

Optimizers

Review. ★ Backpropagation

- Output activation and loss function

Task	# of output nodes	Output activation	Loss function
Binary classification	1	Sigmoid $\hat{y} = 1/(1 + \exp(-z))$	$y \in \{0, 1\}$ Binary cross-entropy $L(y, \hat{y}) = -y \log \hat{y} - (1 - y) \log(1 - \hat{y})$
Multiclass classification	k	Softmax $\hat{y}_i = \frac{\exp\{z_i\}}{\sum_{i=1}^k \exp\{z_i\}}$	Cross-entropy $y = [0, 0, \dots, 1, 0, \dots, 0]$ $L(y, \hat{y}) = - \sum_{i=1}^k y_i \log(\hat{y}_i)$
Regression	k	None $\hat{y}_i = z_i$	Mean squared error $L(y, \hat{y}) = \frac{1}{2} \sum_{i=1}^k (\hat{y}_i - y_i)^2$

Review. Backpropagation

- Refer to [Chap. 2 in Nielsen's Neural Networks and Deep Learning](#)

- Notations: $w_{jk}^l, b_j^l, a_j^l, z_j^l, \delta_j^l = \frac{\delta(C=\text{Loss})}{\delta z_j^l}$
bias act pre-act

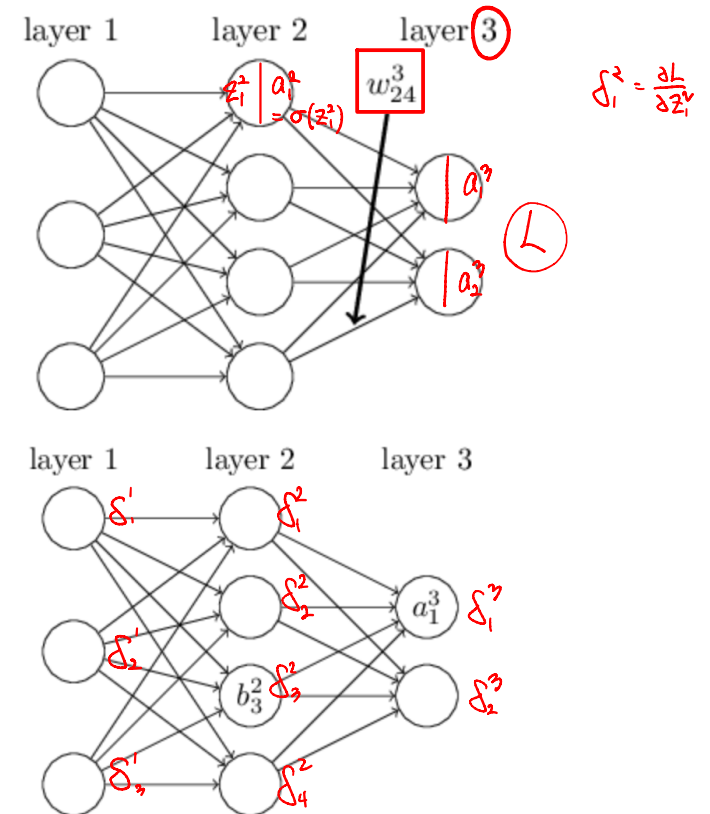
*l for layer
j for destination
k for source*

- The Hadamard product $\odot \rightarrow$ element-wise multiplication.
- Four fundamental equations behind backpropagation

L for last layer

- (BP1) $\delta^L = \nabla_a C \odot \sigma'(z^L)$
- (BP2) $\delta^l = (w^{l+1})^T \delta^{l+1} \odot \sigma'(z^l)$
- (BP3) $\frac{\partial C}{\partial b_j^l} = \delta_j^l$
- (BP4) $\frac{\partial C}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l$

$\delta_i^L = \frac{\partial C}{\partial z_i^L}$ $z_i^L \xrightarrow{\sigma} a_i^L \rightarrow C$
 $\delta_i^l = \frac{\partial C}{\partial z_i^l}$



Review. Backpropagation

$$= \frac{\partial C}{\partial z_i^L}, \quad z_i^L \rightarrow \hat{y}_i \rightarrow C \quad (i \in [1, k])$$

- (BP1) $\delta^L = \nabla_a C \odot \sigma'(z^L)$ for multiclass classification

$$\delta_j^L = \frac{\partial C}{\partial z_j^L} = \sum_{i=1}^{\text{\# output } k} \frac{\partial C}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial z_j^L} = \hat{y}_j - y_j$$

$$\frac{\partial C}{\partial \hat{y}_i} = -\frac{y_i}{\hat{y}_i} \qquad \frac{\partial \hat{y}_i}{\partial z_j^L} = \hat{y}_i (I_{ij} - \hat{y}_j)$$

3-class classification. $y = [1, 0, 0]$

$\odot \hat{y}_1$	0.3	$\delta_1^L = 0.3 - 1$
$\odot \hat{y}_2$	0.3	$\delta_2^L = 0.3 - 0$
$\odot \hat{y}_3$	0.4	$\delta_3^L = 0.4 - 0$

$I = 1$

Review. Backpropagation

- Refer to [Chap. 2 in Nielsen's Neural Networks and Deep Learning](#)

1. **Input x :** Set the corresponding activation a^1 for the input layer.

2. **Feedforward:** For each $l = 2, 3, \dots, L$ compute $z^l = w^l a^{l-1} + b^l$ and $a^l = \sigma(z^l)$.

3. **Output error δ^L :** Compute the vector $\delta^L = \nabla_a C \odot \sigma'(z^L)$.

4. **Backpropagate the error:** For each $l = L - 1, L - 2, \dots, 2$ compute $\delta^l = ((w^{l+1})^T \delta^{l+1}) \odot \sigma'(z^l)$.

5. **Output:** The gradient of the cost function is given by

$$\frac{\partial C}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l \text{ and } \frac{\partial C}{\partial b_j^l} = \delta_j^l.$$

backpropagation



Basic Gradient Descent Algorithms

- Batch gradient descent

Algorithm 1 Batch Gradient Descent at Iteration k

Require: Learning rate ϵ_k

Require: Initial Parameter θ

- 1: **while** stopping criteria not met **do**
 - 2: Compute gradient estimate over N examples:
 - 3: $\hat{\mathbf{g}} \leftarrow +\frac{1}{N} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$
 - 4: Apply Update: $\theta \leftarrow \theta - \epsilon \hat{\mathbf{g}}$
 - 5: **end while**
-

Basic Gradient Descent Algorithms

- Stochastic gradient descent

Algorithm 8.1 Stochastic gradient descent (SGD) update

Require: Learning rate schedule $\epsilon_1, \epsilon_2, \dots$

Require: Initial parameter θ

$k \leftarrow 1$

while stopping criterion not met **do**

 Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$.

 Compute gradient estimate: $\hat{\mathbf{g}} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

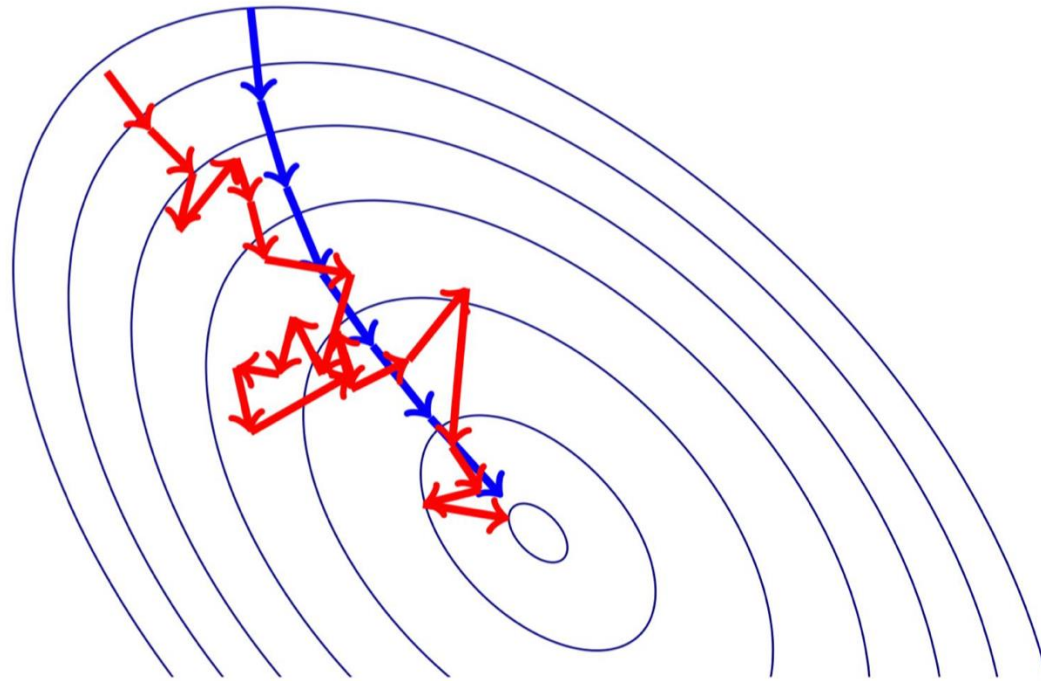
 Apply update: $\theta \leftarrow \theta - \epsilon_k \hat{\mathbf{g}}$

$k \leftarrow k + 1$

end while

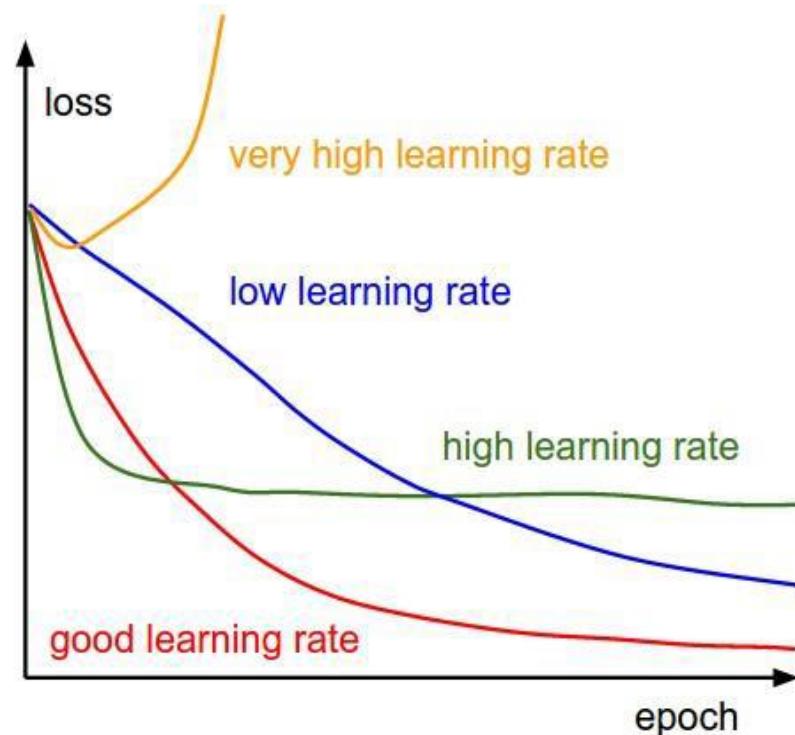
Basic Gradient Descent Algorithms

- Batch gradient descent VS Stochastic gradient descent



Basic Gradient Descent Algorithms

- How to determine the learning rate?
 - If it is too large, the objective value may explode.
 - If it is too small, a lot of iterations are needed.



Basic Gradient Descent Algorithms

- Momentum: accumulates an exponentially decaying moving average of past gradients and continues to move in that direction

previous direction
(accumulation of past gradients)

$$\mathbf{v} \leftarrow \alpha \mathbf{v} - \epsilon \nabla_{\theta} \left(\frac{1}{m} \sum_{i=1}^m L(\mathbf{f}(x^{(i)}; \theta), y^{(i)}) \right)$$

current direction

current gradient

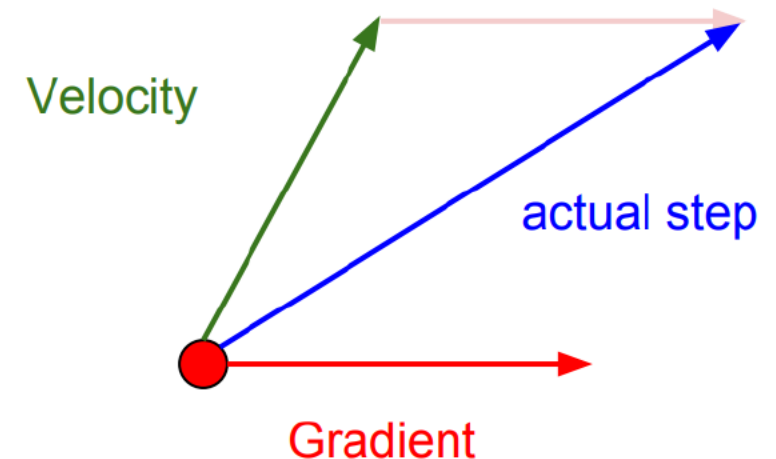


Figure Credit: Prof. Kang (SKKU)

Basic Gradient Descent Algorithms

- Stochastic gradient descent with momentum

Algorithm 8.2 Stochastic gradient descent (SGD) with momentum

Require: Learning rate ϵ , momentum parameter α

Require: Initial parameter θ , initial velocity v

while stopping criterion not met **do**

 Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$.

 Compute gradient estimate: $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$.

 Compute velocity update: $\mathbf{v} \leftarrow \alpha \mathbf{v} - \epsilon \mathbf{g}$.

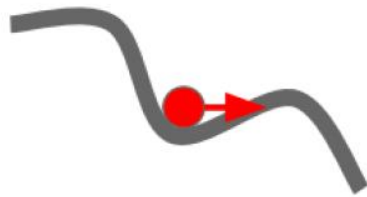
 Apply update: $\theta \leftarrow \theta + \mathbf{v}$.

end while

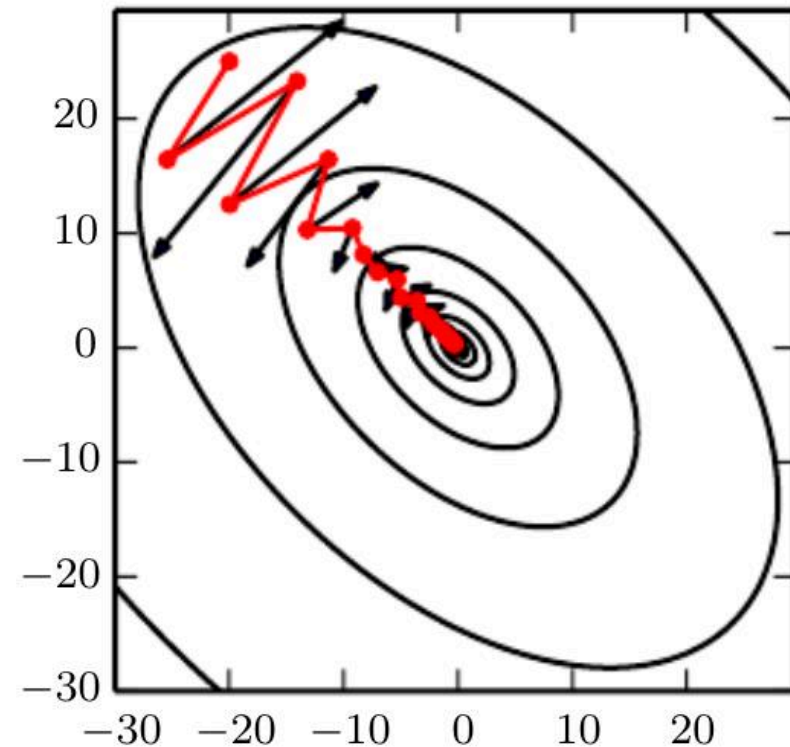
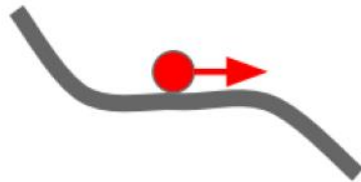
Basic Gradient Descent Algorithms

- Stochastic gradient descent with momentum

Local Minima

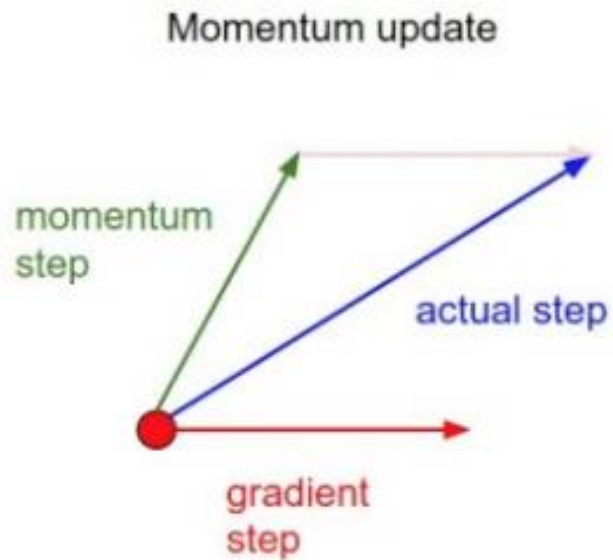


Saddle points



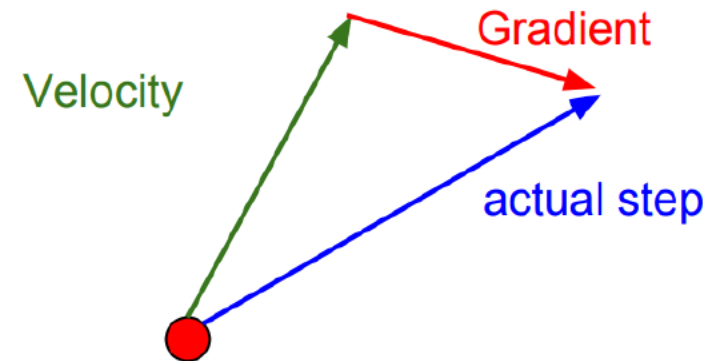
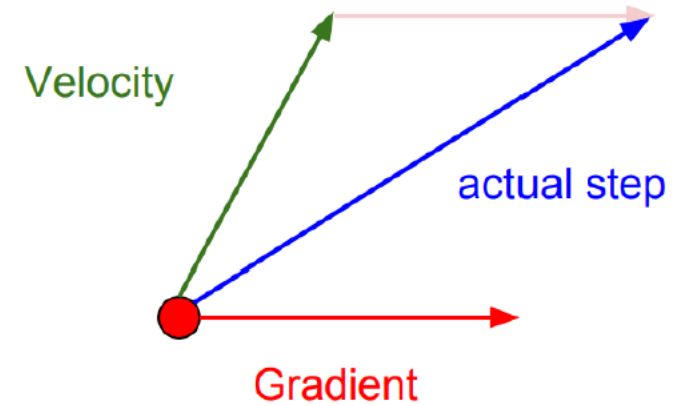
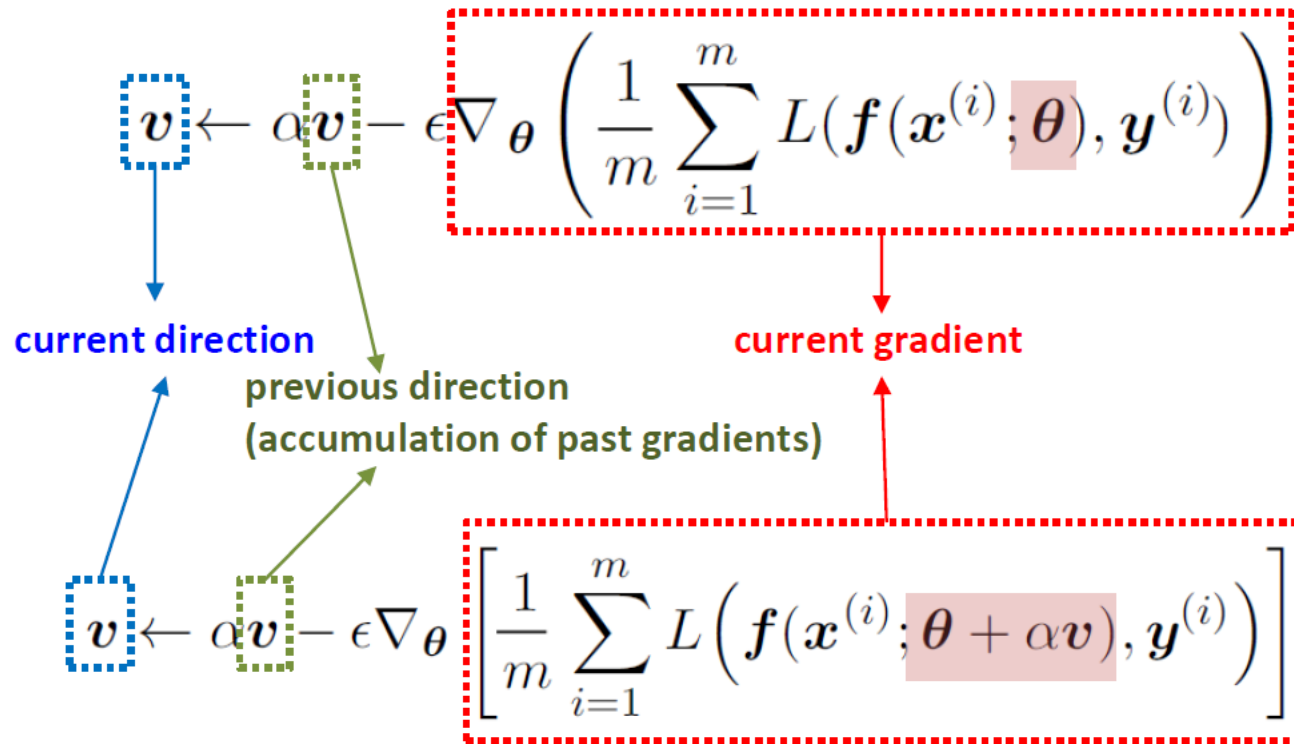
Basic Gradient Descent Algorithms

- Nesterov Momentum: the gradient is evaluated after the current velocity is applied.



Basic Gradient Descent Algorithms

Momentum



Nesterov Momentum

Figure Credit: Prof. Kang (SKKU)

Basic Gradient Descent Algorithms

- Stochastic gradient descent with Nesterov's momentum

Algorithm 8.3 Stochastic gradient descent (SGD) with Nesterov momentum

Require: Learning rate ϵ , momentum parameter α

Require: Initial parameter θ , initial velocity v

while stopping criterion not met **do**

 Sample a minibatch of m examples from the training set $\{x^{(1)}, \dots, x^{(m)}\}$ with corresponding labels $y^{(i)}$.

 Apply interim update: $\tilde{\theta} \leftarrow \theta + \alpha v$.

 Compute gradient (at interim point): $g \leftarrow \frac{1}{m} \nabla_{\tilde{\theta}} \sum_i L(f(x^{(i)}; \tilde{\theta}), y^{(i)})$.

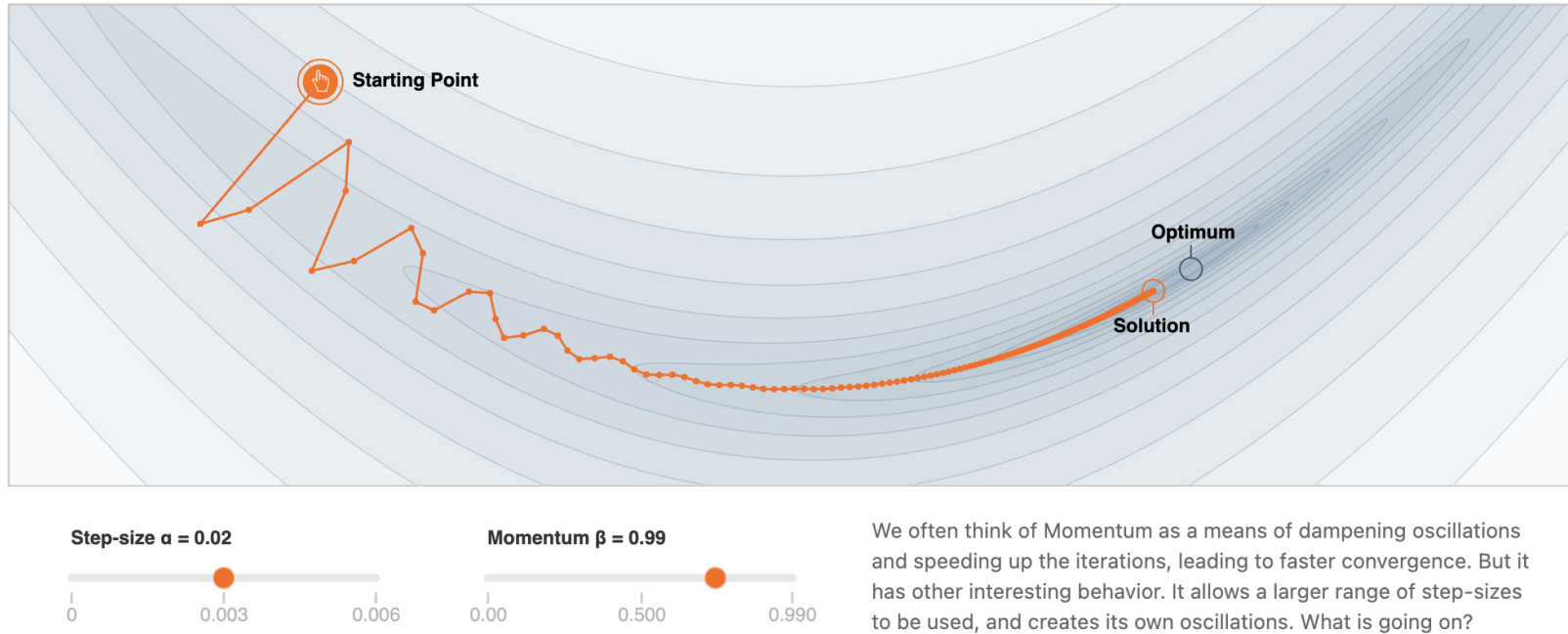
 Compute velocity update: $v \leftarrow \alpha v - \epsilon g$.

 Apply update: $\theta \leftarrow \theta + v$.

end while

Basic Gradient Descent Algorithms

Why Momentum Really Works



<https://distill.pub/2017/momentum/>

Algorithms with Adaptive Learning Rates

- The learning rate significantly affects the performance.
- How to determine the learning rate?
- Adaptive learning rate
 - To use a separate learning rate for each parameter and automatically adapt these learning rates throughout the course of learning

Algorithms with Adaptive Learning Rates

- AdaGrad (2011)
 - Large gradients \rightarrow small update, small gradients \rightarrow large update
 - Code snippet

$$\text{cache} = \begin{bmatrix} 100 \\ 1 \end{bmatrix}$$

```
# Assume the gradient dx and parameter vector x
cache += dx**2
x += - learning_rate * dx / (np.sqrt(cache) + eps)
```

- Cache keeps growing during training \rightarrow learning rate shrinks too quickly

$$0.1 / \begin{bmatrix} 100 \\ 1 \end{bmatrix} = \begin{bmatrix} 0.001 \\ 0.1 \end{bmatrix}$$

Algorithms with Adaptive Learning Rates

- AdaGrad (2011)

Algorithm 8.4 The AdaGrad algorithm

Require: Global learning rate ϵ

Require: Initial parameter θ

Require: Small constant δ , perhaps 10^{-7} , for numerical stability

Initialize gradient accumulation variable $\mathbf{r} = \mathbf{0}$

while stopping criterion not met **do**

Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$.

Compute gradient: $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

Accumulate squared gradient: $\mathbf{r} \leftarrow \mathbf{r} + \mathbf{g} \odot \mathbf{g}$

Compute update: $\Delta \theta \leftarrow -\frac{\epsilon}{\delta + \sqrt{\mathbf{r}}} \odot \mathbf{g}$. (Division and square root applied element-wise)

Apply update: $\theta \leftarrow \theta + \Delta \theta$

end while

Algorithms with Adaptive Learning Rates

- RMSProp (2012)
 - Extension of AdaGrad which solves AdaGrad's aggressive, monotonically decreasing learning rate problem
 - The gradient accumulation → exponentially weighted moving average

$$r \leftarrow r + g \odot g \quad \blacktriangleright \quad r \leftarrow \rho r + (1 - \rho)g \odot g$$

```
cache = decay_rate * cache + (1 - decay_rate) * dx**2  
x += - learning_rate * dx / (np.sqrt(cache) + eps)
```

Algorithms with Adaptive Learning Rates

- RMSProp (2012)

Algorithm 8.5 The RMSProp algorithm

Require: Global learning rate ϵ , decay rate ρ .

Require: Initial parameter θ

Require: Small constant δ , usually 10^{-6} , used to stabilize division by small numbers.

Initialize accumulation variables $\mathbf{r} = 0$

while stopping criterion not met **do**

 Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$.

 Compute gradient: $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

 Accumulate squared gradient: $\mathbf{r} \leftarrow \rho \mathbf{r} + (1 - \rho) \mathbf{g} \odot \mathbf{g}$

 Compute parameter update: $\Delta \theta = -\frac{\epsilon}{\sqrt{\delta + \mathbf{r}}} \odot \mathbf{g}$. ($\frac{1}{\sqrt{\delta + \mathbf{r}}}$ applied element-wise)

 Apply update: $\theta \leftarrow \theta + \Delta \theta$

end while

Algorithms with Adaptive Learning Rates

- Adam (2015) *popular.*
 - Adaptive moment estimation
 - It uses the first moment and the second moment.
 - A variant on the combination of RMSProp and momentum

```
m = beta1*m + (1-beta1)*dx
v = beta2*v + (1-beta2)*(dx**2)
x += - learning_rate * m / (np.sqrt(v) + eps)
```

Algorithms with Adaptive Learning Rates

- Adam (2015)

1st moment
⇒ velocity.
momentum.
2nd moment
⇒ RMSProp.
adaptive LR.

Algorithm 8.7 The Adam algorithm

Require: Step size ϵ (Suggested default: 0.001)

Require: Exponential decay rates for moment estimates, ρ_1 and ρ_2 in $[0, 1)$.
(Suggested defaults: 0.9 and 0.999 respectively)

Require: Small constant δ used for numerical stabilization. (Suggested default: 10^{-8})

Require: Initial parameters θ

Initialize 1st and 2nd moment variables $\mathbf{s} = \mathbf{0}$, $\mathbf{r} = \mathbf{0}$

Initialize time step $t = 0$

while stopping criterion not met **do**

Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$.

Compute gradient: $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

$t \leftarrow t + 1$

Update biased first moment estimate: $\mathbf{s} \leftarrow \rho_1 \mathbf{s} + (1 - \rho_1) \mathbf{g}$

Update biased second moment estimate: $\mathbf{r} \leftarrow \rho_2 \mathbf{r} + (1 - \rho_2) \mathbf{g} \odot \mathbf{g}$

Correct bias in first moment: $\hat{\mathbf{s}} \leftarrow \frac{\mathbf{s}}{1 - \rho_1^t}$

Correct bias in second moment: $\hat{\mathbf{r}} \leftarrow \frac{\mathbf{r}}{1 - \rho_2^t}$

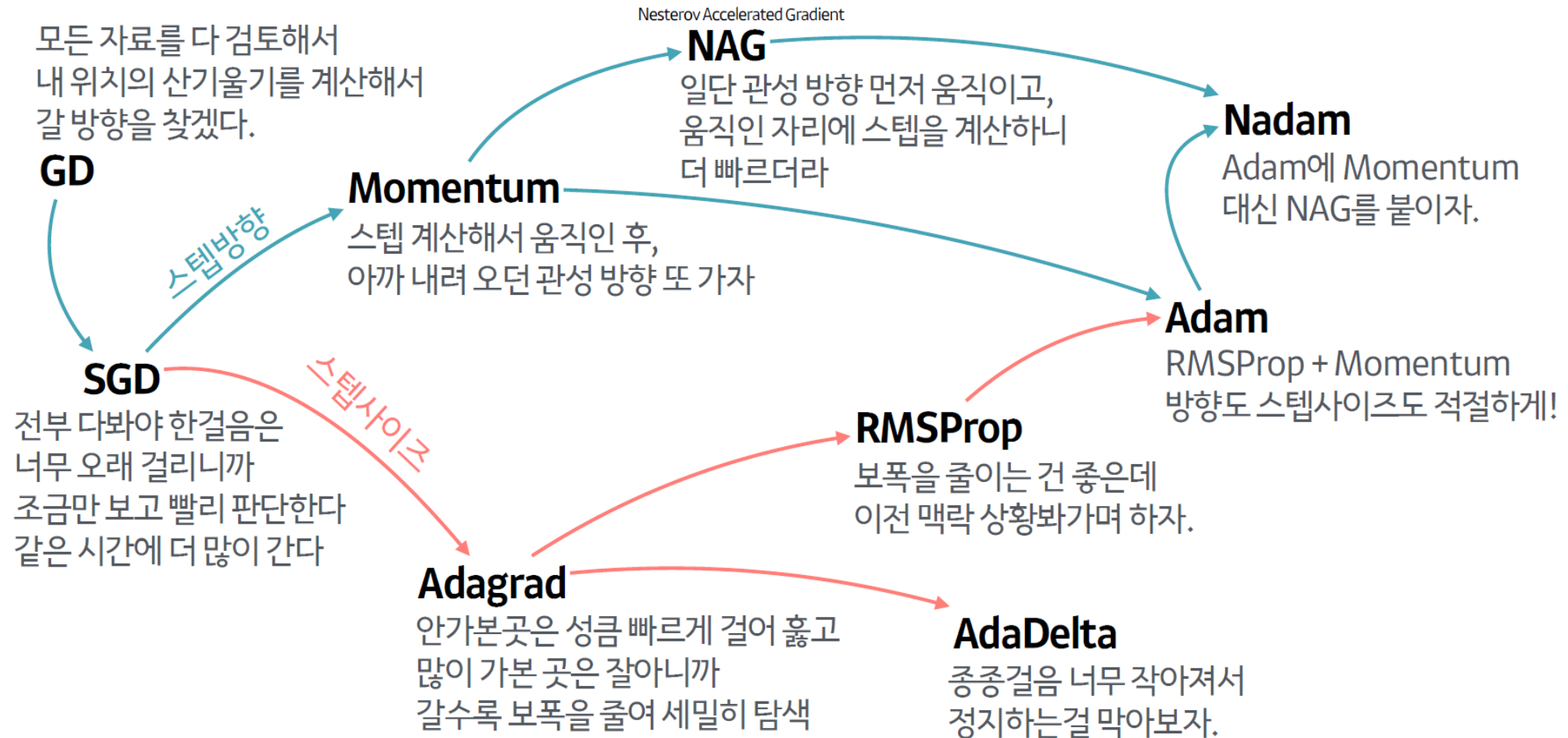
$\rho_1 = \rho_2 = 0.99$

Compute update: $\Delta \theta = -\epsilon \frac{\hat{\mathbf{s}}}{\sqrt{\hat{\mathbf{r}} + \delta}}$ (operations applied element-wise)

Apply update: $\theta \leftarrow \theta + \Delta \theta$

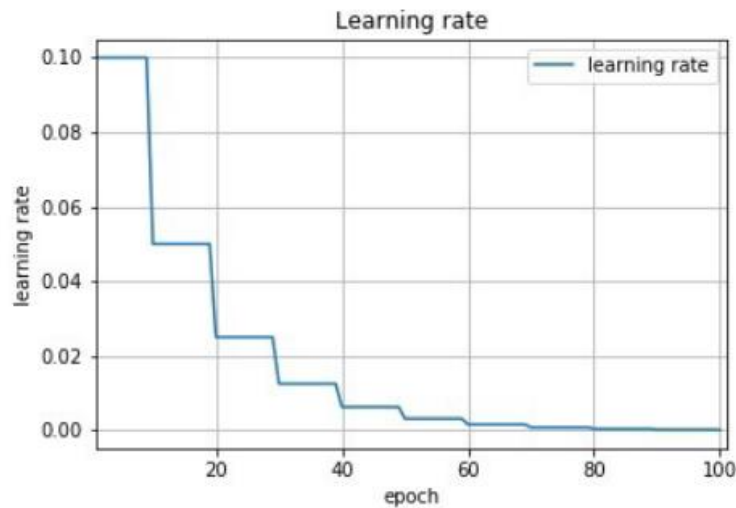
end while

And Many More

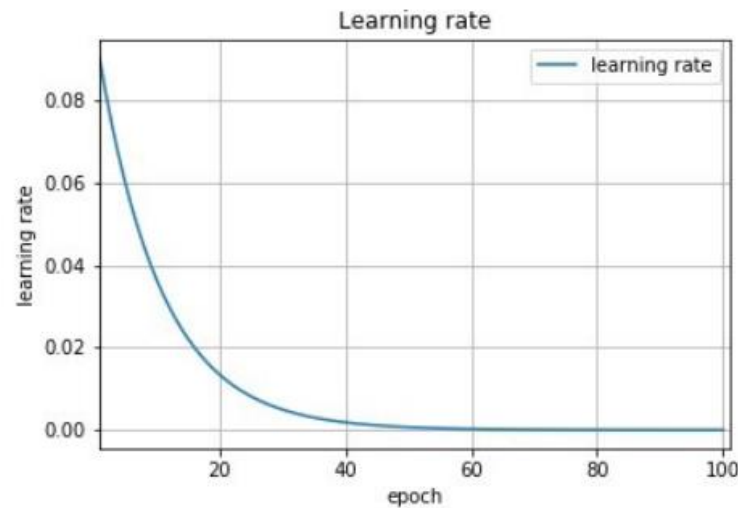


Learning Rate Scheduling

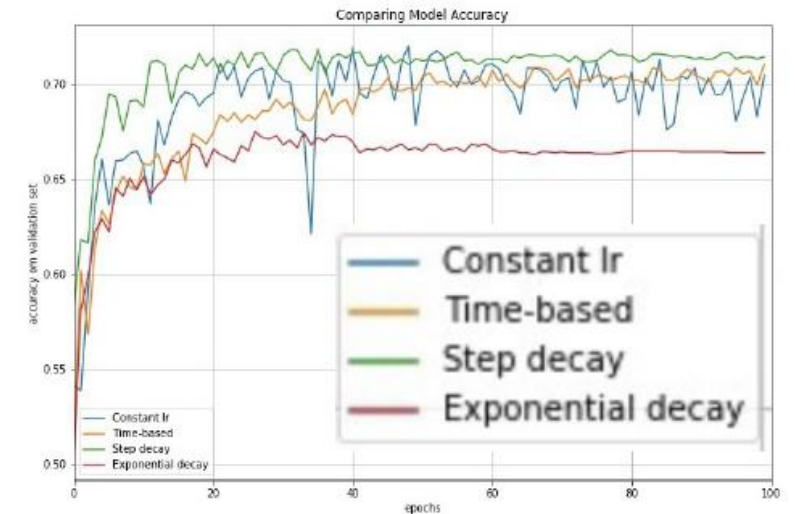
- Learning rate decaying
 - Constant learning rate often prevents convergence.



Step decay

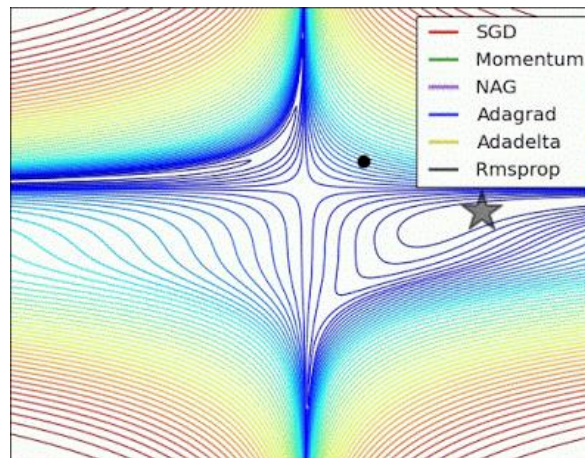
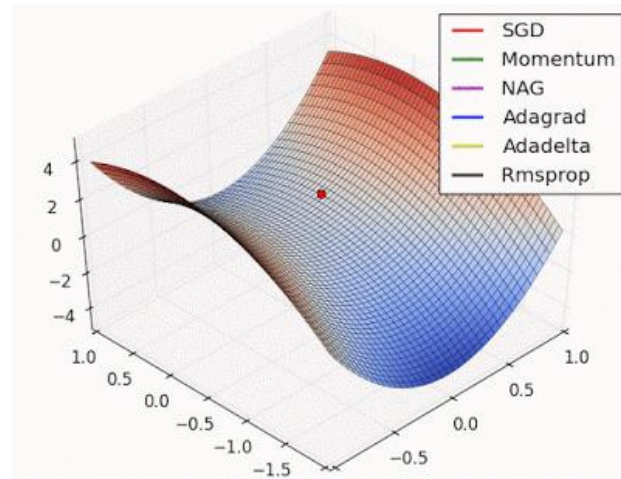
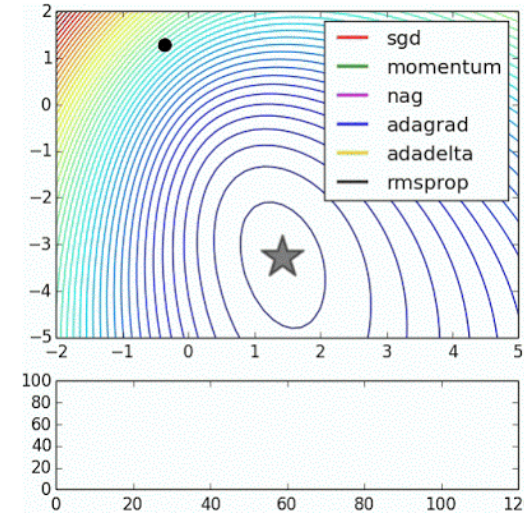
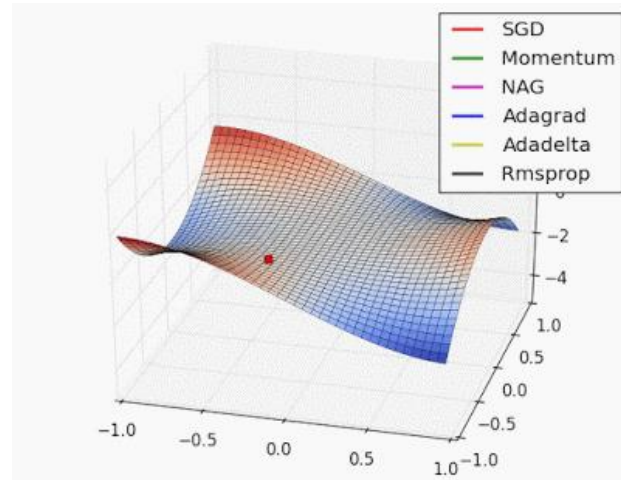


Exponential decay



Test accuracy

Animations for Optimization Algorithms



Reading assignments

- “Understanding deep learning”
 - Chapter 6
- “Dive into deep learning”
 - Section 12