Last time: Inverse RL

| Given | Find |
|-------|------|
| $S, A, T, \{\tau_i\}$ | $R$ |

This Time: Exploration ~~Hierarchical~~



Montezuma's Revenge    Atari

Bandit              ? ? ?

$\varepsilon$-greedy

softmax              $e^{TQ}$

$\rightarrow$ UCB              $Q + C\sqrt{\frac{2\ln T}{N}}$          $\longleftarrow$ Optimistic

$\rightarrow$ Thompson Sampling       distribution over models          Posterior Sampling

                              sample, take optimal
                                        w.r.t. sample

$*$ Bayes Optimal       BAMDP

                        POMDP where unknown model params
                                part of state

$O(\log n)$ Regret

$$IG(z, y \mid a) = E_y\left[H(\hat{p}(z)) - H(\hat{p}(z) \mid y) \mid a\right]$$

↑ params   ↖ action

$$g(a) = IG(z, y \mid a) \qquad \Delta(a) = E[r(a^*) - r(a)]$$

$\rightarrow$ argmin $\frac{\Delta(a)^2}{g(a)}$

                        Russo + Van Roy    "Learning to
                                            optimize with
                                            info directed
                                            sampling"

Optimistic
- new state = good state
- Exploration bonus

$$R^+(s,a) = R(s,a) + B\left(N(s,a)\right)$$

$\uparrow N(s)$

$N\uparrow \Rightarrow B\downarrow$

Large/ Continuous $S$

$\phi_\theta(s)$  "psedo-count"

Small MDP

$$P(s) = \frac{N(s)}{n}$$

after seeing a new state $s$

$$\rightarrow P'(s) = \frac{N(s)+1}{n+1}$$

model that fits these dynamics

$\rightarrow$ fit $p_\theta(s)$  to all seen states $D$

fit $p_{\theta'}(s)$  to $D \cup \{s\}$

$p_\theta(s), p_\theta'(s) \longrightarrow \hat{N}(s)$

$R^+_{(s,a)} = R(s,a) + B\left(\hat{N}(s)\right)$

$$\hat{N}(s) = \hat{n}\, p_\theta(s) \qquad \hat{n} = \frac{1-p_{\theta'}(s)}{p_{\theta'}(s) - p_\theta(s)} \, p_\theta(s)$$

What B?

UCB

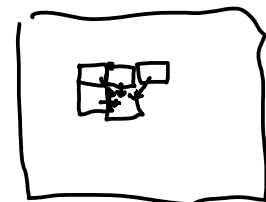$\rightarrow \sqrt{\frac{1}{N(s)}}$   $\leftarrow$ Bellomare   "Unifying Count-Based Exploration"

$\frac{1}{N(s)}$

How to model $p_\theta'(s)$?   CTS
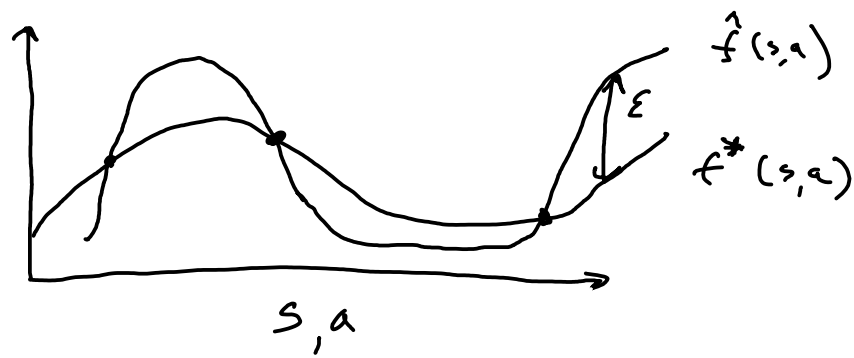
← "hash"
- #Exploration.

compress $s$ into k-bit code w/ $\phi(s)$, then count
$$N(\phi(s))$$

---

- Use a classifier "EX2"

$$p_\theta(s) = \frac{1 - D_s(s)}{D_s(s)} \leftarrow \text{probability of positive}$$

---



$$S, a$$

Use $\varepsilon(s,a) = \| \hat{f}(s,a) - f^*(s,a) \|$ as bonus

What is $f^*$
- $f^*(s,a) = s'$      "Curiosity"

- $f^*(s,a) = f_\phi(s,a)$    where $\phi$ is random
                             random Neural Network

RND - Random Network Distillation

---

# Thompson Sampling Style

$p(Q)$

1. Sample $\hat{Q}$ from $p(Q)$
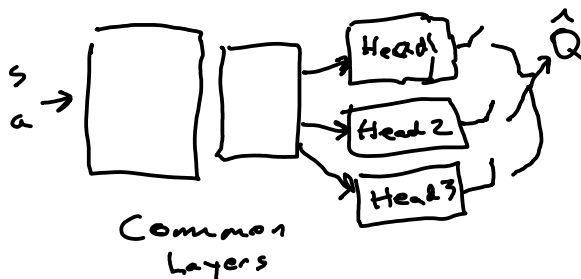2. act according to $\hat{Q}$ for 1 episode
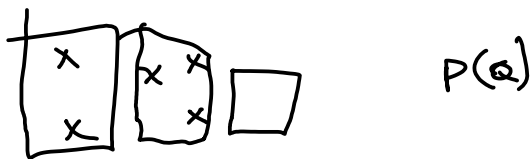
Works with off-policy

## How to maintain $p(\theta)$

- Bootstrapping
    - Resample $D$ $N$ times to get $D_1, \ldots D_N$
    - train $f_{\theta_i}$ on $D_i$
    - Sample from $p(\theta)$ by sampling $i \in 1..N$
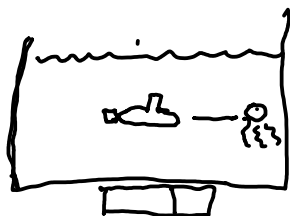      using $f_{\theta_i}$



Common Layers

Osband et al. "Deep exploration via Bootstrapped DQN"

- Dropout



$p(Q)$

+ Don't have to modify Reward
- Doesn't work as well as optimistic

# Information Gain

$$IG(z, y \mid a)$$

about what?

$R(s,a)$ ?  bad for sparse

$P(s)$  state density

$P(s' \mid s, a)$  dynamics

Generally Intractible $\longrightarrow$ Approximations

## Approximations

- prediction gain  $\log p_{\theta'}(s) - \log p_\theta(s)$

  justification  RND

- Variational Inference

  $q_\theta(s) \approx p(s \mid h)$  $\boxed{VIME} \leftarrow$

  IG is like  $D_{KL}\big(p(z \mid y) \,\|\, p(z)\big)$

  $p_\theta(s_{t+1} \mid s_t, a_t)$  $z = \theta$

  $q(\theta \mid \phi) \approx p(\theta \mid h)$  $y = (s_t, a_t, s_{t+1}) \leftarrow$

  $D_{KL}\big(p(\theta \mid h, s_t, a_t, s_{t+1}) \,\|\, p(\theta \mid h)\big)$

  specifically optimize variational lower bound

  $D_{KL}\big(q(\theta \mid \phi) \,\|\, p(h \mid \theta)\, p(\theta)\big)$

  $\uparrow$ product of independent

  Gaussian param.
  distributions



Every step, update $\phi$ to $\phi'$

Use $D_{KL}\big(q(\theta \mid \phi') \,\|\, q(\theta \mid \phi)\big)$ as bonus

Review:

Optimistic: RND ← Best

Thompson: Bootstrapping with many Q networks

theoretically justified →

IG : VIME

---

# Using Expert Demonstration

### Imitation Learning

+ Simple, stable, supervised

- Demos
- Unseen "Distributional shift"
- as good as expert

### RL

+ Exceed human perf

- Reward
- Exploration
- Unstable

Simplest ⤵

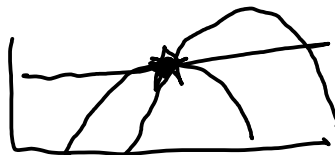- pre-train + fine tune with RL

  Flaws: - could bias
  - could forget

- Use off-policy RL with data from expert

  - Policy Gradient with importance sampling

    $\boxed{\text{Guided Policy Search}}$

➤ Q-learning: Drop into Replay Buffer

$$Q(s,a) \leftarrow r(s,a) + E_{a' \sim \pi_{new}}[Q(s',a')]$$

# Hybrid objective

imitation
$$\sum_{(s,a) \in Demo} \log \pi_\theta (a | s)$$

RL
$$\mathbb{E}_{\pi_\theta} [R(s,a)]$$

Hybrid
$$\mathbb{E}_{\pi_\theta} [R(s,a)] + \lambda \sum_{(s,a) \in Demo} \log \pi_\theta (a|s)$$

Flaws: − choose weight
− domain − specific