# WELCOME AND RECAP

- This is our third session on reinforcement learning

- Sessions I and II focused on some code samples and terminology using python notebooks mimicking the contextual bandit problem (finding the best slot machine to play from a collection of slot machines).

- Tonight's session presents a conceptual foundation for some key principles used in larger scale and more complex learning scenarios

- Prior code and presentations can be found at:
  - https://github.com/jimwill3/NY-AZML-Meetup/tree/CNTK/RL

- In June we will take a first look at an Azure Reinforcement Learning Service called "Personalizer"

https://aischool.microsoft.com/en-us/machine-learning/learning-paths (free self-paced training)

# Markov Decision Processes

Michael H. Evangelista - michael@powerof.to

2019-05-23

Guest WiFi: MSFTGUEST -> Event Code -> "msevent242zc"

Jupyter Notebook:
**https://notebooks.azure.com/mhe500/projects/gridworld**

**(Optional if you want to follow along)**

# What We'll Cover

- Reinforcement learning overview

- MDP basics & describing a simple game as an MDP

- Finding the best actions for a given MDP using Value Iteration

- Value Iteration implementation in Python, solving the simple game

- Extending to more complex problems (like Atari games)

# Reinforcement Learning

*"Reinforcement learning **is learning what to do**— how to map situations to actions—so as to maximize a numerical reward signal. The learner is not told which actions to take, but instead must discover which actions yield the most reward by **trying them**. In the most interesting and challenging cases, actions may affect not only the immediate reward but also the next situation and, through that, **all subsequent rewards**. These two characteristics— **trial-and-error search** and **delayed reward**—are the two most important distinguishing features of reinforcement learning."*

Sutton, Richard S., and Andrew G. Barto. "Reinforcement Learning: an Introduction." *Reinforcement Learning: an Introduction*, The MIT Press., 2018, pp. 1–2.
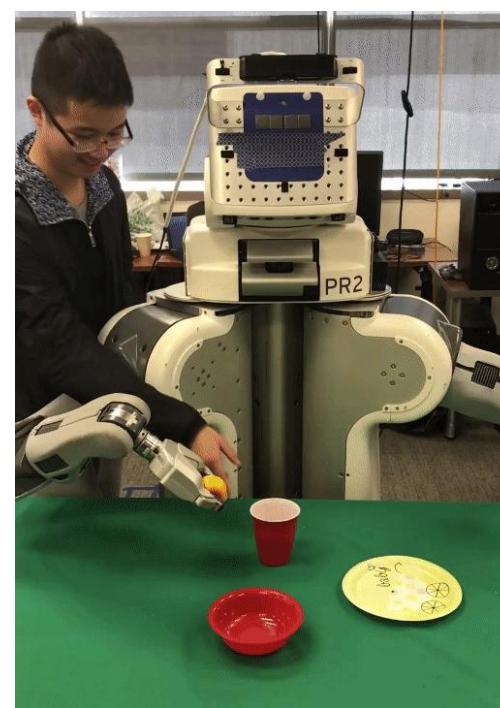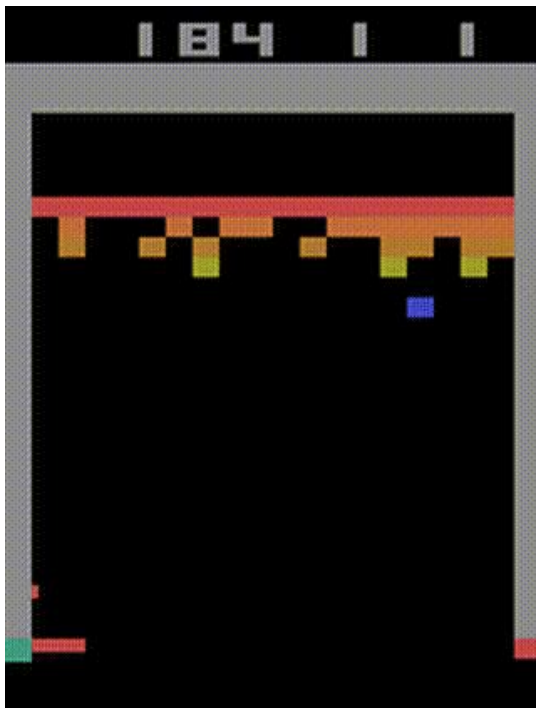Available as free PDF at http://incompleteideas.net/book/the-book-2nd.html

# Why Bother with This?
# Let's get to the CNTK/Tensorflow/MXNet code!

# The Agent-Environment Interface

1. Agent receives state info $S_t$ from Environment at time $t$

2. Agent chooses action $A_t$

3. Environment transitions to and emits $S_{t+1}$

4. Environment emits numerical reward $R_{t+1}$

5. $t = t + 1$

6. Repeat (until task ends or forever)

**Agent's goal:** Maximize Return (cumulative reward from time $t$ and continuing):

$$G_t = \sum_{k=0..\infty} \gamma^k R_{t+k+1}$$

# Markov Decision Process

- Defines sequential decision making problem
- Possible environment states
- Allowed actions
- *Dynamics* – how an action leads from one state to another and how the agent is rewarded

Andrey Markov (1856-1922)

# Simple MDP Illustration

# State vs. Observation

**State**

- A Markov state $s$ "fully characterizes" the future (Markov Property)

- No matter what happened before $s$, the expected outcome (rewards, transitions) are the same going forward

- More simply: "the past doesn't affect the future"

**Observation**

- What you can "see"

- Different states may produce the same observation

- If you don't know ("see") true state, your problem is Partially Observed MDP (POMDP)

- **<u>Not dealing with POMDP today (assume we know true state)</u>**

Agent


Observation*

**\* From human POV. Luckily, for Robocop this is an MDP, not POMDP!**


State


Choose Best Action

**Image Source:** https://www.youtube.com/watch?v=Md14H_qD8iQ

# Robot Grid World

- Robot can move to any adjacent square

- Environment is stochastic (meaning, with some probability the robot may "slip" and go to the left/right instead of forward in desired direction)

- Attempt to move to the rock or past the edge (whether intentionally or by slip) leaves you in current square

- Get the diamond for a score of +1 (win the game)

- Fall into the fiery pit, get a score of –1 (and die)

# Grid World as an MDP

- Set of States - $S$

- Set of Actions - $A$

- Transition Probabity Function - $p(s' \mid s, a)$
  - Probability of ending up in state $s'$, when taking action a from state $s$.

- Reward Function - $r(s, a, s')$
  - Reward gained when action a causes transition from state $s$ to state $s'$.

- Initial State - $s_0$

- Time Horizon - $H$ (can be $\infty$)

- Discount Factor - $\gamma \in [0, 1]$

$$G_t = \sum_{k=0..H} \gamma^k R_{t+k+1}$$

Image Credit: Abbeel, 2017

# Exercise 1: Grid World Transition Probabilities

| s' | p( s' \| s=(2,2), a=up, slip=0.20) |
|---|---|
| (0,0) | 0 |
| (0,1) | 0 |
| (0,2) | 0 |
| (0,3) | 0 |
| (1,0) | 0 |
| (1,1) | 0 |
| (1,2) | 0.80 |
| (1,3) | 0 |
| (2,0) | 0 |
| (2,1) | 0.10 |
| (2,2) | 0 |
| (2,3) | 0.10 |

|   | 0 | 1 | 2 | 3 |
|---|---|---|---|---|
| 0 | (0,0) | (0,1) | (0,2) | (0,3) R=+1 |
| 1 | (0,1) | Rock | (1,2) | (1,3) R=-1 |
| 2 | (0,2) Start | (2,1) | (2,2) | (2,3) |

Attempt to move UP from here

# Remember Bandits (from sessions 1 & 2)?

- A bandit is a one-state MDP...
- ... with stochastic (random) rewards



State $s_0$

Pull Handle #1

Pull Handle #2

Move to $s_0$ with probabilty $p(s_0 \mid S=s_0, A=\text{pull \#1}) = 1.0$
Reward $= r$ with probability $p(R=r \mid S=s_0, A=\text{pull \#1}, S'=s_0)$

Move to $s_0$ with probabilty $p(s_0 \mid S=s_0, A=\text{pull \#2}) = 1.0$
Reward $= r$ with probability $p(R=r \mid S=s_0, A=\text{pull \#2}, S'=s_0)$

# Solving the MDP

Find a policy, $\pi^*$, that yields the best return:



$\pi$

- $\pi^* = \text{argmax}_\pi \, E_\pi[G_t \mid \pi, s_0]$
- $\phantom{\pi^*} = \text{argmax}_\pi \, E_\pi[\sum_{t=0..H} \gamma^t \, r(s_t, a_t, s_{t+1}) \mid \pi, s_0]$

Wait – what the heck does this mean?

- A policy, $\pi$, defines an agent's behavior
- $\pi^*$ is the optimal policy. When followed, $\pi^*$ should, in expectation (i.e., on average), produce the highest possible return starting from any state

# The State- and Action-Value Function Definitions

**State-Value Function under policy $\pi$:**

Intuition: Expected return starting in state $s$ and following policy $\pi$ thereafter.

$$v_\pi(s) = E_\pi[G_t \mid S_t=s]$$
$$= E_\pi[\textstyle\sum_{k=0..\infty} \gamma^k R_{t+k+1} \mid S_t=s] \qquad\qquad \text{(By definition of } G_t)$$

**Action-Value Function under policy $\pi$:**

Intuition: Expected return starting in state $s$, <u>taking action $a$ as the first step</u>, then following policy $\pi$ thereafter.

$$q_\pi(s, a) = E_\pi[G_t \mid S_t=s, A_t=a]$$
$$= E_\pi[\textstyle\sum_{k=0..\infty} \gamma^k R_{t+k+1} \mid S_t=s, A_t=a] \qquad\qquad \text{(By definition of } G_t)$$

***Calculating $q_\pi(s)$ from $v_\pi(s)$:***

Intuition: If we know $v_\pi$ we can easily get $q_\pi$.

$$q_\pi(s, a) = E_\pi[\textstyle\sum_{k=0..\infty} \gamma^k R_{t+k+1} \mid S_t=s, A_t=a]$$
$$= E_\pi[R_{t+1} + \gamma\textstyle\sum_{k=1..\infty} \gamma^{k-1} R_{t+k+1} \mid S_t=s, A_t=a] \qquad \text{(Split summation)}$$
$$= E_\pi[R_{t+1} + \gamma v_\pi(s') \mid S_t=s, A_t=a] \qquad\qquad \text{(By definition of } v_\pi(s'))$$
$$= \textstyle\sum_{s'} p(s' \mid s, a) \, [r(s, a, s') + \gamma v_\pi(s')] \qquad\qquad \text{(By definition of expectation \& } R_{t+1})$$

# <u>Optimal</u> State- and Action-Value Function

**State-Value Function under <u>optimal policy</u> $\pi^*$:**

Intuition: Expected return starting in state $s$ and following the <u>optimal policy</u> $\pi^*$ thereafter.

$$v^*(s) = \max_{\pi} E_{\pi}[G_t \mid S_t{=}s]$$
$$= \max_{\pi} E_{\pi}[\textstyle\sum_{k=0..\infty} \gamma^k R_{t+k+1} \mid S_t{=}s]$$

**Action-Value Function under <u>optimal policy</u> $\pi^*$:**

Intuition: Expected return starting in state $s$, <u>taking action $a$ as the first step</u>, then following the <u>optimal policy</u> $\pi^*$ thereafter. **<u>If we know $q^*(s, a)$, we effectively have our policy</u>**.

$$q^*(s, a) = \max_{\pi} E_{\pi}[G_t \mid S_t{=}s, A_t{=}a]$$
$$= \max_{\pi} E_{\pi}[\textstyle\sum_{k=0..\infty} \gamma^k R_{t+k+1} \mid S_t{=}s, A_t{=}a]$$

# Optimal State- and Action-Value Function

**State-Value Function under optimal policy $\pi^*$:**

Intuition: Expected return starting in state $s$ and following the optimal policy $\pi^*$ thereafter.

$$v^*(s) = \max_{\pi} \mathrm{E}_{\pi}[G_t \mid S_t{=}s]$$

$$= \max_{\pi} \mathrm{E}_{\pi}[\textstyle\sum_{k=0..\infty} \gamma^k R_{t+k+1} \mid S_t{=}s]$$

**Action-Value Function under optimal policy $\pi^*$:**

Intuition: Expected return starting in state $s$, taking action $a$ at first step, then following the optimal policy $\pi^*$ thereafter. **If we know $q^*(s, a)$, we effectively have our policy**.

$$q^*(s, a) = \max_{\pi} \mathrm{E}_{\pi}[G_t \mid S_t{=}s, A_t{=}a]$$

$$= \max_{\pi} \mathrm{E}_{\pi}[\textstyle\sum_{k=0..\infty} \gamma^k R_{t+k+1} \mid S_t{=}s, A_t{=}a]$$

Key: If we can find v*(s), we can get q*(s) and problem solved!

Same As Prior Slide

Just Maximizing Over all Possible Policies

# Exercise 2.1: Find $v^*(s)$

Remember:

$$v^*(s) = \max_{\pi} \mathrm{E}_{\pi}[G_t \mid S_t{=}s]$$
$$= \max_{\pi} \mathrm{E}_{\pi}[\textstyle\sum_{k=0..\infty} \gamma^k R_{t+k+1} \mid S_t{=}s]$$
$$= \max_{\pi} \mathrm{E}_{\pi}[\textstyle\sum_{k=0..\infty} \gamma^k r(s_{t+k}, a_{t+k}, s_{t+k+1}) \mid S_t{=}s]$$

| γ = 1 (no discounting) Slip = 0 (deterministic dynamics) | |
| --- | --- |
| $v^*(0,3) =$ | 1 |
| $v^*(0,2) =$ | 1 |
| $v^*(0,1) =$ | 1 |
| $v^*(2,0) =$ | 1 |
| $v^*(1,3) =$ | -1 |

# Exercise 2.2: Find $v^*(s)$

$$v^*(s) = \max_\pi E_\pi[G_t \mid S_t = s]$$
$$= \max_\pi E_\pi[\textstyle\sum_{k=0..\infty} \gamma^k R_{t+k+1} \mid S_t = s]$$
$$= \max_\pi E_\pi[\textstyle\sum_{k=0..\infty} \gamma^k r(s_{t+k}, a_{t+k}, s_{t+k+1}) \mid S_t = s]$$

| $\gamma$ = **0.9** <br> **Slip = 0 (deterministic dynamics)** | |
|---|---|
| $v^*(0,3) =$ | 1 |
| $v^*(0,2) =$ | $0.9 \times 1 = 0.9$ |
| $v^*(0,1) =$ | $0.9 \times 0.9 \times 1 = 0.81$ ...or... <br> $0.9 \times v^*(0,2) = 0.81$ |
| $v^*(2,0) =$ | $0.9 \times 0.9 \times 0.9 \times 0.9 \times 0.9 \times 1 = 0.59$ ...or... <br> $0.9 \times \max(v^*(1,0), v^*(2,1)) = 0.59$ |
| $v^*(1,3) =$ | -1 |

**Image Credit:** Abbeel, 2017

# Exercise 2.3: Find $v^*(s)$

$$v^*(s) = \max_\pi E_\pi[G_t \mid S_t{=}s]$$
$$= \max_\pi E_\pi[\textstyle\sum_{k=0..\infty} \gamma^k R_{t+k+1} \mid S_t{=}s]$$
$$= \max_\pi E_\pi[\textstyle\sum_{k=0..\infty} \gamma^k r(s_{t+k}, a_{t+k}, s_{t+k+1}) \mid S_t{=}s]$$

| | |
|---|---|
| $\gamma = 0.9$ <br> Slip = **0.2** | |
| $v^*(0,3) =$ | 1 |
| $v^*(0,2) =$ | # 80% chance of success moving right... <br> $0.8 \times 0.9 \times 1 +$ <br> # 10% chance of slip left, bounce off top edge <br> $0.1 \times 0.9 \times v^*(0,2) +$ <br> # 10% chance of slip right.. <br> $0.1 \times 0.9 \times v^*(1,2) = ???$ |
| $v^*(0,1) =$ | I dunno, let Value Iteration figure it out |
| $v^*(2,0) =$ | I dunno, let Value Iteration figure it out |
| $v^*(1,3) =$ | -1 |

**Image Credit:** Abbeel, 2017

# Value Iteration Algorithm

$V_0^*[s] \leftarrow 0 \; \forall s, \; k \leftarrow 0$

Repeat until convergence:

    For each state $s$ in $S$:

        $V_k^*[s] \leftarrow \max_a \sum_{s'} p(s' \mid s, a) \, [r(s, a, s') + \gamma V_{k-1}^*[s']]$    Bellman Optimality Equation*

        $k \leftarrow k+1$

$\pi(s) \leftarrow \mathrm{argmax}_a \, Q_K(s, a) = \mathrm{argmax}_a \sum_{s'} p(s' \mid s, a) \, [r(s, a, s') + \gamma V_K^*[s']]$

**Intuition**:

$V_k^*[s]$ is best possible value of state $s$ if I can take at most $k$ steps. For for each state $s$ we ask:

- For each possible action $a$ I can take from here (state $s$)...

- ...in what states could I possibly end up ($s'$)? And with what probability?

- ...how much immediate reward would I get for landing in sate $s'$?

- ...using my current table of values for $V_{k-1}$, how much is that next state $s'$ worth (after discounting by $\gamma$)?

**The highest value across all the action possiblities is new estimate of $V^*[s]$.**

* Technically, this isn't *exactly* the BOE becuase we're relating iteration k to iteration k-1, but it's derived from BOE and is basically same thing.

# Exercises: Value Iteration Notebook

https://notebooks.azure.com/mhe500/projects/gridworld

# Exercise 3: Effect of Gamma and Stochasticity ("Slip Chance")

|   | 0 | 1 | 2 | 3 | 4 |
|---|---|---|---|---|---|
| 0 |   |   | ? |   |   |
| 1 |   | Rock |   |   |   |
| 2 |   | Rock | +1 | Rock | +10 |
| 3 | s0 |   | ? |   |   |
| 4 | -10 | -10 | -10 | -10 | -10 |

- When will the agent prefer the lower path (go near the cliff)?

- The upper path (stay away from the cliff)?

- When will the agent prefer the +1 reward? +10?

# Now We Can Solve Atari Games!

- Now we can calculate $q*$ for Atari 2600 games

$$Q\left[\ \text{[Pac-Man game image]}\ ,\ \text{Go Up}\ \right] = +50$$

- Just check $q*(s)$ to find the best action during game play

$$Q\left[\ \text{[Pac-Man game image]}\ ,\ \text{Go Right}\ \right] = \text{-}100$$

- **Will this work?**

… other state, action Q values …

# Limitations of Value Iteration

1. Not appropraite for large state spaces
   - Atari VCS has 160x192 pixels, 128 colors (NTSC version)*
   - $128^{(160*192)} = 128^{30,400}$ screen states
   - Too big to hold $V[s]$ or $Q[s, a]$
   - Too big to iterate over all states

2. Requires full knowledge of dynamics
   - No need for trial-and-error
   - Assumption not reasonable for complex environments



*Due to limitations and technical details of Atari VCS graphics implementation, this is an approximation.

# Other Approaches

Addresses unkown dynamics problem

| | Rely on <u>Known</u> Dynamics | Use <u>Samples</u> to Estimate Dynamics (Transitions & Rewards) |
|---|---|---|
| **Tabular Methods**<br><br>Store Full Table of $V^*[s]$ or $Q^*[s, a]$ | **Exact Methods**<br><br>• Value Iteration<br>• Policy Iteration | **Monte Carlo** (learn from sampled full episodes)<br><br>**Temporal Difference Methods** (learn from individually sampled transitions)<br><br>• TD(0) / TD(λ) (learn $V^\pi[s]$ on-policy from $(s, a, r, s')$ samples)<br>• Sarsa (learn $Q^*[s, a]$ on-policy from $(s, a, r, s', a)$ samples)<br>• Q-Learning (learn $Q^*[s, a]$ off-policy from $(s, a, r, s')$ samples) |
| **Non-Tabluar Function Approximators**<br><br>Approximate $q_*(s, a)$ with a non-tabluar function approximator (e.g., neural network) | | **Deep Q-Learning**<br><br>• Train neural network to approximate to $q_*(s, a)$ off-policy<br>• First method to solve Atari 2600 Games (Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D., & Riedmiller, M. (2013). Playing Atari with Deep Reinforcement Learning. Retrieved from http://arxiv.org/abs/1312.5602) (cited > 999 times) |

Addresses large state space problem

# Tabular Q-Learning

$Q[s, a] \leftarrow 0 \; \forall s$

Repeat:

      Choose action $a$ using some exploration policy (e.g., $\varepsilon$-greedy policy)

      Execute $a$ in the environment to get next state $s'$ and reward $R$

      $target \leftarrow R + \gamma \max_{a'} Q[s', a']$ # $Q[s', a']$ considered 0 if episode ended (why?)

      $Q[s, a] \leftarrow Q[s, a] + \alpha \, (target - Q[s, a])$    ← TD Error ($\delta_t$)

      $S \leftarrow S'$

**Intuition/Explanation:**

- $\varepsilon$-greedy means, with $\varepsilon$ probability take a random action, else act according to your policy.

- We collect many state transition samples by running $\varepsilon$-greedy policy and use *actual received reward* plus *next state's (discounted) best action value* (from current $Q[...]$ values) as estimate of $Q[s, a]$

- If $Q[s, a]$ differs from our just-sampled estimate, adjust $Q[s, a]$ by adding a fraction ($\alpha$) of the difference. Eventually, with enough samples, $Q[s, a]$ will converge to $q^*(s, a)$.

- **No use of dynamics (transition probabilities) & no "sweep" (foreach) over all states.**

# Deep Q-Learning Intuition

$Q[s, a] \leftarrow 0 \text{ for all } s \text{ in } S$ ~~(struck through)~~

> Replace $Q$ table with neural network $Q_\Theta$ estimates $q(s, a)$.

While not conveged:

　　Choose action $a$ using some exploration policy (e.g., $\varepsilon$-greedy)

　　Execute $a$ in the environment to get next state $s'$ and reward $R$

　　$target \leftarrow R + \gamma \max_{a'} Q_\Theta(s', a')$ # $Q_\Theta(s', a')$ considered 0 if episode ended (why?)

　　$Q[s, a] \leftarrow Q[s, a] + \alpha\,(target - Q[s, a])$ ~~(struck through)~~

> Train $Q_\Theta$ on batches of X=(s, a), Y=target using "supervised learning" with MSE loss.

　　$S \leftarrow S'$

**Intuition**:

- Transform pixels into lower-dimensional learned internal representation using convolutional network (CNN)
- Replace the $Q$ table with a neural network
- Gather and store experience, just like before
- Use the experience to generate target Q values, just like before
- Train the neural network as you would with supervised learning (what's odd here?)

# Pong $q*(s, a)$-Values



Action Values on Pong

# Atari Frames are <u>Not</u> Markovian States (why not?)



**Action Values on Pong**

Velocities are not represented in frame <u>observations</u>, but are key part of <u>state</u>. Hence, the frame doesn't have the Markov Property.

Mnih et al approximate Markovian State by feeding multiple successive frames into DQN.

**Image Credit:** Mnih, Volodymyr (2017). *Deep Q-Networks (2017 Deep RL Bootcamp)*.
https://sites.google.com/view/deep-rl-bootcamp/lectures

# More Notes on DQN

- Prior slides are *simplified gist &* omit many key contributions of Mnih et al to making DQN work on Atari games

- DQN not the only approach to playing games

- More recent gaming advances (e.g., Starcraft, DotA) using different algorithms (variations of Policy Gradient)



A non-exhaustive, but useful taxonomy of algorithms in modern RL.

# Summary

- Agent-Environment Loop defines how agent & environment interact
- Markov Decision Processes formally define an environment/problem
- MDP assumes state known (else, POMDP)
- MDP can be solved (find $v^*$, $q^*$) through Value Iteration if state space is small enough and dynamics are known
- Sample-based methods (Monte Carlo, TD, Q-Learning) can approximatly solve the MDP without knowing dynamics if state space is small enough
- Neural Networks can be used to replace $Q$ table becuase they generalize and propvide estimate $Q$ value for any state (including unseen states)

# Resources

- **Free Online RL Courses**
  - *OpenAI RL Bootcamp* (https://sites.google.com/view/deep-rl-bootcamp/lectures)
    - Lecture 1 – Pieter Abbeel on MDPs and Exact Solution Methods
    - Lecture 2 – Rocky (Yan) Duan on Temporal Difference (TD) Methods
    - Lecture 3 – Vlad Mnih on Deep Q-Learning / Atari (see also http://arxiv.org/abs/1312.5602)
  - *David Silver 2015 RL Course* (heavier foundations focus (e.g., MDPs, etc. Based heavily on earlier edition of Sutton & Barto) (http://www0.cs.ucl.ac.uk/staff/D.Silver/web/Teaching.html)
  - *Sergey Levine Berkeley Fall 2018 CS294-112* (heavier DNN, robotics, and current trends focus) (http://rail.eecs.berkeley.edu/deeprlcourse/)

- **Books**
  - Sutton, Richard S.., and Andrew G.. Barto. *Reinforcement Learning: an Introduction*. The MIT Press., 2018. Available as free PDF at http://incompleteideas.net/book/the-book-2nd.html
  - Géron, Aurélien. *Hands-on Machine Learning with Scikit-Learn and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems*. O'Reilly, 2018. (See RL chapter for DQN explanation; Note: don't purchase old edition; wait for 2019 second ed. covering Tensorflow 2.0)

# Thank you

(And thank you Microsoft for hosting)