When gradient is small ...

Hung-yi Lee 李宏毅



Warning of Math

Tayler Series Approximation

 $L(\boldsymbol{\theta}) \text{ around } \boldsymbol{\theta} = \boldsymbol{\theta}' \text{ can be approximated below}$ $L(\boldsymbol{\theta}) \approx L(\boldsymbol{\theta}') + \left(\boldsymbol{\theta} - \boldsymbol{\theta}'\right)^T \boldsymbol{g} + \frac{1}{2} \left(\boldsymbol{\theta} - \boldsymbol{\theta}'\right)^T \boldsymbol{H} \left(\boldsymbol{\theta} - \boldsymbol{\theta}'\right)$

Gradient g is a vector

$$\boldsymbol{g} = \nabla L(\boldsymbol{\theta}') \quad \boldsymbol{g}_i = \frac{\partial L(\boldsymbol{\theta}')}{\partial \boldsymbol{\theta}_i}$$

Hessian *H* is a *matrix*

$$\boldsymbol{H}_{ij} = \frac{\partial^2}{\partial \boldsymbol{\theta}_i \partial \boldsymbol{\theta}_j} L(\boldsymbol{\theta'})$$



Hessian



At critical point:
$$v^T H v$$

Hessian $L(\theta) \approx L(\theta') + \frac{1}{2}(\theta - \theta')^T H(\theta - \theta')$
For all v
 $v^T H v > 0 \longrightarrow \text{Around } \theta': L(\theta) > L(\theta') \longrightarrow \text{Local minima}$
 $= H \text{ is positive definite } = \text{All eigen values are positive.}$
For all v
 $v^T H v < 0 \longrightarrow \text{Around } \theta': L(\theta) < L(\theta') \longrightarrow \text{Local maxima}$
 $= H \text{ is negative definite } = \text{All eigen values are negative.}$
Sometimes $v^T H v > 0$, sometimes $v^T H v < 0 \implies \text{Saddle point}$
Some eigen values are positive, and some are negative.



$$x \xrightarrow{w_1} \underbrace{w_2}_{=1} \underbrace{(\hat{y} - w_1 w_2 x)^2}_{=(1 - w_1 w_2)^2} \underbrace{w_1 w_2}_{=1} \underbrace{(1 - w_1 w_2)(-w_2)}_{=0} \underbrace{critical point: w_1 = 0, w_2 = 0}_{H = \begin{bmatrix} 0 & -2 \\ -2 & 0 \end{bmatrix}} \underbrace{\lambda_1 = 2, \lambda_2 = -2}_{A_1 = 2, \lambda_2 = -2} \underbrace{A_1 = 2, \lambda_2 = -2}_{A_1 = 2, \lambda_2 = -2} \underbrace{A_2 = -2, \lambda_2 = -2}_{A_2 = -2} \underbrace{A_2 = -2, \lambda_2 = -2, \lambda_2 = -2, \lambda_2 = -2}_{A_2 = -2, \lambda_2 = -2, \lambda_2$$

2.0

Don't afraid of saddle point? $v^T H v$ At critical point: $L(\theta) \approx L(\theta') + \frac{1}{2}(\theta - \theta')^T H(\theta - \theta')$ Sometimes $v^T H v > 0$, sometimes $v^T H v < 0$ \Rightarrow Saddle point *H* may tell us parameter update direction!

 $m{u}$ is an eigen vector of $m{H}$ λ is the eigen value of $m{u}$ $\lambda < 0$

$$u^T H u = u^T (\lambda u) = \lambda ||u||^2$$

< 0
< 0

$$L(\boldsymbol{\theta}) \approx L(\boldsymbol{\theta}') + \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}')^T \boldsymbol{H} (\boldsymbol{\theta} - \boldsymbol{\theta}') \implies L(\boldsymbol{\theta}) < L(\boldsymbol{\theta}')$$

 $\theta - \theta' = u$ $\theta = \theta' + u$ Decrease L

$$x \xrightarrow{w_{1}} (\hat{y} - w_{1}w_{2}x)^{2} \rightarrow y \Rightarrow \hat{y} = 1$$

$$L = (\hat{y} - w_{1}w_{2}x)^{2} = (1 - w_{1}w_{2})^{2}$$

$$\frac{\partial L}{\partial w_{1}} = 2(1 - w_{1}w_{2})(-w_{2})$$
Critical point: $w_{1} = 0, w_{2} = 0$

$$\frac{\partial U}{\partial w_2} = 2(1 - w_1 w_2)(-w_1) \qquad H = \begin{bmatrix} 0 & -2 \\ -2 & 0 \end{bmatrix} \lambda_1 = 2, \lambda_2 = -2$$

$$\frac{\partial L}{\partial w_2} = 2(1 - w_1 w_2)(-w_1) \qquad \frac{Saddle point}{2}$$

20

$$\lambda_2 = -2$$
 Has eigenvector $\boldsymbol{u} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$
Update the parameter along the direction of \boldsymbol{u}
You can escape the saddle point and decrease the loss.

(this method is seldom used in practice)

0.0

0.5

1.0

1.5

2.0

End of Warning

Saddle Point v.s. Local Minima

• A.D. 1543



來源《三體Ⅲ·死神永生》

Saddle Point v.s. Local Minima

• The Magician Diorena (魔法師狄奧倫娜)

From 3 dimensional space, it is sealed.

It is not in higher dimensions.



Source of image: https://read01.com/mz2DBPE.html#.YECz22gzbIU

Source of image: https://arxiv.org/abs/1712.09913

Saddle Point v.s. Local Minima





higher dimension?

When you have lots of parameters, perhaps local minima is rare?



Source:https://docs.google.com/presentation/d/1siUFXARYRpNiMeSRwgFbt7mZVjkMPhR5od09w0Z8xaU/edit#slide=id.g3 1470fd33a_0_33

Small Gradient ...



The value of a network parameter w

Tips for training: Batch and Momentum

Batch

Review: Optimization with Batch



1 epoch = see all the batches once - Shuffle after each epoch

Consider 20 examples (N=20)

Batch size = N (Full batch)

Update after seeing all the 20 examples



Batch size = 1

Update for each example Update 20 times in an epoch



oldest slides: http://speech.ee.ntu.edu.tw/~tlkagk/courses/MLDS_2015_2/Lecture/DNN%20(v4).pdf old slides: http://speech.ee.ntu.edu.tw/~tlkagk/courses/ML_2017/Lecture/Keras.pdf

Small Batch v.s. Large Batch

• Larger batch size does **not** require longer time to compute gradient (unless batch size is too large)



• Smaller batch requires longer time for one epoch (longer time for seeing all data once)

Time for one **update**



Time for one **epoch**

Consider 20 examples (N=20)

Batch size = N (Full Batch)

Update after seeing all the 20 examples



Batch size = 1

Update for each example Update 20 times in an epoch



MNIST





Smaller batch size has better performance

What's wrong with large batch size? Optimization Fails

- Smaller batch size has better performance
- "Noisy" update is better for training



On Large-Batch Training for Deep Learning: Generalization Gap and Sharp Minima https://arxiv.org/abs/1609.04836

Small Batch v.s. Large Batch

• Small batch is better on testing data?

| | Name | Network Type | Data set |
|----------------|-------|-------------------------|---------------------------------------|
| | F_1 | Fully Connected | MNIST (LeCun et al., 1998a) |
| 3B = 256 | F_2 | Fully Connected | TIMIT (Garofolo et al., 1993) |
| | C_1 | (Shallow) Convolutional | CIFAR-10 (Krizhevsky & Hinton, 2009) |
| LB = | C_2 | (Deep) Convolutional | CIFAR-10 |
| 0 1 x data set | C_3 | (Shallow) Convolutional | CIFAR-100 (Krizhevsky & Hinton, 2009) |
| | C_4 | (Deep) Convolutional | CIFAR-100 |

| | Training Accuracy | | Testing Accuracy | |
|-------|----------------------|----------------------|----------------------|----------------------|
| Name | SB | LB | SB | LB |
| F_1 | $99.66\% \pm 0.05\%$ | $99.92\% \pm 0.01\%$ | $98.03\% \pm 0.07\%$ | $97.81\% \pm 0.07\%$ |
| F_2 | $99.99\% \pm 0.03\%$ | $98.35\% \pm 2.08\%$ | $64.02\% \pm 0.2\%$ | $59.45\% \pm 1.05\%$ |
| C_1 | $99.89\% \pm 0.02\%$ | $99.66\% \pm 0.2\%$ | $80.04\%\pm 0.12\%$ | $77.26\% \pm 0.42\%$ |
| C_2 | $99.99\% \pm 0.04\%$ | $99.99\% \pm 0.01\%$ | $89.24\%\pm 0.12\%$ | $87.26\%\pm 0.07\%$ |
| C_3 | $99.56\% \pm 0.44\%$ | $99.88\% \pm 0.30\%$ | $49.58\% \pm 0.39\%$ | $46.45\% \pm 0.43\%$ |
| C_4 | $99.10\% \pm 1.23\%$ | $99.57\% \pm 1.84\%$ | $63.08\% \pm 0.5\%$ | $57.81\% \pm 0.17\%$ |

On Large-Batch Training for Deep Learning: Generalization Gap and Sharp Minima https://arxiv.org/abs/1609.04836

Small Batch v.s. Large Batch



| | Small | Large |
|---|----------|----------------------|
| Speed for one update (no parallel) | Faster | Slower |
| Speed for one update (with parallel) | Same | Same (not too large) |
| Time for one epoch | Slower | Faster |
| Gradient | Noisy | Stable |
| Optimization | Better 💥 | Worse |
| Generalization | Better | Worse |

Batch size is a hyperparameter you have to decide.

Have both fish and bear's paws?

- Large Batch Optimization for Deep Learning: Training BERT in 76 minutes (https://arxiv.org/abs/1904.00962)
- Extremely Large Minibatch SGD: Training ResNet-50 on ImageNet in 15 Minutes (https://arxiv.org/abs/1711.04325)
- Stochastic Weight Averaging in Parallel: Large-Batch Training That Generalizes Well (https://arxiv.org/abs/2001.02312)
- Large Batch Training of Convolutional Networks (https://arxiv.org/abs/1708.03888)
- Accurate, large minibatch sgd: Training imagenet in 1 hour (https://arxiv.org/abs/1706.02677)

Momentum

Small Gradient ...

Consider the physical world ... Loss How about put this phenomenon in gradient descent?

The value of a network parameter w

(Vanilla) Gradient Descent



Starting at θ^0 Compute gradient g^0 Move to $\theta^1 = \theta^0 - \eta g^0$ Compute gradient g^1 Move to $\theta^2 = \theta^1 - \eta g^1$

Gradient Descent + Momentum

Movement: **movement of last step** minus **gradient** at present



Starting at θ^0 Movement $m^0 = 0$ Compute gradient g^0 Movement $m^1 = \lambda m^0 - \eta g^0$ Move to $heta^1 = heta^0 + m^1$ Compute gradient g^1 Movement $m^2 = \lambda m^1 - \eta g^1$ Move to $\theta^2 = \theta^1 + m^2$

Movement not just based on gradient, but previous movement.

Gradient Descent + Momentum

Movement: **movement of last step** minus **gradient** at present

 m^i is the weighted sum of all the previous gradient: g^0 , g^1 , ..., g^{i-1}



Starting at θ^0 Movement $m^0 = 0$ Compute gradient g^0 Movement $m^1 = \lambda m^0 - \eta g^0$ Move to $\theta^1 = \theta^0 + m^1$ Compute gradient g^1 Movement $m^2 = \lambda m^1 - \eta g^1$ Move to $\theta^2 = \theta^1 + m^2$

Movement not just based on gradient, but previous movement.

Gradient Descent + Momentum



Concluding Remarks

- Critical points have zero gradients.
- Critical points can be either saddle points or local minima.
 - Can be determined by the Hessian matrix.
 - It is possible to escape saddle points along the direction of eigenvectors of the Hessian matrix.
 - Local minima may be rare.
- Smaller batch size and momentum help escape critical points.

Acknowledgement

 感謝作業二助教團隊(陳宣叡、施貽仁、孟妍李 威緒)幫忙跑實驗以及蒐集資料