

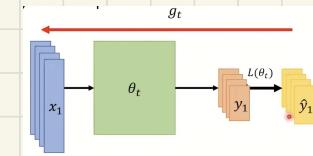
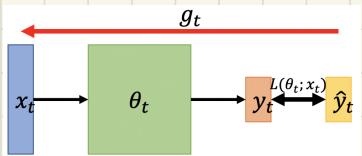
Chapter 4 New Optimization for Deep Learning

概念定义: θ_t : 相应的模型参数

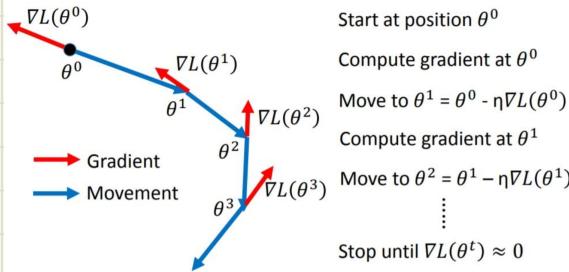
$\nabla L(\theta_t)$ or g_t : 在参数下的梯度, 用来计算 θ_{t+1}

MTR: 从起始时刻到时间 t, 关于梯度的累加值 momentum, 用来计算 θ_{t+1}

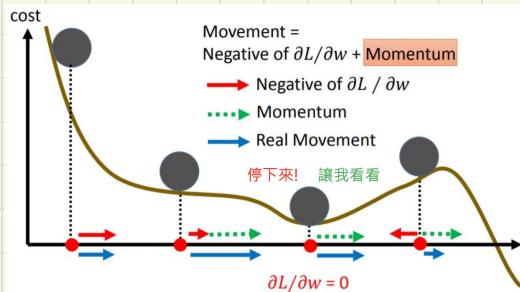
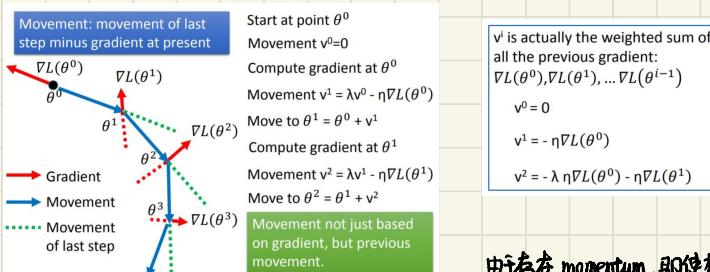
网络模型:



1. Stochastic Gradient Descent (SGD, 1847)



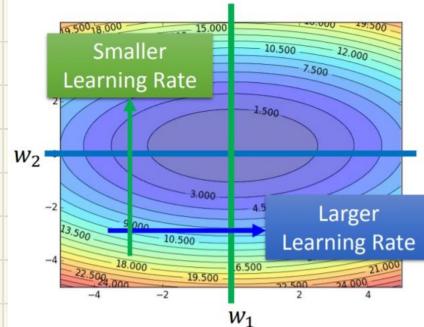
2. Stochastic Gradient Descent with Momentum (SGDM, 1986)



3. Adagrad (2011)

$$\theta_t = \theta_{t-1} - \frac{\eta}{\sqrt{\sum_{i=0}^{t-1} (g_i)^2}} g_{t-1}$$

- 对于不同的参数采用不同的learning rate, 以避免在一些梯度变化过大的维度步进过大



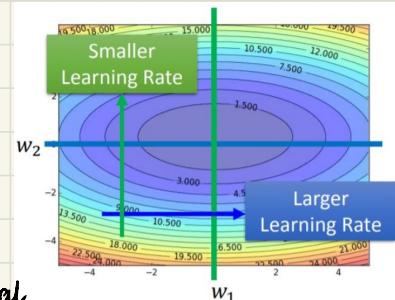
4. RMSProp (2013)

$$\theta_t = \theta_{t-1} - \frac{\eta}{\sqrt{v_t}} g_{t-1}$$

$$v_1 = g_0^2$$

$$v_t = \alpha v_{t-1} + (1 - \alpha)(g_{t-1})^2$$

- 改变了Adagrad中学习率的变化方法, Adagrad中由于分子中不断的累加梯度, 如果梯度过大, 学习率变得很小, 则学习近乎停止
- RMSProp借鉴了momentum的概念, 保证了分母不会无止境的变大 (Exponential moving average (EMA) of squared gradients is not monotonically increasing.)



5. Adam (2014)

• SGDM

$$\theta_t = \theta_{t-1} - \eta m_t$$

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_{t-1}$$



$$\theta_t = \theta_{t-1} - \frac{\eta}{\sqrt{\hat{v}_t + \epsilon}} \hat{m}_t$$

• RMSProp

$$\theta_t = \theta_{t-1} - \frac{\eta}{\sqrt{v_t}} g_{t-1}$$

$$v_1 = g_0^2$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2)(g_{t-1})^2$$

$$\begin{aligned} \hat{m}_t &= \frac{m_t}{1 - \beta_1^t} \\ \hat{v}_t &= \frac{v_t}{1 - \beta_2^t} \end{aligned}$$

de-biasing

$\beta_1 = 0.9$
 $\beta_2 = 0.999$
 $\epsilon = 10^{-8}$

- 学习率的变化部分借鉴了RMSProp
- 梯度部分借鉴了SGDM

除以 $1-\beta$ 的原因是, m_t 经过一发时间会慢下来, 但最初 m_t 比较小, 所以除以一个大于1的值适当放大以靠近近期的稳定值

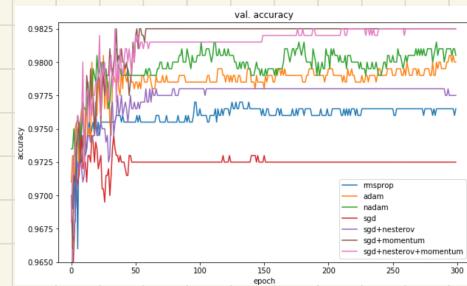
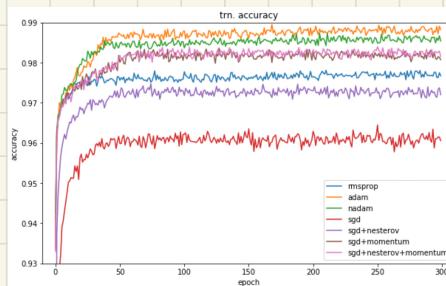
6. Optimizer in Real Application

Adam: Bert, Transformer, Tacotron, Big-GAN, MAML

SGDM: YOLO, ResNet, Mask R-CNN

① Adam vs. SGD

领域内比较有名的 Neural Network 都出自于 SGD 和 Adam，原因可能是两者在 optimizer 中各占了一个极端



Adam: fast training, large generalization gap, unstable

SGDM: stable, little generalization gap, better convergence

因此，SWATs (2017) 年提出了 Begin with Adam, end with SGDM



② 改进 Adam

$$\theta_t = \theta_{t-1} - \frac{\eta}{\sqrt{v_t + \epsilon}} \hat{m}_t$$

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_{t-1}, \beta_1 = 0.$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2)(g_{t-1})^2, \beta_2 = 0.999$$

~~~~~ 每一次乘以 0.999，则此表示的是一步步地将数的相加， $v_t$  受到梯度的影响是  $\frac{1}{1-0.999} = 1000$  步

The "memory" of  $v_t$  keeps roughly 1000 steps!

Adam 的作者发现在训练的后期，大多数梯度值都很小，是 non-informative 的，只有个别 mini-batches 会提供比较大的 informative gradient 以帮助模型下降的方向

| time step | ... | 100000 | 100001 | 100002 | 100003 | ... | 100999            | 101000          |
|-----------|-----|--------|--------|--------|--------|-----|-------------------|-----------------|
| gradient  |     | 1      | 1      | 1      | 1      |     | 100000            | 1               |
| movement  |     | $\eta$ | $\eta$ | $\eta$ | $\eta$ |     | $10\sqrt{10}\eta$ | $10^{-3.5}\eta$ |

对于第 100999 步：

$$\beta_1 = 0 \rightarrow \hat{m}_t = g_{t-1} = 10^5$$

假设之前的 gradient 为 1，则此时  $v_{t-1} = 1$ ， $v_t = 0.999 \times 1 + 0.001 \times 10^5 \approx 10^7$

$$\frac{\eta}{\sqrt{v_t + \epsilon}} \approx \eta \cdot \frac{10^5}{\sqrt{10^7}} = 10\sqrt{10}\eta$$

第 100000 ~ 100998 步并没有提供有意义的梯度，但前进了  $1000\eta$ ；第 100999 步提供了有意义的梯度，却只前进了  $10\sqrt{10}\eta$

↑  
受限于 Adam 中一次更新的最大绝对幅值为  $\sqrt{1-\beta_2}\eta$

假设已经训练了 10 万个 timestamp 且之前的梯度都不大，在这之后碰到了一个很大的梯度

- AMSGrad (ICLR'18)

$$\theta_t = \theta_{t-1} - \frac{\eta}{\sqrt{\hat{v}_t + \epsilon}} m_t$$

$$\hat{v}_t = \max(\hat{v}_{t-1}, v_t)$$

- 通过  $\max$  使  $\hat{v}_t$  记住之前遇到过的最大的梯度，防止在长时间的小梯度进行无效的学习；但只解决了 large learning rate 的情况
- 又有可能像 AdaGard 一样导致分母过大

- AdaBound (ICLR'19)

$$\theta_t = \theta_{t-1} - Clip\left(\frac{\eta}{\sqrt{\hat{v}_t + \epsilon}}\right) \hat{m}_t$$

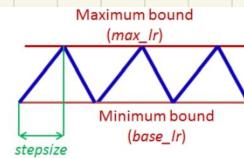
$$Clip(x) = Clip\left(x, 0.1 - \frac{0.1}{(1 - \beta_2)t + 1}, 0.1 + \frac{0.1}{(1 - \beta_2)t}\right)$$

- 通过限制上下界进行裁断的方法防止学习率过大或过小（失去了 Adaptive 的意义）

## ② 改进 SGD

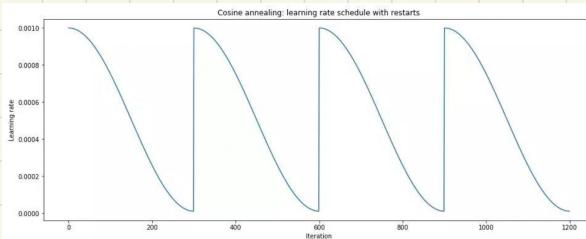
虽然 Adaptive learning rate algorithm 的思想是动态调整学习率，但 SGD 无法做到动态调整，退而求其次寻找一个最优的 “fix learning rate”

- LR Range test [Smith, WACV'17]：学习率过大或过小，performance 都不是最好的；学习率适中效果最好
- Cyclical LR [Smith, WACV'17]：学习率在上界和下界之间反复变换可以取得不错的 performance

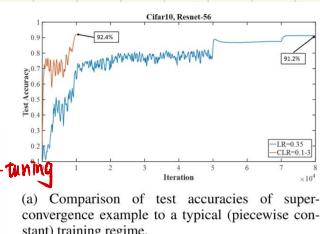
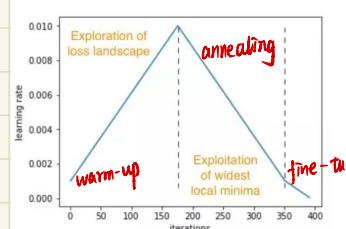


- 遇到 local minima 通过增大学习率跳出局部极小值
- 找到 global minima 后即使增大学习率也不会跳的太久

- SGDR [Loshchilov, ICLR'17]



- One-cycle LR



## 7. Warm-up in Adam

### ① Warm-up 对 Adam 的影响

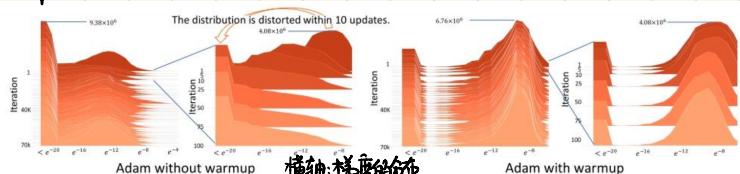
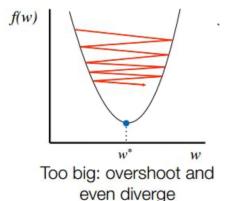


Figure 2: The absolute gradient histogram of the Transformers on the De-En IWSLT' 14 dataset during the training (stacked along the y-axis). X-axis is absolute value in the log scale and the height is the frequency. Without warmup, the gradient distribution is distorted in the first 10 steps.



distorted gradient  
 ↓  
 distorted EMA squared gradient  
 ↓  
 bad learning rate

• 初期混乱的梯度导致 EMA 估计梯度分布不准，进而影响学习率，使得步长过大来回震荡

• 在前期保持一个大的步长，保证即使梯度估计不准也不要大步前进

### ② RAdam

$$\rho_t = \rho_\infty - \frac{2t\beta_2^t}{1-\beta_2^t} \rightarrow \text{effective memory size of EMA}$$

$$\rho_\infty = \frac{2}{1-\beta_2} - 1 \rightarrow \text{max memory size}$$

$$r_t = \sqrt{\frac{\rho_t - 4}{(\rho_\infty - 4)(\rho_\infty - 2)\rho_t}} \rightarrow \text{approximated } \frac{\text{Var}[\frac{1}{\rho_t}]}{\text{Var}[\frac{1}{\rho_\infty}]} \text{ (for } \rho_t > 4\text{)}$$

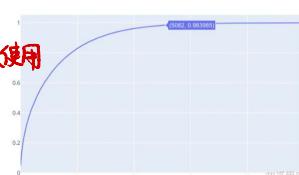
When  $\rho_t \leq 4$  (first few steps of training)

$$\theta_t = \theta_{t-1} - \eta \hat{m}_t \quad (\rho_t < 4 \text{ 使用 Adam, } \rho_t \geq 4 \text{ 使用 SGDM})$$

When  $\rho_t > 4$

$$\theta_t = \theta_{t-1} - \frac{\eta r_t}{\sqrt{\hat{v}_t + \epsilon}} \hat{m}_t$$

$r_t$  is increasing through time!



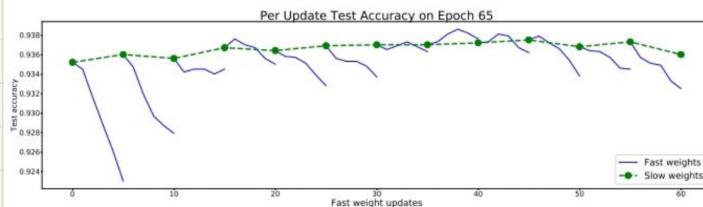
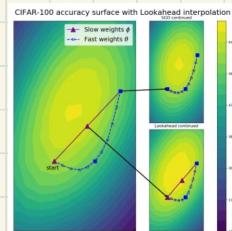
RAdam VS. SWATS :

|              | RAdam                                                                                            | SWATS                                                                 |
|--------------|--------------------------------------------------------------------------------------------------|-----------------------------------------------------------------------|
| Inspiration  | Distortion of gradient at the beginning of training results in inaccurate adaptive learning rate | non-convergence and generalization gap of Adam, slow training of SGDM |
| How?         | Apply warm-up learning rate to reduce the influence of inaccurate adaptive learning rate         | Combine their advantages by applying Adam first, then SGDM            |
| Switch       | SGDM to RAdam                                                                                    | Adam to SGDM                                                          |
| Why switch   | The approximation of the variance of $\hat{v}_t$ is invalid at the beginning of training         | To pursue better convergence                                          |
| Switch point | When the approximation becomes valid                                                             | Some human-defined criteria                                           |

### ③ Lookahead [Zhang, et al., arXiv'19]

核心思想: k step forward, 1 step back "universal wrapper for all optimizer"

For  $t = 1, 2, \dots$  (outer loop)       $\theta$ : fast weight (用浅绿色)       $\phi$ : slow weight  
 $\theta_{t,0} = \phi_{t-1}$   
 For  $i = 1, 2, \dots, k$  (inner loop)  
 $\theta_{t,i} = \theta_{t,i-1} + \text{Optim}(\text{Loss}, \text{data}, \theta_{t,i-1})$       {使用optimizer前进 k 步  
 $\phi_t = \phi_{t-1} + \alpha(\theta_{t,k} - \phi_{t-1})$       向 k 步之前的状态回退一部分}



左图表示一次训练后期的两组权重的变化情况。fast weight不稳定，有可能向坏的方向前进，但在slow weight的控制下，可以保持比较稳定的结果

one step back 的作用是避免过于危险的 exploration，使算法更 stable

Lookahead 算法有效的原因是可以避免进入过于崎岖的地方，尽量在较平坦的地方寻找最优解，使得模型具有更好的 generalization

### ④ Momentum recap (梯度记忆，与 Adam 关系不大): 用未来可能到达的位置梯度参与当前的计算

## 8. Future Position in Current Step

### ① Nesterov accelerated gradient (NAG) Nesterov, ed., SSSR'1983)

$$\text{SGDM} \left\{ \begin{array}{l} \theta_t = \theta_{t-1} - m_t \\ m_t = \lambda m_{t-1} + \eta \nabla L(\theta_{t-1}) \end{array} \right. \rightarrow \left\{ \begin{array}{l} \theta_t = \theta_{t-1} - m_t \\ m_t = \lambda m_{t-1} + \eta \nabla L(\theta_{t-1} - \lambda m_{t-1}) \end{array} \right.$$

使用该项使下次可能到达的位置参与当前的计算中，在保留  $m_t$  的同时，拷贝一份用于计算项，所以需要备份  $\theta_{t-1}$

$$\begin{aligned} \hat{\theta}_t' &= \theta_t - \lambda m_t = \theta_{t-1} - m_t - \lambda m_t \\ &= \theta_{t-1} - \lambda m_t - \lambda m_{t-1} - \eta \nabla L(\theta_{t-1} - \lambda m_{t-1}) \\ &= \theta_{t-1}' - \lambda m_{t-1} - \eta \nabla L(\theta_{t-1}') \end{aligned}$$

$$m_t = \lambda m_{t-1} + \eta \nabla L(\theta_{t-1}')$$

两者相比，NAG 在本次迭代中就用到了即将要到达位置的 momentum (将  $m_t$  超前部署一个 timestamp)

将 SGDM 的二式代入一式可得:  $\theta_t = \theta_{t-1} - \lambda m_{t-1} + \eta \nabla L(\theta_{t-1})$ ,  $m_t = \lambda m_{t-1} + \eta \nabla L(\theta_{t-1})$

## ② Nadam [Dozat, ICLR workshop'16]

将NAG超前部署的梯度加入到Adam中

$$\theta_t = \theta_{t-1} - \frac{\eta}{\sqrt{\hat{v}_t + \epsilon}} \hat{m}_t$$

$$\hat{m}_t = \frac{\beta_1 m_t}{1 - \beta_1^{t+1}} + \frac{(1 - \beta_1) g_{t-1}}{1 - \beta_1^t}$$

SGDM

$$\hat{m}_t = \frac{1}{1 - \beta_1^t} (\beta_1 m_{t-1} + (1 - \beta_1) g_{t-1})$$

$$= \frac{\beta_1 m_{t-1}}{1 - \beta_1^t} + \frac{(1 - \beta_1) g_{t-1}}{1 - \beta_1^t}$$

## ③ AdamW & SGDW with momentum [Loshchilov, arxiv'17]

在2014年至2017年，大家对于正则化和Adam的用途存在误区：

$$L_{l_2}(\theta) = L(\theta) + \gamma ||\theta||^2$$

SGD

$$\theta_t = \theta_{t-1} - \nabla L_{l_2}(\theta_{t-1})$$

$$= \theta_{t-1} - \nabla L(\theta_{t-1}) - \gamma \theta_{t-1}$$

SGDM

$$\theta_t = \theta_{t-1} - \lambda m_{t-1} - \eta (\nabla L(\theta_{t-1}) + \gamma \theta_{t-1})$$

$$m_t = \lambda m_{t-1} + \eta (\nabla L(\theta_{t-1}) + \gamma \theta_{t-1}) ?$$

$$m_t = \lambda m_{t-1} + \eta (\nabla L(\theta_{t-1})) ?$$

Adam

$$m_t = \lambda m_{t-1} + \eta (\nabla L(\theta_{t-1}) + \gamma \theta_{t-1}) ?$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) (\nabla L(\theta_{t-1}) + \gamma \theta_{t-1})^2 ?$$

在计算 momentum 时是否应该加入正则项共同计算，14至17年普遍采用的是加入的方案。

在2017年，Loshchilov提出正则项加入momentum中的计算效果更好，并取得了一定范围的应用（NLP领域）

SGDWM

$$\theta_t = \theta_{t-1} - m_t - \gamma \theta_{t-1}$$

$$m_t = \lambda m_{t-1} + \eta (\nabla L(\theta_{t-1}))$$

AdamW

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) \nabla L(\theta_{t-1})$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) (\nabla L(\theta_{t-1}))^2$$

$$\theta_t = \theta_{t-1} - \eta \left( \frac{1}{\sqrt{\hat{v}_t + \epsilon}} \hat{m}_t + \gamma \theta_{t-1} \right)$$

\* 本章最有用的一个 Optimizer

## 9. Something Helps Optimization

### ① More Exploration:

Shuffling：每一代训练时需要打乱数据集的顺序；每次要重新切分 mini-batch，让不同的样本可以在一个 batch 中

Dropout: 增加梯度随机性

Gradient Noise: 在梯度中增加一个高斯分布的噪声，噪声随着 t 的增大而减小

$$\begin{cases} g_{t,i} = g_{t,i} + N(0, \sigma^2) \\ \sigma_t = \frac{c}{(1+t)^r} \end{cases}$$

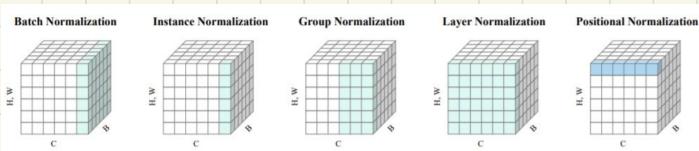
## ② Learning Rate (Teach your model patiently!)

Warm-up: 学习率开始比较小，之后逐渐变大

Curriculum Learning: 开始使用比较简单的数据训练模型，之后再用带有噪声的数据提高模型的泛化能力

Fine-tuning: 使用 pre-trained model 进行验证，避免过多的浪费资源和时间

## ③ Normalization



## ④ Regularization

## 10. Conclusion

| SGD Team                         |                                                                                                                          | Adam Team |                                                                                          |
|----------------------------------|--------------------------------------------------------------------------------------------------------------------------|-----------|------------------------------------------------------------------------------------------|
| 基本算法                             | <ul style="list-style-type: none"> <li>SGM</li> <li>SGDM</li> </ul>                                                      | 基本算法      | <ul style="list-style-type: none"> <li>Adagard</li> <li>RMSProp</li> <li>Adam</li> </ul> |
| 结合两者 : SWATS                     |                                                                                                                          |           |                                                                                          |
| 学习率改进                            | <ul style="list-style-type: none"> <li>LR Range Test</li> <li>Cyclical LR</li> <li>SGDR</li> <li>One-Cycle LR</li> </ul> | 避免极端学习率   | <ul style="list-style-type: none"> <li>AMSGard</li> <li>AdaBound</li> </ul>              |
|                                  |                                                                                                                          | Warm-up   | <ul style="list-style-type: none"> <li>Radam</li> </ul>                                  |
| 未来动量                             | <ul style="list-style-type: none"> <li>NAG</li> </ul>                                                                    | 引入NAG     | <ul style="list-style-type: none"> <li>NAdam</li> </ul>                                  |
| 正则项修正                            | <ul style="list-style-type: none"> <li>SGDW</li> </ul>                                                                   | 正则项修正     | <ul style="list-style-type: none"> <li>AdamW</li> </ul>                                  |
| Wrapper of Optimizer : Lookahead |                                                                                                                          |           |                                                                                          |

| SGDM                                                                                                                                            | Adam                                                                                                                                                                                   |
|-------------------------------------------------------------------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <ul style="list-style-type: none"> <li>Slow</li> <li>Better convergence</li> <li>Stable</li> <li>Smaller generalization gap</li> </ul>          | <ul style="list-style-type: none"> <li>Possibly non-convergence</li> <li>Unstable</li> <li>Larger generalization gap</li> </ul>                                                        |
| SGDM                                                                                                                                            | Adam                                                                                                                                                                                   |
| <ul style="list-style-type: none"> <li>Computer vision</li> <li>image classification</li> <li>segmentation</li> <li>object detection</li> </ul> | <ul style="list-style-type: none"> <li>NLP</li> <li>QA</li> <li>machine translation</li> <li>summary</li> <li>Speech synthesis</li> <li>GAN</li> <li>Reinforcement learning</li> </ul> |