

Computer Engineering Department

Exploration in RL

Mohammad Hossein Rohban, Ph.D.

Spring 2024

Courtesy: Most of slides are adopted from CS 296, UC Berkeley.

What's the problem?

this is easy (mostly)



Why?

this is impossible



Montezuma's revenge

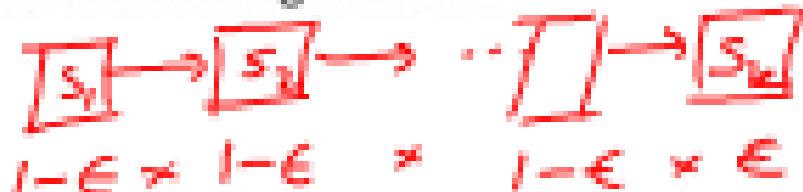


- Getting key = reward
- Opening door = reward
- Getting killed by skull = nothing (is it good? bad?)
- Finishing the game only weakly correlates with rewarding events
- We know what to do because we understand what these sprites mean!

Why exploration can be difficult?

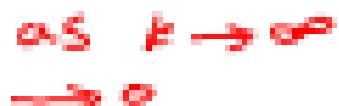
Temporally extended tasks like Montezuma's revenge become increasingly difficult based on

- How extended the task is
- How little you know about the rules



- Let's assume a complex task

- Consisting of multiple sub-task, each is a prerequisite for the next sub-task.
- Each should be solved in a sequence to get a high reward
- Epsilon greedy does not obviously help:
- Suppose you mastered up to the k^{th} sub-task.
- You have to exploit up to the k^{th} task and then explore onwards.
- Now the chance to only explore in the sub-task $(k+1)$ is $(1-\epsilon)^{0(k)} \epsilon^{0(1)}$.
- For $\epsilon_{\text{pt}} = 0.1$, this is $\approx 6\%$. For $\epsilon_{\text{pt}} = 0.5$, this is $\approx 2\%$.



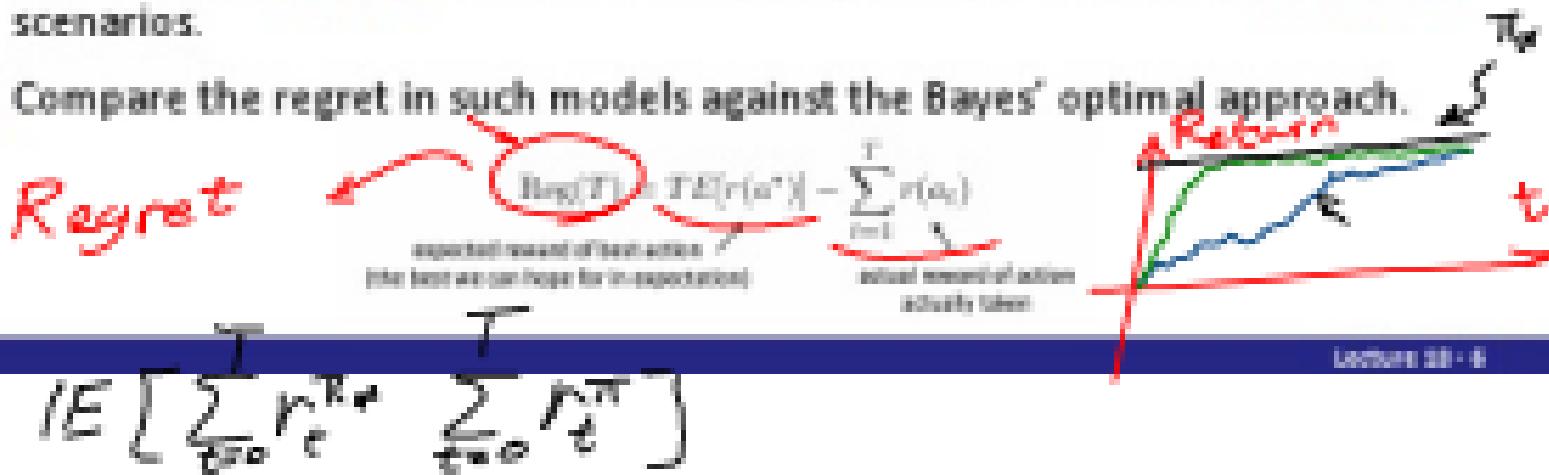
\rightarrow

Exploration and exploitation

- Two potential definitions of exploration problem:
 - How can an agent **discover** high-reward strategies that require a temporally extended sequence of complex behaviors that, individually, are not rewarding?
 - How can an agent decide whether to **attempt new behaviors** (to discover ones with higher reward) or continue to do the best thing it knows so far?

Optimal Exploration?

- Bayesian model of the environment. (POMDP with belief state)
- Optimize the expected reward under all uncertainties.
- Requires knowledge of state dynamic distribution class, the prior, and maintaining the belief state.
- Here we seek simpler solutions which could be extended to more complex scenarios.
- Compare the regret in such models against the Bayes' optimal approach.



Multi-armed Bandit

Bandits

assume $r(a_i) \sim p_{\theta_i}(r_i)$

e.g., $p(r_1 = 1) = \theta_1$ and $p(r_1 = 0) = 1 - \theta_1$

$\theta_i \sim p(\theta)$, but otherwise unknown

this defines a POMDP with $s = [\theta_1, \dots, \theta_n]$

belief state is $p(\theta_1, \dots, \theta_n)$



- solving the POMDP yields the optimal exploration strategy
- but that's overkill: belief state is huge!
- we can do very well with much simpler strategies

how do we measure goodness of exploration algorithm?

regret: difference from optimal policy at time step T:

$$\text{Reg}(T) = TE[r(a^*)] - \sum_{t=1}^T r(a_t)$$

expected reward of best action ✓
 (the best we can hope for in expectation)
 actual reward of action
 actually taken

Optimistic exploration

UCB

keep track of average reward $\hat{\mu}_a$ for each action a

exploitation: pick $a = \arg \max \hat{\mu}_a$

optimistic estimate: $a = \arg \max \hat{\mu}_a + \sigma_a$

some sort of variance estimate

$$\frac{\text{std}(t_j - \bar{t}_{n(a)})}{\sqrt{n(a)}}$$

$$\dots$$

$$k$$

intuition: try each arm until you are sure it's not great

example (Auer et al. Finite-time analysis of the multiarmed bandit problem):

$$a = \arg \max \hat{\mu}_a + \sqrt{\frac{2 \ln T}{N(a)}}$$

number of times we
picked this action

$$\hat{\mu}_a = \frac{\sum_{i=1}^t r_i}{n(a)}$$

$\text{Reg}(T)$ is $O(\log T)$, probably as good as any algorithm

$$\hat{\mu} \quad 1000 \quad vs. \quad 5$$

50 49

$$\hat{P}(a) = \underbrace{\frac{1}{T} \sum_{t=1}^T R_t}_{P(a)}$$

$$\hat{\text{Regret}} = \frac{1}{T} \sum_{t=1}^T Q(\hat{a}^*) - Q(\hat{a}_t)$$

(WNB) \overbrace{t} Subtracted

$$Q(a_i) = \mathbb{E}(r(a_i)) \quad a^* = \operatorname{arg\,max}_a Q(a)$$

$$\text{if } x \geq \epsilon : P(X \geq \epsilon) \leq \frac{E(X)}{\epsilon}$$

Markov Inequality

proof:

$$\begin{aligned} &= \frac{\int_{\epsilon}^{\infty} \epsilon f_X(x) dx}{\epsilon} \\ &\leq \frac{\int_{\epsilon}^{\infty} x f_X(x) dx}{\epsilon} \leq \frac{\int_0^{\infty} x f_X(x) dx}{\epsilon} \\ &= \frac{E(X)}{\epsilon} \end{aligned}$$

$$X \rightarrow e^{tX}$$

$$t > 0 \quad Y$$

$$\forall \epsilon. P(X \geq \epsilon) = P(e^{tX} \geq e^{t\epsilon})$$
$$\leq \underbrace{E(Y)/e^t}_{= \frac{1}{e^t}} = \underbrace{\frac{E(e^{tX})}{e^{t\epsilon}}}_{= \frac{E(e^{tX})}{e^{t\epsilon}}}$$

$$\Rightarrow P(X \geq \epsilon) \leq \inf_{t > 0} \frac{E(e^{tX})}{e^{t\epsilon}}$$

Chernoff Bound

Suppose $E(X) = \alpha$, $X \in [a, b]$

$$\underline{E(e^{tX})} \rightarrow X = \alpha b + (1-\alpha)a$$
$$0 \leq \alpha \leq 1$$

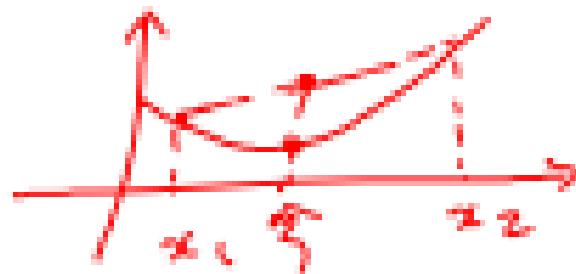
$$\alpha(b-a) = X - a$$

$$\alpha = \frac{X-a}{b-a}$$

$$E(\alpha) = \frac{-a}{b-a}$$

for any convex function f

$$f(\alpha x_1 + (1-\alpha)x_2) \leq \alpha f(x_1) + (1-\alpha)f(x_2)$$



let $\text{for } e^{tx} \quad x = \alpha b + (1-\alpha)a$

$$e^{tx} \leq \alpha e^{tb} + (1-\alpha)e^{ta}$$

$$\rightarrow g''(u) = \frac{(-\gamma e^u)(-\gamma e^u + u) + \gamma e^u(-\gamma e^u)}{(-\gamma e^u + 1 + \gamma)^2} = P(1/\gamma)$$

$$\rightarrow E(e^{tX}) \leq E(\alpha e^{tb} + (1-\alpha)e^{ta})$$

$$= \frac{-\alpha}{b-a} e^{tb} + \frac{b}{b-a} e^{ta}$$

$$= e^{ta} \underbrace{\left(\frac{-\alpha}{b-a} e^{t(b-a)} + \frac{b}{b-a} \right)}$$

$$u := t(b-a)$$

$$e^{g(u)}, \quad \gamma := \frac{a}{b-a}$$

$$g(u) = \gamma u + \log(-\gamma e^u + 1 + \gamma)$$

$$g'(u) = \gamma + \frac{-\gamma e^u}{-\gamma e^u + 1 + \gamma} \rightarrow g'(u) = 0$$

$$g(u) = 0 + \frac{g''(0)u}{1!} + \frac{g'''(\xi)u^2}{2!}$$

$$\xi \in [0, u]$$

$$\underbrace{g'''(\xi)}_{P(1-P)} \leq \frac{1}{4}$$

$$g(u) \leq \frac{1}{8}u^2$$

$$\rightarrow \underline{E(e^{tX})} \leq e^{g(u)} \leq e^{\frac{1}{8}u^2} e^{t^2(b-a)^2} = \underline{e}$$

$$\begin{aligned}
 P(X > \epsilon) &\leq \inf_{t \geq 0} \frac{\mathbb{E}(e^{tX})}{e^{t\epsilon}} \\
 &\leq \inf_{t \geq 0} \frac{e^{\frac{1}{2}t^2(b-a)^2}}{e^{t\epsilon}}
 \end{aligned}$$

$$\frac{1}{4}t(b-a)^2 - \epsilon = 0 \Rightarrow t = \frac{4\epsilon}{(b-a)^2}$$

$$\geq e^{-\frac{2\epsilon^2}{(b-a)^2}}$$

Hoeffding
Bound

$$P\left(\frac{x_1 + \dots + x_n}{n} - E(x) \geq \epsilon\right)$$

$$\inf_t E(e^{tZ}) = e^{tE}$$

$$C = (b-a)^2$$

$$n \cdot h \cdot P \cdot 1 - \frac{1}{2} t^2$$

$$\hat{\mu} \leq E(x) + \sqrt{\frac{2 \log \frac{1}{\delta}}{n}}$$

$$\epsilon = \sqrt{\frac{2 \log \frac{1}{\delta}}{n}}$$

$$U(a_t) := \hat{Q}(a_t) + \sqrt{\frac{2 \log t}{n(a_t)}}$$

$$\begin{aligned} \text{Regret} &= \sum_{t=1}^T Q(a^*) - Q(a_t) \\ &= \sum_{t=1}^T \overbrace{Q(a^*) - U(a_t)}^{< 0 \text{ w.h.p}} + \overbrace{U(a_t) - \hat{Q}(a_t)}^{\geq 0} \end{aligned}$$

$$a_t = \underset{a}{\operatorname{argmax}} U(a) \quad a^* = \underset{a}{\operatorname{argmax}} Q(a)$$

Assumption $\forall t, \forall a, Q(a) - \hat{Q}(a) \leq \sqrt{\frac{2\log \frac{1}{\delta}}{nC\alpha}}$

(1) $\hat{a}_t = a^*$ $Q(a^*) \leq U(a^*)$ w.h.p

$$Q(a^*) - \hat{Q}(a^*) \leq \epsilon$$

$$\Rightarrow Q(a^*) \leq \underbrace{\hat{Q}(a^*) + \epsilon}_{U(a^*)}$$

(2) $a_t \neq a^*$ $U(a_t) \geq U(a^*) \geq Q(a^*)$

w.h.p

$$\Rightarrow U(a_t) \geq Q(a^*)$$

$$\rightarrow \text{Regret} \leq \sum_{t=1}^T \overline{Q(a_t) - Q_{\text{last}}}$$

$$\hat{Q}(a_t) + \sqrt{\frac{2 \log \frac{T}{\delta}}{N(a_t)}}$$

$$\hat{Q}(a_t) - Q(a_t) \leq \sqrt{\frac{2 \log \frac{T}{\delta}}{N(a_t)}} \quad \text{with}$$

prob. of at least $1 - \delta/4$

$$\rightarrow \text{Regret} \leq \sum_{t=1}^T 2 \sqrt{\frac{2 \log \frac{T}{\delta}}{N(a_t)}}$$

w.p. of at least $P(\text{Assumptions})$
 $(= 2^{-m\delta})$

$P(\text{Assumptions : } \forall t \forall a \text{ Confid.}$
 $\text{Bounds are satisfied})$

$= 1 - P(\exists t \exists a \text{ Conf. Bound}$
 $\text{is violated})$

$$= 1 - P\left(\bigcup_{t=1}^T \bigcup_{a \in \mathcal{A}} Q(a_i^{(t)}) - \hat{Q}(a_i^{(t)}) \geq \sqrt{\frac{2 \log \delta / \epsilon}{n_t(a_i)}}\right)$$

Union Bound $P(A \cup B) \leq P(A) + P(B)$

$$P(\text{Assumption}) \geq 1 - \sum_{t,a} P(a - \hat{a} \geq t)$$

$$\sum_{t=1}^T \gamma_t^2 \leq \sum_{t=1}^{\infty} \gamma_t^2 = \pi^2/6 < 2$$

$$\Rightarrow P(\text{Assumption}) \geq \frac{1 - 2m\delta}{m}$$

$$m = |A|$$

$$\begin{aligned}
 \text{Regret} &\leq \sum_{t=1}^T 2 \sqrt{\frac{2 \log t}{n(a_t)}} \\
 &\leq 2 \sqrt{2 \log T} \cdot \sum_{t=1}^T \sqrt{\frac{1}{n(a_t)}} \\
 &\leq \frac{2 \sqrt{2 \log T}}{\sqrt{\sum_{a \in \mathcal{A}} n(a) \sum_{i=1}^{|A|} \sqrt{\frac{1}{i}}}}
 \end{aligned}$$

maximized when $n(a)$
 are all equal.

$$\leq 2 \sqrt{3T\delta} \sum_{a=a_1}^{\infty} \sum_{i=1}^{n_a} \sqrt{\frac{1}{j_i}}$$

$$\sum_{i=1}^m \frac{1}{j_i} \leq 2\sqrt{m}$$

$$2\sqrt{T\delta}$$

$$\leq 4 \underbrace{\sqrt{2m \log T/\delta}}_{\text{simp}}$$

→

Sublinear

Probability matching/posterior sampling

assume $r(a_i) \sim p_{\theta_i}(r_i)$

this defines a POMDP with $s = [\theta_1, \dots, \theta_n]$

belief state is $p(\theta_1, \dots, \theta_n)$

this is a model of our bandit

Idea: sample $\theta_1, \dots, \theta_n \sim p(\theta_1, \dots, \theta_n)$

pretend the model $\theta_1, \dots, \theta_n$ is correct

take the optimal action

update the model

$$(1 - \theta_{a_{t+1}}) \theta_{a_{t+1}}$$



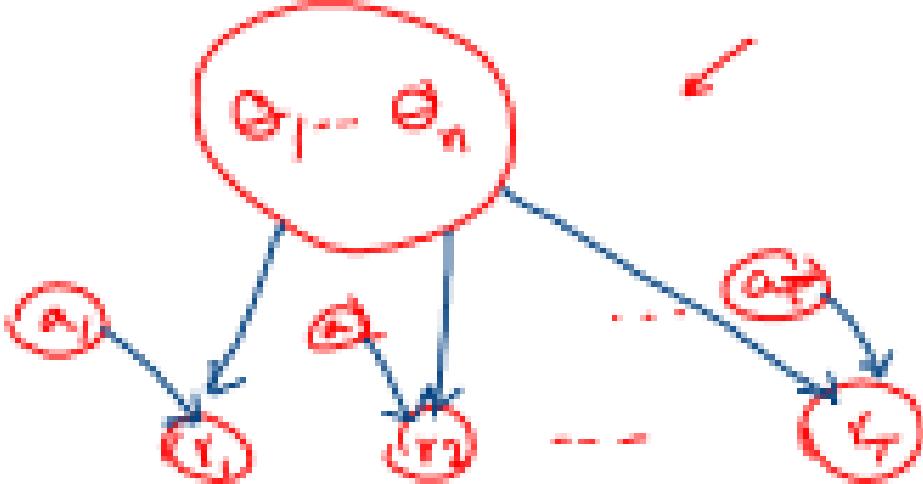
- This is called posterior sampling or Thompson sampling
- Harder to analyze theoretically
- Can work very well empirically
- See: Chapelle & Li, "An Empirical Evaluation of Thompson Sampling."

Belief

$$P(\theta_1, \dots, \theta_n | r_1, \dots, r_t, r_{t+1})$$

$$= P(r_{t+1} | \theta_1, \dots, \theta_n, r_1, \dots, r_t).$$

$$\propto P(\theta_1, \dots, \theta_n | r_1, \dots, r_t) / K$$



Information gain

VIME

Bayesian experimental design:

say we want to determine some latent variable z

(e.g., z might be the optimal action, or its value)

which action do we take?

let $\mathcal{H}(\hat{p}(z))$ be the current entropy of our z estimate

let $\mathcal{H}(\hat{p}(z)|y)$ be the entropy of our z estimate after observation y

(e.g., y might be $r(a)$)

the lower the entropy, the more precisely we know z

$$S_{\text{old}} \xrightarrow{\text{act } a} S_{\text{new}}$$

$$\text{IG}(z, y) = E_y[\mathcal{H}(\hat{p}(z)) - \mathcal{H}(\hat{p}(z)|y)]$$

with a

typically depends on action, so we have $\text{IG}(z, y|a)$

Information gain example

$$\text{IG}(z, y|a) = E_y[\mathbb{H}(\rho(z)) - \mathbb{H}(\rho(z)|y)|a]$$

how much we learn about z from action a , given current beliefs

Example bandit algorithm:

Russo & Van Roy "Learning to Optimize via Information-Directed Sampling"

$y = r_a$, $z = \theta_a$ (parameters of model $p(r_a)$)

$g(a) = \text{IG}(\theta_a, r_a|a)$ – information gain of a

$\Delta(a) = E[r(a^*)] - r(a)$ – expected suboptimality of a

choose a according to $\arg \min_a \frac{\Delta(a)^2}{g(a)}$

don't bother taking actions if you won't learn anything

$$\min_a \frac{-E(r(a))^2}{g(a)}$$

don't take actions that you're sure are suboptimal

General themes

UCB:

$$a = \arg \max \hat{\mu}_a + \sqrt{\frac{2 \ln T}{N(a)}}$$

Thompson sampling:

$$\theta_1, \dots, \theta_n \sim \beta(\theta_1, \dots, \theta_n)$$

$$a = \arg \max_a E_{\theta_a}[r(a)]$$

Info gain:

$$IG(z, y|a)$$

- Most exploration strategies require some kind of uncertainty estimation (even if it's naive)
- Usually assumes some value to new information
 - Assume unknown = good (optimism)
 - Assume sample = truth
 - Assume information gain = good

Optimistic exploration in RL

$$\text{UCB: } a = \arg \max \hat{\mu}_a + \sqrt{\frac{2 \ln T}{N(a)}}$$

"exploration bonus"

lots of functions work, as long as they decrease with $N(a)$

can we use this idea with MDPs?

count-based exploration: use $N(s, a)$ or $N(s)$ to add exploration bonus

$$\text{use } r^+(s, a) = r(s, a) + B(N(s))$$

$$B(x) = \sqrt{x} \quad \text{or} \quad \sqrt{\frac{1}{x}}$$

bonus that decreases with $N(s)$

use $r^+(s, a)$ instead of $r(s, a)$ with any model-free algorithm

+ simple addition to any RL algorithm

- need to tune bonus weight

The trouble with counts

use $r^+(s, a) = r(s, a) + \beta(N(s))$

But wait... what's a count?



Uh oh... we never see the same thing twice!

But some states are more similar than others

Fitting generative models

idea: fit a density model $p_\theta(s)$ or $p_\theta(s, a)$

$p_\theta(s)$ might be high even for a new s

if s is similar to previously seen states

can we use $p_\theta(s)$ to get a "pseudo-count"?

if we have small MDPs

the true probability is:

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

↑
probability/density ↓
count
total states visited

after we see s , we have:

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

Exploring with pseudo-counts

- fit model $p_\theta(s)$ to all states D seen so far
- take a step i and observe s_i
- fit new model $p_{\theta'}(s)$ to $D \cup s_i$
- use $p_\theta(s_i)$ and $p_{\theta'}(s_i)$ to estimate $\hat{N}(s_i)$
- set $r_i^+ = r_i + \delta(\hat{N}(s_i))$ — “pseudo-count”

$$p_\theta(s)$$

$$p_{\theta'}(s)$$

how to get $\hat{N}(s_i)$? use the equations

$$p_\theta(s_i) = \frac{\hat{N}(s_i)}{n}$$

$$p_{\theta'}(s_i) = \frac{\hat{N}(s_i) + 1}{n + 1}$$

two equations and two unknowns!

$$\hat{N}(s_i) = n p_\theta(s_i)$$

$$n = \frac{1 - p_\theta(s_i)}{p_{\theta'}(s_i) - p_\theta(s_i)}$$

What kind of generative modeling to use?



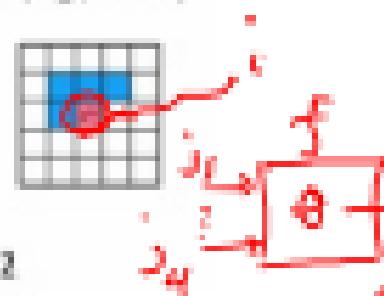
need to be able to output densities, but doesn't necessarily need to produce great samples



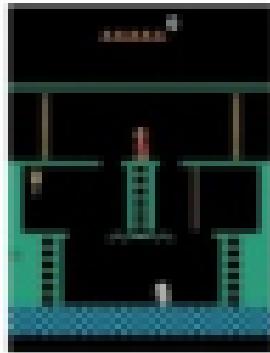
opposite considerations from many popular generative models in the literature (e.g., GANs)

Bellman et al. "CTG" model:
condition each pixel on its top-left neighborhood

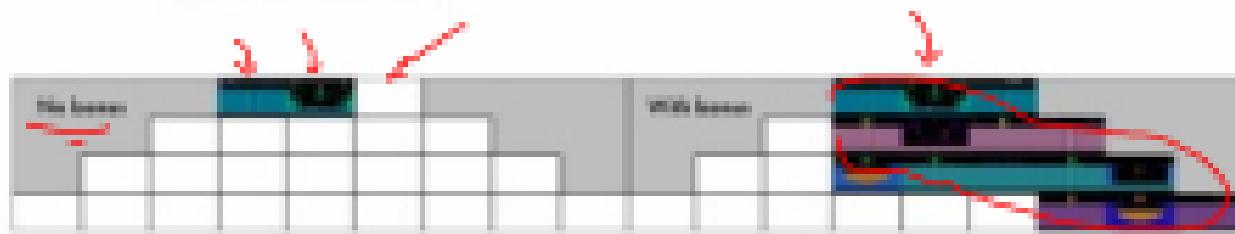
Other models: stochastic neural networks, compression length, etc.



Does it work?



VS



Bellman et al. "Identifying Count-Based Exploration..."

Counting with hashes

What if we still count states, but in a different space?

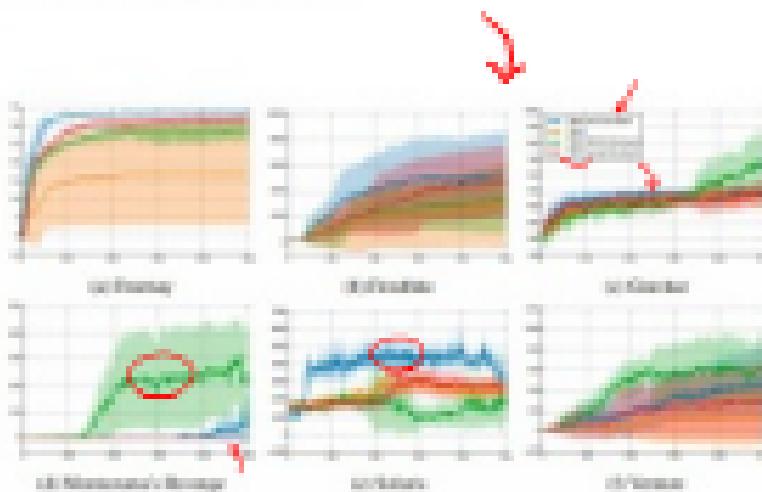
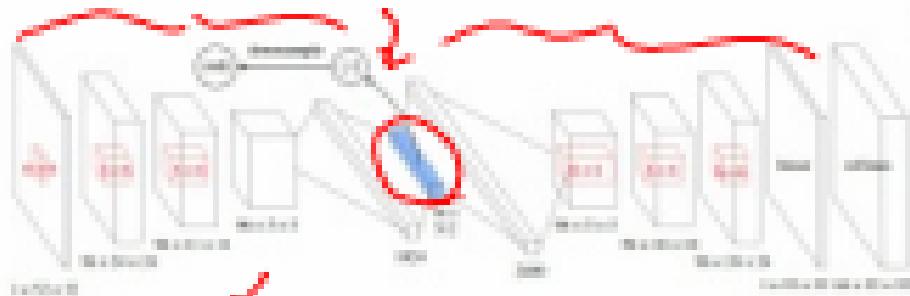
Ideal: compress s into a d -bit code via $\phi(s)$, then count $N(\phi(s))$

shorter codes = more hash collisions

similar states get the same hash? maybe



improve the odds by learning a compression:



Implicit density modeling with exemplar models

$p_\theta(s)$

need to be able to output densities, but doesn't necessarily need to produce great samples

Can we explicitly compare the new state to past states?

Intuition: the state is **novel** if it is **easy** to distinguish from all previous seen states by a classifier

for each observed state s , fit a classifier to classify that state against all past states D , use classifier error to obtain density

$$p_\theta(s) = \frac{1 - D_\theta(s)}{D_\theta(s)}$$

probability that classifier assigns that s is "positive"
positives: $\{s\}$
negatives: D

previously
seen
states

Implicit density modeling with exemplar models

$$\underline{P(j=1|s)}$$

hang on... aren't we just checking if $s = s^*$?

If $s \in \mathcal{D}$, then the optimal $D_{\theta}(s) \neq 1$

In fact: $D_{\theta}(s) = \frac{1}{1 + p(s)}$



$$= \frac{\hat{P}(s|j=1)\hat{P}(j=1)}{\hat{P}(s|j=1)\hat{P}(j=1) + \hat{P}(s|j=0)\hat{P}(j=0)}$$
$$\hat{p}_e(s) = \frac{1 - D_{\theta}(s)}{D_{\theta}(s)}$$

in reality, each state is unique, so we regularize the classifier

isn't one classifier per state a bit much?

train one amortized model: single network that takes in exemplar as input



Figure 1: Density Response

Fu et al. 2012: Exploration with Exemplar Models...

Heuristic estimation of counts via errors

$p_\theta(s)$

need to be able to output densities, but doesn't necessarily need to produce great samples

...and doesn't even need to output great densities

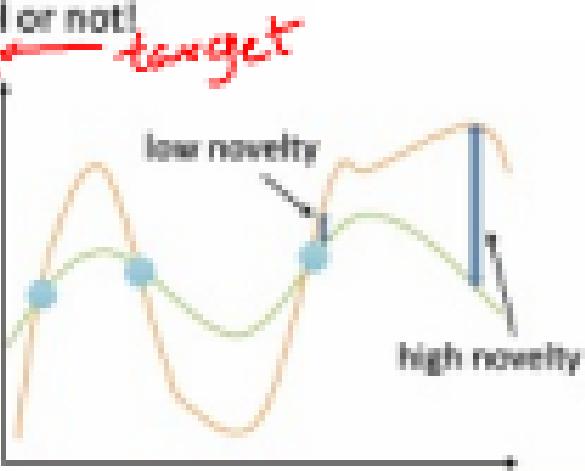
...just need to tell if state is novel or not!

let's say we have some target function $f^*(s, a)$

given our buffer $\mathcal{D} = \{(s_t, a_t)\}$, fit $f_\theta(s, a)$

use $\mathcal{L}(s, a) = \|f_\theta(s, a) - f^*(s, a)\|^2$ as bonus

$$f_\theta(s_t, a_t) = S_{t+1}$$



Heuristic estimation of counts via errors

what should we use for $f^*(s, a)$?

one common choice: set $f^*(s, a) = s'$ – i.e., next state prediction

even simpler: $f^*(s, a) = f_\phi(s, a)$, where ϕ is a random parameter vector

Posterior sampling in deep RL

Thompson sampling:

$$\theta_1, \dots, \theta_n \sim p(\theta_1, \dots, \theta_n)$$

$$\theta = \arg \max_{\theta} E_{\theta} [r(a)]$$

$$P(Q(s, a))$$

What do we sample?

How do we represent the distribution?

bandit setting: $p(\theta_1, \dots, \theta_n)$ is distribution over rewards

$$\theta_1, \dots, \theta_n \rightarrow Q_1, \dots, Q_n$$

MDP analog is the Q -function!



1. sample Q -function Q from $p(Q)$
2. act according to Q for one episode
3. update $p(Q)$

since Q-learning is off-policy, we don't care which Q -function was used to collect data

how can we represent a distribution over functions?

Bootstrap

given a dataset \mathcal{D} , resample with replacement N times to get $\mathcal{D}_1, \dots, \mathcal{D}_N$

train each model f_{θ_i} on \mathcal{D}_i ,

to sample from $p(\theta)$, sample $i \in [1, \dots, N]$ and use f_{θ_i} .

training N big neural nets is expensive, can we avoid it?



Osband et al. "Deep Exploration via Bootstrapped DQN"

Bootstrap

Q_1, \dots, Q_K are samples from the posterior list.

Algorithm 1 Bootstrapped DQN

- 1: Inputs: Value function networks Q with K outputs $\{Q_k\}_{k=1}^K$, Masking distribution M .
- 2: Let \mathcal{H} be a replay buffer storing experience for training.
- 3: for each episode do
- 4: Obtain initial state from environment
- 5: Pick a value function to act using $k = \text{Uniform}[1, \dots, K]$
- 6: for step $t = 1, \dots$ until end of episode do
- 7: Pick an action according to $a_t \in \arg\max_a Q_k(s_t, a)$
- 8: Receive state s_{t+1} and reward r_t from environment, having taking action a_t
- 9: Sample bootstrap mask $m_t \sim M$
- 10: Add $(s_t, a_t, r_t, s_{t+1}, m_t)$ to replay buffer \mathcal{H}
- 11: end for
- 12: end for

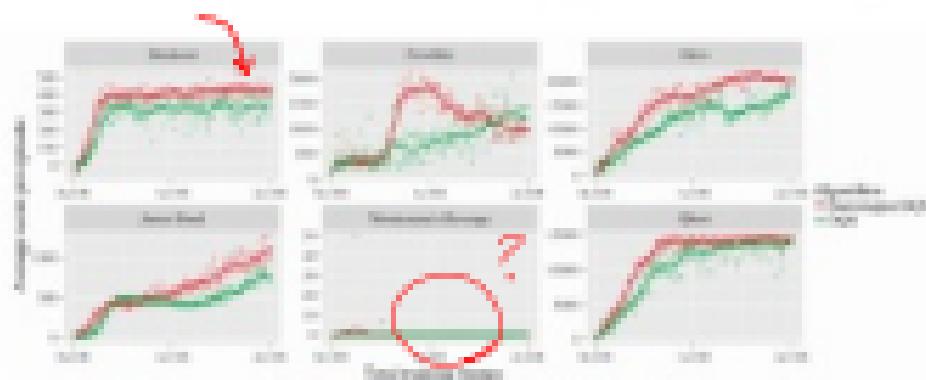
Osband et al. "Deep Exploration via Bootstrapped DQN"

$$h = (s_0, a_0, s_1, r_1, s_2, \dots, s_T)$$

Why does this work?

Exploring with random actions (e.g., epsilon-greedy): oscillate back and forth, might not go to a coherent or interesting place

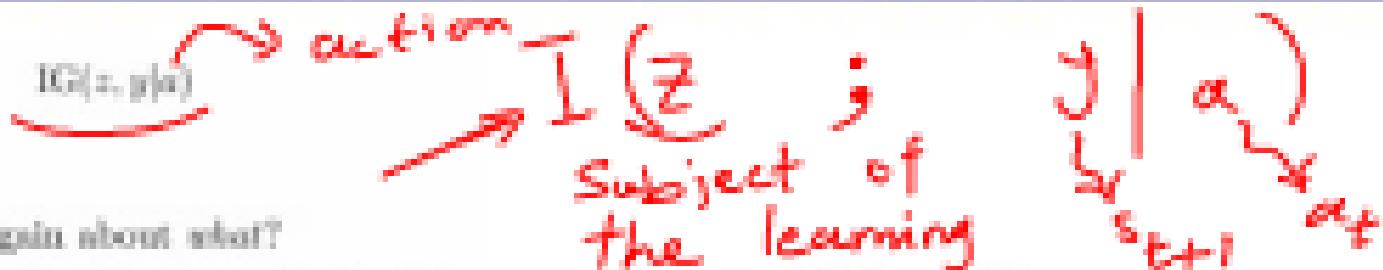
Exploring with random Q-functions: commit to a randomized but internally consistent strategy for an entire episode



- + no change to original reward function
- very good bonuses often do better

Reasoning about information gain (approximately)

Info gain: $IG(z, y|a)$



information gain about what?

information gain about reward $r(s, a)$?

state density $p(s)$? $P_{\pi^*}(s)$

information gain about dynamics $p(s'|s, a)$?

not very useful if reward is sparse

a bit strange, but somewhat makes sense!

good proxy for learning the MDP, though still heuristic

$$z = P_a(s'|s, a)$$

$$z = \theta$$

Generally intractable to use exactly, regardless of what is being estimated!

Reasoning about information gain (approximately)

A few approximations:

$$I(X; T) \approx D_{KL}(P(x|T) \| P(x))$$

prediction gain: $\log p_{\theta^*}(s) - \log p_{\theta}(s)$

(Schmidhuber '91, Bellmire '16)

Intuition: if density changed a lot, the state was novel

MAP?

variational inference

IG can be equivalently written as $D_{KL}(p(z|x) \| p(z))$

learn about transitions $p_{\theta}(s_{t+1}|s_t, a_t), z = \theta$



Intuition: a transition is more informative if it causes belief over θ to change

Idea: use variational inference to estimate $q(\theta|\phi) \rightarrow p(\theta|A) \rightarrow$
given new transition (s, a, s') , update ϕ to get ϕ'

Variational
Inference

$$s_t \xrightarrow{a_t} \boxed{\theta} \rightarrow s_{t+1}, \underline{\theta} = (\underline{\theta}_1, \dots, \underline{\theta}_d)$$

$$P(\underline{\theta} | h) = \prod_{i=1}^d q(\theta_i | \phi_i)$$

$$q(\cdot | (\mu_i, \sigma_i)) = \mathcal{N}(\mu_i, \sigma_i^2) \frac{e^{-\frac{1}{2}\phi_i^2}}{\sqrt{2\pi}}$$



Reasoning about information gain (approximately)

VIME implementations

IG can be equivalently written as $D_{KL}(p(\theta|h, s_t, a_t, s_{t+1}) \| p(\theta|h))$

model parameters θ : $p(\theta|s_{0:t}, a_{0:t})$

newly observed transition

history of all prior transitions



$$q(\theta|\phi) \approx p(\theta|h)$$

specifically, optimize variational lower bound $D_{KL}(q(\theta|\phi) \| p(h|\theta)p(\theta))$

represent $q(\theta|\phi)$ as product of independent Gaussian parameter distributions

with mean ϕ (see Blundell et al. "Weight uncertainty in neural networks")

given new transition (s, a, s') , update ϕ to get ϕ'

i.e., update the network weight means and variances

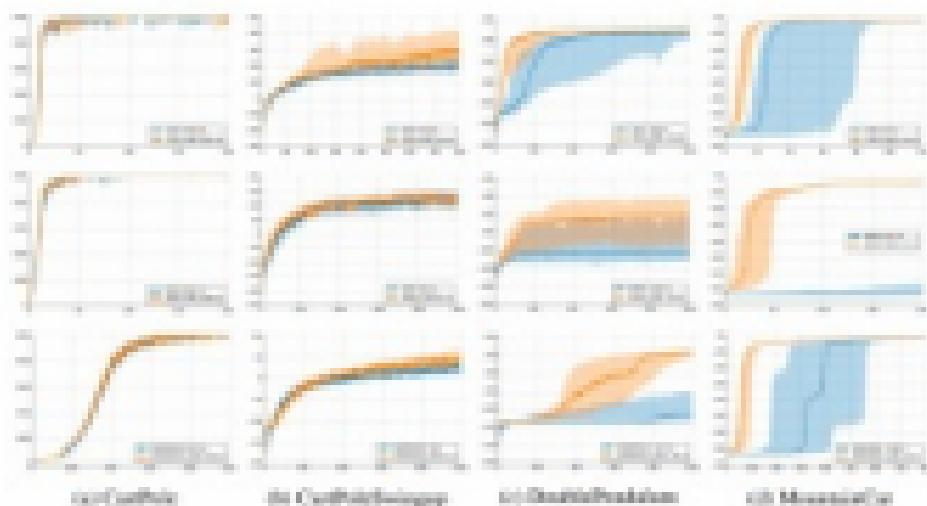
use $D_{KL}(q(\theta|\phi') \| q(\theta|\phi))$ as approximate bonus

$$p(\theta|\mathcal{D}) = \prod_i p(\theta_i|\mathcal{D})$$

$$p(\theta_i|\mathcal{D}) = \mathcal{N}(\mu_i, \sigma_i)$$

$$\frac{1}{\phi}$$

Reasoning about information gain (approximately)



Approximate IG:

- + appealing mathematical formalism
- models are more complex, generally harder to use effectively