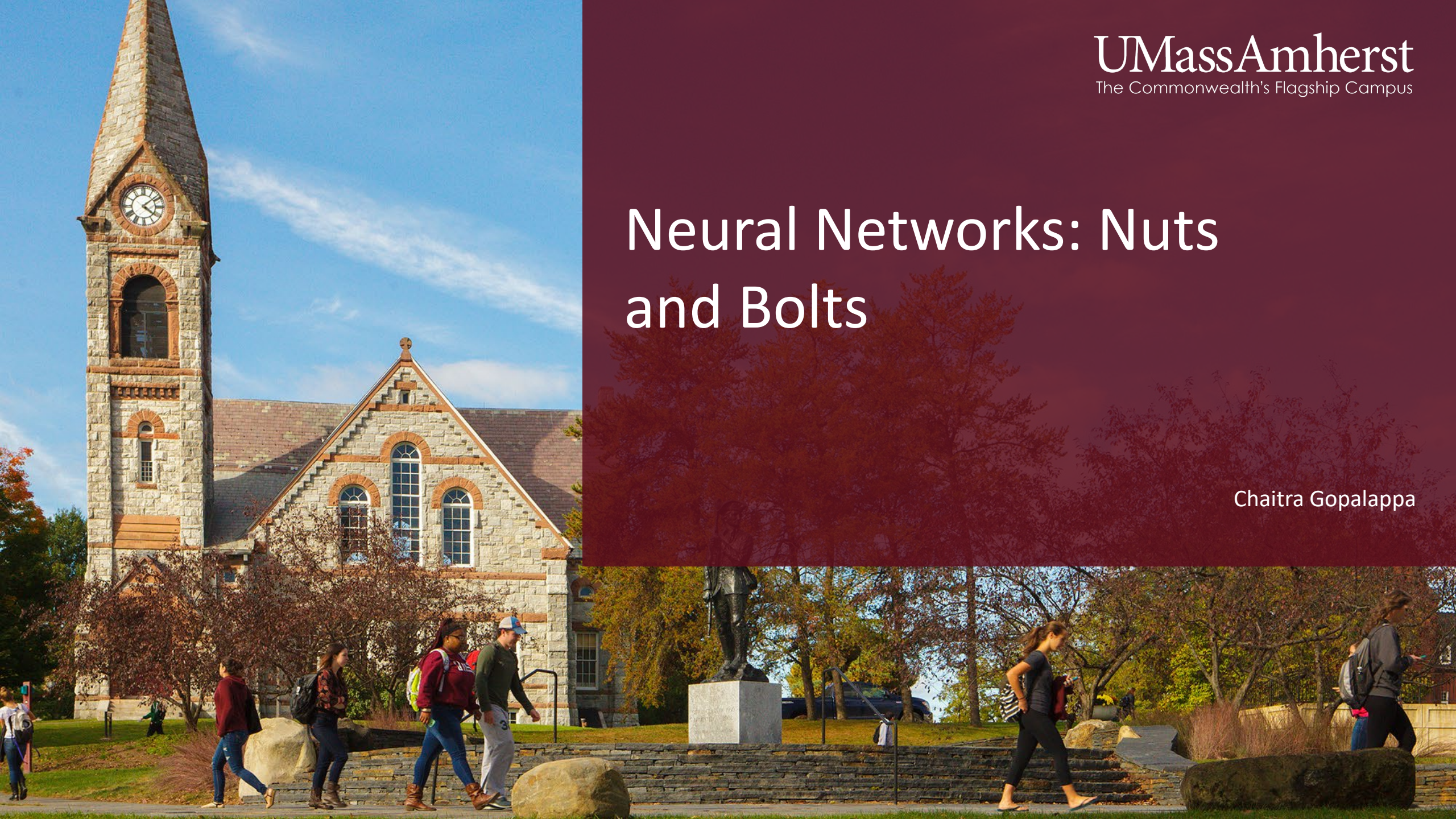


# Neural Networks: Nuts and Bolts

Chaitra Gopalappa





# Batch, mini-batch, incremental (online)

Source: Bengio, Practical recommendations for gradient-based training of deep architectures, 2012, <https://arxiv.org/abs/1206.5533>

# Machine learning

- Typical to split data into train and test sets
- Use 'train' set to train the data
- Test the trained model on the 'test' set

# Batch, mini-batch, incremental (online)

- Batch: use the full train set to train the model
- Mini-batch: divide train set into small mini-batches
  - Typically  $2^n$ : 32, 64, 128, 256 (corresponding to CPU /GPU architecture)
- Incremental/ online learning: single sample at a time
- Bengio, Practical recommendations for gradient-based training of deep architectures, 2012, <https://arxiv.org/abs/1206.5533>

# Batch, mini-batch, incremental (online)

- Batch: use the full train set to train the model

Gradient descent

- Mini-batch: divide train set into small mini-batches

- Typically  $2^n$ : 32, 64, 128, 256 (corresponding to CPU /GPU architecture)

Stochastic Gradient descent

- Incremental/ online learning: single sample at a time

- Bengio, Practical recommendations for gradient-based training of deep architectures, 2012, <https://arxiv.org/abs/1206.5533>

# Batch, mini-batch, incremental (online)

- Batch: use the full train set to train the model
- Mini-batch: divide train set into small mini-batches
  - Typically  $2^n$ : 32, 64, 128, 256 (corresponding to CPU /GPU architecture)
  - Batch-size then becomes a hyperparameter
- Incremental/ online learning: single sample at a time

Gradient descent

Stochastic Gradient  
descent

- Bengio, Practical recommendations for gradient-based training of deep architectures, 2012, <https://arxiv.org/abs/1206.5533>



# Neural Network Optimizers

Source:

Chapter 8: Ian Goodfellow and Yoshua Bengio and Aaron Courville, Deep Learning

<https://www.deeplearningbook.org/contents/optimization.html>

## Recollect line search algorithms

*Transformation*

$$\vec{x}_{m+1} \leftarrow \vec{x}_m - \mu_m \vec{p}_m$$

$\vec{p}_m$ : search direction

$\mu_m$ : step size

## Neural Network Optimizers

- SDG
- SDG with momentum
- SDG with Nesterov

Gradient as  
search direction

- AdaGrad
- RMSProp
- Adam

Adaptive learning  
rate/ step size

- Newtons method
- Conjugate  
gradients
- BFGS

Hessian as  
search direction

Line Search  
methods: identify  
search direction  
and step in that  
direction



Recollect slides from first two weeks  
Generalized: Stochastic gradient  
Descent (SDG) algorithm

- *Transformation*

$$\vec{x}_{m+1} \leftarrow \vec{x}_m - \mu_m Y_m(\vec{x}_m)$$

$$Y_m(\vec{x}_m) = \nabla f(\vec{x}_m) + \text{noise}$$

## SDG in machine learning

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

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

**Require:** Initial parameter  $\theta$

$k \leftarrow 1$

**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**

Sufficient conditions for convergence  $\sum_{k=1}^{\infty} \epsilon_k = \infty$ , and  $\sum_{k=1}^{\infty} \epsilon_k^2 < \infty$ .

Typical to decay learning rate linearly until some iteration

$\tau$ , with  $\alpha = \frac{k}{\tau}$ . After iteration  $\tau$ , leave  $\epsilon$  constant.

$\epsilon_0, \epsilon_{\tau}, \tau$  are hyperparameters

$$\epsilon_k = (1 - \alpha)\epsilon_0 + \alpha\epsilon_{\tau}$$

If  $m = B$  (batch size/ full training dataset) it is general GD

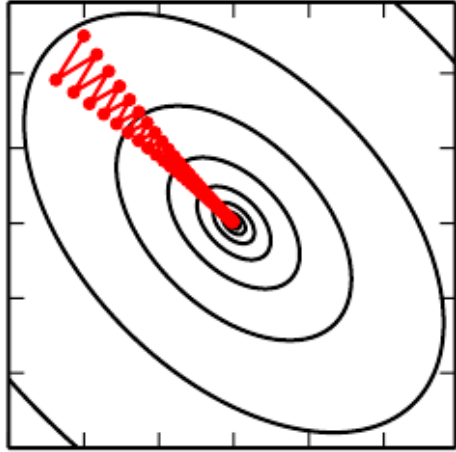
If  $m = 1$  it is incremental or online SGD

If  $1 < m < B$  it is mini-batch SDG; typically  $m = 2^n$  (some exponent of 2)

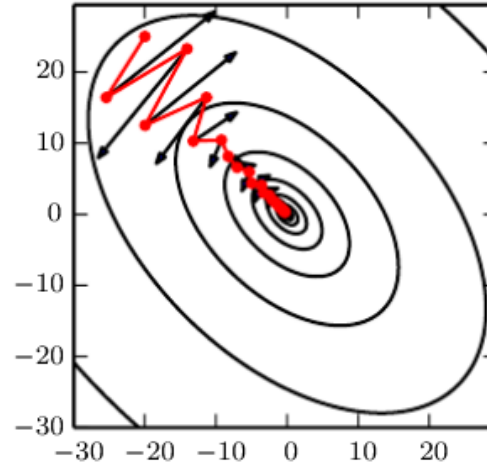
Source: Chapter 8:

<https://www.deeplearningbook.org/contents/optimization.html>

# SDG



# SDG with momentum



The difference between Nesterov momentum and standard momentum is where the gradient is evaluated. With Nesterov momentum, the gradient is evaluated after the current velocity is applied.

---

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

---

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

**Require:** Initial parameter  $\theta$ , 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 targets  $y^{(i)}$ .

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

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

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

**end while**

---

---

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

---

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

**Require:** Initial parameter  $\theta$ , 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**

---

Source: Chapter 8:

<https://www.deeplearningbook.org/contents/optimization.html>

# Algorithms with Adaptive Learning Rates

---

## Algorithm 8.5 The RMSProp algorithm

---

**Require:** Global learning rate  $\epsilon$ , decay rate  $\rho$

**Require:** Initial parameter  $\theta$

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

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

**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:  $\mathbf{r} \leftarrow \rho \mathbf{r} + (1 - \rho) \mathbf{g} \odot \mathbf{g}$ .

**[** 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**

---

---

## Algorithm 8.4 The AdaGrad algorithm

---

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

**Require:** Initial parameter  $\theta$

**Require:** Small constant  $\delta$ , 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**

---

---

## Algorithm 8.6 RMSProp algorithm with Nesterov momentum

---

**Require:** Global learning rate  $\epsilon$ , decay rate  $\rho$ , momentum coefficient  $\alpha$

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

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

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

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

**[** Compute velocity update:  $\mathbf{v} \leftarrow \alpha \mathbf{v} - \frac{\epsilon}{\sqrt{\mathbf{r}}} \odot \mathbf{g}$ . ( $\frac{1}{\sqrt{\mathbf{r}}}$  applied element-wise)

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

**end while**

---

Source: Chapter 8:

<https://www.deeplearningbook.org/contents/optimization.html>

# Algorithms with Adaptive Learning Rates

---

**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  $\theta$

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:  $\hat{\mathbf{r}} \leftarrow \frac{\mathbf{r}}{1 - \rho_2^t}$

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**

---

Source: Chapter 8:

<https://www.deeplearningbook.org/contents/optimization.html>

# SDG function in Pytorch

- <https://pytorch.org/docs/stable/generated/torch.optim.SGD.html#torch.optim.SGD>
- Bengio, Practical recommendations for gradient-based training of deep architectures, 2012, <https://arxiv.org/abs/1206.5533>

## SGD

```
CLASS torch.optim.SGD(params, lr=<required parameter>, momentum=0, dampening=0,  
    weight_decay=0, nesterov=False, *, maximize=False, foreach=None,  
    differentiable=False) [SOURCE]
```

Docs > torch.optim > SGD

---

**input** :  $\gamma$  (lr),  $\theta_0$  (params),  $f(\theta)$  (objective),  $\lambda$  (weight decay),  
 $\mu$  (momentum),  $\tau$  (dampening), **nesterov**, *maximize*

---

**for**  $t = 1$  **to**  $\dots$  **do**

$g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$

**if**  $\lambda \neq 0$

$g_t \leftarrow g_t + \lambda \theta_{t-1}$

**if**  $\mu \neq 0$

**if**  $t > 1$

$\mathbf{b}_t \leftarrow \mu \mathbf{b}_{t-1} + (1 - \tau) g_t$

**else**

$\mathbf{b}_t \leftarrow g_t$

**if** *nesterov*

$g_t \leftarrow g_t + \mu \mathbf{b}_t$

**else**

$g_t \leftarrow \mathbf{b}_t$

**if** *maximize*

$\theta_t \leftarrow \theta_{t-1} + \gamma g_t$

**else**

$\theta_t \leftarrow \theta_{t-1} - \gamma g_t$

---

**return**  $\theta_t$

---



# More variants added over time

- Best way to keep up is to look at NN libraries (Keras, Pytorch, Tensorflow) on available options
  - <https://keras.io/api/optimizers/>
  - <https://pytorch.org/docs/stable/optim.html>
  - [https://www.tensorflow.org/api\\_docs/python/tf/keras/optimizers](https://www.tensorflow.org/api_docs/python/tf/keras/optimizers)

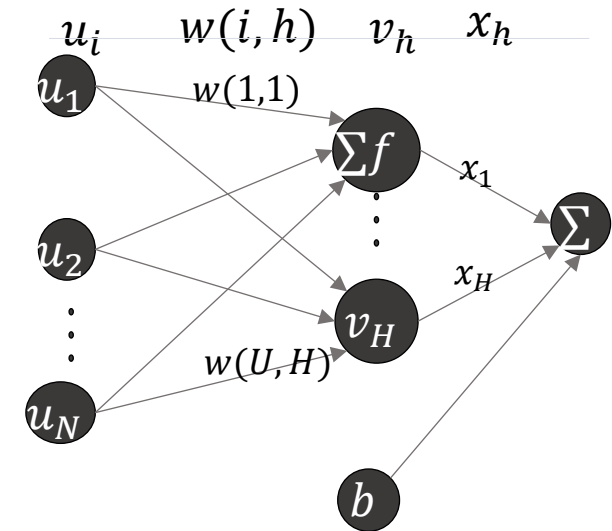
# Tensors and AutoDiff

- [AutoDiff](#): Baydan et. al., Journal of Machine Learning Research 18 (2018) 1-43
- <https://dlsyscourse.org/slides/4-automatic-differentiation.pdf>

# Activation functions

# Challenges with Sigmoid- vanishing gradient

- Sigmoid is a saturation function
- $v_h = \frac{1}{1+e^{-\vec{w}_h \vec{u}}}$  ; if weights are initialized to be very large or small,  $v_h$  saturates at 0 or 1
- Gradient then becomes very small (stops learning)  $\rightarrow$  vanishing gradient problem



# Other activation functions

- Hyperbolic tangent
- Rectified Linear Unit (ReLU)
- Leaky ReLU
- Exponential Linear unit
- Numerous other AF: <https://arxiv.org/abs/1811.03378>



Additional components useful for RL

# Output layer

- Softmax with multiple nodes in output layer → converts output to a probability distribution
  - AF may then be used in output layer as well,
  - In RL, for policy gradient, we will use multiple nodes in output layer with softmax AF.
  - Lapan Ch3-Module1 code uses softmax
- Dropout in output layer :
  - It is a form of regularization in supervised learning, especially when data is sparse; <https://arxiv.org/pdf/1207.0580.pdf>
  - It is not used that much in reinforcement learning, but there is some research in its use in policy gradient <https://doi.org/10.48550/arXiv.2202.11818>
  - Lapan Ch3-Module1 code uses dropout

UMassAmherst  
The Commonwealth's Flagship Campus