## Machine Learning Soundings (機器學習深測)



#### Lecture 3: Optimization in Deep Learning

Hsuan-Tien Lin (林軒田) htlin@csie.ntu.edu.tw

Department of Computer Science & Information Engineering

National Taiwan University (國立台灣大學資訊工程系)



## Roadmap

#### 1 Deep Learning Foundations

Lecture 1: Neural Network

automatic pattern feature extraction from layers of neurons with backprop for GD/SGD

Lecture 3: Optimization in Deep Learning

• Difficulty of Deep Learning Optimization

2 Deep Learning Models

## Difficulty of Deep Learning Optimization

#### error surface complicated

- local minima: not as bad as imagined
- saddle points/local maxima: easily escapable (especially with SGD)
- plateau: need larger learning rate  $\eta$
- ravines: need to avoid oscillation

### stability <> computation trade-off

slow computation of gradient (backprop)

- $\Rightarrow$  SGD on minibatch
- $\Rightarrow$  'instable' estimate of gradient

#### getting more stable estimate? averaging

Optimization in Deep Learning Difficulty of Deep Learning Optimization Running Average Estimate of Gradient

gradient descent:  $\mathbf{w}_t \leftarrow \mathbf{w}_{t-1} - \eta \cdot \mathbf{v}_t$ 

original minibatch SG

gradient estimate  $\mathbf{v}_t = \Delta_t$  from one minibatch SG

#### averaging by multiple SG

if minibatch SG for *M* times at *t*-th iteration, each getting  $\Delta_t^{(m)}$ , more stable gradient estimate by uniform averaging  $\mathbf{v}_t = \frac{1}{M} \sum_{m=1}^{M} \Delta_t^{(m)}$ —needing *M* times more computation than original minibatch SGD

speedup by reusing each  $\Delta_t = \Delta_t^{(1)}$ 

 $\mathbf{v}_t = \frac{1}{M} \sum_{m=1}^{M} \Delta_{t-m+1}$  —'moving window' average of SG

### issue with 'moving window' average:

#### uniformly weighted

Hsuan-Tien Lin (NTU CSIE)

Machine Learning Soundings

Optimization in Deep Learning

Difficulty of Deep Learning Optimization

## Averaging SG Non-uniformly

**Running Average** 

$$\mathbf{v}_t = \beta \mathbf{v}_{t-1} + (1 - \beta) \Delta_t$$

with  $0 \le \beta < 1$  to control how much history to take  $\beta = 0$ : original SGD

$$\mathbf{v}_t = \sum_{m=1}^t \beta^{t-m} (1-\beta) \Delta_t$$

-size-t window, exponentially-decreasing aeveraging

SGD with momentum: optimization direction = current SG ( $\Delta_t$ ) + historical inertia ( $\mathbf{v}_{t-1}$ )

## Benefits of SGD with Momentum

$$\mathbf{v}_t = \beta \mathbf{v}_{t-1} + (1-\beta) \Delta_t$$
$$\mathbf{w}_t = \mathbf{w}_{t-1} - \eta \mathbf{v}_t$$

- some variance in SG canceled out
- oscilliation across ravine dampened
- shallow local optima/saddle points escaped

# SGD with momentum: 'stablize' SG with running average

Difficulty of Deep Learning Optimization

## Per-Component Learning Rate

fixed learning rate :  $\mathbf{w}_t = \mathbf{w}_{t-1} - \eta \mathbf{v}_t$ per-component learning rate :  $\mathbf{w}_t = \mathbf{w}_{t-1} - \eta_t \odot \mathbf{v}_t$ 

intuition: scales error surface

want: smaller step for larger gradient component

## Running Average of Gradient Magnitude

want: smaller step for larger gradient component, say

$$\boldsymbol{\eta}_t = \frac{1}{\|\nabla \boldsymbol{E}(\mathbf{w}_t)\|}$$

- full gradient  $\nabla E$  not available, SG only
- using  $\|\Delta\|$  not very stable

idea: running average of  $\Delta_t \odot \Delta_t$ 

Difficulty of Deep Learning Optimization

RMSProp

$$\mathbf{u}_{t} = \beta \mathbf{u}_{t-1} + (1 - \beta) \Delta_{t} \odot \Delta_{t}$$
  
$$\boldsymbol{\eta}_{t} = \boldsymbol{\eta}_{t} / \sqrt{\mathbf{u}_{t} + \epsilon} \mathbf{w}_{t} = \mathbf{w}_{t-1} - \boldsymbol{\eta}_{t} \odot \Delta_{t}$$

RMSProp: SGD + per-component learng rate using running average of magnitude

## Adam: Adaptive Moment Estimation

#### Adam $\approx$ momentum + RMSProp + global decay

$$\mathbf{v}_t = \beta_1 \mathbf{v}_{t-1} + (1 - \beta_1) \Delta_t$$
  

$$\mathbf{u}_t = \beta_2 \mathbf{u}_{t-1} + (1 - \beta_2) \Delta_t \odot \Delta_t$$
  

$$\eta_t = \eta_t / \sqrt{\mathbf{u}_t + \epsilon} / \sqrt{t/N}$$
  

$$\mathbf{w}_t = \mathbf{w}_{t-1} - \eta_t \odot \mathbf{v}_t$$

- momentum in v<sub>t</sub>
- RMSProp in **u**<sub>t</sub>
- global decay by  $\sqrt{t/N}$
- (some minor correction of estimation)

# Adam usually more aggressive than original SGD (but can also overfit faster)

Hsuan-Tien Lin (NTU CSIE)