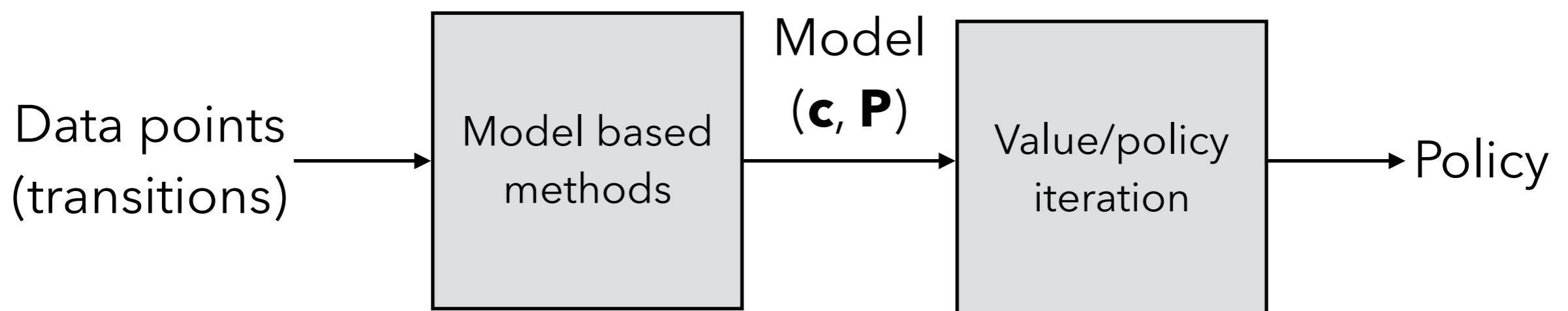


Planning, Learning and Decision Making

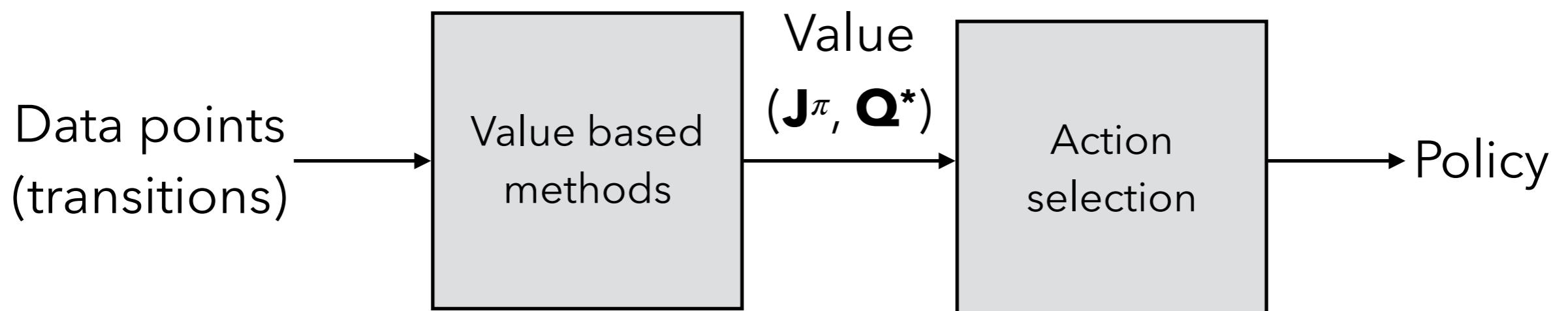
Lecture 20. Reinforcement learning: Policy-gradient

Model based RL



Value based RL

- Value-based methods:



Q-learning

- Given a sample (x_t, a_t, c_t, x_{t+1}) ,
- Compute

$$Q_{t+1}(x_t, a_t) = Q_t(x_t, a_t) + \alpha_t \left[c_t + \gamma \min_{a' \in \mathcal{A}} Q_t(x_{t+1}, a') - Q_t(x_t, a_t) \right]$$

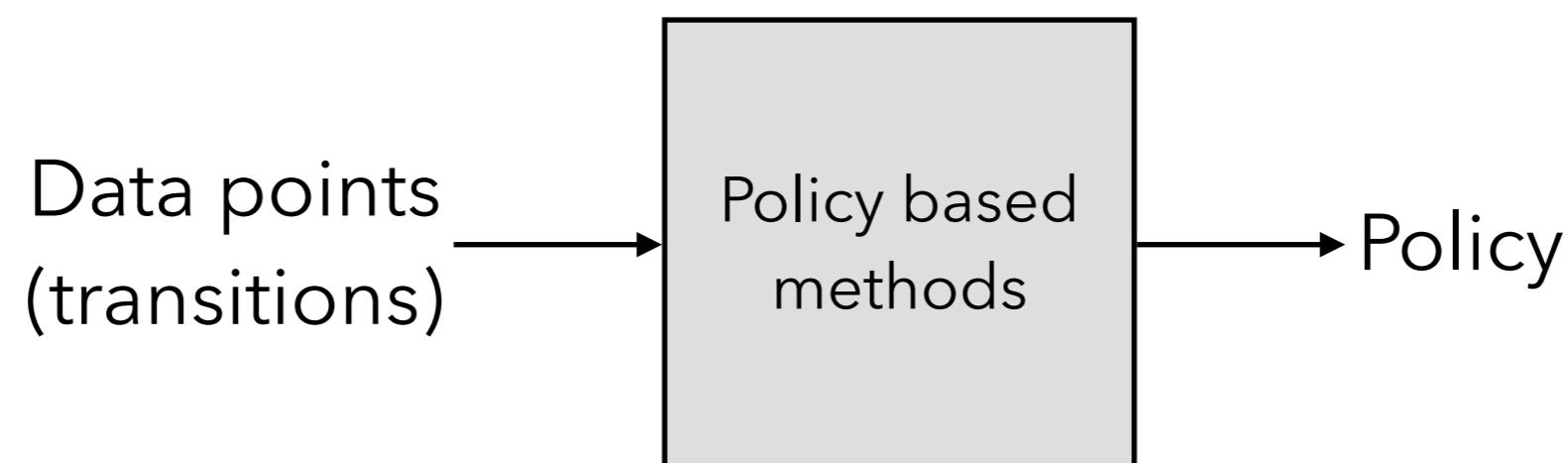
SARSA

- Given a sample $(x_t, a_t, c_t, x_{t+1}, a_{t+1})$,
- Compute

$$Q_{t+1}(x_t, a_t) = Q_t(x_t, a_t) + \alpha_t [c_t + \gamma Q_t(x_{t+1}, a_{t+1}) - Q_t(x_t, a_t)]$$

Policy-based RL

- Policy-based methods:





Policy-based methods

Why policy-based RL?

- As seen in previous classes,
 - ... value-based methods are not guaranteed to converge with approximations
 - ... quality of the result depends on the quality of approximation



Even if approximation is good,
policy may be bad

Example

- In some MDP, the Q-values in state x are

a_1	a_2
5	7

- Which is the optimal action?
 - a_1 - minimum cost

Example

- In some MDP, the Q-values in state x are

a_1	a_2
5	6

- Consider the two approximations

a_1	a_2
6	5.5

a_1	a_2
0	1

Which one has
smaller error?

Example

- In some MDP, the Q-values in state x are

a_1	a_2
5	6

- Consider the two approximations

a_1	a_2
6	5.5

a_1	a_2
0	1

Which one
is better?

Why policy-based RL?

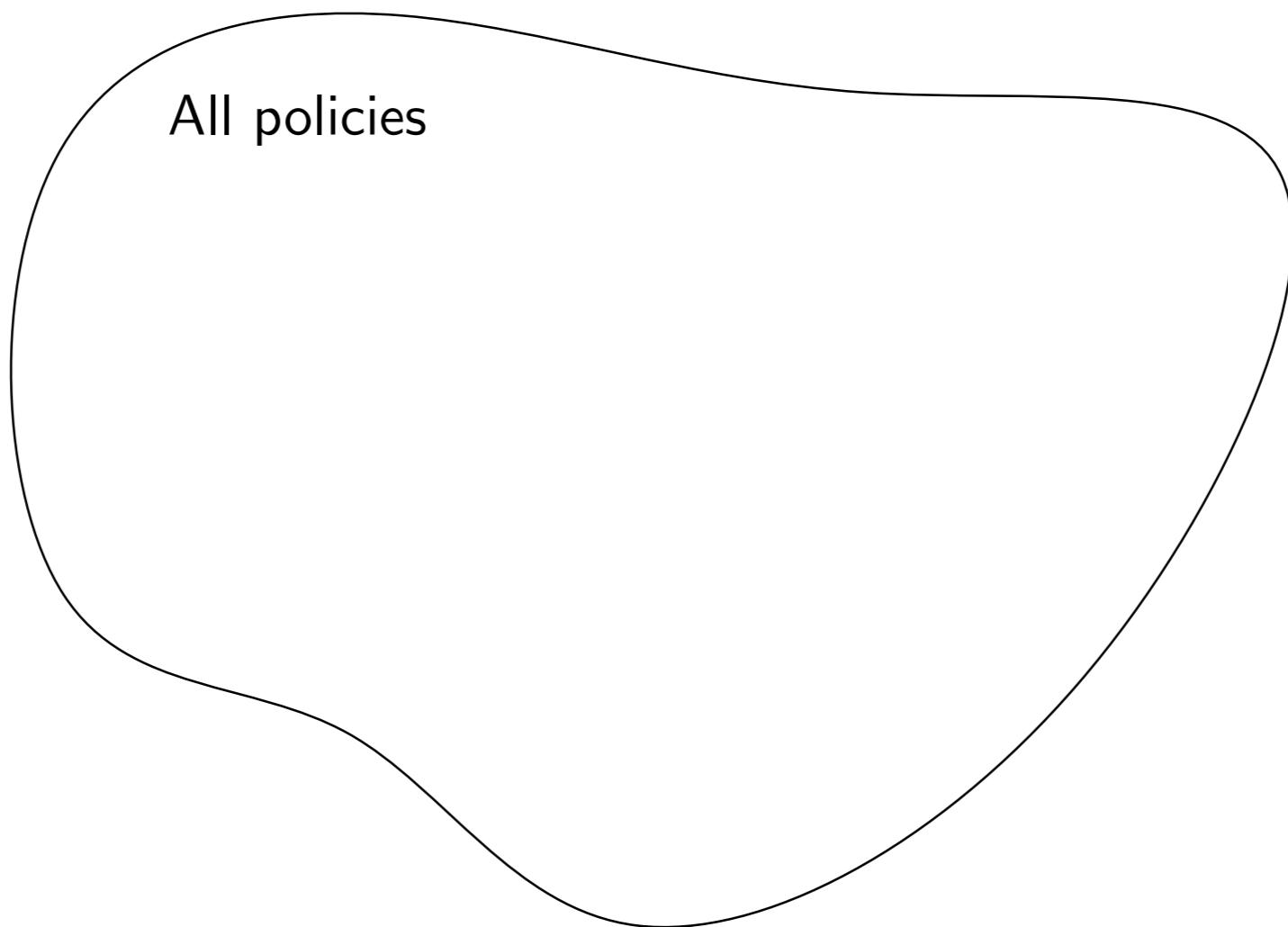
- As seen in previous classes,
 - ... value-based methods are not guaranteed to converge with approximations
 - ... quality of the result depends on the quality of approximation
 - **... even if the approximation is good, the resulting policy may be bad**

What is policy-based RL?

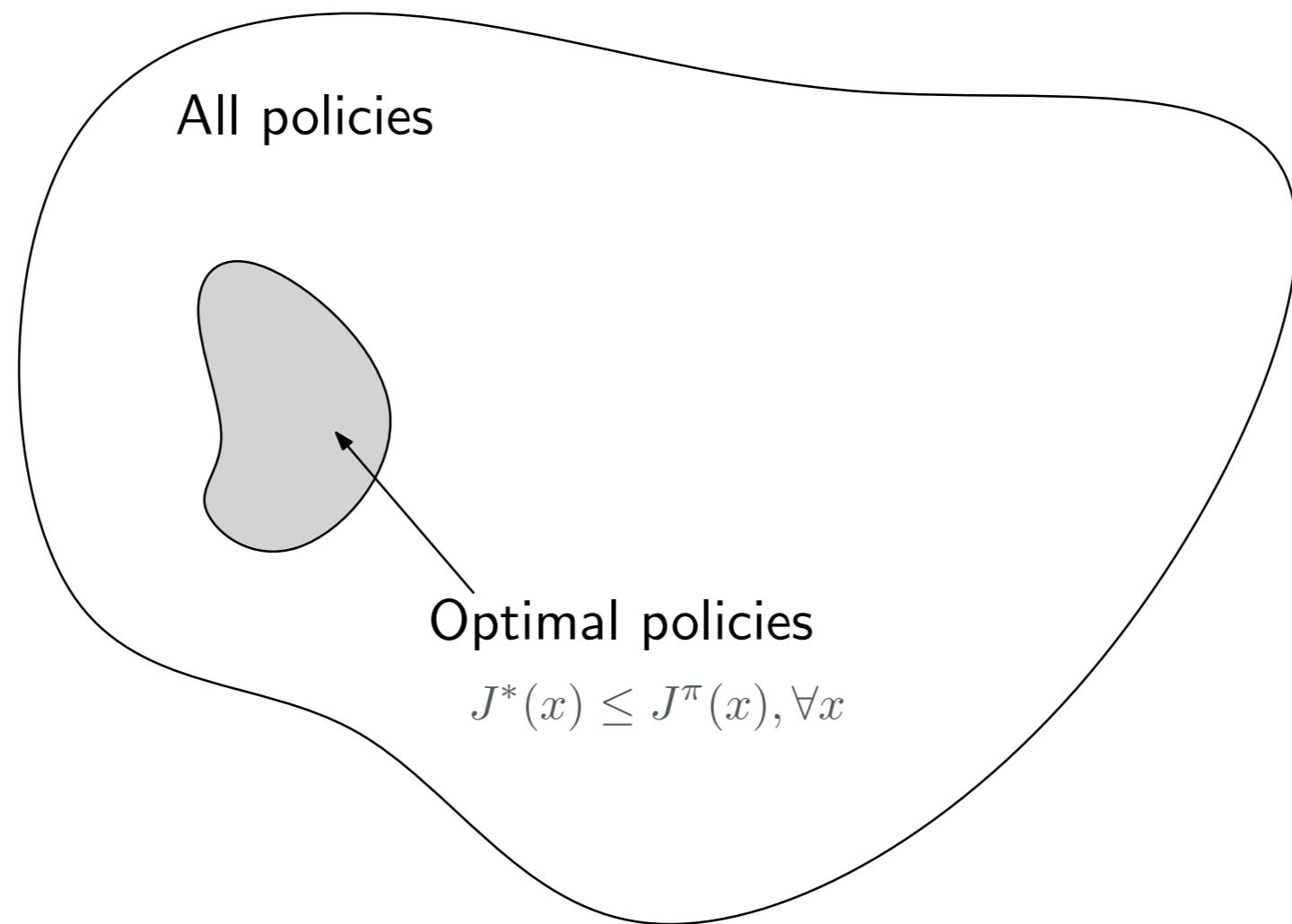
- Instead of selecting an approximation for the value, we select an approximation for the policy
- Policy-based methods select the best policy within the approximation

What is the “best policy”?

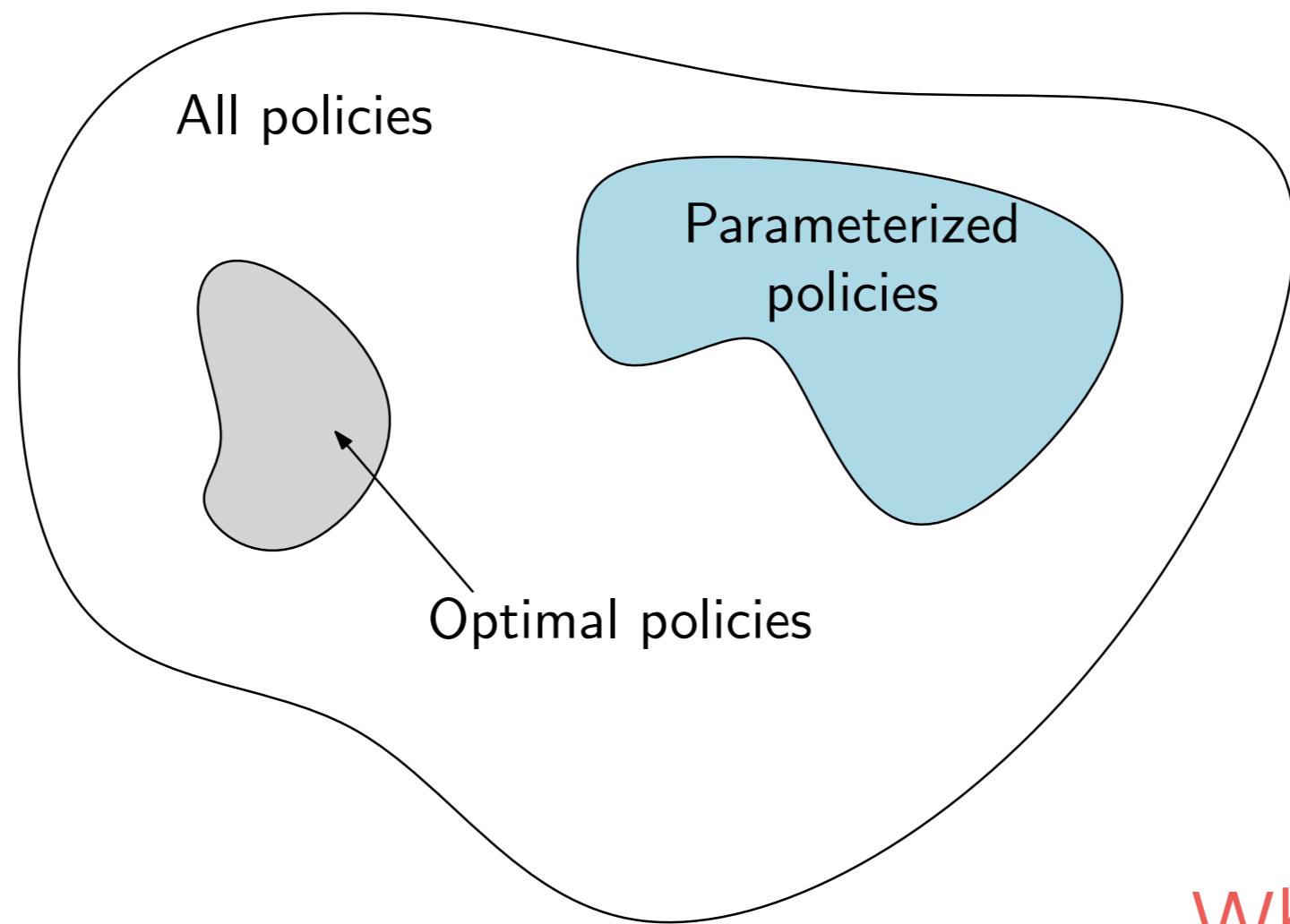
Revisiting optimality



Revisiting optimality



Revisiting optimality



What is “best”?

Example

- Suppose, that, for two policies π_1 and π_2

$$J^{\pi_1} = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \qquad J^{\pi_2} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

- Which one is “best”?

Initial distribution

- We consider the initial distribution μ_0 and set

$$V(\pi) = \mathbb{E}_{\mu_0} [J^\pi(x_0)]$$

$$= \sum_{x \in \mathcal{X}} \mu_0(x) J^\pi(x)$$



Single number

Initial distribution

- We consider the initial distribution μ_0 and set

$$\begin{aligned} V(\pi) &= \mathbb{E}_{\mu_0} [J^\pi(x_0)] \\ &= \sum_{x \in \mathcal{X}} \mu_0(x) J^\pi(x) \end{aligned}$$

- Then, π_1 is better than π_2 if $V(\pi_1) < V(\pi_2)$

Policy improvement

- Several approaches:
 - Hill climbing
 - Evolutionary algorithms
 - ...
 - We focus on **gradient descent**:
 - More efficient
 - Exploits MDP structure better
- BBO
(Black box optimization)

Policy gradient

- We consider a family of parameterized policies $\{\pi_\theta\}$
- θ is the parameter of the policy
- Example:

$$\pi_{\boldsymbol{\theta}}(a \mid x) = \frac{e^{\boldsymbol{\theta}^\top \phi(x, a)}}{\sum_{a' \in \mathcal{A}} e^{\boldsymbol{\theta}^\top \phi(x, a')}}$$

... many other possibilities

- V is a function of θ

How to compute the gradient?

Policy gradient

- **Approach 1:** Finite differences

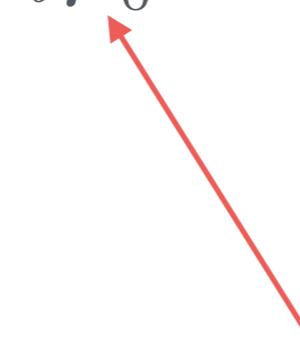
$$\frac{\partial V}{\partial \theta_n} \approx \frac{V(\theta + \varepsilon_n) - V(\theta)}{\varepsilon_n}$$

- Simple to compute
- Works even if policy is not differentiable
- Noisy and slow learning
- Requires evaluation of V

Policy gradient

- **Approach 2:** Analytical
 - Assume that we know $\nabla_{\theta}\pi$
 - Then,

$$\nabla_{\theta}V(\theta) = \nabla_{\theta}\mu_0 J^{\pi_{\theta}}$$

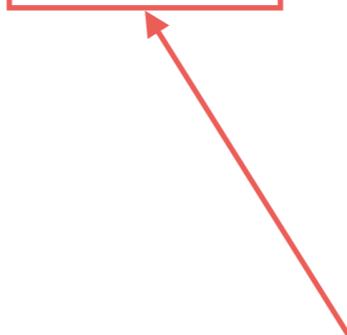


Does not
depend on θ

Policy gradient

- **Approach 2:** Analytical
 - Assume that we know $\nabla_{\theta}\pi$
 - Then,

$$\nabla_{\theta}V(\theta) = \mu_0 \boxed{\nabla_{\theta}J^{\pi_{\theta}}}$$

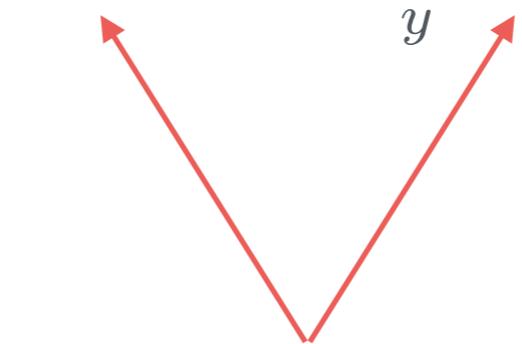


Let us consider
this term

Policy gradient

- Cost-to-go function verifies

$$\begin{aligned} J^{\pi_\theta}(x) &= \sum_a \pi_\theta(a \mid x) Q^\pi(x, a) \\ &= \sum_a \pi_\theta(a \mid x) \left[c(x, a) + \gamma \sum_y P(y \mid x, a) J^{\pi_\theta}(y) \right] \end{aligned}$$



Do not
depend on θ

Policy gradient

- Cost-to-go function verifies

$$\begin{aligned}
 J^{\pi_\theta}(x) &= \sum_a \pi_\theta(a \mid x) Q^\pi(x, a) \\
 &= \sum_a \pi_\theta(a \mid x) \left[c(x, a) + \gamma \sum_y P(y \mid x, a) J^{\pi_\theta}(y) \right]
 \end{aligned}$$

- Then (derivative of the product),

$$\nabla_\theta J^{\pi_\theta}(x) = \sum_a \nabla_\theta \pi_\theta(a \mid x) Q^{\pi_\theta}(x, a) + \gamma \sum_y P_{\pi_\theta}(y \mid x) \nabla_\theta J^{\pi_\theta}(y)$$

Linear
equation

Policy gradient

- We can solve it!

$$\nabla_{\theta} J^{\pi_{\theta}} = \nabla_{\theta} \pi_{\theta} Q^{\pi_{\theta}} + \gamma P_{\pi_{\theta}} \nabla_{\theta} J^{\pi_{\theta}}$$

Statewise
dot-product



Policy gradient

- We can solve it!

$$(I - \gamma P_{\pi_\theta}) \nabla_\theta J^{\pi_\theta} = \nabla_\theta \pi_\theta Q^{\pi_\theta}$$

Policy gradient

- We can solve it!

$$\nabla_{\theta} J^{\pi_{\theta}} = (\mathbf{I} - \gamma \mathbf{P}_{\pi_{\theta}})^{-1} \nabla_{\theta} \pi_{\theta} \mathbf{Q}^{\pi_{\theta}}$$

- We finally get:

$$\nabla_{\theta} V(\theta) = \boxed{\mu_0 (\mathbf{I} - \gamma \mathbf{P}_{\pi_{\theta}})^{-1} \nabla_{\theta} \pi_{\theta} \mathbf{Q}^{\pi_{\theta}}} \mu_{\theta}$$

Policy gradient theorem

- Continuing the computations, we get the following result.

Theorem: The policy gradient is given by

$$\nabla_{\theta} V(\theta) = \sum_x \mu_{\theta}(x) \sum_a \boxed{\nabla_{\theta} \pi_{\theta}(a | x)} Q^{\pi_{\theta}}(x, a).$$

where μ_{θ} is a long-term distribution over \mathcal{X} .

$$\nabla_{\theta} \pi_{\theta}(a | x) = \pi_{\theta}(a | x) \frac{\nabla_{\theta} \pi_{\theta}(a | x)}{\pi_{\theta}(a | x)}$$

Policy gradient theorem

- Continuing the computations, we get the following result.

Theorem: The policy gradient is given by

$$\nabla_{\theta} V(\theta) = \sum_x \mu_{\theta}(x) \sum_a \boxed{\nabla_{\theta} \pi_{\theta}(a | x)} Q^{\pi_{\theta}}(x, a).$$

where μ_{θ} is a long-term distribution over \mathcal{X} .

$$\nabla_{\theta} \pi_{\theta}(a | x) = \pi_{\theta}(a | x) \nabla_{\theta} \log \pi_{\theta}(a | x)$$

Policy gradient theorem

- Continuing the computations, we get the following result.

Theorem: The policy gradient is given by

$$\nabla_{\theta} V(\theta) = \mathbb{E}_{\mu_{\theta}, \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a | x) Q^{\pi_{\theta}}(x, a)].$$

where μ_{θ} is a long-term distribution over \mathcal{X} .



Must still
compute Q^{π}

REINFORCE Algorithm

- Initialize θ
- Collect a trajectory $\{x_0, a_0, c_0, x_1, a_1, c_1, \dots, x_T, a_T, c_T\}$ using π_θ
- For each $t = 0, \dots, T$
 - Update $\theta \leftarrow \theta - \alpha \nabla \log \pi(a_t | x_t) \sum_{\tau=t}^T \gamma^{\tau-t} c_\tau$

$$\sum_{\tau=t}^T \gamma^{\tau-t} c_\tau$$

Estimate of $\gamma^t Q^\pi$

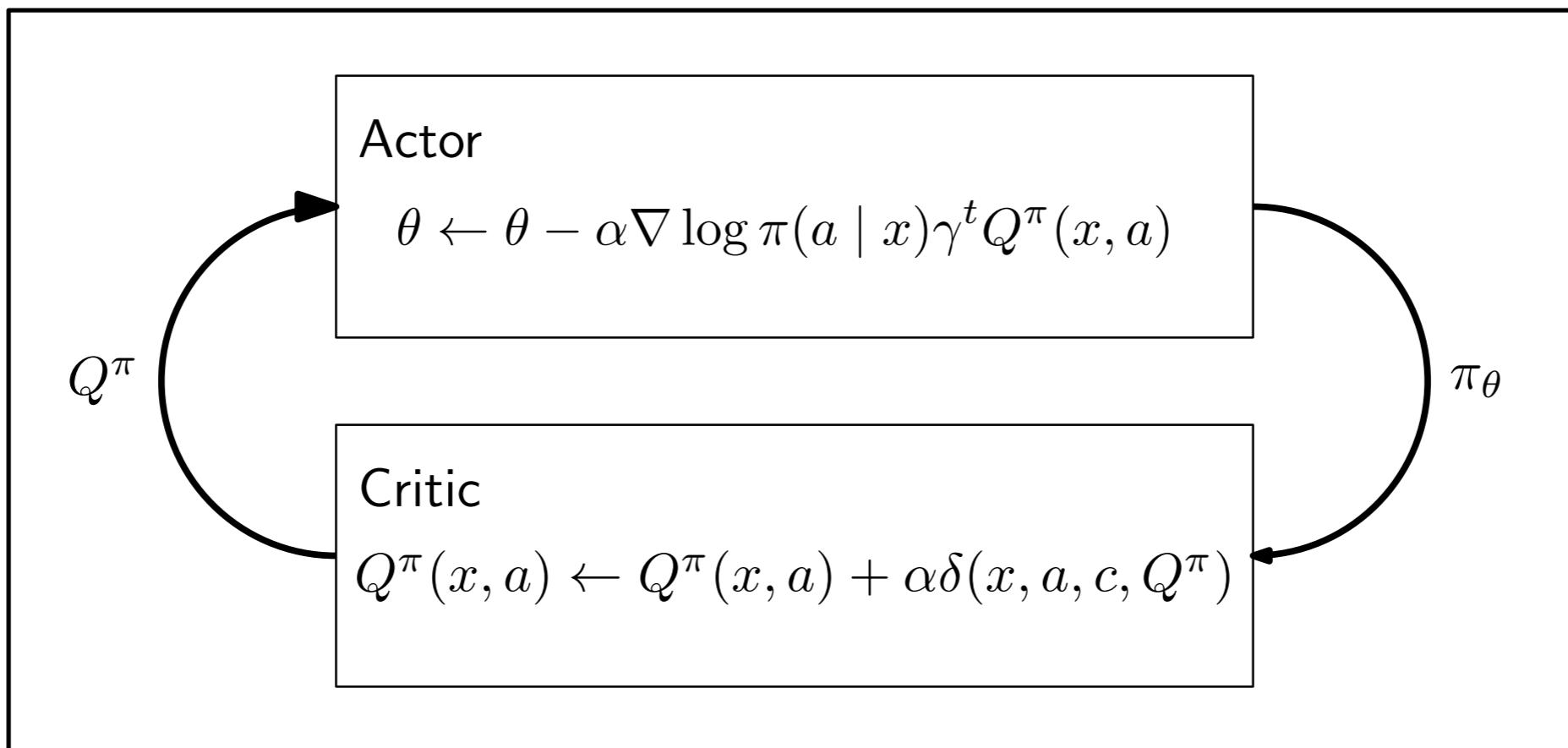
REINFORCE Algorithm

- The REINFORCE algorithm is significantly better than the finite difference approach
- However, gradient estimates have large variance
 - Why not use better algorithm to estimate Q^π ?

Actor-critic architecture

- Split algorithm in two parts:
 - One part maintains/runs the policy (actor)
 - One part computes Q (critic)

Actor-critic architecture



Actor-critic architecture

- However,
 - ... we need to estimate Q^π
 - ... in large domains, we are back to the problems of selecting good approximations for Q^π

... are we?

Policy gradient revisited

- The formula for the gradient gives the answer:

$$\nabla_{\theta} V(\theta) = \mathbb{E}_{\mu_{\theta}, \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a | x) Q^{\pi_{\theta}}(x, a)] .$$

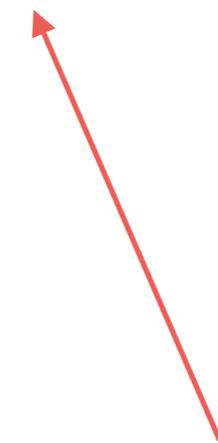


Don't really
need Q^{π}

Policy gradient revisited

- The formula for the gradient gives the answer:

$$\nabla_{\theta} V(\theta) = \mathbb{E}_{\mu_{\theta}, \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a | x) Q^{\pi_{\theta}}(x, a)] .$$



Inner
product

Compatible features

- We let

$$\phi(x, a) = \nabla_\theta \log \pi_\theta(a \mid x)$$

- The functions ϕ are called **compatible features**
- If we approximate Q^π using the features ϕ , we recover the gradient

Policy gradient theorem II

Theorem: Given the compatible basis functions

$$\phi(x, a) = \nabla_{\theta} \log \pi_{\theta}(a \mid x),$$

let

$$Q_{\mathbf{w}}(x, a) = \mathbf{w}^{\top} \phi(x, a),$$

with

$$\mathbf{w} = \operatorname{argmin}_{\mu_{\theta}, \pi_{\theta}} \mathbb{E}_{\mu_{\theta}, \pi_{\theta}} \left[(Q^{\pi_{\theta}}(x, a) - Q_{\mathbf{w}}(x, a))^2 \right]$$

Then,

$$\nabla_{\theta} V(\theta) = \mathbb{E}_{\mu_{\theta}, \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a \mid x) Q_{\mathbf{w}}(x, a)].$$

Actor-critic methods

- They provide satisfactory solutions both for the actor and the critic
- Much used in robotics
- Many variations exist:
 - Critic using batch learning
 - Actor using adjusted/natural gradients
 - Multi time scales
 - ...

Examples in practice

- Gait learning: initial policy



Examples in practice

- Gait learning: final policy



Examples in practice

- Learning to grasp:



Examples in practice

- Learning to intercept a ping-pong ball:

