

# Introduction to policy search in Reinforcement Learning

Djalel Benbouzid – Data:Lab Munich, Volkswagen Group

June 27<sup>th</sup> 2018



# quick outline

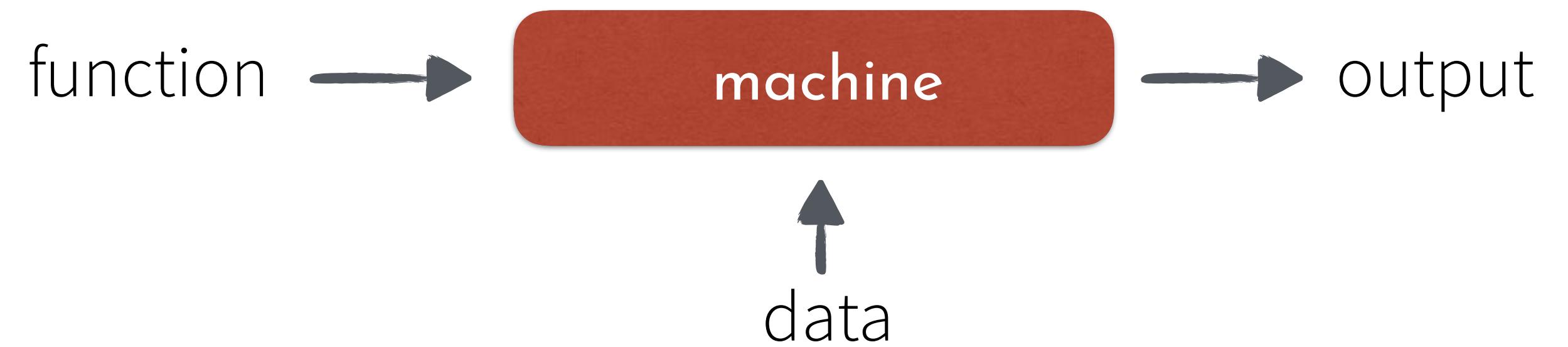
**introduction** and some reminders

introduce a **different way to RL**, wrt. first lecture

**gradients-based** methods

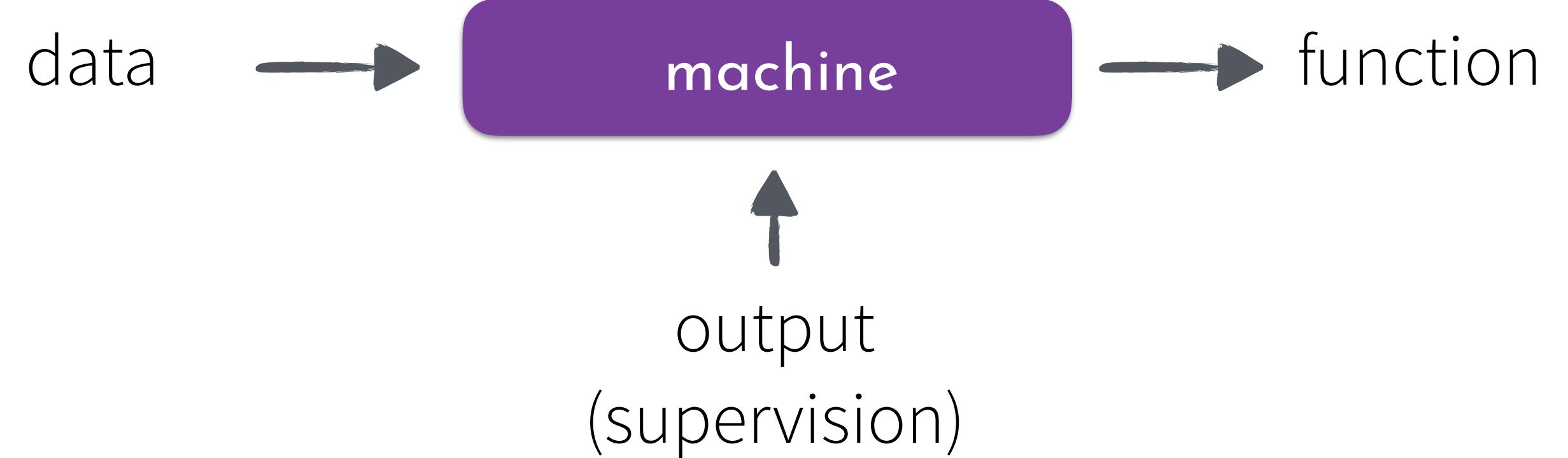
**gradients-free** methods

“traditional”  
programming



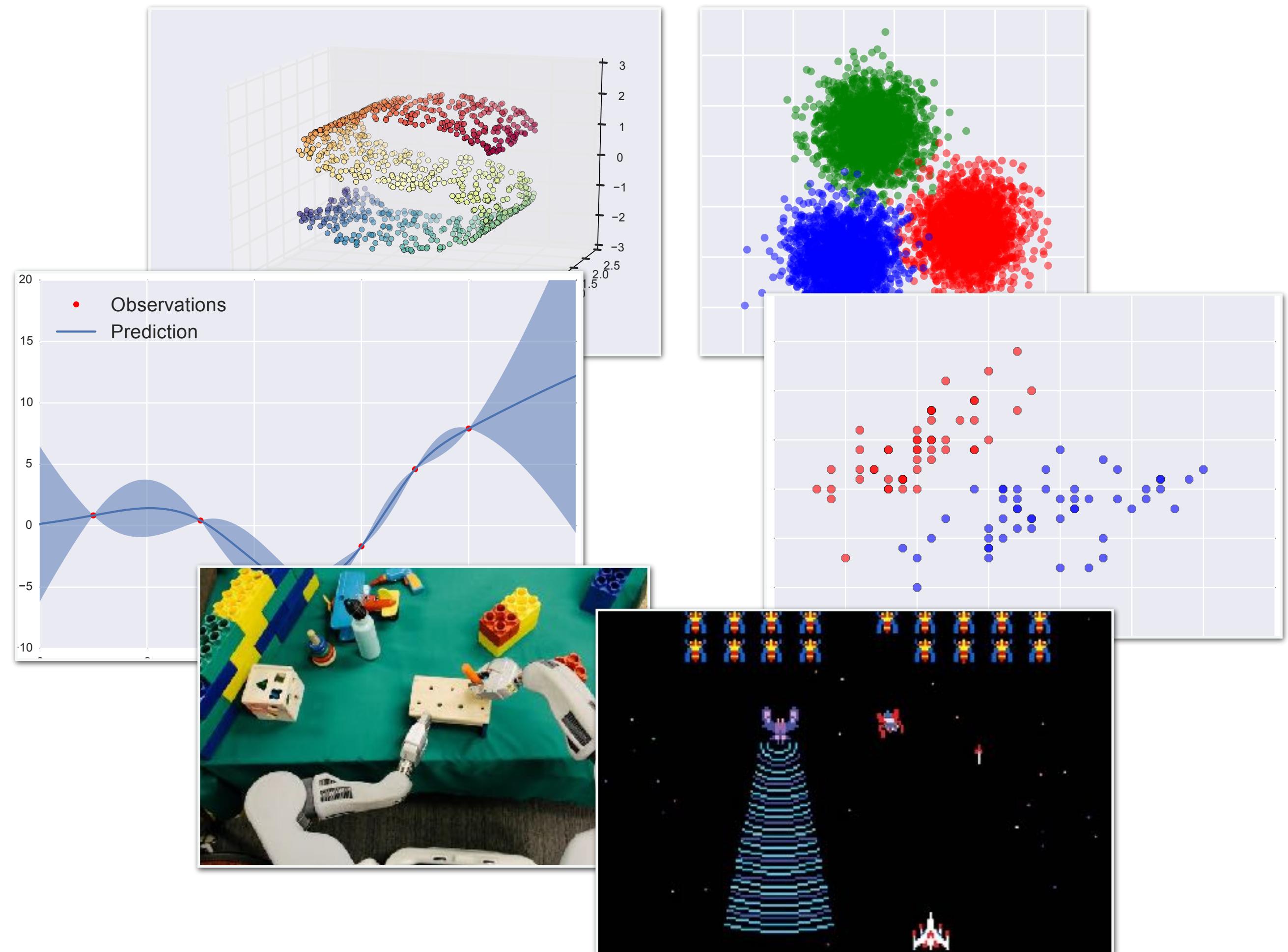
.....

machine learning

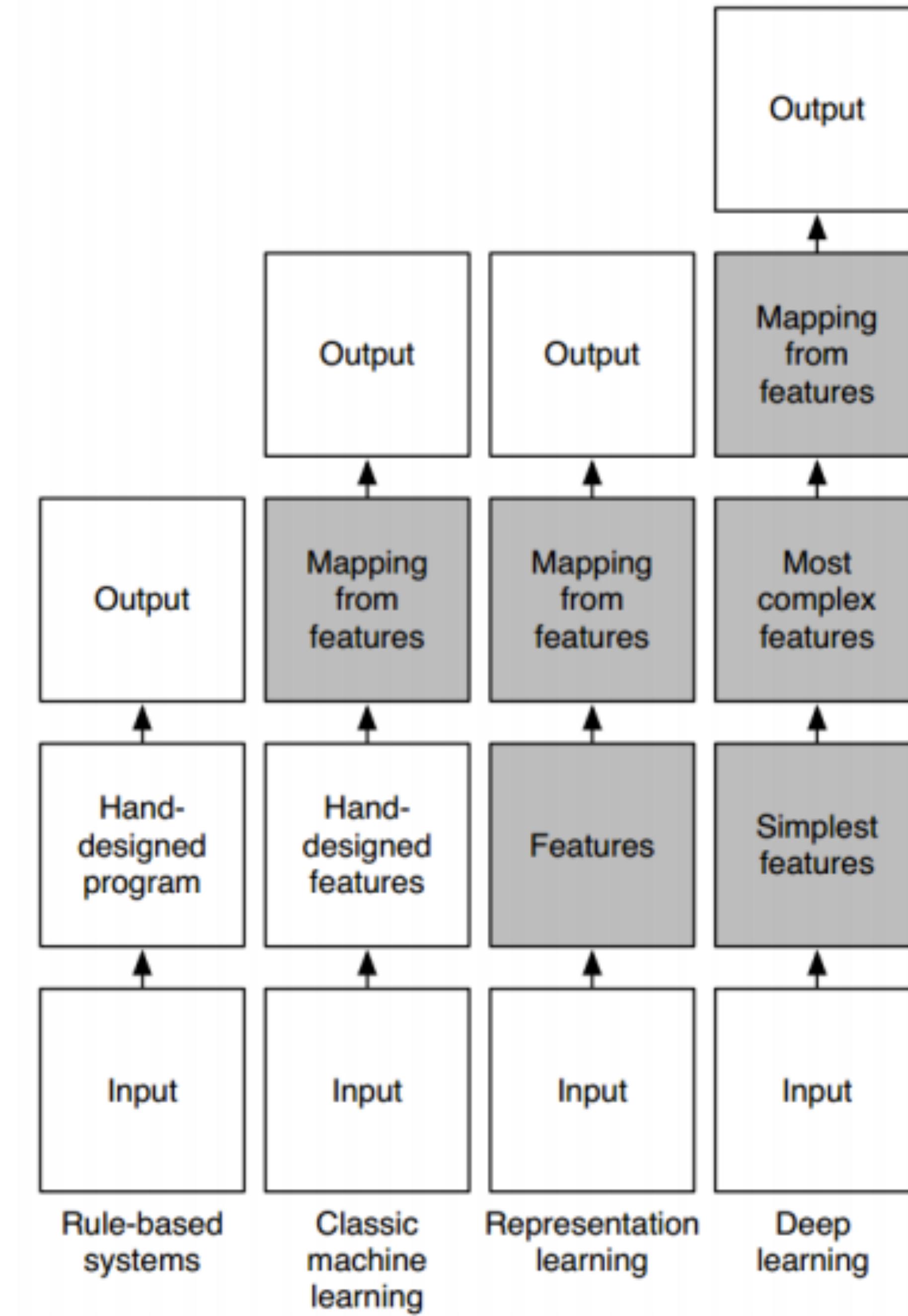


# Machine Learning paradigms

- Supervised Learning
- Unsupervised Learning
- Reinforcement Learning



# what about deep learning?



source: Bengio et. al

# it's all about inferring functions and abstracting problems away

**supervised learning:**  
function mapping, I/O

$$f : \mathcal{X} \rightarrow \mathcal{Y} \qquad f \sim p(y | \mathbf{x})$$

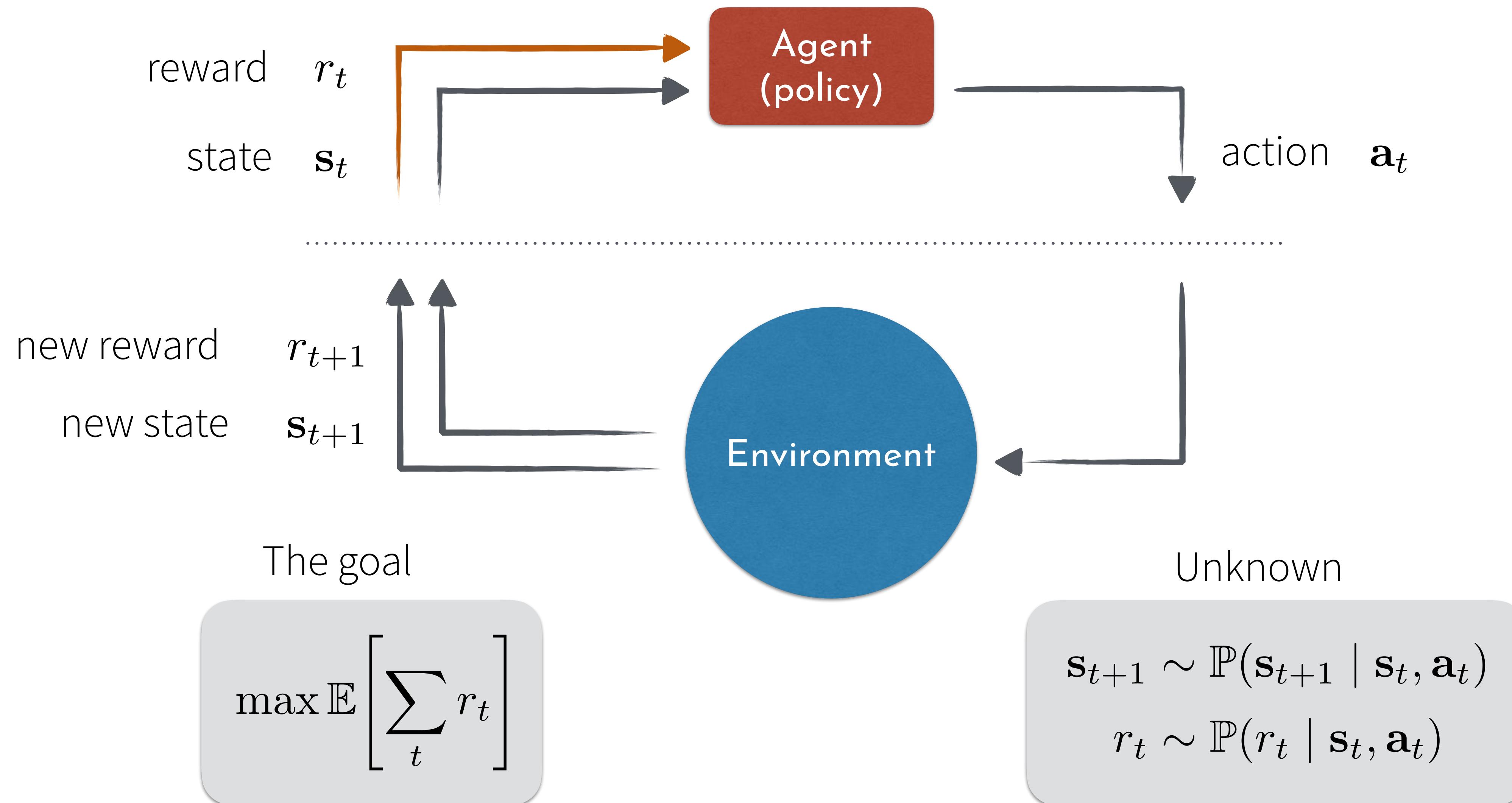
**unsupervised learning:**  
learning densities, and/or representations

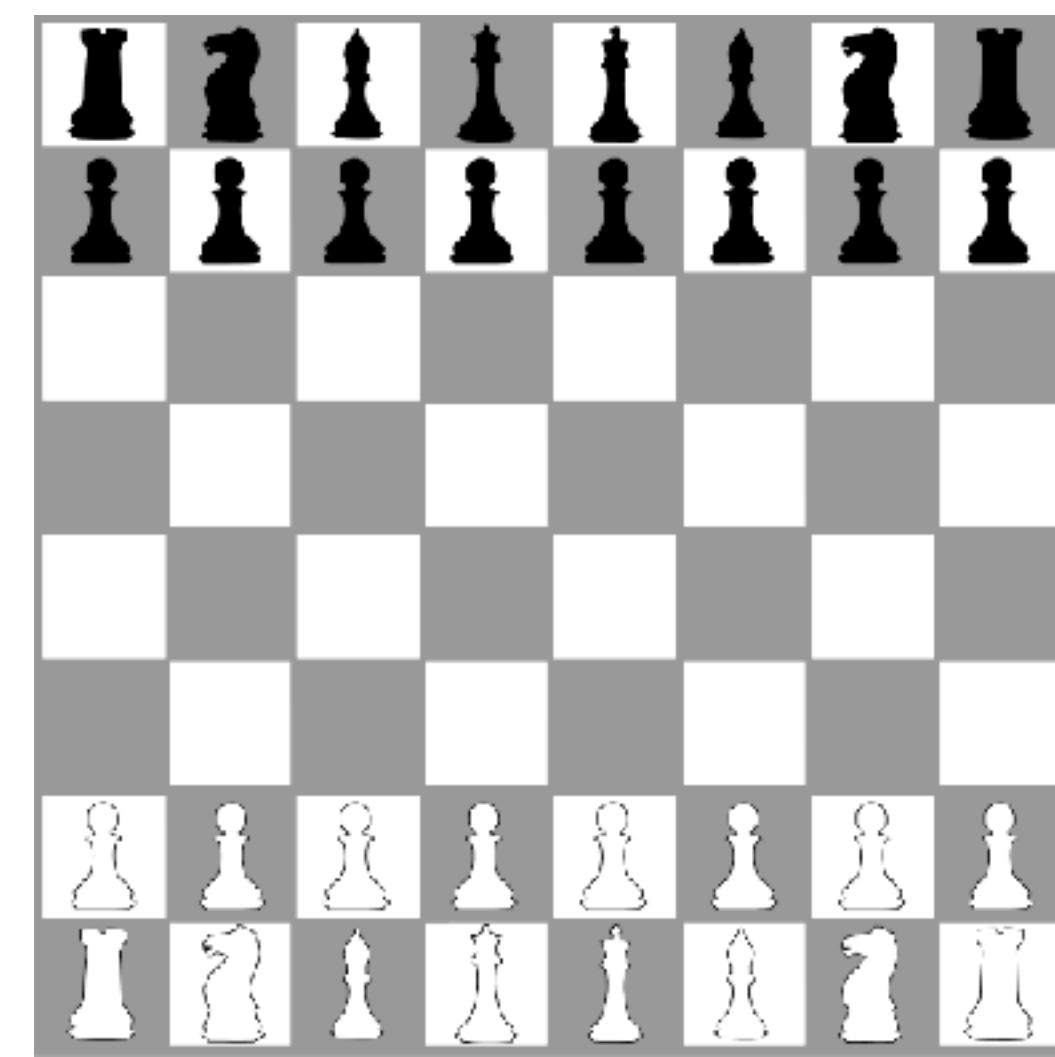
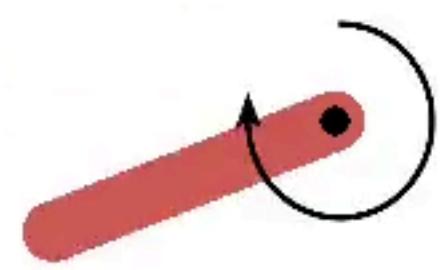
$$p(\mathbf{x}) \qquad \text{or} \qquad p(\mathbf{x}, \mathbf{z})$$

**reinforcement learning:**  
learn a series of actions, a policy

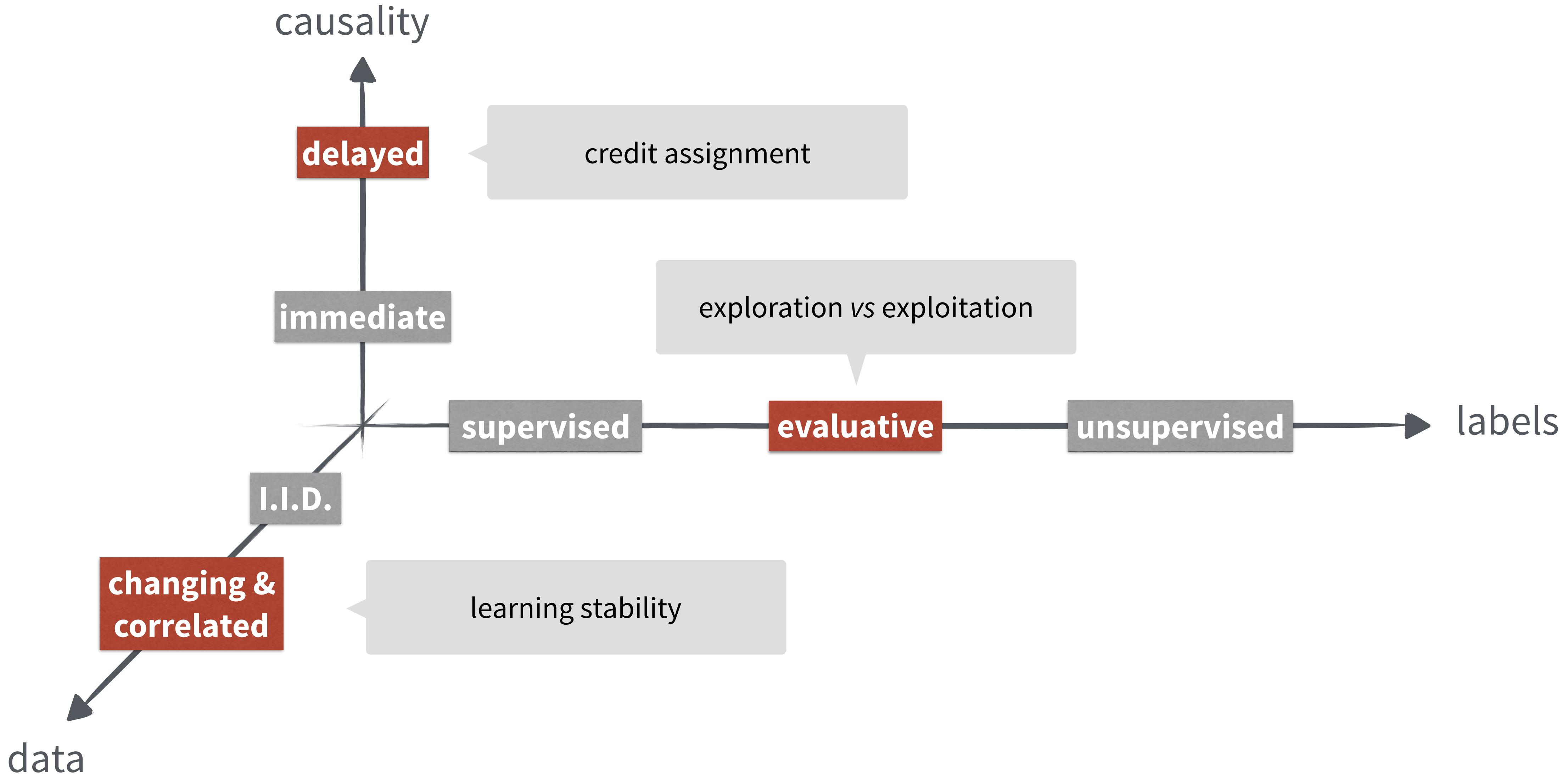
$$\pi(\mathbf{s}) \mapsto a_1, a_2, a_3, \dots$$

# the RL paradigm





# what makes **RL** harder than **SL**?



# two main approaches for Reinforcement Learning

## Value-based

$$V(s) = \mathbb{E}_t[R_t \mid S_t = s]$$

$$Q(s, a) = \mathbb{E}_t[R_t \mid S_t = s, A_t = a]$$

$$\pi^*(s) = \arg \max_a Q^*(s, a)$$

- dynamic programming
- temporal-difference (TD) learning
- indirect method in some sense
- not for large / continuous action space  
(until recently)

## Policy-based

$$\pi_\theta(s, a) = \mathbb{P}(a \mid s, \theta)$$

optimization problem  
Loss function  $J(\theta)$   
directly search in  $\theta$  space

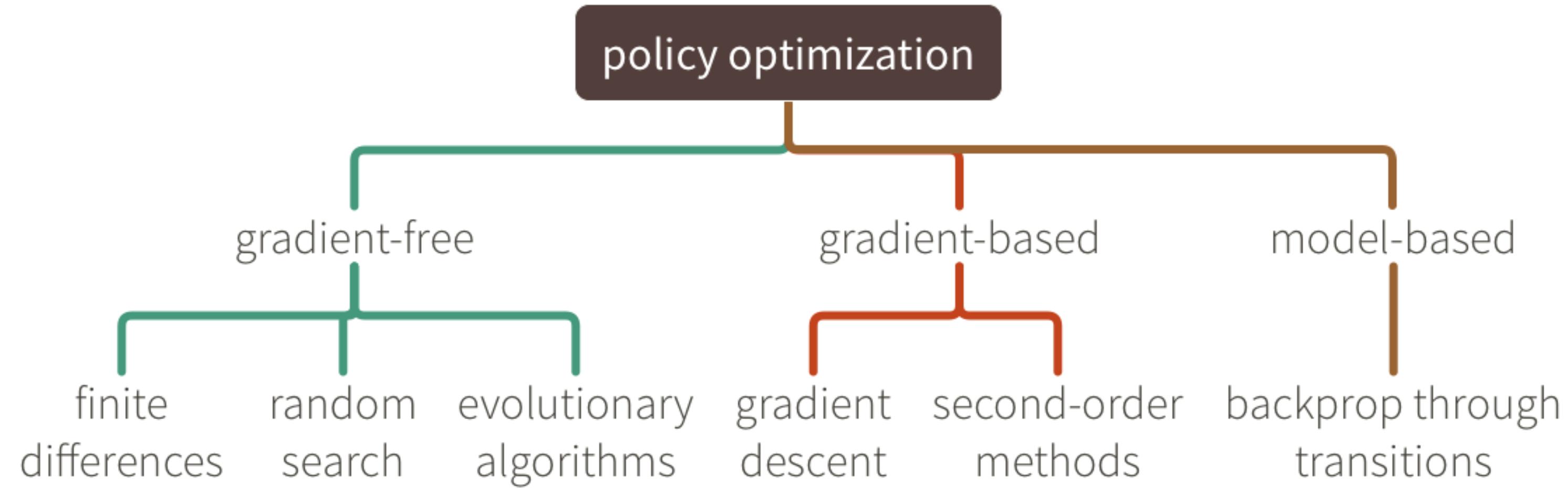
- stochastic policies
- handles continuous action space
- but... suffers from local optima and high variance

# policy search (or optimization)

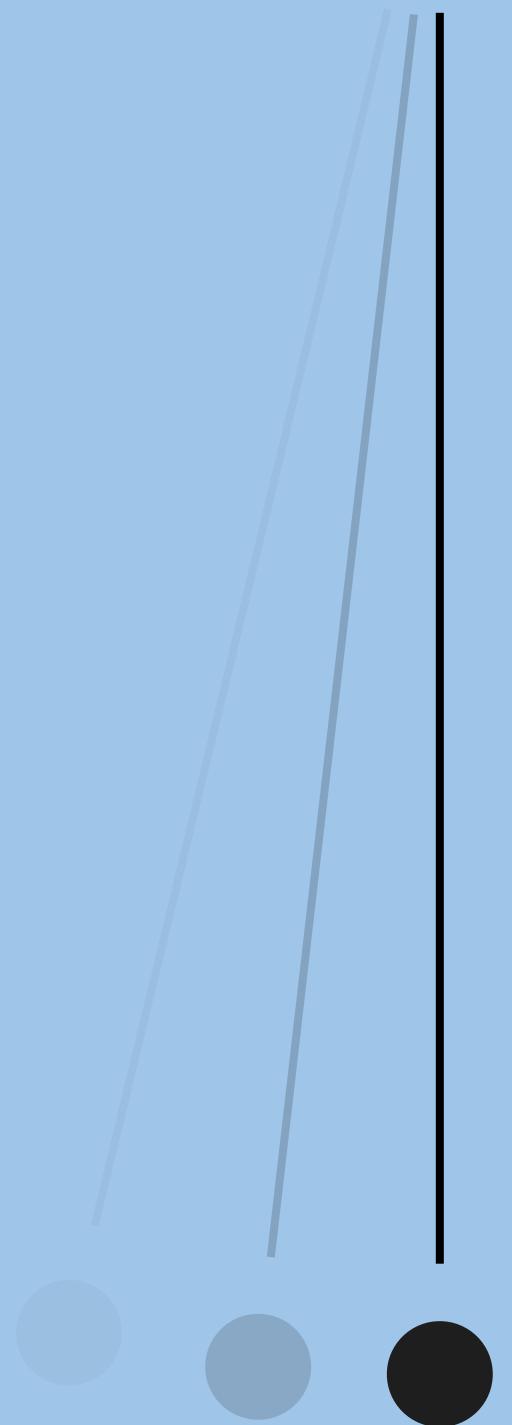
Parametrized policy       $a \sim \pi_\theta(s) = \mathbb{P}(a \mid s, \theta) \quad \text{or} \quad a = \pi_\theta(s) = f(a \mid s, \theta)$

Objective function       $J(\theta) = \mathbb{E}_{p_\theta(\mathbf{s}_{1:t}, \mathbf{a}_{1:t})} \left[ \sum_t r(\mathbf{s}_t, \mathbf{a}_t) \right]$

Find  $\theta$  that maximizes  $J(\theta)$



# Policy Gradient



# policy gradient: general formulation

actions

$$a \sim \pi_\theta(s) = \mathbb{P}(a \mid s, \theta)$$

transitions

$$\mathbf{s}_{t+1} \sim p(\mathbf{s}_{t+1} \mid \mathbf{s}_t, \mathbf{a})$$

trajectories

$$\tau = (\mathbf{s}_1, \mathbf{a}_1, \dots, \mathbf{s}_T, \mathbf{a}_T)$$

$$p_\theta(\tau) = p(\mathbf{s}_1) \prod_{t=1}^T \pi_\theta(\mathbf{a}_t \mid \mathbf{s}_t) p(\mathbf{s}_{t+1} \mid \mathbf{s}_t, \mathbf{a}_t)$$

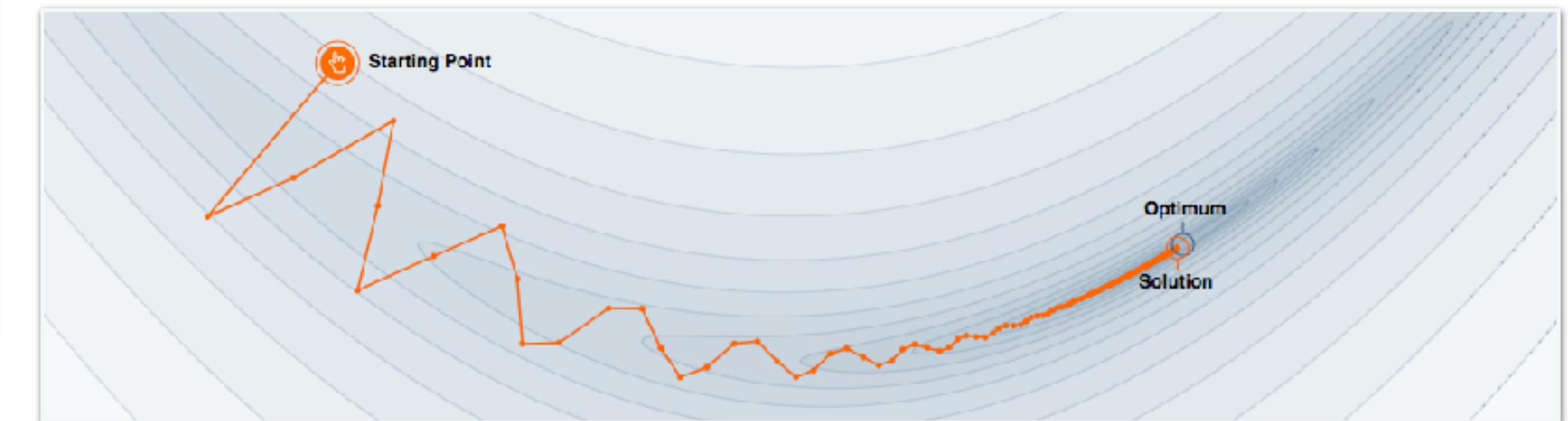
cumulative reward

$$r(\tau) = \sum_t r(\mathbf{s}_t, \mathbf{a}_t)$$

$$\text{Objective function } J(\theta) = \mathbb{E}_{\tau \sim p_\theta(\tau)} [r(\tau)]$$

$$= \int r(\tau) p_\theta(\tau) d\tau$$

$$\theta^\star = \arg \max_\theta \mathbb{E}_{\tau \sim p_\theta(\tau)} \left[ \sum_t r(\mathbf{s}_t, \mathbf{a}_t) \right]$$



non-differentiable objective function!

⇒ likelihood-ratio trick

$$\begin{aligned}
\nabla_{\theta} J(\theta) &= \nabla_{\theta} \int r(\tau) p_{\theta}(\tau) d\tau \\
&= \int r(\tau) \nabla_{\theta} p_{\theta}(\tau) d\tau \\
&= \int r(\tau) \nabla_{\theta} p_{\theta}(\tau) \left( \frac{p_{\theta}(\tau)}{p_{\theta}(\tau)} \right) d\tau \\
&= \int r(\tau) p_{\theta}(\tau) \left( \frac{\nabla_{\theta} p_{\theta}(\tau)}{p_{\theta}(\tau)} \right) d\tau \\
&= \int r(\tau) p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau) d\tau \\
&= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [r(\tau) \nabla_{\theta} \log p_{\theta}(\tau)]
\end{aligned}$$

scale by the trajectory evaluation and follow the gradient of the trajectories logprop  
(in expectation)

many names in the literature: score function estimator, likelihood-ratio trick, etc.

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [r(\tau) \nabla_{\theta} \log p_{\theta}(\tau)]$$

$$p_{\theta}(\tau) = p(\mathbf{s}_1) \prod_{t=1}^T \pi_{\theta}(\mathbf{a}_t \mid \mathbf{s}_t) p(\mathbf{s}_{t+1} \mid \mathbf{s}_t, \mathbf{a}_t)$$

~~$$\nabla_{\theta} \log p_{\theta}(\tau) = \log p(\mathbf{s}_1) + \sum_{t=1}^T \pi_{\theta}(\mathbf{a}_t \mid \mathbf{s}_t) + p(\mathbf{s}_{t+1} \mid \mathbf{s}_t, \mathbf{a}_t)$$~~

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[ \sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t \mid \mathbf{s}_t) \right]$$

$$\approx \frac{1}{N} \sum_{i=1}^N \left( \sum_{t=1}^T r(\mathbf{s}_t^{(i)}, \mathbf{a}_t^{(i)}) \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^{(i)} \mid \mathbf{s}_t^{(i)}) \right)$$

⇒ the REINFORCE algorithm

unbiased estimator

- 1 - interact with the environment (generate  $N$  trajectories)
- 2 - compute their cumulative reward
- 3 - compute an estimate of loss gradient  $g \approx \nabla_{\theta} J(\theta)$
- 4 - update the policy parameters  $\theta \leftarrow \theta + \alpha \times g$

Williams, Ronald J. "Simple statistical gradient-following algorithms for connectionist reinforcement learning." Reinforcement Learning. Springer, Boston, MA, 1992.

Sutton, Richard S., et al. "Policy gradient methods for reinforcement learning with function approximation." Advances in neural information processing systems. 2000.

## REINFORCE, A Monte-Carlo Policy-Gradient Method (episodic)

Input: a differentiable policy parameterization  $\pi(a|s, \boldsymbol{\theta})$ ,  $\forall a \in \mathcal{A}, s \in \mathcal{S}, \boldsymbol{\theta} \in \mathbb{R}^d$

Initialize policy parameter  $\boldsymbol{\theta}$

Repeat forever:

    Generate an episode  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$ , following  $\pi(\cdot| \cdot, \boldsymbol{\theta})$

    For each step of the episode  $t = 0, \dots, T - 1$ :

$G \leftarrow$  return from step  $t$

$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha \gamma^t G \nabla_{\boldsymbol{\theta}} \log \pi(A_t | S_t, \boldsymbol{\theta})$

$$\pi_\theta(\mathbf{a}_t \mid \mathbf{s}_t) \sim \mathcal{N}(\mu(\mathbf{s}_t), \Sigma)$$



often a neural network

## Gaussian policies

$$\nabla_\theta \log \pi_\theta(\mathbf{a}_t \mid \mathbf{s}_t) = \frac{1}{2} \Sigma^{-1} (\mathbf{a}_t - \mu(\mathbf{s}_t)) \nabla_\theta \mu(\mathbf{s}_t)$$

## Softmax policies

$$\pi_\theta(\mathbf{a}_t \mid \mathbf{s}_t) = \frac{e^{f(\mathbf{s}_t, \mathbf{a}_t)}}{\sum_{\mathbf{a}'_t} e^{f(\mathbf{s}_t, \mathbf{a}'_t)}}$$

$$\nabla_\theta \log \pi_\theta(\mathbf{a}_t \mid \mathbf{s}_t) = f(\mathbf{s}_t, \mathbf{a}_t) - \mathbb{E}_{\pi_\theta}[f(\mathbf{s}_t, \cdot)]$$

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [r(\tau) \nabla_{\theta} \log p_{\theta}(\tau)]$$

$$p_{\theta}(\tau) = p(\mathbf{s}_1) \prod_{t=1}^T \pi_{\theta}(\mathbf{a}_t \mid \mathbf{s}_t) p(\mathbf{s}_{t+1} \mid \mathbf{s}_t, \mathbf{a}_t)$$

~~$$\nabla_{\theta} \log p_{\theta}(\tau) = \log p(\mathbf{s}_1) + \sum_{t=1}^T \pi_{\theta}(\mathbf{a}_t \mid \mathbf{s}_t) + p(\mathbf{s}_{t+1} \mid \mathbf{s}_t, \mathbf{a}_t)$$~~

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[ \sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t \mid \mathbf{s}_t) \right]$$

$$\approx \frac{1}{N} \sum_{i=1}^N \left( \sum_{t=1}^T r(\mathbf{s}_t^{(i)}, \mathbf{a}_t^{(i)}) \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^{(i)} \mid \mathbf{s}_t^{(i)}) \right)$$

⇒ the REINFORCE algorithm

unbiased estimator

1 - interact with the environment (generate  $N$  trajectories)

high variance!

2 - compute their cumulative reward

3 - compute an estimate of loss gradient  $g \approx \nabla_{\theta} J(\theta)$

4 - update the policy parameters  $\theta \leftarrow \theta + \alpha \times g$

Williams, Ronald J. "Simple statistical gradient-following algorithms for connectionist reinforcement learning." Reinforcement Learning. Springer, Boston, MA, 1992.

Sutton, Richard S., et al. "Policy gradient methods for reinforcement learning with function approximation." Advances in neural information processing systems. 2000.

# reducing the variance with the temporal structure

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[ \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t \mid \mathbf{s}_t) \sum_{t'=t}^{\textcolor{red}{T}} r(\mathbf{s}_{\textcolor{red}{t}'}, \mathbf{a}_{\textcolor{red}{t}'}) \right] \quad \begin{matrix} \text{actions should only be} \\ \text{influenced by future rewards} \end{matrix}$$

$$= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[ \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t \mid \mathbf{s}_t) Q^{\pi_{\theta}}(\mathbf{s}_t, \mathbf{a}_t) \right]$$

the return  $R_t$  is a sample of the Q value function  $R_t = \sum_{t'=t}^T r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \sim Q^{\pi_{\theta}}(\mathbf{s}_t, \mathbf{a}_t)$

# reducing the variance with a baseline

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[ \sum_{t=1}^T \nabla_{\theta} \log p_{\theta}(\tau) \left( \sum_{t'=t}^T r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) - b(\mathbf{s}_t) \right) \right]$$

$$\begin{aligned} \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} \log p_{\theta}(\tau) B] &= \int p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau) B \, d\tau \\ &= \int p_{\theta}(\tau) \frac{\nabla_{\theta} p_{\theta}(\tau)}{p_{\theta}(\tau)} B \, d\tau \\ &= \int \nabla_{\theta} p_{\theta}(\tau) B \, d\tau \\ &= \nabla_{\theta} \int p_{\theta}(\tau) B \, d\tau \\ &= B \, \nabla_{\theta} \int p_{\theta}(\tau) \, d\tau \\ &= B \times 0 \end{aligned}$$

⇒ still unbiased estimator

# which baseline?

constant, or moving average     $b(\mathbf{s}_t) = b \approx \frac{1}{N} \sum_t r(\mathbf{s}_t, \mathbf{a}_t)$

estimate the value function     $b(\mathbf{s}_t) = \hat{V}_\nu(\mathbf{s}_t) = \mathbb{E}_t[R_t \mid S_t = s]$

**intuition:** only increase probability of actions that improve over the average

*More substance in*

Greensmith, Evan, Peter L. Bartlett, and Jonathan Baxter. "Variance reduction techniques for gradient estimates in reinforcement learning." *Journal of Machine Learning Research*, 2004

# let's recap

$$\theta^\star = \arg \max_{\theta} \mathbb{E}_{\tau \sim p_\theta(\tau)} \left[ \sum_t r(\mathbf{s}_t, \mathbf{a}_t) \right]$$

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim p_\theta(\tau)} \left[ \sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t \mid \mathbf{s}_t) \right]$$

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim p_\theta(\tau)} \left[ \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t \mid \mathbf{s}_t) \ Q^{\pi_{\theta}}(\mathbf{s}_t, \mathbf{a}_t) \right]$$

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim p_\theta(\tau)} \left[ \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t \mid \mathbf{s}_t) \ Q^{\pi_{\theta}}(\mathbf{s}_t, \mathbf{a}_t) - V^{\pi_{\theta}}(\mathbf{s}_t) \right]$$

# actor-critic algorithms

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[ \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) Q^{\pi_{\theta}}(\mathbf{s}_t, \mathbf{a}_t) \right]$$

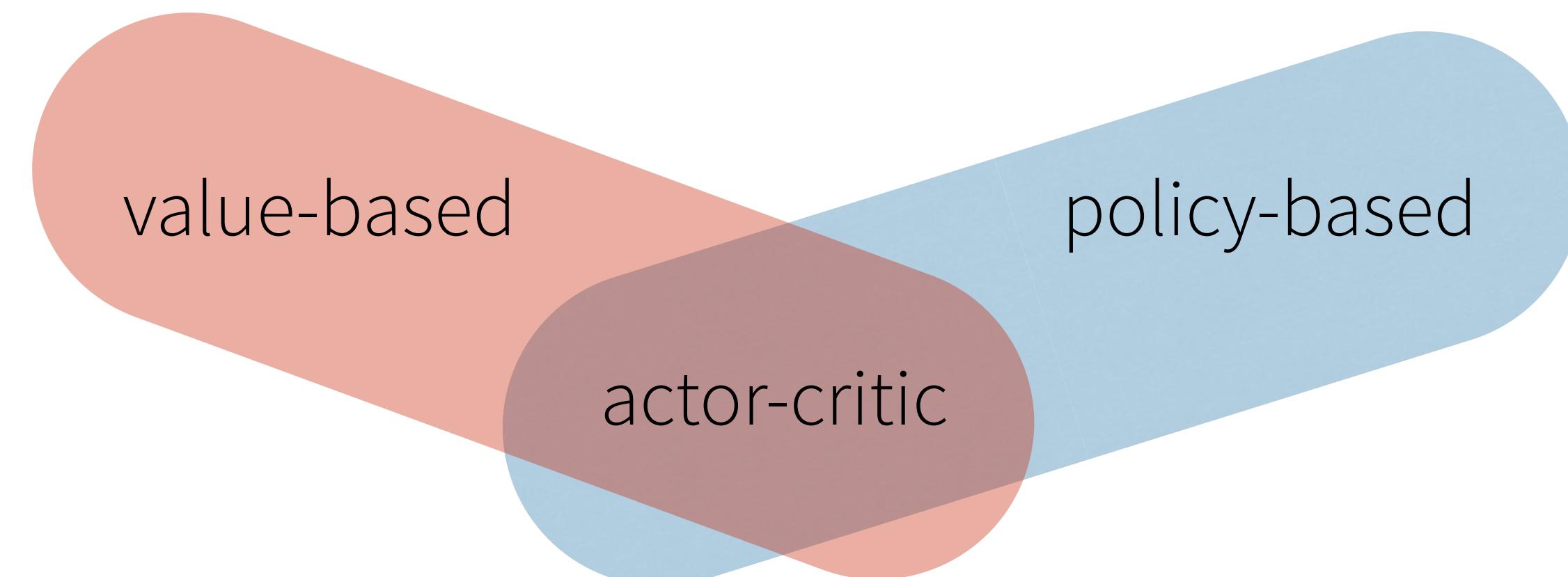
déjà-vu?

Q-learning!

$$\approx \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[ \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \hat{Q}_{\omega}(\mathbf{s}_t, \mathbf{a}_t) \right]$$

instead of sampling  $Q^{\pi_{\theta}}(\mathbf{s}_t, \mathbf{a}_t)$   
one can approximate it  $\hat{Q}_{\omega}(\mathbf{s}_t, \mathbf{a}_t)$

we now have two parameters to fit  
-  $\theta$  for the policy (**the actor**)  
-  $\omega$  for the value function (**the critic**)



# One-step Actor–Critic (episodic)

Input: a differentiable policy parameterization  $\pi(a|s, \theta)$ ,  $\forall a \in \mathcal{A}, s \in \mathcal{S}, \theta \in \mathbb{R}^d$

Input: a differentiable state-value parameterization  $\hat{v}(s, \mathbf{w})$ ,  $\forall s \in \mathcal{S}, \mathbf{w} \in \mathbb{R}^m$

Parameters: step sizes  $\alpha > 0$ ,  $\beta > 0$

Initialize policy parameter  $\theta$  and state-value weights  $\mathbf{w}$

# Repeat forever:

Initialize  $S$  (first state of episode)

$$I \leftarrow 1$$

While  $S$  is not terminal:

$$A \sim \pi(\cdot | S, \theta)$$

Take action  $A$ , observe  $S', R$

$$\delta \leftarrow R + \gamma \hat{v}(S', \mathbf{w}) - \hat{v}(S, \mathbf{w})$$

(if  $S'$  is terminal, then  $\hat{v}(S', \mathbf{w}) \doteq 0$ )

$$\mathbf{w} \leftarrow \mathbf{w} + \beta \delta \nabla_{\mathbf{w}} \hat{v}(S, \mathbf{w})$$

$$\theta \leftarrow \theta + \alpha I \delta \nabla_{\theta} \log \pi(A|S, \theta)$$

$$I \leftarrow \gamma I$$

$$S \leftarrow S'$$

# actor-critic algorithms with a baseline

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[ \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t \mid \mathbf{s}_t) \underbrace{Q^{\pi_{\theta}}(\mathbf{s}_t, \mathbf{a}_t) - V^{\pi_{\theta}}(\mathbf{s}_t)}_{A^{\pi_{\theta}}(\mathbf{s}_t, \mathbf{a}_t)} \right]$$

advantage function

TD error     $\delta_t^{\pi_{\theta}} = r_t + \gamma V^{\pi_{\theta}}(\mathbf{s}_{t+1}) - V^{\pi_{\theta}}(\mathbf{s}_t, )$

is an unbiased sample of the Advantage function     $\delta_t^{\pi_{\theta}} \sim A^{\pi_{\theta}}(\mathbf{s}_t, \mathbf{a}_t)$

# family of policy gradient algorithms

$\nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t \mid \mathbf{s}_t) R_t$	vanilla policy gradient
$\nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t \mid \mathbf{s}_t) (R_t - \hat{V}_{\nu})$	REINFORCE
$\nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t \mid \mathbf{s}_t) \hat{Q}_{\omega}(\mathbf{s}_t, \mathbf{a}_t)$	actor-critic
$\nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t \mid \mathbf{s}_t) \hat{A}_{\omega}(\mathbf{s}_t, \mathbf{a}_t)$	advantage actor-critic
$\nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t \mid \mathbf{s}_t) \delta_t$	TD actor-critic

# example: Asynchronous Advantage Actor-Critic (A3C)

---

**Algorithm S3** Asynchronous advantage actor-critic - pseudocode for each actor-learner thread.

---

// Assume global shared parameter vectors  $\theta$  and  $\theta_v$  and global shared counter  $T = 0$

// Assume thread-specific parameter vectors  $\theta'$  and  $\theta'_v$

Initialize thread step counter  $t \leftarrow 1$

**repeat**

    Reset gradients:  $d\theta \leftarrow 0$  and  $d\theta_v \leftarrow 0$ .

    Synchronize thread-specific parameters  $\theta' = \theta$  and  $\theta'_v = \theta_v$

$t_{start} = t$

    Get state  $s_t$

**repeat**

        Perform  $a_t$  according to policy  $\pi(a_t|s_t; \theta')$

        Receive reward  $r_t$  and new state  $s_{t+1}$

$t \leftarrow t + 1$

$T \leftarrow T + 1$

**until** terminal  $s_t$  **or**  $t - t_{start} == t_{max}$

$$R = \begin{cases} 0 & \text{for terminal } s_t \\ V(s_t, \theta'_v) & \text{for non-terminal } s_t // \text{Bootstrap from last state} \end{cases}$$

**for**  $i \in \{t - 1, \dots, t_{start}\}$  **do**

$R \leftarrow r_i + \gamma R$

    Accumulate gradients wrt  $\theta'$ :  $d\theta \leftarrow d\theta + \nabla_{\theta'} \log \pi(a_i|s_i; \theta')(R - V(s_i; \theta'_v))$

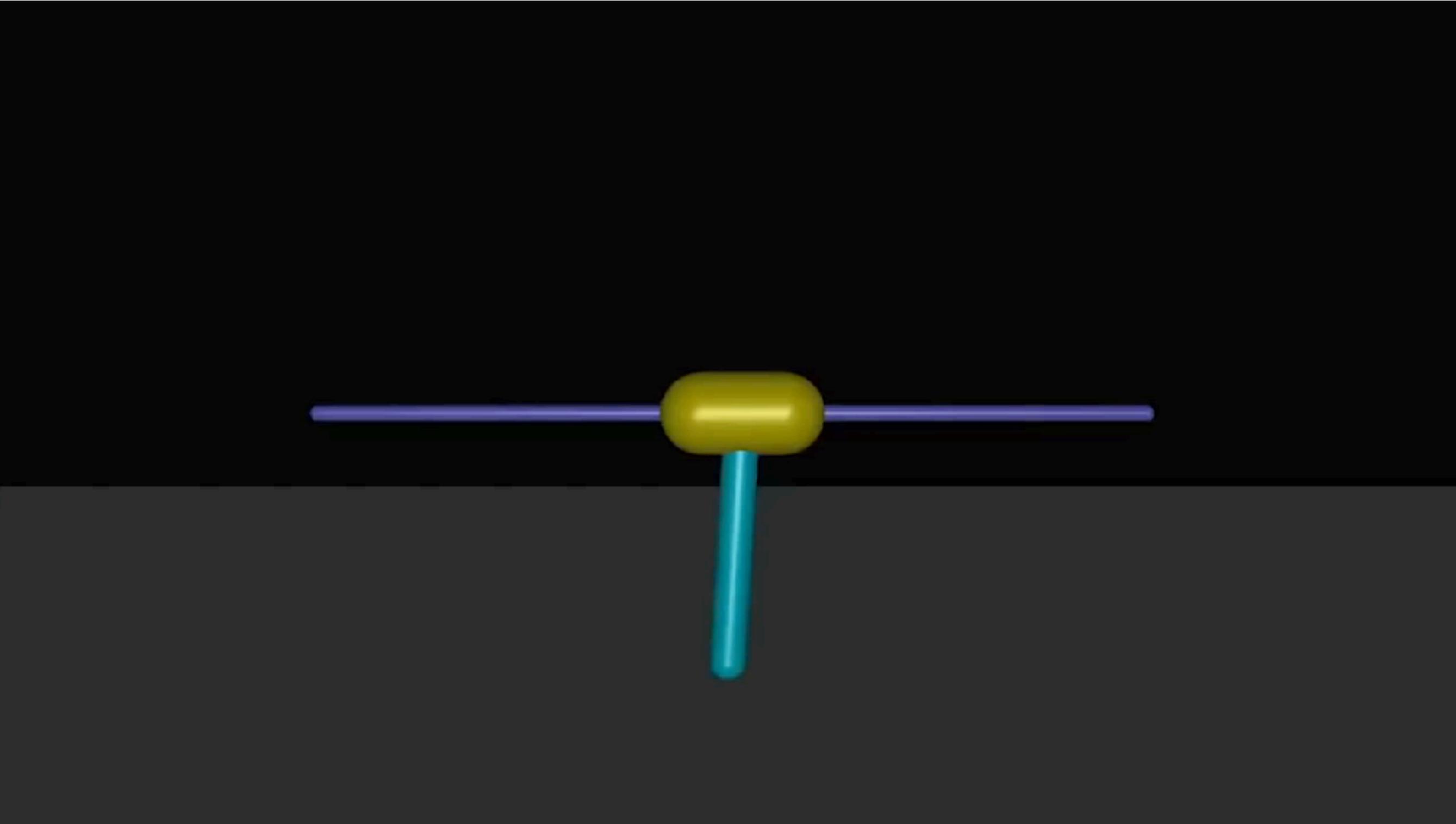
    Accumulate gradients wrt  $\theta'_v$ :  $d\theta_v \leftarrow d\theta_v + \partial(R - V(s_i; \theta'_v))^2 / \partial \theta'_v$

**end for**

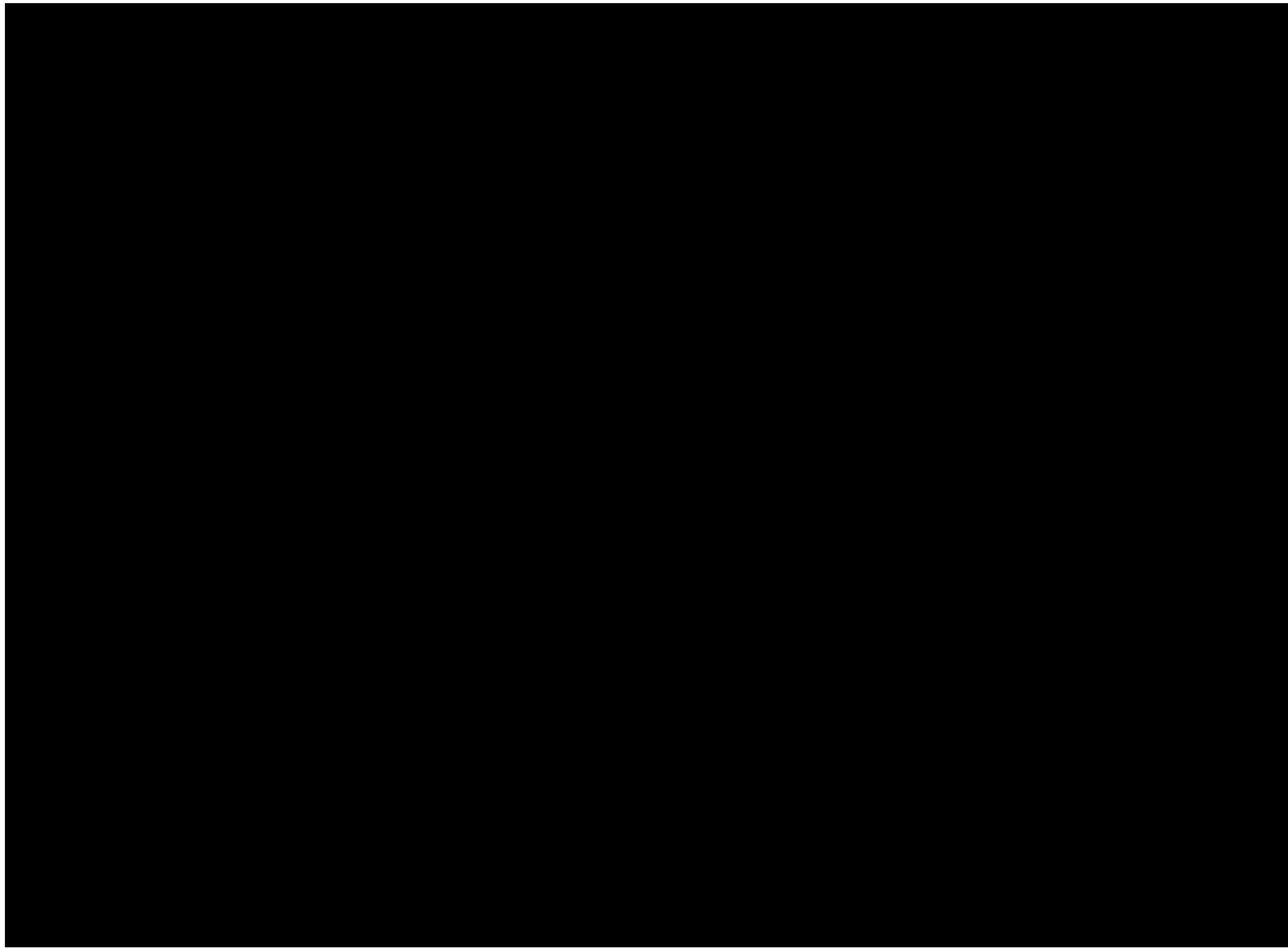
    Perform asynchronous update of  $\theta$  using  $d\theta$  and of  $\theta_v$  using  $d\theta_v$ .

**until**  $T > T_{max}$

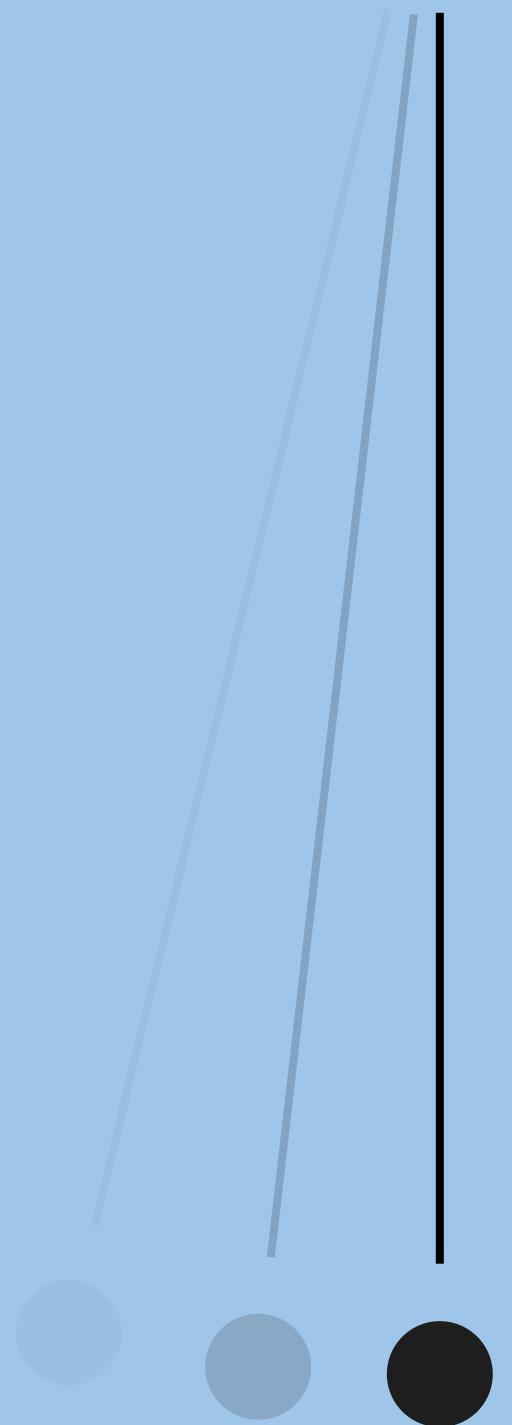
---



using the MuJoCo library



gradient-free methods



# finite differences

- approximate the gradient with numerical methods
- perturb the  $k^{\text{th}}$  component of  $\theta$  by a small step size  $\nu$
- sample  $\delta$  (eg. from a standard Gaussian distribution)
- $$\frac{\partial J(\theta)}{\partial \theta_k} \approx \frac{J(\theta + \nu\delta) - J(\theta - \nu\delta)}{2\nu}$$

---

**Algorithm 1** Basic Random Search (BRS)

- 1: **Hyperparameters:** step-size  $\alpha$ , number of directions sampled per iteration  $N$ , standard deviation of the exploration noise  $\nu$
- 2: **Initialize:**  $\theta_0 = \mathbf{0}$ , and  $j = 0$ .
- 3: **while** ending condition not satisfied **do**
- 4:     Sample  $\delta_1, \delta_2, \dots, \delta_N$  of the same size as  $\theta_j$ , with i.i.d. standard normal entries.
- 5:     Collect  $2N$  rollouts of horizon  $H$  and their corresponding rewards using the policies

$$\pi_{j,k,+}(x) = \pi_{\theta_j + \nu \delta_k}(x) \quad \text{and} \quad \pi_{j,k,-}(x) = \pi_{\theta_j - \nu \delta_k}(x),$$

with  $k \in \{1, 2, \dots, N\}$ .

- 6:     Make the update step:

$$\theta_{j+1} = \theta_j + \frac{\alpha}{N} \sum_{k=1}^N [r(\pi_{j,k,+}) - r(\pi_{j,k,-})] \delta_k .$$

- 7:      $j \leftarrow j + 1$ .

- 8: **end while**

.

---

augmented random search

- 1: **Hyperparameters:** step-size  $\alpha$ , number of directions sampled per iteration  $N$ , standard deviation of the exploration noise  $\nu$ , number of top-performing directions to use  $b$  ( $b < N$  is allowed only for **V1-t** and **V2-t**)
- 2: **Initialize:**  $M_0 = \mathbf{0} \in \mathbb{R}^{p \times n}$ ,  $\mu_0 = \mathbf{0} \in \mathbb{R}^n$ , and  $\Sigma_0 = \mathbf{I}_n \in \mathbb{R}^{n \times n}$ ,  $j = 0$ .
- 3: **while** ending condition not satisfied **do**
- 4:   Sample  $\delta_1, \delta_2, \dots, \delta_N$  in  $\mathbb{R}^{p \times n}$  with i.i.d. standard normal entries.
- 5:   Collect  $2N$  rollouts of horizon  $H$  and their corresponding rewards using the  $2N$  policies

$$\begin{aligned}\mathbf{V1:} \quad & \begin{cases} \pi_{j,k,+}(x) = (M_j + \nu\delta_k)x \\ \pi_{j,k,-}(x) = (M_j - \nu\delta_k)x \end{cases} \\ \mathbf{V2:} \quad & \begin{cases} \pi_{j,k,+}(x) = (M_j + \nu\delta_k) \text{diag}(\Sigma_j)^{-1/2} (x - \mu_j) \\ \pi_{j,k,-}(x) = (M_j - \nu\delta_k) \text{diag}(\Sigma_j)^{-1/2} (x - \mu_j) \end{cases}\end{aligned}$$

for  $k \in \{1, 2, \dots, N\}$ .

- 6:   Sort the directions  $\delta_k$  by  $\max\{r(\pi_{j,k,+}), r(\pi_{j,k,-})\}$ , denote by  $\delta_{(k)}$  the  $k$ -th largest direction, and by  $\pi_{j,(k),+}$  and  $\pi_{j,(k),-}$  the corresponding policies.
- 7:   Make the update step:

$$M_{j+1} = M_j + \frac{\alpha}{b\sigma_R} \sum_{k=1}^b [r(\pi_{j,(k),+}) - r(\pi_{j,(k),-})] \delta_{(k)},$$

where  $\sigma_R$  is the standard deviation of the  $2b$  rewards used in the update step.

- 8:   **V2 :** Set  $\mu_{j+1}$ ,  $\Sigma_{j+1}$  to be the mean and covariance of the  $2NH(j+1)$  states encountered from the start of training.<sup>2</sup>
- 9:    $j \leftarrow j + 1$
- 10: **end while**

augmented  
random  
search

# cross-entropy method

- generate random trajectories according to  $\mathbb{P}(\cdot | \theta)$
- rank them according to their performance (eg. sum of rewards)
- use the best performing ones to improve the sampling process

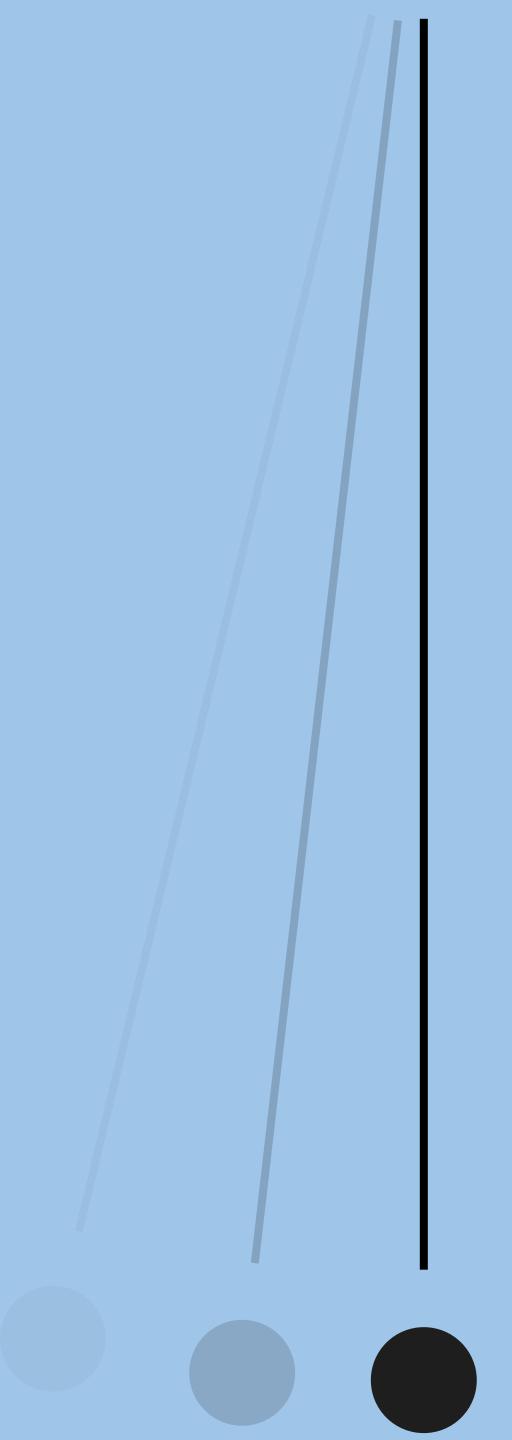
CEM:

```
for iter i = 1, 2, ...
    for population member e = 1, 2, ...
        sample  $\theta^{(e)} \sim P_{\mu^{(i)}}(\theta)$ 
        execute roll-outs under  $\pi_{\theta^{(e)}}$ 
        store  $(\theta^{(e)}, U(e))$ 
    endfor
     $\mu^{(i+1)} = \arg \max_{\mu} \sum_{\bar{e}} \log P_{\mu}(\theta^{(\bar{e})})$ 
    where  $\bar{e}$  indexes over top p %
endfor
```

# bottomline of gradient-free methods

- **easy** to implement
- good to have aside, as a toolbox for **comparison**
- not always suitable for **huge** parameter vectors
- don't **bootstrap**, so they need to know the full trajectory
- no particular assumptions, ie., **blackbox** functions, except fast evaluation of the rewards
- common sense states: when you have a **gradient**, exploit it! (do you agree?)

closing



# (some) research questions

**model-based** reinforcement learning

**efficient** reinforcement learning

**exploration** and intrinsic motivation

**hierarchical** reinforcement learning

# going further

- **Reinforcement Learning: An Introduction.** Richard S. Sutton & Andrew G. Barto  
<http://incompleteideas.net/book/the-book-2nd.html> (PDF freely available)
- **Algorithms for Reinforcement Learning.** Csaba Szepesvári  
<https://sites.ualberta.ca/~szepesva/RLBook.html> (PDF freely available)
- **Approximate Dynamic Programming: Solving the curses of dimensionality .** Warren B. Powell  
<http://adp.princeton.edu/>
- David Silver's excellent course (on video)  
<http://www0.cs.ucl.ac.uk/staff/d.silver/web/Teaching.html>
- Pieter Abbeel & John Schulman's tutorial @NIPS2016 <https://nips.cc/Conferences/2016/Schedule?showEvent=6198>

# libraries & frameworks

- **Showed during the talk**
  - javascript library [https://cs.stanford.edu/people/karpathy/reinforcejs/gridworld\\_dp.html](https://cs.stanford.edu/people/karpathy/reinforcejs/gridworld_dp.html)
- **OpenAI Gym** <https://gym.openai.com/>
- **Keras-RL** <https://github.com/matthiasplappert/keras-rl>
- **TensorForce** <https://reinforce.io/blog/introduction-to-tensorforce/>
- **OpenAI Baselines** <https://github.com/openai/baselines>
- **Coach** <http://coach.nervanasys.com/>

questions?