

Week 4 CS 1.07.2018

Approximate Value Based Methods

Limitations of Tabular Methods

- in tabular SARSA

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)]$$

- any single cell update does not affect any other cell

Limitation

- # of parameters > # of state

Memory Data Time

Function approximation in RL

of parameters independent of # of state

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

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

Any single parameter affects

values of all states.

It is wise to obtain the following mappings:

$$s \mapsto V_{\pi}(s)$$

$$s, a \mapsto Q_{\pi}(s, a)$$

Reduction to Supervised Learning problem - MC

Monte Carlo
ideal world goals

$$s \mapsto E_{\pi} [G_t | S_t = s]$$

$$s, a \mapsto E_{\pi} [G_t | S_t = s, A_t = a]$$

$$(G_t \stackrel{\Delta}{=} \sum_{k=0}^{\infty} \gamma^k R_{t+k})$$

Sample based estimates of goals
 $g(s)$, $g(s, a)$

$$\begin{aligned} s &\mapsto R(s, \pi(s)) + \gamma G_{t+1} \\ s, a &\mapsto R(s, a) + \gamma G_{t+1} \end{aligned} \quad] \text{ numbers known at the end of an episode}$$

Temporal difference TD

Ideal world goals

$$s \mapsto E_{\pi} [R_{t+2} + \gamma v_{\pi}(S_{t+2}) \mid S_t = s]$$

$$s, a \mapsto [E_{\pi} [R_{t+2} + \gamma v_{\pi}(S_{t+2}) \mid S_t = s, A_t = a]]$$

Sample based estimates of goals $g(s)$, $g(s, a)$

$$s \mapsto R(s, \pi(s)) + \gamma \hat{v}_{\pi}(S_{t+1}, w) \quad \text{functions}$$

$$s, a \mapsto R(s, a) + \hat{v}_{\pi}(S_{t+1}, w) \quad \text{of parameters } w.$$

Loss is the same as for regression problem.

$$L(w) = \frac{1}{2} \sum_{s, a} p_{\pi}(s, a) [g(s, a) - \hat{q}_{\pi}(s, a, w)]^2$$

Weight of "importance" $p_{\pi}(s, a)$ can be thought of as

- fraction of time

- policy encounters state s

- and in that state makes action a

Recap: online, off line, on-policy, off-policy

Off-policy - more general, very difficult

- Behaviour policy - collects data (makes actions)
- Target policy - subject to evaluation & improvement

On-policy - less general, easier

- Behaviour policy equals target policy

Usually, but not always

- Online - update policy during the episode (e.g. TD)

- Offline - update policy after an episode ends (e.g. MC)

Loss functions in value based RL

minimizing loss with gradient loss

$$L(w) = \frac{1}{2} \sum_{s,a} p_\pi(s,a) \overbrace{\left[g(s,a) - \tilde{q}_\pi(s,a,w) \right]^2}^{L_{s,a}(w)}$$

Gradient descent (GD)

$$w \leftarrow w - \alpha \nabla_w L(w)$$

Stochastic gradient descent (SGD)

On-policy: $s, a \sim p_\pi$

Off-policy: $s, a \sim p_b$ (behaviour)

$$w \leftarrow w - \alpha \nabla_w L_{s,a}(w)$$

From gradient to semi-gradient

$$L(w) = \frac{1}{2} \sum_{s,a} p_\pi(s,a) [g(s,a) - \hat{q}_\pi(s,a,w)]^2$$

goal is a function of w
in case of TD

(in MC is a simple number)

Stochastic semi-gradient update treats
goals as fixed:

$$\nabla_w g(s,a) = 0$$

This assumption simplifies math a lot.

$$w \leftarrow w - \alpha \nabla_w L_{s,a}(w)$$

$$w \leftarrow w + \alpha [g(s,a) - \hat{q}_\pi(s,a,w)] \nabla_w \hat{q}_\pi(s,a,w)$$

Semi-gradient update

Treat goals $g(s,a)$ as fixed number

Changes parameters to move estimates
closer to targets.

Ignores effect update on the target

Is not a proper gradient

- No SGD's theoretical convergence properties
- Converges reliably in most cases
- More computationally efficient than true SGD
- Meaningful thing to do.

Target $g(s,a)$ define and what and how
do we learn

SARSA

$$g(s,a) = R(s,a) + \gamma \hat{q}_\pi(s_{t+1}, a_{t+1}, m)$$

Expected SARSA

$$g(s,a) = R(s,a) + \gamma \sum_a \pi(a|s_{t+1}) \hat{q}_\pi(s_{t+1}, a, m)$$

Q-learning

$$g(s,a) = R(s,a) + \max_a \hat{q}_\pi(s_{t+1}, a, m)$$

↑
all s are random variables

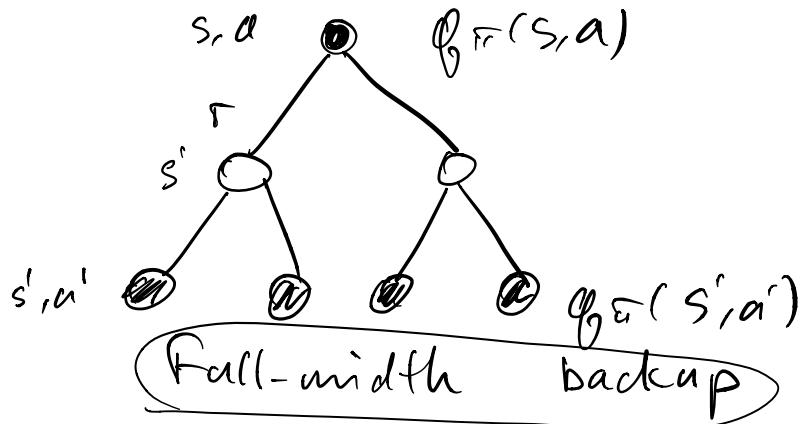
Semi-gradient SARSA (on-policy)

Tabular SARSA

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) +$$

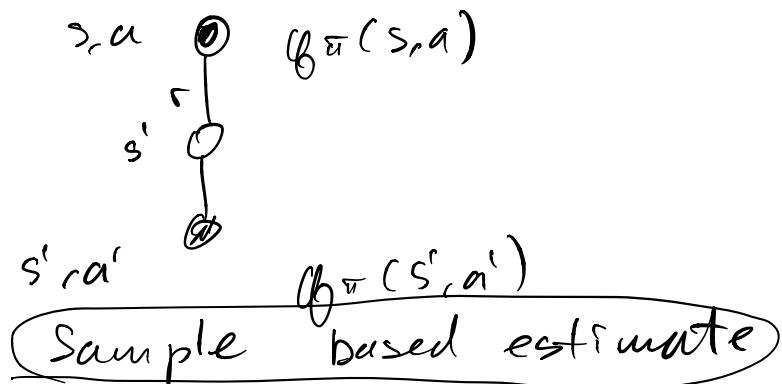
$\leftarrow \gamma [R_{t+1} + \mathbb{E} Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)]$

Bellman expectation equation for $Q(S_t, A_t)$



Approximate SARSA

$$w \leftarrow w + \gamma [R_{t+1} + \hat{q}(S_{t+1}, A_{t+1}, w) - \hat{q}(S_t, A_t, w)] \triangleright \hat{q}(S_t, A_t, w)$$



Semi-gradient expected SARSA (off-policy)

Tabular expected SARSA

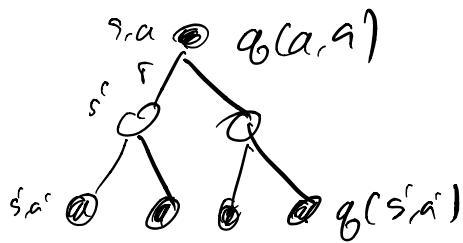
$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) +$$

$$+ \lambda \left[R_{t+2} + \gamma \sum_a \pi(a|S_{t+1}) Q(S_{t+1}, a) - Q(S_t, A_t) \right]$$

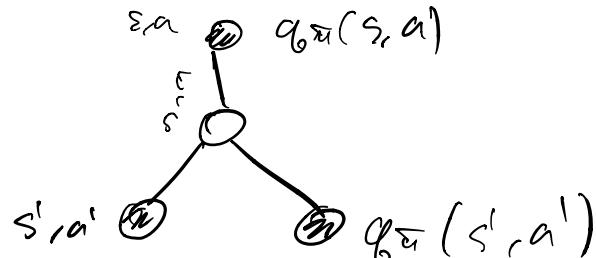
Approximate expected SARSA

$$w \leftarrow w + \lambda \left[R_{t+2} + \gamma \sum_a \pi(a|S_{t+1}) \hat{q}(S_{t+1}, a, w) - \hat{q}(S_t, A_t, w) \right] \cdot \nabla \hat{q}(S_t, A_t, w)$$

Bellman expectation equation for $q(s, a)$



full-width
backup



sample based
estimate

Semi-gradient Q-learning (off-policy)

Tabular Q-learning

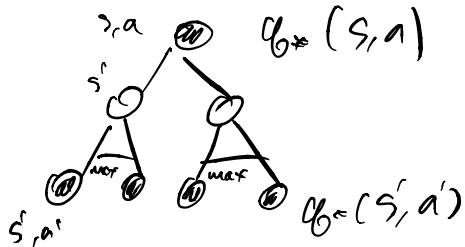
$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \lambda \left[R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t) \right]$$

Approximate Q-learning

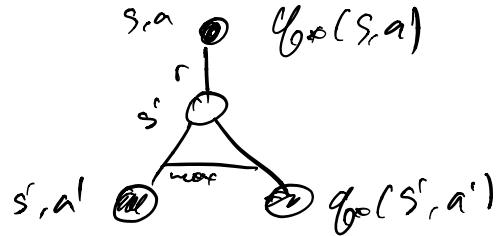
$$w \leftarrow w + \lambda \left[R_{t+2} + \gamma \max_a \hat{q}(S_{t+2}, a, w) - \hat{q}(S_t, A_t, w) \right]$$

$$-\hat{q}(s_t, a_t, m) \} \triangleright \hat{q}(s_t, a_t, m)$$

Bellman optimality equation $q_{\pi^*}(s, a)$



full-width backup



sample based
estimate

Difficulties with Approximate methods

Round (0): general well-known boring problems

- Too many state-action pairs

$$|A| \cdot |S| = 18 \cdot 2^{56} \stackrel{2^{10 \times 160 \times 3}}{\text{(atari)}} \rightarrow 10^{83}$$

- Methods should be sample efficient

- Meaningful decisions need accurate value estimates

- Methods should be very flexible

- Any model has finite # of params

- Even universal approximators are limited.

- Bias-variance tradeoff - overfitting/underfitting
 - Cross validation on scarce data is unreliable

Problems with Supervised Learning in RL

① Unusual data

1. Correlation
2. Dependence on policy
3. Proximity in space and time.

② Nonstationarity

1. Policy improvement
2. Nonstationary task

③ The deadly triad (off-policy learning)

Round (3.2): data-correlated samples

Sequences of highly correlated non-iid data

- Forgetting useful facts and features
- Much of SL rely on iid assumption
- Learning can be inefficient
- SGD loses convergence guarantees

Round (1.2) data-dependence on policy

- Unseen data comes as agent learns new things.
- An agent can learn the fatal behavior
 - New data could be insufficient to unlearn it
- Agent experiments with data stream by exploration

Round (3.3) data-proximity in space and time

- $Q(s, a)$ often changes abruptly in s and a
 - Close states could be arbitrary far in value
 - Successive states could be arbitrary far in value.
 - Unstable gradients, much data needed for SGD
 - Numerical problems.
- TD-specific. Errors in estimates propagate

Approximate Q-learning update

$$w \leftarrow w + \alpha [R_{t+1} + \gamma \max_a \hat{q}(S_{t+1}, a, w) - \hat{q}(S_t, A_t, w)] \triangleright \hat{q}(S_t, A_t, w)$$

Makes $\hat{q}(S_t, A_t, w)$ closer to $R_{t+1} + \gamma \max_a \hat{q}(S_{t+1}, a, w)$

Round (2): inherent nonstationarity - GDT

- Cannot assume fixed training data distribution
 - TD targets are invalidated
 - MC targets no longer apply
- Numeric problems (oscillating behaviour)
 - Small change in Q-values
 - drastic change in policy
 - drastic change in training data

- large gradients
- large update in Q-values.
- The environment can itself be nonstationary

Roual(3): The deadly triad-model divergence

1. Off-policy learning
 - E.g. learning target π_t while following behaviour π_b
2. Bootstrapping
 - Updating a guess towards another guess (TD,DP)
3. Function approximation
 - Using model with # of params smaller than # of states

Divergence is not connected with

- Sampling, exploration, greediness, control
- Complexity of the model (even linear models diverge)

Programming assignments

Approximate Q Learning - Cart Pole

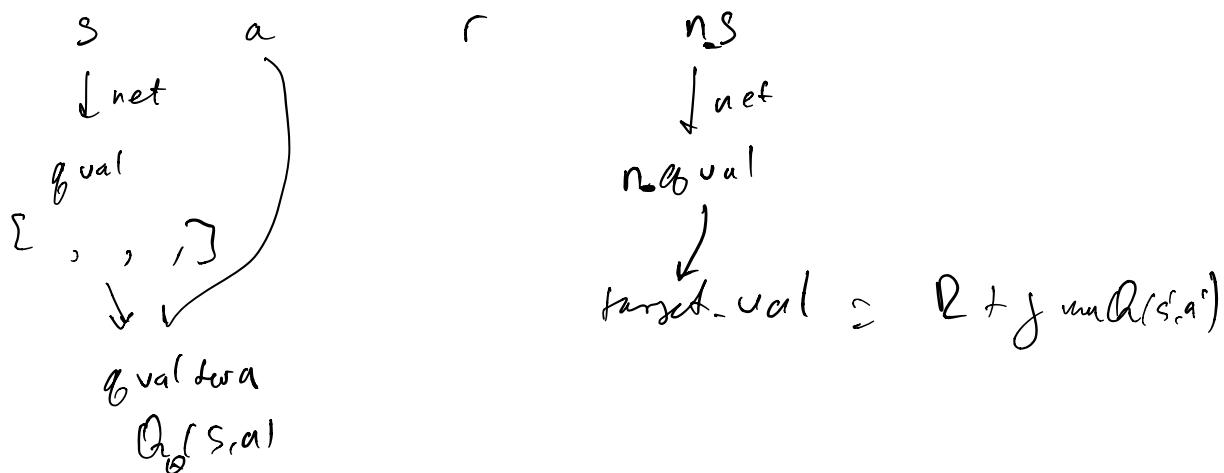
$$L = \frac{1}{N} \sum_i \left(Q_i(s, a) - [r(s, a) + \max_{a'} Q_i(s', a')] \right)^2$$

$Q_\theta \equiv Q_*$ but don't propagate through $Q_*(s, a)$

tabular

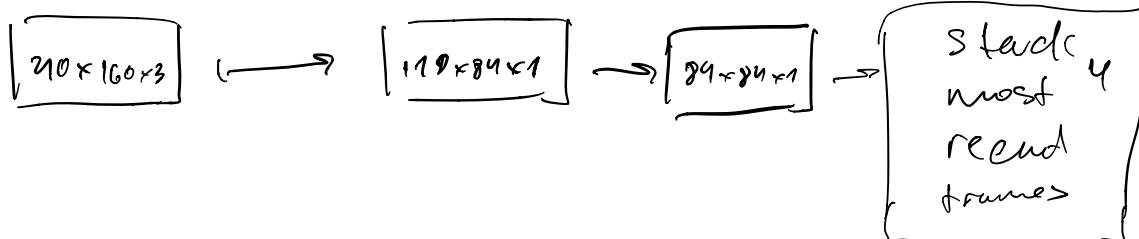
$$V(s) = \max_a Q(s, a)$$

$$Q(s, a) \leftarrow (1-\lambda)Q(s, a) + \lambda \cdot (r + \gamma V(s'))$$



Case Study: Deep Q-Network.

- Epsilon-greedy exploration



Stack n frames = fourth order markov assumption

$$P(R_{t+1}, S_{t+1} | S_0, a_0, \dots, R_t, S_t, a_t) = P(R_{t+1}, S_{t+1} | S_t, a_t)$$

$$P(R_{t+1}, S_{t+1} | S_0, a_0, \dots, R_t, S_t, a_t) = P(R_{t+1}, S_{t+1} | S_t, a_t, \dots, S_{t-3}, a_{t-3})$$

DQN fixes some instability problems

1. Sequential correlated data
 - may hurt convergence and performance
2. Instability of data distribution due to policy changes
 - policy may diverge and/or oscillate
3. Unstable gradients
 - Absolute value of $Q(s, a)$ may vary much across states

Solution \Downarrow

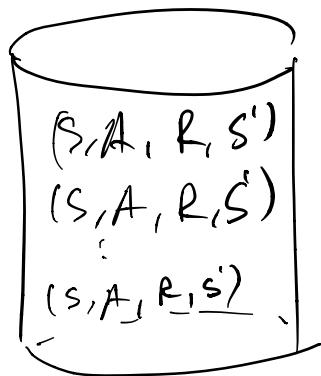
1. Experience replay
2. Target networks
3. Reward clipping

(1) Sequential data \rightarrow Experience replay

Semi-gradient update for approximate Q-learning

$$w \leftarrow w + \alpha [R_{t+1} + \gamma \max_a \hat{q}(S_{t+1}, a, w) - \hat{q}(S_t, A_t, w)] \triangleright \hat{q}(S_t, A_t, w)$$

1. Store tuples (S, A, R, S') in a pool
 2. Sample tuples from pool at random
 3. Update the model of Q-func
 4. Act epsilon-greedy w.r.t Q-func
 5. Add new samples to pool
 6. Go to 2
- ↗ just decorrelate data
if pool is large



only for off-policy learning

Pros / Cons

- + Help against correlations
- + Increase sample efficiency
- Reduce variance of updates
- Computations are easy to parallel
- Analog of sample based model of the world
- Memory intensive (DQN stored 10^6 interactions)
- Random sampling from pool could be improved
- Older interactions were obtained under meater policy

(2) Policy oscillation \rightarrow Target networks

Unwanted source of instability

- Targets are functions of parameters
- Errors in estimates propagate into other estimates

Standard Q-learning update

$$w \leftarrow w + \alpha \{ R_{t+1} + \gamma \max_a \hat{q}(S_{t+1}, a, w) - \hat{q}(S_t, A_t, w) \} \triangleright \hat{q}(S_t, A_t, w)$$

Q-learning update with target networks

$$w \leftarrow w + \alpha \{ R_{t+1} + \gamma \max_a \hat{q}(S_{t+1}, a, w^-) - \hat{q}(S_t, A_t, w) \} \triangleright \hat{q}(S_t, A_t, w)$$

copy parameters to target network

- (hard) Copy w to w^- once in awhile
- (soft) Maintain w^- as moving average of w

(3) Clippable gradients \leftrightarrow Reward clipping

- Don't know the scale of reward beforehand
- Don't want numeric problems with ~~large~~ Q-values.

- Clip the reward to $[-3; 1]$
 - + less peaky Q-values
 - + good gradients
- cannot distinguish between good and great

Honor

DQN: Statistical issues

Recap: Approximate Q-learning

$$Q(s, a_0), Q(s, a_1), Q(s, a_2) \quad \hat{Q}(s_t, a_t) = \underbrace{r + \gamma \max_{a'} Q(s_{t+1}, a')}_{\text{Q-values}}$$

model

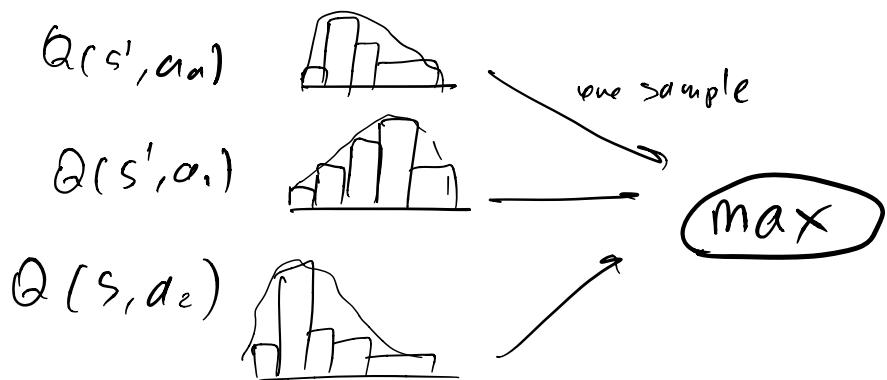
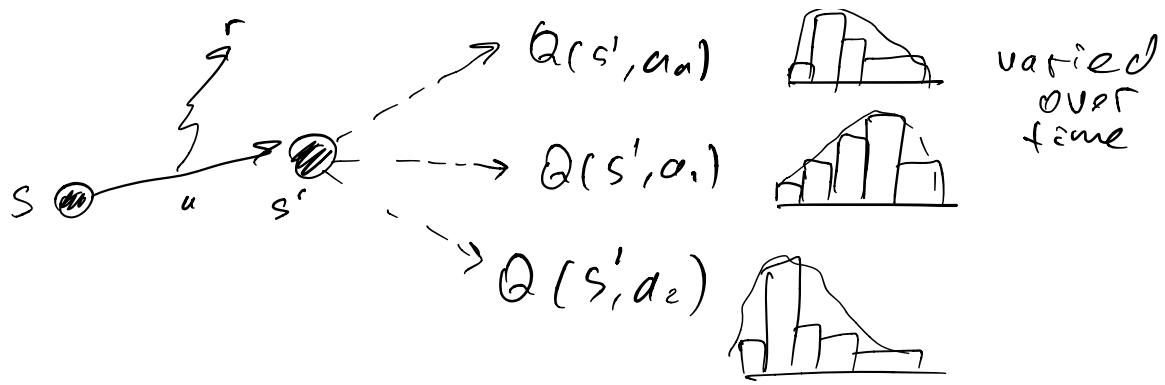
$$\mathcal{L} = \left(Q(s_t, a_t) - \left[r + \gamma \max_{a'} Q(s_{t+1}, a') \right] \right)^2 \quad \text{Objective}$$

W-parameters

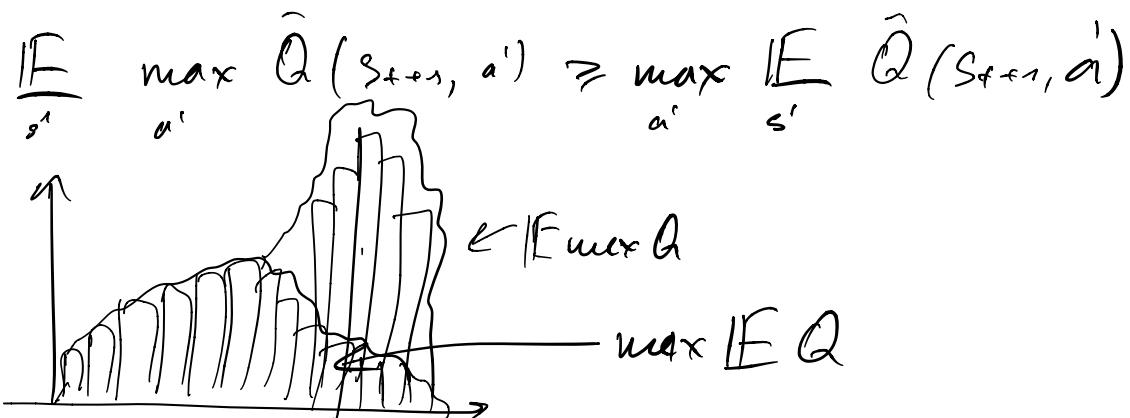
$$w_{t+1} = w_t - \alpha \frac{\partial \mathcal{L}}{\partial w} \quad \text{Gradient step}$$

image

$$\mathcal{L} = \left(Q(s_t, a_t) - \left[r + \gamma \underbrace{\max_{a'} Q(s_{t+1}, a')}_{\text{consider const}} \right] \right)^2$$



so



$$E \max(\xi, \zeta) \geq \max(E\xi, E\zeta)$$

Double Q-learning

Idea: train two Q-functions separately.

Q_1 and Q_2

Objective for Q_1

$$\hat{Q}_1(s_t, a_t) = r_t + \gamma \cdot Q_2(s_{t+1}, \underset{a^*}{\operatorname{argmax}} Q_2(s_{t+1}, a^*))$$

Objective for Q_2

$$\hat{Q}_2(s_t, a_t) = r_t + \gamma Q_1(s_{t+1}, \underset{a^*}{\operatorname{argmax}} Q_1(s_{t+1}, a^*))$$

Algorithm:

Initialize Q_1, Q_2

forever:

sample s_t, a_t, r_t, s_{t+1} from env

if flip_coin() == 'heads':

$$\tilde{Q}_1(s_t, a_t) = r_t + \gamma Q_2(s_{t+1}, \underset{a^*}{\operatorname{argmax}} Q_2(s_{t+1}, a^*))$$

$$Q_1(s_t, a_t) \leftarrow Q_1(s_t, a_t) + \frac{1}{2} [\tilde{Q}_1(s_t, a_t) - Q_1(s_t, a_t)]$$

else:

$$\tilde{Q}_2(s_t, a_t) = r_t + \gamma Q_1(s_{t+1}, \underset{a^*}{\operatorname{argmax}} Q_1(s_{t+1}, a^*))$$

$$Q_2(s_t, a_t) \leftarrow Q_2(s_t, a_t) + \frac{1}{2} [\tilde{Q}_2(s_t, a_t) - Q_2(s_t, a_t)]$$

Now consider DQN

- Neural net policy
- Experience replay
- Target networks

Idea 1: train two networks separately.

Idea 2: use target network instead ($Q_2 =$)

Classic DQN:

$$\hat{Q}(s_t, a_t) = r_t + \gamma \max_{a^*} Q^{\text{old}}(s_{t+1}, a^*)$$

Double DQN:

$$\hat{Q}(s_t, a_t) = r_t + \gamma \cdot Q^{\text{old}}(s_{t+1}, \arg \max_{a^*} Q(s_{t+1}, a^*))$$

step-by-step

$$Q(s_t, a_t) = r + \gamma \max_{a'} \hat{Q}(s_{t+1}, a')$$



$$\hat{Q}(s_t, a_t) = r + \gamma \cdot \hat{Q}(s_{t+1}, \arg \max_a \hat{Q}(s_{t+1}, a))$$



$$\tilde{Q}(s_t, a_t) = r + \gamma \cdot Q^{\text{old}}(s_{t+1}, \arg \max_{a'} \hat{Q}(s_{t+1}, a'))$$

More DQN tricks

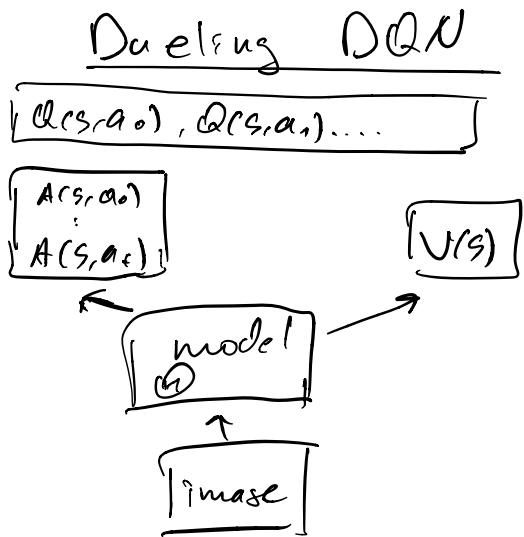
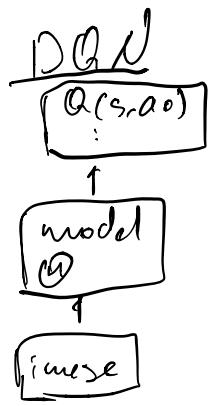
$$\text{Decomposition: } Q(s, a) = V(s) + A(s, a)$$

where $Q(s, a)$ are action values, as usual

$V(s)$ are state values ($V^*(s)$ or $V_\pi(s)$)

$A(s, a)$ is the "advantage"

$$A(s, a) = Q(s, a) - V(s)$$



Estimating $A(s, a)$

option 1:

$$A(s, a_i) = nn(s)[i] - \max_j nn(s)[j]$$

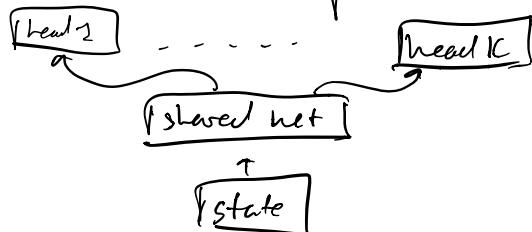
option 2:

$$A(s, a_i) = nn(s)[i] - \frac{1}{|A|} \sum_j nn(s)[j]$$

where $nn(s)$ is a neural network activation
for advantage (left branch)

Bootstrap DQN

Each head predicts all q -values



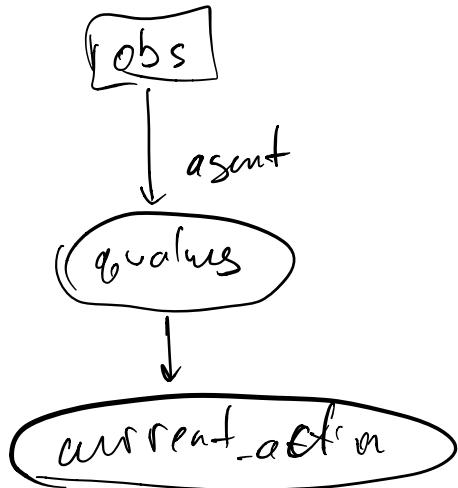
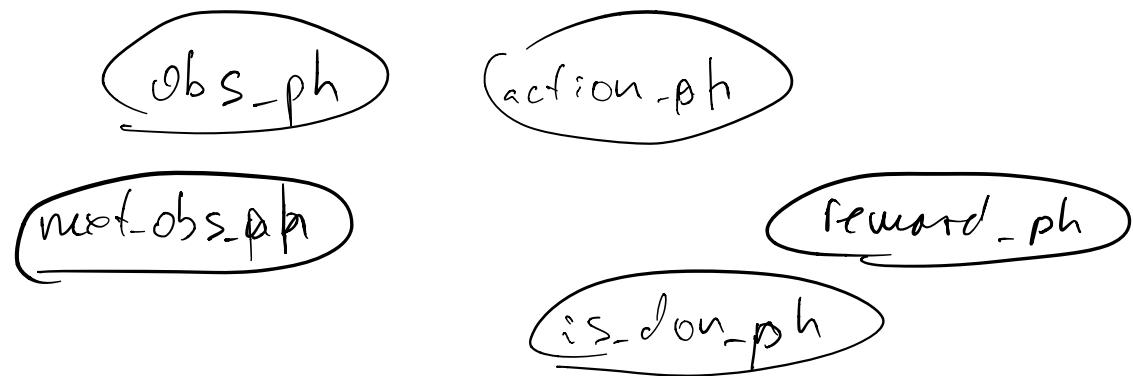
Why: make exploration great again!

Plan: maintain K separate heads
(for example, last 2 layers) and
shared body weights.

↑
pick random head from K.

DQN - Breakout

$$Q_{\text{reference}}(s, a) = r + \gamma \max_{a'} Q_{\text{target}}(s', a')$$



$$L = \frac{1}{N} \sum_i \left\{ Q_\theta(s, a) - Q_{\text{target}}(s, a) \right\}^2$$

$$Q_{\text{target}}(s, a) = r(s, a) + \gamma \max_{a'} Q_{\text{target}}(s', a')$$

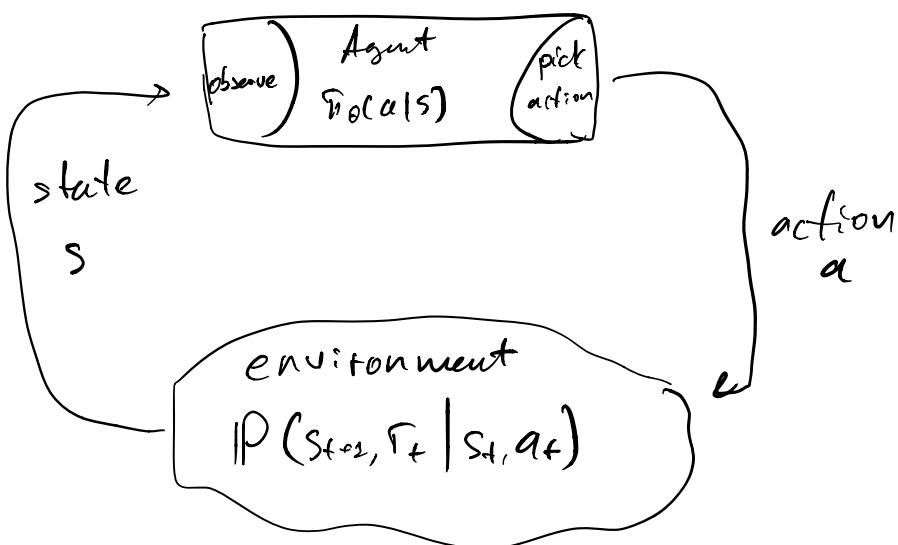
Partial observability

Problem: In most real environments agent observation does not hold all information about system state

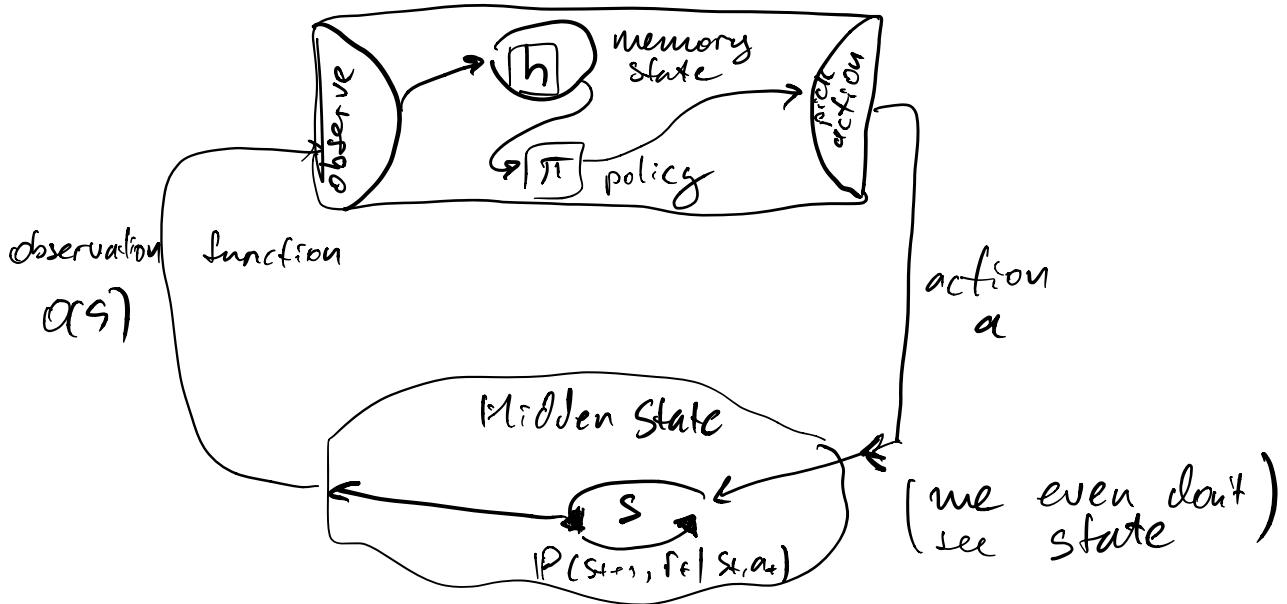
Examples:

- Robots: field of view
- Trading bots: market state
- Atari Breakout: ball velocity

Recap: MDP



Partial Observable MDP



N-step heuristic

Idea:

$$s_t \neq O(s_t)$$

$$s_t \approx (O(s_{t-n}), a_{t-n}, \dots O(s_{t-1}), a_{t-1}, O(s_t))$$

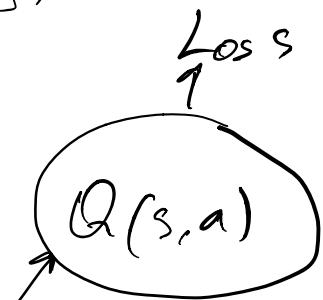
e.g. ball movement in breakout

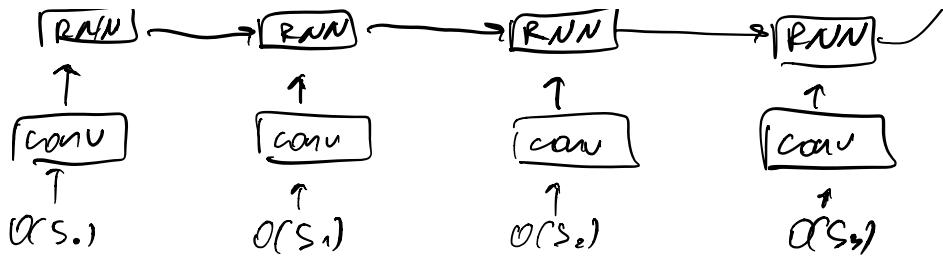
Recurrent neural network

$$L = (Q(s_t, a_t) - \sum \gamma^j \max_a Q(s_{t+j}, a'))^2$$

$$Q(s_t, a_t) \approx Q(h_t(h_{t-2}, O(s_t)))$$

$$h_{t+2} = \sigma(W_{\text{mid}} \cdot h_t + W_{\text{inp}} \cdot O(s_t) + b)$$





char n rule unroll

$$\frac{\partial \text{Loss}}{\partial P(f_a)} \cdot \frac{\partial P(f_a)}{\partial h_3} \cdot \frac{\partial h_3}{\partial h_n} \cdot \frac{\partial h_n}{\partial h_3} \cdot \frac{\partial h_3}{\partial h_2} \cdot \frac{\partial h_2}{\partial h_1} \dots$$

DRQN

- use LSTM on top of DQN body
- predict Q-values from hidden state
- train on partial trajectories
- uses other existing DQN tricks

Week 5 CS 11.07.2018

Policy-based methods

Policy-based RL vs Value-based RL

Intuition

Approximation error

$$L = \mathbb{E} \left[r + \gamma \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t) \right]^2$$

counter-example

	True	A	B
$Q(s_0, a_0)$	3	1	2
$Q(s_0, a_1)$	2	2	1
$Q(s_1, a_0)$	3	3	3
$Q(s_1, a_1)$	100	50	100

A is better →
→ optimal policy

B is worse →
→ but L^2 smaller.

All kinds of Policies

Policies

Deterministic policy $\alpha = \pi_\theta(a|s)$

Stochastic policy $\alpha \sim \pi_\theta(a|s)$
(sampling takes care of exploration)

Two approaches

Value based

Learn value function $Q_\theta(s, a)$ or $V_\theta(s)$

Infer policy $\pi_\theta(a|s) = \#\left[a = \underset{a}{\operatorname{argmax}} Q_\theta(s, a)\right]$

Policy based

Explicitly learn policy $\pi_\theta(s, a)$ or

$\pi_\theta(s) \mapsto a$. Implicitly maximize reward over policy

Policy gradient formalism

Main idea

Objective

• Expected reward

$$J = \mathbb{E}_{\substack{s \sim p(s) \\ a \sim \pi_\theta(s|a) \\ \dots}} R(s, a, s', a', \dots)$$

• Expected discounted reward

$$J = \mathbb{E}_{\substack{s \sim p(s) \\ a \sim \pi_\theta(s|a)}} G(s, a)$$

$$G(s, a) = r + \gamma \cdot G(s', a')$$

• Consider an n-step process for simplicity

$$J = \mathbb{E}_{\substack{s \sim p(s) \\ a \sim \pi_\theta(s|a)}} R(s, a) =$$

$$= \int_s p(s) \cdot \int_a \pi_\theta(s|a) \cdot R(s, a) da ds$$

$$J \approx \frac{1}{N} \sum_{i=0}^N \sum_{s, a \in \mathcal{Z}^i} R(s, a) \leftarrow \begin{array}{l} \text{Reward for agent's action} \\ \text{Sample } N \text{ sessions} \\ \text{by following } \pi_\theta(a|s) \end{array}$$

The log-derivative trick

$$\nabla \log \pi(z) = \frac{1}{\pi(z)} \cdot \nabla \pi(z)$$

$$\pi \cdot \nabla \log \pi(z) = \nabla \pi(z)$$

Analytical inference

$$J = \int_s p(s) \int_a \pi_\theta(a|s) R(s, a) da ds$$

$$\nabla_\theta J = \nabla_\theta \int_s p(s) \int_a \pi_\theta(a|s) R(s, a) da ds =$$

$$\begin{aligned}
 &= \int p(s) \sum_a \nabla_{\theta} \pi_{\theta}(a|s) R(s, a) da ds = \\
 &= \int_s p(s) \cdot \underbrace{\int_a \nabla \log \pi_{\theta}(a|s) R(s, a) da}_{\text{that's expectation}} ds
 \end{aligned}$$

Reinforce

Policy gradient

• Policy gradient

$$\nabla J = \mathbb{E}_{\substack{s \sim p(\cdot) \\ a \sim \pi_{\theta}(s|a)}} \nabla \log \pi_{\theta}(a|s) \cdot Q(s, a)$$

• Approximate with sampling

$$\nabla J \approx \frac{1}{N} \sum_{i=0}^N \sum_{s_i, a_i \in z_i} \nabla \log \pi_{\theta}(a_i|s_i) \cdot Q(s_i, a_i)$$

REINFORCE algorithm [on-policy]

• Initialize NN weights $\theta_0 \leftarrow \text{random}$

• Loop

- Sample N sessions z under current $\pi_{\theta}(a|s)$

- Evaluate policy gradient

$$\nabla J \approx \frac{1}{N} \sum_{i=0}^N \sum_{s,a \in z_i} \nabla \log \pi_\theta(s) \cdot Q(s, a)$$

Ascend $\theta_{i+1} \leftarrow \theta_i + \lambda \cdot \nabla J$

$$Q(s, a) = V(s) + A(s, a)$$

Action s influence $A(s, a)$ only, so $V(s)$ is irrelevant

$$\nabla J \approx \frac{1}{N} \sum_{i=0}^N \sum_{s,a \in z_i} \nabla \log \pi_\theta(s) \cdot (Q(s, a) - b(s))$$

Anything that doesn't depend on action is really $b(s) = V(s)$
 $b(s)$ - baseline.

Programming Assignment: Reinforce

$$E(p) = - \sum_x p(x) \log p(x)$$

$$E(p) = - \int p(x) \log p(x)$$

$$G_t \triangleq R_t + \gamma R_{t+1} + \gamma^2 \cdot R_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k} = \\ = R_t + \gamma G_{t+1}$$

$$T: \emptyset \dots T$$

$$T, \dots, T$$

$$G_0 = \sum_{k=0}^T \gamma^k R_k \quad G_T = \sum_{k=0}^T \gamma^k R_{T+k} = \left\{ k=0 \right\} =$$

$$\left[0, 0, 1, 0, 0, 1, 0 \right] = \gamma^0 \cdot R_T = R_T$$

$G_0 = 0$
 $G_1 = 1 + \gamma \cdot G_0$
 $G_2 = 0 + \gamma \cdot G_1$

Actor-critic

Advantage actor-critic

- Learn both $V(s)$ and $\pi_\theta(a|s)$
- Hope for best of both worlds

- ideal: learn both $\pi_\theta(a|s)$ and $V_\theta(s)$
- Use $V_\theta(s)$ to learn $\pi_\theta(a|s)$ faster!

Non-trivial: how can we estimate $A(s, a)$ from (s, a, r, s') and V-function?

$$A(s, a) = Q(s, a) - V(s)$$

$$Q(s, a) = r + \gamma \cdot V(s')$$

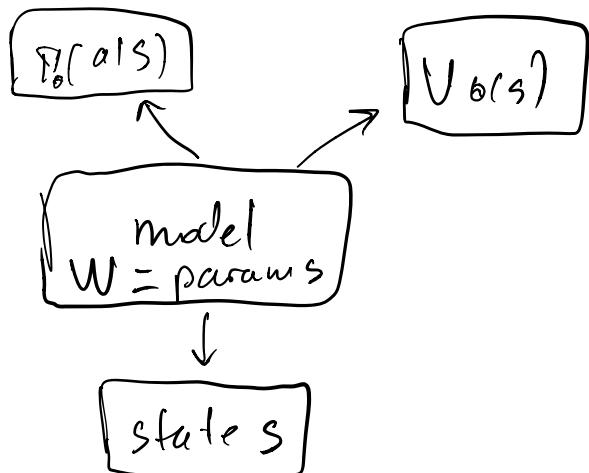
$$A(s, a) = r + \gamma \cdot V(s') - V(s)$$



$$A(s, a) = r + \gamma \cdot V(s') - V(s)$$

$$Q(s, a) = r + \gamma \cdot V(s')$$

$$\triangleright J_{\text{actor}} \approx \frac{1}{N} \sum_{i=0}^N \sum_{s,a \in z_i} \triangleright \log \pi_\theta(a|s) \cdot \underbrace{A(s, a)}_{\text{consider const}}$$



$$\triangleright J_{\text{actor}} \approx \frac{1}{N} \sum_{i=0}^N \sum_{s,a \in z_i} \triangleright \log \pi_\theta(a|s) \cdot A(s, a)$$

$$L_{\text{critic}} \approx \frac{1}{N} \sum_{i=0}^N \sum_{s,a \in z_i} (V_\theta(s) - [r + \gamma \cdot V(s')])^2$$

Duct tape zone

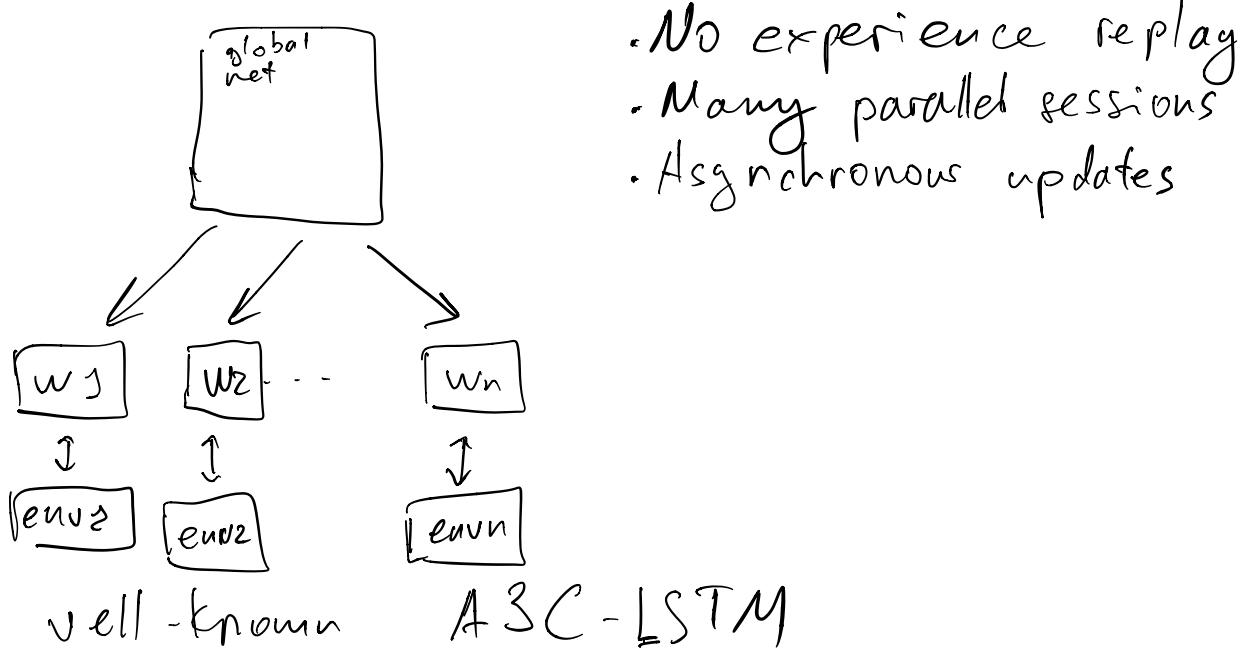
- $V(s)$ errors less important than in Q-learning
- actor still learns even if critic is random, just slower

- Regularize with entropy
 - to prevent premature convergence
- Learn on parallel sessions
 - or super-small experience replay
- Use log softmax for numerical stability

Policy-based vs Value-based

Value base	Policy base + Actor-critic
Q-learning SARSA value-iteration	REINFORCE Advantage Actor-Critic Crossentropy Method
Off policy easier from same	R Stochastic policy from same
<ul style="list-style-type: none"> Solves harder problem Explicit exploration Evaluate states and action Easier to train off-policy 	<ul style="list-style-type: none"> Solves easier problem Innate exploration and stochasticity Easier for continuous actions Compatible with supervised learning

Case study: A3C



Combining supervised & reinforcement learning

Supervised pre-training

Reinforcement learning usually takes long to find optimal policy completely from scratch

We can use existing knowledge to help it.

- Human experience
- Known heuristic
- Previous system

so

- Supervised learning: $\left[\log \pi_\theta(a_{\text{opt}}|s) \right]$

$$\triangleright \ell_h = \underset{x, s_{\text{opt}} \sim D}{\mathbb{E}} \log P_\theta(a_{\text{opt}}|x)$$

- Policy gradient:

$$\triangleright J = \underset{\substack{s \sim d(s) \\ a \sim \pi(a|obs(s))}}{\mathbb{E}} \log \pi_\theta(a|s) \cdot Q(a, s)$$

main difference:
D - reference
 $s \sim d(s)$ - generated

More out there

This domain is huge!

- Trust Region Policy optimization
 - policy gradient on steroids
- Deterministic policy gradients
 - off-policy & less variance
- Many many more:
 - Research happens as you read this line

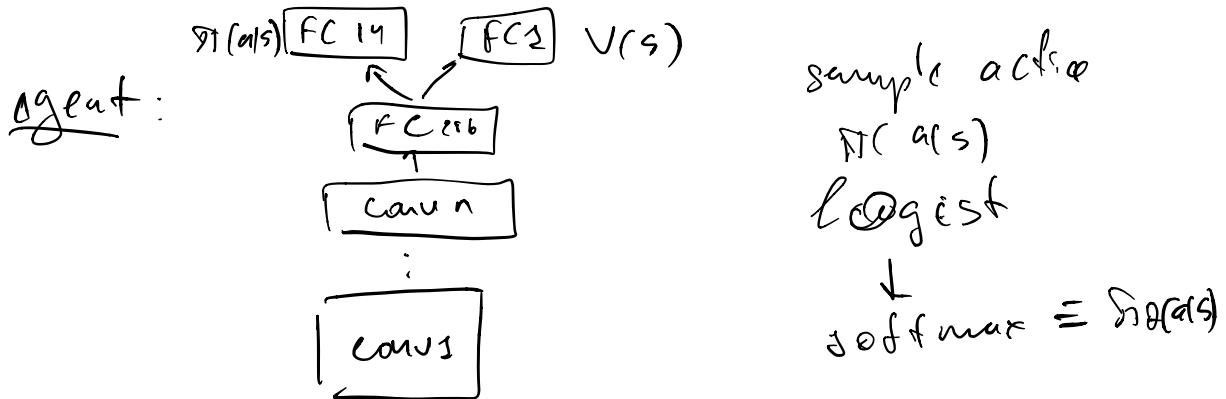
A policy-based agent

$G(s, a)$ - discounted reward

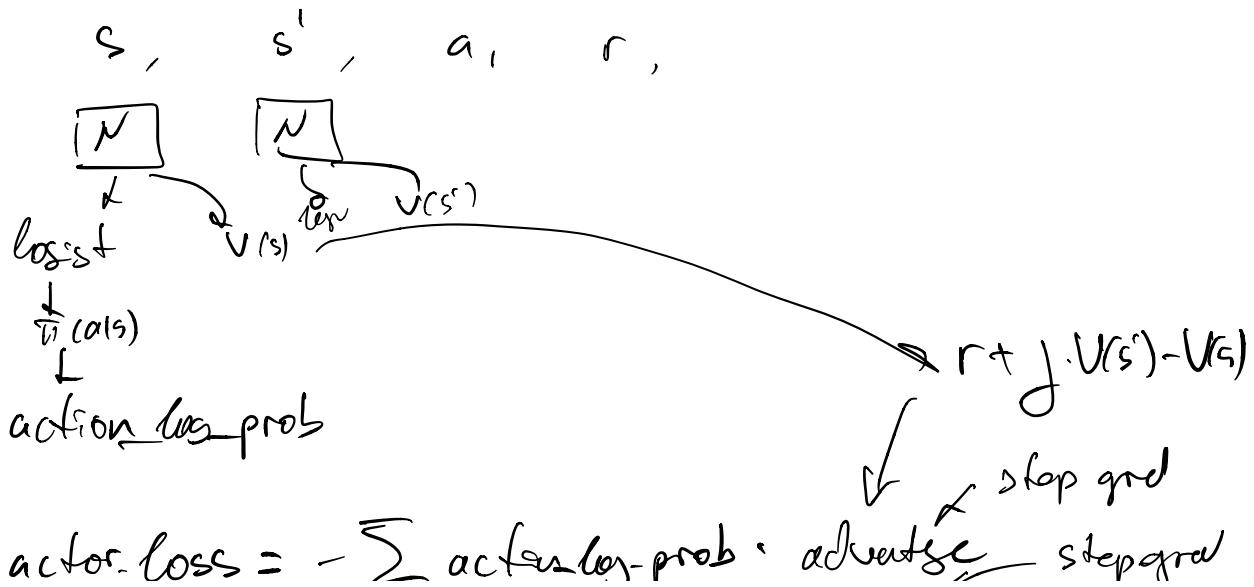
$r(s, a)$ - immediate reward
gamma

Practice Actor-Critic

env: $\begin{pmatrix} \text{uses RL theory} \\ \text{non-action IY} \end{pmatrix}$



algorithm:



$$\text{critic-loss} = -\sum \left\{ V(s) - \left[r + \gamma \cdot V(s') \right]^2 \right\}$$

Week 6

22.07.2018

Measuring exploration

Recap: bandits

Exploration & exploitation

previously

Value-based methods:

We know $Q(s, a) \rightarrow$ we know
optimal policy

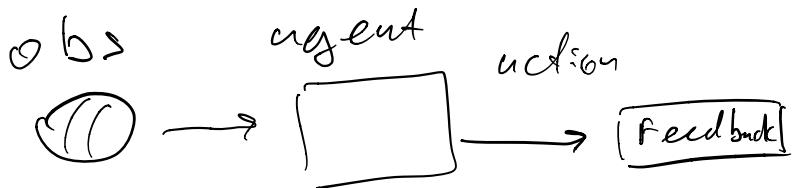
Model-free setting: don't have $p(s', r|s, a)$

Q-learning, SARSA, Expected Value SARSA

Large states spaces

Approximate Q-learning, DQN,
Double DQN.

Multi-armed Bandits



Examples

- banner ads (RTB)
- recommendations
- medical treatment

Basically it's a-step MDP where

- $g(s, a) = r(s, a)$
- $Q(s, a) = \mathbb{E} r(s, a)$
- 50% shorter formulae

Regret: measuring the quality of exploration

Regret of policy $\pi(a|s)$:

Consider an optimal policy, $\pi^*(a|s)$

Regret per tick = optimal - yours

$$h(t) = \sum_{\substack{s,a \sim \pi^* \\ s,a \text{ sample}}} [\mathbb{E} r(s,a) - \mathbb{E}_{\pi} r(s,a)]$$

finite horizon: $t < \max_t$

infinite horizon: $t \rightarrow \infty$

The message just repeat: Regret, Regret, Regret

Strategies:

- ϵ -greedy

With probability ϵ take a uniformly random action: Otherwise take optimal action

- Boltzman

Pick action proportionally to transformed Q-values

$$\pi(a|s) = \text{softmax}(Q(s,a)/\epsilon)$$

- Policy based: add entropy

Uncertainty-based exploration

Intuitive explanation

Thompson Sampling

Uncertainty in rewards

We want to try actions if we believe there's a chance they turn out optimal.

Idea: let's model how certain we are that $Q(s, a)$ is what we predicted.

• Policy:

$P(Q|data)$

- sample once from each Q distribution
- take argmax over samples
- which actions will be taken?

Optimism in face of uncertainty

Idea:

Prioritize actions with uncertain outcomes!

More uncertain = better

Greater expected value = better

Math: try until upper confidence bound is small enough.

Policy:

- Compute 95% upper confidence bound for each a
- Take action with highest confidence bound
- What can we tune here to explore more/less?

UCB-1 (upper confidence bound)

Frequentist approach

There's a number of inequalities that bound

$\Pr(X > t) < \text{something}$

E.g. Hoeffding inequality (arbitrary)
 $t \in [0, 1]$

$$\Pr(X - \mathbb{E}X \geq t) \geq e^{-2nt^2}$$

UCB for bandits

Take actions in proportion to

$$\tilde{\mu}_a = \bar{\mu}_a + \sqrt{\frac{2 \log N}{n_a}}$$

- N number of time-steps so far
- n_a times action a is taken

Upper conf. bound for $r \in [0, 1]$
If not - divide by r_{\max}

Bayesian UCB

The usual way:

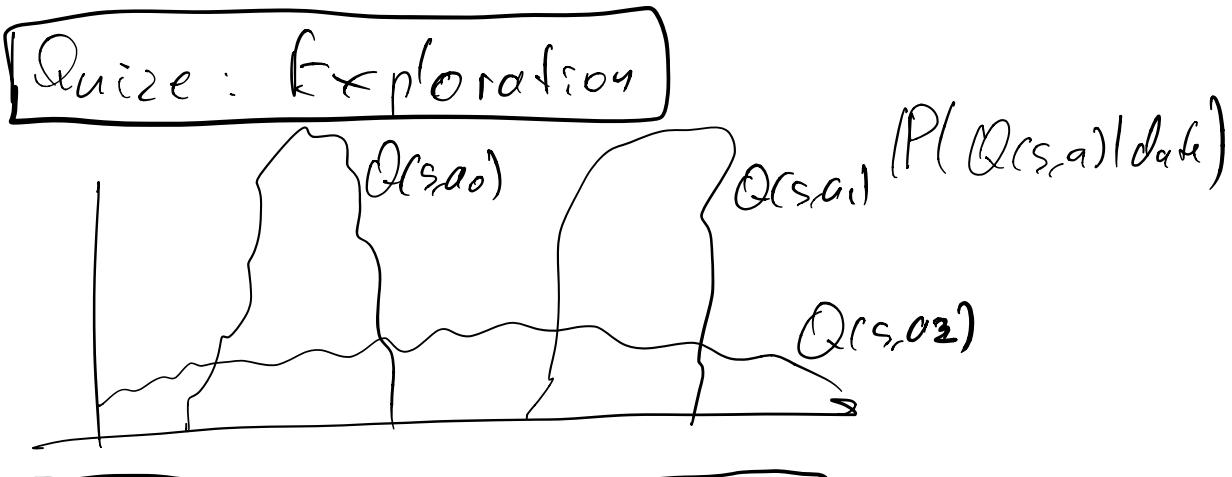
- Start from prior $P(Q)$
- Learn posterior $P(Q | \text{data})$
- Take q -th percentile

What models can learn that?

Approach 1: learn parametric $P(Q)$ e.g. normal

Approach 2: use bayesian neural networks

- UCB vs Bayesian UCB
 - choose any distribution you want
 - learn actual distribution, not upper bound
 - Only as good as your prior.
- Regret: "money you could have saved"
 - if you knew optimal actions from the get go
 - Less regret = better exploration
- ε-greedy and boltzmann never behave optimally
 - Make them greedy in the limit
- Using uncertainty
 - Thompson Sampling, UCB-ε, Bayesian UCB



Programming Assignment: Bandits

K actions

$$r \in \mathbb{R}$$

$$0 \leq \theta_k \leq 1$$

$$J(T) = T \cdot \theta^* - \sum_{t=1}^T r_t$$

$$\theta^* = \max_k \{ \theta_k \}$$

K-pillows T-patients

Bernoulli: n actions: n

$$p_1, \dots, p_n$$

$$\text{optimum reward } \max(p_i)$$

Epsilon-greedy agent:

for $k = 1 \dots K$:

$$\hat{\theta}_k = d_k / (d_k + \beta_k)$$

$x_t \leftarrow \arg\max \hat{\theta}_k$ with prob $1-\epsilon$ else random.

apply x_t, r

$$(d_{k+1}, \beta_{k+1}) \leftarrow (d_k, \beta_k) + (r_t, 1-r_t)$$

done ■

Planning with Monte Carlo Tree Search

Introduction to planning

Monte Carlo Tree Search

Exploration reduces uncertainty

Model-free setting

- Unknown environment dynamics $p(s', r | s, a)$

Model-based settings

- Known model of the world
- Distribution model: $p(s', r | s, a)$
- Sample model: $(s', r) \sim G(s, a)$
Still don't know the global values
of any actions! (we know only immediate r)

We don't know global value estimates

Given the model of the world, what is the best action?

1. Immediate rewards are insufficient
2. Time constraints on recovering sum of future rewards

Planning

Trans form

- know immediate rewards & dynamics into
- value / action-value estimates / explicit policy

Types of planning

Model-free - learn the policy by trial and error

Model-base - plan to obtain the policy

- **Background** planning (DP)

- learning from simulated experience

- improves estimates in arbitrary states

- **Decision-time** planning (heuristic search, MCTS)

- focuses on current state

- plans action selection for current state only

- commit an action only after planning is finished

Guided planning - heuristic search

1. Exhaustive

2. Full-width with fixed depth

3. Full-width with fixed depth, approximate at leaves.

4. Heuristic.

Heuristic search - an umbrella term for many algorithms

- lookahead planning guided by heuristic function
 - select only "valuable" nodes to expand

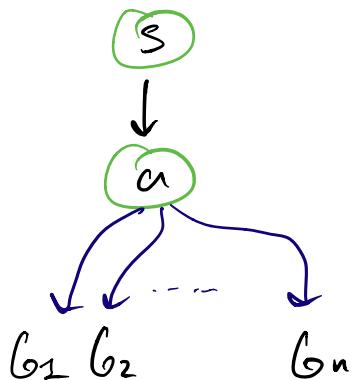
Heuristic search for RL

- + resources focused only on valuable paths
- + nearest states contribute the most to the return
- bad planning results when the model is unreliable
- dependence on the quality of heuristic

Obtain value estimates with rollouts

1. While within computational budget:
 - a. From state s make action a
 - b. Follow **rollout** policy until episode terminates
 - c. Obtain return g_i
2. Output $\hat{q}(s, a) = \frac{1}{n} \sum_{i=1}^n g_i$

- No function approximation
- + Trials are independent \rightarrow parallelization
- Dependence on # of returns averaged



Using rollouts for action selection

Greedy policy - policy improvements as in
Value Iteration

$$\tilde{\pi}(s) \leftarrow \arg \max_a \hat{q}(s, a)$$

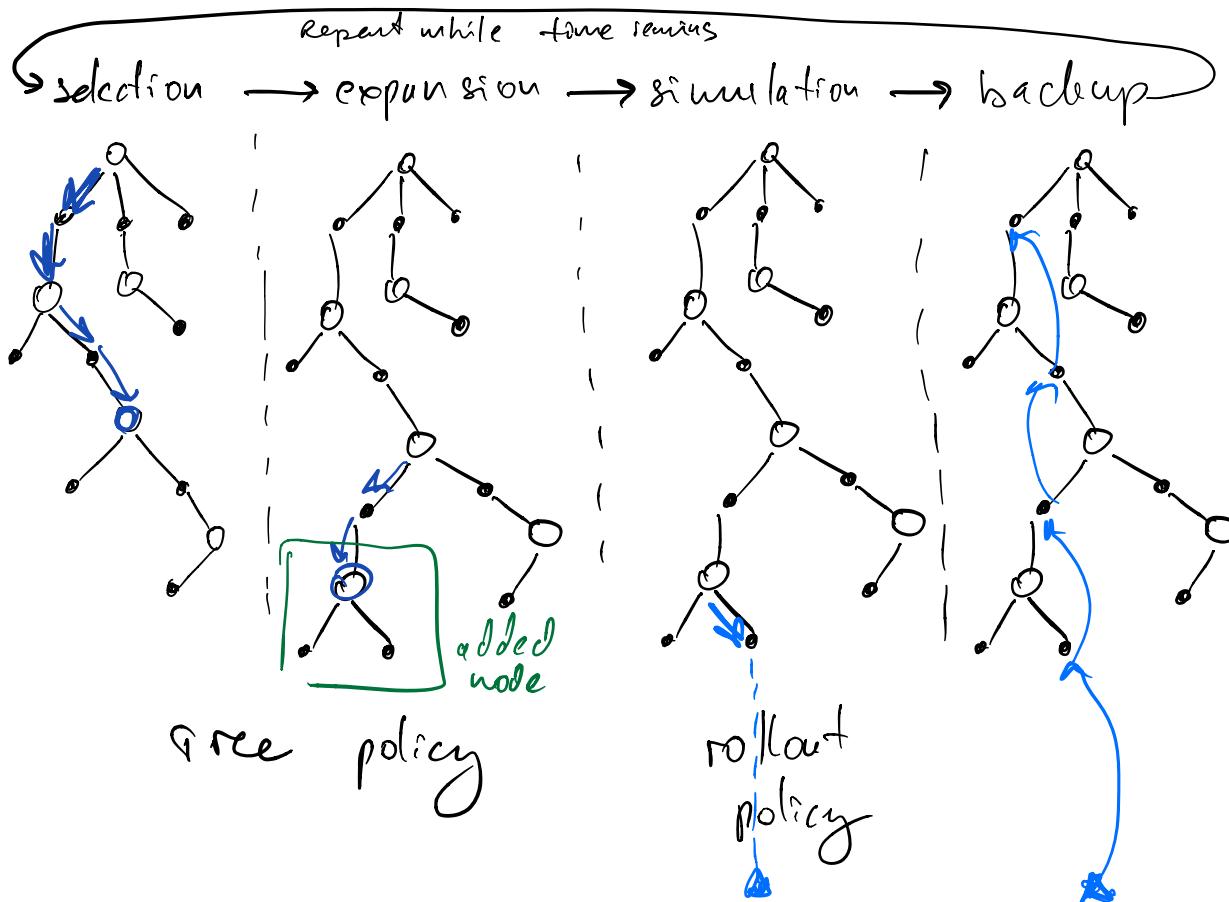
- No estimation of q^* } Good rollout
- No further policy improvement } policy is needed!
- Good rollout policy may require much of valuable time
 - May result in less reliable estimates of $\hat{q}(s, a)$
- Dependence on quality of rollout policy
- + Quite good results even for random rollout policy
- No estimates preserved between successive states
- Does not respect uncertainty in action-value estimates.

Monte Carlo Tree Search

Monte Carlo Tree Search at a glance

- Maintain two policies:
 - tree policy - respects uncertainty, guides the search

- rollout policy estimates the return



- Exploitation - actions with existing good estimates

- Exploration - rarely tested actions

Treating action selection as
multi-armed bandit (UCB)

$$\text{tree}(s) = \arg \max_a \hat{q}(s, a) + 2C \sqrt{\frac{\log N(s)}{N(s, a)}}$$

QUIZ: MCTS

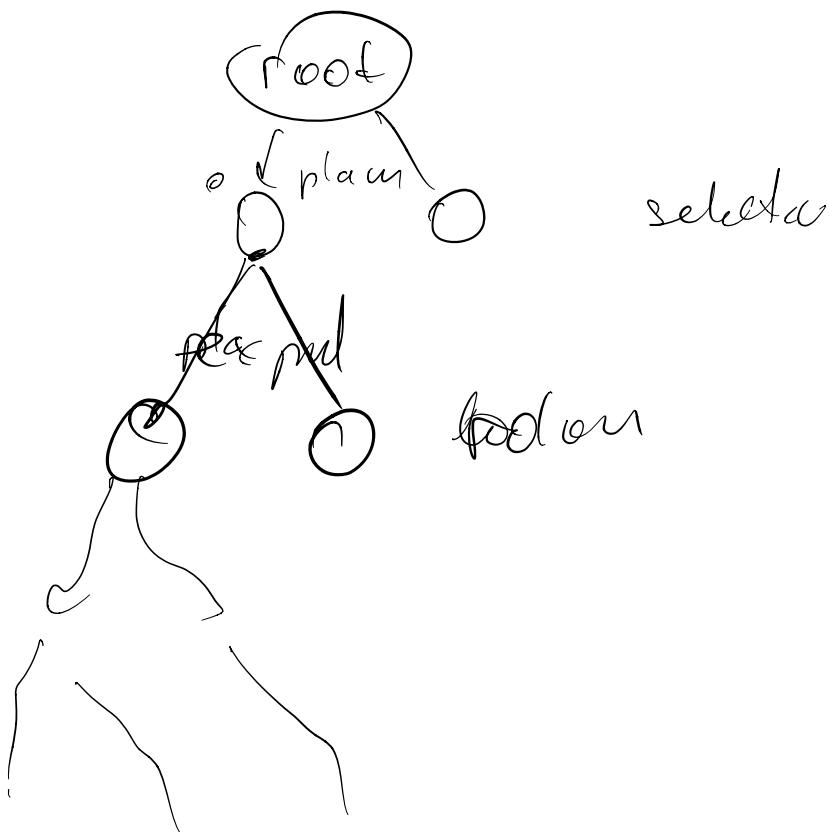
UCB score

$$\frac{x}{y} - \text{numerator}$$
$$\frac{y}{z} - \text{denominator}$$

value - sum
times visited

Node:

is_leaf
is_root
get_mean_value = $\frac{\text{value_sum}}{\text{value_count}}$
ucb_score: mu + ucb



■ RL outside games

General formalism

- Maximize $J = \mathbb{E} \sum_{s \sim p(s)} R(s, a) \text{ over } \pi_{\text{obs}}$
- $R(s, a)$ or $R(\text{session})$ is a black box
 - special case: $R(s, a) = r(s, a) + \gamma R(s', a')$
- Markov property: $P(s' | s, a, \pi) = P(s' | s, a)$
- Special case: $\text{obs}(s) = s$, fully observable

Reinforcement Learning for seq2seq

Encoder - decoder a_i^{pi} / translate
/ score

self-critical policy gradient / weight

$$\nabla J = \mathbb{E}_{x \sim p(x)} \mathbb{E}_{y \sim \pi/g(x)} \log \pi(g|x) \cdot (R(x, s) - b(x))$$

R - negative levenshtein distance

$$b(x) = R(x, \text{greedy}(x))$$

