

Explainable Models via Data-Driven Optimization

Howard Heaton

Typal Academy

Setting

Solve problems where approximate optimization models can be hand-crafted

Learning to Optimize (L2O)

Make parameterized optimization model and use training data to tune it, *i.e.*

(model output) , argmin (prior knowledge) + (data-driven terms)

Goal for Today

Outline the tools needed for audience to create their own L2O models

Machine Learning

$$N(d) = (W^m + b^m) \quad (W^1 d + b^1)$$

- ✓ adapt to available data
- ✗ satisfy constraints / optimality
- ✓ expressive capacity
- ✓ flexible architectures

Traditional Optimization

$$\operatorname{argmin}_{x \in \mathcal{C}} f(x)$$

- ✗ adapt to available data
- ✓ guaranteed optimality
- ✓ interpretable models
- ✓ scalable first-order algorithms

Learning fast approximations of sparse coding by Gregor and LeCun in 2010

Sparse codes can be modeled as solutions to the ℓ_1 regularized problem

$$\min_x \frac{1}{2} \|Ax - d\|_2^2 + \lambda \|x\|_1$$

Proximal gradient updates look like

$$x^{k+1} = \text{shrink} \left(x^k + A^T(Ax^k - d); \lambda \right); \quad = \text{shrink} \left(\underbrace{W_1 x^k + W_2 d}_{\text{feed forward layer}}; \lambda \right);$$

where $W_1 = I - A^T A$, $W_2 = A^T$, and $\text{shrink}(z; \lambda) = \text{sign}(z) \max(|z| - \lambda, 0)$

LISTA Idea: Treat x^k as network output and tune W_1 and W_2 from training data

- 1 Make a network by parameterizing an optimization problem to get

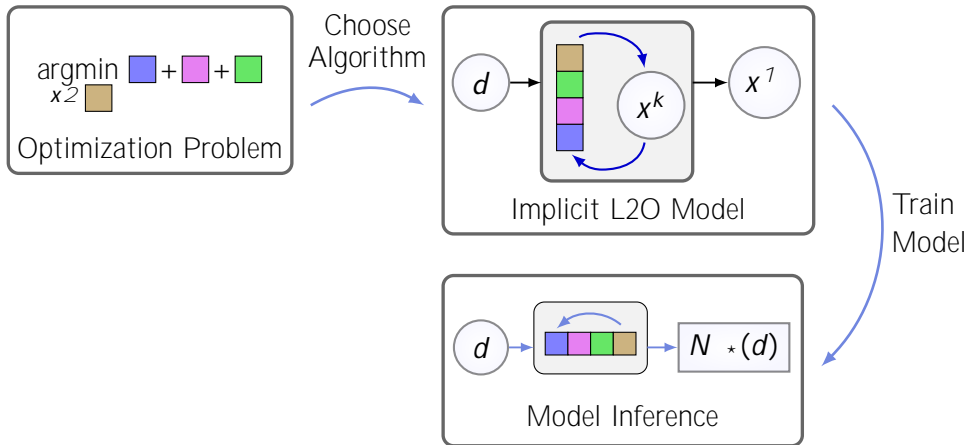
$$N(d), \operatorname{argmin}_x f(x; d)$$

Note: Constraints can be included in this formulation

Note: More general L2O models can use black box and optimization layers

- 2 Forward prop consists of applying an apt first-order algorithm *until convergence*
- 3 Backprop consists of using built-in autograd *on last step of forward prop*

L2O via Building Blocks



Task

Recover a signal $x_d^?$ from linear measurements $d = Ax_d^?$

Key Knowledge

Signal $x_d^?$ has low dimensional structure (but is *not* sparse)

L2O Model

For a “sparsifying matrix” K , we can estimate

$$x_d^? = \underset{x}{\operatorname{argmin}} \|Kx\|_1 \quad \text{s.t. } Ax = d$$

Question: Even if this model is “right” for some choice of K , how do we find K ?

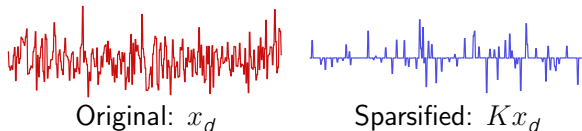


Figure 1: Applying the learned K sparsifies x_d^* (shown for test data d)

Fix weights $K \in \mathbb{R}^{250 \times 250}$, noting $x_d^* \in \mathbb{R}^{250}$ and $d \in \mathbb{R}^{100}$, and let model N be

$$N(d) = \operatorname{argmin}_x \|Kx\|_1 \text{ s.t. } Ax = d$$

For a distribution of measurement/signal pairs $(d; x_d^*)$, train model by minimizing

$$\min \mathbb{E}_d \|Kx_d^* - N(d)\|_2^2$$

This ensures the weights K are tuned as well as possible for the task at hand

Toy Example – Results

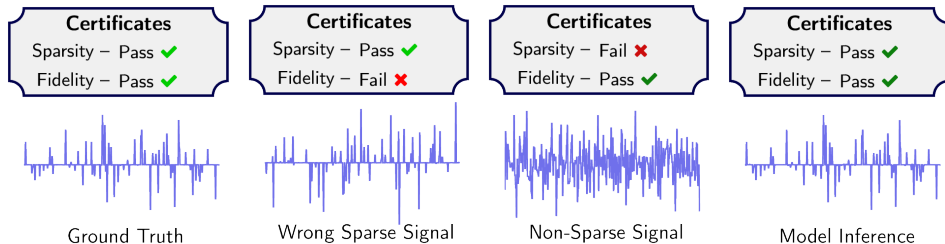


Figure 2: Example inferences for test data d with sparsified Kx of each inference x shown

```
x_fxd_pt = find_fixed_point(d)
pred = apply_opt_update(x_fxd_pt, d)
loss = criterion(pred, labels)
loss.backward()
optimizer.step()
```

Figure 3: Sample PyTorch code for backpropagation. Here `x_fxd_pt` is the *same* as `pred`, except that `pred` has gradients attached and `x_fxd_pt` does not

Informal Theorem:¹ Backpropping through the final step of a fixed point optimization algorithm (as shown above) yields a *preconditioned* gradient

¹Wu Fung, et al. *JFB: Jacobian-Free Backpropagation for Implicit Networks*. 2022.

Prior Knowledge Embedding

Can use regularizer to promote sparsity, low-rank, smooth, laying on manifold, ...

Hard Constraints

Constraints can be put into model (*e.g.* simplex, linear system) and forward propagation can continue until these are met to given tolerance

L O model can inherit all guarantees desired from traditional optimization

1 L2O Overview

2 Numerical Examples

3 Appendix

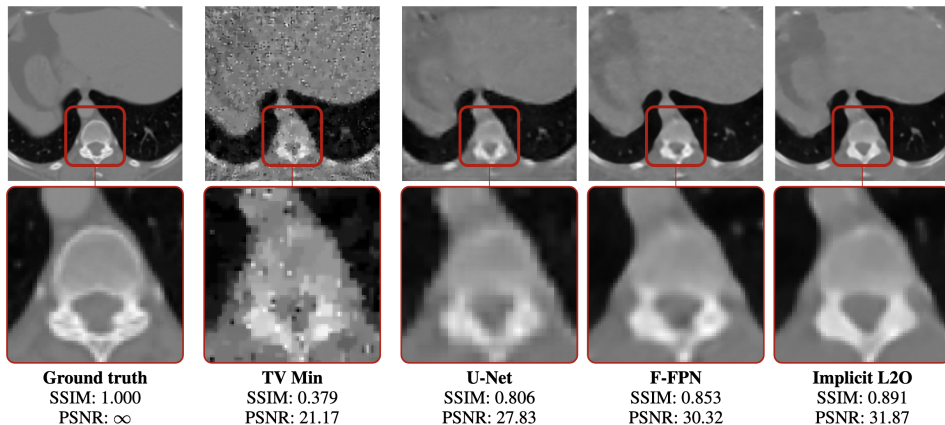


Figure 4: Comparison of techniques, ranging from traditional to fully data-driven

$$(\text{L2O Model}) = N(d), \quad \underset{x \in [0;1]^n}{\operatorname{argmin}} f(Kx) \quad \text{s.t.} \quad KAx = dk$$

Task

Given contextual data d (e.g. weather, construction), predict traffic distribution

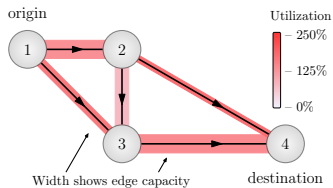
L2O Model

Assume drivers are self-interested and most likely behavior is a Nash equilibria, where game gradient is *unknown*

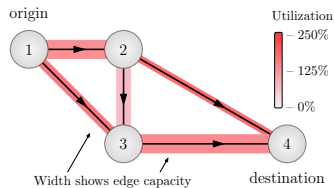
Constraints

Vertex-arc incidence matrix characterizes road network, yielding large linear constraints to get valid traffic flows

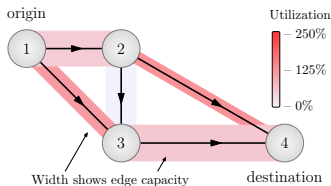
Traffic Routing – Toy Example



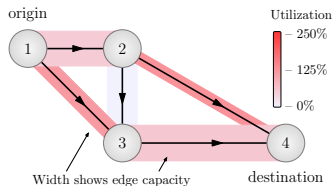
(a) Rainy Day Prediction



(b) Rain Day True



(c) Sunny Day Prediction



(d) Sunny Day True

Figure 5: Predictions of toy traffic from context data d

- Heaton, Wu Fung. *Explainable AI via Learning to Optimize*. 2022.
- McKenzie, Wu Fung, Heaton. *Faster Predict-and-Optimize with Three-Operator Splitting*. 2023.
- Wu Fung, et al. *JFB: Jacobian-free Backpropagation for Implicit Networks*. 2022
- Heaton, Wu Fung, Gibali. *Feasibility-based Fixed Point Networks*. 2021.
- Heaton, et al. *Learn to Predict Equilibria via Fixed Point Networks*. 2021.
- Gilton, Ongie, Willet. *Deep Equilibrium Architectures for Inverse Problems in Imaging*. 2022.
- Heaton, et al. *Safeguarded Learned Convex Optimization*. 2023.
- de Avila Belbute-Peres, et al. *End-to-End Differentiable Physics for Learning and Control*. 2018.

Setting

Solve problems where approximate optimization models can be hand-crafted

Learning to Optimize (L2O)

- Make parameterized optimization model

(model output) , $\operatorname{argmin}(\text{prior knowledge}) + (\text{data-driven terms})$

- Training is “easy” using Jacobian-Free Backpropagation (JFB)
- Outputs leverage available data, have strong guarantees, and are interpretable

1 L2O Overview

2 Numerical Examples

3 Appendix

The toy problem can be generalized to a model of the form

$$\min_x f(Kx) + h(x) \quad \text{s.t.} \quad \|Mx - d\| \leq \kappa;$$

with proximal f and g , can be rewritten as

$$\min_{x; w; p} f(p) + h(x) + \mathbb{1}_{B(d; \kappa)}(w) \quad \text{s.t.} \quad \begin{matrix} 2 & 3 & 2 & 3 \\ 4 & K & 5 & x \\ & M & & 4 & p & 5 \\ & & & & w & = 0; \end{matrix}$$

where p and w are auxiliary variables and $\mathbb{1}_{B(d; \kappa)}$ is the indicator function that is 0 inside the Euclidean ball $B(d; \kappa)$ and ∞ elsewhere

Using Linearized ADMM updates yields the iterative updates

$$p^{k+1} = \text{prox}_f \left(p^k + \left(\frac{k}{1} + (Kx^k - p^k) \right) \right)$$

$$w^{k+1} = \text{proj}_{B(d;)} \left(w^k + \left(\frac{k}{2} + (Mx^k - w^k) \right) \right)$$

$$\frac{k+1}{1} = \frac{k}{1} + (Kx^k - p^{k+1})$$

$$\frac{k+1}{2} = \frac{k}{2} + (Mx^k - w^{k+1})$$

$$r^k = K \frac{k+1}{1} - \frac{k}{1} + M \frac{k+1}{2} - \frac{k}{2}$$

$$x^{k+1} = \text{prox}_h \left(x^k + r^k \right);$$

which can be coded into the `apply_opt_update` function