# Optimizers

# Review. Backpropagation
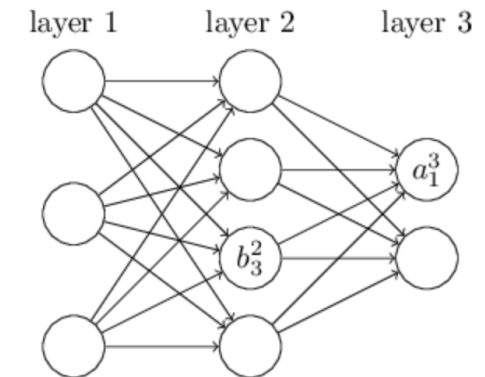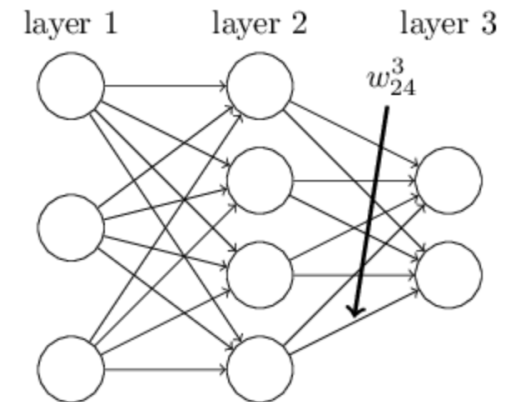
- Output activation and loss function

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

Sangheum Hwang

# Review. Backpropagation

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

  - Notations: $w_{jk}^l, b_j^l, a_j^l, z_j^l, \delta_j^l = \frac{\delta C}{\delta z_j^l}$

  - The Hadamard product $\odot$
  - Four fundamental equations behind backpropagation
    - (BP1) $\delta^L = \nabla_a C \odot \sigma'(z^L)$
    - (BP2) $\delta^l = \left(w^{l+1}\right)^T \delta^{l+1} \odot \sigma'(z^l)$
    - (BP3) $\frac{\partial C}{\partial b_j^l} = \delta_j^l$
    - (BP4) $\frac{\partial C}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l$

# Review. Backpropagation

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

$$\frac{\delta C}{\delta z_j^L} = \sum_{i=1}^{k} \frac{\delta C}{\delta \hat{y}_i} \frac{\delta \hat{y}_i}{\delta z_j} = \hat{y}_j - y_j$$

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

# 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 $l = 2, 3, \ldots, L$ compute
   $z^l = w^l a^{l-1} + b^l$ and $a^l = \sigma(z^l)$.

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

4. **Backpropagate the error:** For each $l = L-1, L-2, \ldots, 2$ 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
   $\frac{\partial C}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l$ and $\frac{\partial C}{\partial b_j^l} = \delta_j^l$.

# Basic Gradient Descent Algorithms

- Batch gradient descent

**Algorithm 1** Batch Gradient Descent at Iteration $k$

**Require:** Learning rate $\epsilon_k$

**Require:** Initial Parameter $\theta$

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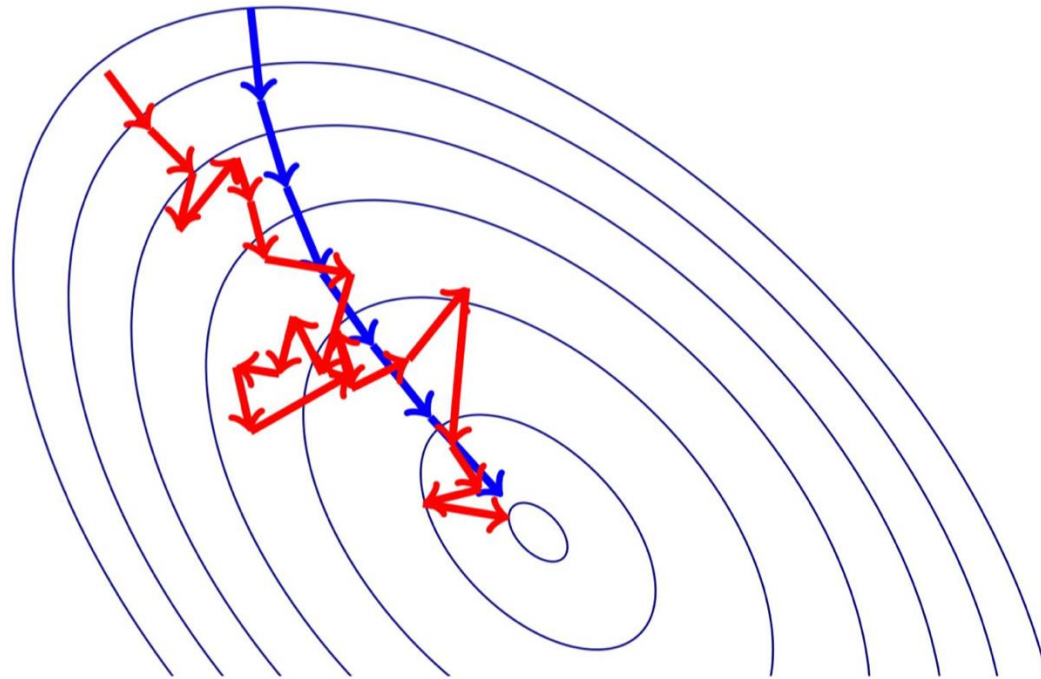3:     $\hat{\mathbf{g}} \leftarrow +\frac{1}{N}\nabla_\theta \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$
4:     Apply Update: $\theta \leftarrow \theta - \epsilon\hat{\mathbf{g}}$
5: **end while**

# Basic Gradient Descent Algorithms

- Stochastic gradient descent

---

**Algorithm 8.1** Stochastic gradient descent (SGD) update

---

**Require:** Learning rate schedule $\epsilon_1, \epsilon_2, \ldots$

**Require:** Initial parameter $\boldsymbol{\theta}$

    $k \leftarrow 1$

    **while** stopping criterion not met **do**

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

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

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

        $k \leftarrow k + 1$
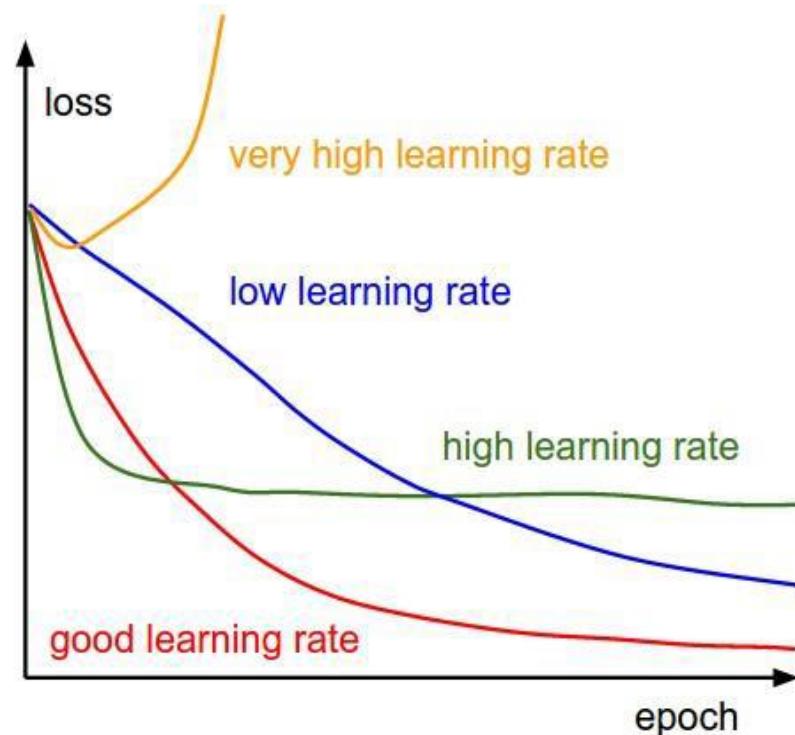
    **end while**

---

# 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)

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

current direction

current gradient

Velocity
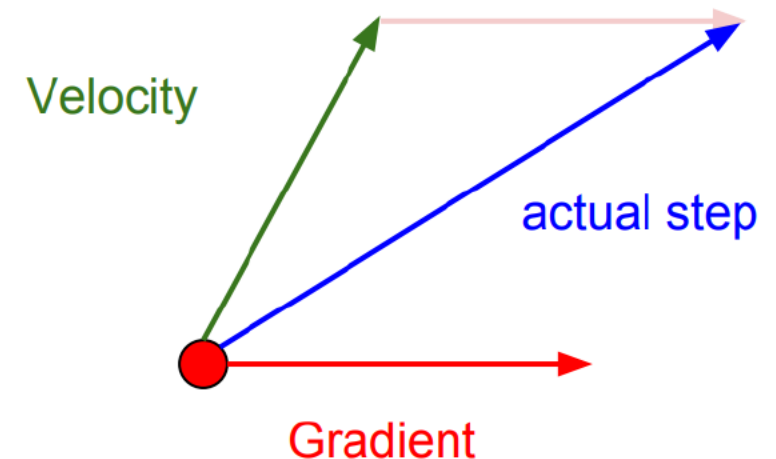
actual step

Gradient

Figure Credit: Prof. Kang (SKKU)

서울과학기술대학교

# Basic Gradient Descent Algorithms

- Stochastic gradient descent with momentum

---

**Algorithm 8.2** Stochastic gradient descent (SGD) with momentum

---

**Require:** Learning rate $\epsilon$, momentum parameter $\alpha$

**Require:** Initial parameter $\boldsymbol{\theta}$, initial velocity $\boldsymbol{v}$

    **while** stopping criterion not met **do**

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

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

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

        Apply update: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \boldsymbol{v}$.
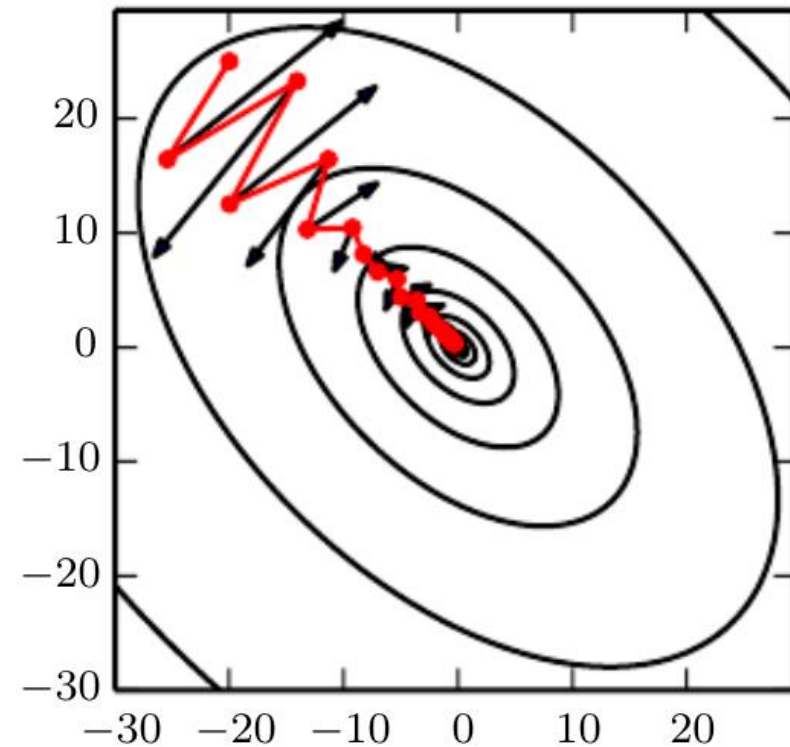
    **end while**

---

# Basic Gradient Descent Algorithms

- Stochastic gradient descent with momentum
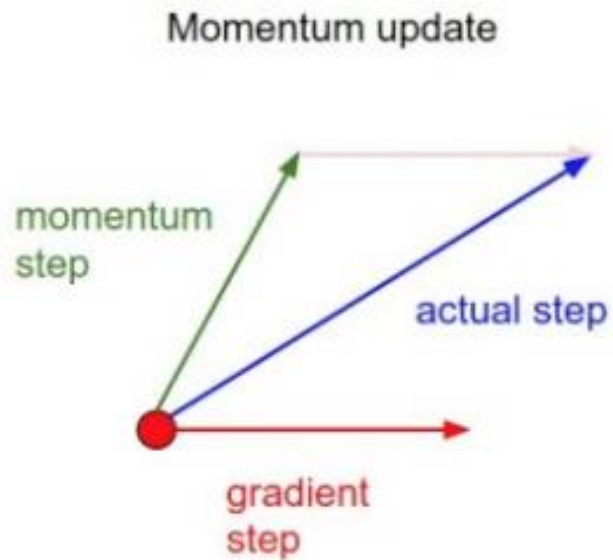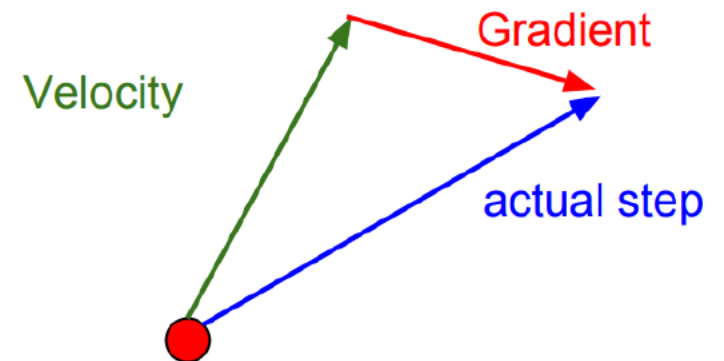
# Basic Gradient Descent Algorithms

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

# Basic Gradient Descent Algorithms

**Momentum**

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

**current direction**

**previous direction**
**(accumulation of past gradients)**

**current gradient**

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

**Nesterov Momentum**

Figure Credit: Prof. Kang (SKKU)

Velocity

actual step

Gradient

Gradient

Velocity

actual step

...g...um Hwang

# Basic Gradient Descent Algorithms

- Stochastic gradient descent with Nesterov's momentum

**Algorithm 8.3** Stochastic gradient descent (SGD) with Nesterov momentum

**Require:** Learning rate $\epsilon$, momentum parameter $\alpha$
**Require:** Initial parameter $\boldsymbol{\theta}$, initial velocity $\boldsymbol{v}$
   **while** stopping criterion not met **do**
      Sample a minibatch of $m$ examples from the training set $\{\boldsymbol{x}^{(1)}, \ldots, \boldsymbol{x}^{(m)}\}$ with corresponding labels $\boldsymbol{y}^{(i)}$.
      Apply interim update: $\tilde{\boldsymbol{\theta}} \leftarrow \boldsymbol{\theta} + \alpha \boldsymbol{v}$.
      Compute gradient (at interim point): $\boldsymbol{g} \leftarrow \frac{1}{m} \nabla_{\tilde{\boldsymbol{\theta}}} \sum_i L(f(\boldsymbol{x}^{(i)}; \tilde{\boldsymbol{\theta}}), \boldsymbol{y}^{(i)})$.
      Compute velocity update: $\boldsymbol{v} \leftarrow \alpha \boldsymbol{v} - \epsilon \boldsymbol{g}$.
      Apply update: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \boldsymbol{v}$.
   **end while**

# Basic Gradient Descent Algorithms



## Why Momentum Really Works

Step-size α = 0.02

Momentum β = 0.99

We often think of Momentum as a means of dampening oscillations and speeding up the iterations, leading to faster convergence. But it has other interesting behavior. It allows a larger range of step-sizes to be used, and creates its own oscillations. What is going on?

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

Sangheum Hwang

# 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 → small update, small gradients → large update
  - Code snippet

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

  - Cache keeps growing during training → learning rate shrinks too quickly

# Algorithms with Adaptive Learning Rates

- AdaGrad (2011)

---

**Algorithm 8.4** The AdaGrad algorithm

**Require:** Global learning rate $\epsilon$

**Require:** Initial parameter $\boldsymbol{\theta}$

**Require:** Small constant $\delta$, perhaps $10^{-7}$, for numerical stability

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

**while** stopping criterion not met **do**

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

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

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

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

Apply update: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \Delta \boldsymbol{\theta}$

**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

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

```
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 $\epsilon$, decay rate $\rho$.

**Require:** Initial parameter $\boldsymbol{\theta}$

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

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

**while** stopping criterion not met **do**

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

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

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

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

    Apply update: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \Delta \boldsymbol{\theta}$

**end while**

# Algorithms with Adaptive Learning Rates

- Adam (2015)
    - 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)

---

**Algorithm 8.7** The Adam algorithm

---

**Require:** Step size $\epsilon$ (Suggested default: 0.001)

**Require:** Exponential decay rates for moment estimates, $\rho_1$ and $\rho_2$ in $[0, 1)$. (Suggested defaults: 0.9 and 0.999 respectively)

**Require:** Small constant $\delta$ used for numerical stabilization. (Suggested default: $10^{-8}$)

**Require:** Initial parameters $\boldsymbol{\theta}$

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

  Initialize time step $t = 0$

  **while** stopping criterion not met **do**

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

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

  $t \leftarrow t + 1$

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

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

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

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

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

  Apply update: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \Delta\boldsymbol{\theta}$

  **end while**

---

# And Many More



모든 자료를 다 검토해서
내 위치의 산기울기를 계산해서
갈 방향을 찾겠다.

**GD**

**SGD**
전부 다봐야 한걸음은
너무 오래 걸리니까
조금만 보고 빨리 판단한다
같은 시간에 더 많이 간다

스텝방향

스텝사이즈

**Momentum**
스텝 계산해서 움직인 후,
아까 내려 오던 관성 방향 또 가자

Nesterov Accelerated Gradient
**NAG**
일단 관성 방향 먼저 움직이고,
움직인 자리에 스텝을 계산하니
더 빠르더라

**Adagrad**
안가본곳은 성큼 빠르게 걸어 훑고
많이 가본 곳은 잘아니까
갈수록 보폭을 줄여 세밀히 탐색

**RMSProp**
보폭을 줄이는 건 좋은데
이전 맥락 상황봐가며 하자.

**AdaDelta**
종종걸음 너무 작아져서
정지하는걸 막아보자.

**Adam**
RMSProp + Momentum
방향도 스텝사이즈도 적절하게!

**Nadam**
Adam에 Momentum
대신 NAG를 붙이자.

# Learning Rate Scheduling

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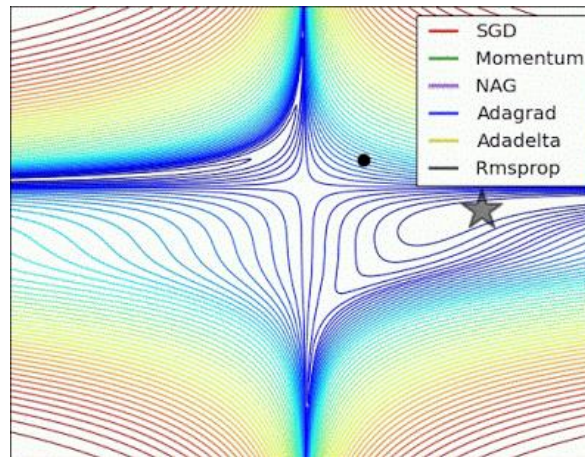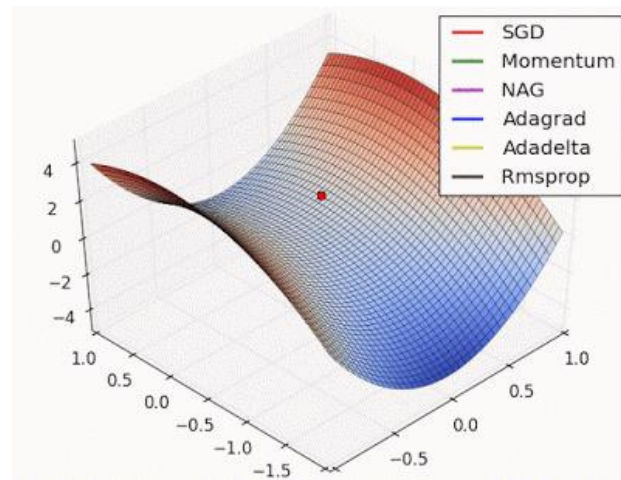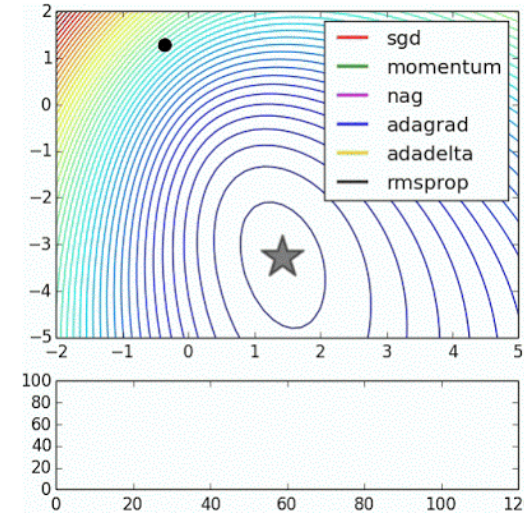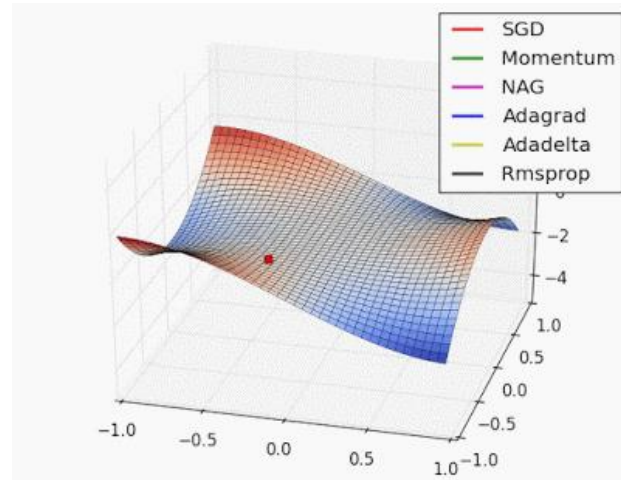


**Step decay**

**Exponential decay**

**Test accuracy**

# Animations for Optimization Algorithms

# Reading assignments

- "Understanding deep learning"
  - Chapter 6
- "Dive into deep learning"
  - Section 12