International Conference on Machine Learning (ICML) 2018

On the Optimization of Deep Networks: Implicit Acceleration by Overparameterization

Sanjeev Arora^{*†} Nadav Cohen[†] Elad Hazan^{*†}

*Princeton University



[†]Institute for Advanced Study



[‡]Google Brain



Conventional wisdom:

Depth boosts expressive power



Conventional wisdom:

Depth boosts expressive power



But complicates optimization



Conventional wisdom:

Depth boosts expressive power



But complicates optimization



This work:

Depth can accelerate optimization!

Problem:

Expressiveness can interfere with our study – deeper nets may seem to optimize faster per being able to reach lower training err

Problem:

Expressiveness can interfere with our study – deeper nets may seem to optimize faster per being able to reach lower training err

Resolution:

 We focus on models whose expressiveness is oblivious to depth – linear neural networks

Problem:

Expressiveness can interfere with our study – deeper nets may seem to optimize faster per being able to reach lower training err

Resolution:

- We focus on models whose expressiveness is oblivious to depth – linear neural networks
- Adding layers amounts to replacing matrix param by product of matrices – overparameterization

Problem:

Expressiveness can interfere with our study – deeper nets may seem to optimize faster per being able to reach lower training err

Resolution:

- We focus on models whose expressiveness is oblivious to depth – linear neural networks
- Adding layers amounts to replacing matrix param by product of matrices – overparameterization

Linear neural networks were studied extensively, cf. [Saxe et al. 2013; Kawaguchi 2016; Hardt & Ma 2016]

Training objective for linear regression with ℓ_p loss:

$$L(\boldsymbol{w}) = \sum_{(\boldsymbol{x},\boldsymbol{y})\in S} \frac{1}{p} (\boldsymbol{x}^T \boldsymbol{w} - \boldsymbol{y})^p$$

Training objective for linear regression with ℓ_p loss:

$$L(\boldsymbol{w}) = \sum_{(\boldsymbol{x},\boldsymbol{y})\in S} \frac{1}{p} (\boldsymbol{x}^T \boldsymbol{w} - \boldsymbol{y})^p \quad \text{convex}$$

Training objective for linear regression with ℓ_p loss:

$$L(\boldsymbol{w}) = \sum_{(\boldsymbol{x},\boldsymbol{y})\in S} \frac{1}{p} (\boldsymbol{x}^T \boldsymbol{w} - \boldsymbol{y})^p \quad \text{convex}$$

Overparameterize – replace vector w by vector w_1 times scalar w_2 :

$$\tilde{L}(\boldsymbol{w_1}, \boldsymbol{w_2}) = \sum_{(\boldsymbol{x}, \boldsymbol{y}) \in S} \frac{1}{p} (\boldsymbol{x}^T \boldsymbol{w_1} \boldsymbol{w_2} - \boldsymbol{y})^p$$

Training objective for linear regression with ℓ_p loss:

$$L(\boldsymbol{w}) = \sum_{(\boldsymbol{x},\boldsymbol{y})\in S} \frac{1}{p} (\boldsymbol{x}^T \boldsymbol{w} - \boldsymbol{y})^p \quad \text{convex}$$

Overparameterize – replace vector w by vector w_1 times scalar w_2 :

$$\tilde{L}(\boldsymbol{w_1}, \boldsymbol{w_2}) = \sum_{(\boldsymbol{x}, \boldsymbol{y}) \in S} \frac{1}{p} (\boldsymbol{x}^T \boldsymbol{w_1} \boldsymbol{w_2} - \boldsymbol{y})^p \quad \text{non-convex}$$

Training objective for linear regression with ℓ_p loss:

$$L(\boldsymbol{w}) = \sum_{(\boldsymbol{x},\boldsymbol{y})\in S} \frac{1}{p} (\boldsymbol{x}^T \boldsymbol{w} - \boldsymbol{y})^p \quad \text{convex}$$

Overparameterize – replace vector w by vector w_1 times scalar w_2 :

$$\tilde{L}(\boldsymbol{w_1}, \boldsymbol{w_2}) = \sum_{(\boldsymbol{x}, \boldsymbol{y}) \in S} \frac{1}{p} (\boldsymbol{x}^T \boldsymbol{w_1} \boldsymbol{w_2} - \boldsymbol{y})^p \quad \text{non-convex}$$

Claim:

• Gradient descent (GD) over $\tilde{L}(w_1, w_2)$ induces on $w = w_1 w_2$:

$$\boldsymbol{w}^{(t+1)} \leftarrow \boldsymbol{w}^{(t)} - \rho^{(t)} \cdot \nabla L(\boldsymbol{w}^{(t)}) - \sum_{\tau=1}^{t-1} \mu^{(t,\tau)} \cdot \nabla L(\boldsymbol{w}^{(\tau)})$$

Training objective for linear regression with ℓ_p loss:

$$L(\boldsymbol{w}) = \sum_{(\boldsymbol{x},\boldsymbol{y})\in S} \frac{1}{p} (\boldsymbol{x}^T \boldsymbol{w} - \boldsymbol{y})^p \quad \text{convex}$$

Overparameterize – replace vector w by vector w_1 times scalar w_2 :

$$\tilde{L}(\boldsymbol{w_1}, \boldsymbol{w_2}) = \sum_{(\boldsymbol{x}, \boldsymbol{y}) \in S} \frac{1}{p} (\boldsymbol{x}^T \boldsymbol{w_1} \boldsymbol{w_2} - \boldsymbol{y})^p \quad \text{non-convex}$$

Claim:

• Gradient descent (GD) over $\tilde{L}(w_1, w_2)$ induces on $w = w_1 w_2$:

$$\boldsymbol{w}^{(t+1)} \leftarrow \boldsymbol{w}^{(t)} - \rho^{(t)} \cdot \nabla L(\boldsymbol{w}^{(t)}) - \sum_{\tau=1}^{t-1} \mu^{(t,\tau)} \cdot \nabla L(\boldsymbol{w}^{(\tau)})$$

adaptive learning rate

Training objective for linear regression with ℓ_p loss:

$$L(\boldsymbol{w}) = \sum_{(\boldsymbol{x},\boldsymbol{y})\in S} \frac{1}{p} (\boldsymbol{x}^T \boldsymbol{w} - \boldsymbol{y})^p \quad \text{convex}$$

Overparameterize – replace vector w by vector w_1 times scalar w_2 :

$$\tilde{L}(\boldsymbol{w_1}, \boldsymbol{w_2}) = \sum_{(\boldsymbol{x}, \boldsymbol{y}) \in S} \frac{1}{p} (\boldsymbol{x}^T \boldsymbol{w_1} \boldsymbol{w_2} - \boldsymbol{y})^p \quad \text{non-convex}$$

Claim:

adap

• Gradient descent (GD) over $\tilde{L}(w_1, w_2)$ induces on $w = w_1 w_2$:

$$\boldsymbol{w}^{(t+1)} \leftarrow \boldsymbol{w}^{(t)} - \rho^{(t)} \cdot \nabla L(\boldsymbol{w}^{(t)}) - \sum_{\tau=1}^{t-1} \mu^{(t,\tau)} \cdot \nabla L(\boldsymbol{w}^{(\tau)})$$

tive learning rate *momentum*

Training objective for linear regression with ℓ_p loss:

$$L(\boldsymbol{w}) = \sum_{(\boldsymbol{x},\boldsymbol{y})\in S} \frac{1}{p} (\boldsymbol{x}^T \boldsymbol{w} - \boldsymbol{y})^p \quad \text{convex}$$

Overparameterize – replace vector w by vector w_1 times scalar w_2 :

$$\tilde{L}(\boldsymbol{w_1}, \boldsymbol{w_2}) = \sum_{(\boldsymbol{x}, \boldsymbol{y}) \in S} \frac{1}{p} (\boldsymbol{x}^T \boldsymbol{w_1} \boldsymbol{w_2} - \boldsymbol{y})^p \quad \text{non-convex}$$

Claim:

• Gradient descent (GD) over $\tilde{L}(w_1, w_2)$ induces on $w = w_1 w_2$:

$$\boldsymbol{w}^{(t+1)} \leftarrow \boldsymbol{w}^{(t)} - \rho^{(t)} \cdot \nabla L(\boldsymbol{w}^{(t)}) - \sum_{\tau=1}^{t-1} \mu^{(t,\tau)} \cdot \nabla L(\boldsymbol{w}^{(\tau)})$$

adaptive learning rate \checkmark "momentum"

• For p > 2, this can speed up optimization

Training objective for linear regression with ℓ_p loss:

$$L(\boldsymbol{w}) = \sum_{(\boldsymbol{x},\boldsymbol{y})\in S} \frac{1}{p} (\boldsymbol{x}^T \boldsymbol{w} - \boldsymbol{y})^p \quad \textbf{convex}$$

Overparameterize – replace vector w by vector w_1 times scalar w_2 :

$$\tilde{L}(\boldsymbol{w_1}, \boldsymbol{w_2}) = \sum_{(\boldsymbol{x}, \boldsymbol{y}) \in S} \frac{1}{p} (\boldsymbol{x}^T \boldsymbol{w_1} \boldsymbol{w_2} - \boldsymbol{y})^p \quad \text{non-convex}$$

Claim:

• Gradient descent (GD) over $\tilde{L}(w_1, w_2)$ induces on $w = w_1 w_2$:

$$\boldsymbol{w}^{(t+1)} \leftarrow \boldsymbol{w}^{(t)} - \rho^{(t)} \cdot \nabla L(\boldsymbol{w}^{(t)}) - \sum_{\tau=1}^{t-1} \mu^{(t,\tau)} \cdot \nabla L(\boldsymbol{w}^{(\tau)})$$

adaptive learning rate \checkmark "momentum"

• For p > 2, this can speed up optimization (cf. [Saxe et al. 2013])

Linear neural network:

$$x \rightarrow W_1 \rightarrow W_2 \rightarrow \cdots \rightarrow W_N \rightarrow y = W_N \cdots W_2 \cdot W_1 x$$

Linear neural network:



Linear neural network:

$$x \rightarrow W_{1} \rightarrow W_{2} \rightarrow \cdots \rightarrow W_{N} \rightarrow y = W_{N} \cdots W_{2} \cdot W_{1} x$$

$$W_{e}$$

end-to-end weight matrix

Given loss *L* over linear model, we have the **overparameterized loss**: $\tilde{L}(W_1, \dots, W_N) = L(W_e)$

Linear neural network:

$$x \rightarrow W_{1} \rightarrow W_{2} \rightarrow \cdots \rightarrow W_{N} \rightarrow y = W_{N} \cdots W_{2} \cdot W_{1} x$$

$$W_{e}$$

end-to-end weight matrix

Given loss *L* over linear model, we have the **overparameterized loss**:

$$\tilde{L}(W_1, \cdots, W_N) = L(W_e)$$

Question:

How does W_e behave during GD over W_1, \dots, W_N ?

Theorem:

GD over W_1, \dots, W_N with small learning rate and near-zero init, induces on W_e the **end-to-end update rule**:

$$W_e^{(t+1)} \longleftrightarrow W_e^{(t)}$$

$$- \eta \sum_{j=1}^{N} \left[W_e^{(t)} \left(W_e^{(t)} \right)^T \right]^{\frac{j-1}{N}} \nabla L \left(W_e^{(t)} \right) \left[\left(W_e^{(t)} \right)^T W_e^{(t)} \right]^{\frac{N-j}{N}}$$

Theorem:

GD over W_1, \dots, W_N with small learning rate and near-zero init, induces on W_e the end-to-end update rule:

$$W_e^{(t+1)} \longleftrightarrow W_e^{(t)} - \eta \sum_{j=1}^{N} \left[W_e^{(t)} \left(W_e^{(t)} \right)^T \right]^{\frac{j-1}{N}} \nabla L \left(W_e^{(t)} \right) \left[\left(W_e^{(t)} \right)^T W_e^{(t)} \right]^{\frac{N-j}{N}}$$

which is a:

preconditioner promoting movement in directions already taken

Theorem:

GD over W_1, \dots, W_N with small learning rate and near-zero init, induces on W_e the end-to-end update rule:

$$W_e^{(t+1)} \longleftrightarrow W_e^{(t)} - \eta \sum_{j=1}^{N} \left[W_e^{(t)} \left(W_e^{(t)} \right)^T \right]^{\frac{j-1}{N}} \nabla L \left(W_e^{(t)} \right) \left[\left(W_e^{(t)} \right)^T W_e^{(t)} \right]^{\frac{N-j}{N}}$$

which is a:

- preconditioner promoting movement in directions already taken
- certain combination of adaptive learning rate and "momentum"

Theorem:

GD over W_1, \dots, W_N with small learning rate and near-zero init, induces on W_e the **end-to-end update rule**:

$$W_e^{(t+1)} \longleftrightarrow W_e^{(t)} - \eta \sum_{j=1}^N \left[W_e^{(t)} \left(W_e^{(t)} \right)^T \right]^{\frac{j-1}{N}} \nabla L \left(W_e^{(t)} \right) \left[\left(W_e^{(t)} \right)^T W_e^{(t)} \right]^{\frac{N-j}{N}}$$

which is a:

- preconditioner promoting movement in directions already taken
- certain combination of adaptive learning rate and "momentum"

Overparameterization with deep linear net induces on GD a certain acceleration scheme! No dependence on layer widths!

Theorem:

There exists no objective func (of W_e) over which GD gives the end-to-end update rule

Theorem:

There exists no objective func (of W_e) over which GD gives the end-to-end update rule

Proof sketch:

Fundamental theorem for line integrals:

$$\oint_{\Gamma} \nabla g = 0 \quad \forall \text{ func } g \text{ , closed curve } \Gamma$$

Construct curve on which line integral of end-to-end updates $\neq 0$

Theorem:

There exists no objective func (of W_e) over which GD gives the end-to-end update rule

Proof sketch:

Fundamental theorem for line integrals:

$$\oint_{\Gamma} \nabla g = 0 \quad \forall \text{ func } g \text{ , closed curve } \Gamma$$

Construct curve on which line integral of end-to-end updates $\neq 0$

Effect of overparameterization by depth cannot be attained via any modification of the objective!

Regression problem from UCI ML Repository; ℓ_4 loss

Regression problem from UCI ML Repository; ℓ_4 loss



The Effect of Depth

Regression problem from UCI ML Repository; ℓ_4 loss



The Effect of Depth

Depth can speed up GD, even w/o any change in expressiveness, and despite introducing non-convexity!

Regression problem from UCI ML Repository; ℓ_4 loss



Depth can speed up GD, even w/o any change in expressiveness, and despite introducing non-convexity!

The Effect of Depth

Depth vs. Explicit Accelerators

1000000

Regression problem from UCI ML Repository; ℓ_4 loss



Depth can speed up GD, even w/o any change in expressiveness, and despite introducing non-convexity!

This speed up can outperform explicit acceleration methods designed for convex problems!

800000

1000000

Depth vs. Explicit Accelerators

TensorFlow CNN tutorial for MNIST

Overparameterization: fully-connected layers \rightarrow depth-2 linear nets

TensorFlow CNN tutorial for MNIST

Overparameterization: fully-connected layers \rightarrow depth-2 linear nets



With +15% in params, and no change in expressiveness, overparameterization accelerated by orders of magnitude!

• Depth radically changes obj landscape, turns it highly non-convex

- Depth radically changes obj landscape, turns it highly non-convex
- Conventional wisdom: this complicates optimization

- Depth radically changes obj landscape, turns it highly non-convex
- Conventional wisdom: this complicates optimization
- We show:
 - For linear nets, depth induces on GD a preconditioning scheme

- Depth radically changes obj landscape, turns it highly non-convex
- Conventional wisdom: this complicates optimization
- We show:
 - For linear nets, depth induces on GD a preconditioning scheme
 - Preconditioning combines adaptive learning rate and "momentum"

- Depth radically changes obj landscape, turns it highly non-convex
- Conventional wisdom: this complicates optimization
- We show:
 - For linear nets, depth induces on GD a preconditioning scheme
 - Preconditioning combines adaptive learning rate and "momentum"
 - Effect cannot be attained via any modification of original objective

- Depth radically changes obj landscape, turns it highly non-convex
- Conventional wisdom: this complicates optimization
- We show:
 - For linear nets, depth induces on GD a preconditioning scheme
 - Preconditioning combines adaptive learning rate and "momentum"
 - Effect cannot be attained via any modification of original objective
 - Can lead to significant speed ups, despite no change in expressiveness

Thank You