

Machine Learning from Scratch: Stochastic Gradient Descent and Adam Optimizer

JAMES GABBARD AND
DANIEL MILLER

18.0851 FINAL
PROJECT

Introduction

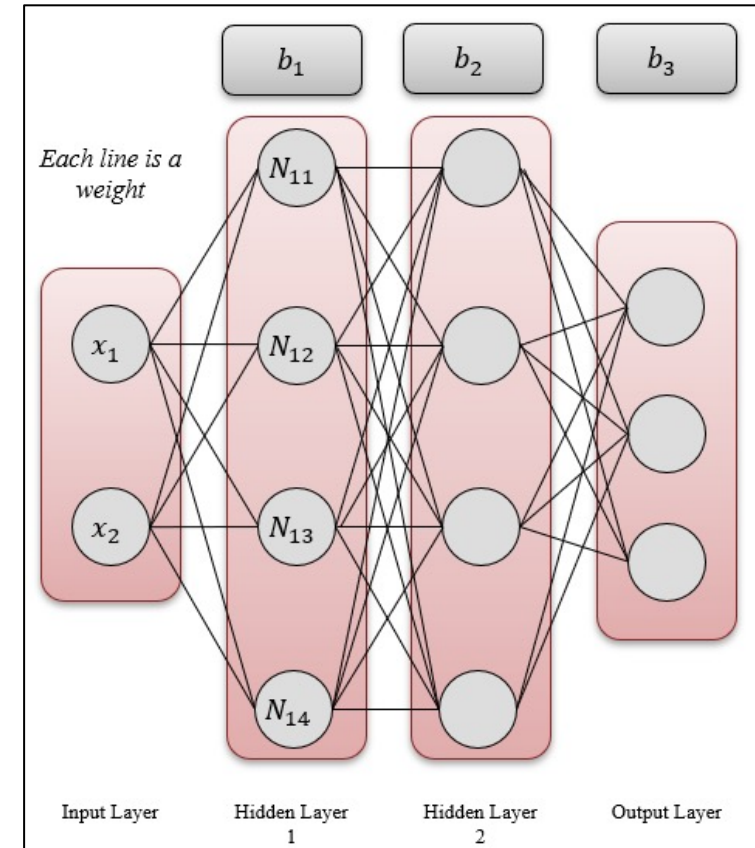
- Neural networks are an increasingly important part of research and daily life.
- While their applications are widely discussed, the algorithms that allow them to learn are not.
- How do they actually learn?
 - Optimization Algorithms:
 - Stochastic Gradient Descent and Backpropagation
 - Adam Optimizer
 - Application test:
 - Tensorflow Playground-style classification problems
 - Reinforcement Learning: CartPole
- All coding completed in MATLAB

Neural Network Overview

- Composed of layers d deep and n neurons wide
 - Network shown here: 2 layers
- Connected via weights and biases
- Activation functions add nonlinearities

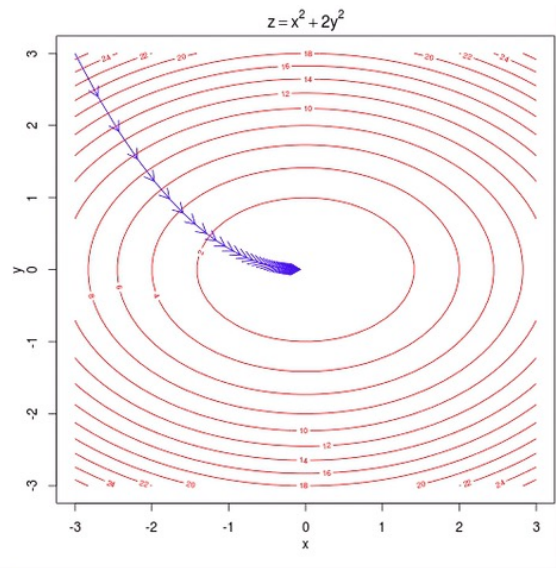
$$N_1 = \text{Activation}(W_1 x^T + b_1)$$

- In order for a network to learn, a method must be devised to shape these weights and biases to produce an accurate output.
 - Minimize a cost function.
- Classification problem: Cost function is the error between data labels and neural network classification



Stochastic Gradient Descent (SGD)

- With a set of data X , SGD optimizes minimizing an objective function $J(\theta)$ using a randomly (STOCHASTIC) selected minibatch of data by via GRADIENT DESCENT, $-\nabla_{\theta}J(\theta)$.



Algorithm 2 Stochastic Gradient Descent (SGD)

```
1: procedure SGD TRAINING( $X$ )
2:   initialize  $\theta = (W, b)$  to small random numbers
3:   randomize order of training examples in  $\mathbf{X}$ 
4:   while not converged do
5:     Create minibatch of length  $B$  from  $X$ 
6:      $a_1 = \text{minibatch}_j$ 
7:      $\nabla_{\theta} = \text{Backpropagation}(X, \theta)$ 
8:     for  $i \leftarrow [2, L]$  do
9:        $\theta_i \leftarrow \theta_i + \frac{\eta}{B} \nabla_{\theta}$ 
10:    end for
11:  end while
12: end procedure
```

Backpropagation (In 5 Minutes)

W^1, \dots, W^L and b^1, \dots, b^L are weight matrices and bias vectors

σ^ℓ is the activation function at layer ℓ

$z^\ell = W^\ell a^{\ell-1} + b^\ell$ is the weighted input to layer ℓ

a^ℓ is the activation at layer ℓ

C is our cost function

$\delta_\ell = \frac{\partial C}{\partial z_\ell}$ is the “error”, a quantity central to backpropagation

For elementwise multiplication (“dot star”), we’ll use \odot

$$c = a \odot b \quad \leftrightarrow \quad c_i = a_i b_i$$

Plan of Attack

Plan of attack:

Forward pass: compute a 's (activations) and z 's (weighted inputs) for each layer

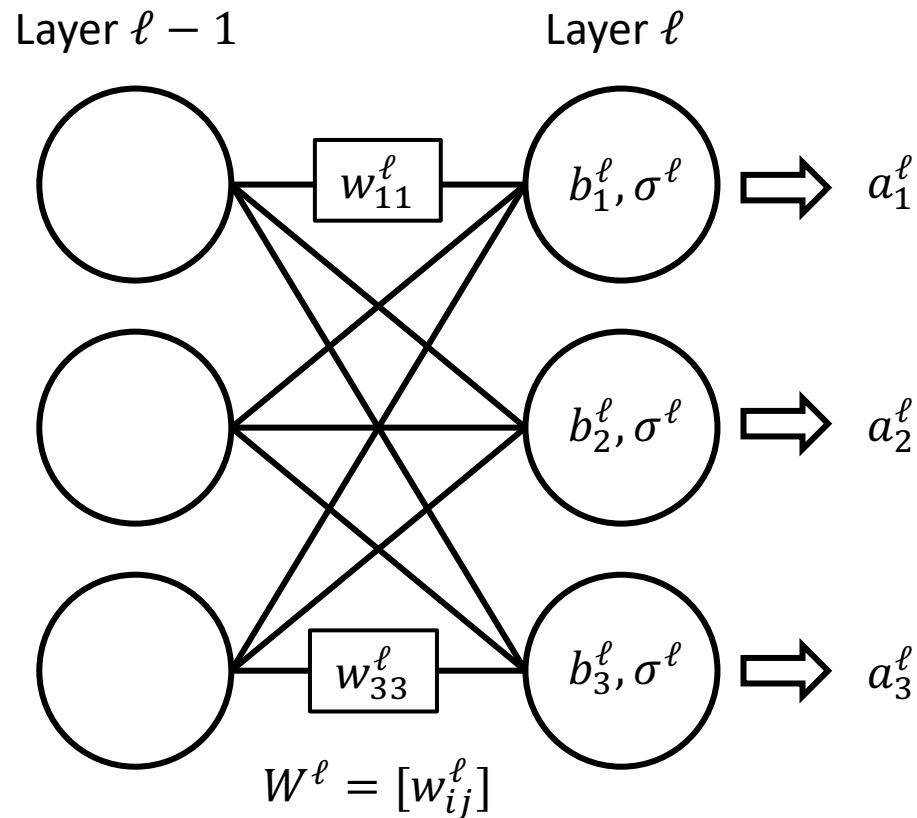
Backward pass:

$$\delta^\ell = \frac{\partial C}{\partial z_\ell} = \frac{\partial C}{\partial a_L} \times \frac{\partial a_L}{\partial z_L} \times \frac{\partial z_L}{\partial a_{L-1}} \times \frac{\partial a_{L-1}}{\partial z_{L-1}} \times \frac{\partial z_{L-1}}{\partial a_{L-2}} \times \dots$$

Compute Gradients:

$$\frac{\partial C}{\partial b}, \frac{\partial C}{\partial W}$$

Feed Forward: Compute a^ℓ and z^ℓ



Feed Forward:

$$z_i^\ell = \sum_j w_{ij}^\ell a_j^{\ell-1} + b_i^\ell$$

$$a_i^\ell = \sigma^\ell(z_i^\ell)$$

Derivatives:

$$\frac{\partial z_i^\ell}{\partial a_j^{\ell-1}} = w_{ij}^\ell$$

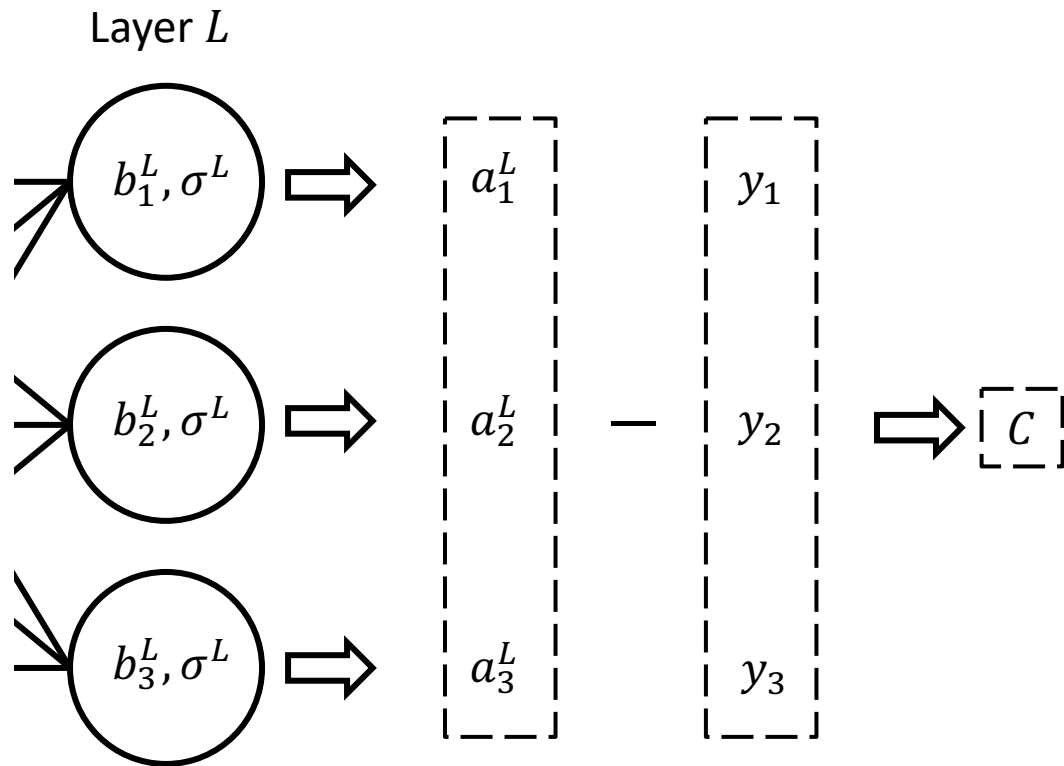
$$\frac{da_i^\ell}{dz_i^\ell} = \sigma'(z_i^\ell)$$

Vector Notation:

$$z^\ell = W^\ell a^{\ell-1} + b^\ell$$

$$a^\ell = \sigma^\ell(z^\ell)$$

The Last Layer: Compute $\delta^L = \partial C / \partial z_L$



Quadratic Cost Function:

$$C = \frac{1}{2} \sum_i (a_i^L - y_i)^2$$
$$= \frac{1}{2} \sum_i [\sigma^L(z_i^L) - y_i]^2$$

Taking Derivatives:

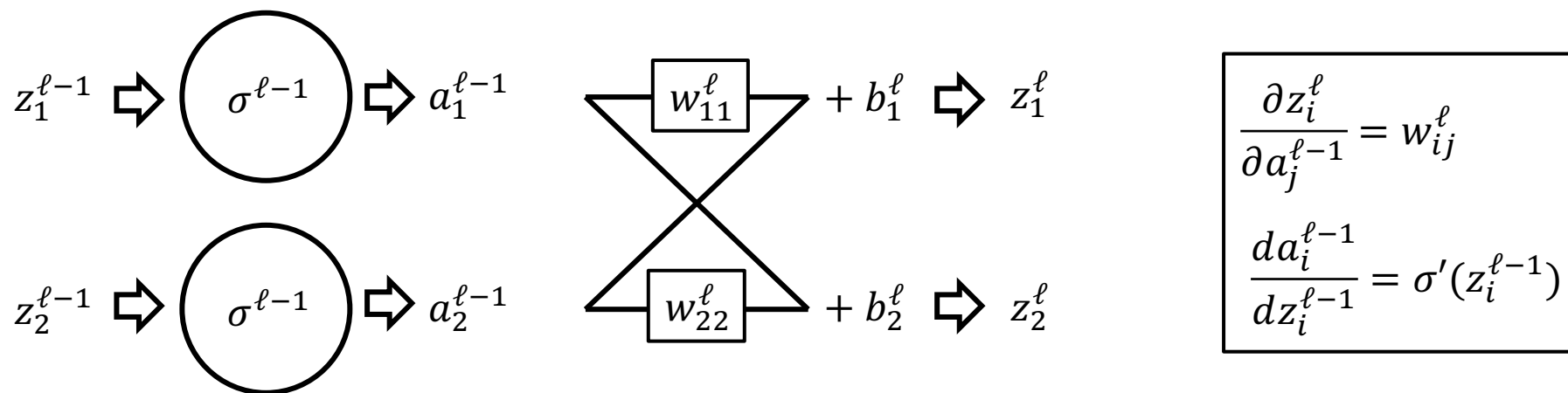
$$\frac{\partial C}{\partial z_i^L} = [\sigma^L(z_i^L) - y_i] \sigma'(z_i^L)$$

Substitute:

Vector Notation:

$$\delta_i^L = (a_i^L - y_i) \sigma'(z_i^L) \quad \Rightarrow \quad \delta^L = (a^L - y) \odot \sigma'(z^L)$$

Propagating Error: Given δ^ℓ , Compute $\delta^{\ell-1}$



Chain Rule:

Substitute:

Vector Notation:

$$\frac{\partial C}{\partial a_i^{\ell-1}} = \sum_j \frac{\partial C}{\partial z_j^\ell} \times \frac{\partial z_j^\ell}{\partial a_i^{\ell-1}} \Rightarrow \frac{\partial C}{\partial a_i^{\ell-1}} = \sum_j \delta_j^\ell w_{ji}^\ell \Rightarrow \frac{\partial C}{\partial a^{\ell-1}} = (W^\ell)^T \delta^\ell$$

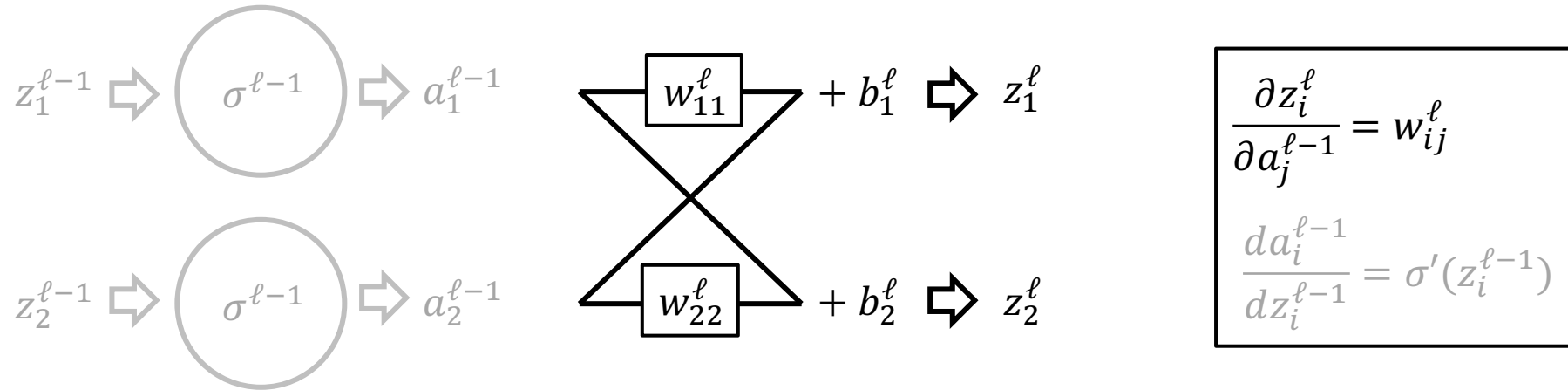
Chain Rule:

Substitute:

Vector Notation:

$$\frac{\partial C}{\partial z_i^{\ell-1}} = \frac{\partial C}{\partial a_i^{\ell-1}} \times \frac{da_i^{\ell-1}}{dz_i^{\ell-1}} \Rightarrow \frac{\partial C}{\partial z_i^{\ell-1}} = \frac{\partial C}{\partial a_i^{\ell-1}} \sigma'(z_i^{\ell-1}) \Rightarrow \delta_i^{\ell-1} = (W^\ell)^T \delta^\ell \odot \sigma'(z^{\ell-1})$$

Propagating Error: Given δ^ℓ , Compute $\delta^{\ell-1}$



Chain Rule:

Substitute:

Vector Notation:

$$\frac{\partial C}{\partial a_i^{\ell-1}} = \sum_j \frac{\partial C}{\partial z_j^\ell} \times \frac{\partial z_j^\ell}{\partial a_i^{\ell-1}} \quad \Rightarrow \quad \frac{\partial C}{\partial a_i^{\ell-1}} = \sum_j \delta_j^\ell w_{ji}^\ell \quad \Rightarrow \quad \frac{\partial C}{\partial a^{\ell-1}} = (W^\ell)^T \delta^\ell$$

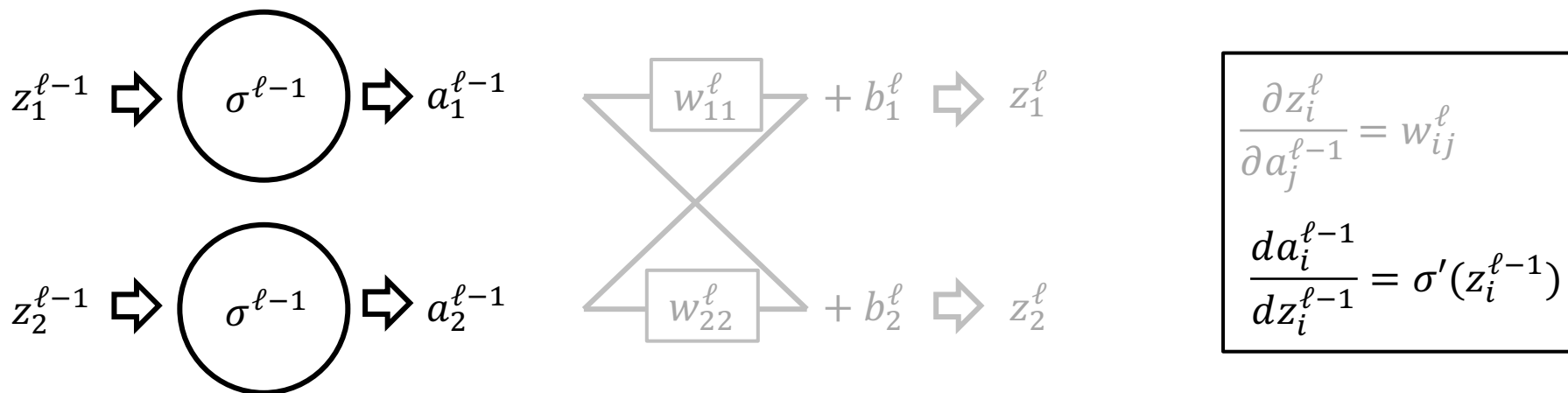
Chain Rule:

Substitute:

Vector Notation:

$$\frac{\partial C}{\partial z_i^{\ell-1}} = \frac{\partial C}{\partial a_i^{\ell-1}} \times \frac{da_i^{\ell-1}}{dz_i^{\ell-1}} \quad \Rightarrow \quad \frac{\partial C}{\partial z_i^{\ell-1}} = \frac{\partial C}{\partial a_i^{\ell-1}} \sigma'(z_i^{\ell-1}) \quad \Rightarrow \quad \delta_i^{\ell-1} = (W^\ell)^T \delta^\ell \odot \sigma'(z^{\ell-1})$$

Propagating Error: Given δ^ℓ , Compute $\delta^{\ell-1}$



Chain Rule:

Substitute:

Vector Notation:

$$\frac{\partial C}{\partial a_i^{\ell-1}} = \sum_j \frac{\partial C}{\partial z_j^\ell} \times \frac{\partial z_j^\ell}{\partial a_i^{\ell-1}} \Rightarrow \frac{\partial C}{\partial a_i^{\ell-1}} = \sum_j \delta_j^\ell w_{ji}^\ell \Rightarrow \frac{\partial C}{\partial a^{\ell-1}} = (W^\ell)^T \delta^\ell$$

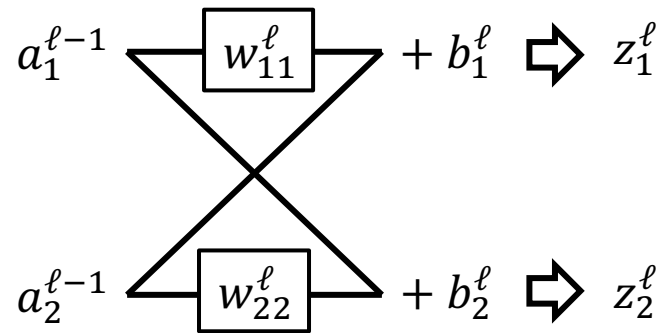
Chain Rule:

Substitute:

Vector Notation:

$$\frac{\partial C}{\partial z_i^{\ell-1}} = \frac{\partial C}{\partial a_i^{\ell-1}} \times \frac{da_i^{\ell-1}}{dz_i^{\ell-1}} \Rightarrow \frac{\partial C}{\partial z_i^{\ell-1}} = \frac{\partial C}{\partial a_i^{\ell-1}} \sigma'(z_i^{\ell-1}) \Rightarrow \delta_i^{\ell-1} = (W^\ell)^T \delta^\ell \odot \sigma'(z^{\ell-1})$$

Gradients: Given δ^ℓ , Compute $\partial C / \partial b^\ell$ and $\partial C / \partial W^\ell$



Weighted Input:

$$z_i^\ell = \sum_j w_{ij}^\ell a_j^{\ell-1} + b_i^\ell$$

Derivatives:

$$\frac{\partial z_j^\ell}{\partial b_i^\ell} = Id_{ij}$$

$$\frac{\partial z_i^\ell}{\partial w_{ij}^\ell} = a_j^{\ell-1}$$

Apply Chain Rule:

$$\frac{\partial C}{\partial b_i^\ell} = \frac{\partial C}{\partial z_j^\ell} \times \frac{\partial z_j^\ell}{\partial b_i^\ell}$$

$$\frac{\partial C}{\partial w_{ij}^\ell} = \frac{\partial C}{\partial z_i^\ell} \times \frac{\partial z_i^\ell}{\partial w_{ij}^\ell}$$

Substitute for derivatives:

$$\frac{\partial C}{\partial b_i^\ell} = \delta_j^\ell \times Id_{ij} = \delta_i^\ell$$

$$\frac{\partial C}{\partial w_{ij}^\ell} = \delta_i^\ell \times a_j^{\ell-1}$$

Vector Notation:

$$\frac{\partial C}{\partial b^\ell} = \delta^\ell$$

$$\frac{\partial C}{\partial W^\ell} = \delta^\ell (a^{\ell-1})^T$$

The Complete Backpropagation Algorithm

Feed Forward:

$$z^\ell = W^\ell a^{\ell-1} + b^\ell$$

$$a^\ell = \sigma^\ell(z^\ell)$$

Backpropagate Error:

$$\delta^L = (a^L - y) \odot \sigma'(z^L)$$

$$\delta^{\ell-1} = (W^\ell)^T \delta^\ell \odot \sigma'(z^{\ell-1})$$

Calculate Gradients:

$$\frac{\partial C}{\partial b^\ell} = \delta^\ell$$

$$\frac{\partial C}{\partial W^\ell} = \delta^\ell (a^{\ell-1})^T$$

Translating to MATLAB

Setting up the network

```
% Initializing network by hand: size is [2,3,3,2];
W2 = 0.5*ones(3,2);
b2 = 0.5*ones(3,1);
W3 = 0.5*ones(3,3);
b3 = 0.5*ones(3,1);
W4 = 0.5*ones(2,3);
b4 = 0.5*ones(2,1);

% Cell arrays for each layer
W = {[],W2, W3, W4};
b = {[],b2, b3, b4};
sigma = {[],@ReLU, @ReLU, @sigmoid};
grad = {[],@grad_ReLU, @grad_ReLU, @grad_sigmoid};

% Training data here; one data point for this example
features = {[2;2]};
labels = {[0.3;0.7]};

% Initialize input and activation arrays
z = {}; % Weighted input to each layer
a = {}; % Activation at each layer

% Max Index, (For readability)
L = length(W);
```

Backpropagation

```
% Set input layer
a{1} = features{1}; %features(1) is a cell

% Forward pass: apply the network to the data
for i = 2:L % Layer 1 is the input layer
    z{i} = W{i} * a{i-1} + b{i};
    a{i} = sigma{i}(z{i});
end

% Calculate error for output layer
delta{L} = (labels{1} - a{L}) .* grad{L}(z{L});
grad_b{L} = delta{L};
grad_W{L} = delta{L} * a{L-1}' ;

% Backward Pass: the chain rule
for i = (L - 1):-1:2 % Looping backwards
    % Backpropagate error
    delta{i} = (W{i+1}'*delta{i+1}).*grad{i}(z{i});

    % Compute gradient for weights and bias
    grad_b{i} = delta{i};
    grad_W{i} = delta{i} * a{i-1}' ;
end
```

Stochastic Gradient Descent (SGD)

- Network update:

$$\theta_i \leftarrow \theta_i + \frac{\eta}{B} \nabla_{\theta}$$

where η is the learning rate and B is the batch size

- Repeat process until converged

Algorithm 2 Stochastic Gradient Descent (SGD)

```
1: procedure SGD TRAINING( $X$ )
2:   initialize  $\theta = (W, b)$  to small random numbers
3:   randomize order of training examples in  $\mathbf{X}$ 
4:   while not converged do
5:     Create minibatch of length  $B$  from  $X$ 
6:      $a_1 = \text{minibatch}_j$ 
7:      $\nabla_{\theta} = \text{Backpropagation}(X, \theta)$ 
8:     for  $i \leftarrow [2, L]$  do
9:        $\theta_i \leftarrow \theta_i + \frac{\eta}{B} \nabla_{\theta}$ 
10:    end for
11:  end while
12: end procedure
```

Adaptive Moment Estimation (Adam) Optimizer

- Modification of SGD
- Introduce m and v : the first and second moments of inertia of the gradient
 - First moment = running mean
 - Second moment = running uncentered variance
- Update Rule:

$$m_i \leftarrow \beta_1 m_i + (1 - \beta_1) \nabla_{\theta}$$

$$v_i \leftarrow \beta_2 v_i + (1 - \beta_2) \nabla_{\theta}^2$$

$$\beta_1, \beta_2 \sim 1$$

These build momentum!

Algorithm 3 Adam

```
1: procedure ADAM( $X$ )
2:   initialize  $\theta = (W, b)$  to small random numbers
3:   randomize order of training examples in  $\mathbf{X}$ 
4:   initialize  $m \leftarrow 0, v \leftarrow 0$ 
5:   while not converged do
6:     Create minibatch of length  $B$  from  $X$ 
7:      $a_1 = \text{minibatch}_j$ 
8:      $\nabla_{\theta} = \text{Backpropagation}(X, \theta)$ 
9:     for  $i \leftarrow [2, L]$  do
10:       $m_i \leftarrow \beta_1 m_i + (1 - \beta_1) \nabla_{\theta}$ 
11:       $v_i \leftarrow \beta_2 v_i + (1 - \beta_2) \nabla_{\theta}^2$ 
12:       $\hat{m}_i = \frac{m_i}{1 - \beta_1^t}$ 
13:       $\hat{v}_i = \frac{v_i}{1 - \beta_2^t}$ 
14:       $\theta_i \leftarrow \theta_i - \frac{\alpha}{\sqrt{\hat{v}_i} + \epsilon} \hat{m}_i$ 
15:    end for
16:  end while
17: end procedure
```

Adam Optimizer

- Update network parameters:

$$\theta_i \leftarrow \theta_i - \frac{\alpha}{\sqrt{\hat{v}_i} + \epsilon} \hat{m}_i$$

- To avoid biasing the first few parameter updates, a correction step must be completed.

$$\hat{m}_i = \frac{m_i}{1 - \beta_1^t}$$
$$\hat{v}_i = \frac{v_i}{1 - \beta_2^t}$$

Algorithm 3 Adam

```
1: procedure ADAM( $X$ )
2:   initialize  $\theta = (W, b)$  to small random numbers
3:   randomize order of training examples in  $\mathbf{X}$ 
4:   initialize  $m \leftarrow 0, v \leftarrow 0$ 
5:   while not converged do
6:     Create minibatch of length  $B$  from  $X$ 
7:      $a_1 = \text{minibatch}_j$ 
8:      $\nabla_{\theta} = \text{Backpropagation}(X, \theta)$ 
9:     for  $i \leftarrow [2, L]$  do
10:       $m_i \leftarrow \beta_1 m_i + (1 - \beta_1) \nabla_{\theta}$ 
11:       $v_i \leftarrow \beta_2 v_i + (1 - \beta_2) \nabla_{\theta}^2$ 
12:       $\hat{m}_i = \frac{m_i}{1 - \beta_1^t}$ 
13:       $\hat{v}_i = \frac{v_i}{1 - \beta_2^t}$ 
14:       $\theta_i \leftarrow \theta_i - \frac{\alpha}{\sqrt{\hat{v}_i} + \epsilon} \hat{m}_i$ 
15:    end for
16:  end while
17: end procedure
```

Adam Optimizer

Why is the bias correction required? Consider the equation for the second moment.

$$v_i \leftarrow \beta_2 v_i + (1 - \beta_2) \nabla_{\theta}^2$$

At some time step t , this moving average can be written as a sum of previous second moments.

$$v_t = (1 - \beta_2) \sum_{i=1}^t \beta_2^{t-i} \cdot \nabla_{\theta}^2$$

How does this compare against the true second moment, $\mathbb{E}[\nabla_{\theta}^2]$?

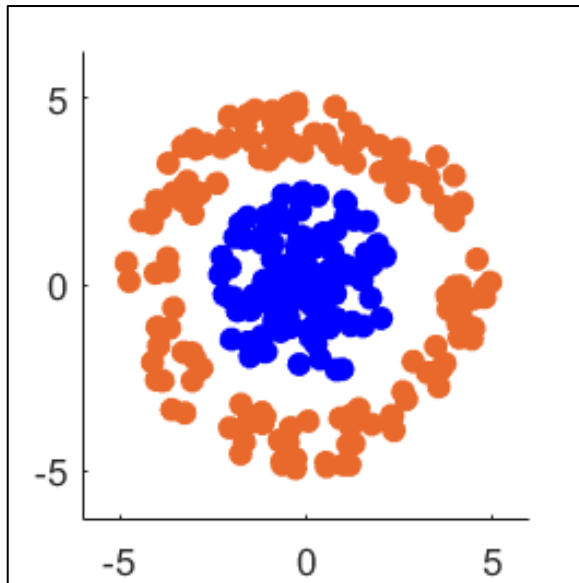
Adam Optimizer

$$\begin{aligned}\mathbb{E}[v_t] &= \mathbb{E}\left[(1 - \beta_2) \sum_{i=1}^t \beta_2^{t-i} \cdot \nabla_{\theta}^2\right] \\ &= \mathbb{E}[\nabla_{\theta}^2] \cdot (1 - \beta_2) \sum_{i=1}^t \beta_2^{t-i} \\ &= \mathbb{E}[\nabla_{\theta}^2] \cdot (1 - \beta_2^t)\end{aligned}$$

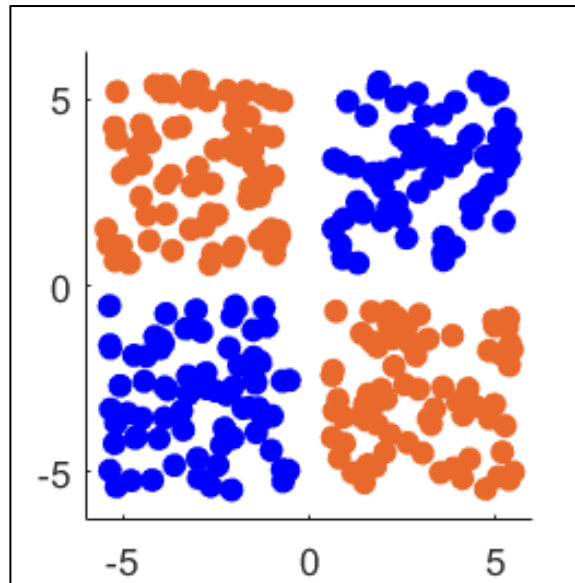
Therefore, to find the true value $\mathbb{E}[\nabla_{\theta}^2]$,

$$\mathbb{E}[\nabla_{\theta}^2] = \hat{v}_i = \frac{v_i}{1 - \beta_2^t}$$

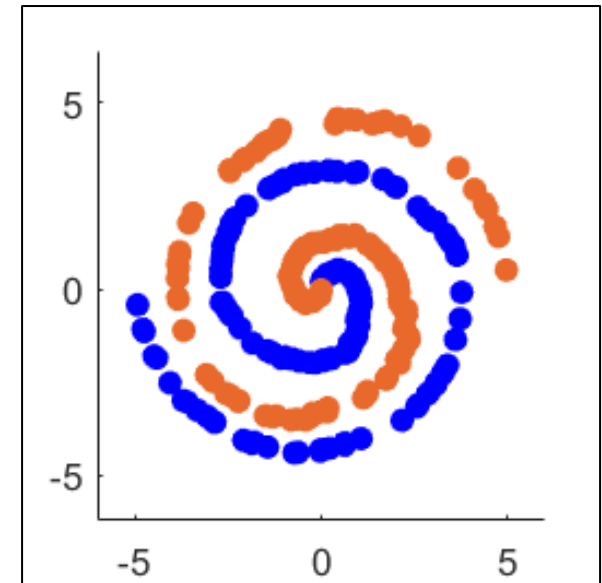
Tensorflow Playground



Circles

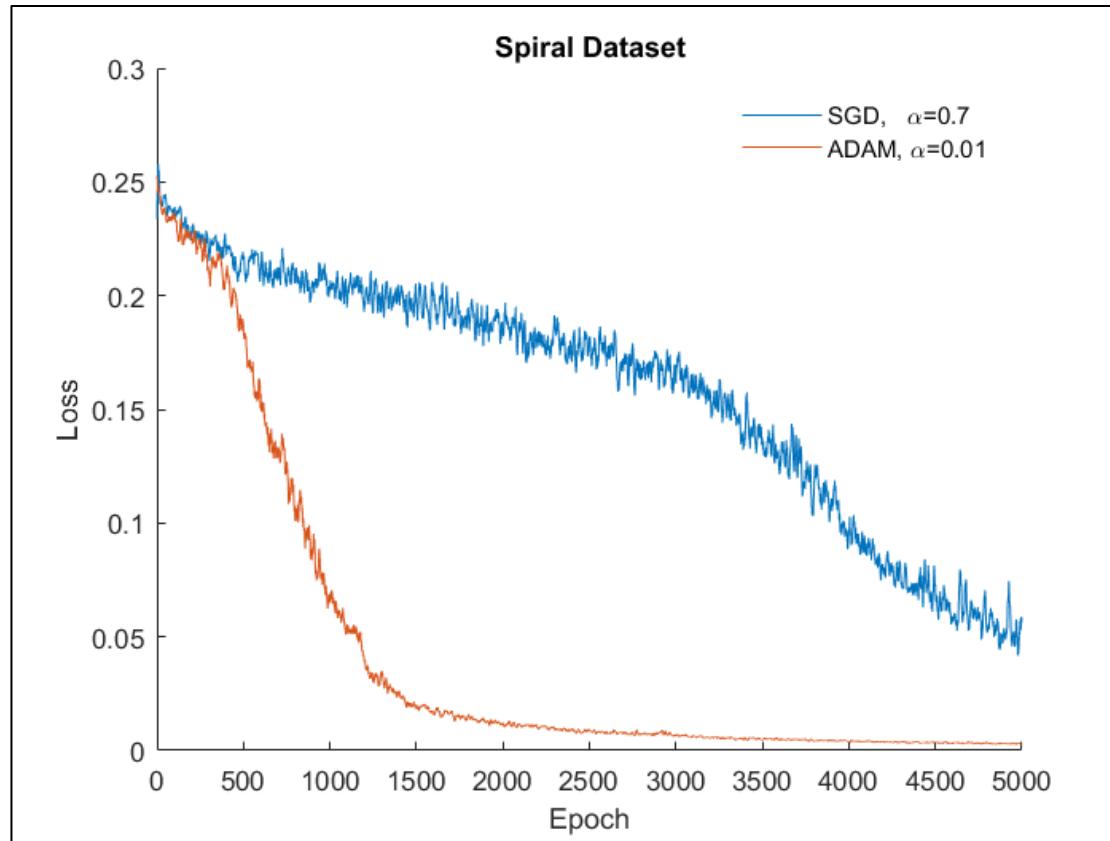


Square Regions

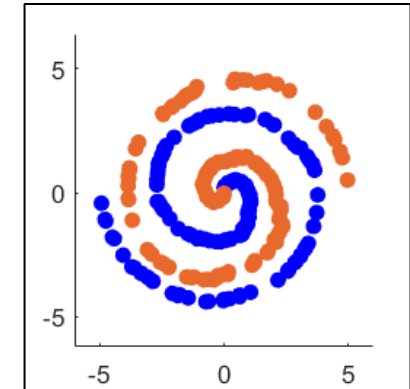


Spirals

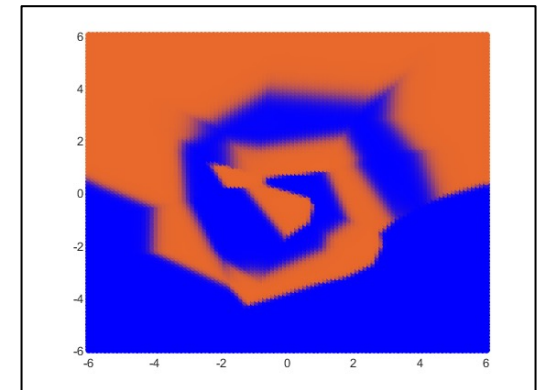
Tensorflow Playground Results



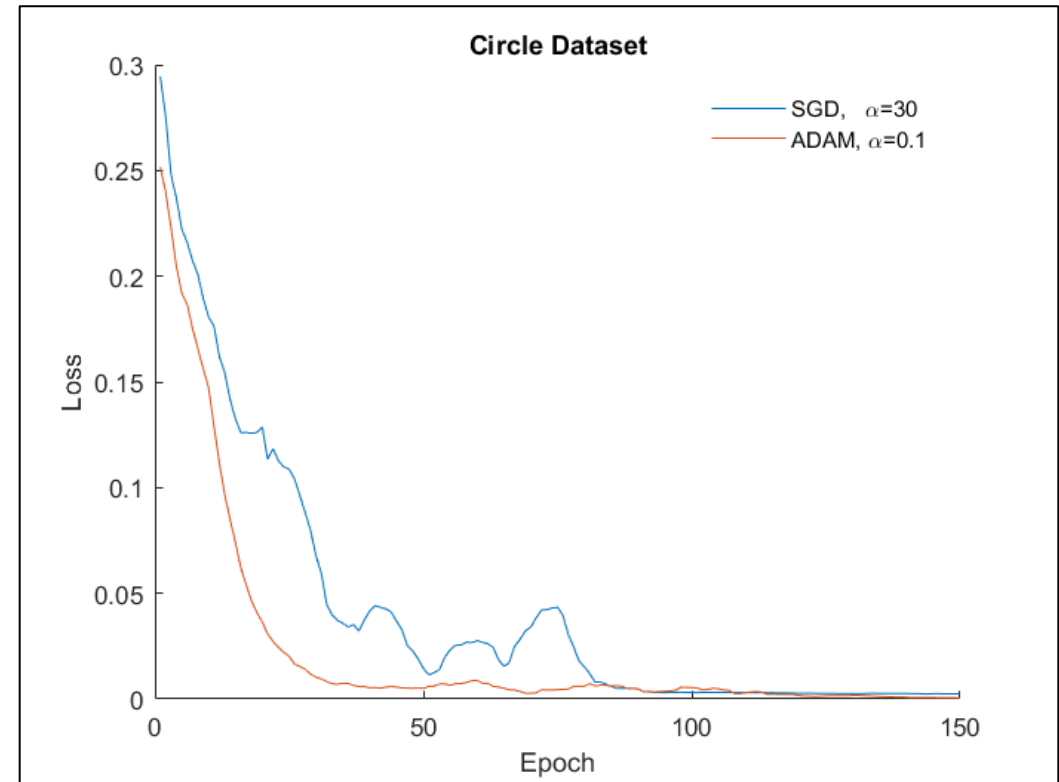
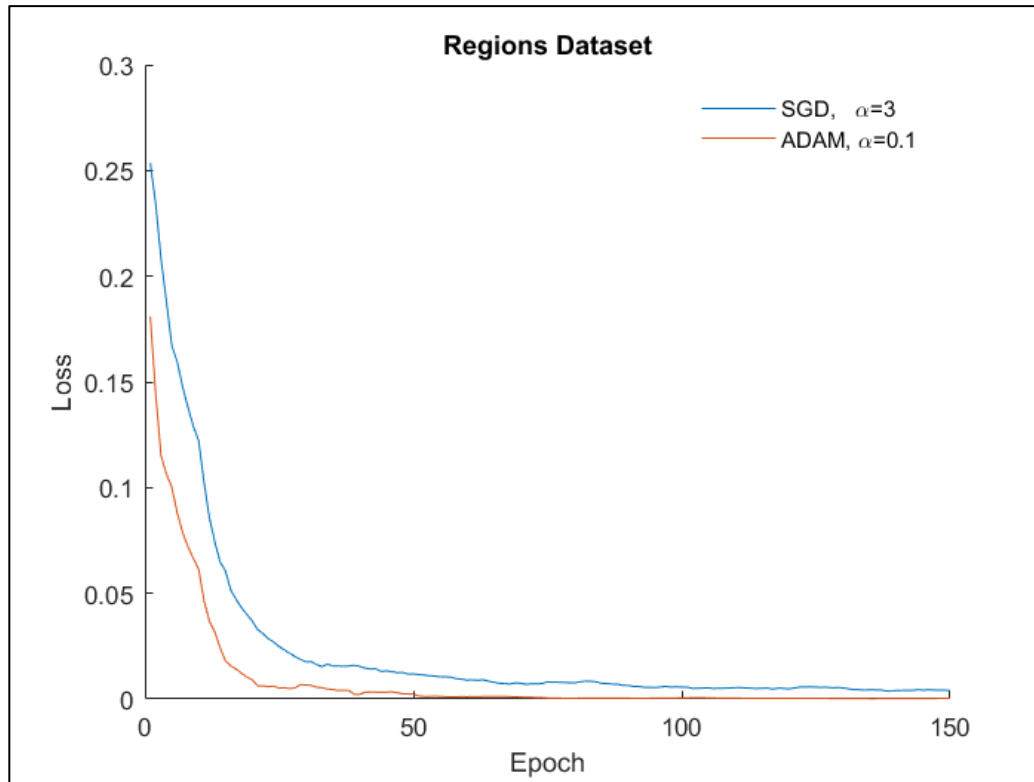
← SGD vs Adam
Spiral data set
(moving avg.)



Adam
Classification of
Spiral data set →

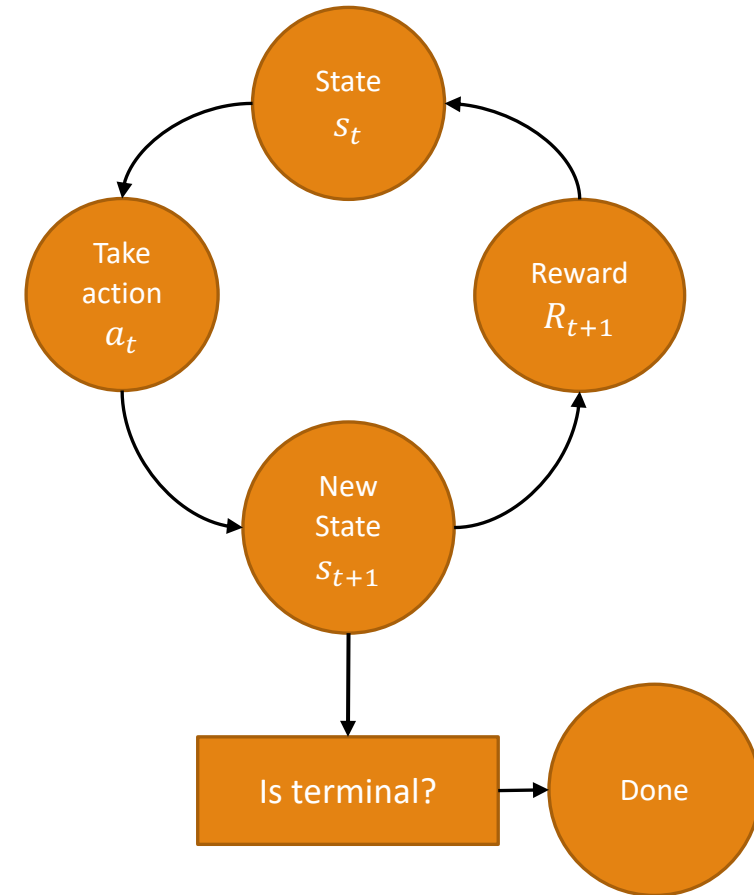


Tensorflow Playground Results



Reinforcement Learning (In 5 Minutes)

- Training an agent (the AI) to learn by giving positive or negative feedback
- Modeling a decision process:
 - Begin in a **state** s_t at time t
 - Choose an **action** a_t from a list of possible actions
 - Transition to a **new state** s_{t+1}
 - Receive appropriate **reward**, R_{t+1}
 - Repeat until entering a terminal state – success or failure
- Goal: Choose actions (policy) to maximize total rewards

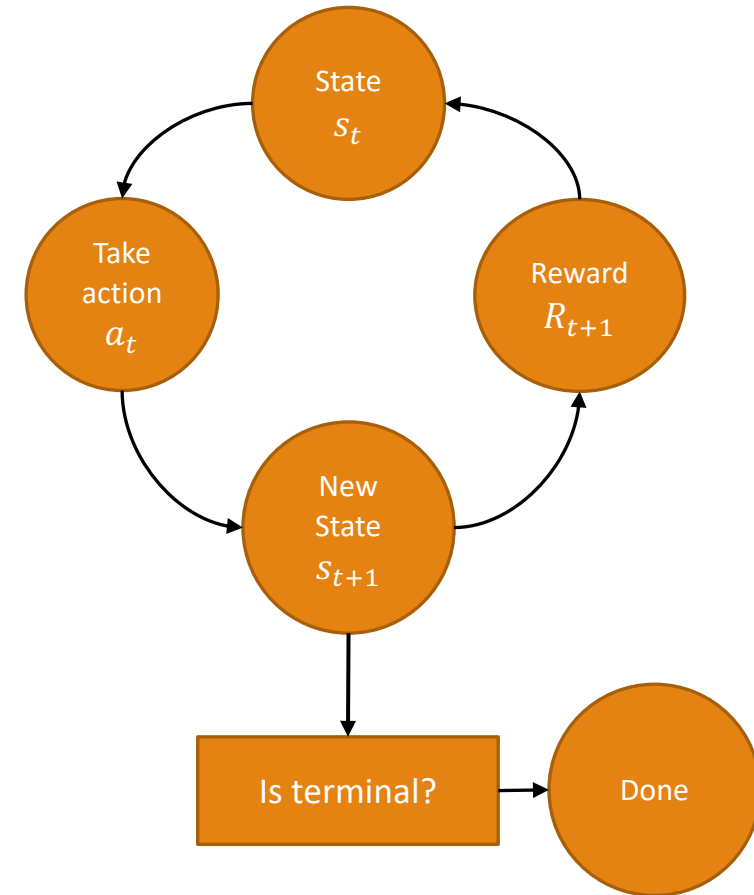


Expected Rewards

- What is total reward?
- Discount factor: $0 < \gamma < 1$
- Total discounted rewards:
$$G_t = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots$$

The “quality” of an action:

$$Q(a, s) = \mathbb{E}(G_t \mid a_t = a, s_t = s)$$



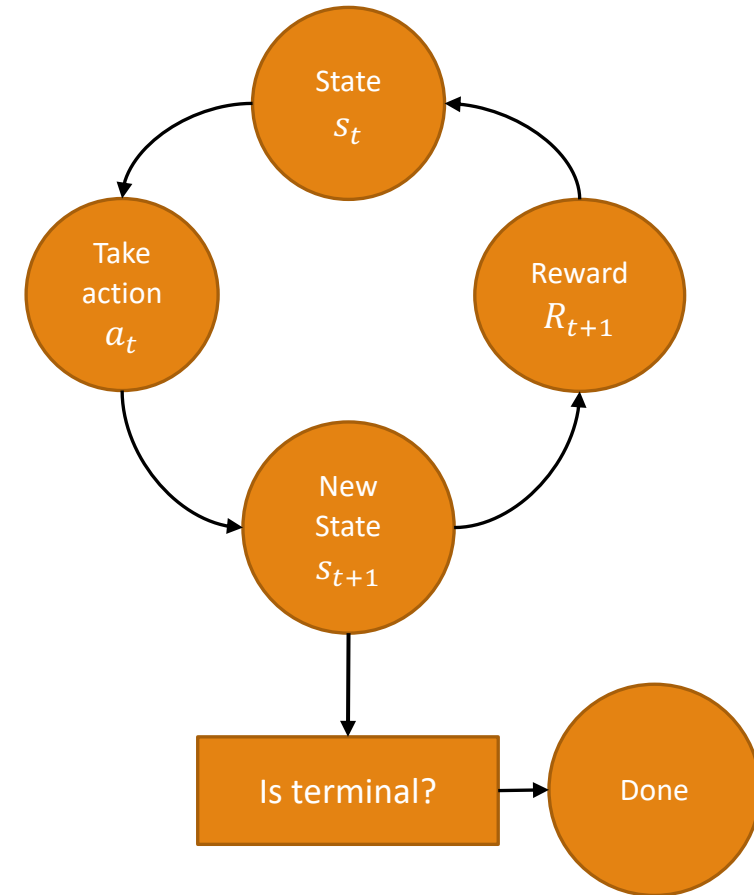
Results on MDP's

(Markov Decision Processes)

1. An optimal policy is to always choose the highest quality action

2. When following an optimal policy, Q obeys the Bellman Equation:

$$Q(s_t, a_t) = R_{t+1} + \gamma \max_a Q(s_{t+1}, a)$$



Reinforcement Learning

- DQN approximates the quality function using a neural network.

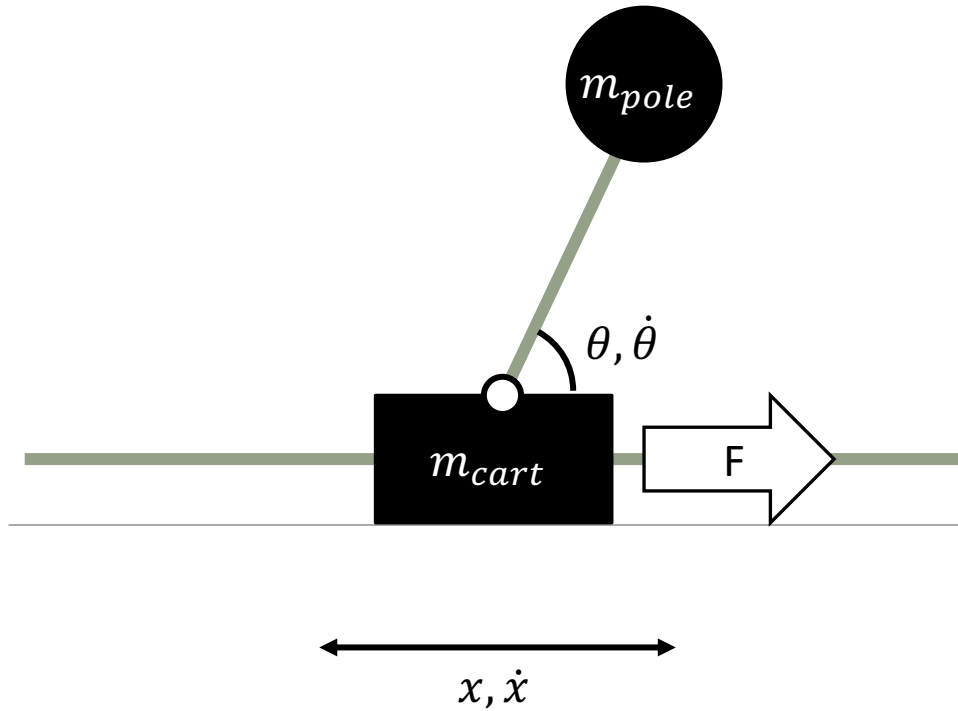
- The network is trained with a cost function based on the Bellman Equation:

$$C = \left[Q^*(s_t, a_t) - R_{t+1} - \gamma \max_a Q^*(s_{t+1}, a) \right]^2$$

- Every action becomes training data for the network (improves with time)

```
1: procedure DQN( $X$ )
2:   Initialize Neural Network to represent action-value function  $Q_\theta$ 
3:   Initialize Experience Replay memory,  $M$ 
4:   Initialize target action-value function  $\hat{Q}_{\theta-}$ 
5:   for  $ep = [1, N]$  do
6:      $s_t \leftarrow s_0$ 
7:     while  $s_t$  is not a terminal state do
8:       With probability  $\epsilon$  select a random action  $a_t$ 
9:       Otherwise, select  $a_t = \operatorname{argmax}_a Q(s_t, a)$ 
10:      Receive reward  $R_{t+1}$  and new state  $s_{t+1}$ 
11:      Store transition  $\{s_t, a_t, s_{t+1}, R_{t+1}\}$  in  $M$ 
12:      Train  $Q_\theta$  using random experiences from  $M$ 
13:       $s_t \leftarrow s_{t+1}$ 
14:    end while
15:     $\hat{Q}_{\theta-} \leftarrow Q_\theta$ 
16:  end for
17: end procedure
```

Sample Problem: Cart Pole



State $s_t = \{x, \dot{x}, \theta, \dot{\theta}\}$

Terminal states:

$$|x| > x_{max} \quad \text{or} \\ |\theta| > \theta_{max}$$

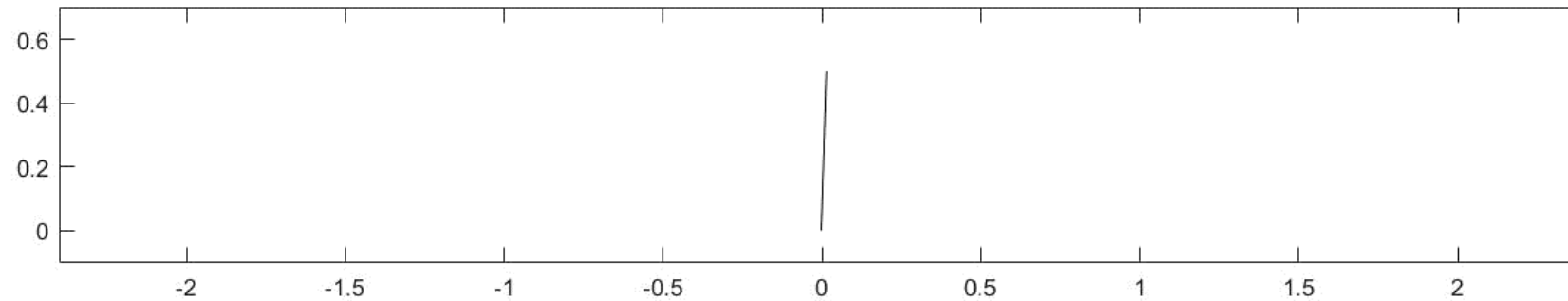
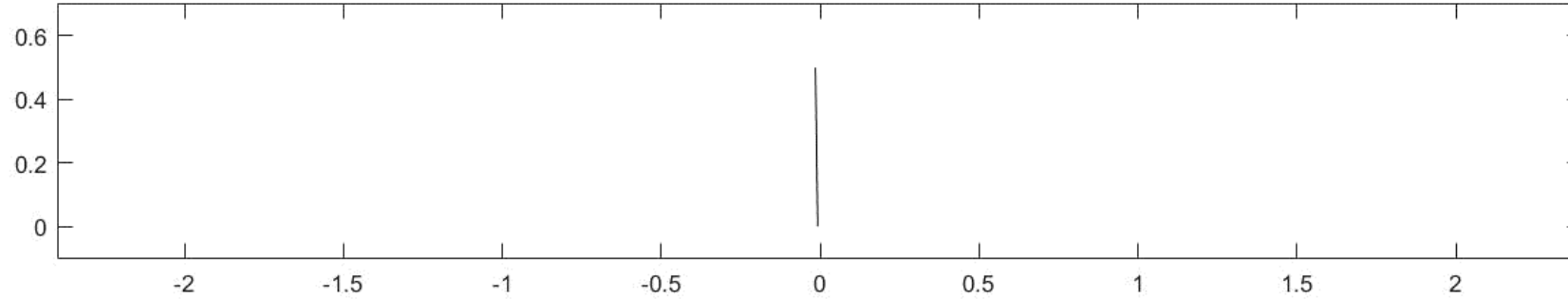
Actions:

$$a_1 \rightarrow F = 10 \text{ N} \\ a_2 \rightarrow F = -10 \text{ N}$$

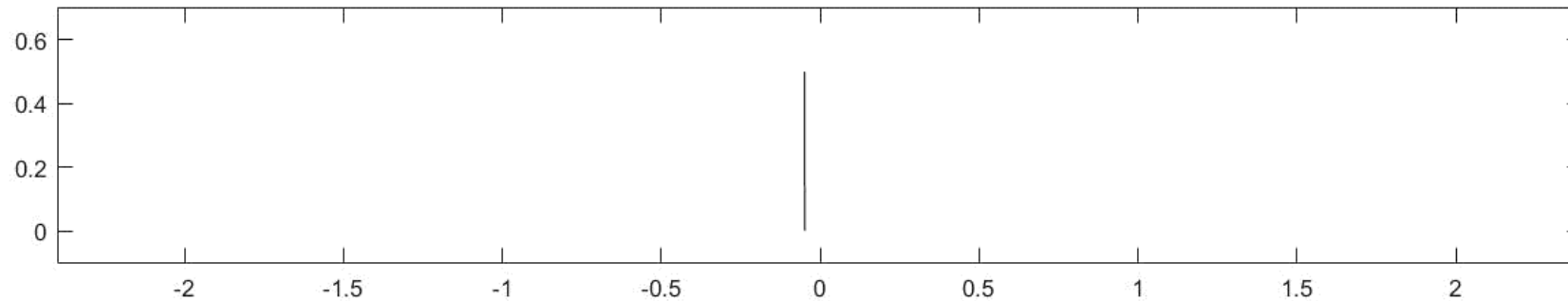
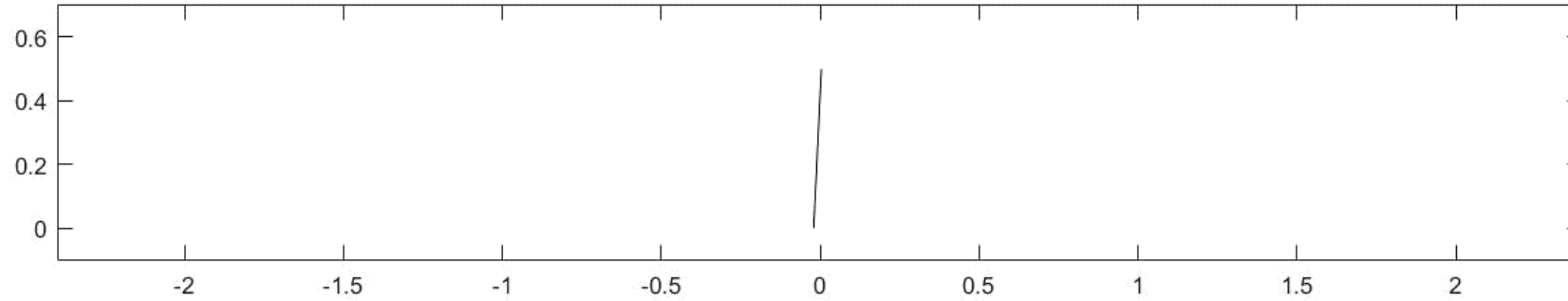
Reward:

$$R_t = 0 \text{ for a terminal state} \\ R_t = 1 \text{ for all other states}$$

Results: Getting Started



Results: Mastery



Thank you!

Any questions?

Kingma, Diederik P, and Jimmy Ba. “ADAM: A METHOD FOR STOCHASTIC OPTIMIZATION.” *CoRR*, abs/1412.6980, 2014, arxiv.org/abs/1412.6980.

Mnih, Volodymyr, et al. “Human-Level Control through Deep Reinforcement Learning.” *Nature*, vol. 518, 26 Feb. 2015, pp. 529–533., doi:<https://doi.org/10.1038/nature14236>.

Nielson, Michael A. *Neural Networks and Deep Learning*. Determination Press, 2015, neuralnetworksanddeeplearning.com/chap2.html.

Ruder, Sebastian. “An Overview of Gradient Descent Optimization Algorithms.” Sebastian Ruder, Sebastian Ruder, 19 Jan. 2016, ruder.io/optimizing-gradient-descent.

Smilkov, Daniel, and Shan Carter. “Tensorflow - Neural Network Playground.” *A Neural Network Playground*, playground.tensorflow.org/.

Increasing Stability: Target Network

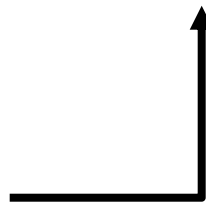
- Q-value: The quality of a state – the discounted reward expected from an action at a given state
- Bellman Equation:

$$Q^*(s_t, a_t) = R_{t+1} + \gamma \max_a Q^*(s_{t+1}, a)$$


- Framed as a cost function:

$$C = Q^*(s_t, a_t) - R_{t+1} + \gamma \max_a Q^*(s_{t+1}, a)$$

Working Network
Optimized every step



Target Network
Copied every ~500 steps



Increasing Stability: Experience Recall

Unstable: Current action immediately becomes training data

Improved: Current action is stored in memory, “recalled” later

