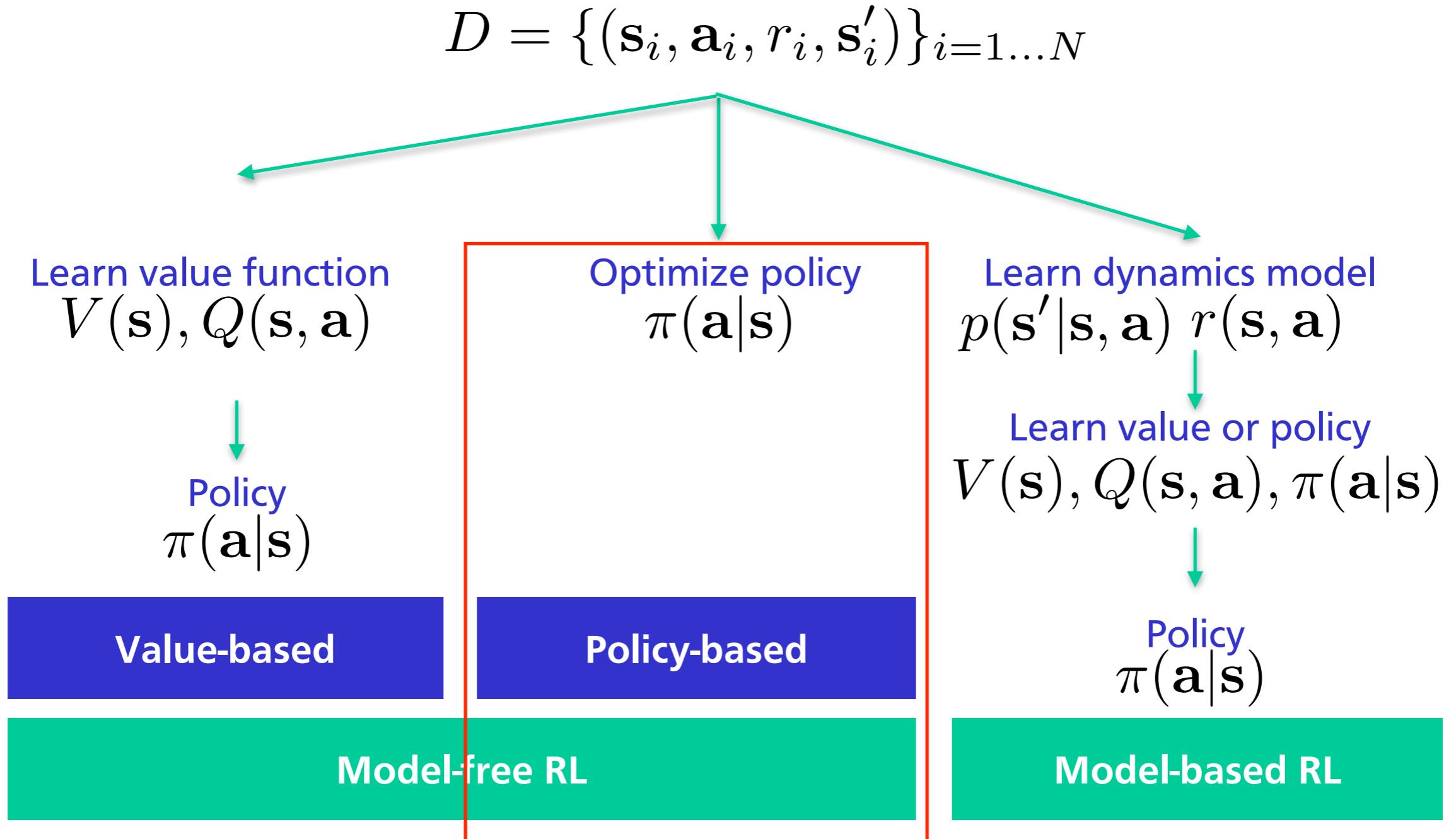


POLICY SEARCH METHODS 2

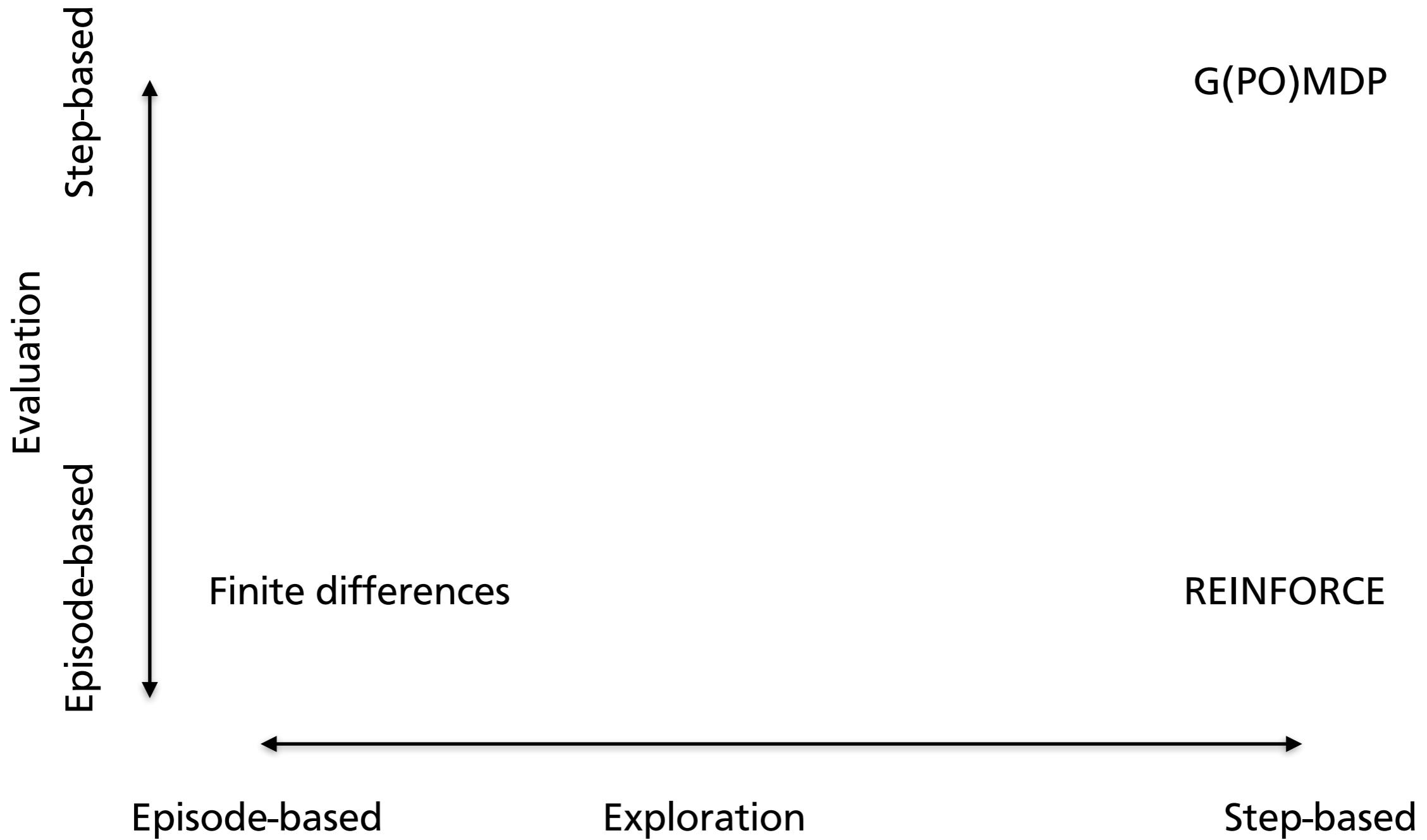
Herke van Hoof

Big picture: How to learn policies



Thanks to Jan Peters

Big picture: How to learn policies



Policy Gradient Theorem

By shuffling the terms around, we can reformulate Reinforce / G(PO)MDP

$$\begin{aligned} & \mathbb{E}_\tau \left[\sum_{t=1}^T r_t \sum_{t'=1}^t \nabla \log p(\mathbf{a}_{t'} | \mathbf{s}_{t'}) \right] \\ &= \mathbb{E}_\tau \left[\sum_{t'=1}^T \nabla \log p(\mathbf{a}_{t'} | \mathbf{s}_{t'}) \sum_{t=t'}^T r_t \right] \end{aligned}$$

	a_1	a_2	a_3
r_1	→		
r_2	→	→	
r_3	→	→	→

	a_1	a_2	a_3
r_1	↓		
r_2	↓	↓	
r_3	↓	↓	↓

Policy Gradient Theorem

By shuffling the terms around, we can reformulate Reinforce / G(PO)MDP

$$\mathbb{E}_\tau \left[\sum_{t=1}^T r_t \sum_{t'=1}^t \nabla \log p(\mathbf{a}_{t'} | \mathbf{s}_{t'}) \right]$$

$$= \mathbb{E}_\tau \left[\sum_{t'=1}^T \nabla \log p(\mathbf{a}_{t'} | \mathbf{s}_{t'}) \sum_{t=t'}^T r_t \right]$$

value func q_{π} is defined as (or should)
= Expectation of return.

Replace return by its expected

sum of imm reward = true value
 $q_{\pi}(s, a)$ =estimated value func, it brings some error

$$= \mathbb{E}_\tau \left[\sum_{t'=1}^T \nabla \log p(\mathbf{a}_{t'} | \mathbf{s}_{t'}) q_{\pi}(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \right]$$

	a_1	a_2	a_3
r_1	→		
r_2	→	→	
r_3	→	→	→

3

Policy Gradient Theorem

$$\nabla J = \mathbb{E}_\tau \left[\sum_{t'=1}^T \nabla \log \pi(\mathbf{a}_{t'} | \mathbf{s}_{t'}) q_\pi(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \right]$$

State-action pairs at each time step will contribute equally to the total gradient. So this is equivalent to an expectation over the on-pol

what is the advantage that we replace true reward as estimated reward:
it is good for continuing case: when the episode is really long, when the episode does not terminate.

$$\nabla J = \mathbb{E}_{\mu(\mathbf{s})\pi(\mathbf{a}|\mathbf{s})} \nabla \log \pi(\mathbf{a}|\mathbf{s}) q_\pi(\mathbf{s}, \mathbf{a})$$

Formal proof, also for the continuing case:

Sutton et al., Policy Gradient Methods for Reinforcement Learning with Function Approximation

PGT Actor Critic

has bias
reduce the variance
compared with actor only policy based method

$$\nabla J = \mathbb{E}_{\mu(s)\pi(a|s)} \nabla \log \pi(a|s) q_\pi(s, a)$$

Maybe we can replace q_π by an estimate (q_w or $v_w(s_{t+1}) + r_{t+1}$)

Plug in samples to estimate the expected value

$$\theta_{t+1} = \theta_t + \alpha (R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w})) \nabla \log \pi(\mathbf{a}_t | \mathbf{s}_t, \theta_t)$$

Advantage: break the high variance of Monte-Carlo returns

Disadvantage: possibly include bias?

PGT Actor Critic

$$\theta_{t+1} = \theta_t + \alpha (R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w})) \nabla \log \pi(\mathbf{a}_t | \mathbf{s}_t, \theta_t)$$

This algorithm has a parametrised policy *and* a parametrised action-value function, with parameters θ and w , respectively.

The policy is often called an actor while the value function is called critic. A method that uses both is an *actor-critic* method.

PGT Actor Critic

it is biased in this case. because it always converge to the target. the target is an estimated return, not true return, so it is always biased. but it has low variance. because we dont do MC

Note the difference between a critic and a baseline

Reinforce w baseline

$$\theta_{t+1} = \theta_t + \alpha (G_t - \hat{v})(\mathbf{s}_t, \mathbf{w}) \nabla \log \pi(\mathbf{a}_t | \mathbf{s}_t, \theta_t)$$

(Removes variance by centering targets, always unbiased)

Actor-Critic

$$\theta_{t+1} = \theta_t + \alpha (R_{t+1} + \gamma \hat{v}(\mathbf{s}_{t+1}, \mathbf{w})) \nabla \log \pi(\mathbf{a}_t | \mathbf{s}_t, \theta_t)$$

(Removes variance by removing long-term dependencies. Bias?)

Actor-Critic + baseline

$$\begin{aligned} \theta_{t+1} &= \theta_t + \alpha (R_{t+1} + \gamma \hat{v}(\mathbf{s}_{t+1}, \mathbf{w}) - \hat{v}(\mathbf{s}_t, \mathbf{w})) \nabla \log \pi(\mathbf{a}_t | \mathbf{s}_t, \theta_t) \\ &= \delta \end{aligned}$$

TD error
= delta

PGT Actor critic

One-step Actor–Critic (episodic), for estimating $\pi_\theta \approx \pi_*$

Input: a differentiable policy parameterization $\pi(a|s, \theta)$

Input: a differentiable state-value function parameterization $\hat{v}(s, w)$

Parameters: step sizes $\alpha^\theta > 0$, $\alpha^w > 0$

Initialize policy parameter $\theta \in \mathbb{R}^{d'}$ and state-value weights $w \in \mathbb{R}^d$ (e.g., to 0)

Loop forever (for each episode):

 Initialize S (first state of episode)

$I \leftarrow 1$

 Loop while S is not terminal (for each time step):

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

 Take action A , observe S', R

$\delta \leftarrow R + \gamma \hat{v}(S', w) - \hat{v}(S, w)$ (if S' is terminal, then $\hat{v}(S', w) \doteq 0$)

$w \leftarrow w + \alpha^w \delta \nabla \hat{v}(S, w)$

$\theta \leftarrow \theta + \alpha^\theta I \delta \nabla \ln \pi(A|S, \theta)$ With baseline!

$I \leftarrow \gamma I$

$S \leftarrow S'$

If we discount, future time steps
less important

Compatible function approximation thm.

Disadvantage of actor-critic: possibly include bias?

Question: are there functions q_π that give an unbiased estimate?

$$\nabla J = \mathbb{E}_{\mu(\mathbf{s}), a} [\nabla \log \pi(\mathbf{a}_{t'} | \mathbf{s}_{t'}) q_{\pi}(\mathbf{s}_{t'}, \mathbf{a}_{t'})] \stackrel{?}{=} \mathbb{E}_{\mu(\mathbf{s}), a} [\nabla \log \pi(\mathbf{a}_{t'} | \mathbf{s}_{t'}) \hat{q}_{\mathbf{w}}(\mathbf{s}_{t'}, \mathbf{a}_{t'})]$$

[Sutton 2000]

Compatible function approximation thm.

Disadvantage of actor-critic: possibly include bias?

Question: are there functions q_π that give an unbiased estimate?

$$\nabla J = \mathbb{E}_{\mu(\mathbf{s}), a} [\nabla \log \pi(\mathbf{a}_{t'} | \mathbf{s}_{t'}) q_{\pi}(\mathbf{s}_{t'}, \mathbf{a}_{t'})] \stackrel{?}{=} \mathbb{E}_{\mu(\mathbf{s}), a} [\nabla \log \pi(\mathbf{a}_{t'} | \mathbf{s}_{t'}) \hat{q}_{\mathbf{w}}(\mathbf{s}_{t'}, \mathbf{a}_{t'})]$$

Answer: Yes, if value function is 'compatible':

$$\nabla_{\mathbf{w}} \hat{q}_{\mathbf{w}} = \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{s}, \mathbf{a}) \quad \text{e.g. } \hat{q}_{\mathbf{w}} = \mathbf{w}^T \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{s}, \mathbf{a})$$

and

$$\mathbb{E}_{\mathbf{s} \sim \mu, \mathbf{a} \sim \pi(\cdot, \mathbf{s})} [(q_{\pi}(\mathbf{s}, \mathbf{a}) - \hat{q}_{\mathbf{w}}(\mathbf{s}, \mathbf{a})) \nabla_{\mathbf{w}} \hat{q}_{\mathbf{w}}(\mathbf{s}, \mathbf{a})] = 0$$

$$(\text{e.g. } \mathbf{w} = \arg \min_{\mathbf{w}} \mathbb{E}_{\mathbf{s} \sim \mu, \mathbf{a} \sim \pi(\cdot, \mathbf{s})} (q_{\pi}(\mathbf{s}, \mathbf{a}) - \hat{q}_{\mathbf{w}}(\mathbf{s}, \mathbf{a}))^2)$$

Compatible function approximation

Proof:

$$1) \quad \nabla_{\mathbf{w}} \hat{q}_{\mathbf{w}} = \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{s}, \mathbf{a})$$

$$2) \quad \mathbb{E}_{\mathbf{s} \sim \mu, \mathbf{a} \sim \pi(\cdot, \mathbf{s})} [(q_{\pi}(\mathbf{s}, \mathbf{a}) - \hat{q}_{\mathbf{w}}(\mathbf{s}, \mathbf{a})) \nabla_{\mathbf{w}} \hat{q}_{\mathbf{w}}(\mathbf{s}, \mathbf{a})] = 0$$

Substitute (1) into (2)

$$\mathbb{E}_{\mu(\mathbf{s}), \mathbf{a}} [(q_{\pi}(\mathbf{s}, \mathbf{a}) - \hat{q}_{\mathbf{w}}(\mathbf{s}, \mathbf{a})) \nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{a}|\mathbf{s})] = 0$$



[Sutton 2000]

Compatible function approximation

Proof:

$$1) \quad \nabla_{\mathbf{w}} \hat{q}_{\mathbf{w}} = \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{s}, \mathbf{a})$$

$$2) \quad \mathbb{E}_{\mathbf{s} \sim \mu, \mathbf{a} \sim \pi(\cdot, \mathbf{s})} [(q_{\pi}(\mathbf{s}, \mathbf{a}) - \hat{q}_{\mathbf{w}}(\mathbf{s}, \mathbf{a})) \nabla_{\mathbf{w}} \hat{q}_{\mathbf{w}}(\mathbf{s}, \mathbf{a})] = 0$$

Substitute (1) into (2)

$$\mathbb{E}_{\mu(\mathbf{s}), \mathbf{a}} [(q_{\pi}(\mathbf{s}, \mathbf{a}) - \hat{q}_{\mathbf{w}}(\mathbf{s}, \mathbf{a})) \nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{a}|\mathbf{s})] = 0$$

$$\mathbb{E}_{\mu(\mathbf{s}), \mathbf{a}} [q_{\pi}(\mathbf{s}, \mathbf{a}) \nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{a}|\mathbf{s})] = \mathbb{E}_{\mu(\mathbf{s}), \mathbf{a}} [\hat{q}_{\mathbf{w}}(\mathbf{s}, \mathbf{a}) \nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{a}|\mathbf{s})]$$

$$\xrightarrow{} = \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$$

Again yielding a **unbiased and consistent estimate**

This can again be combined with a **baseline**

Compatible function approximation

Compatibility:

$$\nabla_{\mathbf{w}} \hat{q}_{\mathbf{w}} = \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{s}, \mathbf{a}) \quad \text{e.g. } \hat{q}_{\mathbf{w}} = \mathbf{w}^T \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{s}, \mathbf{a})$$

Example: Softmax policy

$$\pi(a|s) = \frac{e^{\boldsymbol{\theta}^T \phi_{sa}}}{\sum_b e^{\boldsymbol{\theta}^T \phi_{sb}}}$$

Compatible function approximation

Compatibility:

$$\nabla_{\mathbf{w}} \hat{q}_{\mathbf{w}} = \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{s}, \mathbf{a}) \quad \text{e.g. } \hat{q}_{\mathbf{w}} = \mathbf{w}^T \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{s}, \mathbf{a})$$

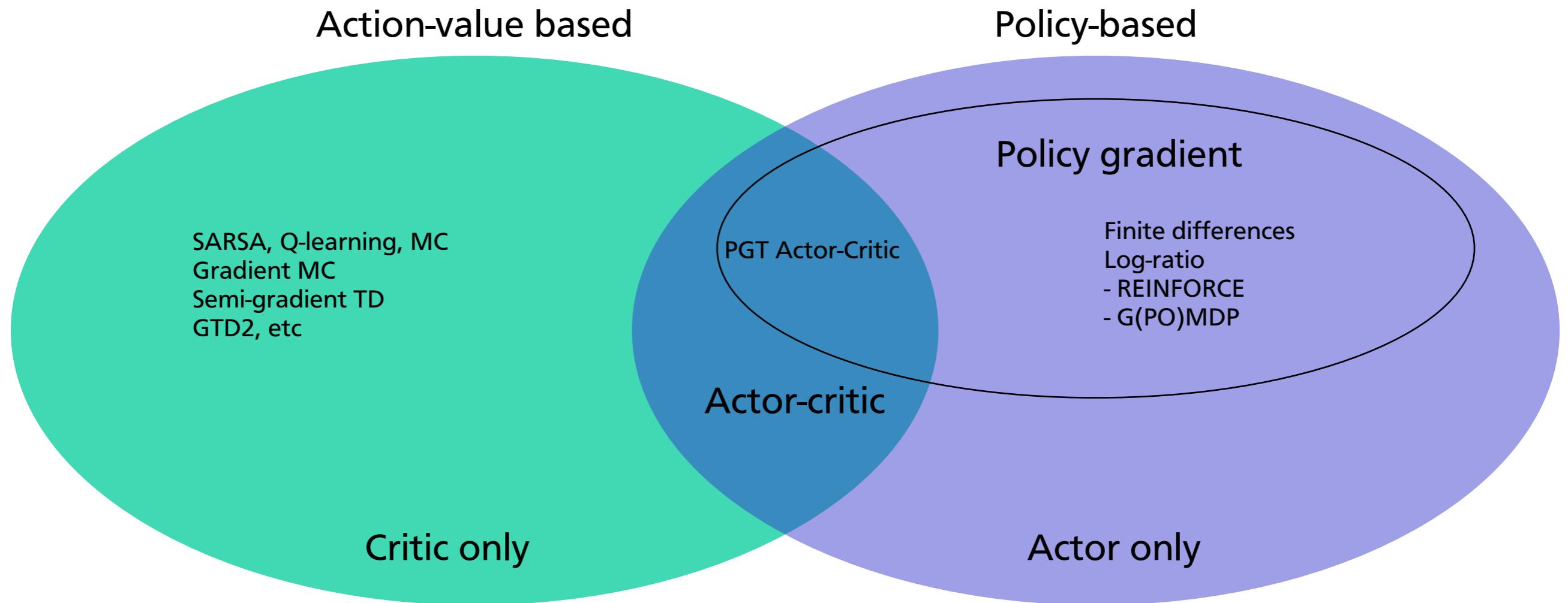
Example: Softmax policy

$$\pi(a|s) = \frac{e^{\boldsymbol{\theta}^T \phi_{sa}}}{\sum_b e^{\boldsymbol{\theta}^T \phi_{sb}}}$$

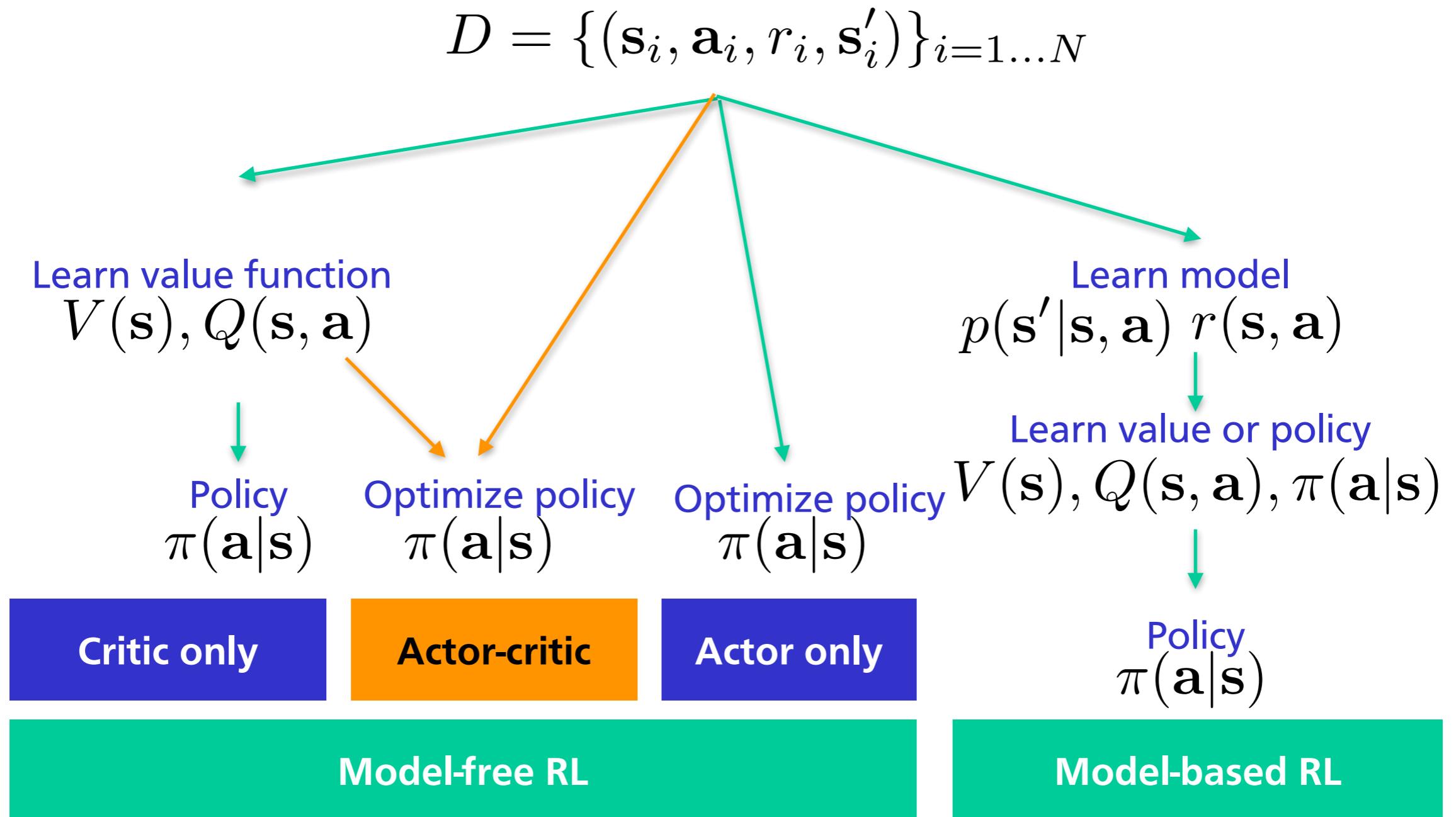
$$\begin{aligned} \nabla \log \frac{e^{\boldsymbol{\theta}^T \phi_{sa}}}{\sum_b e^{\boldsymbol{\theta}^T \phi_{sb}}} &= \nabla \boldsymbol{\theta}^T \phi_{sa} - \nabla \log \sum_b e^{\boldsymbol{\theta}^T \phi_{sb}} \\ &= \phi_{sa} - \frac{1}{\sum_b e^{\boldsymbol{\theta}^T \phi_{sb}}} \sum_b e^{\boldsymbol{\theta}^T \phi_{sb}} \phi_{sb} &= \phi_{sa} - \underbrace{\sum_b \frac{e^{\boldsymbol{\theta}^T \phi_{sb}}}{\sum_b e^{\boldsymbol{\theta}^T \phi_{sb}}} \phi_{sb}}_{=\pi(b|s)} \end{aligned}$$

$$\hat{q}_{\mathbf{w}}(s, a) = \mathbf{w}^T \left(\phi_{sa} - \sum_b \pi(b|s) \phi_{sb} \right)$$

Policies and action-values

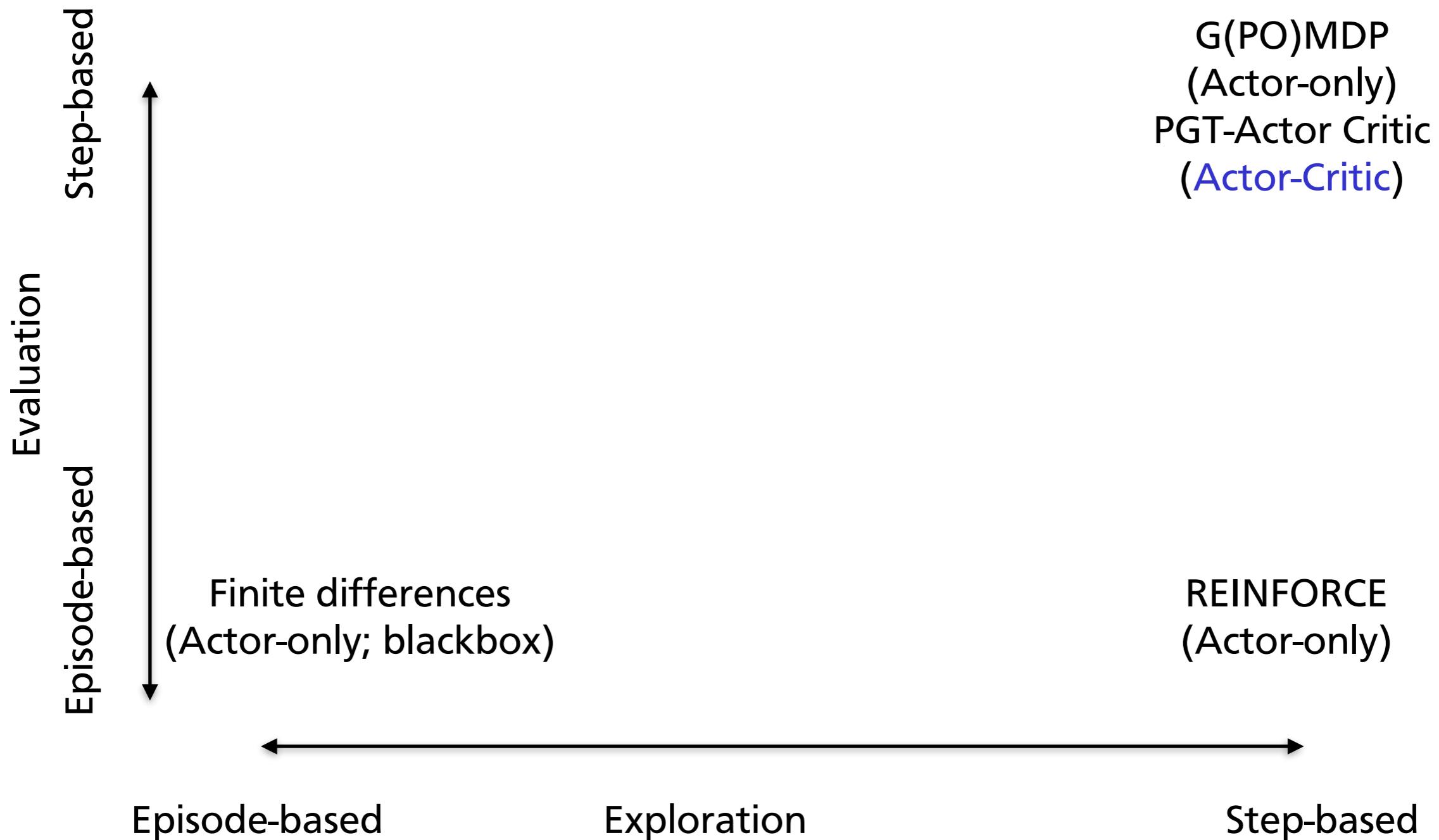


Big picture: How to learn policies



Thanks to Jan Peters

Comparison



Before we look at more advanced policy search methods, let's look at when we can use policy search methods

Advantages of policy-based methods

Policy search methods typically preferred in any of the below cases:

Problems with continuous actions

If we need to learn stochastic policies

If we have prior knowledge about the type of policy

If it is important to have ‘small’ policy updates between subsequent time steps

Advantages of policy-based methods

Policy search methods typically preferred in any of the below cases:

Problems with continuous actions

If we need to learn stochastic policies

If we have prior knowledge about the type of policy

If it is important to have ‘small’ policy updates between subsequent time steps

Many of these aspects present with physical systems like robots...

Actor critic addresses some weaknesses

Actor-only methods have high variance from Monte-Carlo

- Actor-critic lowers variance using critic

A lot of the methods we discussed are specific to episodic setting

- Actor-critic can be formulated for continuing setting

Actor-critic can be ‘fiddly’, many moving parts

Requires stochastic policies, what if deterministic is optimal?

- If amount of randomness is learned, can get close to deterministic
- We will also see a policy gradient method to learn deterministic policies

A CLOSER LOOK AT TARGET VALUES

- The target values are the part that differs most between policy gradient methods
- ‘Vanilla’ policy gradients can be written as:

$$g = \mathbb{E} \left[\sum_{t=0}^{\infty} \Psi_t \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right]$$

- With different choices for Ψ , e.g:

- G(PO)MPD:
$$\Psi_t = \sum_{t'=t}^{\infty} r_{t'}$$

- PGT-AC:
$$\Psi_t = q_{\pi}(s_t, a_t)$$

A CLOSER LOOK AT TARGET VALUES

- The target values are the part that differs most between policy gradient methods
- ‘Vanilla’ policy gradients can be written as:

$$g = \mathbb{E} \left[\sum_{t=0}^{\infty} \Psi_t \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right]$$

- With different choices for Ψ , e.g:

➤ G(PO)MPD: $\Psi_t = \sum_{t'=t}^{\infty} r_{t'}$

➤ PGT-AC: $\Psi_t = q_{\pi}(s_t, a_t)$

Less variance with value function baseline:

$$\Psi_t = \sum_{t'=t}^{\infty} r_{t'} - V(s_t)$$

$$\Psi_t = q_{\pi}(s_t, a_t) - V(s_t)$$

A CLOSER LOOK AT TARGET VALUES

- Different target values have their own problems:

- $$\Psi_t = \sum_{t'=t}^{\infty} r_{t'} - V(\mathbf{s}_t)$$

Although the baseline reduces variance somewhat,
REINFORCE still has a very high variance

- $$\Psi_t = q_{\pi}(\mathbf{s}_t, \mathbf{a}_t) - V(\mathbf{s}_t)$$

An approximation of the q-function might introduce bias
(except when using compatible functions and Q has converged)

- What does this remind you of?

REMEMBER N-STEP METHODS?

- We can also use n-step returns!
- Let's first define the **target value** as the **advantage**

$$A(s_t, a_t) = Q(s_t, a_t) - V(s_t)$$

advanced estimator of target reward
= sum over TD errors
=subtract baseline + TD(n) estimator

- Now an n-step advantage estimator is given by

$$\hat{A}_t^n = \sum_{l=0}^{n-1} \delta_{t+l} = -V(s_t) + r_t + \gamma r_{t+1} + \cdots + \gamma^{n-1} r_{t+n-1} + \gamma^n V(s_{t+n})$$

if n is large, then we update more states in the past. but variance is higher
if n is small, then reward propagate slower. but lower variance

- Again (like with n-)

- High-n: update more states +variance

????? 1h 29min

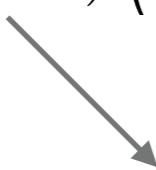
- Low-n: only 1 step back-up, bad when V is bad +bias

REMEMBER N-STEP METHODS?

- Rather than choosing 1 value of n, we can combine many:

$$\hat{A}_t^{\text{GAE}(\gamma, \lambda)} := (1 - \lambda) \left(\hat{A}_t^{(1)} + \lambda \hat{A}_t^{(2)} + \lambda^2 \hat{A}_t^{(3)} + \dots \right)$$

Why 1-λ?



REMEMBER N-STEP METHODS?

- Rather than choosing 1 value of n, we can combine many:

$$\hat{A}_t^{\text{GAE}(\gamma, \lambda)} := (1 - \lambda) \left(\hat{A}_t^{(1)} + \lambda \hat{A}_t^{(2)} + \lambda^2 \hat{A}_t^{(3)} + \dots \right)$$

=

=

=

=

REMEMBER N-STEP METHODS?

- Rather than choosing 1 value of n, we can combine many:

$$\begin{aligned}\hat{A}_t^{\text{GAE}(\gamma, \lambda)} &:= (1 - \lambda) \left(\hat{A}_t^{(1)} + \lambda \hat{A}_t^{(2)} + \lambda^2 \hat{A}_t^{(3)} + \dots \right) \\ &= (1 - \lambda) \left(\delta_t^V + \lambda (\delta_t^V + \gamma \delta_{t+1}^V) + \lambda^2 (\delta_t^V + \gamma \delta_{t+1}^V + \gamma^2 \delta_{t+2}^V) + \dots \right) \\ &= (1 - \lambda) \left(\delta_t^V (1 + \lambda + \lambda^2 + \dots) + \gamma \delta_{t+1}^V (\lambda + \lambda^2 + \lambda^3 + \dots) \right. \\ &\quad \left. + \gamma^2 \delta_{t+2}^V (\lambda^2 + \lambda^3 + \lambda^4 + \dots) + \dots \right) \\ &= (1 - \lambda) \left(\delta_t^V \left(\frac{1}{1 - \lambda} \right) + \gamma \delta_{t+1}^V \left(\frac{\lambda}{1 - \lambda} \right) + \gamma^2 \delta_{t+2}^V \left(\frac{\lambda^2}{1 - \lambda} \right) + \dots \right) \\ &= \sum_{l=0}^{\infty} (\gamma \lambda)^l \delta_{t+l}^V\end{aligned}$$

skip推导过程。重要的是结果：sum of discounted TD errors
gamma: discounting factor

(Generalized Advantage Estimates, Schulman [2016])

REMEMBER N-STEP METHODS?

- Sidenote:
- A similar trick can also be used to learn value functions
- In that case, the trick is called $\text{TD}(\lambda)$
- 1-step TD is equal to $\text{TD}(\lambda)$ when choosing $\lambda=0$
(in that case longer returns have a weight of 0).
1-step TD is also called $\text{TD}(0)$
- $\text{TD}(\lambda)$ can be implemented without storing the whole episode
by using a trick called “eligibility traces”.
(won’t be covered in this course...)

this trace is not covered in this course

REMEMBER N-STEP METHODS?

- 2 special cases:

$$\hat{A}_t^{\text{GAE}(\gamma, \lambda)} := (1 - \lambda) \left(\hat{A}_t^{(1)} + \lambda \hat{A}_t^{(2)} + \lambda^2 \hat{A}_t^{(3)} + \dots \right) = \sum_{l=0}^{\infty} (\gamma \lambda)^l \delta_{t+l}^V$$

- GAE($\gamma, 0$) :

GAE($\gamma, 1$) :

REMEMBER N-STEP METHODS?

- 2 special cases:

$$\hat{A}_t^{\text{GAE}(\gamma, \lambda)} := (1 - \lambda) \left(\hat{A}_t^{(1)} + \lambda \hat{A}_t^{(2)} + \lambda^2 \hat{A}_t^{(3)} + \dots \right) = \sum_{l=0}^{\infty} (\gamma \lambda)^l \delta_{t+l}^V$$

- GAE($\gamma, 0$) : $\hat{A}_t := \delta_t = r_t + \gamma V(s_{t+1}) - V(s_t)$
- GAE($\gamma, 1$) : $\hat{A}_t := \sum_{l=0}^{\infty} \gamma^l \delta_{t+l} = \sum_{l=0}^{\infty} \gamma^l r_{t+l} - V(s_t)$

GENERALIZED ADVANTAGE ESTIMATION

- Notice that we were writing GAE as a function of λ and γ ?
- Normally γ is specified by the problem, and it defines how much we really care about getting a reward now or later
- But we can plan with a smaller γ , treating it as a hyperparameter
(we thus have a ‘problem γ ’ and a ‘solution γ ’)
- What do you think is the effect on bias? On variance?

GENERALIZED ADVANTAGE ESTIMATION

Initialize policy parameter θ_0 and value function parameter ϕ_0 .

for $i = 0, 1, 2, \dots$ **do**

 Simulate current policy π_{θ_i} until N timesteps are obtained.

 Compute δ_t^V at all timesteps $t \in \{1, 2, \dots, N\}$, using $V = V_{\phi_i}$.

 Compute $\hat{A}_t = \sum_{l=0}^{\infty} (\gamma \lambda)^l \delta_{t+l}^V$ at all timesteps.

 Compute θ_{i+1} with TRPO update, Equation (31).

 Compute ϕ_{i+1} with Equation (30).

end for

*could use any value learning
algorithm, here a new method
that uses trust regions
(not covered)*

*This implementation
requires full episodes
(disadvantages?)*

*could use any actor-critic
algorithm, here: TRPO (see next week!)*

Figure from Schulman [2016]

RESULTS: DIFFERENT λ

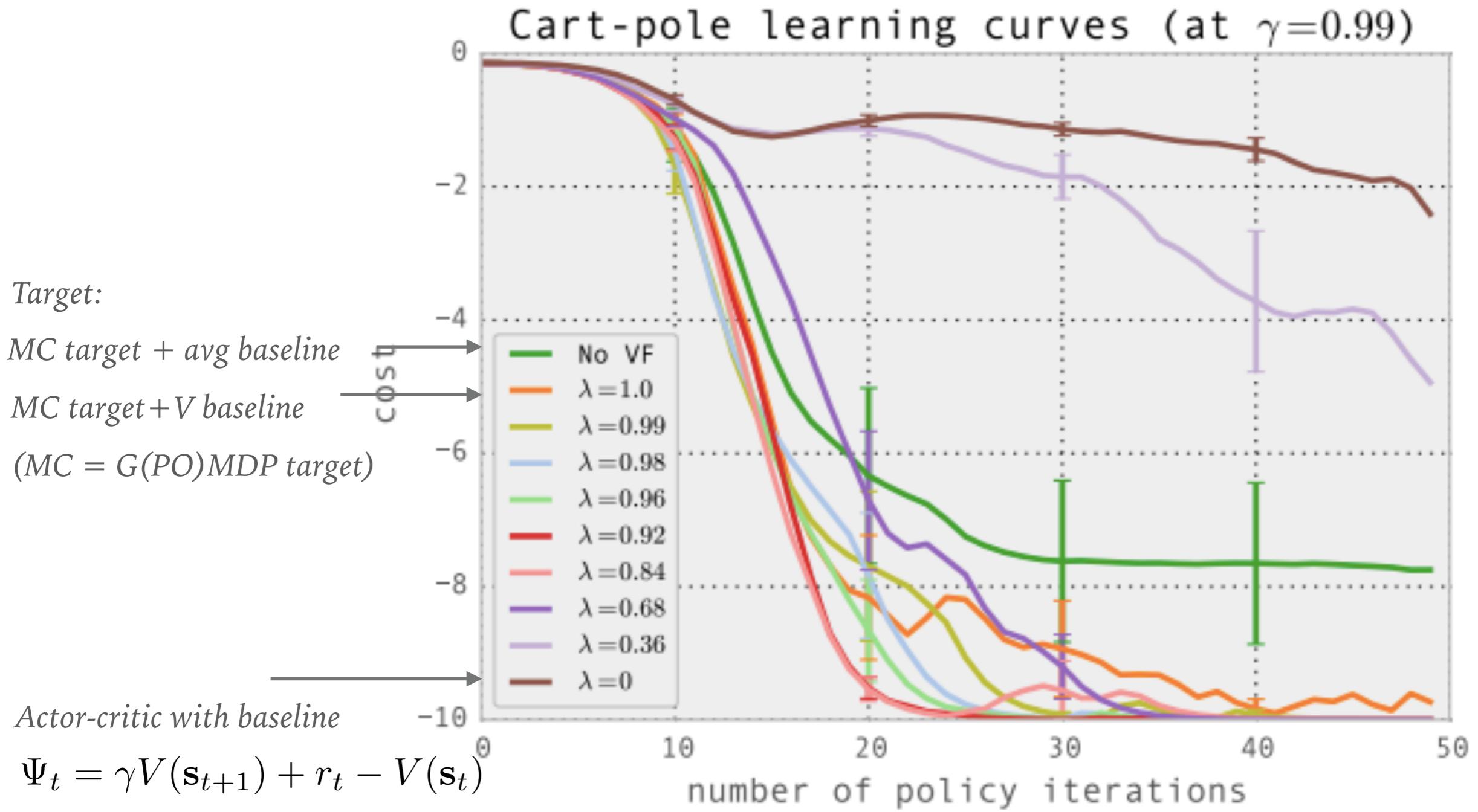
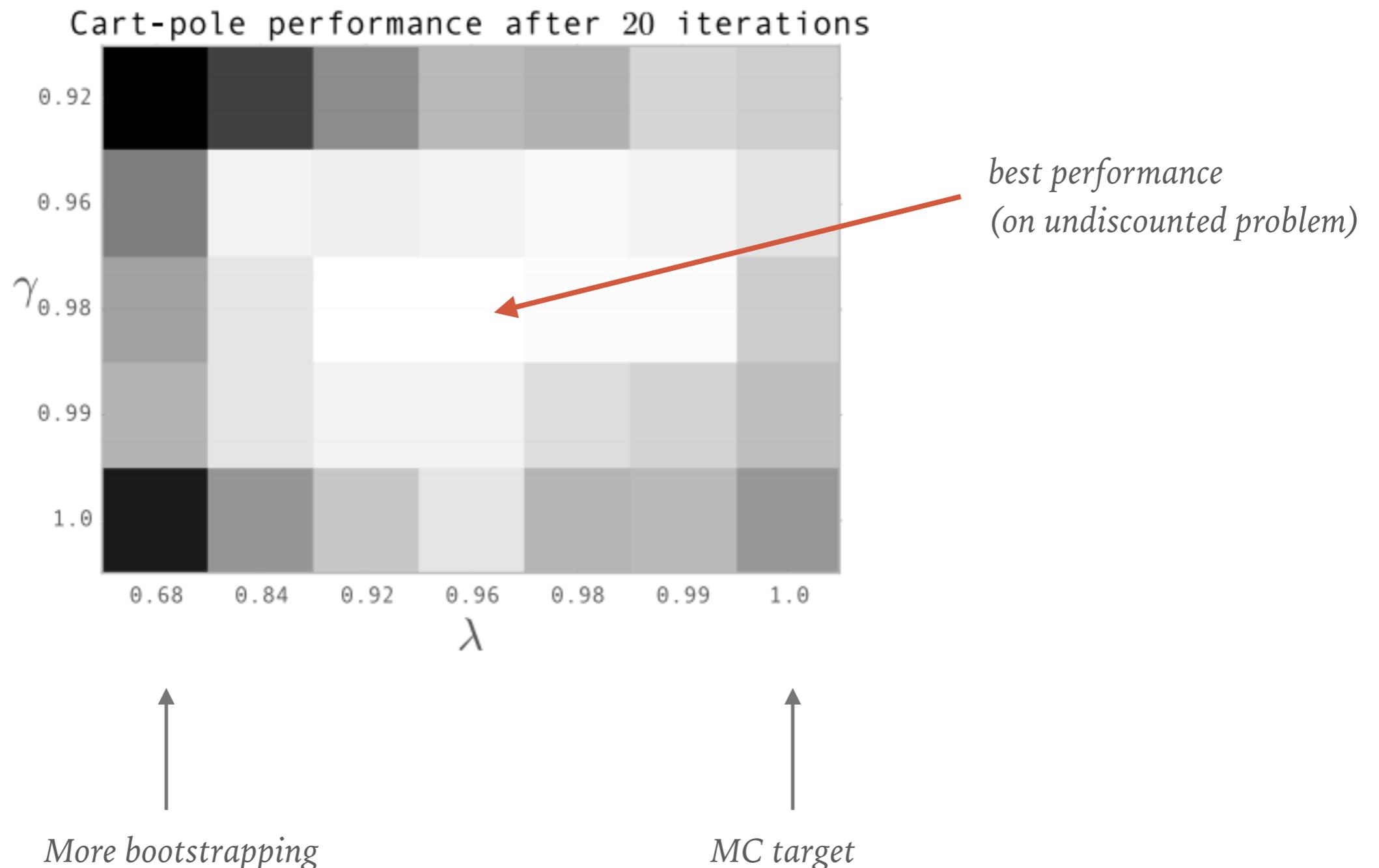


Figure from Schulman [2016]

RESULTS: DIFFERENT λ AND γ



GENERALIZED ADVANTAGE ESTIMATION

- Generalized advantage estimation (GAE) can be used with many policy gradient methods
- E.g. well-known deep-RL method ‘A3C’ [Mnih, 2016] uses an ‘advantage actor critic’ with policy update proportional to:

$$\theta_{t+1} \leftarrow \theta_t + \alpha A_t \nabla_\theta \log \pi_\theta(a_t | s_t)$$

- This is often used with GAE

CONCLUSIONS ABOUT GAE

- GAE can interpolate between:
 - MC (∞ -step return) used by REINFORCE
 - TD (1-step return) used by PGT-AC
- n-step methods are comparable, but GAE averages over many different n
- Can be tuned to problem at hand using λ and γ
- Can yield huge improvement compared to either TD or MC
- Can be used with many policy-based algorithms

Thanks for your attention!

Feedback?

h.c.vanhoof@uva.nl

REFERENCES

- [Schulman 2016] Schulman, J., Moritz, P., Levine, S., Jordan, M.I. and Abbeel, P. High-dimensional continuous control using generalised advantage estimates. In ICLR.