

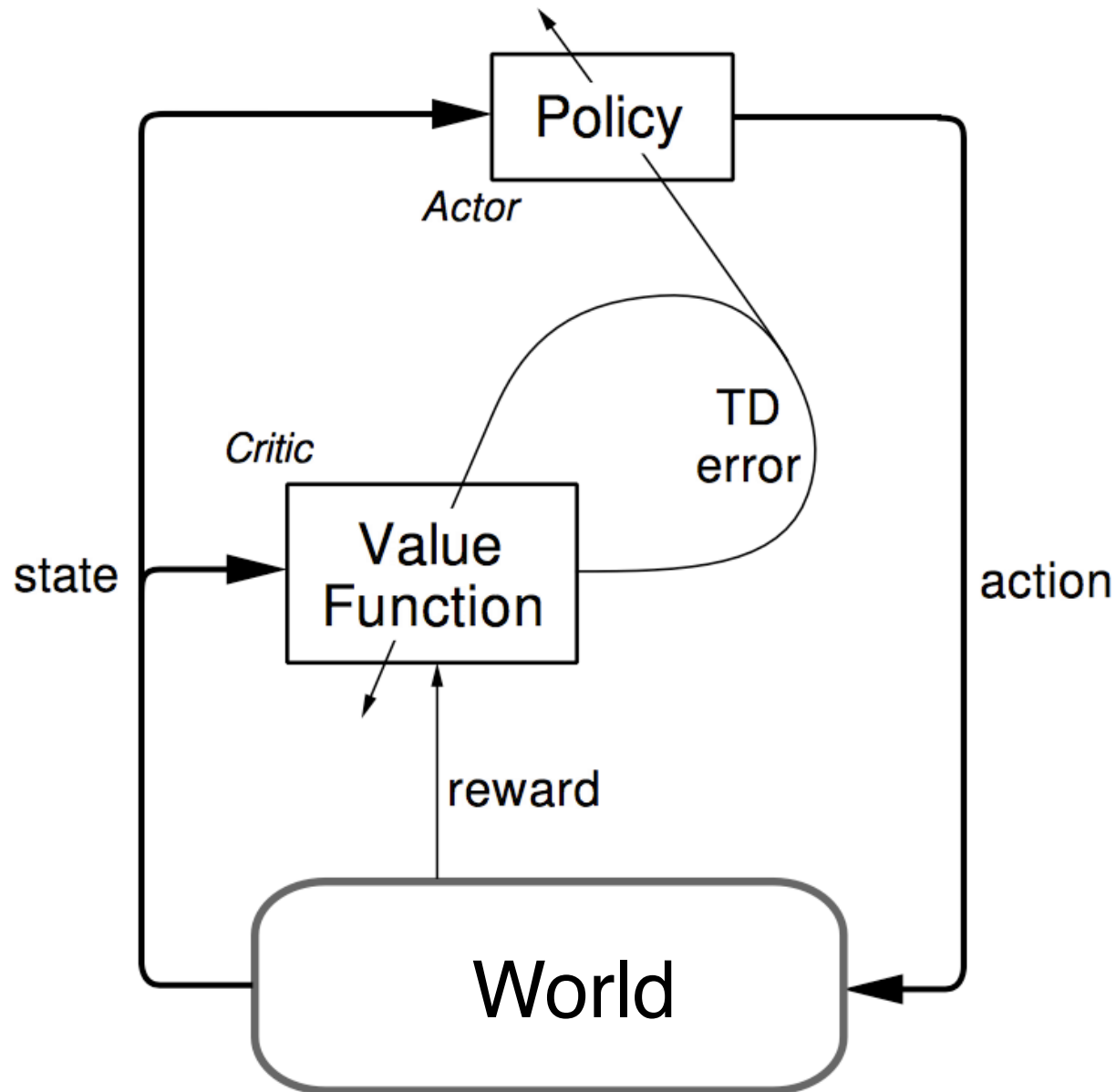
Policy Approximation

- Policy = a function from state to action
 - How does the agent select actions?
 - In such a way that it can be affected by learning?
 - In such a way as to assure exploration?
- Approximation: there are too many states and/or actions to represent all policies
 - To handle large/continuous action spaces

What is learned and stored?

1. *Action-value methods*: learn the value of each action; pick the max (usually)
2. *Policy-gradient methods*: learn the parameters \mathbf{u} of a stochastic *policy*, update by $\nabla_{\mathbf{u}} \text{Performance}$
 - including *actor-critic methods*, which learn both value and policy parameters
3. *Dynamic Policy Programming*
4. *Drift-diffusion models* (Psychology)

Actor-critic architecture



Action-value methods

- The *value of an action in a state given a policy* is the expected future reward starting from the state taking that first action, then following the policy thereafter

$$q_{\pi}(s, a) = \mathbb{E} \left[\sum_{t=1}^{\infty} \gamma^{t-1} R_t \mid S_0 = s, A_0 = a \right]$$

- **Policy:** pick the max most of the time

$$A_t = \arg \max_a \hat{Q}_t(S_t, a)$$

but sometimes pick at random (ϵ -greedy)

We should never discount
when approximating policies!



γ is ok if there is a
start state/distribution

Average reward setting

- All rewards are compared to the average reward

$$q_{\pi}(s, a) = \mathbb{E} \left[\sum_{t=1}^{\infty} R_t - \bar{r}(\pi) \mid S_0 = s, A_0 = a \right]$$

- where

$$\bar{r}(\pi) = \lim_{t \rightarrow \infty} \frac{1}{t} \mathbb{E} [R_1 + R_2 + \cdots + R_t \mid A_{0:t-1} \sim \pi]$$

- and we learn an approximation

$$\bar{r}_t \approx \bar{r}(\pi_t)$$

Why approximate policies rather than values?

- In many problems, the policy is simpler to approximate than the value function
- In many problems, the optimal policy is stochastic
 - e.g., bluffing, POMDPs
- To enable smoother change in policies
- To avoid a search on every step (the max)
- To better relate to biology

Policy-gradient methods

- The policy itself is learned and stored
 - the policy is parameterized by $\mathbf{u} \in \mathbb{R}^n$
 - we learn and store \mathbf{u}

$$\mathbb{P}r [A_t = a] = \pi_{\mathbf{u}_t}(a|S_t)$$

- \mathbf{u} is updated by approximate gradient ascent

$$\mathbf{u}_{t+1} = \mathbf{u}_t + \alpha \widehat{\nabla_{\mathbf{u}} \bar{r}(\pi_{\mathbf{u}})}$$

eg, linear-exponential policies (discrete actions)

- The “preference” for action a in state s is linear in \mathbf{u}

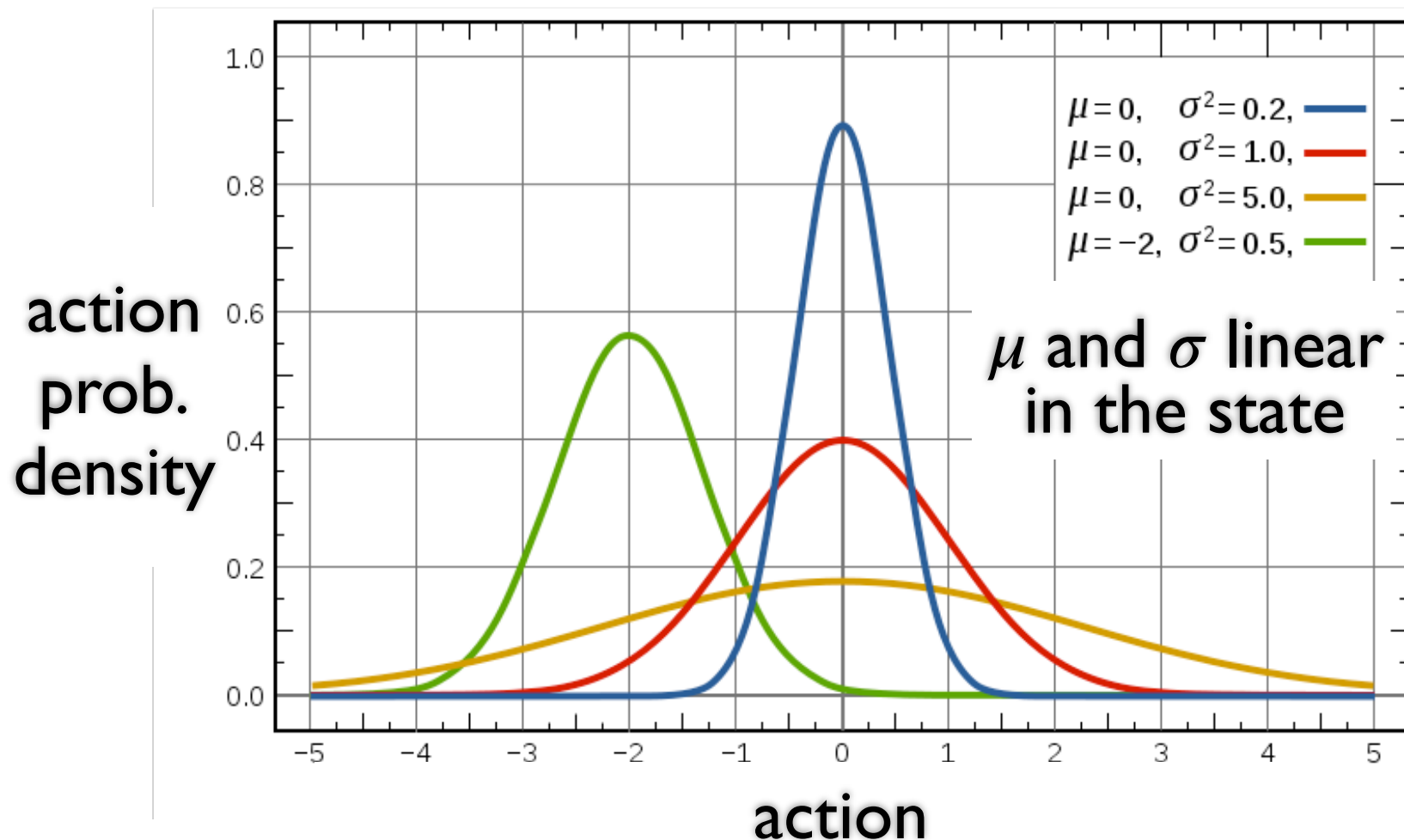
$$\mathbf{u}^\top \mathbf{x}_{sa} \equiv \sum_i u(i) \mathbf{x}_{sa}(i)$$

feature vector $\in \mathbb{R}^n$

- The probability of action a in state s is exponential in its preference

$$\pi_{\mathbf{u}}(a|s) = \frac{e^{\mathbf{u}^\top \mathbf{x}_{sa}}}{\sum_b e^{\mathbf{u}^\top \mathbf{x}_{sb}}}$$

eg, linear-gaussian policies (continuous actions)



eg, linear-gaussian policies (continuous actions)

- The mean and std. dev. for the action taken in state s are linear and linear-exponential in $\mathbf{u}_\mu, \mathbf{u}_\sigma$

$$\mu(s) = \mathbf{u}_\mu^\top \phi_s \quad \sigma(s) = e^{\mathbf{u}_\sigma^\top \phi_s}$$

- The probability density function for the action taken in state s is gaussian

$$\pi_{\mathbf{u}}(a|s) = \frac{1}{\sigma(s)\sqrt{2\pi}} e^{-\frac{(a-\mu(s))^2}{2\sigma(s)^2}}$$

Policy Approximation

Average Reward per step objective

No discounting over the on-policy distribution ^{& proof}

→ The policy gradient theorem and its proof

→ Gradient Ascent (stochastic)

Approximating the gradient

Eligibility ^{f_{ns}} ~~test~~ for discrete & continuous actions

A final algorithm

$$V_{\pi}^{\delta}(s) \equiv \sum_{t=0}^{\infty} E_{\pi} [\gamma^t R_{t+1} | S_0 = s]$$

$$= \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \delta V_{\pi}^{\delta}(s')]$$

$$d_{\pi}(s) \equiv \lim_{t \rightarrow \infty} P_r(S_t = s | A_{0:t} \sim \pi)$$

exists &
assume does not
depend on S_0
ergodic

$$d_{\pi}(s') = \sum_s d_{\pi}(s) \sum_a \pi(a|s) p(s' | s, a)$$

$$\bar{r}(\pi) \equiv \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T R_t$$

assume exists and
does not depend on S_0

$$= \sum_s d_{\pi}(s) \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) r$$

$$\tilde{V}_{\pi}(s) \equiv \sum_{t=0}^{\infty} (E_{\pi} [\gamma^t R_{t+1} | S_0 = s] - \bar{r}(\pi))$$

$$= \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r - \bar{r}(\pi) + \tilde{V}_{\pi}(s')]$$

$$\tilde{Q}_{\pi}(s, a) \equiv \text{similar}$$

Is $J(\pi) \doteq \sum_s d_\pi(s) V_\pi^\delta(s)$ a good objective?

Or is it the same as $\bar{r}(\pi)$?

$$\begin{aligned} J(\pi) &\doteq \sum_s d_\pi(s) V_\pi^\delta(s) \\ &= \sum_s d_\pi(s) \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) [r + \gamma V_\pi^\delta(s')] \\ &= \bar{r}(\pi) + \gamma \sum_s d_\pi(s) \sum_a \pi(a|s) \sum_{s'} p(s'|s,a) V_\pi^\delta(s') \\ &= \bar{r}(\pi) + \gamma \sum_{s'} V_\pi^\delta(s') \sum_s d_\pi(s) \sum_a \pi(a|s) p(s'|s,a) \\ &= \bar{r}(\pi) + \gamma \sum_{s'} V_\pi^\delta(s') d_\pi(s') \\ &= \bar{r}(\pi) + \gamma J(\pi) \\ &= \bar{r}(\pi) + \gamma \bar{r}(\pi) + \gamma^2 \bar{r}(\pi) + \gamma^3 \bar{r}(\pi) + \dots \\ &= \frac{\bar{r}(\pi)}{1-\gamma} \end{aligned}$$

i.e., a scaled $\bar{r}(\pi)$, thus no effect of γ
on the ordering of any policies

~~St~~

Gradient Ascent

Let $u \in \mathbb{R}^n$ be the policy parameter

$$u_{t+1} = u_t + \alpha \nabla_u r(\pi_u)$$

Stochastic Gradient Ascent

$$u_{t+1} = u_t + \alpha \nabla_u r(\pi_u)$$

where $E_{\pi}[\nabla_u r(\pi_u)] = \nabla_u r(\pi_u)$

Policy Gradient Theorem

$$\nabla_u r(\pi_u) = \sum_s d_{\pi}(s) \sum_a \tilde{q}_{\pi}(s,a) \nabla_u \pi_u(s,a)$$

$$= E_{\pi} \left[\tilde{q}_{\pi}(s_t, A_t) \frac{\nabla_u \pi_u(s_t, A_t)}{\pi_u(s_t, A_t)} \right]$$

$$u_{t+1} = u_t + \alpha \left[Q(s_t, A_t) - b(s_t) \right] \frac{\nabla_u \pi_u(s_t, A_t)}{\pi_u(s_t, A_t)}$$

Proof of P.G.T.

$$\begin{aligned}
 \nabla_n \tilde{V}_\pi(s) &= \nabla_n \sum_a \pi_n(a|s) \tilde{Q}_\pi(s,a) \\
 &= \sum_a \left[\nabla \pi_n(a|s) \tilde{Q}_\pi(s,a) + \pi_n(a|s) \nabla \tilde{Q}_\pi(s,a) \right] \\
 &= \sum_a \left[\nabla \pi_n(a|s) \tilde{Q}_\pi(s,a) + \pi_n(a|s) \nabla \sum_{s',r} p(s',r|s,a) [r - \bar{r}(\pi_n) + \tilde{V}_\pi(s')] \right] \\
 &= \sum_a \left[\nabla \pi_n(a|s) \tilde{Q}_\pi(s,a) + \pi_n(a|s) \left[\nabla \bar{r}(\pi_n) + \sum_{s'} p(s'|s,a) \nabla \tilde{V}_\pi(s') \right] \right]
 \end{aligned}$$

Re-arranging terms

$$\nabla \bar{r}(\pi_n) = \sum_a \left[\nabla \pi_n(a|s) \tilde{Q}_\pi(s,a) + \pi_n(a|s) \sum_{s'} p(s'|s,a) \nabla \tilde{V}_\pi(s') \right] - \nabla \tilde{V}_\pi(s)$$

Summing over d_π :

$$\begin{aligned}
 \sum_s d_\pi(s) \nabla \bar{r}(\pi_n) &= \sum_s d_\pi(s) \sum_a \nabla \pi_n(a|s) \tilde{Q}_\pi(s,a) + \sum_s d_\pi(s) \pi_n(a|s) \sum_{s'} p(s',r|s,a) \nabla \tilde{V}_\pi(s') \\
 &\quad - \sum_s d_\pi(s) \nabla \tilde{V}_\pi(s)
 \end{aligned}$$

$$\nabla \bar{r}(\pi_n) = \sum_s d_\pi(s) \sum_a \tilde{Q}_\pi(s,a) \nabla \pi_n(a|s) \quad \text{QED}$$

Approximating the Gradient

PGT

$$\nabla_u F(\pi_u) = \sum_s d_\pi(s) \sum_a \nabla_u \pi_u(a|s) q_{\pi}(s,a)$$

note $\sum_a \nabla_u \pi(a|s) = 0 \quad \forall s$

$$= \sum_s d_\pi(s) \sum_a \nabla_u \pi_u(a|s) [q_{\pi}(s,a) - b(s)] \quad \forall b: S \rightarrow \mathbb{R}$$

$$= \sum_s d_\pi(s) \sum_a \pi_u(a|s) [q_{\pi}(s,a) - b(s)] \underbrace{\frac{\nabla_u \pi_u(a|s)}{\pi_u(a|s)}}_{g(s,a)}$$

$$= E_{\pi} \left[(q_{\pi}(S_t, A_t) - b(S_t)) g(S_t, A_t) \right]$$

$$\begin{aligned} u_{t+1} &= u_t + \alpha \left[\hat{q}(S_t, A_t) - b(S_t) \right] \frac{\nabla \pi_u(A_t|S_t)}{\pi_u(A_t|S_t)} \\ &= u_t + \alpha \left[\hat{q}(S_t, A_t) - b(S_t) \right] \frac{\nabla \pi_u(A_t|S_t)}{\pi_u(A_t|S_t)} \end{aligned}$$

e.g.

$$= u_t + \alpha \left[R_{t+1} - \bar{R}_t + \hat{v}(S_{t+1}, w_t) - \hat{v}(S_t, w_t) \right] \frac{\nabla \pi(A_t|S_t)}{\pi(A_t|S_t)}$$

Elig traces:

$$e_t^u = e_{t-1}^u + \frac{\nabla_u \pi_u(A_t|S_t)}{\pi(A_t|S_t)}$$

\uparrow
 e_t^u

Final Complete Algorithm

$$u, e^u \in \mathbb{R}^n$$

$$v, e^v \in \mathbb{R}^m$$

$$\bar{R} \in \mathbb{R}$$

actor
critic

init
all
to 0

On each step, in state S :

Choose A according to $\pi_u(\cdot | S)$

Take A ; observe S', R

$$\delta \leftarrow R - \bar{R} + v^T x(S') - v^T x(S)$$

$$\bar{R} \leftarrow \bar{R} + \alpha_R \delta$$

$$e^v \leftarrow \lambda e^v + x(S)$$

$$v \leftarrow v + \alpha_v \delta e^v$$

$$e^u \leftarrow \lambda e^u + \frac{\nabla_u \pi_u(A | S)}{\pi_u(A | S)}$$

$$u \leftarrow u + \alpha_u \delta e^u$$

The algorithms mentioned above are independent of the structure of the policy distribution used in the policy. For discrete actions, the Gibbs distribution is often used. In this paper, for continuous actions, we define the policy such that actions are taken according to a normal distribution, as suggested by Williams [10], with a probability density function defined as $\mathcal{N}(s, a) = \frac{1}{\sqrt{2\pi\sigma^2(s)}} \exp\left(-\frac{(a-\mu(s))^2}{2\sigma^2(s)}\right)$ where $\mu(s)$ and $\sigma(s)$ are respectively the mean and standard deviation of the distribution $\pi(\cdot|s)$.

In our policy parameterization, the scalars $\mu(s) = \mathbf{u}_\mu^\top \mathbf{x}_\mu(s)$ and $\sigma(s) = \exp(\mathbf{u}_\sigma^\top \mathbf{x}_\sigma(s))$ are defined as linear combinations, where the parameters of the policy are $\mathbf{u} = (\mathbf{u}_\mu^\top, \mathbf{u}_\sigma^\top)^\top$, and the features vector in state s is $\mathbf{x}_u(s) = (\mathbf{x}_\mu(s)^\top, \mathbf{x}_\sigma(s)^\top)^\top$.

The compatible features $\frac{\nabla_{\mathbf{u}} \pi(a|s)}{\pi(a|s)}$ depend on the structure of the probability density of the policy. Given that our policy density is a normal distribution, the compatible features for the mean and the standard deviation are [10]:

$$\frac{\nabla_{\mathbf{u}_\mu} \pi(a|s)}{\pi(a|s)} = \frac{1}{\sigma(s)^2} (a - \mu(s)) \mathbf{x}_\mu(s) \quad (9)$$

$$\frac{\nabla_{\mathbf{u}_\sigma} \pi(a|s)}{\pi(a|s)} = \left(\frac{(a - \mu(s))^2}{\sigma(s)^2} - 1 \right) \mathbf{x}_\sigma(s) \quad (10)$$

where $\frac{\nabla_{\mathbf{u}} \pi(a|s)}{\pi(a|s)} = \left(\frac{\nabla_{\mathbf{u}_\mu} \pi(a|s)}{\pi(a|s)}, \frac{\nabla_{\mathbf{u}_\sigma} \pi(a|s)}{\pi(a|s)} \right)^\top$.

The compatible feature in (9), used to update the parameters \mathbf{u}_μ of the policy, has a $\frac{1}{\sigma(s)^2}$ factor: the smaller the standard deviation is, the larger the norm of $\frac{\nabla_{\mathbf{u}_\mu} \pi(a_t|s_t)}{\pi(a_t|s_t)}$ is, and vice-versa. We observed that such an effect can cause instability, particularly because $\lim_{\sigma \rightarrow 0} \frac{(a-\mu(s))}{\sigma(s)^2} = \infty$.

The generality of the policy-gradient strategy

- Can be applied whenever we can compute the effect of parameter changes on the action probabilities, $\nabla_{\mathbf{u}}\pi(a|s)$
- E.g., has been applied to spiking neuron models
- There are many possibilities other than linear-exponential and linear-gaussian
 - e.g., mixture of random, argmax, and fixed-width gaussian; learn the mixing weights
 - drift/diffusion models?

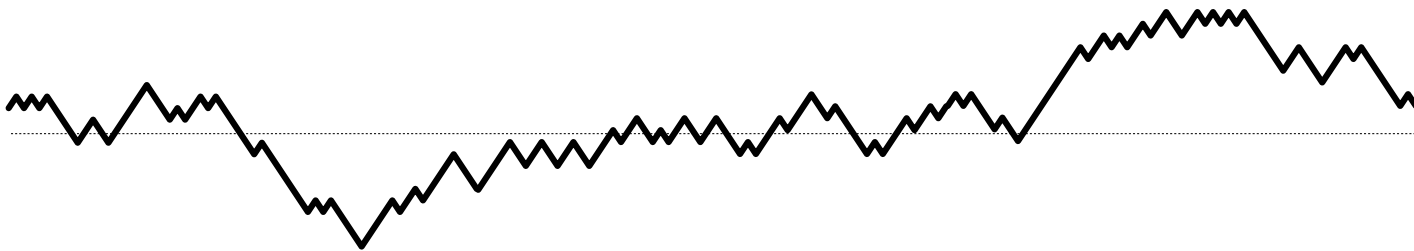
Can policy gradient methods solve the *twitching problem*?

(the problem of decisiveness in adaptive behavior)

- The problem:
 - we need stochastic policies to get exploration
 - but all of our policies have been i.i.d. (independent, identically distributed)
 - if the time step is small, the robot just twitches! 🙄
 - really, no aspect of behavior should depend on the length of the time step

Can we design a parameterized policy whose stochasticity is independent of the time step?

- let a “noise” variable take a random walk, drifting but tending back to zero



- *add* it to the action, and adapt its parameters by the PG algorithm (or have several such noise variables)

The generality of the policy-gradient strategy (2)

- Can be applied whenever we can compute the effect of parameter changes on the action probabilities, $\nabla_{\mathbf{u}}\pi(a|s)$
- Can we apply PG when outcomes are viewed as action?
 - e.g., the action of “turning on the light” or the action of “going to the bank”
 - is this an alternate strategy for temporal abstraction?
- We would need to learn—not compute—the gradient of these states w.r.t. the policy parameter

Have we eliminated action?

- If any state can be an action, then what is still special about actions?
- The parameters/weights are what we can really, directly control
- We have always, in effect, “sensed” our actions (even in the ϵ -greedy case)
- Perhaps actions are just sensory signals that we can usually control easily
- Perhaps there is no longer any need for a special concept of action in the RL framework