# Lecture 6: Model-free RL with Value Function Approximation Continued [1]

Emma Brunskill

CS234 Reinforcement Learning.

Winter 2023

---

# Class Structure

- Last time: Model-free value function approximation control and Deep Q-learning
- This time: Model-free value function approximation and more DQN
- Next time: Policy search in large spaces / policy gradient methods

# Refresh Your Understanding: Modified AB Example: (Ex. 6.4, Sutton & Barto, 2018)

*[handwritten: TD → cert inty equivalent MDP soln for $V^\pi$]*

- Two states $A, B$ with $\gamma = 1$. Policy evaluation.
- Given 8 episodes of experience:
    - $A, 1, B, 0$ (observed 2 times) (state, reward, next state, next reward)
    - $B, 1$ (observed 4 times) (state, reward)
    - $B, 0$ (observed 2 times) (state, reward)
- Imagine run TD updates over data infinite number of times, and (separately) MC over data an infinite number of times?
- What is $V^{TD}(B)$ and $V^{TD}(A)$? What is $V^{MC}(B)$ and $V^{MC}(A)$?

*[handwritten below:]*

.5  1.5  .5  1

$r + \gamma V(B)$
$= 1 + \gamma \cdot .5 = 1.5$

$1 + 0 = 1$
$1 + 0 = 1$

# Refresh Your Understanding: Modified AB Example: (Ex. 6.4, Sutton & Barto, 2018). Solution

- Two states $A, B$ with $\gamma = 1$
- Given 8 episodes of experience:
  - $A, 1, B, 0$ (observed 2 times)
  - $B, 1$ (observed 4 times)
  - $B, 0$ (observed 2 times)
- Imagine run TD updates over data infinite number of times, and (separately) MC over data an infinite number of times?
- What is $V^{TD}(B)$ and $V^{TD}(A)$? What is $V^{MC}(B)$ and $V^{MC}(A)$?
  V(B) = 0.5 for TD and MC. V(A) = 1.5 for TD. V(A) = 1.0 for MC.

# Table of Contents

# Table of Contents

# Convergence Guarantees for Linear Value Function Approximation for Policy Evaluation

- Define the mean squared error of a linear value function approximation for a particular policy $\pi$ relative to the true value as

$$MSVE_\mu(\boldsymbol{w}) = \sum_{s \in S} \mu(s)(V^\pi(s) - \hat{V}^\pi(s; \boldsymbol{w}))^2$$

- where
  - $\mu(s)$: probability of visiting state $s$ under policy $\pi$. Note $\sum_s \mu(s) = 1$
  - $\hat{V}^\pi(s; \boldsymbol{w}) = \boldsymbol{x}(s)^T \boldsymbol{w}$, a linear value function approximation

# Convergence Guarantees for Linear Value Function Approximation for Policy Evaluation

- Define the mean squared error of a linear value function approximation for a particular policy $\pi$ relative to the true value as

$$MSVE_\mu(\boldsymbol{w}) = \sum_{s \in S} \mu(s)(V^\pi(s) - \hat{V}^\pi(s; \boldsymbol{w}))^2$$



- where
  - $\mu(s)$: probability of visiting state $s$ under policy $\pi$. Note $\sum_s \mu(s) = 1$
  - $\hat{V}^\pi(s; \boldsymbol{w}) = \boldsymbol{x}(s)^T \boldsymbol{w}$, a linear value function approximation
- Monte Carlo policy evaluation with VFA converges to the weights $\boldsymbol{w}_{MC}$ which has the <mark>minimum mean squared error</mark> possible with <mark>respect to the distribution $\mu$</mark>:

  *no Markov assumptions*

$$MSVE_\mu(\boldsymbol{w}_{MC}) = \min_{\boldsymbol{w}} \sum_{s \in S} \mu(s)(V^\pi(s) - \hat{V}^\pi(s; \boldsymbol{w}))^2$$

# Convergence Guarantees for TD Linear VFA for Policy Evaluation: Preliminaries

- For infinite horizon, the Markov Chain defined by a MDP with a particular policy will eventually converge to a probability distribution over states $d(s)$
- $d(s)$ is called the stationary distribution over states of $\pi$
- $\sum_s d(s) = 1$
- $d(s)$ satisfies the following balance equation:

$$d(s') = \sum_s \sum_a \underbrace{\pi(a|s)}\, \underbrace{p(s'|s,a)d(s)}_{\text{Markov}}$$

# Convergence Guarantees for Linear Value Function Approximation for Policy Evaluation

- Define the mean squared error of a linear value function approximation for a particular policy $\pi$ relative to the true value given the distribution $d$ as

$$MSVE_d(\boldsymbol{w}) = \sum_{s \in S} d(s)(V^\pi(s) - \hat{V}^\pi(s; \boldsymbol{w}))^2$$

- where
  - $d(s)$: stationary distribution of $\pi$ in the true decision process
  - $\hat{V}^\pi(s; \boldsymbol{w}) = \boldsymbol{x}(s)^T \boldsymbol{w}$, a linear value function approximation
- TD(0) policy evaluation with VFA converges to weights $\boldsymbol{w}_{TD}$ which is within a <u>constant</u> factor of the min mean squared error possible given distribution $d$:

$$MSVE_d(\boldsymbol{w}_{TD}) \leq \underline{\frac{1}{1-\gamma}} \min_{\boldsymbol{w}} \sum_{s \in S} d(s)(V^\pi(s) - \hat{V}^\pi(s; \boldsymbol{w}))^2$$

- TD(0) policy evaluation with VFA converges to weights $\boldsymbol{w}_{TD}$ which is within a constant factor of the min mean squared error possible for distribution $d$:

$$MSVE_d(\boldsymbol{w}_{TD}) \leq \frac{1}{1-\gamma} \min_{\boldsymbol{w}} \sum_{s \in S} d(s)(V^\pi(s) - \hat{V}^\pi(s; \boldsymbol{w}))^2$$

- If the VFA is a tabular representation (one feature for each state), what is the $MSVE_d$ for TD?

1. Depends on the problem
2. MSVE $= 0$ for TD
3. Not sure

- TD(0) policy evaluation with VFA converges to weights $\mathbf{w}_{TD}$ which is within a constant factor of the min mean squared error possible for distribution $d$:

$$MSVE_d(\mathbf{w}_{TD}) \leq \frac{1}{1-\gamma} \min_{\mathbf{w}} \sum_{s \in S} d(s)(V^\pi(s) - \hat{V}^\pi(s; \mathbf{w}))^2$$

- If the VFA is a tabular representation (one feature for each state), what is the $MSVE_d$ for TD?

MSVE = 0 for TD

$s. \ t. \quad s^f d \ cond$
$\quad \quad vn \ \propto$

# Table of Contents

# Recall Incremental Model-Free Control Approaches

- Similar to policy evaluation, true state-action value function for a state is unknown and so substitute a target value

- In Monte Carlo methods, use a return $G_t$ as a substitute target

$$\Delta \boldsymbol{w} = \alpha(G_t - \hat{Q}(s_t, a_t; \boldsymbol{w}))\nabla_{\boldsymbol{w}}\hat{Q}(s_t, a_t; \boldsymbol{w})$$

- For SARSA instead use a TD target $r + \gamma\hat{Q}(s', a'; \boldsymbol{w})$ which leverages the current function approximation value
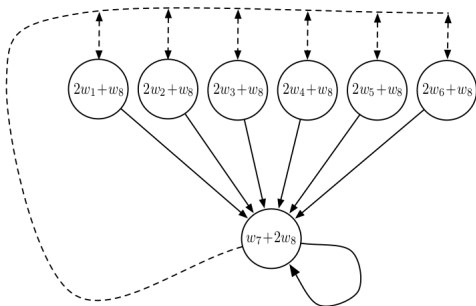
$$\Delta \boldsymbol{w} = \alpha(r + \gamma\hat{Q}(s', a'; \boldsymbol{w}) - \hat{Q}(s, a; \boldsymbol{w}))\nabla_{\boldsymbol{w}}\hat{Q}(s, a; \boldsymbol{w})$$

- For Q-learning instead use a TD target $r + \gamma\max_{a'}\hat{Q}(s', a'; \boldsymbol{w})$ which leverages the max of the current function approximation value

$$\Delta \boldsymbol{w} = \alpha(r + \gamma\max_{a'}\hat{Q}(s', a'; \boldsymbol{w}) - \hat{Q}(s, a; \boldsymbol{w}))\nabla_{\boldsymbol{w}}\hat{Q}(s, a; \boldsymbol{w})$$

# Challenges of Off Policy Control: Baird Example [1]



$\sim$1994/
1995

$\pi(\text{solid}|\cdot) = 1$
$\mu(\text{dashed}|\cdot) = 6/7$
$\mu(\text{solid}|\cdot) = 1/7$
$\gamma = 0.99$

- Behavior policy and target policy are not identical
- Value can diverge

bootstrapping
func approx
off policy

- Informally, updates involve doing an (approximate) Bellman backup followed by best trying to fit underlying value function to a particular feature representation
- Bellman operators are contractions, but value function approximation fitting can be an expansion
- Geoff Gordon 1995.

$$\| B V_1 - B V_2 \|_\infty \leq \| V_1 - V_2 \|_\infty$$

$$\| O B V_1 - O B V_2 \| \nleq \| V_1 - V_2 \|_\infty$$

$\uparrow$ func proj

# Convergence of Policy Evaluation and Control Methods with VFA

$\longrightarrow$ convergence

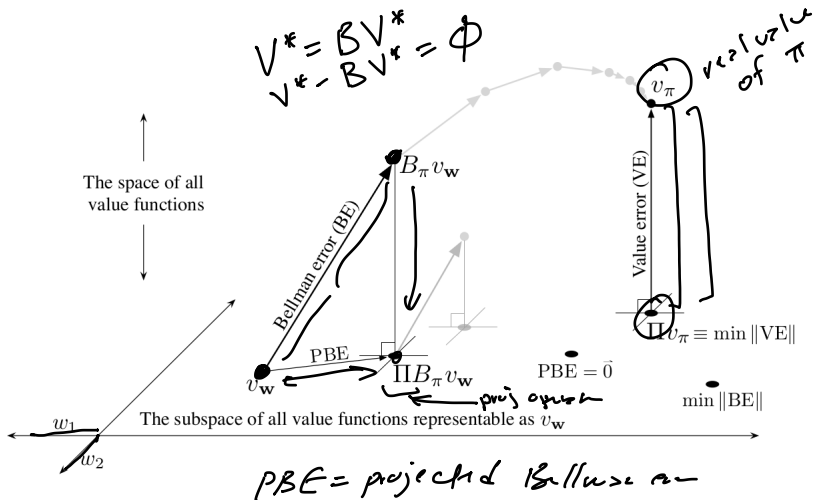| Algorithm | Tabular | Linear VFA | General VFA |
|:---:|:---:|:---:|:---:|
| Monte-Carlo Control | ✓ | MC explore starts 2022 | ✗ |
| Sarsa | ✓ | chatter | ✗ |
| Q-learning | ✓ | ✗ | ✗ |

$\epsilon$

$\alpha$

# Active Area: Off Policy Learning with Function Approximation

- Extensive work in better TD-style algorithms with value function approximation, some with convergence guarantees: see Chp 11 SB
- Will come up further later in this course

$$V^* = BV^* = \phi$$
$$V^* = BV^* = \phi$$

real value of $\pi$

The space of all value functions

$B_\pi v_{\mathbf{w}}$

Bellman error (BE)

Value error (VE)

$v_\pi$

PBE

$\Pi B_\pi v_{\mathbf{w}}$

proj option

$v_{\mathbf{w}}$

$\Pi v_\pi \equiv \min \|\text{VE}\|$

$\text{PBE} = \vec{0}$

$\min \|\text{BE}\|$

$w_1$

$w_2$

The subspace of all value functions representable as $v_{\mathbf{w}}$

PBE = projected Bellman er

$$PBE = \text{projected Bellman er}$$

[1] Figure from Sutton and Barto 2018

# Table of Contents

# Table of Contents

# Maximization Bias[2]

*episodes are 1 step*

- Consider single-state MDP ($|S| = 1$) with 2 actions, and both actions have 0-mean **random** rewards, ($\mathbb{E}(r|a = a_1) = \mathbb{E}(r|a = a_2) = 0$).
- Then $Q(s, a_1) = Q(s, a_2) = 0 = V(s)$
- Assume there are prior samples of taking action $a_1$ and $a_2$

  $a_1 \quad -.1 \quad .5 \quad .1 \cdots$
  $a_2 \quad -.5 \quad .3 \quad .2 \cdot .05 \cdots$

- Let $\hat{Q}(s, a_1), \hat{Q}(s, a_2)$ be the finite sample estimate of $Q$
- Use an unbiased estimator for $Q$: e.g. $\hat{Q}(s, a_1) = \frac{1}{n(s,a_1)} \sum_{i=1}^{n(s,a_1)} r_i(s, a_1)$
- Let $\hat{\pi} = \arg\max_a \hat{Q}(s, a)$ be the greedy policy w.r.t. the estimated $\hat{Q}$

$$\hat{V}^{\hat{\pi}} = E[\max \hat{Q}(s,a_1), \hat{Q}(s,a_2)]$$
$$\geq \max E[\hat{Q}(s,a_1) \hat{Q}(s,a_2)] \qquad \text{Jensen's inequality}$$
$$= \max(0,0)$$
$$= 0 = V^{\pi}$$

---

[2]Example from Mannor, Simester, Sun and Tsitsiklis. Bias and Variance Approximation in Value Function Estimates. Management Science 2007

# Maximization Bias[3] Proof

- Consider single-state MDP ($|S| = 1$) with 2 actions, and both actions have 0-mean random rewards, ($\mathbb{E}(r|a = a_1) = \mathbb{E}(r|a = a_2) = 0$).
- Then $Q(s, a_1) = Q(s, a_2) = 0 = V(s)$
- Assume there are prior samples of taking action $a_1$ and $a_2$
- Let $\hat{Q}(s, a_1), \hat{Q}(s, a_2)$ be the finite sample estimate of $Q$
- Use an unbiased estimator for $Q$: e.g. $\hat{Q}(s, a_1) = \frac{1}{n(s,a_1)} \sum_{i=1}^{n(s,a_1)} r_i(s, a_1)$
- Let $\hat{\pi} = \arg\max_a \hat{Q}(s, a)$ be the greedy policy w.r.t. the estimated $\hat{Q}$
- *Even though each estimate of the state-action values is unbiased*, the estimate of $\hat{\pi}$'s value $\hat{V}^{\hat{\pi}}$ can be biased:
  $\hat{V}^{\hat{\pi}}(s) = \mathbb{E}[\max \hat{Q}(s, a_1), \hat{Q}(s, a_2)]$
  $\geq \max[\mathbb{E}[\hat{Q}(s, a_1)], [\hat{Q}(s, a_2)]]$
  $= max[0, 0] = V^{\pi}$,
  where the inequality comes from Jensen's inequality.

---

[3]Example from Mannor, Simester, Sun and Tsitsiklis. Bias and Variance Approximation in Value Function Estimates. Management Science 2007

# Table of Contents

# Double Q-Learning

- The greedy policy w.r.t. estimated $Q$ values can yield a maximization bias during finite-sample learning
- Avoid using max of estimates as estimate of max of true values
- Instead split samples and use to create two independent unbiased estimates of $Q_1(s_1, a_i)$ and $Q_2(s_1, a_i)$ $\forall a$.
  - Use one estimate to select max action: $a^* = \arg\max_a Q_1(s_1, a)$
  - Use other estimate to estimate value of $a^*$: $Q_2(s, a^*)$
  - Yields unbiased estimate: $\mathbb{E}(Q_2(s, a^*)) = Q(s, a^*)$

# Double Q-Learning

- The greedy policy w.r.t. estimated $Q$ values can yield a maximization bias during finite-sample learning
- Avoid using max of estimates as estimate of max of true values
- Instead split samples and use to create two independent unbiased estimates of $Q_1(s_1, a_i)$ and $Q_2(s_1, a_i) \; \forall a$.
  - Use one estimate to select max action: $a^* = \arg\max_a Q_1(s_1, a)$
  - Use other estimate to estimate value of $a^*$: $Q_2(s, a^*)$
  - Yields unbiased estimate: $\mathbb{E}(Q_2(s, a^*)) = Q(s, a^*)$
- Why is this an unbiased estimate of the max state-action value?

# Double Q-Learning

- The greedy policy w.r.t. estimated $Q$ values can yield a maximization bias during finite-sample learning
- Avoid using max of estimates as estimate of max of true values
- Instead split samples and use to create two independent unbiased estimates of $Q_1(s_1, a_i)$ and $Q_2(s_1, a_i)$ $\forall a$.
    - Use one estimate to select max action: $a^* = \arg\max_a Q_1(s_1, a)$
    - Use other estimate to estimate value of $a^*$: $Q_2(s, a^*)$
    - Yields unbiased estimate: $\mathbb{E}(Q_2(s, a^*)) = Q(s, a^*)$
- Why is this an unbiased estimate of the max state-action value? Using independent samples to estimate the value

- If acting online, can alternate samples used to update $Q_1$ and $Q_2$, using the other to select the action chosen
- Next slides extend to full MDP case (with more than 1 state)

# Double Q-Learning

1: Initialize $Q_1(s, a)$ and $Q_2(s, a), \forall s \in S, a \in A$ $t = 0$, initial state $s_t = s_0$
2: **loop**
3:    Select $a_t$ using $\epsilon$-greedy $\pi(s) = \arg\max_a Q_1(s_t, a) + Q_2(s_t, a)$
4:    Observe $(r_t, s_{t+1})$
5:    **if** (with 0.5 probability) **then**
6:      $Q_1(s_t, a_t) \leftarrow Q_1(s_t, a_t) + \alpha(r_t + \gamma Q_2(s_{t+1}, \arg\max_a Q_1(s_{t+1}, a)) - Q_1(s_t, a_t))$
7:    **else**
8:      $Q_2(s_t, a_t) \leftarrow Q_2(s_t, a_t) + \alpha(r_t + \gamma Q_1(s_{t+1}, \arg\max_a Q_2(s_{t+1}, a)) - Q_2(s_t, a_t))$
9:    **end if**
10:    $t = t + 1$
11: **end loop**

Compared to Q-learning, how does this change the: memory requirements, $2x$ computation requirements per step, amount of data required? dosn't chge
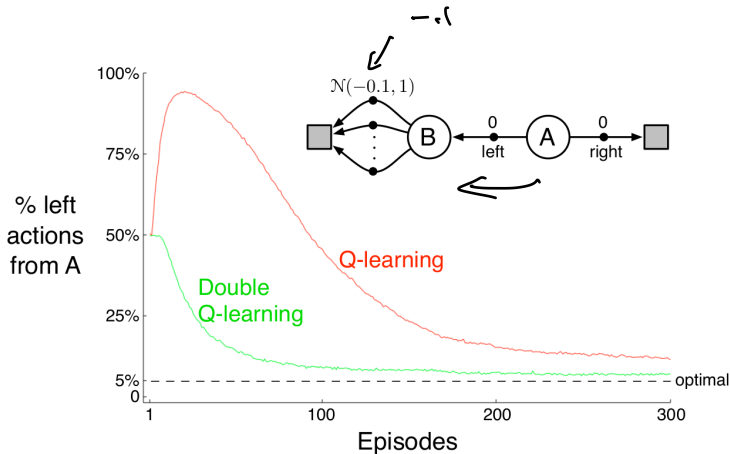
# Double Q-Learning

1: Initialize $Q_1(s,a)$ and $Q_2(s,a), \forall s \in S, a \in A$ $t = 0$, initial state $s_t = s_0$
2: **loop**
3:     Select $a_t$ using $\epsilon$-greedy $\pi(s) = \arg\max_a Q_1(s_t, a) + Q_2(s_t, a)$
4:     Observe $(r_t, s_{t+1})$
5:     **if** (with 0.5 probability) **then**
6:         $Q_1(s_t, a_t) \leftarrow Q_1(s_t, a_t) + \alpha(r_t + \gamma Q_2(s_{t+1}, \arg\max_a Q_1(s_{t+1}, a)) - Q_1(s_t, a_t))$
7:     **else**
8:         $Q_2(s_t, a_t) \leftarrow Q_2(s_t, a_t) + \alpha(r_t + \gamma Q_1(s_{t+1}, \arg\max_a Q_2(s_{t+1}, a)) - Q_2(s_t, a_t))$
9:     **end if**
10:    $t = t + 1$
11: **end loop**

Compared to Q-learning, how does this change the: memory requirements, computation requirements per step, amount of data required?

Doubles the memory, same computation requirements, data requirements are subtle– might reduce amount of exploration needed due to lower bias

Due to the maximization bias, Q-learning spends much more time selecting suboptimal actions than double Q-learning.

# Table of Contents

# Recall DQN

- Deep Q-learning (DQN): Q-learning with deep neural networks **and**
  - Experience replay
  - Fixed Q-targets

$$\Delta \boldsymbol{w} = \alpha(r + \gamma \max_{a'} \hat{Q}(s', a'; \boldsymbol{w}^-) - \hat{Q}(s, a; \boldsymbol{w}))\nabla_{\boldsymbol{w}}\hat{Q}(s, a; \boldsymbol{w})$$

# Recall DQN Pseudocode

1: Input $C$, $\alpha$, $D = \{\}$, Initialize $w$, $w^- = w$, $t = 0$
2: Get initial state $s_0$
3: **loop**
4:     Sample action $a_t$ given $\epsilon$-greedy policy for current $\hat{Q}(s_t, a; w)$
5:     Observe reward $r_t$ and next state $s_{t+1}$
6:     Store transition $(s_t, a_t, r_t, s_{t+1})$ in replay buffer $D$
7:     Sample random minibatch of tuples $(s_i, a_i, r_i, s_{i+1})$ from $D$
8:     **for** $j$ in minibatch **do**
9:         **if** episode terminated at step $i + 1$ **then**
10:             $y_i = r_i$
11:         **else**
12:             $y_i = r_i + \gamma \max_{a'} \hat{Q}(s_{i+1}, a'; w^-)$
13:         **end if**
14:         Do gradient descent step on $(y_i - \hat{Q}(s_i, a_i; w))^2$ for parameters $w$: $\Delta w = \alpha(y_i - \hat{Q}(s_i, a_i; w))\nabla_w \hat{Q}(s_i, a_i; w)$
15:     **end for**
16:     $t = t + 1$
17:     **if** mod(t,C) == 0 **then**
18:         $w^- \leftarrow w$
19:     **end if**
20: **end loop**

## Double DQN

- Double DQN (Deep Reinforcement Learning with Double Q-Learning, Van Hasselt et al, AAAI 2016)
- Extend double Q learning to DQN
- Current Q-network $\boldsymbol{w}$ is used to select actions
- Older Q-network $\boldsymbol{w}^-$ is used to evaluate actions

$$\Delta \boldsymbol{w} = \alpha(r + \gamma \overbrace{\hat{Q}(\underbrace{\arg\max_{a'} \hat{Q}(s', a'; \boldsymbol{w})}_{\text{Action selection: } \boldsymbol{w}}; \boldsymbol{w}^-)}^{\text{Action evaluation } \boldsymbol{w}^-} - \hat{Q}(s, a; \boldsymbol{w}))$$

# Double DQN

- Double DQN (Deep Reinforcement Learning with Double Q-Learning, Van Hasselt et al, AAAI 2016)
- Extend double Q learning to DQN
- Current Q-network $w$ is used to select actions
- Older Q-network $w^-$ is used to evaluate actions

$$\Delta w = \alpha(r + \gamma \overbrace{\hat{Q}(\underbrace{\arg\max_{a'} \hat{Q}(s', a'; w)}_{\text{Action selection: } w}; w^-)}^{\text{Action evaluation: } w^-} - \hat{Q}(s, a; w))$$

- How is this different from fixed target network update used in DQN?

# Double DQN

- Double DQN (Deep Reinforcement Learning with Double Q-Learning, Van Hasselt et al, AAAI 2016)
- Extend double Q learning to DQN
- Current Q-network $w$ is used to select actions
- Older Q-network $w^-$ is used to evaluate actions

$$\Delta w = \alpha(r + \gamma \underbrace{\hat{Q}(\overbrace{\arg\max_{a'} \hat{Q}(s', a'; w)}^{\text{Action evaluation: } w^-}; w^-)}_{\text{Action selection: } w} - \hat{Q}(s, a; w))$$

- How is this different from fixed target network update used in DQN? In DQN the same weights $w^-$ were used to choose the best action at $s'$ and evaluate its value $\hat{Q}(s', a'; w^-)$
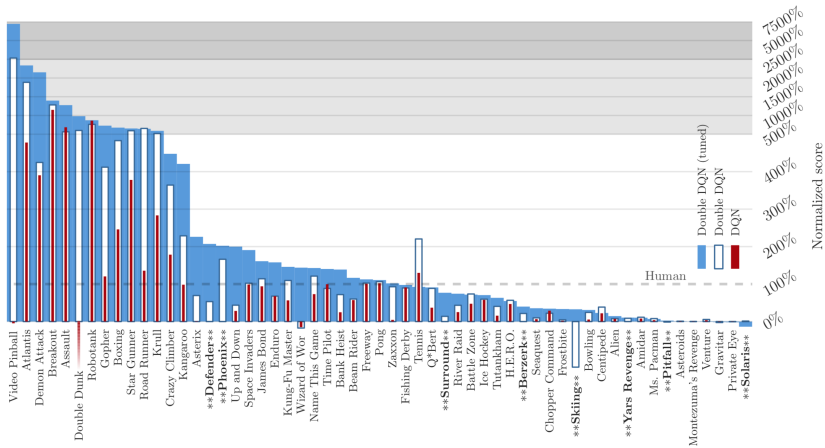
# Double DQN



Figure: van Hasselt, Guez, Silver, 2015

## Double DQN

- Double DQN (Deep Reinforcement Learning with Double Q-Learning, Van Hasselt et al, AAAI 2016)
- Extend double Q learning to DQN
- Current Q-network $w$ is used to select actions
- Older Q-network $w^-$ is used to evaluate actions

$$\Delta w = \alpha(r + \gamma \overbrace{\hat{Q}(\underbrace{\arg\max_{a'} \hat{Q}(s', a'; w)}_{\text{Action selection: } w}; w^-)}^{\text{Action evaluation: } w^-} - \hat{Q}(s, a; w))$$

- **Very small code change, often can lead to significantly improved results**

# Table of Contents
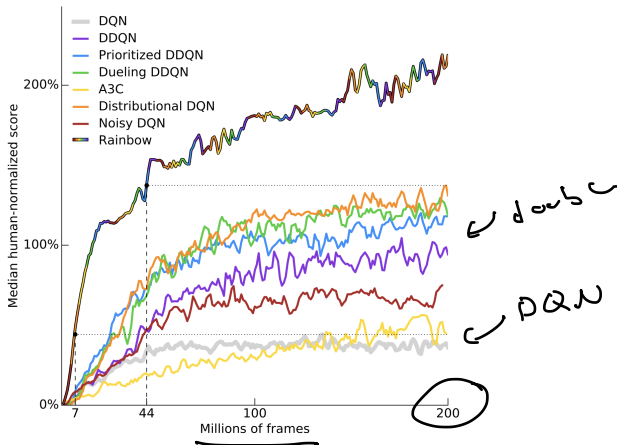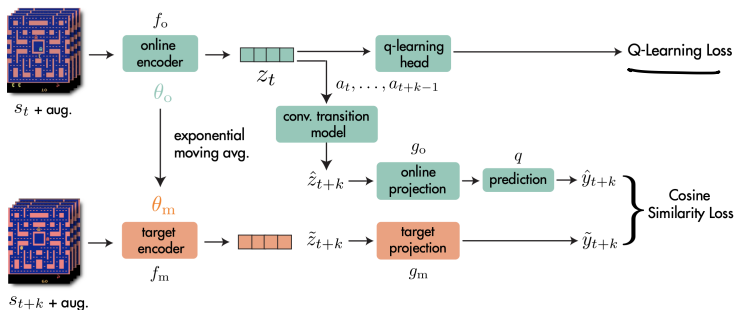
Figure: Median human-normalized performance across 57 Atari games. Curves smoothed with a moving avg over 5 points.

- One (of many) significant ideas: use additional objectives



Figure: Data-efficient reinforcement learning with self-predictive representations. Schwarzer et al. ICLR 2021.

# What is Enabling Progress?

- Benchmark tasks. Atari, Atari 100k, Mujoco, ...
- Standing on the shoulders of giants... : building on past algorithms
  - and code bases for said algorithms

# Model-free value function approximation RL: What You Should Know

- Be able to derive weight update for generic function approximation for $Q/V^\pi$
- Understand various (MC/SARSA/Q-learning) targets used when updating Q function
- Know what TD vs MC converge to for policy evaluation with a linear function approximator
- Be able to implement DQN
- Define the maximization bias and give one tool for alleviating it

# Class Structure

- Last time: Model-free value function approximation control and Deep Q-learning
- This time: Model-free value function approximation and more DQN
- Next time: Policy search in large spaces / policy gradient methods

# Lecture 6: Refresh Your Knowledge

- In TD learning with linear VFA (select all):
  1. $w = w + \alpha(r(s_t) + \gamma x(s_{t+1})^T w - x(s_t)^T w)x(s_t)$
  2. $V(s) = w(s)x(s)$
  3. Asymptotic convergence to the true best minimum MSE linear representable $V(s)$ is guaranteed for $\alpha \in (0, 1)$, $\gamma < 1$.
  4. Not sure

- In TD learning with linear VFA (select all):
  1. $\mathbf{w} = \mathbf{w} + \alpha(r(s_t) + \gamma \mathbf{x}(s_{t+1})^T \mathbf{w} - \mathbf{x}(s_t)^T \mathbf{w})\mathbf{x}(s_t)$
  2. $V(s) = \mathbf{w}(s)\mathbf{x}(s)$
  3. Asymptotic convergence to the true best minimum MSE linear representable $V(s)$ is guaranteed for $\alpha \in (0, 1)$, $\gamma < 1$.
  4. Not sure

Answer: 1 is true. Convergence is not guaranteed to the best, the resulting one may still be worse than the best MSE solution by a factor of $\frac{1}{1-\gamma}$. It is also important to know that this is with respect to the stationary distirbution $d(s)$. Also note the weights do not depend on the state.