# Reinforcement Learning:
# Markov Decision Processes

Matt Gormley
Lecture 15
Oct.14, 2019

# Reminders

- **Homework 5: Neural Networks**
  - **Out: Fri, Oct. 11**
  - **Due: Fri, Oct. 25 at 11:59pm**
- **Recitation:**
  - **Thu, Oct 17th at 7:30pm – 8:30pm in GHC 4401**
  - **(also available on Panopto)**
- **Today's In-Class Poll**
  - **http://p15.mlcourse.org**

# Q&A

# OTHER APPROACHES TO DIFFERENTIATION
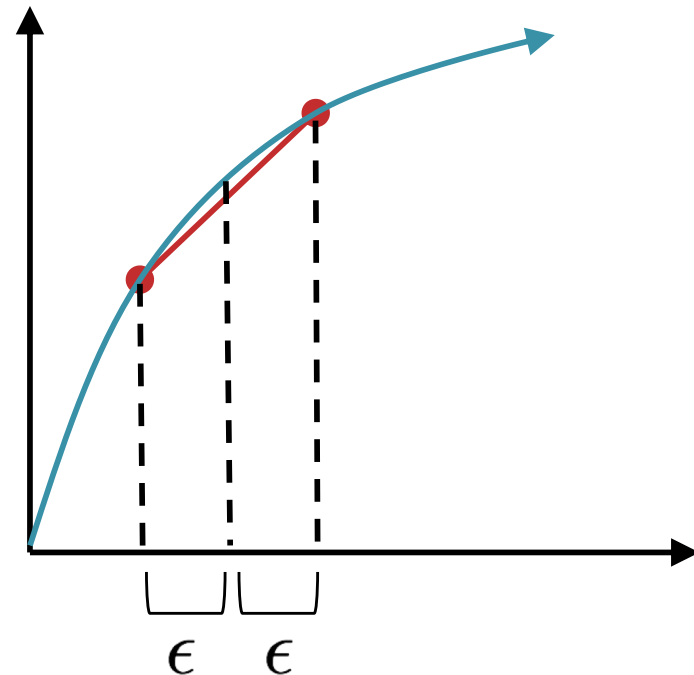
# Finite Difference Method

The *centered* finite difference approximation is:

$$\frac{\partial}{\partial \theta_i} J(\boldsymbol{\theta}) \approx \frac{(J(\boldsymbol{\theta} + \epsilon \cdot \boldsymbol{d}_i) - J(\boldsymbol{\theta} - \epsilon \cdot \boldsymbol{d}_i))}{2\epsilon} \qquad (1)$$

where $\boldsymbol{d}_i$ is a 1-hot vector consisting of all zeros except for the $i$th entry of $\boldsymbol{d}_i$, which has value 1.

**Notes:**

- Suffers from issues of floating point precision, in practice
- Typically only appropriate to use on small examples with an appropriately chosen epsilon

# Symbolic Differentiation

**Differentiation Quiz #1:**

Suppose x = 2 and z = 3, what are dy/dx and dy/dz for the function below? **Round your answer to the nearest integer.**

$$y = \exp(xz) + \frac{xz}{\log(x)} + \frac{\sin(\log(x))}{xz}$$

**Answer:** *Answers below are in the form [dy/dx, dy/dz]*

A.  [42, -72]          E.  [1208, 810]
B.  [72, -42]          F.  [810, 1208]
C.  [100, 127]         G.  [1505, 94]
D.  [127, 100]         H.  [94, 1505]

# Symbolic Differentiation

## Differentiation Quiz #2:
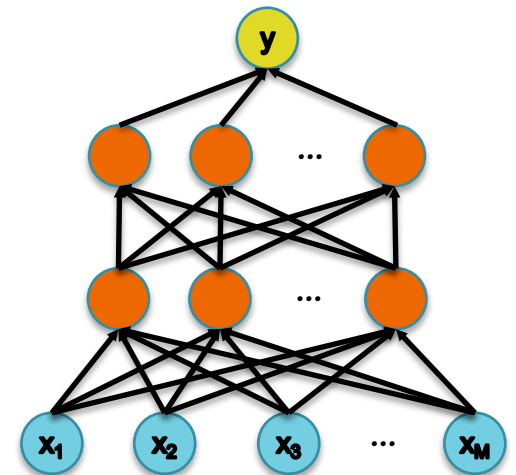
A neural network with 2 hidden layers can be written as:

$$y = \sigma(\boldsymbol{\beta}^T \sigma((\boldsymbol{\alpha}^{(2)})^T \sigma((\boldsymbol{\alpha}^{(1)})^T \mathbf{x})))$$

where $y \in \mathbb{R}$, $\mathbf{x} \in \mathbb{R}^{D^{(0)}}$, $\boldsymbol{\beta} \in \mathbb{R}^{D^{(2)}}$ and $\boldsymbol{\alpha}^{(i)}$ is a $D^{(i)} \times D^{(i-1)}$ matrix. Nonlinear functions are applied elementwise:

$$\sigma(\mathbf{a}) = [\sigma(a_1), \ldots, \sigma(a_K)]^T$$

Let $\sigma$ be sigmoid: $\sigma(a) = \frac{1}{1+exp-a}$

What is $\frac{\partial y}{\partial \beta_j}$ and $\frac{\partial y}{\partial \alpha_j^{(i)}}$ for all $i, j$.

# Summary

1. **Neural Networks**…
   - provide a way of learning features
   - are highly nonlinear prediction functions
   - (can be) a highly parallel network of logistic regression classifiers
   - discover useful hidden representations of the input

2. **Backpropagation**…
   - provides an efficient way to compute gradients
   - is a special case of reverse-mode automatic differentiation

# Backprop Objectives

*You should be able to...*

- Construct a computation graph for a function as specified by an algorithm
- Carry out the backpropagation on an arbitrary computation graph
- Construct a computation graph for a neural network, identifying all the given and intermediate quantities that are relevant
- Instantiate the backpropagation algorithm for a neural network
- Instantiate an optimization method (e.g. SGD) and a regularizer (e.g. L2) when the parameters of a model are comprised of several matrices corresponding to different layers of a neural network
- Apply the empirical risk minimization framework to learn a neural network
- Use the finite difference method to evaluate the gradient of a function
- Identify when the gradient of a function can be computed at all and when it can be computed efficiently

# LEARNING PARADIGMS

# Learning Paradigms

| Paradigm | Data | |
| --- | --- | --- |
| Supervised | $\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^{N}$ | $\mathbf{x} \sim p^*(\cdot)$ and $y = c^*(\cdot)$ |
| $\hookrightarrow$ Regression | $y^{(i)} \in \mathbb{R}$ | |
| $\hookrightarrow$ Classification | $y^{(i)} \in \{1, \ldots, K\}$ | |
| $\hookrightarrow$ Binary classification | $y^{(i)} \in \{+1, -1\}$ | |
| $\hookrightarrow$ Structured Prediction | $\mathbf{y}^{(i)}$ is a vector | |

# Learning Paradigms

| Paradigm | Data | |
| --- | --- | --- |
| Supervised | $\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^{N}$ | $\mathbf{x} \sim p^{*}(\cdot)$ and $y = c^{*}(\cdot)$ |
| $\hookrightarrow$ Regression | $y^{(i)} \in \mathbb{R}$ | |
| $\hookrightarrow$ Classification | $y^{(i)} \in \{1, \ldots, K\}$ | |
| $\hookrightarrow$ Binary classification | $y^{(i)} \in \{+1, -1\}$ | |
| $\hookrightarrow$ Structured Prediction | $\mathbf{y}^{(i)}$ is a vector | |
| Unsupervised | $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^{N}$ | $\mathbf{x} \sim p^{*}(\cdot)$ |

# Learning Paradigms

| Paradigm | Data | |
|---|---|---|
| Supervised | $\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^{N}$ | $\mathbf{x} \sim p^*(\cdot)$ and $y = c^*(\cdot)$ |
| $\hookrightarrow$ Regression | $y^{(i)} \in \mathbb{R}$ | |
| $\hookrightarrow$ Classification | $y^{(i)} \in \{1, \ldots, K\}$ | |
| $\hookrightarrow$ Binary classification | $y^{(i)} \in \{+1, -1\}$ | |
| $\hookrightarrow$ Structured Prediction | $\mathbf{y}^{(i)}$ is a vector | |
| Unsupervised | $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^{N}$ | $\mathbf{x} \sim p^*(\cdot)$ |
| Semi-supervised | $\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^{N_1} \cup \{\mathbf{x}^{(j)}\}_{j=1}^{N_2}$ | |

# Learning Paradigms

| Paradigm | Data |
| --- | --- |
| Supervised | $\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^{N}$ $\quad$ $\mathbf{x} \sim p^*(\cdot)$ and $y = c^*(\cdot)$ |
| $\hookrightarrow$ Regression | $y^{(i)} \in \mathbb{R}$ |
| $\hookrightarrow$ Classification | $y^{(i)} \in \{1, \ldots, K\}$ |
| $\hookrightarrow$ Binary classification | $y^{(i)} \in \{+1, -1\}$ |
| $\hookrightarrow$ Structured Prediction | $\mathbf{y}^{(i)}$ is a vector |
| Unsupervised | $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^{N}$ $\quad$ $\mathbf{x} \sim p^*(\cdot)$ |
| Semi-supervised | $\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^{N_1} \cup \{\mathbf{x}^{(j)}\}_{j=1}^{N_2}$ |
| Online | $\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), (\mathbf{x}^{(3)}, y^{(3)}), \ldots\}$ |

# Learning Paradigms

| Paradigm | Data | |
|---|---|---|
| Supervised | $\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^{N}$ | $\mathbf{x} \sim p^*(\cdot)$ and $y = c^*(\cdot)$ |
| $\hookrightarrow$ Regression | $y^{(i)} \in \mathbb{R}$ | |
| $\hookrightarrow$ Classification | $y^{(i)} \in \{1, \ldots, K\}$ | |
| $\hookrightarrow$ Binary classification | $y^{(i)} \in \{+1, -1\}$ | |
| $\hookrightarrow$ Structured Prediction | $\mathbf{y}^{(i)}$ is a vector | |
| Unsupervised | $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^{N}$ | $\mathbf{x} \sim p^*(\cdot)$ |
| Semi-supervised | $\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^{N_1} \cup \{\mathbf{x}^{(j)}\}_{j=1}^{N_2}$ | |
| Online | $\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), (\mathbf{x}^{(3)}, y^{(3)}), \ldots\}$ | |
| Active Learning | $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^{N}$ and can query $y^{(i)} = c^*(\cdot)$ at a cost | |

# Learning Paradigms

| Paradigm | Data |
|---|---|
| Supervised | $\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^{N}$     $\mathbf{x} \sim p^*(\cdot)$ and $y = c^*(\cdot)$ |
| $\hookrightarrow$ Regression | $y^{(i)} \in \mathbb{R}$ |
| $\hookrightarrow$ Classification | $y^{(i)} \in \{1, \ldots, K\}$ |
| $\hookrightarrow$ Binary classification | $y^{(i)} \in \{+1, -1\}$ |
| $\hookrightarrow$ Structured Prediction | $\mathbf{y}^{(i)}$ is a vector |
| Unsupervised | $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^{N}$     $\mathbf{x} \sim p^*(\cdot)$ |
| Semi-supervised | $\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^{N_1} \cup \{\mathbf{x}^{(j)}\}_{j=1}^{N_2}$ |
| Online | $\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), (\mathbf{x}^{(3)}, y^{(3)}), \ldots\}$ |
| Active Learning | $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^{N}$ and can query $y^{(i)} = c^*(\cdot)$ at a cost |
| Imitation Learning | $\mathcal{D} = \{(s^{(1)}, a^{(1)}), (s^{(2)}, a^{(2)}), \ldots\}$ |

# Learning Paradigms

| Paradigm | Data |
|---|---|
| Supervised | $\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^{N}$     $\mathbf{x} \sim p^*(\cdot)$ and $y = c^*(\cdot)$ |
| $\hookrightarrow$ Regression | $y^{(i)} \in \mathbb{R}$ |
| $\hookrightarrow$ Classification | $y^{(i)} \in \{1, \ldots, K\}$ |
| $\hookrightarrow$ Binary classification | $y^{(i)} \in \{+1, -1\}$ |
| $\hookrightarrow$ Structured Prediction | $\mathbf{y}^{(i)}$ is a vector |
| Unsupervised | $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^{N}$     $\mathbf{x} \sim p^*(\cdot)$ |
| Semi-supervised | $\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^{N_1} \cup \{\mathbf{x}^{(j)}\}_{j=1}^{N_2}$ |
| Online | $\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), (\mathbf{x}^{(3)}, y^{(3)}), \ldots\}$ |
| Active Learning | $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^{N}$ and can query $y^{(i)} = c^*(\cdot)$ at a cost |
| Imitation Learning | $\mathcal{D} = \{(s^{(1)}, a^{(1)}), (s^{(2)}, a^{(2)}), \ldots\}$ |
| Reinforcement Learning | $\mathcal{D} = \{(s^{(1)}, a^{(1)}, r^{(1)}), (s^{(2)}, a^{(2)}, r^{(2)}), \ldots\}$ |

# REINFORCEMENT LEARNING

# Examples of Reinforcement Learning

- How should a robot behave so as to optimize its "performance"? (Robotics)

- How to automate the motion of a helicopter? (Control Theory)

- How to make a good chess-playing program? (Artificial Intelligence)

# Autonomous Helicopter

Video:

https://www.youtube.com/watch?v=VCdxqnofcnE

# Robot in a room



actions: UP, DOWN, LEFT, RIGHT

UP

80%     move UP
10%     move LEFT
10%     move RIGHT

- reward +1 at [4,3], -1 at [4,2]
- reward -0.04 for each step

- what's the strategy to achieve max reward?
- what if the actions were NOT deterministic?

# History of Reinforcement Learning

- Roots in the psychology of animal learning (Thorndike,1911).

- Another independent thread was the problem of optimal control, and its solution using dynamic programming (Bellman, 1957).

- Idea of temporal difference learning (on-line method), e.g., playing board games (Samuel, 1959).

- A major breakthrough was the discovery of Q-learning (Watkins, 1989).

# What is special about RL?

- RL is learning how to map states to actions, so as to <span style="color:red">maximize</span> a numerical <span style="color:red">reward</span> over time.

- Unlike other forms of learning, it is a multistage decision-making process (often <span style="color:red">Markovian</span>).

- An RL agent must learn by <span style="color:red">trial-and-error</span>. (Not entirely supervised, but interactive)

- Actions may affect not only the immediate reward but also subsequent rewards (<span style="color:red">Delayed effect</span>).

# Elements of RL

- A policy
    - A map from state space to action space.
    - May be stochastic.
- A reward function
    - It maps each state (or, state-action pair) to a real number, called reward.
- A value function
    - Value of a state (or, state-action pair) is the total expected reward, starting from that state (or, state-action pair).
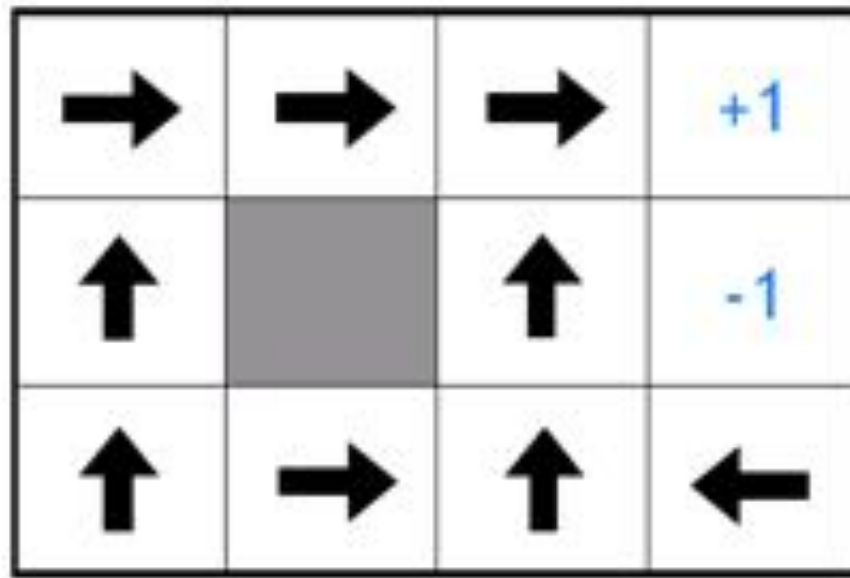
# Policy

# Reward for each step -2

# Reward for each step: -0.1

# The Precise Goal

- To find a policy that maximizes the Value function.
  - transitions and rewards usually not available

- There are different approaches to achieve this goal in various situations.

- Value iteration and Policy iteration are two more classic approaches to this problem. But essentially both are dynamic programming.

- Q-learning is a more recent approaches to this problem. Essentially it is a temporal-difference method.

# MARKOV DECISION PROCESSES

# Markov Decision Process

- For **supervised learning** the **PAC learning framework** provided assumptions about where our data came from:

$$\mathbf{x} \sim p^*(\cdot) \text{ and } y = c^*(\cdot)$$

- For **reinforcement learning** we assume our data comes from a **Markov decision process** (MDP)

# Markov Decision Process

*Whiteboard*

- Components: states, actions, state transition probabilities, reward function
- Markovian assumption
- MDP Model
- MDP Goal: Infinite-horizon Discounted Reward
- deterministic vs. nondeterministic MDP
- deterministic vs. stochastic policy