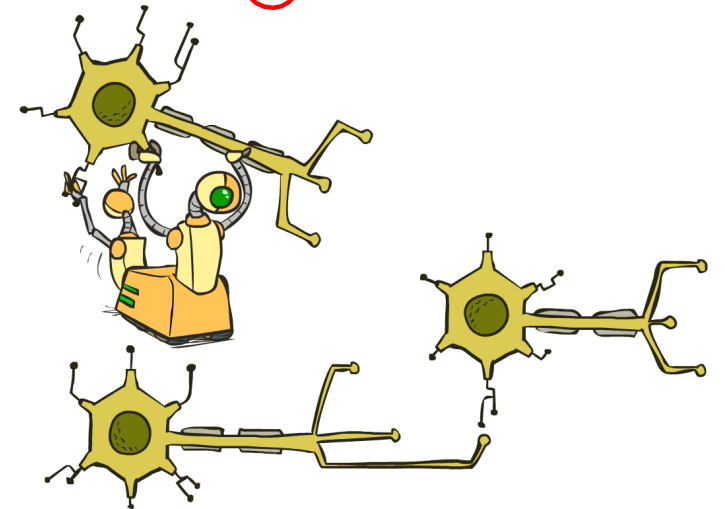# CSE 573 :
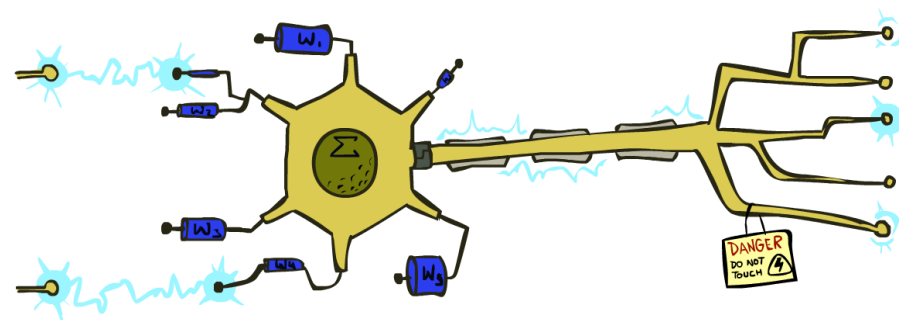# Artificial Intelligence

Hanna Hajishirzi (+ Raj)

Neural Networks

slides adapted from
Dan Klein, Pieter Abbeel ai.berkeley.edu
And Dan Weld, Luke Zettlemoyer

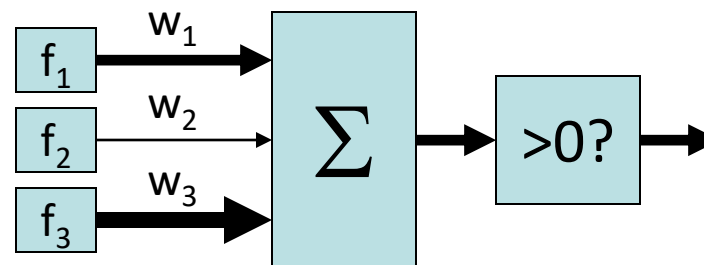# Reminder: Linear Classifiers

- Inputs are feature values
- Each feature has a weight
- Sum is the activation

$$\text{activation}_w(x) = \sum_i w_i \cdot f_i(x) = w \cdot f(x)$$

- If the activation is:
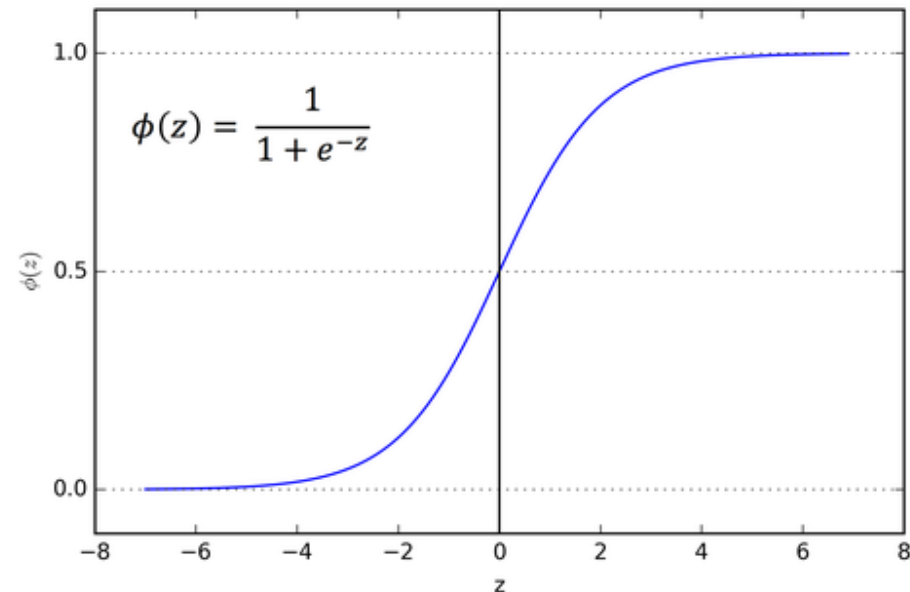  - Positive, output +1
  - Negative, output -1

# Example

- Learning the OR function: 0, 1 inputs. 0, 1 outputs
- Draw a **straight** line separating the greens and reds (two classes)

decision boundary

0,1   1,1

0,0   1,0

(0,1)

1,0

0,0

x > 1

y > 1

# Recap: How to get probabilistic decisions?

- Activation:   $z = w \cdot f(x)$

- If   $z = w \cdot f(x)$   very positive → want probability going to 1

- If   $z = w \cdot f(x)$   very negative → want probability going to 0

- Sigmoid function

$$\phi(z) = \frac{1}{1 + e^{-z}}$$



$\phi(z) = \frac{1}{1 + e^{-z}}$

# Best w?

- Maximum likelihood estimation:

$$\max_{w} \; ll(w) = \max_{w} \sum_{i} \log P(y^{(i)} | x^{(i)}; w)$$
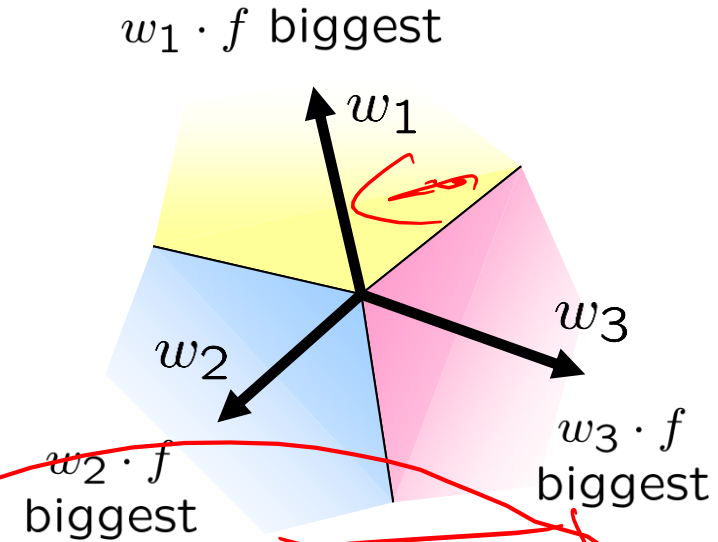
with:

$$P(y^{(i)} = +1 | x^{(i)}; w) = \frac{1}{1 + e^{-w \cdot f(x^{(i)})}}$$

$$P(y^{(i)} = -1 | x^{(i)}; w) = 1 - \frac{1}{1 + e^{-w \cdot f(x^{(i)})}}$$

**= Logistic Regression**

# Recap: Multiclass Logistic Regression

- **Multi-class linear classification**

  - A weight vector for each class:  $w_y$

  - Score (activation) of a class y:  $w_y \cdot f(x)$

  - Prediction w/highest score wins:  $y = \arg\max_y \ w_y \cdot f(x)$

$$w_1 \cdot f \text{ biggest}$$

$$w_2 \cdot f \text{ biggest} \qquad w_3 \cdot f \text{ biggest}$$

- **How to make the scores into probabilities?**

$$z_1, z_2, z_3 \rightarrow \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}, \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}}, \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$\underbrace{\phantom{z_1, z_2, z_3}}_{\text{original activations}} \qquad\qquad\qquad \underbrace{\phantom{\frac{e^{z_3}}{e^{z_1}}}}_{\text{softmax activations}}$$

# Best w?

- Maximum likelihood estimation:

$$\max_w \; ll(w) = \max_w \; \sum_i \log P(y^{(i)}|x^{(i)}; w)$$

with: $\qquad P(y^{(i)}|x^{(i)}; w) = \dfrac{e^{w_{y^{(i)}} \cdot f(x^{(i)})}}{\sum_y e^{w_y \cdot f(x^{(i)})}}$

softmax

$\dfrac{1}{1 + e^{w \cdot f(x^{(i)})}}$

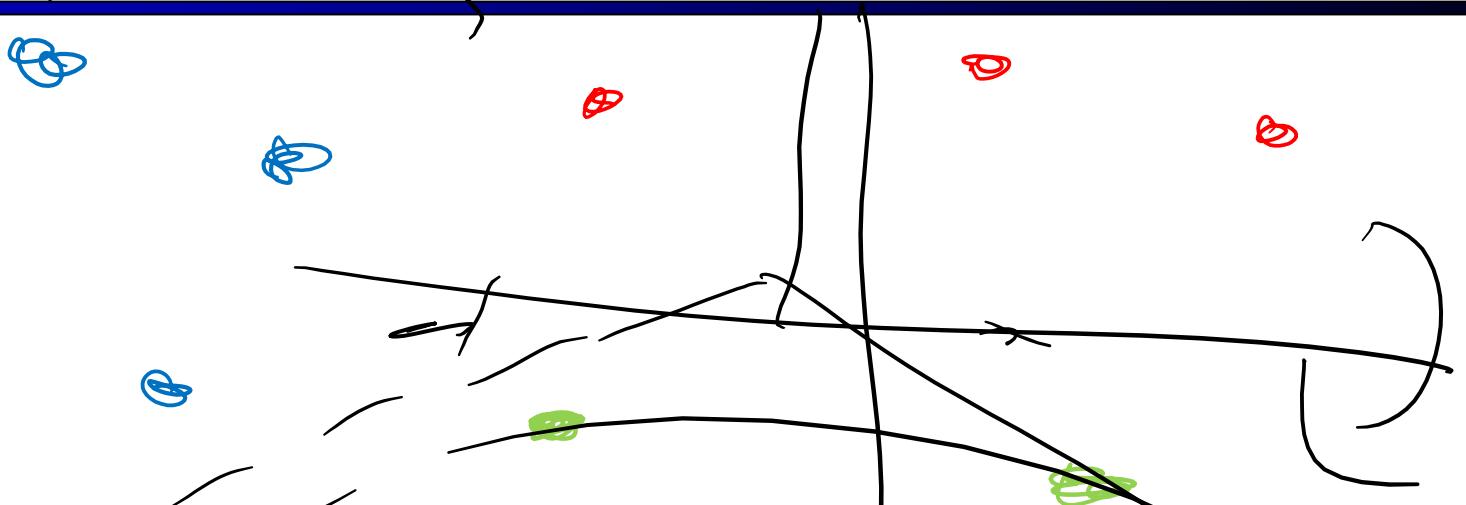**= Multi-Class Logistic Regression**

# Optimization

- Optimization
  - i.e., how do we solve:

$$\max_w \; ll(w) = \max_w \left( \sum_i \log P(y^{(i)} | x^{(i)}; w) \right)$$

$$\begin{bmatrix} m \\ b \end{bmatrix}$$

$$y = mx + b$$

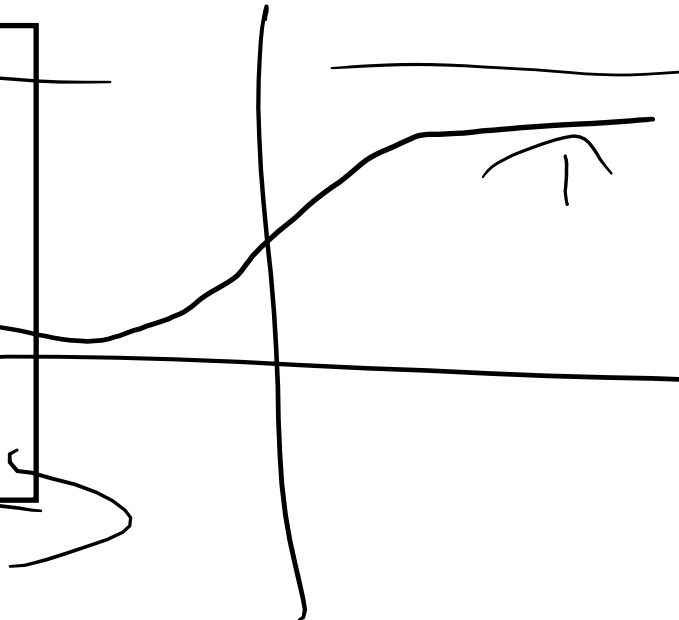# Hill Climbing

- **simple, general idea**
  - Start wherever
  - Repeat: move to the best neighboring state
  - If no neighbors better than current, quit

- **What's particularly tricky when hill-climbing for multiclass logistic regression?**
  - Optimization over a continuous space
    - Infinitely many neighbors!
    - How to do this efficiently?

# Optimization Procedure: Gradient Ascent

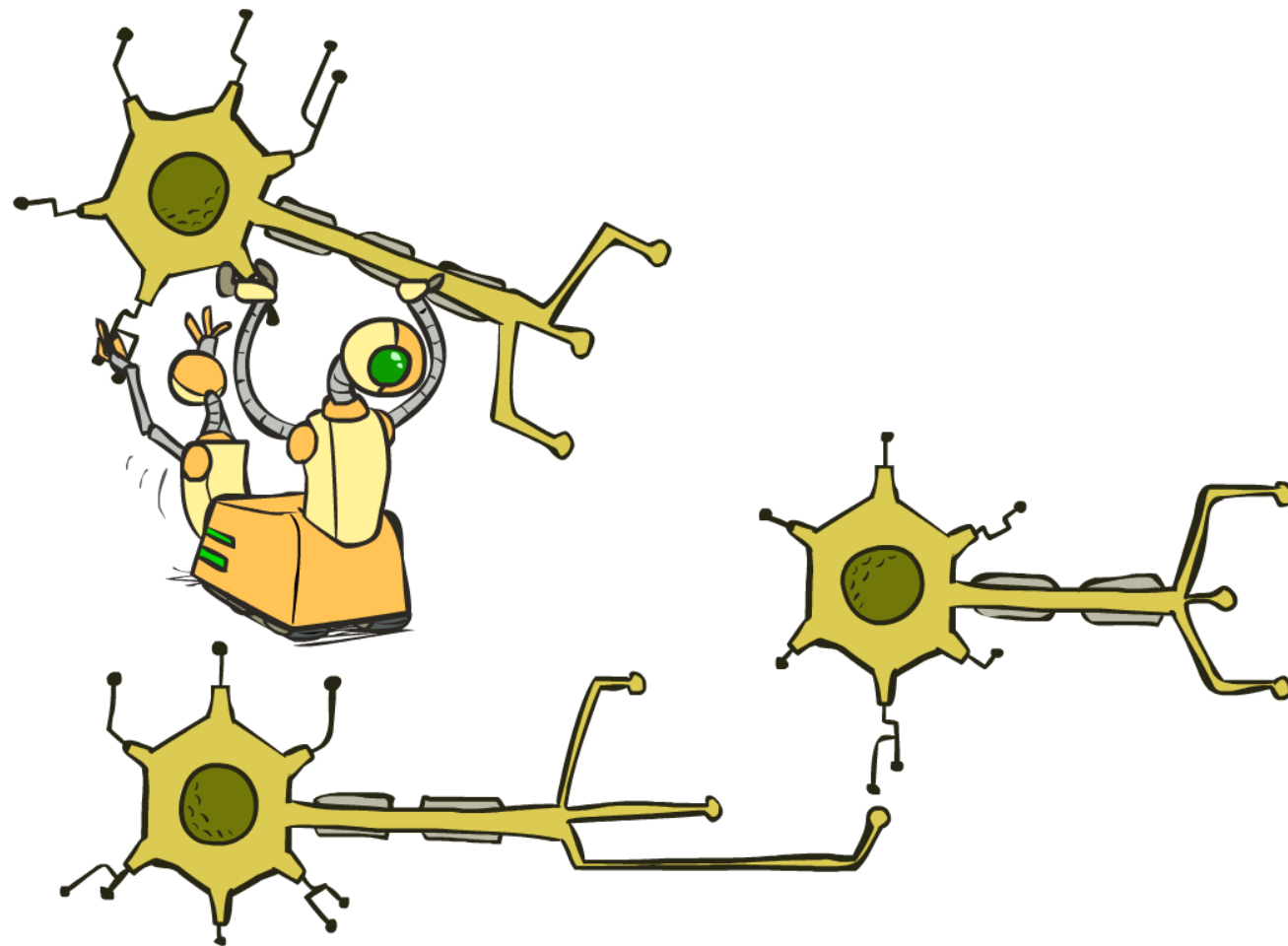- init $w$
- for iter = 1, 2, …

$$w \leftarrow w + \alpha * \nabla g(w)$$

- $\alpha$ : learning rate --- tweaking parameter that needs to be chosen carefully

- How? Try multiple choices

  - Crude rule of thumb: update changes $w$ about 0.1 – 1 %
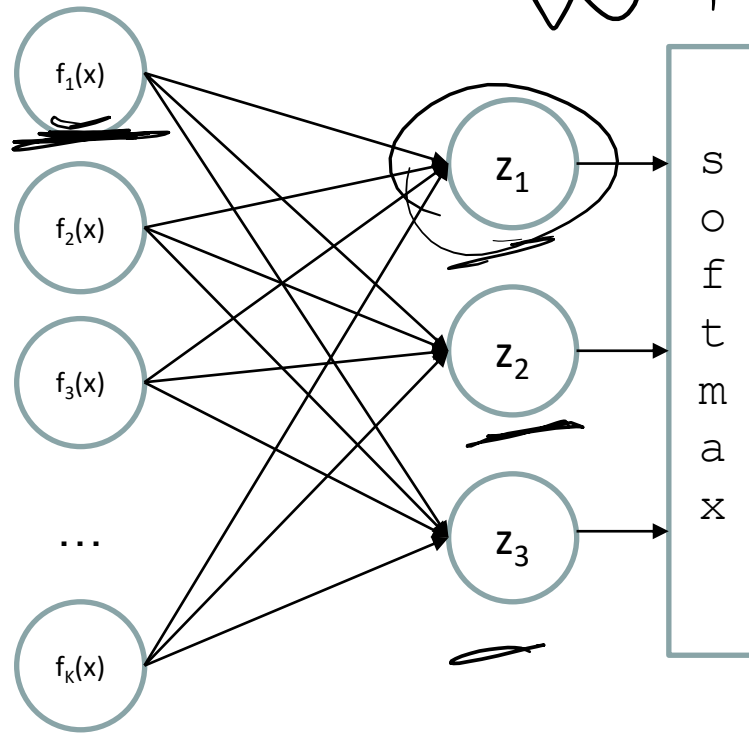
# How about computing all the derivatives?

- We'll talk about that once we covered neural networks, which are a generalization of logistic regression

# Neural Networks

# Multi-class Logistic Regression



■ = special case of neural network

$\begin{bmatrix} b_1(x) \\ \vdots \\ b_n(x) \end{bmatrix} = z_1$

$w_1$

s
o
f
t
m
a
x

$P(y_1|x;w) = \dfrac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}$

$P(y_2|x;w) = \dfrac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}}$

$P(y_3|x;w) = \dfrac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}}$

$y^{(0)}$
$y^{(1)}$
$y^{(2)}$

# Deep Neural Network = Also learn the features!
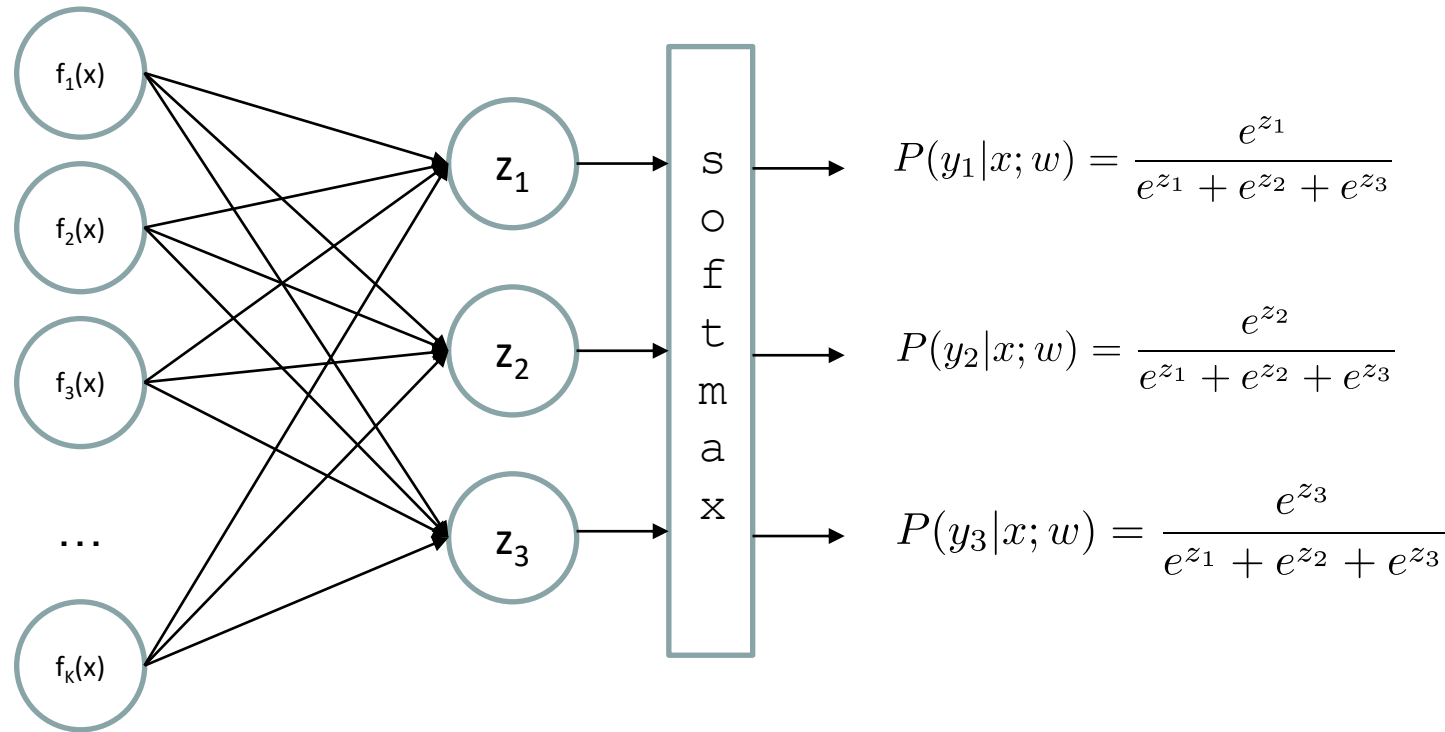


$$P(y_1|x;w) = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$P(y_2|x;w) = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$P(y_3|x;w) = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

# Deep Neural Network = Also learn the features!



$$z_i^{(k)} = g\left(\sum_j W_{i,j}^{(k-1,k)} z_j^{(k-1)}\right)$$

**g = nonlinear activation function**

# Deep Neural Network = Also learn the features!



$$z_i^{(k)} = g(\sum_j W_{i,j}^{(k-1,k)} z_j^{(k-1)})$$

**g = nonlinear activation function**

# Common Activation Functions

## Sigmoid Function



$$g(z) = \frac{1}{1 + e^{-z}}$$

$$g'(z) = g(z)(1 - g(z))$$

## Hyperbolic Tangent



$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$g'(z) = 1 - g(z)^2$$

## Rectified Linear Unit (ReLU)



$$g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 1, & z > 0 \\ 0, & \text{otherwise} \end{cases}$$

# Why non-linear activations?

- To understand, lets try to learn the XOR function
- Draw a **straight** line through the graph such that greens are on one side and reds are on another



| Input 1 | Input 2 | Output |
|---------|---------|--------|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 1 | 0 |
| 1 | 0 | 1 |

# Deep Neural Network: Also Learn the Features!

- Training the deep neural network is just like logistic regression:

$$\max_w \quad ll(w) = \max_w \sum_i \log P(y^{(i)}|x^{(i)}; w)$$

just w tends to be a much, much larger vector ☺

→just run gradient ascent

+ stop when log likelihood of hold-out data starts to decrease

# Neural Networks Properties

- Theorem (Universal Function Approximators).  A two-layer neural network with a sufficient number of neurons can approximate any continuous function to any desired accuracy.

- Practical considerations
  - Can be seen as learning the features

  - Large number of neurons
    - Danger for overfitting
    - (hence early stopping!)

# How about computing all the derivatives?

- But neural net f is never one of those?
  - No problem: CHAIN RULE:

  If $$f(x) = g(h(x))$$

  Then $$f'(x) = g'(h(x))h'(x)$$

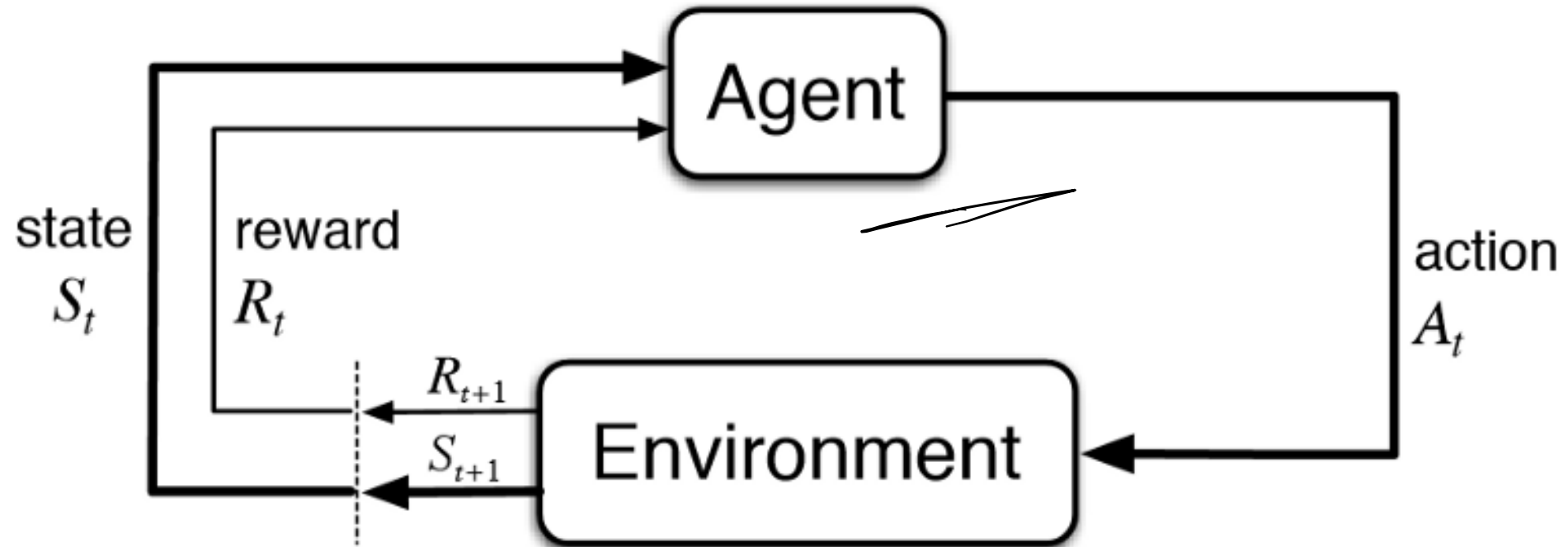  **→ Derivatives can be computed by following well-defined procedures**

# Automatic Differentiation

- Automatic differentiation software
  - e.g. Theano, TensorFlow, PyTorch, Chainer
  - Only need to program the function g(x,y,w)
  - Can automatically compute all derivatives w.r.t. all entries in w

- Need to know this exists
- How this is done?  --  outside of scope of CSE573

# Summary of Key Ideas

- Optimize probability of label given input $\quad \max\limits_{w} \quad ll(w) = \max\limits_{w} \sum\limits_{i} \log P(y^{(i)}|x^{(i)}; w)$

- Continuous optimization
  - Gradient ascent:
    - Compute steepest uphill direction = gradient (= just vector of partial derivatives)
    - Take step in the gradient direction
    - Repeat (until held-out data accuracy starts to drop = "early stopping")

- Deep neural nets
  - Last layer = still logistic regression
  - Now also many more layers before this last layer
    - = computing the features
    - → the features are learned rather than hand-designed
  - Automatic differentiation gives the derivatives efficiently (how? = outside of scope of 573)

# Deep Reinforcement Learning

# Reinforcement Learning = Learning by Interaction

AGENT GOAL: FIND POLICY π: S → A
TO MAXIMIZE REWARD

# Reinforcement Learning
# = Learning by Interaction

# Markov Decision Process

Mathematical formulation of the RL problem

Defined by: $(\mathcal{S}, \mathcal{A}, \mathcal{R}, \mathbb{P}, \gamma)$

$\mathcal{S}$ : set of possible states
$\mathcal{A}$ : set of possible actions
$\mathcal{R}$ : distribution of reward given (state, action) pair
$\mathbb{P}$ : transition probability i.e. distribution over next state given (state, action) pair
$\gamma$ : discount factor

**Markov property**: Current state completely characterizes the state of the world

$$p(r, s'|s, a) = Prob\left[R_{t+1} = r, S_{t+1} = s' \mid S_t = s, A_t = a\right]$$

# Recap: Solving for the optimal policy

**Value iteration** algorithm: Use Bellman equation as an iterative update

$$Q_{i+1}(s, a) = \mathbb{E}\left[r + \gamma \max_{a'} Q_i(s', a')|s, a\right]$$

$Q_i$ will converge to Q* as i -> infinity

# Approximate MDP solvers

What's the problem with this?
1. Not scalable. Must compute Q(s,a) for every state-action pair. If state is e.g. current game state pixels, computationally infeasible to compute for entire state space!
2. Real problems do not give you transition matrix. Again, need to account for all possibilities.

Solution:  use a function approximator to estimate Q(s,a). E.g. a neural network!

# Deep Q-learning

Remember: want to find a Q-function that satisfies the Bellman Equation:

$$Q^*(s, a) = \mathbb{E}[r + \gamma \max_{a'} Q^*(s', a') \mid s, a]$$

<span style="color:blue">Forward Pass</span>

Loss function: $L_i(\theta_i) = \mathbb{E}\left[(y_i - Q(s, a; \theta_i)^2\right]$
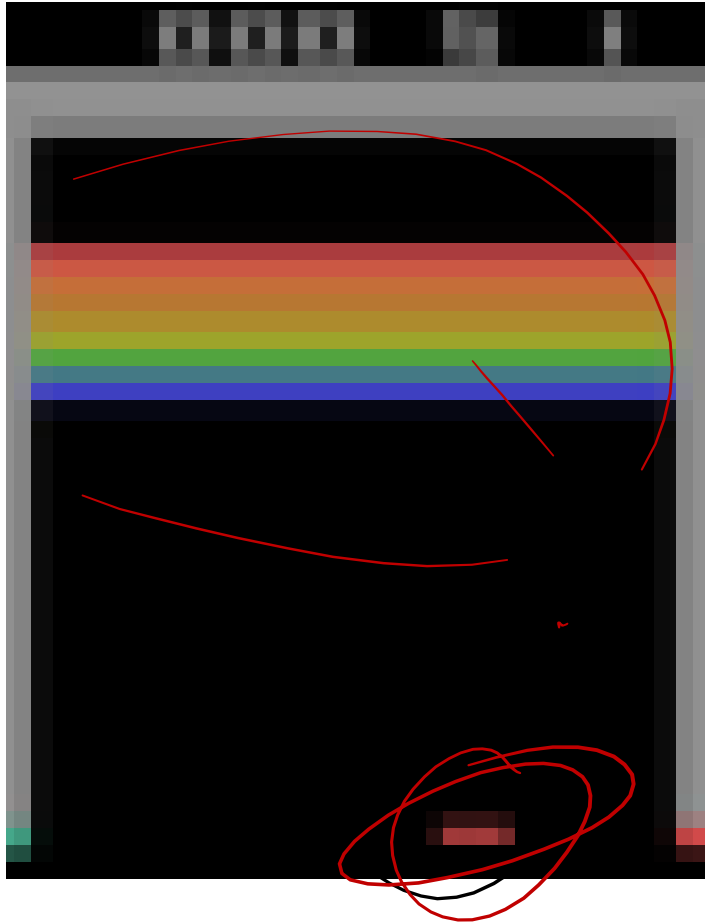
where $\quad y_i = \mathbb{E}[r + \gamma \max_{a'} Q^*(s', a') \mid s, a]$

<span style="color:blue">Backward Pass</span>

Gradient update (with respect to Q-function parameters θ):

$$\nabla_{\theta_i} L_i(\theta_i) = \mathbb{E}\left[r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) - Q(s, a; \theta_i))\nabla_{\theta_i} Q(s, a; \theta_i)\right]$$

# Breakout as an MDP



- **Objective**: Complete the game with the highest score
- **State:** Raw pixel inputs of the game state
- **Action:** Game controls e.g. Start, Left, Right, Stay
- **Reward:** Score increase/decrease at each time step

$$\rightarrow 256 \times 256 \times 3 \times^?$$

# Deep Q-Learning for Breakout

- $\pi(s) = argmax_a\ Q(s, a)$ once trained
- Next questions:
  - What is the structure of this network?
  - How do we train this network?

# Network Structure



Convolution + ReLU non-linearity     Convolution + ReLU non-linearity

To vector

action 1
action 2
action 3
action 4
action 5
action 6
action 7

Convolutional layers

Fully connected layers

# Multi-class Logistic Regression

- = special case of neural network



$$P(y_1|x;w) = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$P(y_2|x;w) = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$P(y_3|x;w) = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

# Pytorch Demo! Open those Colab Notebooks if you want to follow along!

https://colab.research.google.com/drive/1PR8wBdglJ10yikyUkU8fmjl7-WviWC_W?usp=sharing

# https://tinyurl.com/cse573-dqn

# Training the Q-network: Experience Replay

Learning from batches of consecutive samples is problematic:
- Samples are correlated => inefficient learning
- Current Q-network parameters determines next training samples (e.g. if maximizing action is to move left, training samples will be dominated by samples from left-hand size) => can lead to bad feedback loops

# Training the Q-network: Experience Replay

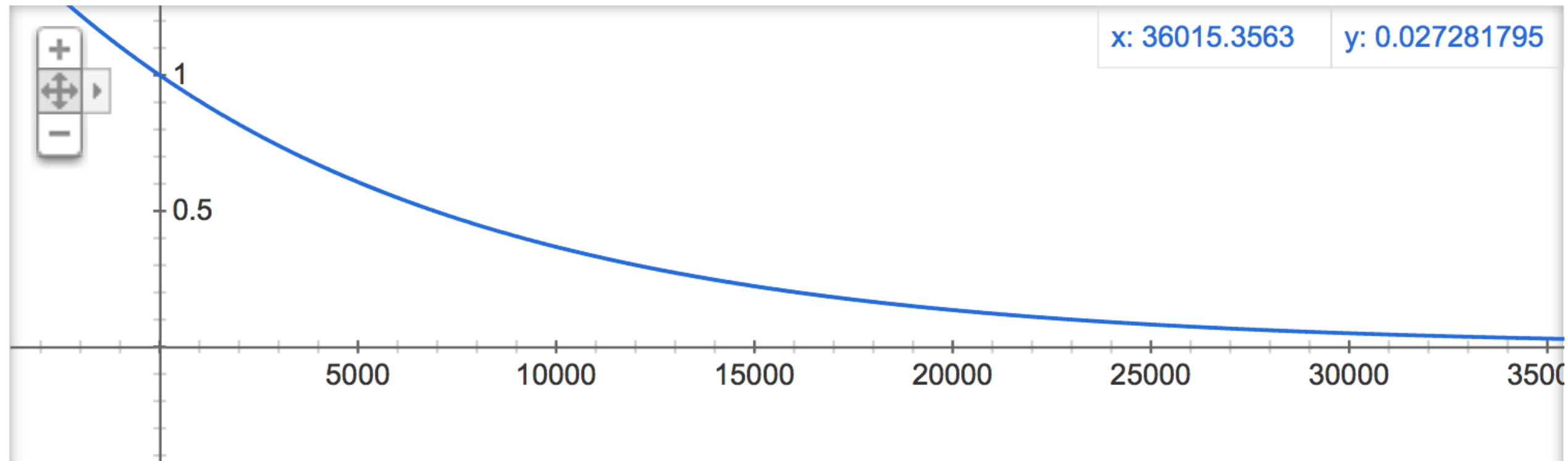Learning from batches of consecutive samples is problematic:
-   Samples are correlated => inefficient learning
-   Current Q-network parameters determines next training samples (e.g. if maximizing action is to move left, training samples will be dominated by samples from left-hand size) => can lead to bad feedback loops

Address these problems using **experience replay**
-   Continually update a **replay memory** table of transitions ($s_t$, $a_t$, $r_t$, $s_{t+1}$) as game (experience) episodes are played
-   Train Q-network on random minibatches of transitions from the replay memory, instead of consecutive samples

37

# Training the Q-network: Epsilon Greedy

- Epsilon is the term that decides how often the agent randomly picks an action

**Algorithm 1: deep Q-learning with experience replay.**

Initialize replay memory $D$ to capacity $N$

Initialize action-value function $Q$ with random weights $\theta$

Initialize target action-value function $\hat{Q}$ with weights $\theta^- = \theta$

**For** episode $= 1, M$ **do**

   Initialize sequence $s_1 = \{x_1\}$ and preprocessed sequence $\phi_1 = \phi(s_1)$

   **For** $t = 1, \text{T}$ **do**

      With probability $\varepsilon$ select a random action $a_t$

      otherwise select $a_t = \text{argmax}_a\, Q(\phi(s_t), a; \theta)$

      Execute action $a_t$ in emulator and observe reward $r_t$ and image $x_{t+1}$

      Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$

      Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in $D$

      Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from $D$

$$
\text{Set } y_j = \begin{cases} r_j & \text{if episode terminates at step } j+1 \\ r_j + \gamma\, \max_{a'} \hat{Q}(\phi_{j+1}, a'; \theta^-) & \text{otherwise} \end{cases}
$$

      Perform a gradient descent step on $\left(y_j - Q(\phi_j, a_j; \theta)\right)^2$ with respect to the network parameters $\theta$

      Every $C$ steps reset $\hat{Q} = Q$

   **End For**

**End For**