Lecture 5: Training Deep Neural Networks II

Madhavan Mukund

https://www.cmi.ac.in/~madhavan

Advanced Machine Learning September–December 2021

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへで

Ill conditioning

 Ill conditioning — small change in input produces a large change in output



Ill conditioning

 Ill conditioning — small change in input produces a large change in output

• Gradient
$$\nabla_{\theta} = \frac{\partial}{\partial \theta_i} J(\theta)$$



- Ill conditioning small change in input produces a large change in output
- Gradient $\nabla_{\theta} = \frac{\partial}{\partial \theta_i} J(\theta)$
- Impact of update $\theta \epsilon \nabla_{\theta}$ on cost $J(\theta)$?



- Ill conditioning small change in input produces a large change in output
- Gradient $\nabla_{\theta} = \frac{\partial}{\partial \theta_i} J(\theta)$
- Impact of update $\theta \epsilon \nabla_{\theta}$ on cost $J(\theta)$?
- Depends on curvature, given by second derivative



- Ill conditioning small change in input produces a large change in output
- Gradient $\nabla_{\theta} = \frac{\partial}{\partial \theta_i} J(\theta)$
- Impact of update $\theta \epsilon \nabla_{\theta}$ on cost $J(\theta)$?
- Depends on curvature, given by second derivative
- Hessian: $H_{\theta} = \frac{\delta^2}{\delta \theta_i \delta \theta_j} J(\theta)$



- Ill conditioning small change in input produces a large change in output
- Gradient $\nabla_{\theta} = \frac{\partial}{\partial \theta_i} J(\theta)$
- Impact of update $\theta \epsilon \nabla_{\theta}$ on cost $J(\theta)$?
- Depends on curvature, given by second derivative
- Hessian: $H_{\theta} = \frac{\delta^2}{\delta \theta_i \delta \theta_j} J(\theta)$
- Using Taylor expansion, impact of update $\theta \epsilon \nabla_{\theta}$,

 $J(heta) -
abla_{ heta}^{ op}
abla_{ heta} + rac{1}{2}
abla_{ heta}^{ op} H_{ heta}
abla_{ heta}$



- Ill conditioning small change in input produces a large change in output
- Gradient $\nabla_{\theta} = \frac{\partial}{\partial \theta_i} J(\theta)$
- Impact of update $\theta \epsilon \nabla_{\theta}$ on cost $J(\theta)$?
- Depends on curvature, given by second derivative
- Hessian: $H_{\theta} = \frac{\delta^2}{\delta \theta_i \delta \theta_j} J(\theta)$
- Using Taylor expansion, impact of update $\theta \epsilon \nabla_{\theta}$,

 $J(heta) -
abla_{ heta}^{ op}
abla_{ heta} + rac{1}{2}
abla_{ heta}^{ op} H_{ heta}
abla_{ heta}$

• Analyze H_{θ} to check for ill conditioning



- Locally steepest direction of descent may be far from the optimum
 - Elliptical contours vs circular contours



- Locally steepest direction of descent may be far from the optimum
 - Elliptical contours vs circular contours
- Gradient changes rapidly along the direction of steepest descent
 - Taking large steps is problematic



- Locally steepest direction of descent may be far from the optimum
 - Elliptical contours vs circular contours
- Gradient changes rapidly along the direction of steepest descent
 - Taking large steps is problematic
- Ill-conditioned Hessian *H* second derivatives
 - Computing Hessian is expensive
 - "Second order" methods are not used in practice



- Locally steepest direction of descent may be far from the optimum
 - Elliptical contours vs circular contours
- Gradient changes rapidly along the direction of steepest descent
 - Taking large steps is problematic
- Ill-conditioned Hessian *H* second derivatives
 - Computing Hessian is expensive
 - "Second order" methods are not used in practice
- Instead, heuristics like momentum and adaptive learning rates



Madhavan Mukund

SGD convergence can be very slow

э

▶ ∢ ⊒

- SGD convergence can be very slow
- Momentum in physics mass × velocity

- 4 西

э

- SGD convergence can be very slow
- Momentum in physics mass × velocity
- Introduce velocity v in SGD assume unit mass
 - Moving average of past gradients, exponential decay
 - If gradient remains steady, velocity increases

- SGD convergence can be very slow
- Momentum in physics mass × velocity
- Introduce velocity v in SGD assume unit mass
 - Moving average of past gradients, exponential decay
 - If gradient remains steady, velocity increases

Update rule • $\mathbf{v} \leftarrow \alpha \mathbf{v} - \epsilon \nabla_{\theta} \left(\frac{1}{m} \sum_{i=1}^{m} \mathcal{O}(\mathbf{x}_{i}; \theta), y_{i}) \right)$ $\blacksquare \theta \leftarrow \theta + v$



 $\theta \geq \theta - \varepsilon \nabla_{\theta} J$

э

4 E

- SGD convergence can be very slow
- Momentum in physics mass × velocity
- Introduce velocity v in SGD assume unit mass
 - Moving average of past gradients, exponential decay
 - If gradient remains steady, velocity increases

x = 0.9

Update rule

$$v \leftarrow \alpha v - \epsilon \nabla_{\theta} \left(\frac{1}{m} \sum_{i=1}^{m} L(f(x_i; \theta), y_i) \right)$$

$$\theta \leftarrow \theta + v$$

- Hyperparameter $\alpha \in [0, 1)$ "friction", exponentially decaying history
 - With constant gradient g, in the limit $\frac{\epsilon g}{1-\alpha}$, geometric progression

Nesterov momentum optimization

 Measure cost function slightly ahead, in direction of momentum

Nesterov momentum optimization

- Measure cost function slightly ahead, in direction of momentum
- Update rule

•
$$v \leftarrow \alpha v - \epsilon \nabla_{\theta} \left(\frac{1}{m} \sum_{i=1}^{m} L(f(x_i; \theta + \beta m), y_i) \right)$$

• $\theta \leftarrow \theta + v$

Nesterov momentum optimization

- Measure cost function slightly ahead, in direction of momentum
- Update rule

•
$$\mathbf{v} \leftarrow \alpha \mathbf{v} - \epsilon \nabla_{\theta} \left(\frac{1}{m} \sum_{i=1}^{m} L(f(\mathbf{x}_i; \theta + \beta \mathbf{m}), \mathbf{y}_i) \right)$$

• $\theta \leftarrow \theta + \mathbf{v}$

 Controls sideways oscillations better



Adjusting the trajectory

If features have different scales, gradient descent is steeper in some dimensions

How can we correct for this?



Adagrad

- Adagrad update rule
 - $g \leftarrow \frac{1}{m} \nabla_{\theta} \sum_{i=1}^{m} L(f(x_i; \theta), y_i)$ • $r \leftarrow r + g \cdot g$ • $\Delta \theta \leftarrow \frac{\epsilon}{\delta + \sqrt{r}} \cdot g$, where $\delta \approx 10^{-7}$, for numerical stability
 - $\blacksquare \ \theta \leftarrow \theta + \Delta \theta$

< 一型

Adagrad

Adagrad update rule



- $\blacksquare \ \theta \leftarrow \theta + \Delta \theta$
- Shrink learning rate in each dimension according to entire history of squared gradient

-

э

Adagrad

Adagrad update rule



- $\mathbf{\theta} \leftarrow \boldsymbol{\theta} + \Delta \boldsymbol{\theta}$
- Shrink learning rate in each dimention according to entire history of squared gradient



э

▶ < ∃ ▶</p>

Adaptive learning rates

RMSProp

- Using entire history shrinks learning rate too much
- Exponentially decaying average, discard extreme past
- Update rule

•
$$g \leftarrow \frac{1}{m} \nabla_{\theta} \sum_{i=1}^{m} L(f(x_i; \theta), y_i)$$

• $r \leftarrow \rho r + (1 - \rho)(g \cdot g)$, where ρ is decay rate
• $\Delta \theta$
• $\frac{\epsilon}{\sqrt{\delta + r}}$, where $\delta \approx 10^{-6}$
• $\theta \leftarrow \theta + \Delta \theta$

• New hyperparameter ρ

э

Adaptive moments — combines RMSProp and moments

э

▶ ∢ ⊒

< □ > < 円

Adam

- Adaptive moments combines RMSProp and moments
- Update rule
 - Step size ϵ ; two decay rates ρ_1 , ρ_2 ; two moments, s = r = 0; time step t = 0





Adam

- Adaptive moments combines RMSProp and moments
- Update rule
 - Step size ϵ ; two decay rates ρ_1 , ρ_2 ; two moments, s = r = 0; time step t = 0
 - $g \leftarrow \frac{1}{m} \nabla_{\theta} \sum_{i=1}^{m} L(f(x_i; \theta), y_i)$
 - $s \leftarrow \rho_1 s + (1 \rho_1)g; r \leftarrow \rho_2 r + (1 \rho_2)(g \cdot g)$
 - Correct bias in first and second moments: $\hat{s} \leftarrow \frac{s}{1-\rho_1^t}$, $\hat{r} \leftarrow \frac{r}{1-\rho_2^t}$

$$\Delta \theta = -\epsilon \frac{\hat{s}}{\sqrt{\hat{r}} + \delta}; \ \theta \leftarrow \theta + \Delta \theta$$

- No clear theoretical justification for combinining momentum and scaling
- Fairly robust with respect to values of hyperparameters

Adaptive learning rates

- Choosing a fixed learning rate is hard
- Make a learning rate a function of iteration number
- Power scheduling, exponential scheduling, piecewise constant scheduling



10/15

Regularization

- ℓ_1 and ℓ_2 regularization, as usual
- Dropout
 - Disable nodes with probability
 - Analogy multifunctional employees
 - Scale weights after training

Origend Lost + Parameter Lest Replanzation



э

< ∃

- The loss function for regression is convex
- Gradient descent converges to global optimum



э

< E

- The loss function for regression is convex
- Gradient descent converges to global optimum
- Loss function for neural networks is not convex
- In general, gradient descent only finds local minima



- The loss function for regression is convex
- Gradient descent converges to global optimum
- Loss function for neural networks is not convex
- In general, gradient descent only finds local minima
- How many local minima are there?



- The loss function for regression is convex
- Gradient descent converges to global optimum
- Loss function for neural networks is not convex
- In general, gradient descent only finds local minima
- How many local minima are there?
- How does it affect gradient descent?



Is the model that fits the data unique?

э

- Is the model that fits the data unique?
- Non-identifiable two or more settings of the parameters are observationally equivalent

э

Is the model that fits the data unique?

- Non-identifiable two or more settings of the parameters are observationally equivalent
- Symmetry
 - Fully connected network, permutations of a layer are indistinguishable



Is the model that fits the data unique?

- Non-identifiable two or more settings of the parameters are observationally equivalent
- Symmetry
 - Fully connected network, permutations of a layer are indistinguishable
- Piecewise linear activation ReLU
 - Scale inputs by k, multiply output by 1/k



Is the model that fits the data unique?

- Non-identifiable two or more settings of the parameters are observationally equivalent
- Symmetry
 - Fully connected network, permutations of a layer are indistinguishable
- Piecewise linear activation ReLU
 - Scale inputs by k, multiply output by 1/k
- Large numbers of local minima!



How to measure the impact of local minima?

э

- How to measure the impact of local minima?
- Training process will see local ups and downs, but "bumpy" surface may not give a good picture

- How to measure the impact of local minima?
- Training process will see local ups and downs, but "bumpy" surface may not give a good picture
- Instead [Goodfellow et al]
 - **Random initialization** θ_i
 - **SGD** finds an optimum value θ_f

- How to measure the impact of local minima?
- Training process will see local ups and downs, but "bumpy" surface may not give a good picture
- Instead [Goodfellow et al]
 - **Random initialization** θ_i
 - **SGD** finds an optimum value θ_f
 - Check loss along the linearly interpolation $\theta_{\alpha} = \alpha \cdot \theta_f + (1 - \alpha) \cdot \theta_i$



- How to measure the impact of local minima?
- Training process will see local ups and downs, but "bumpy" surface may not give a good picture
- Instead [Goodfellow et al]
 - **Random initialization** θ_i
 - **SGD** finds an optimum value θ_f
 - Check loss along the linearly interpolation $\theta_{\alpha} = \alpha \cdot \theta_f + (1 - \alpha) \cdot \theta_i$
 - Are there problematic local minima along the path?



- How to measure the impact of local minima?
- Training process will see local ups and downs, but "bumpy" surface may not give a good picture
- Instead [Goodfellow et al]
 - **Random initialization** θ_i
 - **SGD** finds an optimum value θ_f
 - Check loss along the linearly interpolation $\theta_{\alpha} = \alpha \cdot \theta_f + (1 - \alpha) \cdot \theta_i$
 - Are there problematic local minima along the path?



Typically, no!



- Critical points zero gradient
 - Minimum, maximum or inflection point
 - k critical points $\rightarrow k/3$ are minima



- Critical points zero gradient
 - Minimum, maximum or inflection point
 - k critical points $\rightarrow k/3$ are minima
- In d dimensions
 - Should be minimum in *d* directions
 - k critical points $\rightarrow k/3^d$ are minima



- Critical points zero gradient
 - Minimum, maximum or inflection point
 - k critical points $\rightarrow k/3$ are minima
- In d dimensions
 - Should be minimum in *d* directions
 - k critical points $\rightarrow k/3^d$ are minima
- Large fraction of critical values are saddle points



- Critical points zero gradient
 - Minimum, maximum or inflection point
 - k critical points $\rightarrow k/3$ are minima
- In d dimensions
 - Should be minimum in *d* directions
 - k critical points $\rightarrow k/3^d$ are minima
- Large fraction of critical values are saddle points
- Does not seem to be a problem for SGD



- Critical points zero gradient
 - Minimum, maximum or inflection point
 - k critical points $\rightarrow k/3$ are minima
- In d dimensions
 - Should be minimum in *d* directions
 - k critical points $\rightarrow k/3^d$ are minima
- Large fraction of critical values are saddle points
- Does not seem to be a problem for SGD
- Solving directly for zero gradient is problematic



- Critical points zero gradient
 - Minimum, maximum or inflection point
 - k critical points $\rightarrow k/3$ are minima
- In d dimensions
 - Should be minimum in *d* directions
 - k critical points $\rightarrow k/3^d$ are minima
- Large fraction of critical values are saddle points
- Does not seem to be a problem for SGD
- Solving directly for zero gradient is problematic



- Critical points zero gradient
 - Minimum, maximum or inflection point
 - k critical points $\rightarrow k/3$ are minima
- In d dimensions
 - Should be minimum in *d* directions
 - k critical points $\rightarrow k/3^d$ are minima
- Large fraction of critical values are saddle points
- Does not seem to be a problem for SGD
- Solving directly for zero gradient is problematic

