

DSCI 565: OPTIMIZATION ALGORITHMS

*This content is protected and may not
be shared, uploaded, or distributed.*

Ke-Thia Yao

Lecture 17: 2025-10-29

Sparse Features and Learning Rates

2

- In natural language, some words occur much more often than others, e.g., the words *learning* vs *preconditioning*
- A global learning rate schedule for all words may cause the rate to decrease too slowly for some words and too fast for other words
- How to keep track of which feature is frequent, and which is rare?
- Adagrad proposes maintaining a summary of past gradient variances for all parameters

Adagrad

3

- Use variable \underline{s}_t to accumulate past gradient variance
- $l(y_t, f(\mathbf{x}_t, \mathbf{w}))$ computes the loss wrt to sample (\mathbf{x}_t, y_t) and weight \mathbf{w}
- \mathbf{g}_t is the gradient of the loss
- Initially $\mathbf{s}_0 = \mathbf{0}$
- \mathbf{w}_t is the weight at time t
- For some fixed η , and fixed $\epsilon > 0$ to prevent division by zero

Se shape = weight shape

$$\begin{aligned}\mathbf{g}_t &= \partial_{\mathbf{w}} l(y_t, f(\mathbf{x}_t, \mathbf{w})), \\ \mathbf{s}_t &= \mathbf{s}_{t-1} + \mathbf{g}_t^2, \\ \mathbf{w}_t &= \mathbf{w}_{t-1} - \frac{\eta}{\sqrt{\mathbf{s}_t + \epsilon}} \cdot \mathbf{g}_t.\end{aligned}$$

gradient descent
 $\mathbf{w}_t = \mathbf{w}_{t-1} - \eta \mathbf{g}_t$
Adagrad

$$\mathbf{w}_t = \mathbf{w}_{t-1} - \frac{\eta}{\sqrt{\mathbf{s}_t + \epsilon}} \mathbf{g}_t$$

Adagrad Rationale

4

- Why $s_t = \underline{s_{t-1}} + \underline{g_t^2}$?

- Similar to analysis of Newton's method we want to precondition with the diagonal of the Hessian matrix
- For deep learning even computing the diagonal of second derivatives is difficult due to memory constraints
- But the diagonal elements does depend on the gradients

- Why $-\frac{\eta}{\sqrt{(s_t + \epsilon)}}$?

- For stochastic gradient descent, under convexity assumption the optimal rate schedule is $\mathcal{O}\left(\frac{1}{\sqrt{T}}\right)$
- See section 12.4.3
- Here instead of the number of steps T , we use $\underline{s_t}$

RMSProp

$$s_t = \gamma s_{t-1} + (1 - \gamma) g_t^2$$

5

- For Adagrad the state $s_t = s_{t-1} + \underline{g_t^2}$ can grow without bound
- This can cause adaptive learning rate $\frac{\eta}{\sqrt{(s_t + \epsilon)}}$ to decay too quickly
- RMSProp proposes using leaky averaging $\underline{\gamma}$

$$\begin{aligned}s_t &\leftarrow \underline{\gamma s_{t-1}} + \underline{(1 - \gamma) g_t^2}, \\ x_t &\leftarrow x_{t-1} - \frac{\eta}{\sqrt{s_t + \epsilon}} \odot g_t.\end{aligned}$$

- For some fixed η , and for small $\epsilon > 0$, typically $\epsilon = 10^{-6}$

Adadelta

6

- Enhances Adagrad by
 - Using leaky average for state variable s_t , similar to RMSProp
 - Decoupling per-coordinate scaling with learning rate adjustment
 - Which also removes the need for a global learning rate η
 - PyTorch implementation still has learning rate but defaults to $\eta = 1$

Adadelta

7

- Uses two state variables

- Variable s_t for maintaining the variance of the gradient (with leaky average)

$$s_t = \rho s_{t-1} + (1 - \rho) g_t^2.$$

- Variable Δx_t to store the variance of the parameter (with leaky average)
 - Removing the need to manually select a global rate η (no need)

$$\begin{aligned} x_t &= x_{t-1} - g'_t \\ g'_t &= \frac{\sqrt{\Delta x_{t-1} + \epsilon}}{\sqrt{s_t + \epsilon}} \odot g_t \\ \Delta x_t &= \rho \Delta x_{t-1} + (1 - \rho) g'^2_t \end{aligned}$$

- With $\Delta x_0 = 0$, and small $\epsilon > 0$

Adam

8

- Adam combines all techniques we have seen
 - Minibatch stochastic gradient descent
 - Momentum
 - Per-coordinate scaling (Adagrad, RMSProp)
 - Decoupling per-coordinate scaling from learning rate adjustment (Adadelta)

Adam

Momentum + Variance of the gradient

9

- Uses two state variables with leaky averages

- v_t for the momentum

- s_t for the variance of the gradient

- Initialize $v_0 = s_0 = 0$

- Typically, $\beta_1 = 0.9$ and $\beta_2 = 0.999$

- Variance estimate moves much more slowly than momentum

- At startup biased toward smaller values, rescale:

- Update equation

$$x_t \leftarrow x_{t-1} - g'_t$$

$$g'_t = \frac{\eta \hat{y}_t}{\sqrt{\hat{s}_t} + \epsilon}$$

$$\begin{aligned} v_t &\leftarrow \beta_1 v_{t-1} + (1 - \beta_1) g_t, \\ s_t &\leftarrow \beta_2 s_{t-1} + (1 - \beta_2) g_t^2. \end{aligned}$$

↳ $\hat{v}_t = \frac{v_t}{1 - \beta_1^t}$ and $\hat{s}_t = \frac{s_t}{1 - \beta_2^t}$

AdamW

10

- Adam with decoupled weight decay

$$\mathbf{x}_t \leftarrow \mathbf{x}_{t-1} - \eta \left(\frac{\hat{\mathbf{v}}_t}{\sqrt{\hat{s}_t} + \epsilon} + \boxed{\lambda \mathbf{x}_{t-1}} \right)$$

where $\lambda \in \mathbb{R}$

- Currently, the most popular for training big models
- Loshchilov and Hutter, 2019

Learning Rate Scheduling

11

- So far focused on algorithms for updating weight vectors rather than the **rate** at which they are updated
- Aspect of **rate scheduling**
 - Magnitude
 - Decay
 - Initialization

Toy Problem

12

□ Convolutional network with cross entropy loss on FashionMNIST data

```
def net_fn():
    model = nn.Sequential(
        nn.Conv2d(1, 6, kernel_size=5, padding=2), nn.ReLU(),
        nn.MaxPool2d(kernel_size=2, stride=2),
        nn.Conv2d(6, 16, kernel_size=5), nn.ReLU(),
        nn.MaxPool2d(kernel_size=2, stride=2),
        nn.Flatten(),
        nn.Linear(16 * 5 * 5, 120), nn.ReLU(),
        nn.Linear(120, 84), nn.ReLU(),
        nn.Linear(84, 10))

    return model

loss = nn.CrossEntropyLoss()
```

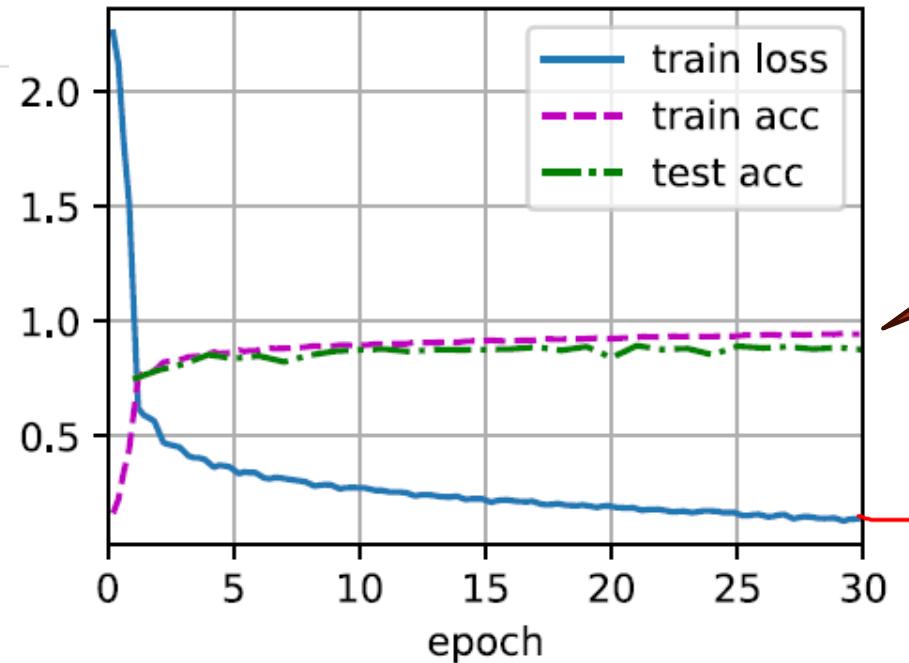
Baseline with No Schedule

13

```
lr, num_epochs = 0.3, 30  
net = net_fn()  
trainer = torch.optim.SGD(net.parameters(), lr=lr)  
train(net, train_iter, test_iter, num_epochs, loss, trainer, device)
```

Constant learning rate

train loss 0.145, train acc 0.944, test acc 0.877

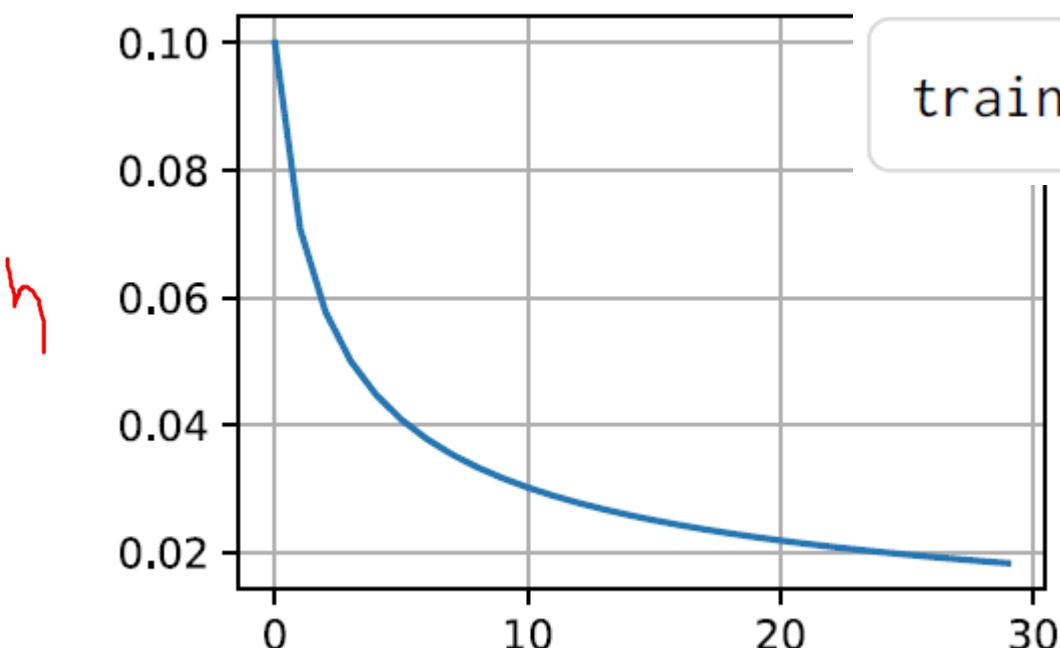


Gap between
train and test

Square Root Scheduler

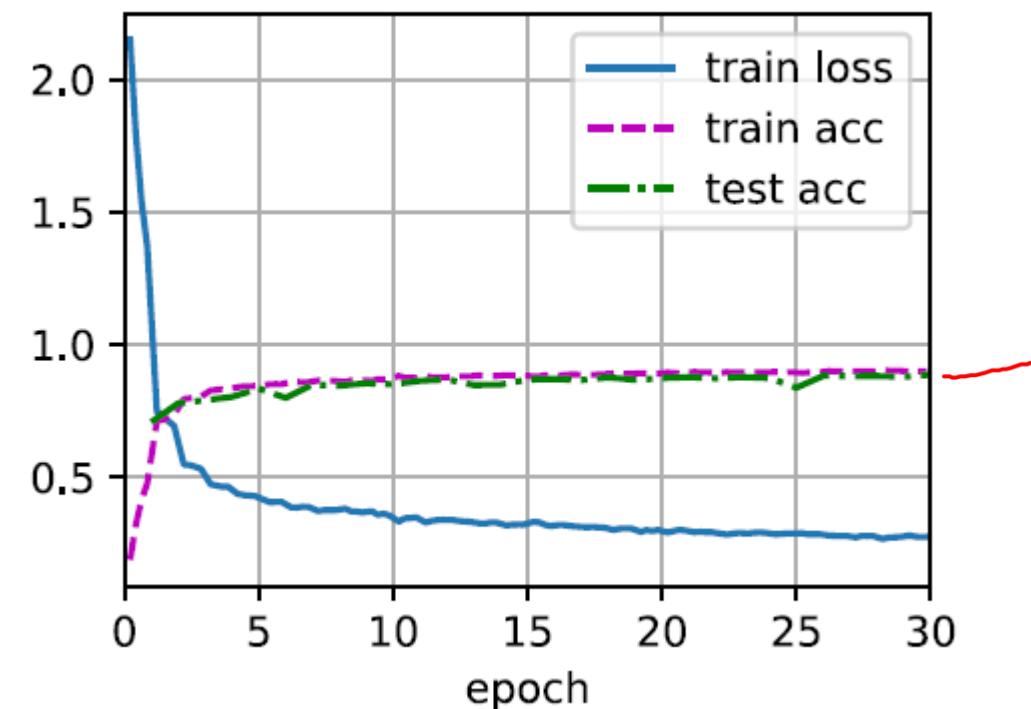
14

- ◻ $\eta = \eta_0(t + 1)^{-\frac{1}{2}}$ with $\eta_0 = 0.1$



train loss 0.273, train acc 0.900, test acc 0.886

- ◻ Good for convex problems
- ◻ Slower decrease for nonconvex

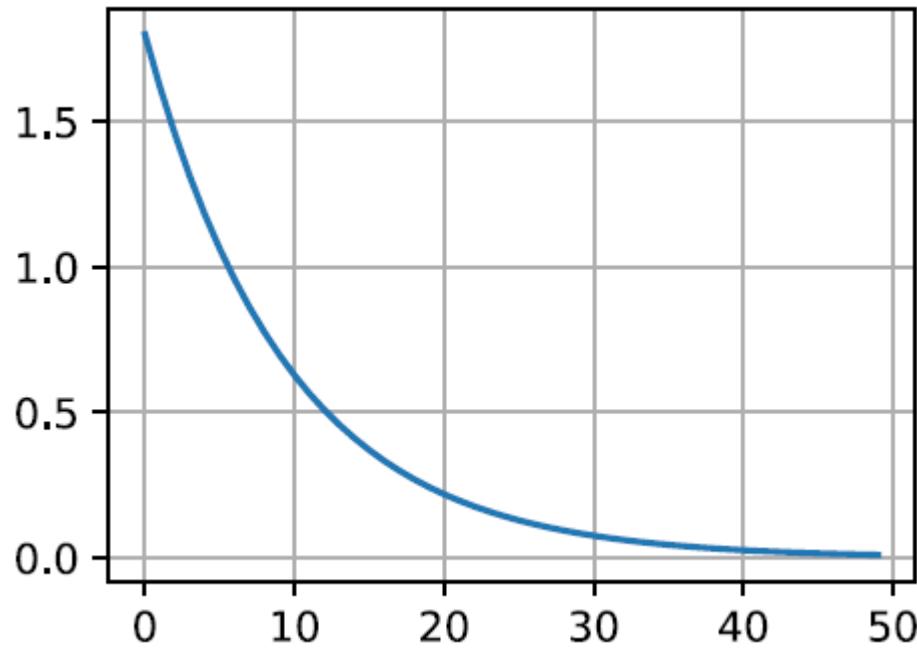


Factor Scheduler

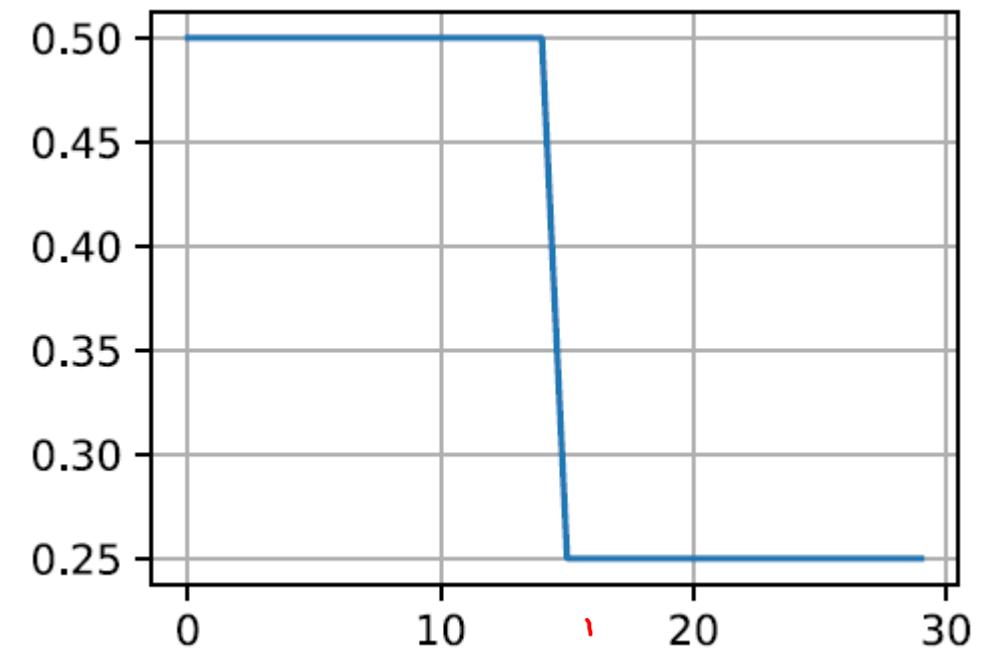
15

- Factor scheduler

$$\eta_{t+1} \leftarrow \alpha \eta_t \text{ for } \alpha \in (0,1)$$



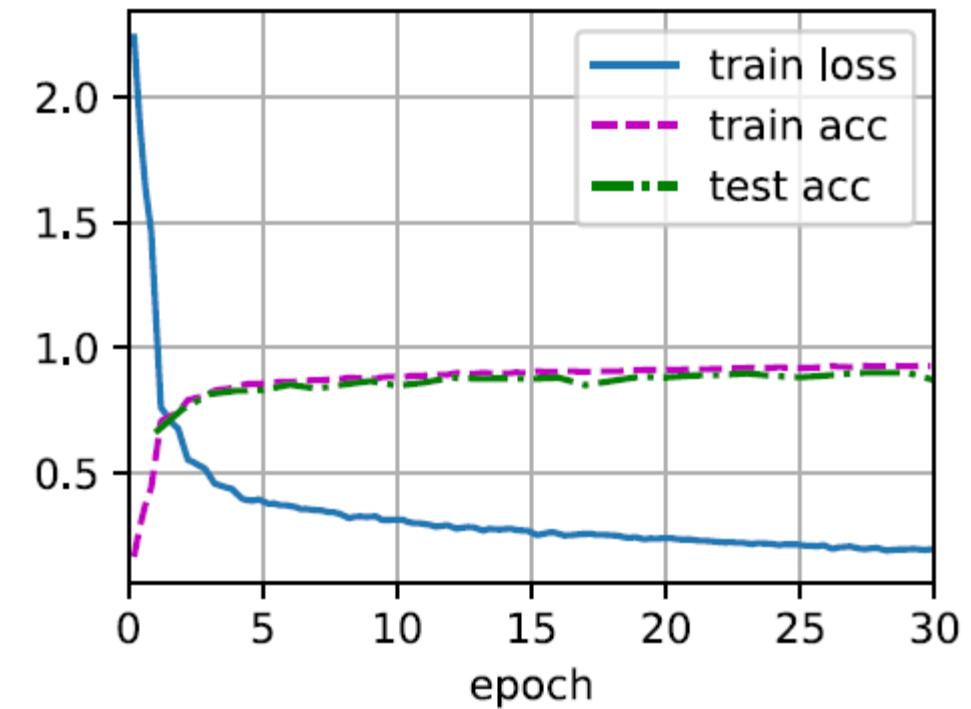
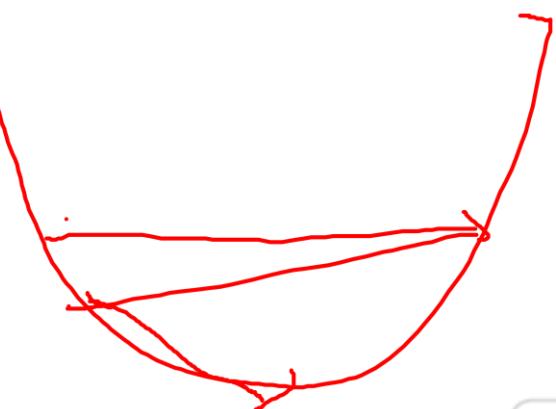
- ↑ multi step
- Multifactor scheduler
 - Apply factor scheduler update only at certain time steps $s = \{15, 30, 45\}$



Multifactor Scheduler

16

- Intuition for piecewise constant learning rate
 - Optimize until near local optima
 - Then decrease rate to get even closer



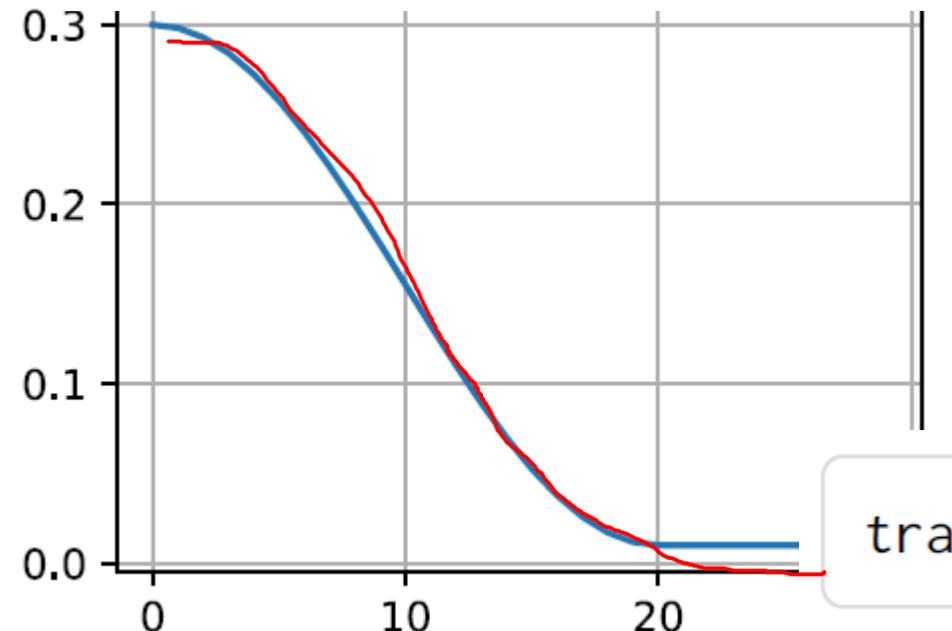
train loss 0.194, train acc 0.927, test acc 0.869

Cosine Scheduler

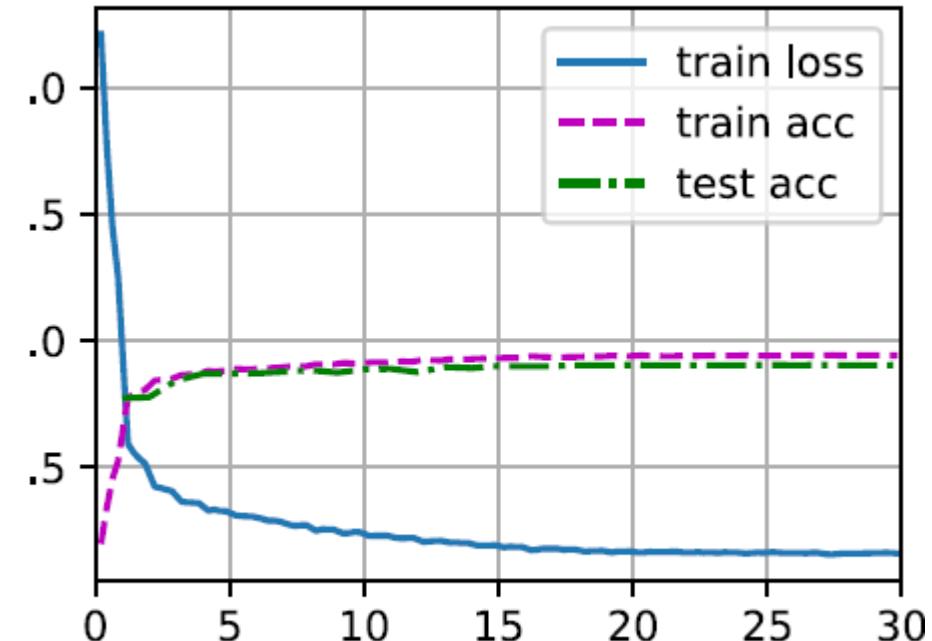
17

- Larger rates in the beginning
- Very small rates to refine solution
- Popular for computer vision

$$\eta_t = \eta_T + \frac{\eta_0 - \eta_T}{2} (1 + \cos(\pi t/T))$$



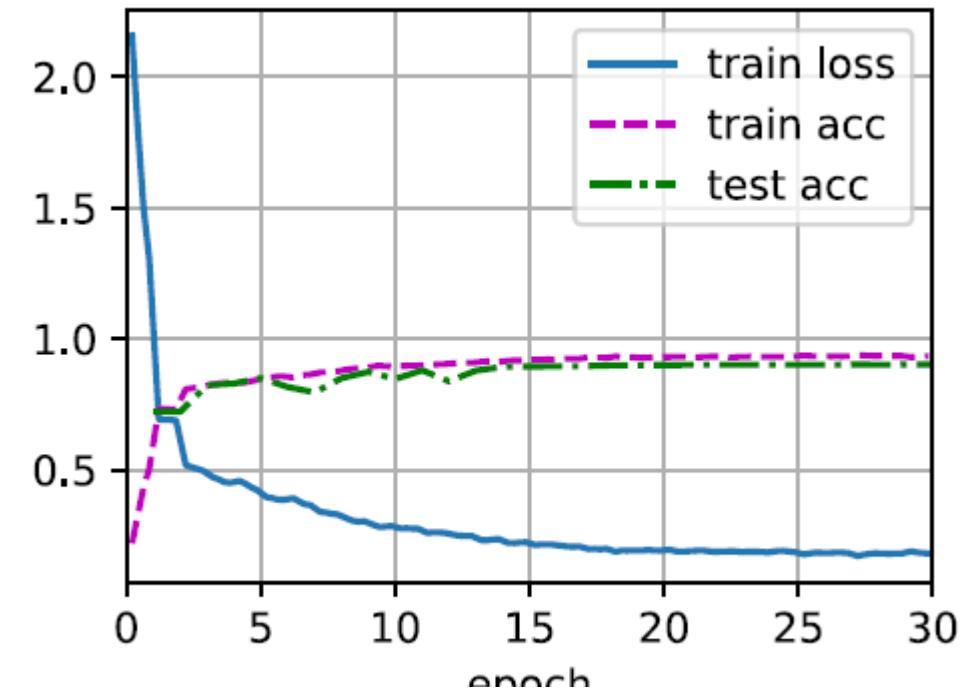
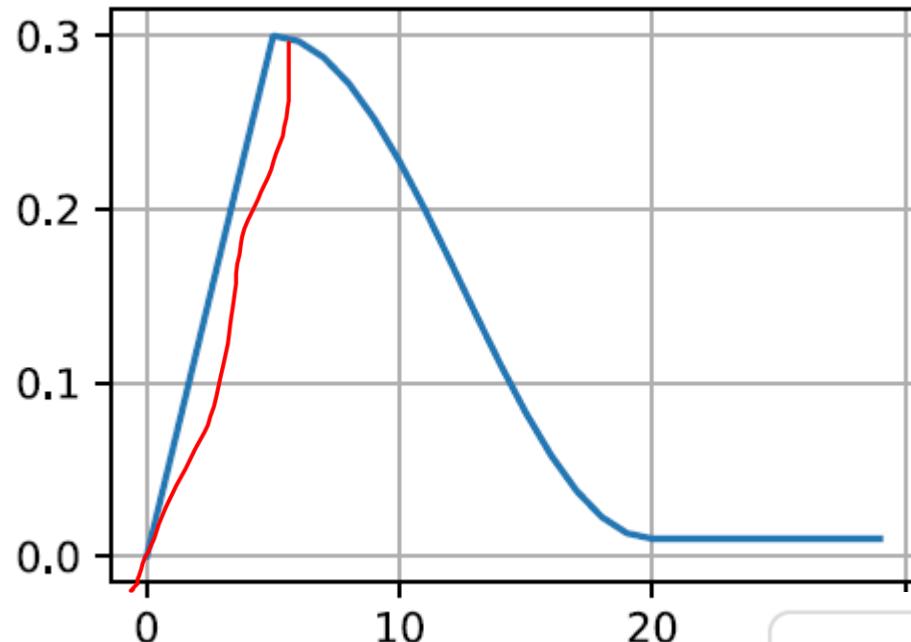
train loss 0.159, train acc 0.942, test acc 0.904



Warmup

18

- Sometimes large rates at startup cause optimization to diverge
- Start with small rates then switch to cosine



train loss 0.181, train acc 0.934, test acc 0.901