

This Week's Plan

- ① Introduction on function approximation
- ② Value function approximation for prediction
- ③ Value function approximation for control
- ④ Batch RL and least square prediction and control
- ⑤ Deep Q-Learning: Entering the modern RL world

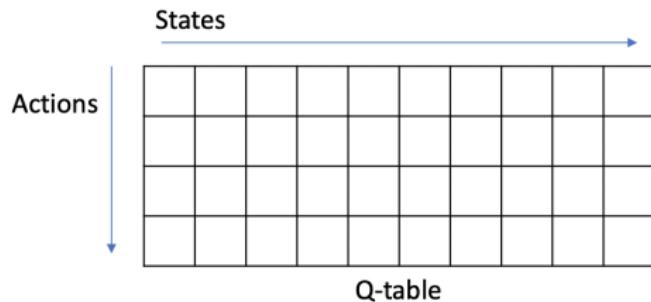
Introduction: Scaling up RL

- ① Previous lectures on small RL problems:
 - ① Cliff walk: 4×16 states
 - ② Mountain car: 1600 states
 - ③ Tic-Tac-Toe: 10^3 states
- ② Large-scale problems:
 - ① Backgammon: 10^{20} states
 - ② Chess: 10^{47} states
 - ③ Game of Go: 10^{170} states
 - ④ Robot Arm and Helicopter have continuous state space
 - ⑤ Number of atoms in universe: 10^{80}
- ③ Challenge: How can we scale up the model-free methods for prediction and control?

Introduction: Scaling up RL

- ① In tabular methods we represent value function by a lookup table:

- ① Every state s has an entry $V(s)$
- ② Every state-action pair s, a has an entry $Q(s, a)$



- ② Challenges with large MDPs:

- ① too many states or actions to store in memory
- ② too slow to learn the value of each state individually

Introduction: Scaling up RL

- ① In tabular methods we represent value function by a lookup table:
 - ① Every state s has an entry $V(s)$
 - ② Every state-action pair s, a has an entry $Q(s, a)$
- ② Challenges with large MDPs:
 - ① too many states or actions to store in memory
 - ② too slow to learn the value of each state individually

Scaling up RL with Function Approximation

- ① How to avoid explicitly learning or storing for every single state:
 - ① Dynamics or reward model
 - ② Value function, state-action function
 - ③ Policy
- ② Solution: Estimate with function approximation

$$\hat{v}(s, \mathbf{w}) \approx v^\pi(s)$$

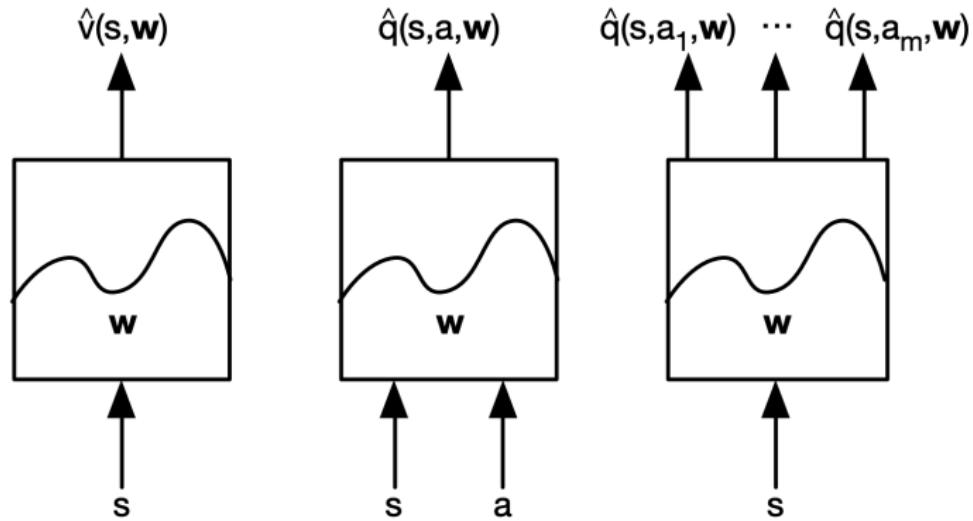
$$\hat{q}(s, a, \mathbf{w}) \approx q^\pi(s, a)$$

$$\hat{\pi}(a, s, \mathbf{w}) \approx \pi(a|s)$$

- ① Generalize from seen states to unseen states
- ② Update the parameter \mathbf{w} using MC or TD learning

Types of value function approximation

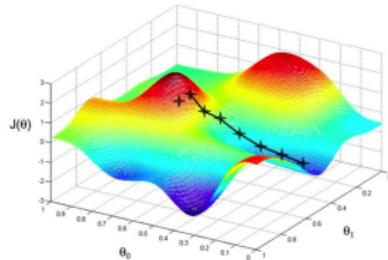
Several function designs:



Function Approximators

- ① Many possible function approximators:
 - ① Linear combinations of features
 - ② Neural networks
 - ③ Decision trees
 - ④ Nearest neighbors
- ② We will focus on function approximators that are differentiable
 - ① Linear feature representations
 - ② Neural networks

Review on Gradient Descend



- ① Consider a function $J(\mathbf{w})$ that is a differentiable function of a parameter vector \mathbf{w}
- ② Goal is to find parameter \mathbf{w}^* that minimizes J
- ③ Define the gradient of $J(\mathbf{w})$ to be
$$\nabla_{\mathbf{w}} J(\mathbf{w}) = \left(\frac{\partial J(\mathbf{w})}{\partial w_1}, \frac{\partial J(\mathbf{w})}{\partial w_2}, \dots, \frac{\partial J(\mathbf{w})}{\partial w_n} \right)^T$$
- ④ Adjust \mathbf{w} in the direction of the negative gradient, where α is step-size

$$\Delta \mathbf{w} = -\frac{1}{2}\alpha \nabla_{\mathbf{w}} J(\mathbf{w})$$

Value Function Approximation with an Oracle

- ① We assume that we have the oracle for knowing the true value for $v^\pi(s)$ for any given state s
- ② Then the objective is to find the best approximate representation of $v^\pi(s)$
- ③ Thus use the mean squared error and define the loss function as

$$J(\mathbf{w}) = \mathbb{E}_\pi \left[(v^\pi(S) - \hat{v}(s, \mathbf{w}))^2 \right]$$

- ④ Follow the gradient descend to find a local minimum

$$\Delta \mathbf{w} = -\frac{1}{2} \alpha \nabla_{\mathbf{w}} J(\mathbf{w})$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \Delta \mathbf{w}$$

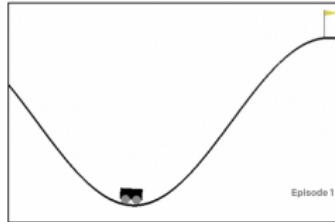
Representing State with Feature Vectors

- ① Represent state using a feature vector

$$\mathbf{x}(S) = (x_1(S), \dots, x_n(S))^T$$

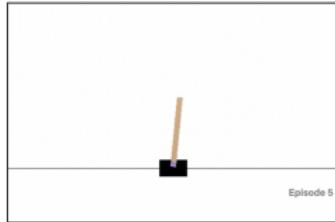
- ② For example:

Mountain Car



Position of car, velocity of car

Cart Pole



Position of cart, velocity of cart,
angle of pole, rotation rate of pole

Game of Go in AlphaGo

Extended Data Table 2 | Input features for neural networks

Feature	# of planes	Description
Stone colour	3	Player stone / opponent stone / empty
Onces	1	A constant plane filled with 1
Turns since	8	How many turns since a move was played
Liberties	8	Number of liberties (empty adjacent points)
Capture size	8	How many opponent stones would be captured
Self-atari size	8	How many of own stones would be captured
Liberties after move	8	Number of liberties after this move is played
Ladder capture	1	Whether a move at this point is a successful ladder capture
Ladder escape	1	Whether a move at this point is a successful ladder escape
Sensibleness	1	Whether a move is legal and does not fill its own eyes
Zeros	1	A constant plane filled with 0
Player color	1	Whether current player is black

48 places of 19x19 feature maps

Linear Value Function Approximation

- ① Represent value function by a linear combination of features

$$\hat{v}(S, \mathbf{w}) = \mathbf{x}(S)^T \mathbf{w} = \sum_{j=1}^n x_j(S) w_j$$

- ② The objective function is quadratic in parameter \mathbf{w}

$$J(\mathbf{w}) = \mathbb{E}_\pi \left[(v^\pi(S) - \mathbf{x}(S)^T \mathbf{w})^2 \right]$$

- ③ Thus the update rule is as simple as

$$\Delta \mathbf{w} = \alpha (v^\pi(S) - \hat{v}(S, \mathbf{w})) \mathbf{x}(S)$$

Update = StepSize × PredictionError × FeatureValue

- ④ Stochastic gradient descent converges to global optimum. Because in the linear case, there is only one optimum, thus local optimum is automatically converge to or near the global optimum.

Linear Value Function Approximation with Table Lookup Feature

- ① Table lookup is a special case of linear value function approximation
- ② Table lookup feature is one-hot vector as follows

$$\mathbf{x}^{table}(S) = (\mathbf{1}(S = s_1), \dots, \mathbf{1}(S = s_n))^T$$

- ③ Then we can see that each element on the parameter vector \mathbf{w} indicates the value of each individual state

$$\hat{v}(S, \mathbf{w}) = (\mathbf{1}(S = s_1), \dots, \mathbf{1}(S = s_n))(\mathbf{w}_1, \dots, \mathbf{w}_n)^T$$

- ④ Thus we have $\hat{v}(s_k, \mathbf{w}) = w_k$

Value Function Approximation for Model-free Prediction

- ① In practice, no access to oracle of the true value $v^\pi(s)$ for any state s
- ② Recall model-free prediction
 - ① Goal is to evaluate v^π following a fixed policy π
 - ② A lookup table is maintained to store estimates v^π or q^π
 - ③ Estimates are updated after each episode (MC method) or after each step (TD method)
- ③ Thus what we can do is to **include the function approximation step in the loop**

Incremental VFA Prediction Algorithms

- ① We assumed that true value function $v^\pi(s)$ given by supervisor/oracle

$$\Delta \mathbf{w} = \alpha \left(v^\pi(S) - \hat{v}(S, \mathbf{w}) \right) \nabla_{\mathbf{w}} \hat{v}(S_t, \mathbf{w})$$

- ② But in RL there is no supervisor, only rewards
- ③ In practice, we substitute the **target** for $v^\pi(s)$
 - ① For MC, the target is the actual return G_t

$$\Delta \mathbf{w} = \alpha \left(G_t - \hat{v}(S_t, \mathbf{w}) \right) \nabla_{\mathbf{w}} \hat{v}(S_t, \mathbf{w})$$

- ② For TD(0), the target is the TD target $R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w})$

$$\Delta \mathbf{w} = \alpha \left(R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w}) - \hat{v}(S_t, \mathbf{w}) \right) \nabla_{\mathbf{w}} \hat{v}(S_t, \mathbf{w})$$

Monte-Carlo Prediction with VFA

- ① Return G_t is an **unbiased**, but **noisy** sample of true value $v^\pi(S_t)$
- ② Why unbiased? $\mathbb{E}[G_t] = v^\pi(S_t)$
- ③ So we have the training data that can be used for supervised learning in VFA:

$$\langle S_1, G_1 \rangle, \langle S_2, G_2 \rangle, \dots, \langle S_T, G_T \rangle$$

- ④ Using linear Monte-Carlo policy evaluation

$$\begin{aligned}\Delta \mathbf{w} &= \alpha \left(G_t - \hat{v}(S_t, \mathbf{w}) \right) \nabla_{\mathbf{w}} \hat{v}(S_t, \mathbf{w}) \\ &= \alpha \left(G_t - \hat{v}(S_t, \mathbf{w}) \right) \mathbf{x}(S_t)\end{aligned}$$

- ⑤ Monte-Carlo prediction converges, in both linear and non-linear value function approximation.

TD Prediction with VFA

- ① TD target $R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w})$ is a **biased** sample of true value $v^\pi(S_t)$
- ② Why biased? It is drawn from our previous estimate, rather than the true value: $\mathbb{E}[R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w})] \neq v^\pi(S_t)$
- ③ We have the training data used for supervised learning in VFA:

$$< S_1, R_2 + \gamma \hat{v}(S_2, \mathbf{w}) >, < S_2, R_3 + \gamma \hat{v}(S_3, \mathbf{w}) >, \dots, < S_{T-1}, R_T >$$

- ④ Using linear TD(0), the stochastic gradient descend update is

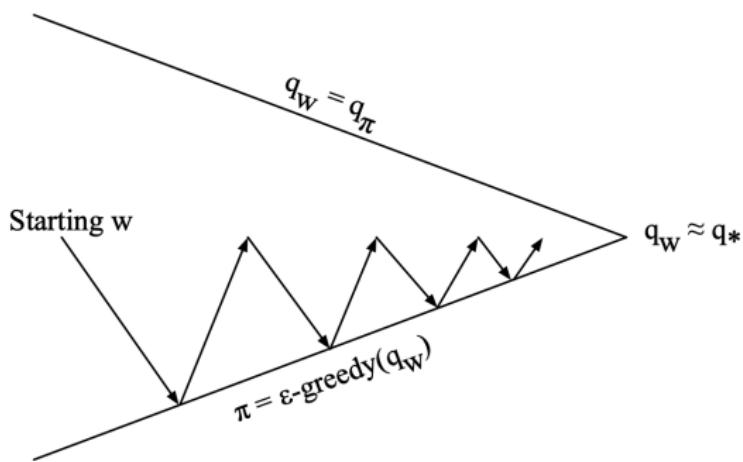
$$\begin{aligned}\Delta \mathbf{w} &= \alpha \left(R + \gamma \hat{v}(S', \mathbf{w}) - \hat{v}(S, \mathbf{w}) \right) \nabla_{\mathbf{w}} \hat{v}(S, \mathbf{w}) \\ &= \alpha \left(R + \gamma \hat{v}(S', \mathbf{w}) - \hat{v}(S, \mathbf{w}) \right) \mathbf{x}(S)\end{aligned}$$

This is also called as semi-gradient, as we ignore the effect of changing the weight vector \mathbf{w} on the target

- ⑤ Linear TD(0) converges(close) to global optimum

Control with Value Function Approximation

Generalized policy iteration



- ① Policy evaluation: **approximate** policy evaluation, $\hat{q}(\cdot, \cdot, \mathbf{w}) \approx q^\pi$
- ② Policy improvement: ϵ -greedy policy improvement

Action-Value Function Approximation

- ① Approximate the action-value function

$$\hat{q}(s, a, \mathbf{w}) \approx q^\pi(s, a)$$

- ② Minimize the MSE (mean-square error) between approximate action-value and true action-value (assume oracle)

$$J(\mathbf{w}) = \mathbb{E}_\pi[(q^\pi(s, a) - \hat{q}(s, a, \mathbf{w}))^2]$$

- ③ Stochastic gradient descend to find a local minimum

$$\Delta \mathbf{w} = \alpha(q^\pi(s, a) - \hat{q}(s, a, \mathbf{w})) \nabla_{\mathbf{w}} \hat{q}(s, a, \mathbf{w})$$

Linear Action-Value Function Approximation

- ① Represent state and action using a feature vector

$$\mathbf{x}(s, a) = (x_1(s, a), \dots, x_n(s, a))^T$$

- ② Represent action-value function by a linear combination of features

$$\hat{q}(s, a, \mathbf{w}) = \mathbf{x}(s, a)^T \mathbf{w} = \sum_{j=1}^n x_j(s, a) w_j$$

- ③ Thus the stochastic gradient descend update

$$\Delta \mathbf{w} = \alpha(q^\pi(s, a) - \hat{q}(s, a, \mathbf{w})) \mathbf{x}(s, a)$$

Incremental Control Algorithm

Same to the prediction, there is no oracle for the true value $q^\pi(s, a)$, so we substitute a target

- ① For MC, the target is the return G_t

$$\Delta \mathbf{w} = \alpha(G_t - \hat{q}(S_t, A_t, \mathbf{w})) \nabla_{\mathbf{w}} \hat{q}(S_t, A_t, \mathbf{w})$$

- ② For Sarsa, the target is the TD target $R_{t+1} + \gamma \hat{q}(S_{t+1}, A_{t+1}, \mathbf{w})$

$$\Delta \mathbf{w} = \alpha(R_{t+1} + \gamma \hat{q}(S_{t+1}, A_{t+1}, \mathbf{w}) - \hat{q}(S_t, A_t, \mathbf{w})) \nabla_{\mathbf{w}} \hat{q}(S_t, A_t, \mathbf{w})$$

- ③ For Q-learning, the target is the TD target

$$R_{t+1} + \gamma \max_a \hat{q}(S_{t+1}, a, \mathbf{w})$$

$$\Delta \mathbf{w} = \alpha(R_{t+1} + \gamma \max_a \hat{q}(S_{t+1}, a, \mathbf{w}) - \hat{q}(S_t, A_t, \mathbf{w})) \nabla_{\mathbf{w}} \hat{q}(S_t, A_t, \mathbf{w})$$

Semi-gradient Sarsa for VFA Control

Episodic Semi-gradient Sarsa for Estimating $\hat{q} \approx q_*$

Input: a differentiable function $\hat{q} : \mathcal{S} \times \mathcal{A} \times \mathbb{R}^d \rightarrow \mathbb{R}$

Initialize value-function weights $\mathbf{w} \in \mathbb{R}^d$ arbitrarily (e.g., $\mathbf{w} = \mathbf{0}$)

Repeat (for each episode):

$S, A \leftarrow$ initial state and action of episode (e.g., ε -greedy)

Repeat (for each step of episode):

Take action A , observe R, S'

If S' is terminal:

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha [R - \hat{q}(S, A, \mathbf{w})] \nabla \hat{q}(S, A, \mathbf{w})$$

Go to next episode

Choose A' as a function of $\hat{q}(S', \cdot, \mathbf{w})$ (e.g., ε -greedy)

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha [R + \gamma \hat{q}(S', A', \mathbf{w}) - \hat{q}(S, A, \mathbf{w})] \nabla \hat{q}(S, A, \mathbf{w})$$

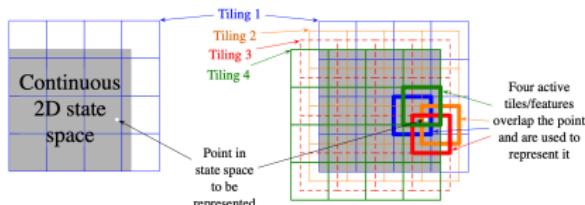
$$S \leftarrow S'$$

$$A \leftarrow A'$$

Mountain Car Example



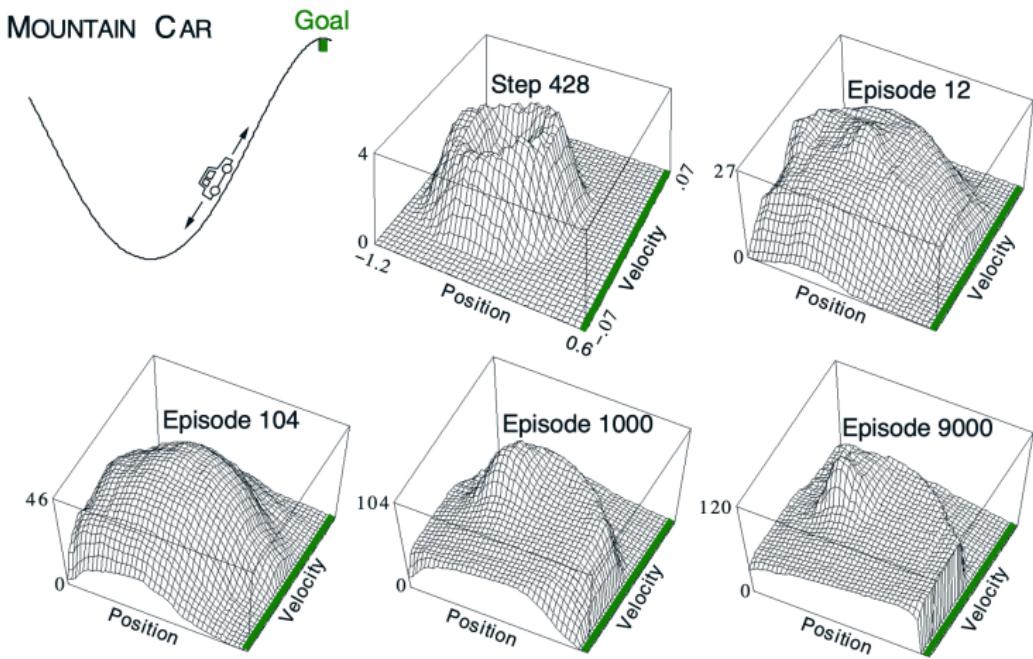
- ① Control task: move the car up the hill with forward and backward acceleration
 - ① Continuous state: [position of the car, velocity of the car]
 - ② Action: [full throttle forward, full throttle backward, zero throttle]
- ② Grid-tiling coding of the continuous 2D state space with 4 tilings



- ③ Q function approximator: $\hat{q}(s, a, \mathbf{w}) = \mathbf{w}^T \mathbf{x}(s, a) = \sum_{i=1}^4 w_i x_i(s, a)$

Mountain Car Example

- ① Visualization of cost-to-go function ($-\max_a \hat{q}(s, a, w)$) learned over episodes



Mountain Car Example

- ① Code example: [https://github.com/cuhkrlcourse/RLexample/
blob/master/modelfree/q_learning_mountaincar.py](https://github.com/cuhkrlcourse/RLexample/blob/master/modelfree/q_learning_mountaincar.py)

Convergence of Control Methods with VFA

- ① For Sarsa,

$$\Delta \mathbf{w} = \alpha (R_{t+1} + \gamma \hat{q}(s_{t+1}, a_{t+1}, \mathbf{w}) - \hat{q}(s_t, a_t, \mathbf{w})) \nabla_{\mathbf{w}} \hat{q}(s_t, a_t, \mathbf{w})$$

- ② For Q-learning,

$$\Delta \mathbf{w} = \alpha (R_{t+1} + \gamma \max_a \hat{q}(s_{t+1}, a, \mathbf{w}) - \hat{q}(s_t, a_t, \mathbf{w})) \nabla_{\mathbf{w}} \hat{q}(s_t, a_t, \mathbf{w})$$

- ③ TD with VFA doesn't follow the gradient of any objective function
- ④ The updates involve doing an approximate Bellman backup followed by fitting the underlying value function
- ⑤ That is why TD can diverge when off-policy or using non-linear function approximation
- ⑥ Challenge for off-policy control: behavior policy and target policy are not identical, thus value function approximation can diverge

The Deadly Triad for the Danger of Instability and Divergence

Potential problematic issues:

- ① **Function approximation**: A scalable way of generalizing from a state space much larger than the memory and computational resources
- ② **Bootstrapping**: Update targets that include existing estimates (as in dynamic programming or TD methods) rather than relying exclusively on actual rewards and complete returns (as in MC methods)
- ③ **Off-policy training**: training on a distribution of transitions other than that produced by the target policy
- ④ See details at Textbook Chapter 11.3

Convergence of Control Methods

Algorithm	Table Lookup	Linear	Non-Linear
Monte-Carlo Control	✓	(✓)	X
Sarsa	✓	(✓)	X
Q-Learning	✓	X	X

(✓) moves around the near-optimal value function

Batch Reinforcement Learning

- ① Incremental gradient descend update is simple
- ② But it is not sample efficient
- ③ Batch-based methods seek to find the best fitting value function for a batch of the agent's experience

Least Square Prediction

- ① Given the value function approximation $\hat{v}(s, \mathbf{w}) \approx v^\pi(s)$
- ② The experience \mathcal{D} consisting of $< \text{state}, \text{value} >$ pairs (may from one episode or many previous episodes)

$$\mathcal{D} = \{ < s_1, v_1^\pi >, \dots, < s_T, v_T^\pi > \}$$

- ③ Objective: To optimize the parameter \mathbf{w} that best fit all the experience \mathcal{D}
- ④ Least squares algorithms are used to minimize the sum-squared error between $\hat{v}(s_t, \mathbf{w})$ and the target values v_t^π

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \mathbb{E}_{\mathcal{D}}[(v^\pi - \hat{v}(s, \mathbf{w}))^2]$$

$$= \arg \min_{\mathbf{w}} \sum_{t=1}^T (v_t^\pi - \hat{v}(s_t, \mathbf{w}))^2$$

Stochastic Gradient Descent with Experience Replay

- ① Given the experience consisting of $\langle \text{state}, \text{value} \rangle$ pairs

$$\mathcal{D} = \{ \langle s_1, v_1^\pi \rangle, \dots, \langle s_T, v_T^\pi \rangle \}$$

- ② Iterative solution could be to repeat the following two steps
 - ① Randomly sample one pair $\langle \text{state}, \text{value} \rangle$ from the experience

$$\langle s, v^\pi \rangle \sim \mathcal{D}$$

- ② Apply the stochastic gradient descent update

$$\Delta \mathbf{w} = \alpha(v^\pi - \hat{v}(s, \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(s, \mathbf{w})$$

- ③ The solution from the gradient descend method converges to the least squares solution

$$\mathbf{w}^{LS} = \arg \min_{\mathbf{w}} \sum_{t=1}^T (v_t^\pi - \hat{v}(s_t, \mathbf{w}))^2$$

Solving Linear Least Squares Prediction

- ① Experience replay finds least squares solution, but it may take many iterations
- ② Using linear value function approximation $\hat{v}(s, \mathbf{w}) = \mathbf{x}(s)^T \mathbf{w}$, we can solve the least squares solution directly

Solving Linear Least Squares Prediction 2

- ① At minimum of $\text{LS}(\mathbf{w})$, the expected update must be zero

$$\mathbb{E}_{\mathcal{D}}[\Delta \mathbf{w}] = 0$$

- ② thus

$$\Delta \mathbf{w} = \alpha \sum_{t=1}^T \mathbf{x}(s_t)(v_t^\pi - \mathbf{x}(x_t)^T \mathbf{w}) = 0$$

$$\sum_{t=1}^T \mathbf{x}(s_t)v_t^\pi = \sum_{t=1}^T \mathbf{x}(s_t)\mathbf{x}(s_t)^T \mathbf{w}$$

$$\mathbf{w} = \left(\sum_{t=1}^T \mathbf{x}(s_t)\mathbf{x}(s_t)^T \right)^{-1} \sum_{t=1}^T \mathbf{x}(s_t)v_t^\pi$$

- ③ For N features, the matrix inversion complexity is $O(N^3)$

Linear Least Squares Prediction Algorithms

- ① Again, we do not know the true values v_t^π
- ② In practice, the training data must use noisy or biased samples of v_t^π
 - ① LSMC: Least squares Monte-Carlo uses return
 $v_t^\pi \approx G_t$
 - ② LSTD: Least squares TD uses TD target
 $v_t^\pi \approx R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w})$

Linear Least Squares Prediction Algorithms

LSMC

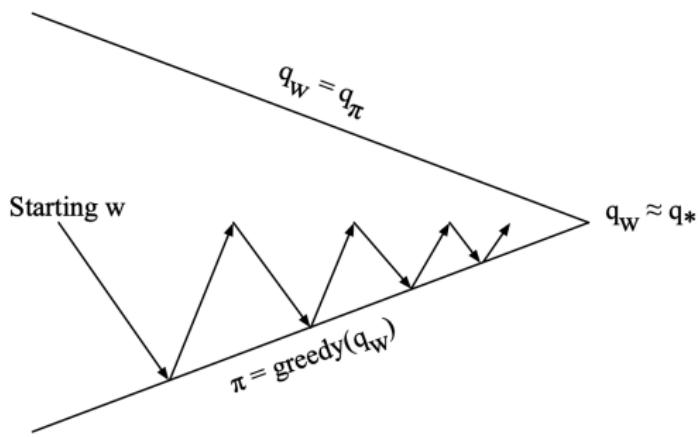
$$0 = \sum_{t=1}^T \alpha(G_t - \hat{v}(S_t, \mathbf{w})) \mathbf{x}(S_t)$$
$$\mathbf{w} = \left(\sum_{t=1}^T \mathbf{x}(S_t) \mathbf{x}(S_t)^T \right)^{-1} \sum_{t=1}^T \mathbf{x}(S_t) G_t$$

LSTD

$$0 = \sum_{t=1}^T \alpha(R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w}) - \hat{v}(S_t, \mathbf{w})) \mathbf{x}(S_t)$$
$$\mathbf{w} = \left(\sum_{t=1}^T \mathbf{x}(S_t) (\mathbf{x}(S_t) - \gamma \mathbf{x}(S_{t+1}))^T \right)^{-1} \sum_{t=1}^T \mathbf{x}(S_t) R_{t+1}$$

Least Square Control with Value Function Approximation

Generalized policy iteration



- ① Policy evaluation: Policy evaluation by least squares Q-learning
- ② Policy improvement: greedy policy improvement

Least Squares Action-Value Function Approximation

- ① Approximate action-value function $q^\pi(s, a)$
- ② using linear combination of features $\mathbf{x}(s, a)$

$$\hat{q}(s, a, \mathbf{w}) = \mathbf{x}(s, a)^T \mathbf{w} \approx q^\pi(s, a)$$

- ③ minimize least squares error between $\hat{q}(s, a, \mathbf{w})$ and $q^\pi(s, a)$
- ④ The experience is generate from using policy π , consisting of
 $< state, action, action - value >$

Least Squares Control

- ① For policy evaluation we want to efficiently use all the experience
- ② For control we want to improve the policy
- ③ The experience is generated from many policies (previous old policies)
- ④ So to evaluate $q^\pi(s, a)$ we must learn **off-policy**
- ⑤ We follow the same idea of Q-Learning:
 - ① Use the experience generated by old policy:

$$S_t, A_t, R_{t+1}, S_{t+1} \sim \pi_{old}$$

- ② Consider the alternative successor action $A' = \pi_{new}(S_{t+1})$ using greedy
- ③ Update $\hat{q}(S_t, A_t, \mathbf{w})$ towards value of alternative action
 $R_{t+1} + \gamma \hat{q}(S_{t+1}, A', \mathbf{w})$

Least Squares Q-Learning

- ① Consider the following linear Q-learning update

$$\delta = R_{t+1} + \gamma \hat{q}(S_{t+1}, \pi(S_{t+1}), \mathbf{w}) - \hat{q}(S_t, A_t, \mathbf{w})$$

$$\Delta \mathbf{w} = \alpha \delta \mathbf{x}(S_t, A_t)$$

- ② LSTDQ algorithm: solve the total sum of the gradient = zero:

$$0 = \sum_{t=1}^T \alpha(R_{t+1} + \gamma \hat{q}(S_{t+1}, \pi(S_{t+1}), \mathbf{w}) - \hat{q}(S_t, A_t, \mathbf{w})) \mathbf{x}(S_t, A_t)$$

$$\mathbf{w} = \left(\sum_{t=1}^T \mathbf{x}(S_t, A_t) (\mathbf{x}(S_t, A_t) - \gamma \mathbf{x}(S_{t+1}, \pi(S_{t+1})))^T \right)^{-1} \sum_{t=1}^T \mathbf{x}(S_t, A_t) R_{t+1}$$

Least Squares Policy Iteration Algorithm

- ① Use LSTDQ for policy evaluation
- ② Repeatedly re-evaluate experience \mathcal{D} with different policies

```
function LSPI-TD( $\mathcal{D}, \pi_0$ )
     $\pi' \leftarrow \pi_0$ 
    repeat
         $\pi \leftarrow \pi'$ 
         $Q \leftarrow \text{LSTDQ}(\pi, \mathcal{D})$ 
        for all  $s \in \mathcal{S}$  do
             $\pi'(s) \leftarrow \underset{a \in \mathcal{A}}{\text{argmax}} Q(s, a)$ 
        end for
    until ( $\pi \approx \pi'$ )
    return  $\pi$ 
end function
```

Convergence of Control Methods

Algorithm	Table Lookup	Linear	Non-Linear
Monte-Carlo Control	✓	(✓)	✗
Sarsa	✓	(✓)	✗
Q-Learning	✓	✗	✗
LSPI	✓	(✓)	-

(✓) moves around the near-optimal value function

Deep Q-Learning

- ① DeepMind's **Nature** paper: Mnih, Volodymyr; et al. (2015).
Human-level control through deep reinforcement learning
 - ① Entering the period of modern RL: deep learning + reinforcement learning

Review on Stochastic Gradient Descend for Function Approximation

- ① Goal: Find the parameter vector \mathbf{w} that minimizes the loss between a true value function $v_\pi(s)$ and its approximation $\hat{v}_\pi(s, \mathbf{w})$ as represented with a particular function approximator parameterized by \mathbf{w}
- ② The mean square error loss function is as

$$J(\mathbf{w}) = \mathbb{E}_\pi \left[(v_\pi(S) - \hat{v}(s, \mathbf{w}))^2 \right]$$

- ③ Follow the gradient descend to find a local minimum

$$\begin{aligned}\Delta \mathbf{w} &= -\frac{1}{2} \alpha \nabla_{\mathbf{w}} J(\mathbf{w}) \\ \mathbf{w}_{t+1} &= \mathbf{w}_t + \Delta \mathbf{w}\end{aligned}$$

Linear Function Approximation

- ① Represent value function by a linear combination of features

$$\hat{v}(S, \mathbf{w}) = \mathbf{x}(S)^T \mathbf{w} = \sum_{j=1}^n x_j(S) w_j$$

- ② Objective function is $J(\mathbf{w}) = \mathbb{E}_\pi \left[(v_\pi(S) - \mathbf{x}(S)^T \mathbf{w})^2 \right]$
- ③ Update is as simple as $\Delta \mathbf{w} = \alpha(v_\pi(S) - \hat{v}(S, \mathbf{w})) \mathbf{x}(S)$
- ④ But there is no oracle for the true value $v_\pi(S)$, we substitute with the target from MC or TD
 - ① for MC policy evaluation,

$$\Delta \mathbf{w} = \alpha \left(G_t - \hat{v}(S_t, \mathbf{w}) \right) \mathbf{x}(S_t)$$

- ② for TD policy evaluation,

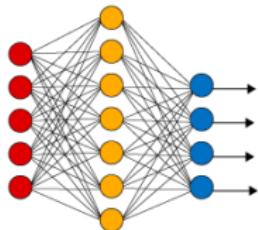
$$\Delta \mathbf{w} = \alpha \left(R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w}) - \hat{v}(S_t, \mathbf{w}) \right) \mathbf{x}(S_t)$$

Linear vs Nonlinear Value Function Approximation

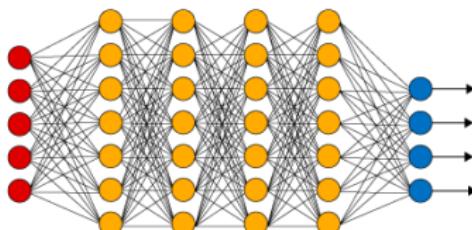
- ① Linear VFA often works well given the right set of features
- ② But it requires manual designing of the feature set
- ③ Alternative is to use a much richer function approximator that is able to directly learn from states without requiring the feature design
- ④ Nonlinear function approximator: Deep neural networks

Deep Neural Networks

Simple Neural Network



Deep Learning Neural Network



● Input Layer

● Hidden Layer

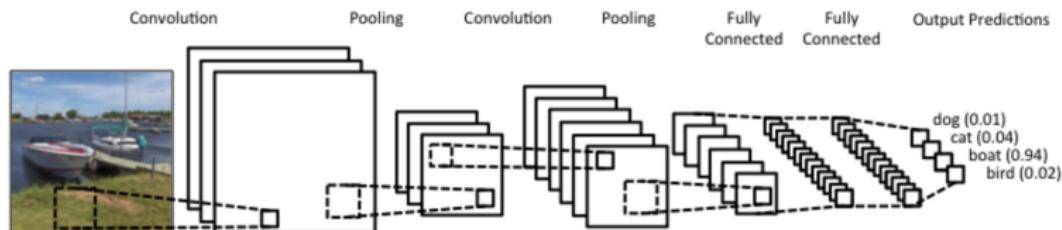
● Output Layer

- ① Multiple layers of linear functions, with non-linear operators between layers

$$f(\mathbf{x}; \theta) = \mathbf{W}_{L+1}^T \sigma(\mathbf{W}_L^T \sigma(\dots \sigma(\mathbf{W}_1^T \mathbf{x} + \mathbf{b}_1) + \dots + \mathbf{b}_{L-1}) + \mathbf{b}_L) + \mathbf{b}_{L+1}$$

- ② The chain rule to backpropagate the gradient to update the weights using the loss function $L(\theta) = \frac{1}{n} \sum_{i=1}^n \left(y^{(i)} - f(\mathbf{x}; \theta) \right)^2$

Convolutional Neural Networks



- ① Convolution encodes the local information in 2D feature map
- ② Layers of convolution, reLU, batch normalization, etc.
- ③ CNNs are widely used in computer vision (more than 70% top conference papers using CNNs)
- ④ A detailed introduction on CNNs:
<http://cs231n.github.io/convolutional-networks/>

Deep Reinforcement Learning

- ① Frontier in machine learning and artificial intelligence
- ② Deep neural networks are used to represent
 - ① Value function
 - ② Policy function (policy gradient methods to be introduced)
 - ③ World model
- ③ Loss function is optimized by stochastic gradient descent (SGD)
- ④ Challenges
 - ① Efficiency: too many model parameters to optimize
 - ② The Deadly Triad for the danger of instability and divergence in training
 - ① Nonlinear function approximation
 - ② Bootstrapping
 - ③ Off-policy training

Deep Q-Networks (DQN)

- ① DeepMind's **Nature** paper: Mnih, Volodymyr; et al. (2015).
Human-level control through deep reinforcement learning
- ② DQN represents the action value function with neural network approximator
- ③ DQN reaches a professional human gaming level across many Atari games using the same network and hyperparameters



4 Atari Games: Breakout, Pong, Montezuma's Revenge, Private Eye

Recall: Action-Value Function Approximation

- ① Approximate the action-value function

$$\hat{q}(S, A, \mathbf{w}) \approx q_{\pi}(S, A)$$

- ② Minimize the MSE (mean-square error) between approximate action-value and true action-value (assume oracle)

$$J(\mathbf{w}) = \mathbb{E}_{\pi}[(q_{\pi}(S, A) - \hat{q}(S, A, \mathbf{w}))^2]$$

- ③ Stochastic gradient descend to find a local minimum

$$\Delta \mathbf{w} = \alpha(q_{\pi}(S, A) - \hat{q}(S, A, \mathbf{w})) \nabla_{\mathbf{w}} \hat{q}(S, A, \mathbf{w})$$

Recall: Incremental Control Algorithm

Same to the prediction, there is no oracle for the true value $q_\pi(S, A)$, so we substitute a target

- ① For MC, the target is the return G_t

$$\Delta \mathbf{w} = \alpha(G_t - \hat{q}(S_t, A_t, \mathbf{w})) \nabla_{\mathbf{w}} \hat{q}(S_t, A_t, \mathbf{w})$$

- ② For Sarsa, the target is the TD target $R_{t+1} + \gamma \hat{q}(S_{t+1}, A_{t+1}, \mathbf{w})$

$$\Delta \mathbf{w} = \alpha(R_{t+1} + \gamma \hat{q}(S_{t+1}, A_{t+1}, \mathbf{w}) - \hat{q}(S_t, A_t, \mathbf{w})) \nabla_{\mathbf{w}} \hat{q}(S_t, A_t, \mathbf{w})$$

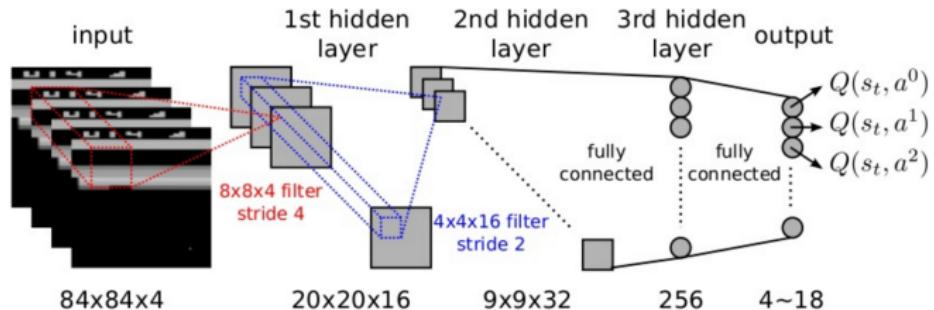
- ③ For Q-learning, the target is the TD target

$$R_{t+1} + \gamma \max_a \hat{q}(S_{t+1}, a, \mathbf{w})$$

$$\Delta \mathbf{w} = \alpha(R_{t+1} + \gamma \max_a \hat{q}(S_{t+1}, a, \mathbf{w}) - \hat{q}(S_t, A_t, \mathbf{w})) \nabla_{\mathbf{w}} \hat{q}(S_t, A_t, \mathbf{w})$$

DQN for Playing Atari Games

- ① End-to-end learning of values $Q(s, a)$ from input pixel frame
- ② Input state s is a stack of raw pixels from latest 4 frames
- ③ Output of $Q(s, a)$ is 18 joystick/button positions
- ④ Reward is the change in score for that step
- ⑤ Network architecture and hyperparameters fixed across all games

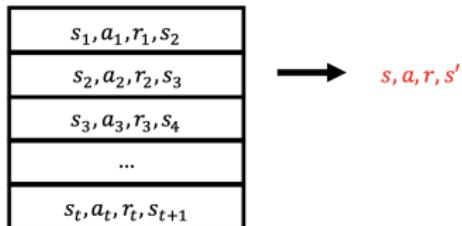


Q-Learning with Value Function Approximation

- ① Two of the issues causing problems:
 - ① Correlations between samples
 - ② Non-stationary targets
- ② Deep Q-learning (DQN) addresses both of these challenges by
 - ① Experience replay
 - ② Fixed Q targets

DQNs: Experience Replay

- ① To reduce the correlations among samples, store transition (s_t, a_t, r_t, s_{t+1}) in replay memory \mathcal{D}



- ② To perform experience replay, repeat the following
 - ① sample an experience tuple from the dataset: $(s, a, r, s') \sim \mathcal{D}$
 - ② compute the target value for the sampled tuple: $r + \gamma \max_{a'} \hat{Q}(s', a', \mathbf{w})$
 - ③ use stochastic gradient descent to update the network weights

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

DQNs: Fixed Targets

- ① To help improve stability, fix the target weights used in the target calculation for multiple updates
- ② Let a different set of parameter \mathbf{w}^- be the set of weights used in the target, and \mathbf{w} be the weights that are being updated
- ③ To perform experience replay with fixed target, repeat the following
 - ① sample an experience tuple from the dataset: $(s, a, r, s') \sim \mathcal{D}$
 - ② compute the target value for the sampled tuple:
 $r + \gamma \max_{a'} \hat{Q}(s', a', \mathbf{w}^-)$
 - ③ use stochastic gradient descent to update the network weights

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

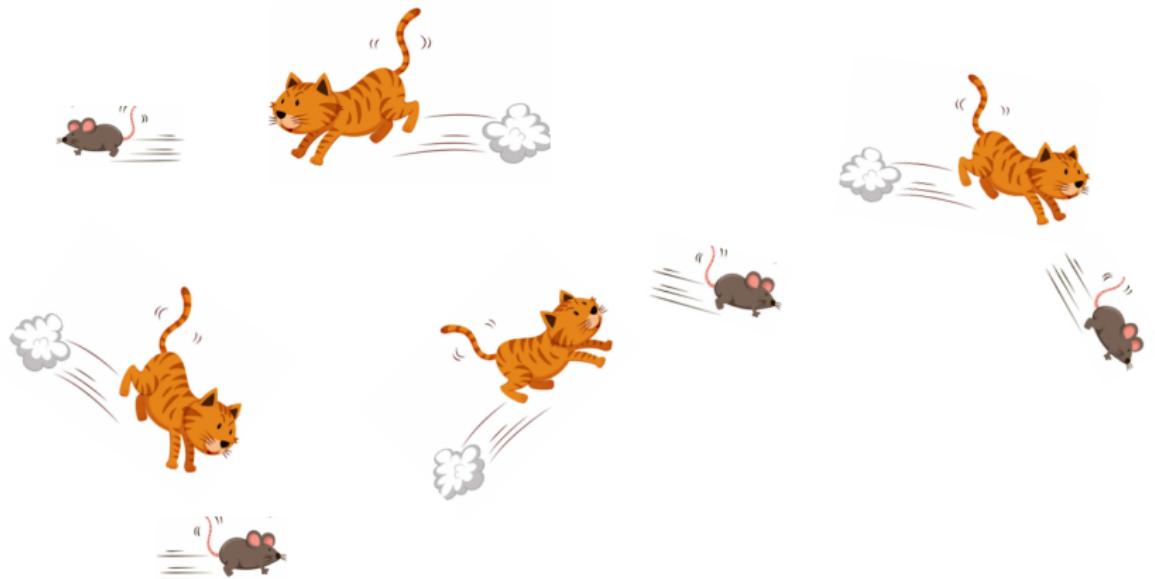
Why fixed target

- ① In the original update, both Q estimation and Q target shifts at each time step
- ② Imagine a cat (Q estimation) is chasing after a mouse (Q target)
- ③ The cat must reduce the distance to the mouse



Why fixed target

- ① Both the cat and mouse are moving,



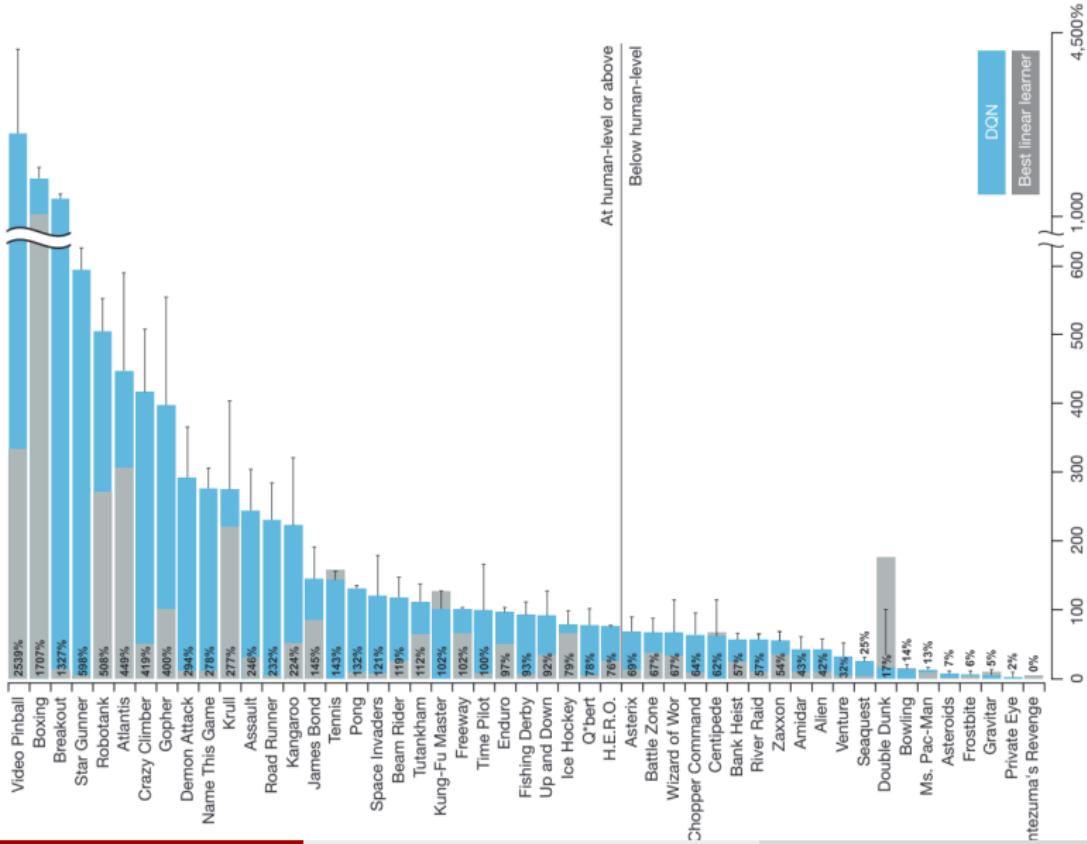
Why fixed target

- ① This could lead to a strange path of chasing (an oscillated training history)



- ② Solution: fix the target for a period of time during the training

Performance of DQNs on Atari



Ablation Study on DQNs

① Game score under difference conditions

	Replay Fixed-Q	Replay Q-learning	No replay Fixed-Q	No replay Q-learning
Breakout	316.81	240.73	10.16	3.17
Enduro	1006.3	831.25	141.89	29.1
River Raid	7446.62	4102.81	2867.66	1453.02
Seaquest	2894.4	822.55	1003	275.81
Space Invaders	1088.94	826.33	373.22	301.99

Demo of DQNs

- ① Demo of deep q-learning for Breakout:

<https://www.youtube.com/watch?v=V1eYniJ0Rnk>

- ② Demo of Flappy Bird by DQN:

<https://www.youtube.com/watch?v=xM62SpKAZHU>

- ③ Code of DQN in PyTorch: <https://github.com/cuhkrlcourse/DeepRL-Tutorials/blob/master/01.DQN.ipynb>

- ④ Code of Flappy Bird:

<https://github.com/xmfbit/DQN-FlappyBird>

Summary of DQNs

- ① DQN uses experience replay and fixed Q-targets
- ② Store transition $(s_t, a_t, r_{t+1}, s_{t+1})$ in replay memory \mathcal{D}
- ③ Sample random mini-batch of transitions (s, a, r, s') from \mathcal{D}
- ④ Compute Q-learning targets w.r.t. old, fixed parameters \mathbf{w}^-
- ⑤ Optimizes MSE between Q-network and Q-learning targets using stochastic gradient descent

Improving DQN

- ① Success in Atari has led to a huge excitement in using deep neural networks to do value function approximation in RL
- ② Many follow-up works on improving DQNs
 - ① **Double DQN**: Deep Reinforcement Learning with Double Q-Learning. Van Hasselt et al, AAAI 2016
 - ② **Dueling DQN**: Dueling Network Architectures for Deep Reinforcement Learning. Wang et al, best paper ICML 2016
 - ③ **Prioritized Replay**: Prioritized Experience Replay. Schaul et al, ICLR 2016
- ③ A nice tutorial on the relevant algorithms:
<https://github.com/cuhkrlcourse/DeepRL-Tutorials>

Improving DQN: Double DQN

- ① Handles the problem of the overestimation of Q-values
- ② Idea: use the two networks to decouple the action selection from the target Q value generation
- ③ Vanilla DQN:

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

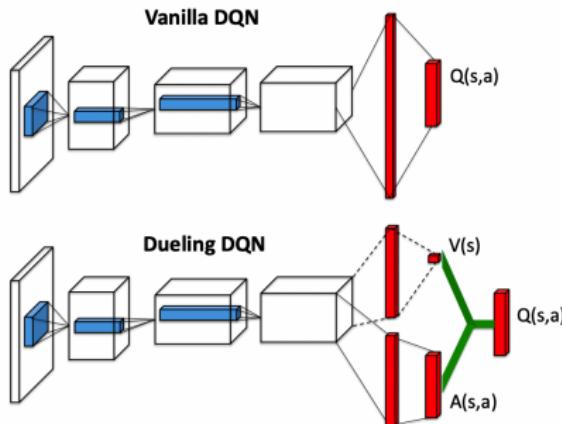
- ④ Double DQN:

$$\Delta \mathbf{w} = \alpha \left(r + \gamma \hat{Q}(s', \arg \max_{a'} Q(s', a', \mathbf{w}), \mathbf{w}^-) - Q(s, a, \mathbf{w}) \right) \nabla_{\mathbf{w}} \hat{Q}(s, a, \mathbf{w})$$

- ⑤ Code: https://github.com/cuhkrlcourse/DeepRL-Tutorials/blob/master/03.Double_DQN.ipynb

Improving DQN: Dueling DQN

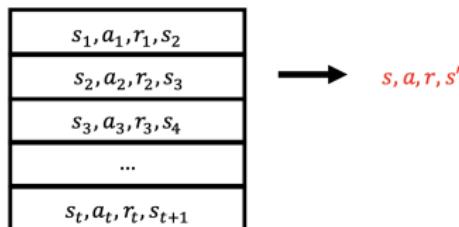
- ① One branch estimates $V(s)$, other branch estimates the advantage for each action $A(s, a)$. Then $Q(s, a) = A(s, a) + V(s)$



- ② By decoupling the estimation, intuitively the DuelingDQN can learn which states are (or are not) valuable without having to learn the effect of each action at each state
- ③ Code: [https://github.com/cuhkrlcourse/DeepRL-Tutorials/
blob/master/04.Dueling_DQN.ipynb](https://github.com/cuhkrlcourse/DeepRL-Tutorials/blob/master/04.Dueling_DQN.ipynb)

Improving DQN: Prioritized Experience Replay

- ① Transition $(s_t, a_t, r_{t+1}, s_{t+1})$ is stored in and sampled from the replay memory \mathcal{D}



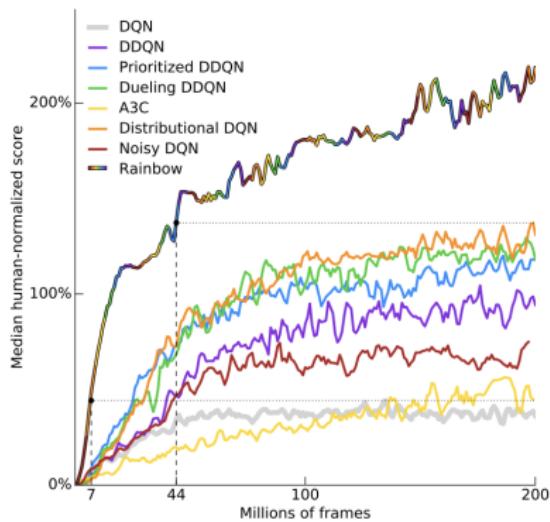
- ② Priority is on the experience where there is a big difference between our prediction and the TD target, since it means that we have a lot to learn about it.
- ③ Define a priority score for each tuple i

$$p_i = |r + \gamma \max_{a'} Q(s_{i+1}, a', \mathbf{w}^-) - Q(s_i, a_i, \mathbf{w})|$$

- ④ Code: [https://github.com/cuhkrlcourse/DeepRL-Tutorials/
blob/master/06.DQN_PriorityReplay.ipynb](https://github.com/cuhkrlcourse/DeepRL-Tutorials/blob/master/06.DQN_PriorityReplay.ipynb)

Improving over the DQN

- ① Rainbow: Combining Improvements in Deep Reinforcement Learning.
Matteo Hessel et al. AAAI 2018.
<https://arxiv.org/pdf/1710.02298.pdf>
- ② It examines six extensions to the DQN algorithm and empirically studies their combination



Optional Homework

- ① Go through the Jupyter tutorial and training your own gaming agent:
<https://github.com/cuhkrlcourse/DeepRL-Tutorials>
- ② Good resource for your course project
 - ① Make sure it works for simple environment such as Pong
- ③ Next week: Policy-based RL