JHU vision lab

Global Optimality in Neural Network Training

Benjamin D. Haeffele and René Vidal

Johns Hopkins University, Center for Imaging Science. Baltimore, USA









Generalization







Are there principled ways to design networks?



- How many layers?
- Size of layers?
- Choice of layer types?
- How does architecture impact expressiveness? [1]



How to train neural networks?





How to train neural networks?

• Problem is non-convex.





How to train neural networks?

• Problem is non-convex.





How to train neural networks?

- Problem is non-convex.
- What does the loss surface look like? [1]
- Any guarantees for network training? [2]
- How to guarantee optimality?
- When will local descent succeed?

[1] Choromanska, et al., "The loss surfaces of multilayer networks." Artificial Intelligence and Statistics. (2015)
[2] Janzamin, et al., "Beating the perils of non-convexity: Guaranteed training of neural networks using tensor methods." arXiv (2015).



Performance Guarantees?



X Complex

- How do networks generalize?
- How should networks be regularized?
- How to prevent overfitting?



Interrelated Problems



- Optimization can impact generalization. [1]
- Architecture has a strong effect on the generalization of networks. [2]
- Some architectures could be easier to optimize than others.





 Are there properties of the network architecture that allow efficient optimization?





- Are there properties of the network architecture that allow efficient optimization?
 - Positive Homogeneity
 - Parallel Subnetwork Structure





- Are there properties of the network architecture that allow efficient optimization?
 - Positive Homogeneity
 - Parallel Subnetwork Structure
- Are there properties of the regularization that allow efficient optimization?





- Are there properties of the network architecture that allow efficient optimization?
 - Positive Homogeneity
 - Parallel Subnetwork Structure
- Are there properties of the regularization that allow efficient optimization?
 - Positive Homogeneity
 - Adapt network architecture to data [1]



Today's Talk: The Results







Today's Talk: The Results





• A local minimum such that one subnetwork is all zero is a global minimum.



Today's Talk: The Results



• Once the size of the network becomes large enough...

Today's Framework

• Local descent can reach a global minimum from any initialization.

Non-Convex Function



Outline



- 1. Network properties that allow efficient optimization
 - Positive Homogeneity
 - Parallel Subnetwork Structure
- 2. Network size from regularization
- 3. Theoretical guarantees
 - Sufficient conditions for global optimality
 - Local descent can reach global minimizers



• Start with a network. Network Outputs W^1 W^2 W^3

Network Weights



• Scale the weights by a non-negative constant.





• Scale the weights by a non-negative constant.





• The network output scales by the constant to some power.





• The network output scales by the constant to some power.



Network Mapping $\Phi(W^1, W^2, W^3) = X$ $\Psi^1 \circ W^2 \circ W^3 = o^p$





• The network output scales by the constant to some power.



Network Mapping $\Phi(W^{1}, W^{2}, W^{3}) = X$ $\Phi(\alpha W^{1}, \alpha W^{2}, \alpha W^{3}) = \alpha^{p} X$

 \boldsymbol{p} - Degree of positive homogeneity






















































• Simple Network





• Simple Network





X

• Simple Network





X

• Simple Network





• Simple Network



• Typically each weight layer increases degree of homogeneity by 1.



Some Common Positively Homogeneous Layers

- ✓ Fully Connected + ReLU
- ✓ Convolution + ReLU
- ✓Max Pooling
- ✓ Linear Layers
- ✓ Mean Pooling
- ✓Max Out
- ✓ Many possibilities...





Some Common Positively Homogeneous Layers





Outline



- 1. Network properties that allow efficient optimization
 - Positive Homogeneity
 - Parallel Subnetwork Structure
- 2. Network regularization
- 3. Theoretical guarantees
 - Sufficient conditions for global optimality
 - Local descent can reach global minimizers



• Subnetworks with identical architecture connected in parallel.





- Subnetworks with identical architecture connected in parallel.
- Simple Example: Single hidden layer network





- Subnetworks with identical architecture connected in parallel.
- Simple Example: Single hidden layer network





- Subnetworks with identical architecture connected in parallel.
- Simple Example: Single hidden layer network



• Subnetwork: One ReLU hidden unit





Any positively homogeneous subnetwork can be used



• Subnetwork: Multiple ReLU layers





• Example: Parallel AlexNets[1]





Outline

Architecture



Generalization/ Regularization





- 1. Network properties that allow efficient optimization
 - Positive Homogeneity
 - Parallel Subnetwork Structure
- 2. Network regularization
- 3. Theoretical guarantees
 - Sufficient conditions for global optimality
 - Local descent can reach global minimizers







 $\Theta(W^1, W^2, W^3) = \|W^1\|_F^2 + \|W^2\|_F^2 + \|W^3\|_F^2$





$$\Theta(W^1, W^2, W^3) = \|W^1\|_F^2 + \|W^2\|_F^2 + \|W^3\|_F^2$$



$$\Theta(\alpha W^1, \alpha W^2, \alpha W^3) = \alpha^2 \Theta(W^1, W^2, W^3)$$



$$\Theta(W^1, W^2, W^3) = \|W^1\|_F^2 + \|W^2\|_F^2 + \|W^3\|_F^2$$



$$\Theta(\alpha W^1, \alpha W^2, \alpha W^3) = \alpha^2 \Theta(W^1, W^2, W^3)$$
$$\Phi(\alpha W^1, \alpha W^2, \alpha W^3) = \alpha^3 \Phi(W^1, W^2, W^3)$$



$$\Theta(W^1, W^2, W^3) = \|W^1\|_F^2 + \|W^2\|_F^2 + \|W^3\|_F^2$$



$$\|W^{-1}\|_{F}^{2} + \|W^{-1}\|_{F}^{2} + \|W^{-1}\|_{F}^{2}$$
$$\Theta(\alpha W^{1}, \alpha W^{2}, \alpha W^{3}) = \alpha^{2} \Theta(W^{1}, W^{2}, W^{3})$$
$$\Phi(\alpha W^{1}, \alpha W^{2}, \alpha W^{3}) = \alpha^{3} \Phi(W^{1}, W^{2}, W^{3})$$



 W^3

$$\Theta(W^1, W^2, W^3) = \|W^1\|_F^2 + \|W^2\|_F^2 + \|W^3\|_F^2$$



$$\Theta(\alpha W^1, \alpha W^2, \alpha W^3) = \alpha^2 \Theta(W^1, W^2, W^3)$$

$$\Phi(\alpha W^1, \alpha W^2, \alpha W^3) = \alpha^3 \Phi(W^1, W^2, W^3)$$

Degrees of positive homogeneity don't match = Bad things happen.

Network Weights

 W^2

 W^1



 W^3

$$\Theta(W^1, W^2, W^3) = \|W^1\|_F^2 + \|W^2\|_F^2 + \|W^3\|_F^2$$



$$\Theta(\alpha W^1, \alpha W^2, \alpha W^3) = \alpha^2 \Theta(W^1, W^2, W^3)$$
$$\Phi(\alpha W^1, \alpha W^2, \alpha W^3) = \alpha^3 \Phi(W^1, W^2, W^3)$$

Degrees of positive homogeneity don't match = Bad things happen.

Network Weights

 W^2

 W^1

Proposition: There will always exist non-optimal local minima.



Start with a positively homogeneous network with parallel structure







• Take the weights of one subnetwork.





$$\theta(W_3^1, W_3^2, W_3^3, W_3^4, W_3^5)$$





$$\theta(W_3^1, W_3^2, W_3^3, W_3^4, W_3^5)$$

- Non-negative.
 - Positively homogeneous with same degree as network mapping.





$$W_{3}^{1} W_{3}^{2} W_{3}^{3} W_{3}^{4} W_{3}^{5}$$

$$\theta(W_3^1, W_3^2, W_3^3, W_3^4, W_3^5)$$

- Non-negative.
- Positively homogeneous with same degree as network mapping.

$$\Phi(\alpha W) = \alpha^p \Phi(W)$$
$$\theta(\alpha W) = \alpha^p \theta(W)$$



$$W_{3}^{1} W_{3}^{2} W_{3}^{3} W_{3}^{4} W_{3}^{5}$$

$$\theta(W_3^1, W_3^2, W_3^3, W_3^4, W_3^5)$$

- Non-negative.
- Positively homogeneous with same degree as network mapping.

$$\Phi(\alpha W) = \frac{\alpha^p}{\alpha^p} \Phi(W)$$
$$\theta(\alpha W) = \frac{\alpha^p}{\alpha^p} \theta(W)$$



• Define a regularization function on the weights.



 $\theta(W_3^1, W_3^2, W_3^3, W_3^4, W_3^5)$

- Non-negative.
- Positively homogeneous with same degree as network mapping.

Example: Product of norms $||W_3^1|||W_3^2|||W_3^3|||W_3^3|||W_3^4|||W_3^5||$



• Sum over all the subnetworks.

θ (Subnetwork 3)





• Sum over all the subnetworks.

θ (Subnetwork 3)+ θ (Subnetwork 2)





• Sum over all the subnetworks.



 θ (Subnetwork 3)+ θ (Subnetwork 2)+ θ (Subnetwork 1)



• Sum over all the subnetworks.



 θ (Subnetwork 3)+ θ (Subnetwork 2)+ θ (Subnetwork 1) $\Theta(W) = \sum \theta(\text{Subnetwork } i)$ i=1

 $r={\rm \#\,of\,Subnetworks}$



• Allow the number of subnetworks to vary.



$$\Theta(W) = \sum_{i=1}^{r} \theta(\text{Subnetwork } i)$$

r = # of Subnetworks

- Adding a subnetwork is penalized by an additional term in the sum.
- Acts to constrain the number of subnetworks.



Outline



- 1. Network properties that allow efficient optimization
 - Positive Homogeneity
 - Parallel Subnetwork Structure

2. Network regularization

- 3. Theoretical guarantees
 - Sufficient conditions for global optimality
 - Local descent can reach global minimizers



Our problem






Our problem

- The non-convex problem we're interested in $f(W) \equiv \ell(Y, \Phi(W)) + \lambda \Theta(W)$





Our problem

• The non-convex problem we're interested in $f(W) \equiv \ell(Y, \Phi(W)) + \lambda \Theta(W)$ $\Theta(W) = \sum_{i=1}^{r} \theta(\text{Subnetwork } i)$



i=1

Our problem

• The non-convex problem we're interested in $f(W) \equiv \ell(Y, \Phi(W)) + \lambda \Theta(W)$ Labels $\Theta(W) = \sum_{i=1}^{r} \theta(\text{Subnetwork } i)$

Loss Function: $\ell(Y, X)$

Assume convex and once differentiable in \boldsymbol{X}

Examples:

- Cross-entropy
- Least-squares











• Induces a convex function on the network outputs.

$$F(X) \equiv \ell(Y, X) + \lambda \Omega(X)$$





• Induces a convex function on the network outputs.

$$F(X) \equiv \ell(Y, X) + \lambda \Omega(X)$$

Induced Function: $\Omega(X)$

Comes from the regularization



• Induces a convex function on the network outputs.

$$F(X) \equiv \ell(Y, X) + \lambda \Omega(X)$$

Induced Function: $\Omega(X)$

Comes from the regularization

$$\Omega(X) = \min_{W} \Theta(W)$$
 subject to $X = \Phi(W)$



• Induces a convex function on the network outputs.

$$F(X) \equiv \ell(Y, X) + \lambda \Omega(X)$$

Induced Function: $\Omega(X)$

Comes from the regularization

$$\Omega(X) = \min_{W} \Theta(W) \text{ subject to } X = \Phi(W)$$
$$\Omega(X) \le \Theta(W) \qquad \forall X = \Phi(W)$$



• Induces a convex function on the network outputs.

$$F(X) \equiv \ell(Y, X) + \lambda \Omega(X)$$

Induced Function: $\Omega(X)$

Comes from the regularization

$$\Omega(X) = \min_{W} \Theta(W) \text{ subject to } X = \Phi(W)$$
$$\Omega(X) \le \Theta(W) \quad \forall X = \Phi(W)$$

Proposition: $\Omega(X)$ is a convex function.



• Induces a convex function on the network outputs.

$$F(X) \equiv \ell(Y, X) + \lambda \Omega(X)$$

• The convex problem provides an achievable lower bound for the non-convex network training problem.

$$f(W) \equiv \ell(Y, \Phi(W)) + \lambda \Theta(W)$$

$$F(X) \le f(W) \qquad \forall X = \Phi(W)$$



• Induces a convex function on the network outputs.

$$F(X) \equiv \ell(Y, X) + \lambda \Omega(X)$$

• The convex problem provides an achievable lower bound for the non-convex network training problem.

$$f(W) \equiv \ell(Y, \Phi(W)) + \lambda \Theta(W)$$

$$F(X) \le f(W) \qquad \forall X = \Phi(W)$$

• Use the convex function as an analysis tool to study the non-convex network training problem.



Sufficient Conditions for Global Optimality



• **Theorem:** A local minimum such that one subnetwork is all zero is a global minimum.



Sufficient Conditions for Global Optimality



• Theorem: A local minimum such that one subnetwork is all zero is a global minimum.



Sufficient Conditions for Global Optimality



- Theorem: A local minimum such that one subnetwork is all zero is a global minimum.
- Intuition: The local minimum satisfies the optimality conditions for the convex problem.

 $F(X) \le f(W)$



Global Minima from Local Descent

• **Theorem:** If the size of the network is large enough (has enough subnetworks), then a global minimum can always be reached by local descent from any initialization.





Global Minima from Local Descent

• **Theorem:** If the size of the network is large enough (has enough subnetworks), then a global minimum can always be reached by local descent from any initialization.





Global Minima from Local Descent

• **Theorem:** If the size of the network is large enough (has enough subnetworks), then a global minimum can always be reached by local descent from any initialization.

• Meta-Algorithm:

• If not at a local minima, perform local descent

Non-Convex Function

- At local minima, test if first Theorem is satisfied
- If not, add a subnetwork in parallel and continue
- Maximum number of subnetworks guaranteed to be bounded by the dimensions of the network output



Today's Framework

Conclusions

- Network size matters
 - Optimize network weights AND network size
 - Current: Size = Number of parallel subnetworks
 - Future: Size = Number of layers, neurons per layer, etc...
- Regularization design matters
 - Match the degrees of positive homogeneity between network and regularization
 - Regularization can control the size of the network
- Not done yet
 - Several practical and theoretical limitations



Vision Lab @ Johns Hopkins University http://www.vision.jhu.edu

Center for Imaging Science @ Johns Hopkins University http://www.cis.jhu.edu

Work supported by NSF grants 1447822, 1618485 and 1618637



