \forall i = 2, \dots, N \\ & \boldsymbol{\lambda} = \mathbf{W}_{N+1}\mathbf{y}_N + \mathbf{b}_{N+1}, \end{split}
```



```
minimize c(\lambda) convex function subject to \mathbf{y}_1 = \text{ReLU}(\mathbf{W}_1\mathbf{d} + \mathbf{b}_1), \mathbf{y}_i = \text{ReLU}(\mathbf{W}_i\mathbf{y}_{i-1} + \mathbf{b}_i), \forall i = 2, \dots, N, \lambda = \mathbf{W}_{N+1}\mathbf{y}_N + \mathbf{b}_{N+1}, \mathbf{w}_i = \mathbf{w}_{N+1}\mathbf{y}_N + \mathbf{b}_{N+1},
```



- Problem: NP hard.
- Solution: Introduce the complementarity condition as a bilinear penalty term in the objective function.

$$\begin{aligned} & \underset{(\mathbf{d}, \mathbf{y}, \mathbf{v}, \boldsymbol{\lambda}) \in \mathcal{O}}{\text{minimize}} \ c(\boldsymbol{\lambda}) + \rho \sum_{i=1}^{N} \mathbf{y}_{i}^{\top} \mathbf{v}_{i} \\ & \text{subject to} \ \mathbf{y}_{1} = \mathbf{W}_{1} \mathbf{d} + \mathbf{b}_{1} + \mathbf{v}_{1}, \\ & \mathbf{y}_{i} = \mathbf{W}_{i} \mathbf{y}_{i-1} + \mathbf{b}_{i} + \mathbf{v}_{i}, \ \forall i = 2, \dots, N \\ & \boldsymbol{\lambda} = \mathbf{W}_{N+1} \mathbf{y}_{N} + \mathbf{b}_{N+1}, \end{aligned}$$

- Non-negative condition on $\mathbf{v}_i, \mathbf{v}_i \geq \mathbf{0}, \forall i = 1, ..., N$ is included in set \mathcal{O}
- Problem: Non-convex problem.
- Solution: Reformulate the objective function into the difference of two convex functions and solve iteratively using the Difference of Convex Approach!

Difference of Convex Approach



$$c(\pmb{\lambda}) +
ho \sum_{i=1}^N \mathbf{y}_i^ op \mathbf{v}_i$$

$$\underbrace{c(\lambda) + \frac{\rho}{4} \sum_{i=1}^{N} \|\mathbf{y}_i + \mathbf{v}_i\|_2^2 - \underbrace{\frac{\rho}{4} \sum_{i=1}^{N} \|\mathbf{y}_i - \mathbf{v}_i\|_2^2}_{f_2(\mathbf{y}, \mathbf{v}), \text{convex function}}$$

[Jara-Moroni et al., 2018]

Difference of Convex Approach



$$c(\lambda) + \rho \sum_{i=1}^{N} \mathbf{y}_{i}^{\top} \mathbf{v}_{i} \qquad \Longrightarrow \qquad \underbrace{c(\lambda) + \frac{\rho}{4} \sum_{i=1}^{N} \|\mathbf{y}_{i} + \mathbf{v}_{i}\|_{2}^{2} - \underbrace{\frac{\rho}{4} \sum_{i=1}^{N} \|\mathbf{y}_{i} - \mathbf{v}_{i}\|_{2}^{2}}_{f_{2}(\mathbf{y}, \mathbf{v}), \text{convex function}} + \underbrace{\frac{\rho}{4} \sum_{i=1}^{N} \|\mathbf{y}_{i} - \mathbf{v}_{i}\|_{2}^{2}}_{f_{2}(\mathbf{y}, \mathbf{v}), \text{convex function}}$$

Algorithm 1 DCA for solving the bilinear problem

input: feasible guess \mathbf{d}^0 , \mathbf{v}^0 , \mathbf{v}^0 , λ^0 , tolerance $\varepsilon_{\text{tol}} > 0$ output: optimized NN input d* repeat

set $k \leftarrow k + 1$. get $\mathbf{d}^{k+1}, \mathbf{v}^{k+1}, \mathbf{v}^{k+1}, \lambda^{k+1}$ by solving DCA subproblem until $f(\mathbf{y}^k, \mathbf{v}^k) - f(\mathbf{y}^{k+1}, \mathbf{v}^{k+1}) < \varepsilon_{\text{tol}}$ return $\mathbf{d}^{\star} \leftarrow \mathbf{d}^{k+1}$

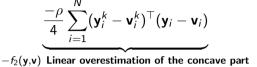
 $f(\mathbf{v}^k, \mathbf{v}^k) = f_1(\mathbf{v}^k, \mathbf{v}^k) - f_2(\mathbf{v}^k, \mathbf{v}^k)$. Here, we write $\mathbf{d}, \mathbf{y}, \mathbf{v}, \lambda$ to represent $\mathbf{d}_i, \mathbf{y}_i, \mathbf{v}_i, \lambda_i \forall i = 1, ..., N$

[Jara-Moroni et al., 2018]

Difference of Convex Approach Solution



$$\underbrace{-\frac{\rho}{4}\sum_{i=1}^{N}||\mathbf{y}_{i}-\mathbf{v}_{i}||_{2}^{2}}_{-f_{2}(\mathbf{y},\mathbf{v}),\text{concave function}} = =$$



Difference of Convex Approach Solution



$$\underbrace{-\frac{\rho}{4}\sum_{i=1}^{N}\lVert\mathbf{y}_{i}-\mathbf{v}_{i}\rVert_{2}^{2}}_{-f_{2}(\mathbf{y},\mathbf{v}),\text{concave function}} \Rightarrow \underbrace{\frac{-\rho}{4}\sum_{i=1}^{N}(\mathbf{y}_{i}^{k}-\mathbf{v}_{i}^{k})^{\top}(\mathbf{y}_{i}-\mathbf{v}_{i})}_{-f_{2}(\mathbf{y},\mathbf{v})} \text{ Linear overestimation of the concave part}$$

DCA subproblem

» XJara-Moroni et al., 2018

$$-\frac{\rho}{2}\sum_{i=1}^{N}(\mathbf{y}_{i}^{k}-\mathbf{v}_{i}^{k})^{\top}(\mathbf{y}_{i}-\mathbf{v}_{i})$$
 subject to $\mathbf{y}_{1}=\mathbf{W}_{1}\mathbf{d}+\mathbf{b}_{1}+\mathbf{v}_{1}$, $\mathbf{v}_{i}=\mathbf{W}_{i}\mathbf{v}_{i-1}+\mathbf{b}_{i}+\mathbf{v}_{i},\ \forall i=2,\ldots,N$

 $\lambda = \mathbf{W}_{N+1}\mathbf{y}_N + \mathbf{b}_{N+1}$

 $\underset{(\mathbf{d}, \mathbf{y}, \mathbf{v}, \boldsymbol{\lambda}) \in \mathcal{O}}{\text{minimize}} c(\boldsymbol{\lambda}) + \frac{\rho}{4} \sum_{i=1}^{N} ||\mathbf{y}_i + \mathbf{v}_i||_2^2$

Linear overestimation of the concave part

Sensitivity to penalty parameter ρ



The solution is extremely sensitive to ρ

- Too Large: Convergence slow
 - Too Small: Violate complimetarity condition

Need to find the smallest ρ that satisfies the complimentarity condition.

General compact formulation minimize $\mathbf{c}^{\top}\mathbf{v}$ d.v.v subject to $Ad \geqslant f$. Vd + Wv + b = v. $0 \leq \mathbf{v}_i \perp \mathbf{v}_i \geqslant \mathbf{0}$. $\forall i = 1, \ldots, N$.

 $\mathsf{d},\mathsf{y},\mathsf{v}$ subject to $\mathsf{Ad}\geqslant \mathsf{f}$: $\lambda,$

 $\mathbf{y}_i, \mathbf{v}_i \geqslant 0,$ $\forall i = 1, \dots, N$

 $Vd + Wv + b = v : \mu$

Part of the derivation draw inspiration from Prop.16 in[Jara-Moroni et al., 2018]

Sensitivity to penalty parameter ρ



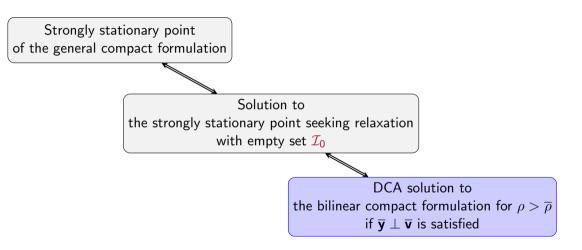
Strongly stationary point seeking relaxation

```
Let (\tilde{\mathbf{d}}, \tilde{\mathbf{v}}, \tilde{\mathbf{v}}) be some feasible point for the general compact formulation
                 minimize \mathbf{c}^{\top}\mathbf{v}
                       d.y.v
             subject to Ad \geqslant f : \lambda.
                                                Vd + Wy + b = v : \mu
                                                                                                                                                                        \mathcal{I}_{\nu}(\tilde{\mathbf{v}},\tilde{\mathbf{v}}) \triangleq \{i \mid \tilde{\mathbf{v}}_i = 0 < \tilde{\mathbf{v}}_i\},
                                                y_i = 0 : \mu_i^{y_0} \quad \forall i \in \mathcal{I}_{v}(\tilde{\mathbf{y}}, \tilde{\mathbf{v}}),
                                                                                                                                                                        \mathcal{I}_{V}(\tilde{\mathbf{y}},\tilde{\mathbf{v}}) \triangleq \{i \mid \tilde{y}_{i} > 0 = \tilde{v}_{i}\},\
                                                y_i \geqslant 0 : \mu_i^{y_+} \quad \forall i \in \mathcal{I}_{\nu}(\tilde{\mathbf{y}}, \tilde{\mathbf{v}}),
                                                                                                                                                                         \mathcal{I}_0(\tilde{\mathbf{v}},\tilde{\mathbf{v}}) \triangleq \{i \mid \tilde{\mathbf{v}}_i = 0 = \tilde{\mathbf{v}}_i\}.
                                                 v_i = 0 : \mu_i^{v_0} \quad \forall i \in \mathcal{I}_{v}(\tilde{\mathbf{y}}, \tilde{\mathbf{v}}),
                                                 v_i \geqslant 0 : \mu_i^{v_+} \quad \forall i \in \mathcal{I}_{v}(\tilde{\mathbf{y}}, \tilde{\mathbf{v}}),
                                                v_i, v_i \geqslant 0, \quad \forall i \in \mathcal{I}_0(\tilde{\mathbf{v}}, \tilde{\mathbf{v}}),
```

The solution $(\overline{\mathbf{d}}, \overline{\mathbf{y}}, \overline{\mathbf{v}})$ is the strongly stationary point of the general compact formulation if $\overline{\mathbf{y}}_i \perp \overline{\mathbf{v}}_i \quad \forall i=1,...,N$, where the strongly stationary point is the point that satisfies the KKTs of the original non-convex optimization.

Connections between the three cases





Obtain feasible point $(\tilde{\mathbf{d}}, \tilde{\mathbf{y}}, \tilde{\mathbf{v}})$





Obtain feasible point $(\tilde{\mathbf{d}}, \tilde{\mathbf{y}}, \tilde{\mathbf{v}})$



Feasibility insight of trained neural network

- $\mathcal{I}_0(\tilde{\mathbf{y}}, \tilde{\mathbf{v}}) = \emptyset$ $\tilde{\mathbf{y}}_i \neq \tilde{\mathbf{v}}_i, \forall i = 1, \dots, N$ $\mathbf{W}_i \mathbf{y}_{i-1} + \mathbf{b}_i \neq \mathbf{0}, \forall i = 1, \dots, N$
 - Often true
 - If not, we can losslessly reduce NN size —- no more zero neuron!

Computing ρ from stationary point



- $(\overline{\mathbf{d}}, \overline{\mathbf{y}}, \overline{\mathbf{v}})$: From strongly stationary point seeking relaxation
- **Obtain** ρ : From partial Karush–Kuhn–Tucker conditions of the strongly stationary point seeking relaxation formulation and the bilinear compact formulation
- Combine lagrangian multipliers $\mu^{v} = [\mu^{v_0 \top} \mu^{v_+ \top}]^{\top}$, $\mu^{y} = [\mu^{y_0 \top} \mu^{y_+ \top}]^{\top}$

Relaxation partial KKT

$$\begin{aligned} \mathbf{A}^{\top} \boldsymbol{\lambda} + \mathbf{V}^{\top} \boldsymbol{\mu}^{\boldsymbol{\nu}} &= \mathbf{0}, \\ \mathbf{W}^{\top} \boldsymbol{\mu}^{\boldsymbol{\nu}} + \boldsymbol{\mu}^{\boldsymbol{y}} &= \mathbf{c}, \\ \mathbf{0} &\leqslant \boldsymbol{\lambda} \perp \mathbf{A} \mathbf{d} - \mathbf{f} \geqslant \mathbf{0}, \\ \mathbf{y}^{\top} \boldsymbol{\mu}^{\boldsymbol{y}} &= 0, \quad \mathbf{v}^{\top} \boldsymbol{\mu}^{\boldsymbol{\nu}} &= 0, \\ \boldsymbol{\mu}_{i}^{\boldsymbol{y}} &\geqslant 0, \quad \forall i \in \mathcal{I}_{\boldsymbol{v}}(\tilde{\mathbf{y}}, \tilde{\mathbf{v}}) \\ \boldsymbol{\mu}_{i}^{\boldsymbol{v}} &\geqslant 0, \quad \forall i \in \mathcal{I}_{\boldsymbol{v}}(\tilde{\mathbf{y}}, \tilde{\mathbf{v}}) \end{aligned}$$

Bilinear partial KKT

$$\mathbf{A}^{ op} oldsymbol{\lambda} + \mathbf{V}^{ op} oldsymbol{\mu} = \mathbf{0},$$

$$\mathbf{0} \leqslant \mathbf{y} \perp - \mathbf{W}^{\top} \boldsymbol{\mu} + \mathbf{c} + \rho \mathbf{v} \geqslant \mathbf{0},$$

$$0 \leqslant \lambda \perp Ad - f \geqslant 0$$
,

$$\mathbf{0} \leqslant \mathbf{v} \perp \rho \mathbf{v} + \boldsymbol{\mu} \geqslant \mathbf{0}.$$

Computation for $\overline{\rho}$



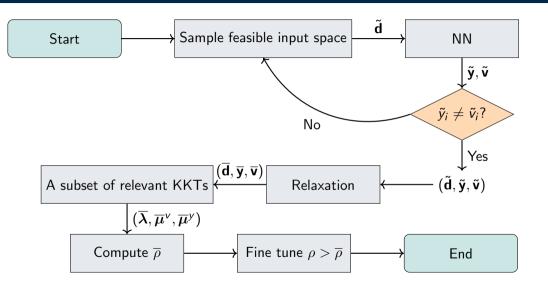
We substitute μ^{ν} from relation partial KKT into μ from bilinear partial KKT. Through some algebraic manipulation, we retrieve the following condition on $\overline{\rho}$

Lower bound for ρ

$$\bar{\rho} \triangleq \max \left\{ 0, \left\{ -\frac{\overline{\mu}_i^y}{\overline{\nu}_i} \middle| \ \overline{\nu}_i > 0 \right\}, \left\{ -\frac{\overline{\mu}_i^v}{\overline{\nu}_i} \middle| \ \overline{y}_i > 0 \right\} \right\}.$$

Computation for DCA penalty term ρ algorithm

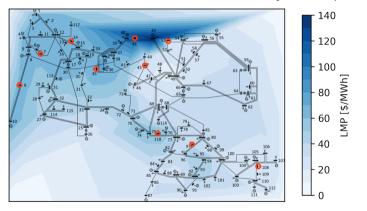




Data center workload allocation in power grids



- Consider a network of spatially distributed data centers.
- Operators allocate workloads to minimize the cost of electricity consumption.



• In practice, we do not have network data or other parameters.

• Use a neural network to map the workload allocation to the cost of electricity. >>> X. Liu

Bilevel optimization to solve data center load allocation



- Goal: Allocate data center load to minimize the cost of electricity consumption.
- **Given:** Forecast of total electricity demand Δ .
- Data center allocation formulated as a bilevel optimization problem where d is the data center loads with upper and lower bounds $\overline{\mathbf{d}}$ and \mathbf{d} . π is the price per demand.
- Challenge: Unknown underlying physical model

minimize
$$oldsymbol{\pi}^{ op} \mathbf{d}$$
 subject to $oldsymbol{\underline{d}} \leqslant \mathbf{d} \leqslant oldsymbol{\overline{d}}$ $oldsymbol{1}^{ op} \mathbf{d} = \Delta$ $oldsymbol{\pi} \in$ Unknown model

Bilevel optimization to solve data center load allocation



- Goal: Allocate data center load to minimize the cost of electricity consumption.
- **Given:** Forecast of total electricity demand Δ .
- Data center allocation formulated as a bilevel optimization problem where d is the data center loads with upper and lower bounds $\overline{\mathbf{d}}$ and \mathbf{d} . π is the price per demand.
- Challenge: Unknown underlying physical model
- Solution: Replace by a neural network trained with a labeled dataset on historical record $\{(\mathbf{d}_1, \lambda_1), \ldots, (\mathbf{d}_n, \lambda_n)\}.$

$$\begin{array}{ll} \underset{\mathbf{d}}{\mathsf{minimize}} & \boldsymbol{\pi}^{\top}\mathbf{d} \\ \mathsf{subject to} & \underline{\mathbf{d}} \leqslant \mathbf{d} \leqslant \overline{\mathbf{d}} \\ & \mathbf{1}^{\top}\mathbf{d} = \Delta \\ & \boldsymbol{\pi} \in \\ & \mathsf{Unknown model} \end{array}$$

minimize $\mathbb{N}\mathbb{N}(\mathbf{d})$ subject to $\mathbf{d} \leqslant \mathbf{d} \leqslant \overline{\mathbf{d}}$

 $\mathbf{1}^{\mathsf{T}}\mathbf{d} = \mathbf{\Lambda}$

Data generation for NN training



• **Method:** Using locational marginal prices (LMPs) as π derived from the optimal power flow (OPF) model. Use this data to train neural network.

Data generation for NN training



• **Method:** Using locational marginal prices (LMPs) as π derived from the optimal power flow (OPF) model. Use this data to train neural network.

$$\label{eq:continuity} \begin{split} \underset{\underline{\mathbf{p}} \leqslant \mathbf{p} \leqslant \overline{\mathbf{p}}}{\text{minimize}} \quad \mathbf{p}^\top \mathbf{C} \mathbf{p} + \mathbf{c}^\top \mathbf{p} \\ \text{subject to} \quad \mathbf{1}^\top (\mathbf{p} - \boldsymbol{\ell} - \mathbf{d}) = 0 \quad : \lambda \\ |\mathbf{F} (\mathbf{p} - \boldsymbol{\ell} - \mathbf{d})| \leqslant \overline{\mathbf{f}} \quad : \underline{\boldsymbol{\mu}}, \overline{\boldsymbol{\mu}} \end{split}$$

- p: generator dispatch
- $[\underline{p}, \overline{p}]$: feasible range of dispatch
- ℓ: conventional loads
- d: data center loads
- **F**: the matrix of power transfer distribution factors, which maps net power injections $(\mathbf{p} \ell \mathbf{d})$ to power flows as $\mathbf{F}(\mathbf{p} \ell \mathbf{d})$
- $\overline{\mathbf{f}}$: line capacity

Data generation for NN training



• **Method:** Using locational marginal prices (LMPs) as π derived from the optimal power flow (OPF) model. Use this data to train neural network.

$$\label{eq:continuity} \begin{split} \underset{\underline{\mathbf{p}} \leqslant \mathbf{p} \leqslant \overline{\mathbf{p}}}{\text{minimize}} \quad \mathbf{p}^\top \mathbf{C} \mathbf{p} + \mathbf{c}^\top \mathbf{p} \\ \text{subject to} \quad \mathbf{1}^\top (\mathbf{p} - \boldsymbol{\ell} - \mathbf{d}) = \mathbf{0} \quad : \lambda \\ |\mathbf{F} (\mathbf{p} - \boldsymbol{\ell} - \mathbf{d})| \leqslant \overline{\mathbf{f}} \quad : \underline{\boldsymbol{\mu}}, \overline{\boldsymbol{\mu}} \end{split}$$

Locational Marginal Price(LMP)

$$oldsymbol{\pi} = \mathbf{1} \lambda^\star - \mathbf{F}^ op \overline{oldsymbol{\mu}}^\star + \mathbf{F}^ op oldsymbol{\mu}^\star$$

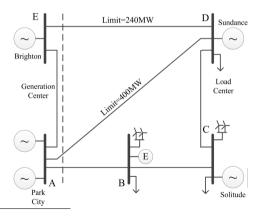
[Chatzivasileiadis, 2018]

- p: generator dispatch
- $[p, \overline{p}]$: feasible range of dispatch
- ℓ: conventional loads
- d: data center loads
- **F**: the matrix of power transfer distribution factors, which maps net power injections $(\mathbf{p} \ell \mathbf{d})$ to power flows as $\mathbf{F}(\mathbf{p} \ell \mathbf{d})$
- $\overline{\mathbf{f}}$: line capacity

Experiment on PJM-5 bus system

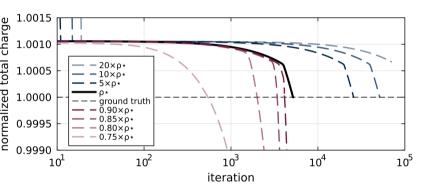


- 3 loads act as data centers with demand in the range of [0.8, 1] of the nominal value.
- NN with 2 hidden layers, 50 neurons each. 10,000 training samples
- Global solution is computed using SOS1 constraints and Gurobi solver.



Experiment on PJM-5 bus system

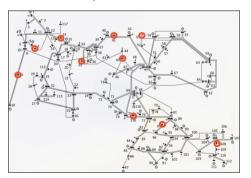




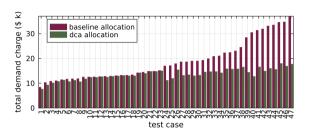
- ρ^{\star} : fast convergence
- $\rho < \rho^{\star}$: Does not satisfy complementarity condition
- $\rho > \rho^{\star}$: Slow convergence



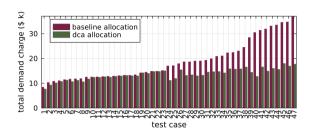
- NN training: NN with 5-hidden layer and 1,000 ReLUs, trained on 12,500 demand samples in the range of [0,100] MWh.
- DCA solution: Computes a new demand allocation to improve the baseline.
- **Samples:** Randomly picked 50 baseline cases, of which DCA could successfully converge and 47 of them roughly match the output of the OPF model.





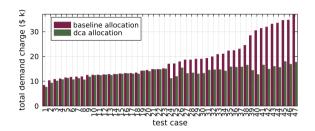


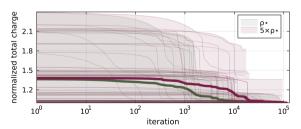




- First 23 cases: The savings are small. Lightly congested with little to no benefit from spatial load redistribution
- Remaining cases: Significant improvement





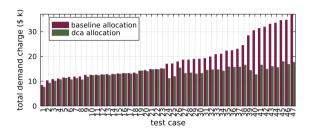


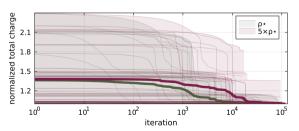
- First 23 cases: The savings are small. Lightly congested with little to no benefit from spatial load redistribution
- Remaining cases: Significant improvement

>>> X. Liu 24 / 38

Experiment on IEEE 118 bus system





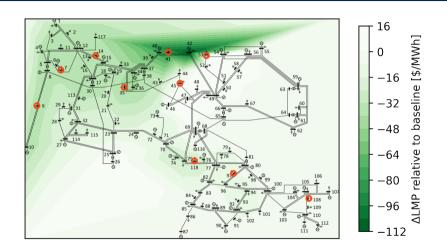


- First 23 cases: The savings are small. Lightly congested with little to no benefit from spatial load redistribution
- Remaining cases: Significant improvement
- The smaller ρ consistently leads to faster convergence.
- Each iteration takes ≈ 0.37 seconds, allowing for convergence in around 1 hour on average.

>>> X. Liu 24 / 3

Experiment on IEEE 118 bus system





- The difference between the LMPs induced at baseline and the DCA solutions.
- LMPs were significantly reduced at the data centers. >>> X. Liu

Conclusions



- Developed a difference-of-convex algorithm to solve decision-making problems with objectives or constraints represented by trained neural networks.
- Avoided the computational bottleneck of ReLU logical constraints through informed penalization of such constraints in the objective function.
- Demonstrated on the data center load allocation application and showed significant reduction in cost compared to baseline.



Questions?



References I



Chatzivasileiadis, S. (2018).

Lecture notes on optimal power flow (opf).

arXiv preprint arXiv:1811.00943.



Enabling certification of verification-agnostic networks via memory-efficient semidefinite programming.

Advances in Neural Information Processing Systems, pages 5318–5331.

Dvorkin, V. (2025). Lecture slides: Computational power systems.

Unpublished lecture slides, accessed on 2025-04-30.

Fazlyab, M., Morari, M., and Pappas, G. J. (2020).

Safety verification and robustness analysis of neural networks via quadratic constraints and semidefinite programming.

IEEE Trans. Autom. Control., 67(1):1–15.

References II



Grimstad, B. and Andersson, H. (2019).

 ${\sf ReLU\ networks\ as\ surrogate\ models\ in\ mixed-integer\ linear\ programs}.$

Comput. Chem. Eng., 131:106580.

Jara-Moroni, F., Pang, J.-S., and Wächter, A. (2018).

A study of the difference-of-convex approach for solving linear programs with complementarity constraints. *Mathematical Programming*, 169:221–254.

Tjeng, V., Xiao, K., and Tedrake, R. (2017).

Evaluating robustness of neural networks with mixed integer programming.

Preprint.

Turner, M. et al. (2024).

 $Pyscipopt-ml:\ Embedding\ trained\ machine\ learning\ models\ into\ mixed-integer\ programs.$

Preprint.

References III





Wang, S. et al. (2021).

Beta-crown: Efficient bound propagation with per-neuron split constraints for neural network robustness verification.

Zhou, A. (2023).

Pjm 5 bus data.

https://dx.doi.org/10.21227/7ex2-1t83.

Advances in neural information processing systems, 34:29909–29921.

Accessed on 2025-04-30.

The End

Backup I: Derivation of LMP



Partial Lagrangian function (dualize the coupling constraints only):

$$egin{aligned} \mathcal{L}(\mathbf{p},\lambda,\overline{m{\mu}},\underline{m{\mu}}) = & c(\mathbf{p}) - \lambda \mathbf{1}^{ op}(\mathbf{p}-\mathbf{d}) \ &+ \overline{m{\mu}}^{ op}(\mathbf{F}(\mathbf{p}-\mathbf{d})-\overline{\mathbf{f}}) + m{\mu}^{ op}(-\mathbf{F}(\mathbf{p}-\mathbf{d})-\overline{\mathbf{f}}) \end{aligned}$$

Group terms corresponding to dispatch \mathbf{p} , demand \mathbf{d} and line limits \mathbf{f} :

$$\mathcal{L} = \mathcal{L}^{f p} + \mathcal{L}^{f d} + \mathcal{L}^{f f}, \quad ext{where}$$

 $\mathcal{L}^{\mathbf{p}}(\mathbf{p}, \lambda, \overline{\mu}, \mu) = c(\mathbf{p}) - (\mathbf{1}\lambda - \mathbf{F}^{\top}\overline{\mu} + \mathbf{F}^{\top}\mu)^{\top}\mathbf{p}$

$$\mathcal{L}^{\mathsf{d}}(\lambda, \overline{\mu}, \underline{\mu}) = (\mathbf{1}\lambda - \mathbf{F}^{\top} \overline{\mu} + \mathbf{F}^{\top} \underline{\mu})^{\top} \mathbf{d}$$
$$\mathcal{L}^{\mathsf{f}}(\overline{\mu}, \underline{\mu}) = -(\underline{\mu} + \overline{\mu})^{\top} \overline{\mathsf{f}}$$

Power dispatch **p** and demand **d** share the same multiplier but with oposite signs

Backup I: Derivation of LMP



$$\pi^{\star}(\lambda^{\star}, \overline{\mu}^{\star}, \underline{\mu}^{\star}) = \underbrace{\mathbf{1}\lambda^{\star}}_{\mathsf{uniform}} - \underbrace{\mathbf{F}^{\top}(\overline{\mu}^{\star} - \underline{\mu}^{\star})}_{\mathsf{congestion}} \in \mathbb{R}^{N}$$

- π_n^{\star} is the cost of supplying the next unit of demand at node n
- ullet in case of congestion $(\overline{\mu}^{\star}>0$ or $\overline{\mu}^{\star}>0)$, electricity price varies across the grid
- price at the reference bus is λ^* (the ref. column of **F** is zero)

Backup II: DC-OPF (PTDF formulation)



- Use matrix $\mathbf{F} \in \mathbb{R}^{E \times N}$ of power transfer distribution factors (PTDF)
 - how the power flow in line e changes w.r.t. to the change of power injection at node n?
 - obtained by manipulating the DC bus admittance matrix B
- Power flows $f = \mathbf{F}(\mathbf{p} \mathbf{d})$ (distribution of net injections across power lines) minimize $c(\mathbf{p})$ generation cost subject to $\mathbf{1}^{\top}(\mathbf{p} \mathbf{d}) = 0$ active power balance $|\mathbf{F}(\mathbf{p} \mathbf{d})| \leqslant \overline{f}$ power flow limits $\mathbf{p} \leqslant \mathbf{p} \leqslant \overline{\mathbf{p}}$ min/max gen p-limits

[Dvorkin, 2025]

Backup III: Big-M



Constraints:

$$z\geqslant 0$$
 $\mu\geqslant 0$
 $z\mu=0$

Big-M reformulation:

$$\left\{\begin{array}{l} \mu \leqslant Mu \\ z \leqslant M(1-u) \end{array}\right.$$

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

Backup IV: Semidefinate Programming



$$egin{aligned} x_{i+1} &= \mathsf{ReLU}\left(\mathbb{L}_i\left(x_i
ight)
ight). \ & x_{i+1} \geqslant 0\left[\lambda_i^\mathrm{a}\right], x_{i+1} \geqslant \mathbb{L}_i\left(x_i
ight)\left[\lambda_i^\mathrm{b}\right] \ & x_{i+1} \odot \left(x_{i+1} - \mathbb{L}_i\left(x_i
ight)
ight) \leqslant 0\left[\lambda_i^\mathrm{c}\right], x_i \odot x_i - \left(\ell_i + u_i
ight) \odot x_i + \ell_i \odot u_i \leqslant 0\left[\lambda_i^\mathrm{d}\right]. \ & x_{i+1}, \lambda_i) = \left(-x_{i+1}
ight)^\top \lambda_i^\mathrm{a} + \left(\mathbb{L}_i\left(x_i
ight) - x_{i+1}
ight)^\top \lambda_i^\mathrm{b} \end{aligned}$$

$$\mathcal{L}\left(x_{i}, x_{i+1}, \lambda_{i}\right) = \left(-x_{i+1}\right)^{\top} \lambda_{i}^{\mathbf{a}} + \left(\mathbb{L}_{i}\left(x_{i}\right) - x_{i+1}\right)^{\top} \lambda_{i}^{\mathbf{b}} \\ + \left(x_{i+1} \odot \left(x_{i+1} - \mathbb{L}_{i}\left(x_{i}\right)\right)\right)^{\top} \lambda_{i}^{\mathbf{c}} + \left(x_{i} \odot x_{i} - \left(\ell_{i} + u_{i}\right) \odot x_{i} + \ell_{i} \odot u_{i}\right)^{\top} \lambda_{i}^{\mathbf{d}} \\ = \underbrace{\left(\ell_{i} \odot u_{i}\right)^{\top} \lambda_{i}^{\mathbf{d}}}_{\text{independent of } x_{i}, x_{i+1}} - \underbrace{x_{i+1}^{\top} \lambda_{i}^{\mathbf{a}} + \left(\mathbb{L}_{i}\left(x_{i}\right)\right)^{\top} \lambda_{i}^{\mathbf{b}} - x_{i+1}^{\top} \lambda_{i}^{\mathbf{b}} - x_{i}^{\top} \left(\left(\ell_{i} + u_{i}\right) \odot \lambda_{i}^{\mathbf{d}}\right)}_{\text{linear in } x_{i}, x_{i+1}} \\ + x_{i+1}^{\top} \operatorname{diag}\left(\lambda_{i}^{\mathbf{c}}\right) x_{i+1} - x_{i+1}^{\top} \operatorname{diag}\left(\lambda_{i}^{\mathbf{c}}\right) \mathbb{L}_{i}\left(x_{i}\right) + x_{i}^{\top} \operatorname{diag}\left(\lambda_{i}^{\mathbf{d}}\right) x_{i} \,.$$

Quadratic in x_i, x_{i+1}

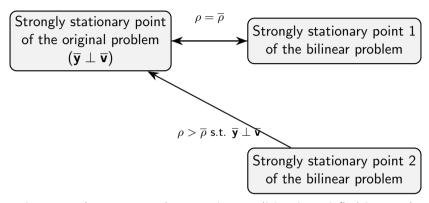
Backup V: Neuron split and bound propagation



- ReLU relaxation: w^{\top} ReLU $(v) \ge w^{\top}$ **D**v + b' where **D** is a diagonal matrix containing free variables $0 \le \alpha_j \le 1$ only when $\mathbf{u}_j > 0 > \mathbf{I}_j$ and $w_j \ge 0$, while its rest values as well as constant b' are determined by $\mathbf{I}, \mathbf{u}, w$.
- $\min_{x \in \mathcal{C}} f(x) \ge \min_{x \in \mathcal{C}} \mathbf{a}_{CROWN}^{\top} x + c_{CROWN}$ where \mathbf{a}_{CROWN} and c_{CROWN} can be computed using $\mathbf{W}^{(i)}, \mathbf{b}^{(i)}, \mathbf{l}^{(i)}, \mathbf{u}^{(i)}$ in polynomial time.

Backup VI: Additional argument for ρ





- The fine tuning to make sure complementarity condition is satisfied is easy because intuitively, if $\overline{\rho}$ works as a penalty parameter in one case, the penalty is penalize the bilinear part enough to satisfy the complementarity condition.
- This result is also proven experimentally.

»» X. Liu