

**Include Optimization as a
Layer in NN**

Outline

- What is OptNet / feedforward pass \Rightarrow
- How to update parameters in the OptNet / backward pass \Leftarrow
- Application using OptNet
 - Representation power
 - MINIST \Leftarrow OptNet
 - Sudoku \Leftarrow

implicit differentiation.

OptNet: Differentiable Optimization as a Layer in Neural Networks

What is OptNet

- Feedforward NN

Form $\underline{z_{i+1}} = \sigma(\underline{W_i z_i} + \underline{b_i})$

Parameters $\underline{W_i}, \underline{b_i}$

- OptNet

$z_{i+1} = \underset{z}{\operatorname{argmin}} \frac{1}{2} z^T Q(z_i) z + q(z_i)^T z$ *quadratic objective*
subject to $A(z_i)z = b(z_i)$ \Rightarrow 等式 (1)
 $G(z_i)z \leq h(z_i)$ \Rightarrow 不等式.

$\Rightarrow \underline{Q(z_i)}, \underline{q(z_i)}, \underline{A(z_i)}, \underline{b(z_i)}, \underline{G(z_i)}, \underline{h(z_i)}$

Can depend on the previous layer z_i

What is OptNet

- Feedforward NN

Form $z_{i+1} = \sigma(W_i z_i + b_i)$

Parameters W_i, b_i

Gradients $\frac{\partial z_{i+1}}{\partial W_i} = \frac{\partial}{\partial W_i}(\text{relu}(W_i z_i + b_i))$

$$\frac{\partial L}{\partial W_i} = \frac{\partial L}{\partial z_{i+1}} \cdot \frac{\partial z_{i+1}}{\partial W_i} = \text{diag}(\text{relu}(\text{sign}(b + W_i z_i)))_i^T \otimes z_i$$

- OptNet

$$z_{i+1} = \underset{z}{\text{argmin}} \frac{1}{2} z^T Q(z_i) z + q(z_i)^T z$$

subject to $A(z_i) z = b(z_i)$

$$G(z_i) z \leq h(z_i)$$

$$Q(z_i) \left[q(z_i), A(z_i), b(z_i), G(z_i), h(z_i) \right]$$

$$\frac{\partial z_{i+1}}{\partial Q(z_i)}, \frac{\partial z_{i+1}}{\partial q(z_i)}, \dots, ?$$

How to update parameters in the OptNet

$$\begin{aligned}
 z_{i+1} = \operatorname{argmin}_z \quad & \frac{1}{2} z^T Q(z_i) z + q(z_i)^T z \\
 \text{subject to} \quad & A(z_i) z = b(z_i) \\
 & G(z_i) z \leq h(z_i)
 \end{aligned} \tag{1}$$

- Output of $i, i + 1$ layer, Variable $z_i, z_{i+1} : \mathbb{R}^n$
- Problem data of optimization problem, Parameter Q, q, A, b, G, h

$$\underline{Q \in \mathbb{R}^{n \times n}}, \underline{q \in \mathbb{R}^n}, \underline{A \in \mathbb{R}^{m \times n}}, \underline{b \in \mathbb{R}^m}, \underline{G \in \mathbb{R}^{p \times n}}, \underline{h \in \mathbb{R}^p}$$

- Gradients: derivative of the solution to its parameter

$$\frac{\partial z_{i+1}}{\partial Q(z_i)}, \frac{\partial z_{i+1}}{\partial q(z_i)}, \dots, ?$$

\Downarrow m 个等式约束 \Downarrow p 个不等式

How to update parameters in the OptNet

$$\begin{aligned}
 z_{i+1} = \operatorname{argmin}_z \quad & \frac{1}{2} z^T Q(z_i) z + q(z_i)^T z \\
 \text{subject to} \quad & A(z_i) z = b(z_i) \\
 & G(z_i) z \leq h(z_i)
 \end{aligned} \tag{1}$$

$Qz \rightarrow dQz + Qdz$

• Write Lagrangian $L(z, \nu, \lambda) = \frac{1}{2} z^T Q z + q^T z + \nu^T (A z - b) + \lambda^T (G z - h)$

• Write KKT conditions

$$\frac{dL}{dz} = 0, \frac{dL}{d\lambda} = 0, \frac{dL}{d\nu} = 0$$

dual variable

$$\begin{aligned}
 Qz^* + q + A^T \nu^* + G^T \lambda^* &= 0 \\
 Az^* - b &= 0 \\
 D(\lambda^*)(Gz^* - h) &= 0,
 \end{aligned}$$

①

z^*, ν^*, λ^* are the optimal primal and dual variables

$$z^* = z_{i+1}$$

How to update parameters in the OptNet

$$\begin{aligned}
 z_{i+1} = \operatorname{argmin}_z \quad & \frac{1}{2} z^T Q(z_i) z + q(z_i)^T z \\
 \text{subject to} \quad & A(z_i) z = b(z_i) \\
 & G(z_i) z \leq h(z_i)
 \end{aligned} \tag{1}$$

• Write derivative of KKT conditions

$$\begin{aligned}
 dQz^* + Qdz + dq + dA^T \nu^* + \\
 A^T d\nu + dG^T \lambda^* + G^T d\lambda = 0 \\
 dAz^* + Adz - db = 0
 \end{aligned}$$

$$D(Gz^* - h)d\lambda + D(\lambda^*)(dGz^* + Gdz - dh) = 0$$

\Rightarrow

$$\begin{bmatrix} Q & G^T & A^T \\ D(\lambda^*)G & D(Gz^* - h) & 0 \\ A & 0 & 0 \end{bmatrix} \begin{bmatrix} dz \\ d\lambda \\ d\nu \end{bmatrix} = \begin{bmatrix} dQz^* + dq + dG^T \lambda^* + dA^T \nu^* \\ D(\lambda^*)dGz^* - D(\lambda^*)dh \\ dAz^* - db \end{bmatrix}$$

②

ab

da · b + a · db

③

$$\Rightarrow \boxed{\frac{\partial z_{i+1}}{\partial Q(z_i)}}, \frac{\partial z_{i+1}}{\partial q(z_i)}, \dots, ?$$

$dq \cdot d\lambda \Rightarrow 0$

How to update parameters in the OptNet

$$\begin{bmatrix} Q & G^T & A^T \\ D(\lambda^*)G & D(Gz^* - h) & 0 \\ A & 0 & 0 \end{bmatrix} \begin{bmatrix} dz \\ d\lambda \\ d\nu \end{bmatrix} = - \begin{bmatrix} dQz^* + dq + dG^T \lambda^* + dA^T \nu^* \\ D(\lambda^*)dGz^* - D(\lambda^*)dh \\ dAz^* - db \end{bmatrix}.$$

$$\frac{\partial z_{i+1}}{\partial Q(z_i)}, \frac{\partial z_{i+1}}{\partial q(z_i)}, \dots, ?$$

$l = c^T x$
 $\downarrow \quad \downarrow \quad \searrow$
 标量 $n \times 1$ $n \times 1$
 $dl = \boxed{dc^T} \cdot x + c^T \boxed{dx}$
 $\frac{dl}{dc_i} = \frac{dc}{dc_i} \cdot x + \cancel{c^T \frac{dx}{dc_i}}$
 $= [0 \dots 0 \underset{\uparrow}{1} 0 \dots 0] x$
 $\frac{dl}{dc} = \begin{bmatrix} 1 & & & \\ & \ddots & & \\ & & 1 & \\ & & & \ddots \end{bmatrix} x$

How to update parameters in the OptNet

$$\overset{K}{\begin{bmatrix} Q & G^T & A^T \\ D(\lambda^*)G & D(Gz^* - h) & 0 \\ A & 0 & 0 \end{bmatrix}} \begin{bmatrix} dz \\ d\lambda \\ d\nu \end{bmatrix} = - \begin{bmatrix} dQz^* + dq + dG^T\lambda^* + dA^T\nu^* \\ D(\lambda^*)dGz^* - D(\lambda^*)dh \\ dAz^* - db \end{bmatrix}.$$

I

$dq, dC_i \Rightarrow 0.$

$$\boxed{\frac{\partial z_{i+1}}{\partial Q(z_i)}, \frac{\partial z_{i+1}}{\partial q(z_i)}, \dots, ?}$$

$$\begin{bmatrix} \frac{dz}{dQ} \\ \frac{d\lambda}{dQ} \\ \frac{d\nu}{dQ} \end{bmatrix}$$

$$= -K^{-1} \begin{bmatrix} z^* \otimes I_{n \times n} \\ 0 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} \frac{dz}{dq} \\ \frac{d\lambda}{dq} \\ \frac{d\nu}{dq} \end{bmatrix}$$

$$= -K^{-1} \begin{bmatrix} I \\ 0 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} \frac{dz}{db} \\ \frac{d\lambda}{db} \\ \frac{d\nu}{db} \end{bmatrix}$$

$$= -K^{-1} \begin{bmatrix} 0 \\ 0 \\ -I \end{bmatrix}$$

$$\begin{bmatrix} \frac{dz}{dh} \\ \frac{d\lambda}{dh} \\ \frac{d\nu}{dh} \end{bmatrix}$$

$$= -K^{-1} \begin{bmatrix} 0 \\ -D(\lambda^*) \\ 0 \end{bmatrix}$$

How to update parameters in the OptNet

- To include previous backward pass vector $\frac{\partial l}{\partial z^*} \in \mathbb{R}^n$

- We actually want $\nabla_b l = \frac{\partial l}{\partial z^*} \frac{\partial z^*}{\partial b}, \nabla_Q l, \nabla_A l, \dots$

$$\begin{bmatrix} Q & G^T & A^T \\ D(\lambda^*)G & D(Gz^* - h) & 0 \\ A & 0 & 0 \end{bmatrix} \begin{bmatrix} dz \\ d\lambda \\ d\nu \end{bmatrix} = - \begin{bmatrix} dQz^* + dG + dG^T \lambda^* + dA^T \nu^* \\ D(\lambda^*)dGz^* - D(\lambda^*)dh \\ dAz^* - db \end{bmatrix}.$$

$$\frac{\partial l}{\partial b} = \begin{bmatrix} \frac{\partial l}{\partial z} & 0 & 0 \end{bmatrix} \begin{bmatrix} \frac{\partial z}{\partial b} \\ \frac{\partial \lambda}{\partial b} \\ \frac{\partial \nu}{\partial b} \end{bmatrix} = \begin{bmatrix} \frac{\partial l}{\partial z} & 0 & 0 \end{bmatrix} K^T \begin{bmatrix} 0 \\ 0 \\ I \end{bmatrix}$$

$$\frac{\partial l}{\partial b} = \begin{bmatrix} \frac{\partial l}{\partial z} & 0 & 0 \end{bmatrix} K^T \begin{bmatrix} \frac{\partial z}{\partial b} \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 & 0 & I \end{bmatrix} (K^T)^T \begin{bmatrix} \frac{\partial l}{\partial z} \\ 0 \\ 0 \end{bmatrix}$$

How to update parameters in the OptNet

• To include previous backward pass vector $\frac{\partial l}{\partial z^*} \in \mathbb{R}^n$

• We actually want $\nabla_b l = \frac{\partial l}{\partial z^*} \frac{\partial z^*}{\partial b}$, $\nabla_Q l$, $\nabla_A l$, ...

the differential matrix

$$\begin{bmatrix} Q & G^T & A^T \\ D(\lambda^*)G & D(Gz^* - h) & 0 \\ A & 0 & 0 \end{bmatrix} \begin{bmatrix} dz \\ d\lambda \\ d\nu \end{bmatrix} = - \begin{bmatrix} Q & G^T D(\lambda^*) & A^T \\ G & D(Gz^* - h) & 0 \\ A & 0 & 0 \end{bmatrix}^{-1} \begin{bmatrix} \left(\frac{\partial \ell}{\partial z^*}\right)^T \\ 0 \\ 0 \end{bmatrix} \quad (7)$$

then the relevant gradients with respect to all the QP parameters can be given by

$$\begin{aligned} \nabla_Q \ell &= \frac{1}{2} (d_z z^T + z d_z^T) & \nabla_q \ell &= d_z \\ \nabla_A \ell &= d_\nu z^T + \nu d_z^T & \nabla_b \ell &= -d_\nu \\ \nabla_G \ell &= D(\lambda^*) (d_\lambda z^T + \lambda d_z^T) & \nabla_h \ell &= -D(\lambda^*) d_\lambda \end{aligned} \quad (8)$$

Application using OptNet

- Representation power: OptNet layer can represent any piecewise linear univariate function

ReLU.

$$f(x) = \sum_{i=1}^k w_i \max(a_i x + b_i, 0)$$

两层神经网络

↑

$$z_{i+1} = \underset{z \in \mathbb{R}, t \in \mathbb{R}^k}{\text{minimize}} \quad \|t\|_2^2 + (z - w^T t)^2$$

subject to $a_i x + b_i \leq t_i, \quad i = 1, \dots, k$

$t_i = 0$
 $t_i = a_i x + b_i \Rightarrow \text{ReLU}(a_i x + b_i)$
 $z = \sum_{i=1}^k w_i \max(a_i x + b_i, 0) = f(x).$

Diagram illustrating the representation of a piecewise linear function $f(x)$ as a sum of ReLU functions. The function is shown as a series of linear segments with slopes a_1, a_2, \dots, a_k and intercepts b_1, b_2, \dots, b_k . The weights w_i are applied to each segment to form the final function $f(x)$.

Application using OptNet

- Representation power: converse is False

$$f'(x) = \max\{a_1^T x, a_2^T x, a_3^T x\}$$

ReLU nature

$$\max\{a_i^T x + b_i, 0\}$$

单层网络无法实现

OptNet

$$\begin{aligned} & \underset{z}{\text{minimize}} \quad z^2 \\ & \text{subject to} \quad a_i^T x \leq z, \quad i = 1, \dots, 3. \end{aligned}$$

||

Application using OptNet

- MINIST:

- Learn constraints and dependencies over the output or latent space of a model.
- Show these layers can be included within existing network architectures and efficiently propagate the gradients through the layer.

- Model

- Fully connected layer: FC600-FC10-FC10-SoftMax
- Fully connected layer with OptNet: FC600-FC10-Optnet10-SoftMax

The learnable parameter is L, G, z_0, s_0



$$\begin{aligned} \min_{z \in \mathbb{R}^{10}} & z^T Q z + z_i^T z \\ \text{s.t. } & Q = LL^T + \epsilon I, \\ & Gz \leq Gz_0 + s_0 \\ & Az = Az_0 \end{aligned}$$

$$z \in \mathbb{R}^{10}.$$

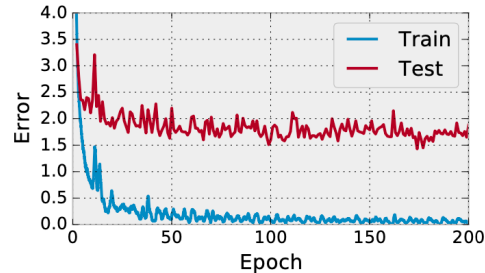
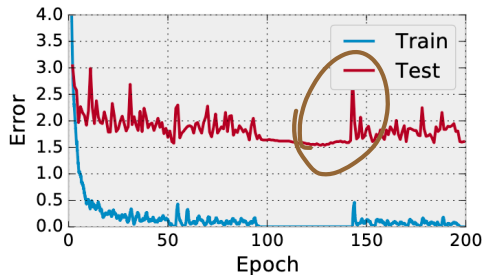
Application using OptNet

- MINIST:

- Learn constraints and dependencies over the output or latent space of a model.
- Show these layers can be included within existing network architectures and efficiently propagate the gradients through the layer.

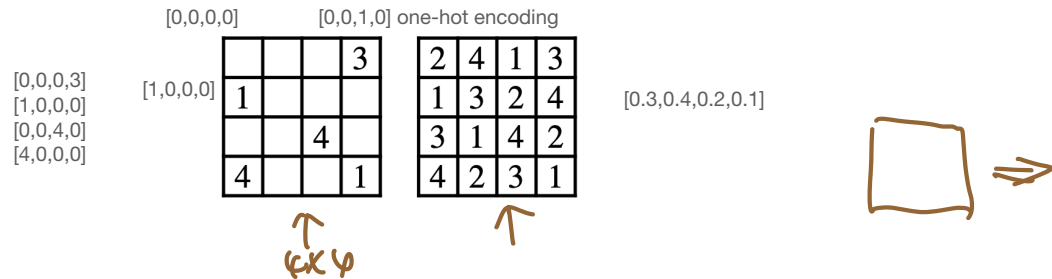
- Result

- slightly lower error and less variance using the OptNet network



Baseline

Application using OptNet



• Sudoku

- a constraint satisfaction problem, nontrivial for computers to know the rules if only given the initial problem and corresponding solved example
- Show OptNet can capture **complex strict relationships** between all input and output variables
- Input: 4×4 tensor, Output: $4 \times 4 \times 4$ tensor (4×4 grid with one-hot encoding)

Handwritten note: 4×4 tensor, one-hot encoding

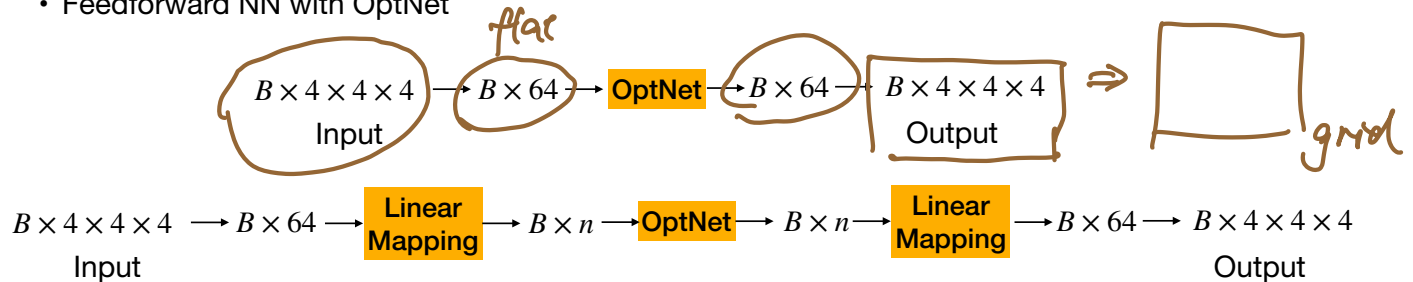
Application using OptNet

- Experiment

- Train on 9000 puzzles, test on 1000 different puzzles
- Error rate: the percentage of puzzles solved incorrectly

- Model

- multilayer feedforward network: 10 convolutional layers with 512 3x3 filters
- Feedforward NN with OptNet



Application using OptNet

- Model

- multilayer feedforward network: 10 convolutional layers with 512 3x3 filters
- Feedforward NN with OptNet

$$B \times 4 \times 4 \times 4 \xrightarrow{\text{Input}} B \times 64 \xrightarrow{\text{OptNet}} B \times 64 \xrightarrow{\text{Output}} B \times 4 \times 4 \times 4$$

$$B \times 4 \times 4 \times 4 \xrightarrow{\text{Input}} B \times 64 \xrightarrow{\text{Linear Mapping}} B \times n \xrightarrow{\text{OptNet}} B \times n \xrightarrow{\text{Linear Mapping}} B \times 64 \xrightarrow{\text{Output}} B \times 4 \times 4 \times 4$$

OptNet

$$z^* = \min_z 0.5 \epsilon z^T z - z_i^T z$$

$$\text{s.t. } \boxed{Az = b}, z \geq 0$$

Application using OptNet

- Result

- the convolutional is able to learn all of the necessary logic for the task and ends up over-fitting to the training data.
- OptNet network learns most of the correct hard constraints within the first three epochs and is able to generalize much better to unseen examples.

