

# SGD or Adam?

Presenter: Wei Li

Advisor: I-Chen Wu

# Outline

- Recap: gradient descent optimization algorithms
  - Gradient descent variants
  - Gradient Descent Optimization Algorithms
- SGD vs. Adam
  - Pros and Cons of Adam
  - SWATS: Adam+SGD
- Choice of optimizer

# Reference

- Keskar, Nitish Shirish, and Richard Socher. "**Improving Generalization Performance by Switching from Adam to SGD.**" arXiv preprint arXiv:1712.07628 (2017).
- Reddi, Sashank J., Satyen Kale, and Sanjiv Kumar. "**On the convergence of adam and beyond.**" International Conference on Learning Representations. 2018.
- Wilson, Ashia C., et al. "**The marginal value of adaptive gradient methods in machine learning.**" Advances in Neural Information Processing Systems. 2017.
- Ruder, Sebastian. "**An overview of gradient descent optimization algorithms.**" *arXiv preprint arXiv:1609.04747*(2016).

# Recap: Gradient descent optimization

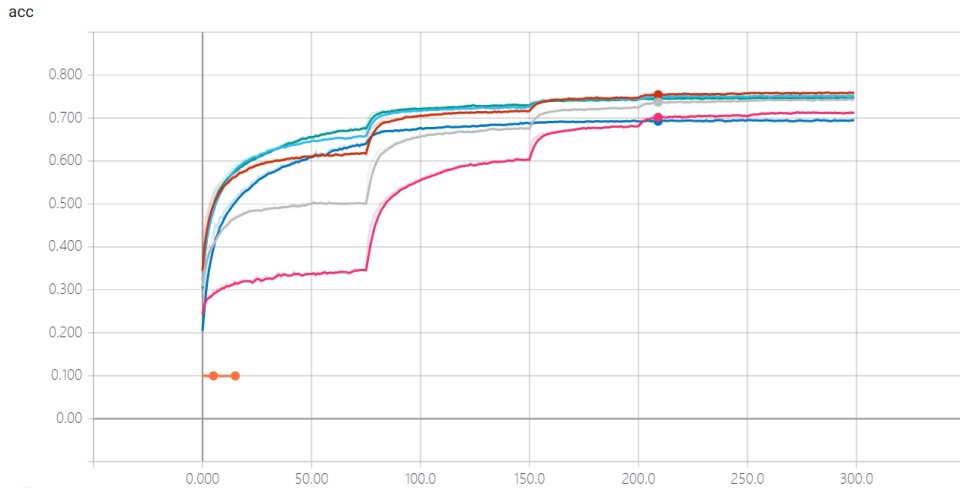
Gradient descent variants

Gradient Descent Optimization Algorithms

# Gradient descent variants

- Batch Gradient Descent(BGD)
- Stochastic Gradient Descent(SGD)
- Mini-Batch Gradient Descent(MBGD)
  - Which is called **SGD** in **Deep learning**
- About batch size:
  - Batch size = 1 (SGD)
    - **Maybe** reach best testing accuracy
    - Sometimes can't converge
    - Training speed: slow
  - Batch size =  $\infty$  (BGD)
    - Can get best training accuracy, bad testing accuracy
    - Training speed: fast
  - Batch size = **a reasonable number** ?
    - can get good performance
    - Training speed: depend on batch size

# About batch size (LeNet on CIFAR-10)



Name	Smoothed Value	Value	Step	Time	Relative
lenet_a_1	0.09986	0.09986	5.000	Fri Mar 9, 02:54:31	11m 8s
lenet_b_10	0.09972	0.09930	15.00	Fri Mar 9, 12:48:45	3m 48s
lenet_c_32	0.7127	0.7136	299.0	Fri Mar 9, 08:01:40	1h 3m 11s
lenet_d_64	0.7436	0.7449	299.0	Fri Mar 9, 09:55:10	58m 13s
lenet_e_128	0.7588	0.7593	299.0	Fri Mar 9, 06:02:26	55m 40s
lenet_f_256	0.7519	0.7533	299.0	Fri Mar 9, 06:58:13	55m 32s
lenet_g_512	0.7467	0.7470	299.0	Fri Mar 9, 08:56:40	54m 46s
lenet_h_1024	0.6945	0.6944	299.0	Fri Mar 9, 03:50:06	53m 32s

# About batch size (LeNet on MNIST)

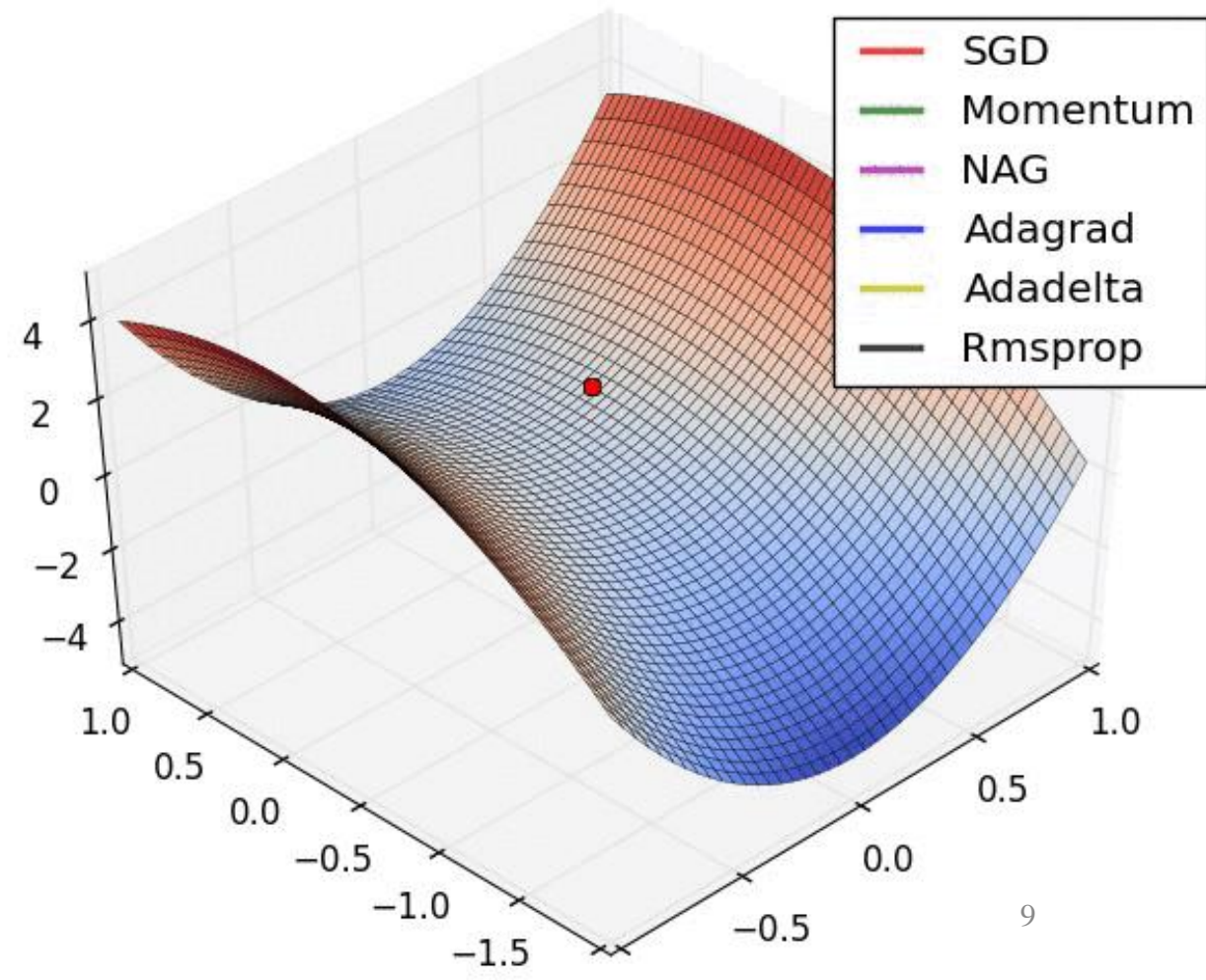
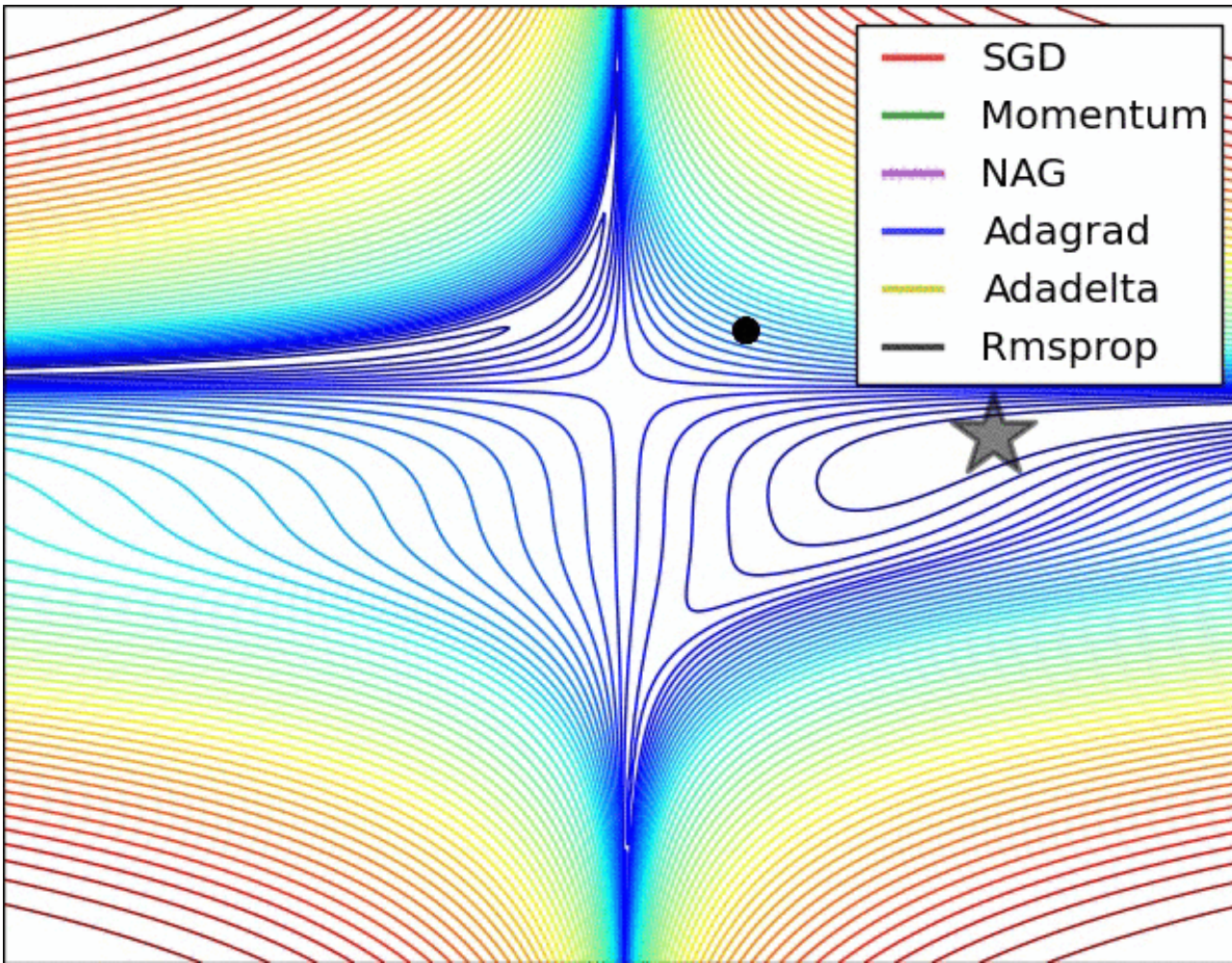
Batch_Size	5000	2000	1000	500	256	100	50	20	10	5	2	1
Total Epoches	200	200	200	200	200	200	200	200	200	200	200	200
Total Iterations	1999	4999	9999	19999	38999	99999	199999	499999	999999	1999999	cannot converge	
Time of 200 Epoches	1	1.068	1.16	1.38	1.75	3.016	5.027	8.513	13.773	24.055		
Achieve 0.99 Accuracy at Epoch	-	-	135	78	41	45	24	9	9	-		
Time of Achieve 0.99 Accuracy	-	-	2.12	1.48	1	1.874	1.7	1.082	1.729	-		
Best Validation Score	0.015	0.011	0.01	0.01	0.01	0.009	0.0098	0.0084	0.01	0.032		
Best Score Achieved at Epoch	182	170	198	100	93	111	38	49	51	17		
Best Test Score	0.014	0.01	0.01	0.01	0.01	0.008	0.0083	0.0088	0.008	0.0262		
Final Test Error (200 epoches)	0.0134	0.01	0.01	0.01	0.01	0.009	0.0082	0.0088	0.008	0.0662		

# Gradient Descent Optimization Algorithms

- SGD
- SGD with Momentum(**SGDM**, default use in many projects)
- SGD with Nesterov Accelerated Gradient (**NAG**)
- **Adaptive optimization algorithms:**
  - AdaGrad
  - AdaDelta / RMSProp
  - Adam
  - Nadam
  - etc.



# Gradient Descent Optimization Algorithms



# SGD

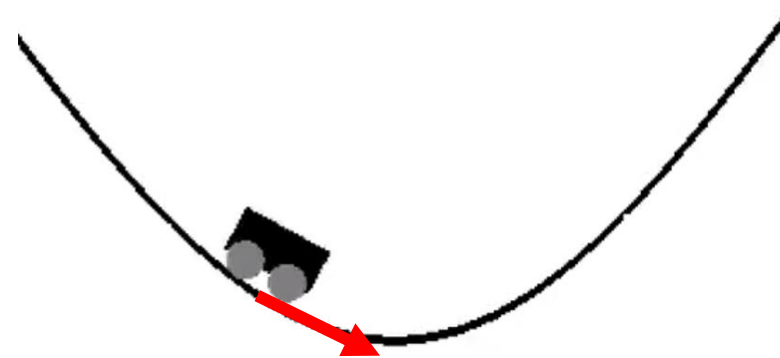
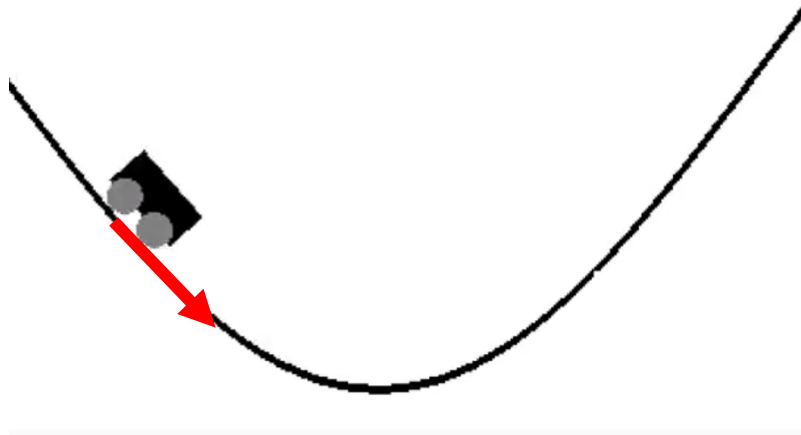
- **SGD**

$$\theta = \theta - \alpha \cdot G(J(\theta))$$

$\theta$ : weights

$\alpha$ : learning rate    $\beta$ : momentum term

$G(\cdot)$ : gradient    $J(\cdot)$ : loss function



# SGD(cont.)

- SGD

$$v = \mathbf{0} + \alpha \cdot G(J(\theta))$$

$$\theta = \theta - v$$

- SGDM

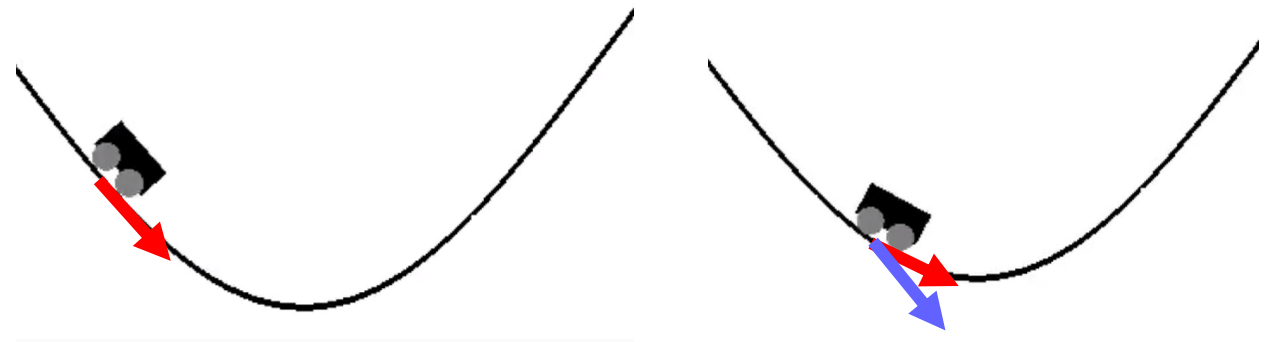
$$v = \beta v + \alpha \cdot G(J(\theta))$$

$$\theta = \theta - v$$

$\theta$ : weights

$\alpha$ : learning rate  $\beta$ : momentum term

$G(\cdot)$ : gradient  $J(\cdot)$ : loss function



# SGD(cont.)

- SGD

$$v = \mathbf{0} + \alpha \cdot G(J(\theta))$$

$$\theta = \theta - v$$

- SGDM

$$v = \beta v + \alpha \cdot G(J(\theta))$$

$$\theta = \theta - v$$

$\theta$ : weights

$\alpha$ : learning rate  $\beta$ : momentum term

$G(\cdot)$ : gradient  $J(\cdot)$ : loss function



# SGD(cont.)

- SGD

$$v = \mathbf{0} + \alpha \cdot G(J(\theta))$$

$$\theta = \theta - v$$

- SGDM

$$v = \beta v + \alpha \cdot G(J(\theta))$$

$$\theta = \theta - v$$

- NAG

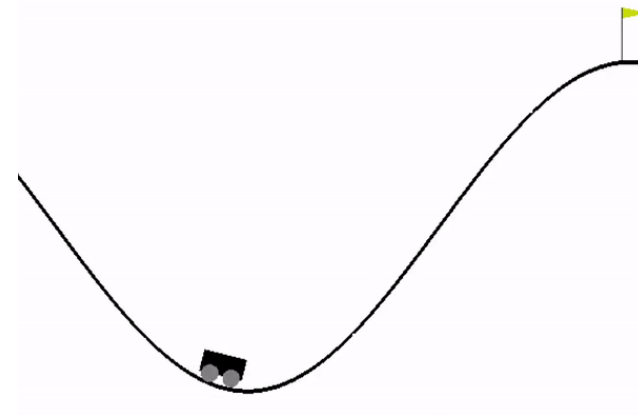
$$v = \beta v + \alpha \cdot G(J(\theta - \beta v))$$

$$\theta = \theta - v$$

$\theta$ : weights

$\alpha$ : learning rate    $\beta$ : momentum term

$G(\cdot)$ : gradient    $J(\cdot)$ : loss function



# SGD(cont.)

- SGD

$$v = \mathbf{0} + \alpha \cdot G(J(\theta))$$

$$\theta = \theta - v$$

- SGDM

$$v = \beta v + \alpha \cdot G(J(\theta))$$

$$\theta = \theta - v$$

- NAG

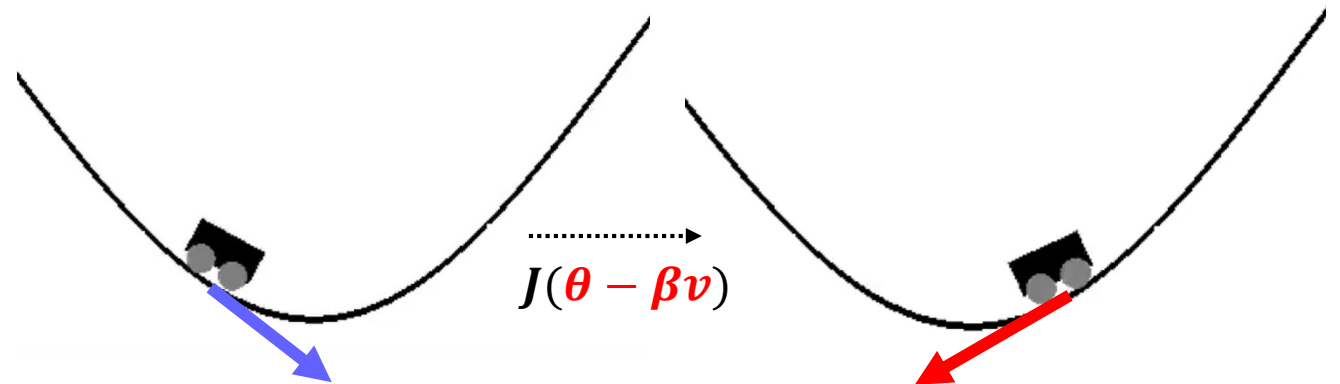
$$v = \beta v + \alpha \cdot G(J(\theta - \beta v))$$

$$\theta = \theta - v$$

$\theta$ : weights

$\alpha$ : learning rate  $\beta$ : momentum term

$G(\cdot)$ : gradient  $J(\cdot)$ : loss function



# SGD(cont.)

- SGD

$$v = \mathbf{0} + \alpha \cdot G(J(\theta))$$

$$\theta = \theta - v$$

- SGDM

$$v = \beta v + \alpha \cdot G(J(\theta))$$

$$\theta = \theta - v$$

- NAG

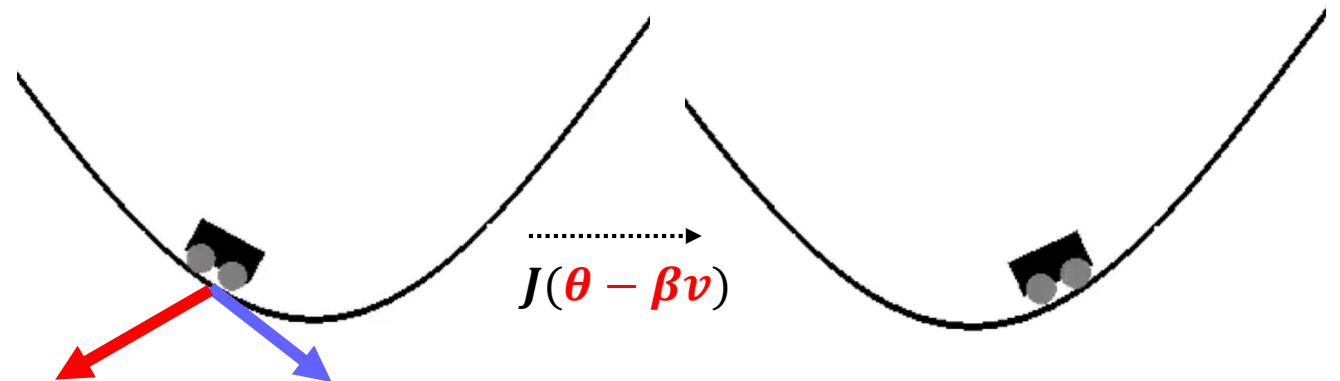
$$v = \beta v + \alpha \cdot G(J(\theta - \beta v))$$

$$\theta = \theta - v$$

$\theta$ : weights

$\alpha$ : learning rate  $\beta$ : momentum term

$G(\cdot)$ : gradient  $J(\cdot)$ : loss function





# SGD(cont.)

- SGD

$$v = \mathbf{0} + \alpha \cdot G(J(\theta))$$

$$\theta = \theta - v$$

- SGDM

$$v = \beta v + \alpha \cdot G(J(\theta))$$

$$\theta = \theta - v$$

- NAG

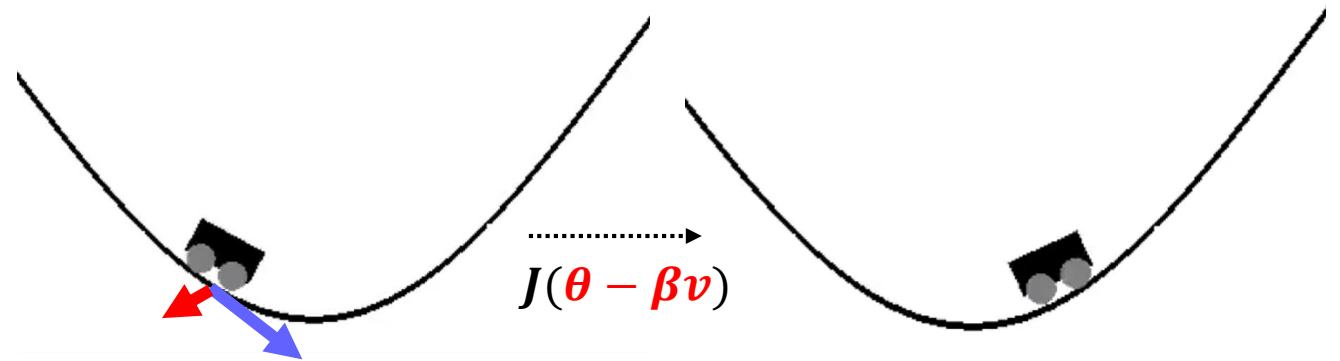
$$v = \beta v + \alpha \cdot G(J(\theta - \beta v))$$

$$\theta = \theta - v$$

$\theta$ : weights

$\alpha$ : learning rate  $\beta$ : momentum term

$G(\cdot)$ : gradient  $J(\cdot)$ : loss function





# SGD(cont.)

- SGD

$$v = \mathbf{0} + \alpha \cdot G(J(\theta))$$

$$\theta = \theta - v$$

- SGDM

$$v = \beta v + \alpha \cdot G(J(\theta))$$

$$\theta = \theta - v$$

- NAG

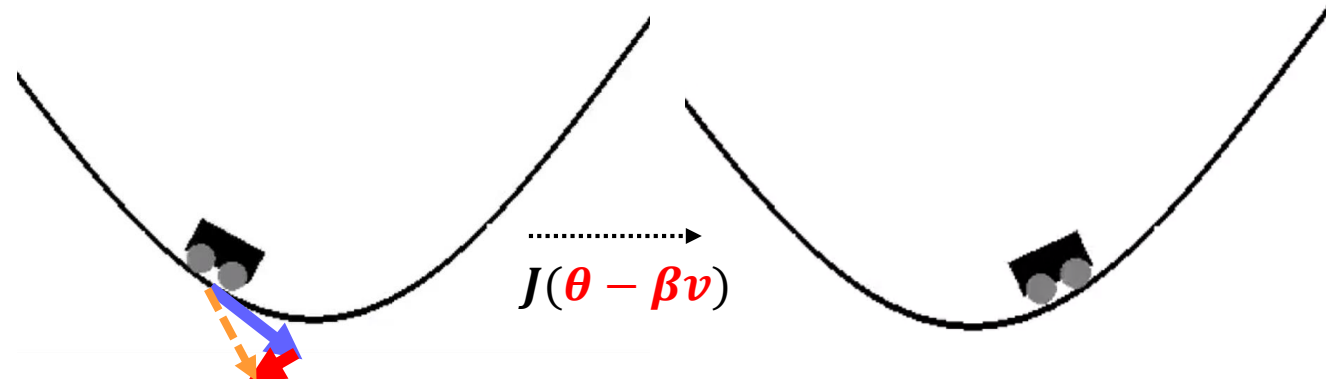
$$v = \beta v + \alpha \cdot G(J(\theta - \beta v))$$

$$\theta = \theta - v$$

$\theta$ : weights

$\alpha$ : learning rate  $\beta$ : momentum term

$G(\cdot)$ : gradient  $J(\cdot)$ : loss function



# Adaptive optimization

- **Adagrad**

$$s = \frac{\alpha}{\sqrt{\sum_{k=1}^t (g_{k,i}^2) + \epsilon}} \cdot g_{t,i}$$
$$\theta = \theta - s$$

$\theta$ : weights

$\alpha$ : learning rate

$\beta$ : momentum term

$G(\cdot)$ : gradient

$J(\cdot)$ : loss function

$$g_{t,i} = G(J(\theta_{t,i}))$$

# Adaptive optimization(cont.)

- Adagrad

$$s = \frac{\alpha}{\sqrt{\sum_{k=1}^t (g_{k,i}^2) + \epsilon}} \cdot g_{t,i}$$
$$\theta = \theta - s$$

$\theta$ : weights

$\alpha$ : learning rate

$\beta$ : momentum term

$G(\cdot)$ : gradient

$J(\cdot)$ : loss function

$$g_{t,i} = G(J(\theta_{t,i}))$$

- Adadelta

$$s = \frac{\alpha}{\sqrt{\mathbf{E}[g^2]_{t,i} + \epsilon}} \cdot g_{t,i}$$
$$\theta = \theta - s$$

# Adaptive optimization(cont.)

- **Adam**

- Adam use estimates of first and second moments of the gradients.

$$m = \beta_1 v + (1 - \beta_1) G^1(J(\theta))$$

$$v = \beta_2 v + (1 - \beta_2) G^2(J(\theta))$$

$$s = \frac{\alpha}{\sqrt{v} + \epsilon} \cdot m$$

$$\theta = \theta - s$$

$\theta$ : weights

$\alpha$ : learnig rate

$\beta_{1,2}$ : momentum term

$G(\cdot)$ : gradient

$J(\cdot)$ : loss function

# SGD vs. Adam

Pros and Cons of Adam

SWATS: Adam+SGD

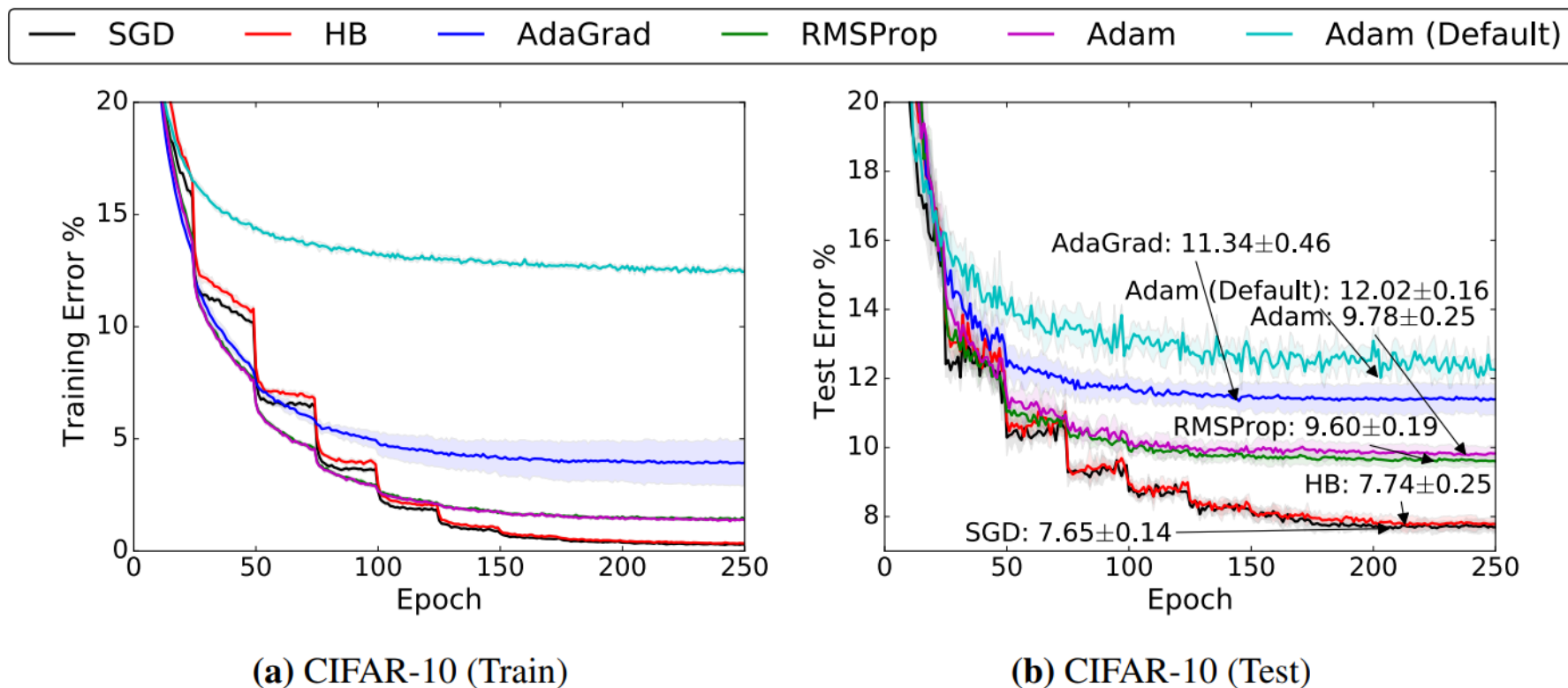
# Pros and Cons of Adam

- Advantages:
  - Fast convergence
  - Good Performance than other adaptive optimization algorithms
  - Needn't to change learning rate by yourself
- Disadvantages:
  - More memory usage
  - Generalize worse (often significantly worse) than SGD, even when these solutions have better training performance
    - **The Marginal Value of Adaptive Gradient Methods in Machine Learning**

# Pros and Cons of Adam

- Advantages:
  - Fast convergence
  - Good Performance than other adaptive optimization algorithms
  - Needn't to change learning rate by yourself
- Disadvantages:
  - More memory usage
  - Generalize worse (often significantly worse) than SGD, even when these solutions have better training performance
    - **The Marginal Value of Adaptive Gradient Methods in Machine Learning**
- **Another Problem:**
  - May not converge to the optimal solution
    - ICLR2018: **On the Convergence of Adam and Beyond**
    - **Many** proofs in appendix

# VGG+BN+Dropout network for CIFAR-10



**Figure 1:** Training (left) and top-1 test error (right) on CIFAR-10. The annotations indicate where the best performance is attained for each method. The shading represents  $\pm$  one standard deviation computed across five runs from random initial starting points. In all cases, adaptive methods are performing worse on both train and test than non-adaptive methods.



# Learning Rate Clip

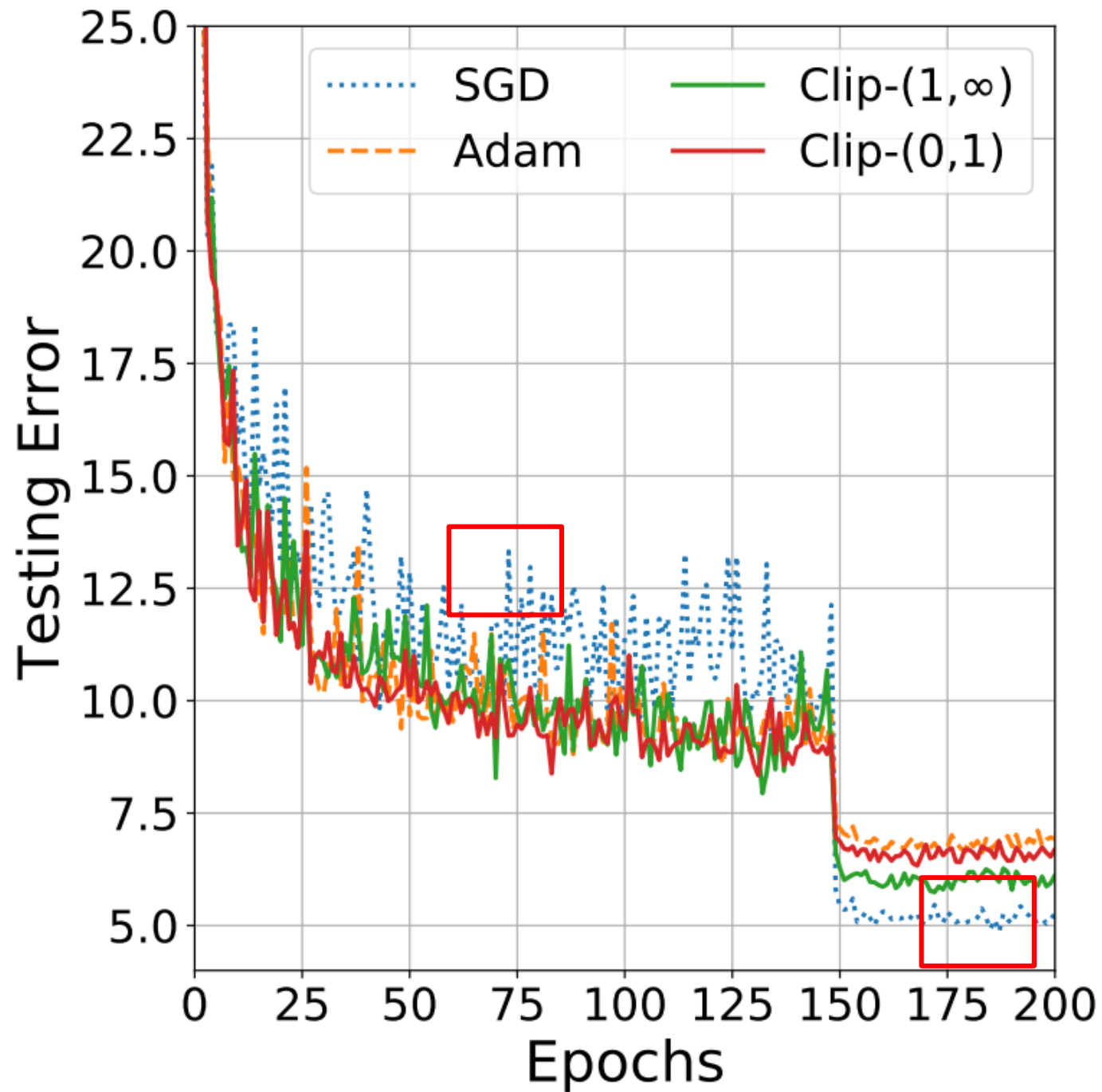
$$w_k = w_{k-1} - \alpha_{k-1} \cdot \frac{\sqrt{1 - \beta_2^k}}{1 - \beta_1^k} \cdot \frac{m_{k-1}}{\sqrt{v_{k-1} + \epsilon}},$$

**clip**( $x, a, b$ )

The function **clip**( $x, a, b$ ) clips the vector  $x$  element-wise such that the output is constrained to be in  $[a, b]$ .

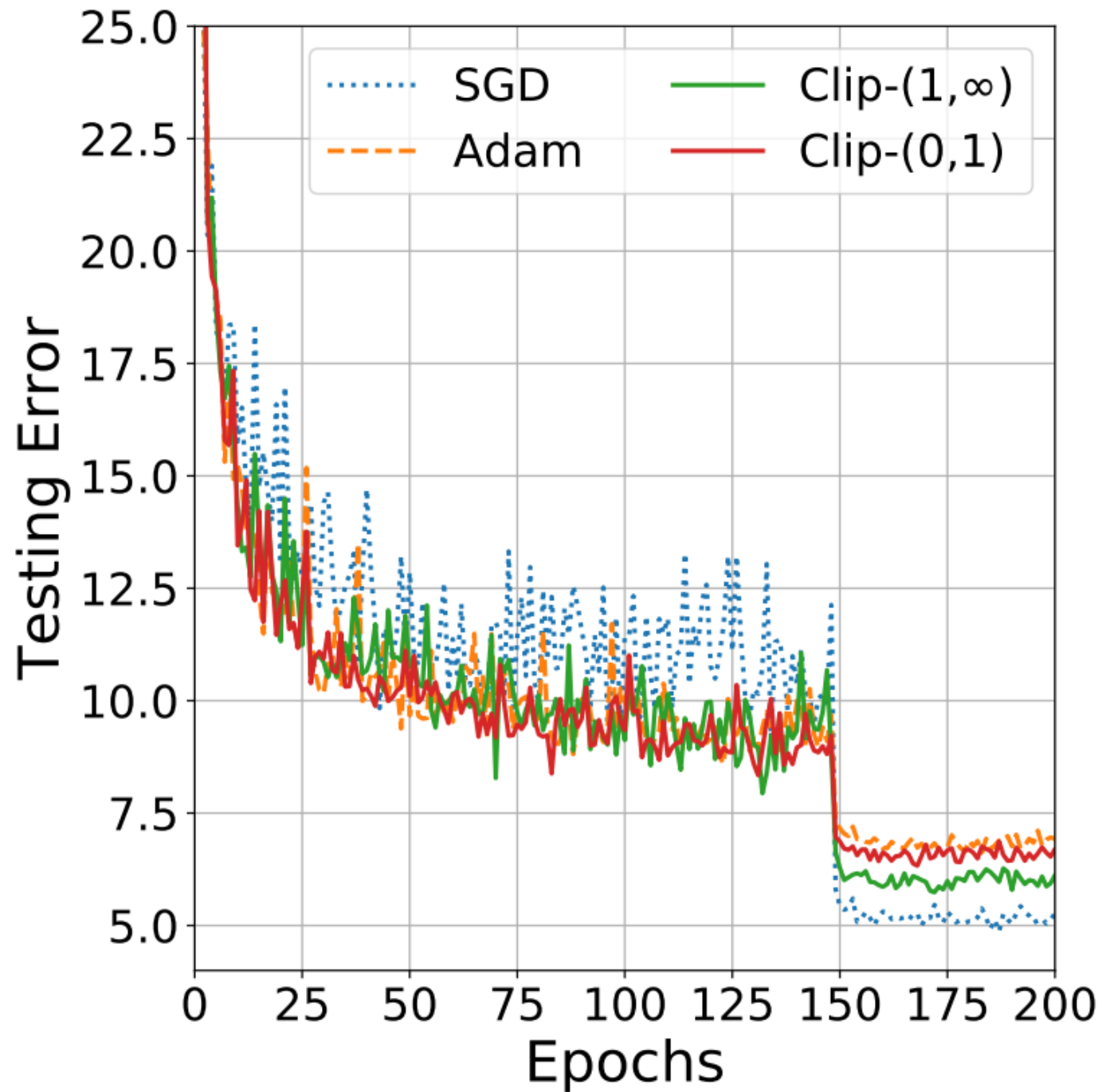
$$w_k = w_{k-1} -$$

$$\text{clip} \left( \frac{\sqrt{1 - \beta_2^k}}{1 - \beta_1^k} \frac{\alpha_{k-1}}{\sqrt{v_{k-1} + \epsilon}}, p \cdot \alpha_{sgd}, q \cdot \alpha_{sgd} \right) m_{k-1}.$$



Training the **DenseNet** architecture on the CIFAR-10 data set with four optimizers

1. Adam's testing error decrease faster than SGD
2. SGD converges to the expected testing error of  $\approx 5\%$  while Adam stagnates in performance at around  $\approx 7\%$  error.
3. Adam(with clip)'s performance are better than Adam(default)



Training the **DenseNet** architecture on the CIFAR-10 data set with four optimizers

1. Adam's testing error decrease faster than SGD
2. SGD converges to the expected testing error of  $\approx 5\%$  while Adam stagnates in performance at around  $\approx 7\%$  error.
3. Adam(with clip)'s performance are better than Adam(default)

# SWATS: Adam+SGD

- Why not combine Adam and SGD?
  - Use Adam first
  - Then switch to SGD
- Questions:
  - **When to switch Adam to SGD?**
  - **What learning rate to choose for SGD after the switch?**
- Answer:
  - **SWATS**

# SWATS: Adam+SGD

## Algorithm 1 SWATS

**Inputs:** Objective function  $f$ , initial point  $w_0$ , learning rate  $\alpha = 10^{-3}$ , accumulator coefficients  $(\beta_1, \beta_2) = (0.9, 0.999)$ ,  $\epsilon = 10^{-9}$ , phase=Adam.

```
1: Initialize  $k \leftarrow 0, m_k \leftarrow 0, a_k \leftarrow 0, \lambda_k \leftarrow 0$ 
2: while stopping criterion not met do
3:    $k = k + 1$ 
4:   Compute stochastic gradient  $g_k = \hat{\nabla} f(w_{k-1})$ 
5:   if phase = SGD then
6:      $v_k = \beta_1 v_{k-1} + g_k$ 
7:      $w_k = w_{k-1} - (1 - \beta_1) \Lambda v_k$ 
8:     continue
9:   end if
10:   $m_k = \beta_1 m_{k-1} + (1 - \beta_1) g_k$ 
11:   $a_k = \beta_2 a_{k-1} + (1 - \beta_2) g_k^2$ 
12:   $p_k = -\alpha_k \frac{\sqrt{1 - \beta_2^k}}{1 - \beta_1^k} \frac{m_k}{\sqrt{a_k + \epsilon}}$ 
13:   $w_k = w_k + p_k$ 
14:  if  $p_k^T g_k \neq 0$  then
15:     $\gamma_k = \frac{p_k^T p_k}{-p_k^T g_k}$ 
16:     $\lambda_k = \beta_2 \lambda_{k-1} + (1 - \beta_2) \gamma_k$ 
17:    if  $k > 1$  and  $|\frac{\lambda_k}{(1 - \beta_2^k)} - \gamma_k| < \epsilon$  then
18:      phase = SGD
19:       $v_k = 0$ 
20:       $\Lambda = \lambda_k / (1 - \beta_2^k)$ 
21:    end if
22:  else
23:     $\lambda_k = \lambda_{k-1}$ 
24:  end if
25: end while
return  $w_k$ 
```

# SWATS: Adam+SGD

- What learning rate to choose for SGD after the switch?

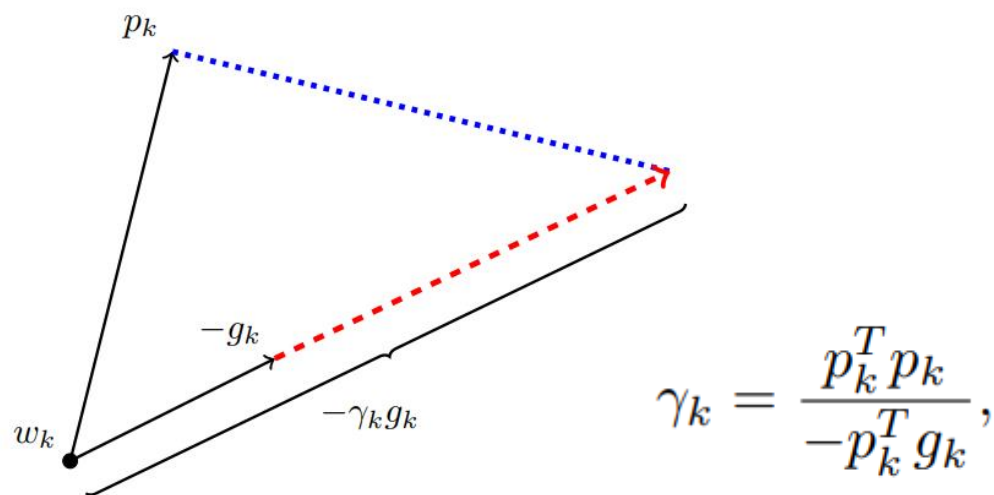


Figure 3. Illustrating the learning rate for SGD ( $\gamma_k$ ) estimated by our proposed projection given an iterate  $w_k$ , a stochastic gradient  $g_k$  and the Adam step  $p_k$ .

## Algorithm 1 SWATS

**Inputs:** Objective function  $f$ , initial point  $w_0$ , learning rate  $\alpha = 10^{-3}$ , accumulator coefficients  $(\beta_1, \beta_2) = (0.9, 0.999)$ ,  $\epsilon = 10^{-9}$ , phase=Adam.

```

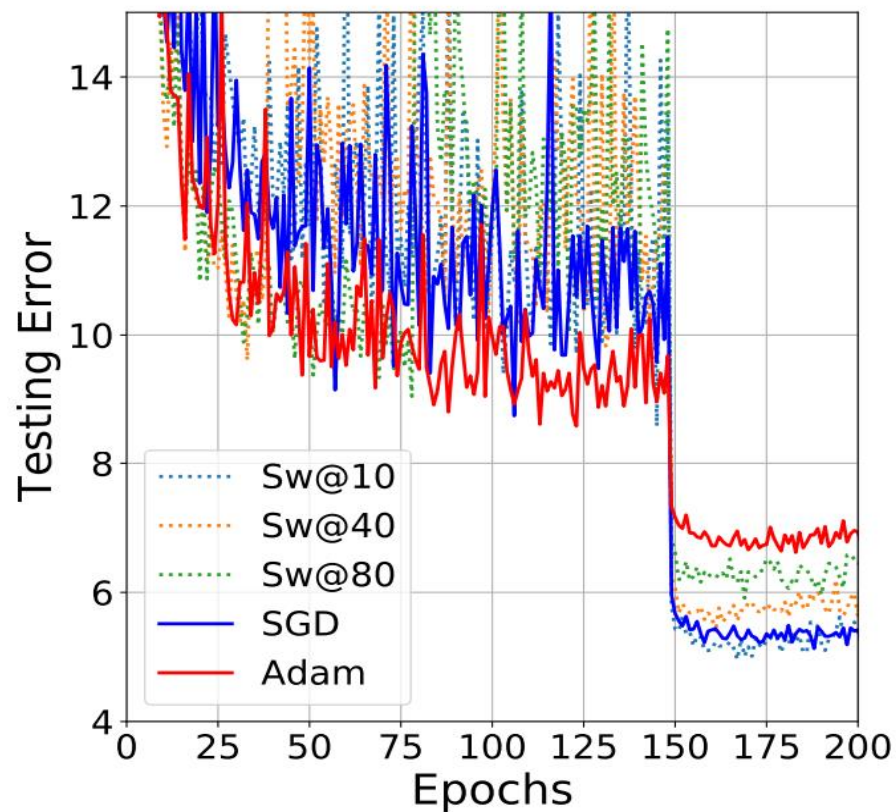
1: Initialize  $k \leftarrow 0, m_k \leftarrow 0, a_k \leftarrow 0, \lambda_k \leftarrow 0$ 
2: while stopping criterion not met do
3:    $k = k + 1$ 
4:   Compute stochastic gradient  $g_k = \hat{\nabla} f(w_{k-1})$ 
5:   if phase = SGD then
6:      $v_k = \beta_1 v_{k-1} + g_k$ 
7:      $w_k = w_{k-1} - (1 - \beta_1) \Lambda v_k$ 
8:     continue
9:   end if
10:   $m_k = \beta_1 m_{k-1} + (1 - \beta_1) g_k$ 
11:   $a_k = \beta_2 a_{k-1} + (1 - \beta_2) g_k^2$ 
12:   $p_k = -\alpha_k \frac{\sqrt{1 - \beta_2^k}}{1 - \beta_1^k} \frac{m_k}{\sqrt{a_k + \epsilon}}$ 
13:   $w_k = w_k + p_k$ 
14:  if  $p_k^T a_k \neq 0$  then
15:     $\gamma_k = \frac{p_k^T p_k}{-p_k^T g_k}$ 
16:     $\lambda_k = \beta_2 \lambda_{k-1} + (1 - \beta_2) \gamma_k$ 
17:    if  $k > 1$  and  $|\frac{\lambda_k}{(1 - \beta_2^k)} - \gamma_k| < \epsilon$  then
18:      phase = SGD
19:       $v_k = 0$ 
20:       $\Lambda = \lambda_k / (1 - \beta_2^k)$ 
21:    end if
22:  else
23:     $\lambda_k = \lambda_{k-1}$ 
24:  end if
25: end while

return  $w_k$ 

```

# SWATS: Adam+SGD

- When to switch Adam to SGD?



## Algorithm 1 SWATS

**Inputs:** Objective function  $f$ , initial point  $w_0$ , learning rate  $\alpha = 10^{-3}$ , accumulator coefficients  $(\beta_1, \beta_2) = (0.9, 0.999)$ ,  $\epsilon = 10^{-9}$ , phase=Adam.

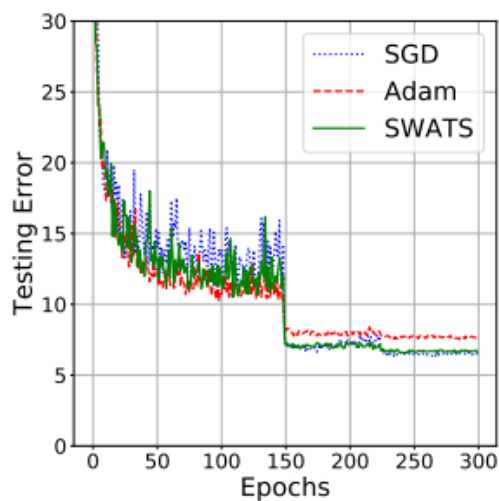
```

1: Initialize  $k \leftarrow 0, m_k \leftarrow 0, a_k \leftarrow 0, \lambda_k \leftarrow 0$ 
2: while stopping criterion not met do
3:    $k = k + 1$ 
4:   Compute stochastic gradient  $g_k = \hat{\nabla} f(w_{k-1})$ 
5:   if phase = SGD then
6:      $v_k = \beta_1 v_{k-1} + g_k$ 
7:      $w_k = w_{k-1} - (1 - \beta_1) \Lambda v_k$ 
8:     continue
9:   end if
10:   $m_k = \beta_1 m_{k-1} + (1 - \beta_1) g_k$ 
11:   $a_k = \beta_2 a_{k-1} + (1 - \beta_2) g_k^2$ 
12:   $p_k = -\alpha_k \frac{\sqrt{1 - \beta_2^k}}{1 - \beta_1^k} \frac{m_k}{\sqrt{a_k + \epsilon}}$ 
13:   $w_k = w_k + p_k$ 
14:  if  $p_k^T g_k \neq 0$  then
15:     $\gamma_k = \frac{p_k^T p_k}{-p_k^T g_k}$ 
16:     $\lambda_k = \beta_2 \lambda_{k-1} + (1 - \beta_2) \gamma_k$ 
17:    if  $k > 1$  and  $|\frac{\lambda_k}{(1 - \beta_2^k)} - \gamma_k| < \epsilon$  then
18:      phase = SGD
19:       $v_k = 0$ 
20:       $\Lambda = \lambda_k / (1 - \beta_2^k)$ 
21:    end if
22:  else
23:     $\lambda_k = \lambda_{k-1}$ 
24:  end if
25: end while

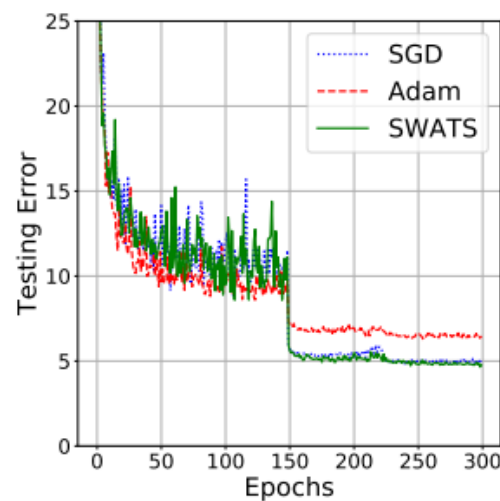
return  $w_k$ 

```

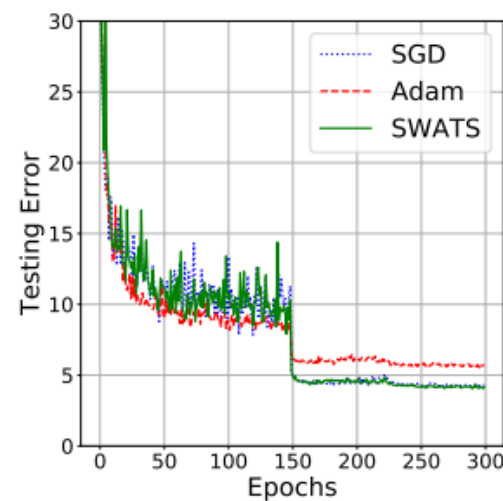




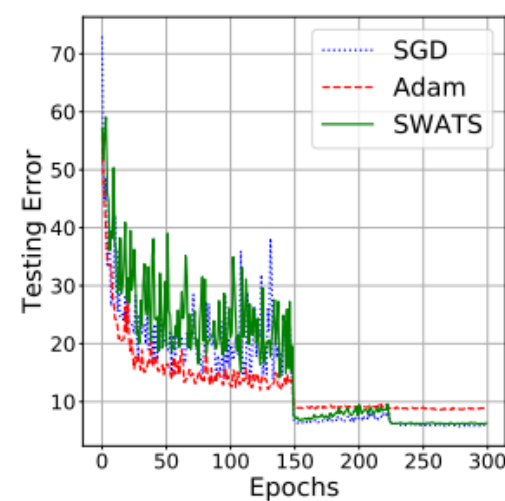
(a) ResNet-32 — CIFAR-10



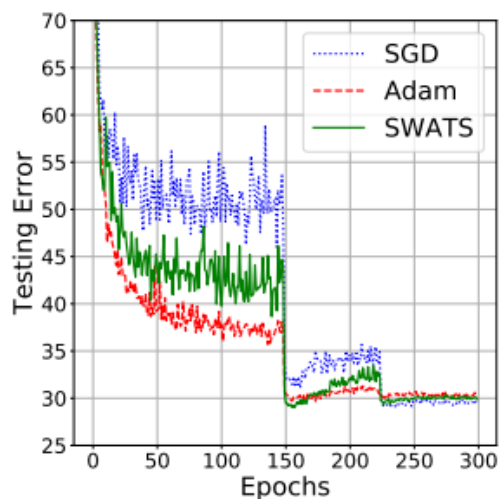
(b) DenseNet — CIFAR-10



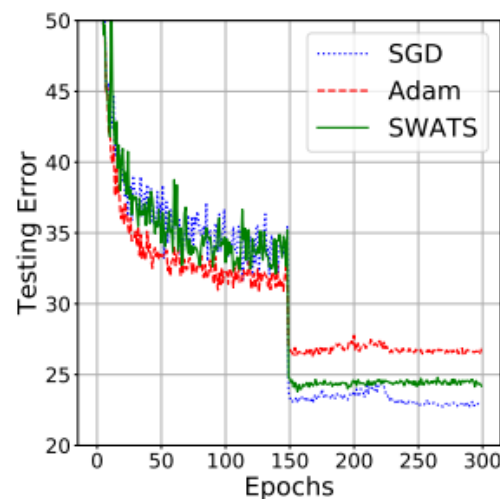
(c) PyramidNet — CIFAR-10



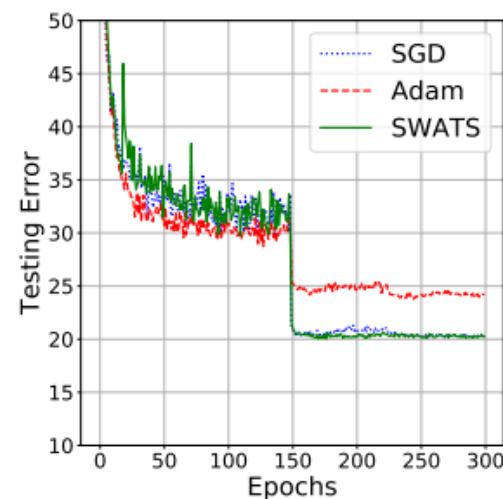
(d) SENet — CIFAR-10



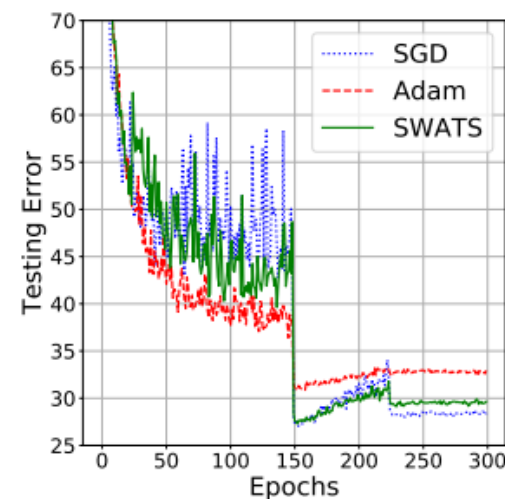
(e) ResNet-32 — CIFAR-100



(f) DenseNet — CIFAR-100



(g) PyramidNet — CIFAR-100



(h) SENet — CIFAR-100

Figure 4. Numerical experiments comparing SGD(M), Adam and SWATS with tuned learning rates on the ResNet-32, DenseNet, PyramidNet and SENet architectures on CIFAR-10 and CIFAR-100 data sets.



# Choice of optimizer

# Choice of optimizer

- The two recommended updates to use are either SGD+Nesterov Momentum or Adam
- Use Adam to check algorithm's correctness.
- Use SGD to training
  - Determine a good learning rate schedule
  - Adam may have lowest training error/loss, but not val.
- Use Adam first, then switch to SGD.(such as SWATS)
  - Consider when to switch and learning rate after switch