

Deep Reinforcement Learning and Control

Pathwise derivatives, DDPG

CMU 10-703

Katerina Fragkiadaki



$\nabla_{\theta} \log \pi_{\theta}(a)$ for Gaussian policy

- ▶ Policy is Gaussian $a \sim \mathcal{N}(\mu(s, \theta), \sigma^2 I)$
- ▶ Variance may be fixed σ^2 , or can also be parameterized

▶ Remember: univariate Gaussian $\mathcal{N}(\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(a-\mu)^2}{\sigma^2}}$

▶ $\nabla_{\theta} \log \pi_{\theta}(a | s) = \text{const.} \cdot \frac{(\mu_{\theta}(s) - a)}{\sigma^2} \nabla_{\theta} \mu_{\theta}(s)$

$$\hat{g} = \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) (G_t^i - b(s_t^{(i)}))$$

$\nabla_{\theta} \log \pi_{\theta}(a)$ for softmax policy

- Policy is the output of a softmax

$$\pi_{\theta}(a | s) = \frac{e^{h_{\theta}(s,a)}}{\sum_b e^{h_{\theta}(s,b)}}$$

$$\nabla_{\theta} \log \pi_{\theta}(a) = \nabla_{\theta} h_{\theta}(s, a) - \nabla_{\theta} \log \sum_b e^{h_{\theta}(s,b)}$$

$$\nabla_{\theta} h_{\theta}(s, a) - \frac{1}{\sum_b e^{h_{\theta}(s,b)}} \nabla_{\theta} \sum_b e^{h_{\theta}(s,b)}$$

$$\nabla_{\theta} h_{\theta}(s, a) - \frac{1}{\sum_b e^{h_{\theta}(s,b)}} \sum_b \nabla_{\theta} e^{h_{\theta}(s,b)}$$

$$\nabla_{\theta} h_{\theta}(s, a) - \frac{1}{\sum_b e^{h_{\theta}(s,b)}} \sum_b e^{h_{\theta}(s,b)} \nabla_{\theta} h_{\theta}(s, b)$$

$$\nabla_{\theta} h_{\theta}(s, a) - \sum_b \frac{e^{h_{\theta}(s,b)}}{\sum_b e^{h_{\theta}(s,b)}} \nabla_{\theta} h_{\theta}(s, b)$$

$$\nabla_{\theta} h_{\theta}(s, a) - \sum_b \pi_{\theta}(s, b) \nabla_{\theta} h_{\theta}(s, b)$$

$$\hat{g} = \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) (G_t^i - b(s_t^{(i)}))$$

Computing Gradients of Expectations

When the variable w.r.t. which we are differentiating appears **in the distribution**:

$$\nabla_{\theta} \mathbb{E}_{x \sim P_{\theta}(x)} f(x) = \mathbb{E}_{x \sim P_{\theta}(x)} \nabla_{\theta} \log P_{\theta}(x) f(x)$$

likelihood ratio gradient estimator

When the variable w.r.t. which we are differentiating appears **inside the expectation**:

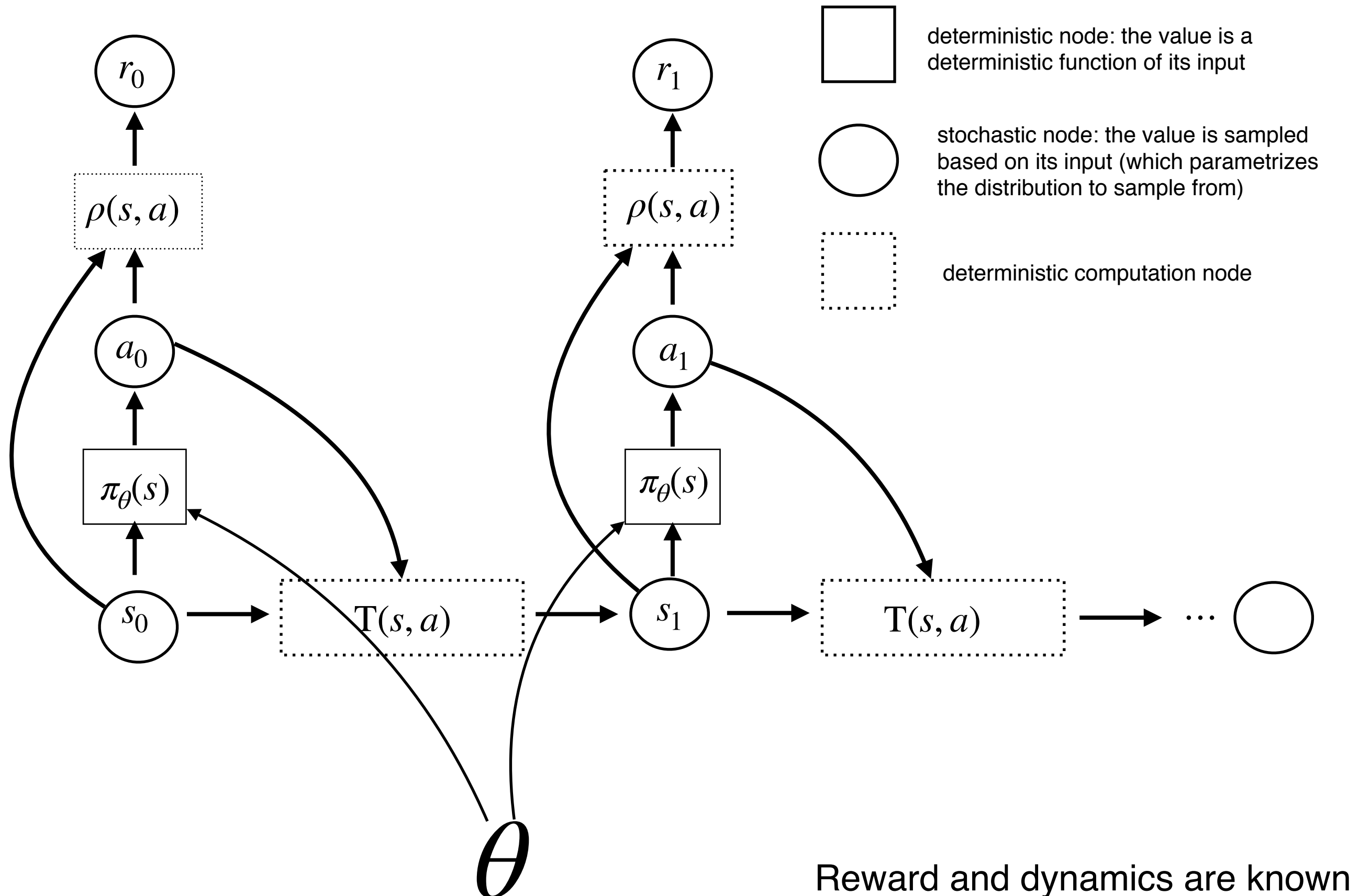
$$\nabla_{\theta} \mathbb{E}_{z \sim \mathcal{N}(0,1)} f(x(\theta), z) = \mathbb{E}_{z \sim \mathcal{N}(0,1)} \nabla_{\theta} f(x(\theta), z) = \mathbb{E}_{z \sim \mathcal{N}(0,1)} \frac{df(x(\theta), z)}{dx} \frac{dx}{d\theta}$$

pathwise derivative

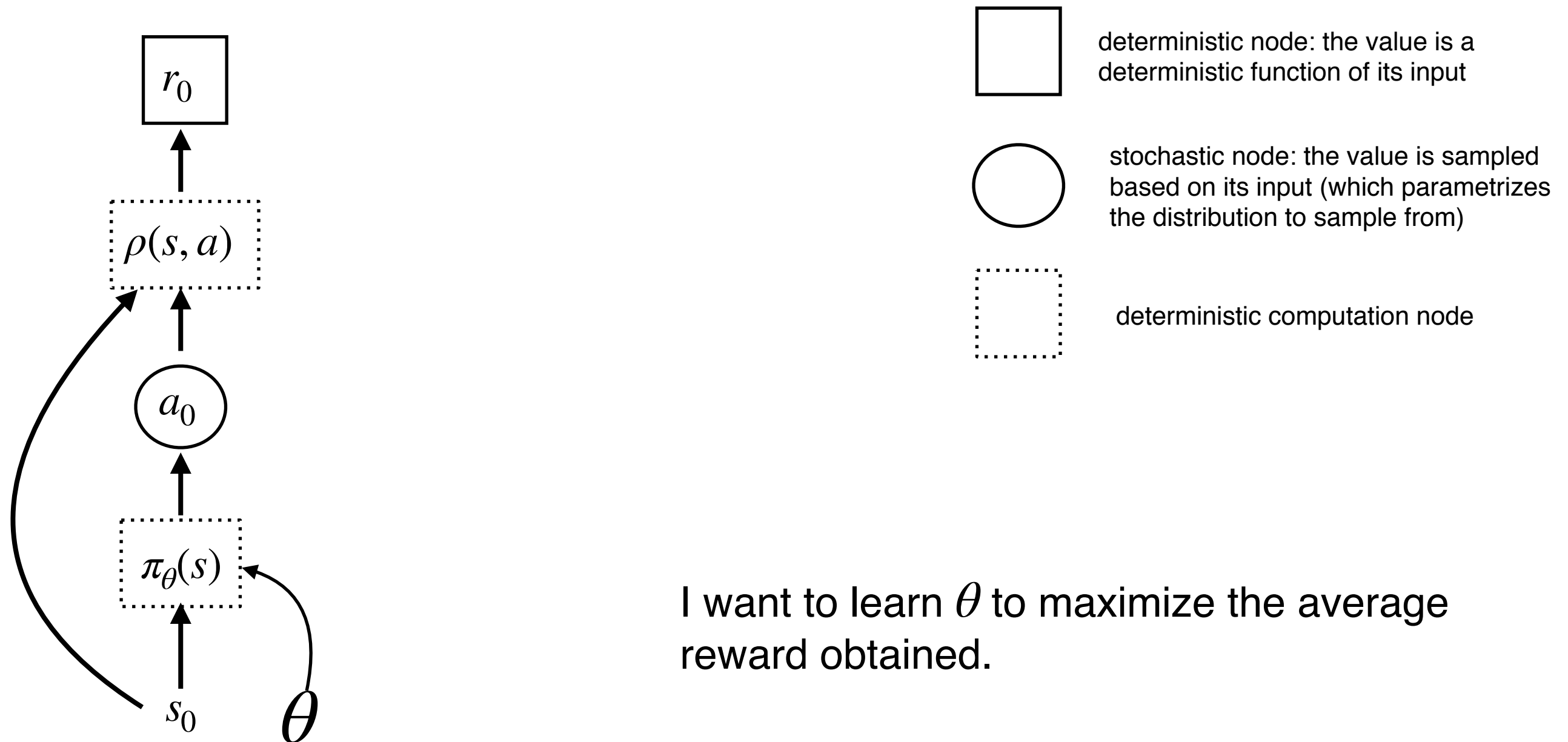
Re-parametrization trick: For some distributions $P_{\theta}(x)$ we can switch from one gradient estimator to the other.

Q: From which to which? Why would we want to do so?

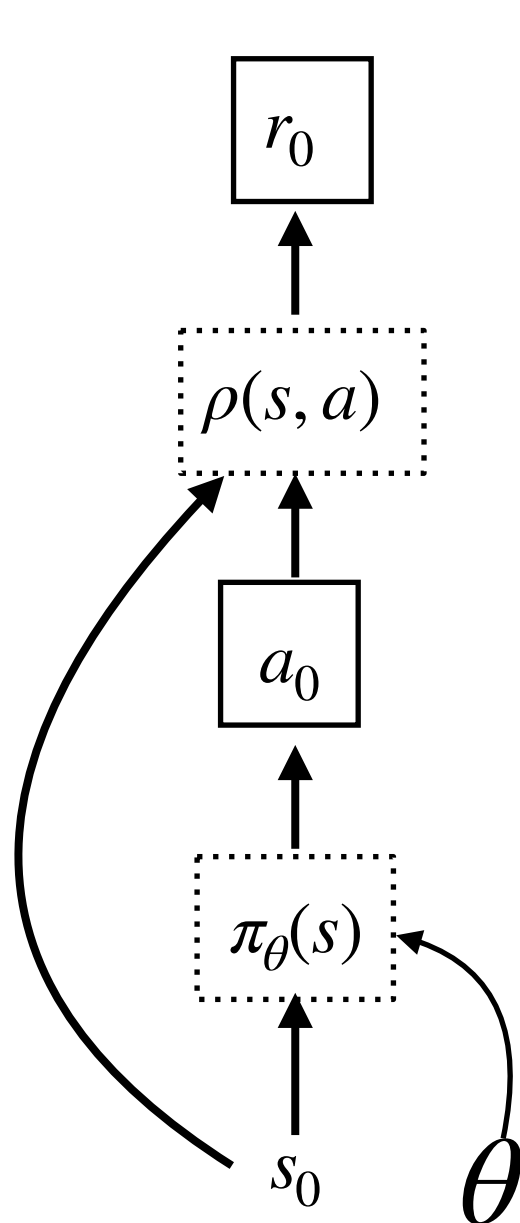
Known MDP (what if that was possible)



Known MDP-let's make it simpler

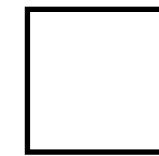


Deterministic policy

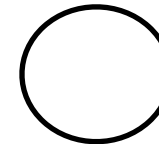


$$a = \pi_\theta(s)$$

$$\max_{\theta} \rho(s_0, a)$$



deterministic node: the value is a deterministic function of its input



stochastic node: the value is sampled based on its input (which parametrizes the distribution to sample from)



deterministic computation node

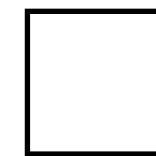
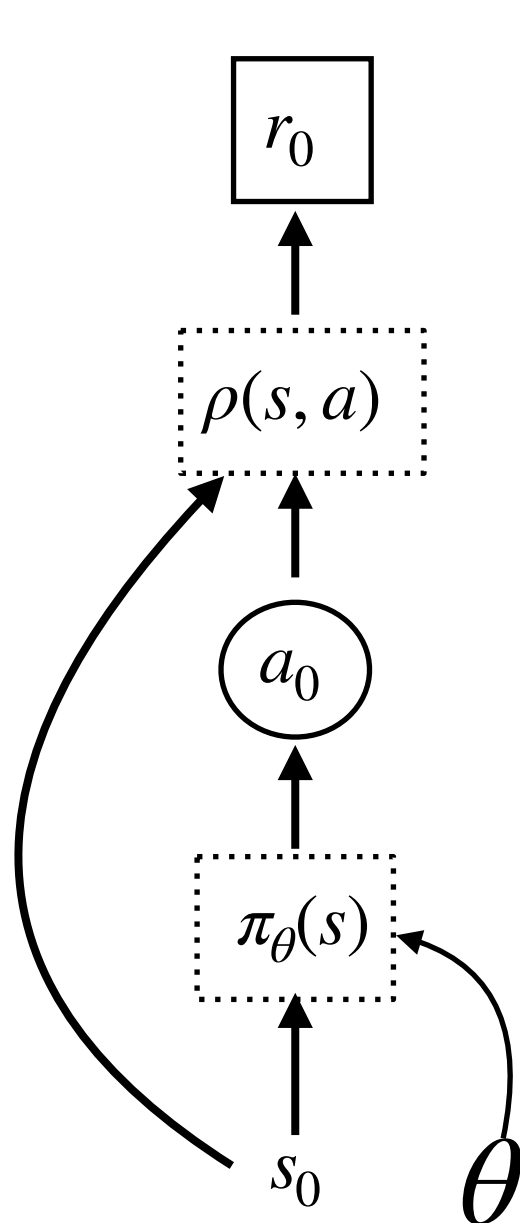
I want to learn θ to maximize the average reward obtained.

I can compute the gradient with backpropagation.

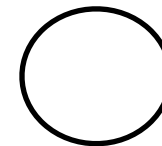
$$\nabla_{\theta} \rho(s, a) = \rho_a \pi_{\theta\theta}$$

Derivative of the *known* reward function w.r.t. the action

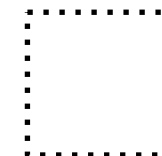
Stochastic policy



deterministic node: the value is a deterministic function of its input



stochastic node: the value is sampled based on its input (which parametrizes the distribution to sample from)



deterministic computation node

I want to learn θ to maximize the **expected** reward obtained.

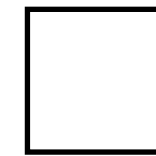
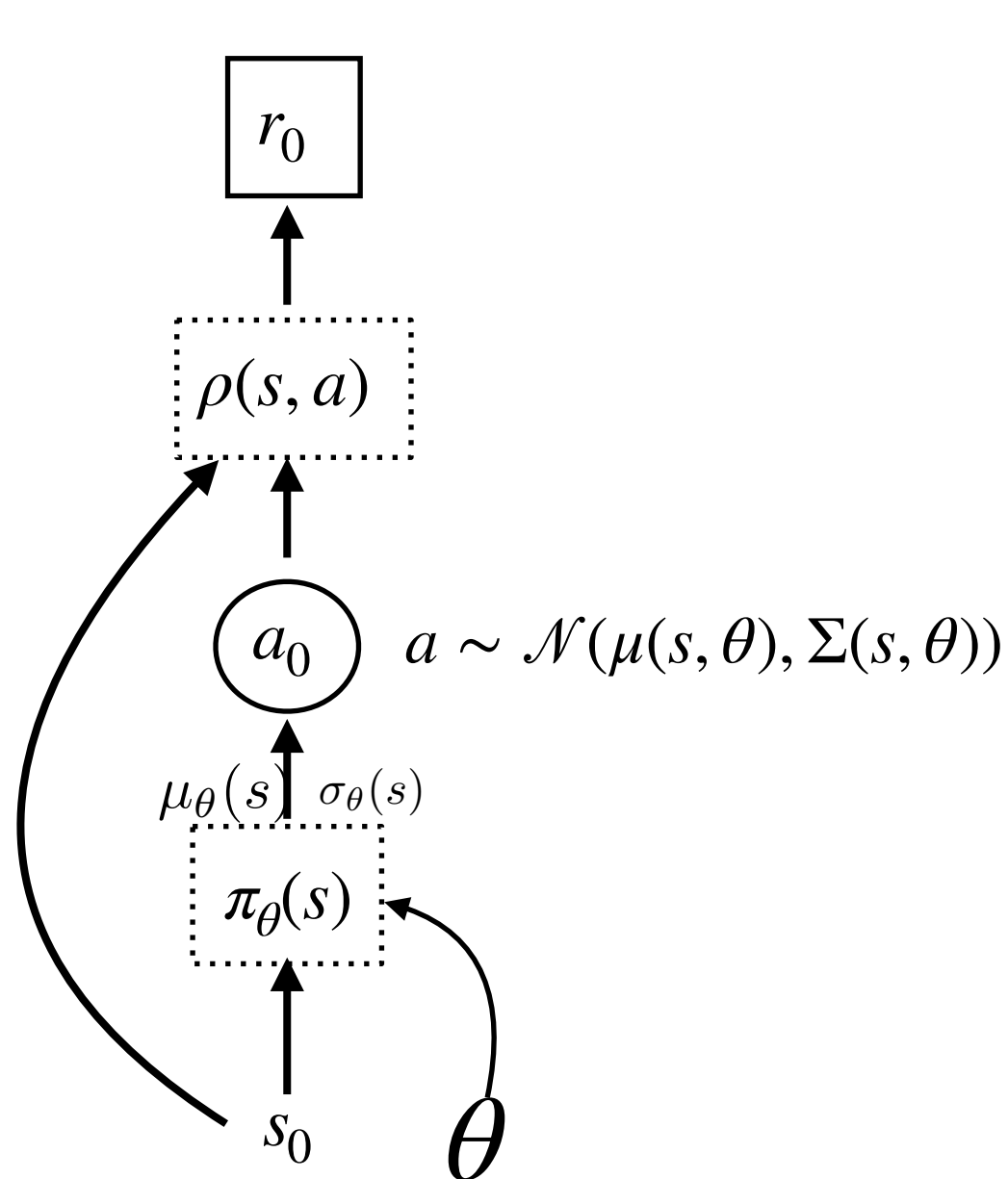
Likelihood ratio estimator, works for both continuous and discrete actions

$$\mathbb{E}_a \nabla_\theta \log \pi_\theta(s) \rho(s, a)$$

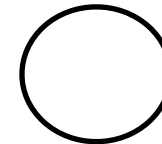
$$\max_{\theta} \mathbb{E}_a \rho(s_0, a)$$

It does not use the derivative of the reward w.r.t. the action.

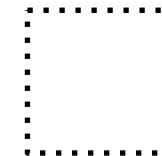
Example: Gaussian policy



deterministic node: the value is a deterministic function of its input



stochastic node: the value is sampled based on its input (which parametrizes the distribution to sample from)



deterministic computation node

I want to learn θ to maximize the average reward obtained.

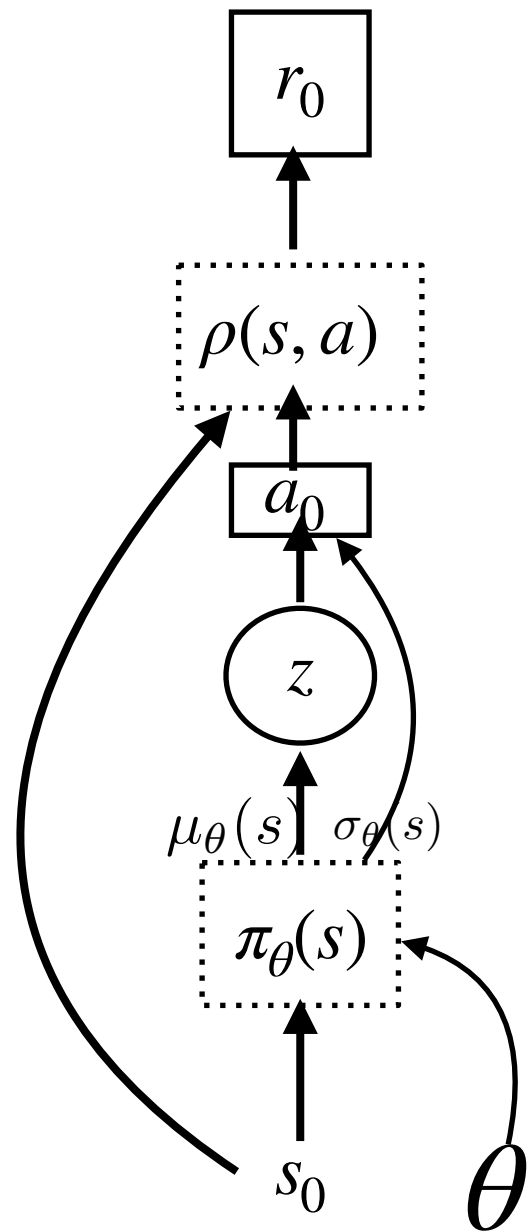
$$\mathbb{E}_a \nabla_\theta \log \pi_\theta(s) \rho(s, a)$$

If σ^2 is constant:

$$\max_{\theta} \mathbb{E}_a \rho(s_0, a)$$

$$\nabla_\theta \log \pi_\theta(s, a) = \frac{(a - \mu(s; \theta)) \frac{\partial \mu(s; \theta)}{\partial \theta}}{\sigma^2}$$

Re-parametrization for Gaussian



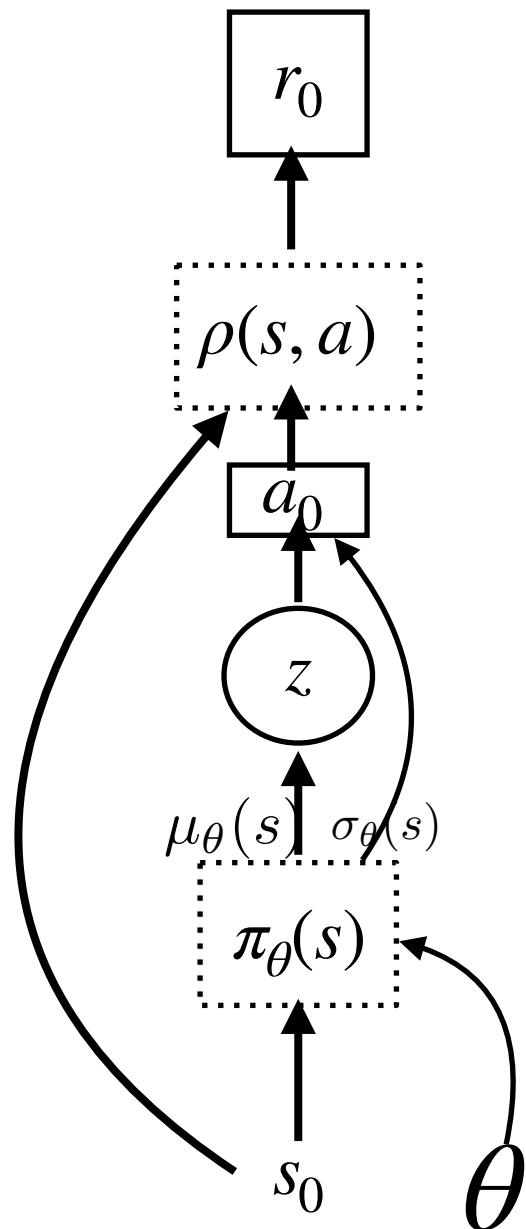
$$a \sim \mathcal{N}(\mu(s, \theta), \Sigma(s, \theta))$$

We can either:

- Assume σ fixed (spherical or isotropic Gaussian)
- Learn $\sigma(s, \theta)$ one value for all action coordinates (spherical or isotropic Gaussian)
- Learn $\sigma^i(s, \theta), i = 1 \cdots n$, (diagonal covariance)
- Learn a full covariance matrix $\Sigma(s, \theta)$

$$\max_{\theta} \mathbb{E}_a \rho(s_0, a)$$

Re-parametrization for Gaussian



Instead of: $a \sim \mathcal{N}(\mu(s, \theta), \Sigma(s, \theta))$

We can write: $a = \mu(s, \theta) + z \odot \sigma(s, \theta) \quad z \sim \mathcal{N}(0, I)$

Why?

Because: $\mathbb{E}_z(\mu(s, \theta) + z\sigma(s, \theta)) = \mu(s, \theta)$
 $\text{Var}_z(\mu(s, \theta) + z\sigma(s, \theta)) = \sigma(s, \theta)^2$

Qs:

- Does a depend on θ ?
- Does z depend on θ ?

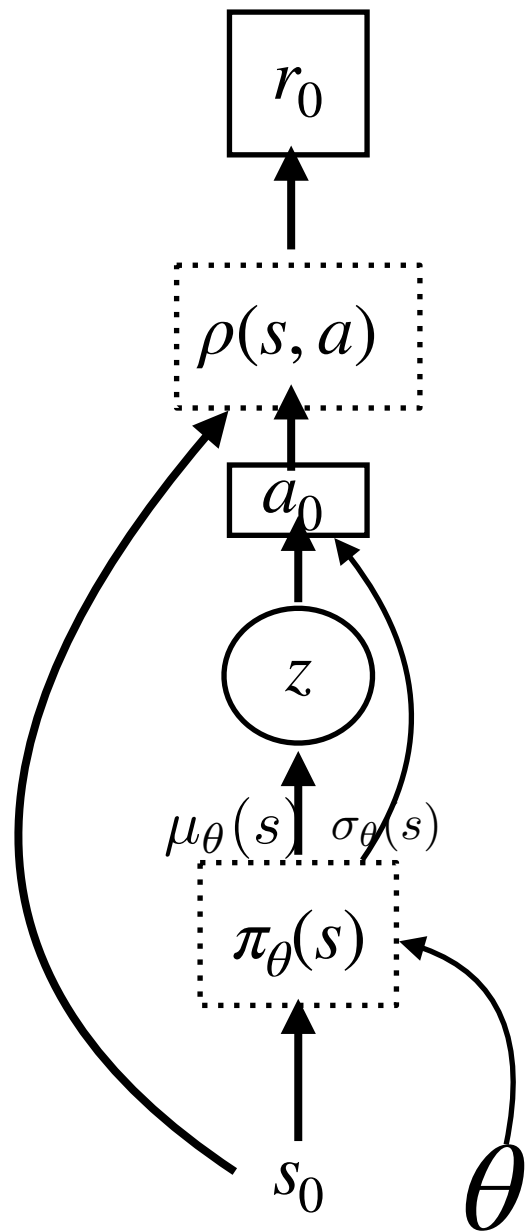
$$\max_{\theta} \mathbb{E}_a \rho(s_0, a)$$



$$\max_{\theta} \mathbb{E}_z \rho(s_0, a(z))$$

$$a = \mu(s, \theta) + L(s, \theta)z, \quad \Sigma = L(s, \theta)L(s, \theta)^{\top}$$

Re-parametrization for Gaussian



Instead of: $a \sim \mathcal{N}(\mu(s, \theta), \Sigma(s, \theta))$

We can write: $a = \mu(s, \theta) + z \odot \sigma(s, \theta) \quad z \sim \mathcal{N}(0, I)$

Why?

Because: $\mathbb{E}_z(\mu(s, \theta) + z\sigma(s, \theta)) = \mu(s, \theta)$

$\text{Var}_z(\mu(s, \theta) + z\sigma(s, \theta)) = \sigma(s, \theta)^2$

What do we gain?

$$\nabla_{\theta} \mathbb{E}_z [\rho(a(\theta, z), s)] = \mathbb{E}_z \frac{d\rho(a(\theta, z), s)}{da} \frac{da(\theta, z)}{d\theta}$$

$$\frac{da(\theta, z)}{d\theta} = \frac{d\mu(s, \theta)}{d\theta} + z \odot \frac{d\sigma(s, \theta)}{d\theta}$$

Sample estimate:

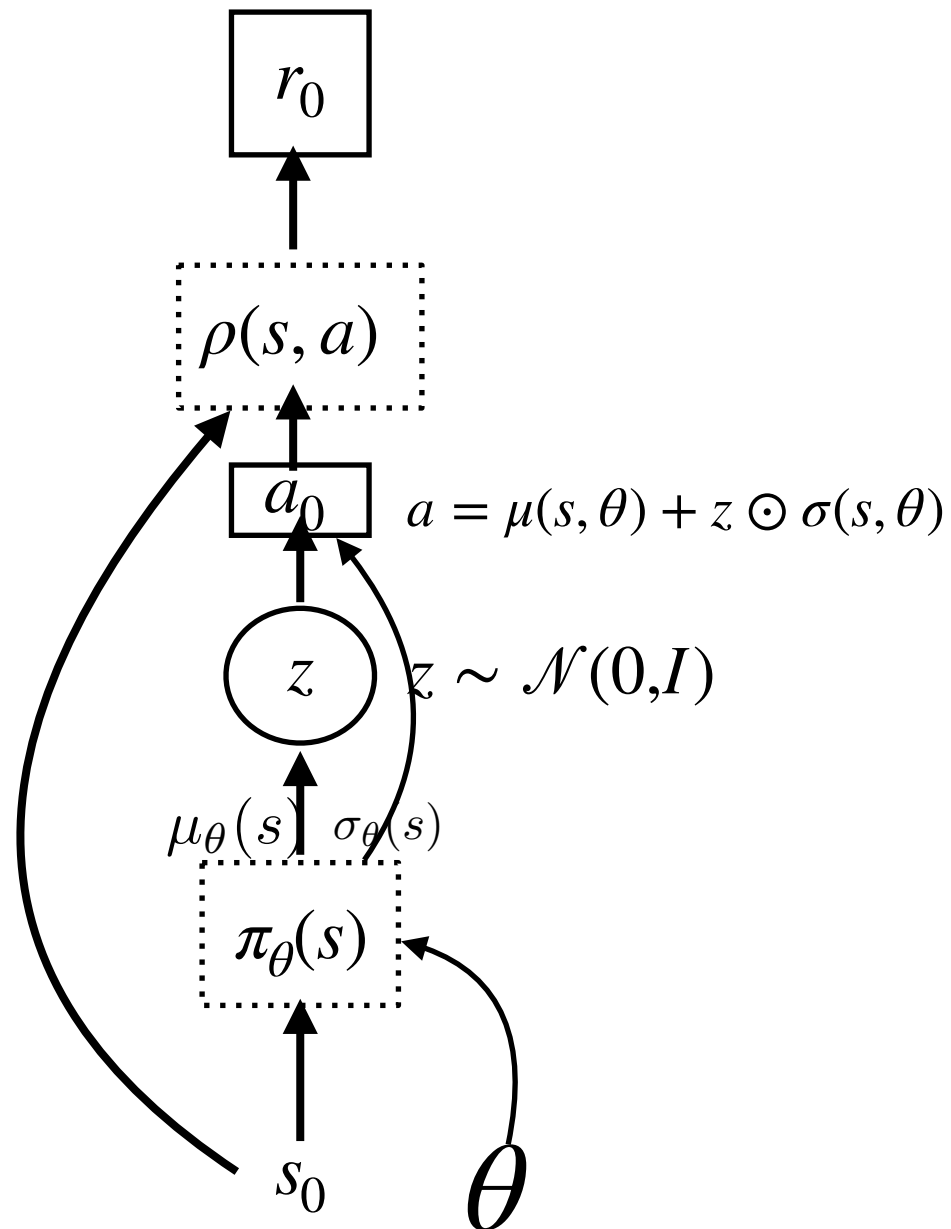
$$\max_{\theta} \mathbb{E}_a \rho(s_0, a)$$



$$\max_{\theta} \mathbb{E}_z \rho(s_0, a(z))$$

$$\nabla_{\theta} \frac{1}{N} \sum_{i=1}^N [\rho(a(\theta, z_i), s)] = \frac{1}{N} \sum_{i=1}^N \frac{d\rho(a(\theta, z_i), s)}{da} \frac{da(\theta, z_i)}{d\theta} \Big|_{z=z_i}$$

Re-parametrization for Gaussian



Likelihood ratio grad estimator:

$$\mathbb{E}_a \nabla_\theta \log \pi_\theta(s, a) \rho(s, a)$$

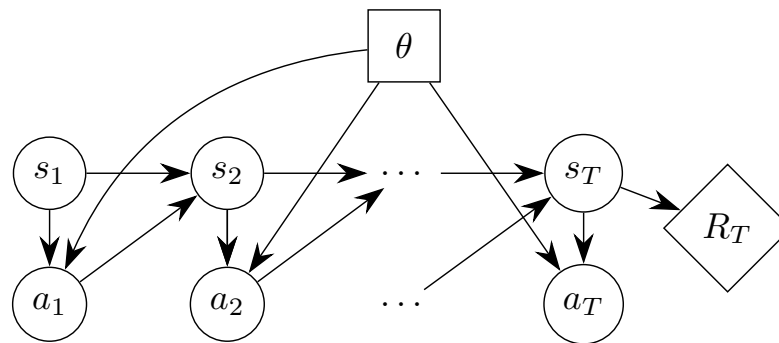
Pathwise derivative:

$$\mathbb{E}_z \frac{d\rho(a(\theta, z), s)}{da} \frac{da(\theta, z)}{d\theta}$$

The pathwise derivative uses the derivative of the reward w.r.t. the action!

Re-parametrized Policy Gradients

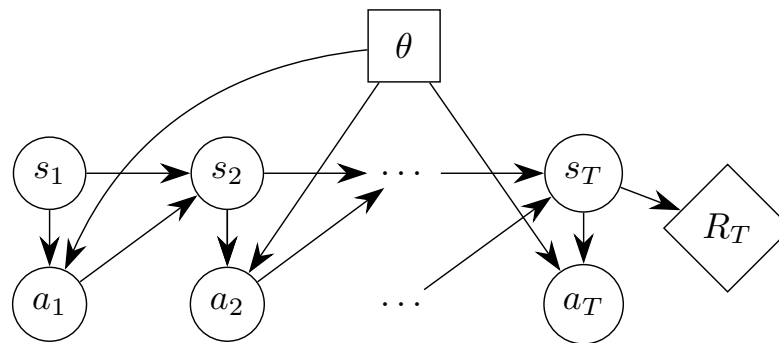
- Episodic MDP:



We want to compute: $\nabla_{\theta} \mathbb{E}[R_T]$

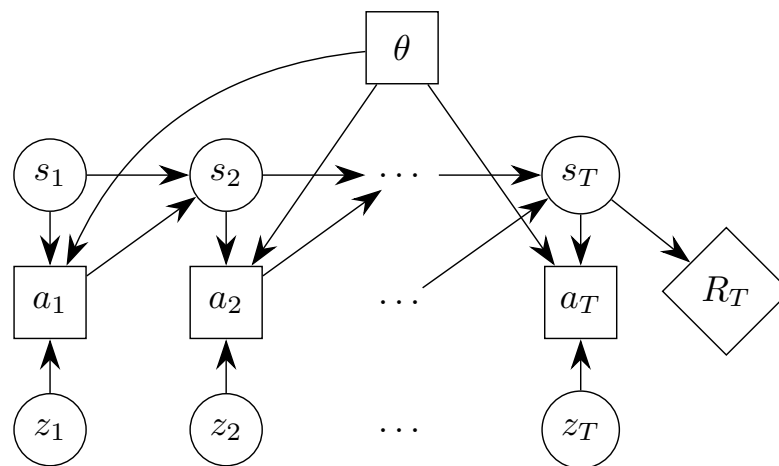
Re-parametrized Policy Gradients

- Episodic MDP:



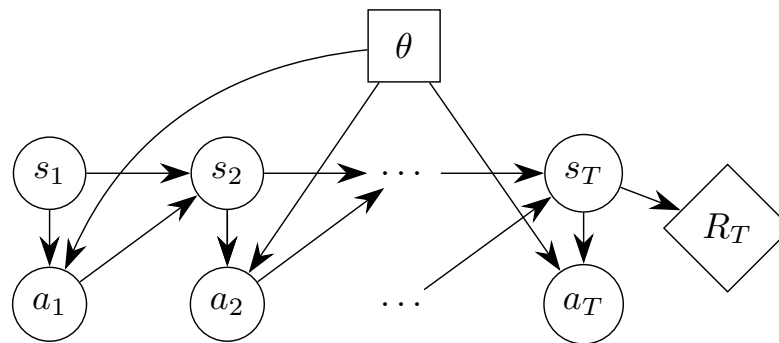
We want to compute: $\nabla_{\theta} \mathbb{E}[R_T]$

- Reparameterize: $a_t = \pi(s_t, z_t; \theta)$. z_t is noise from fixed distribution.



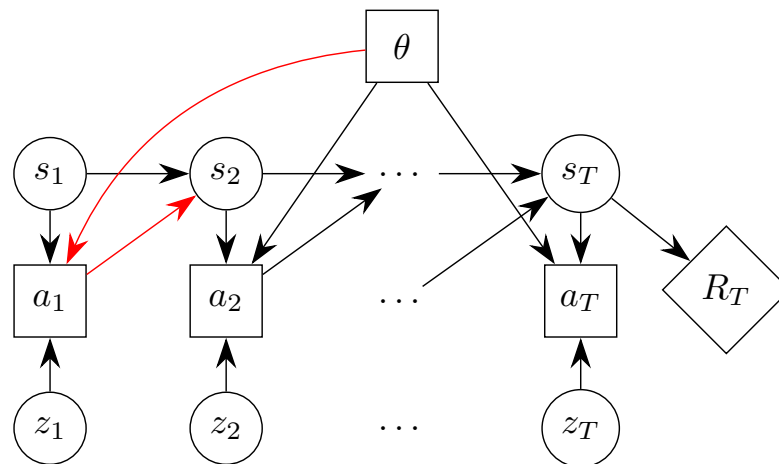
Re-parametrized Policy Gradients

- Episodic MDP:



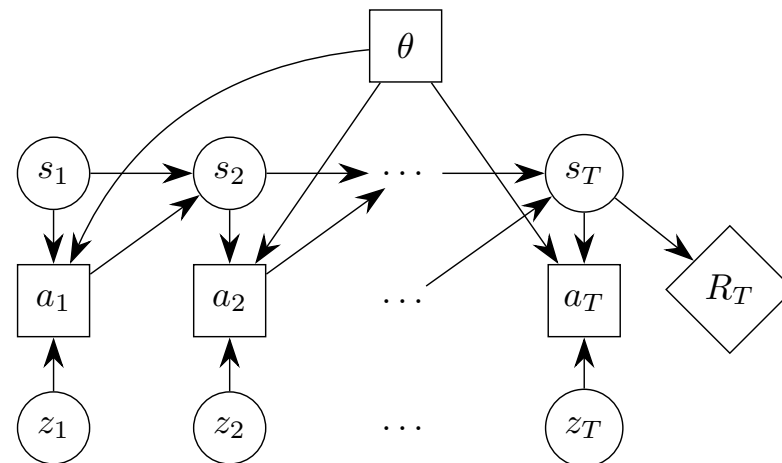
We want to compute: $\nabla_{\theta} \mathbb{E}[R_T]$

- Reparameterize: $a_t = \pi(s_t, z_t; \theta)$. z_t is noise from fixed distribution.



For pathwise derivative to work, we need transition dynamics and reward function to be known.

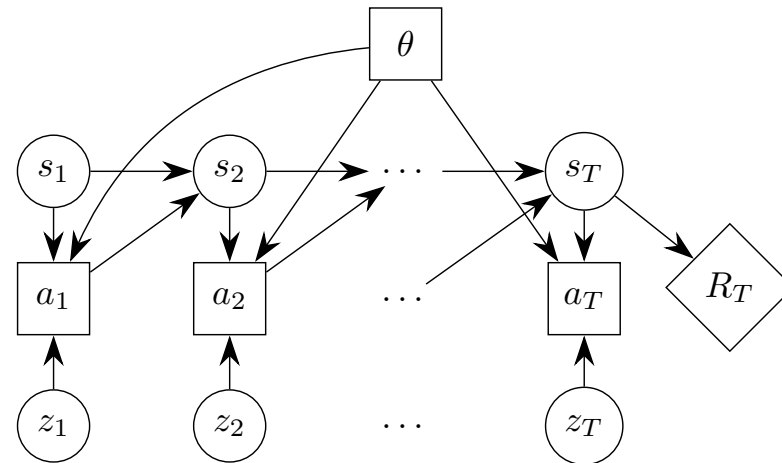
Re-parametrized Policy Gradients



$$\frac{d}{d\theta} \mathbb{E} [R_T] = \mathbb{E} \left[\sum_{t=1}^T \frac{dR_T}{da_t} \frac{da_t}{d\theta} \right] = \mathbb{E} \left[\sum_{t=1}^T \frac{d}{da_t} \mathbb{E} [R_T | a_t] \frac{da_t}{d\theta} \right]$$

For path wise derivative to work, we need transition dynamics and reward function to be known, or...

Re-parametrized Policy Gradients



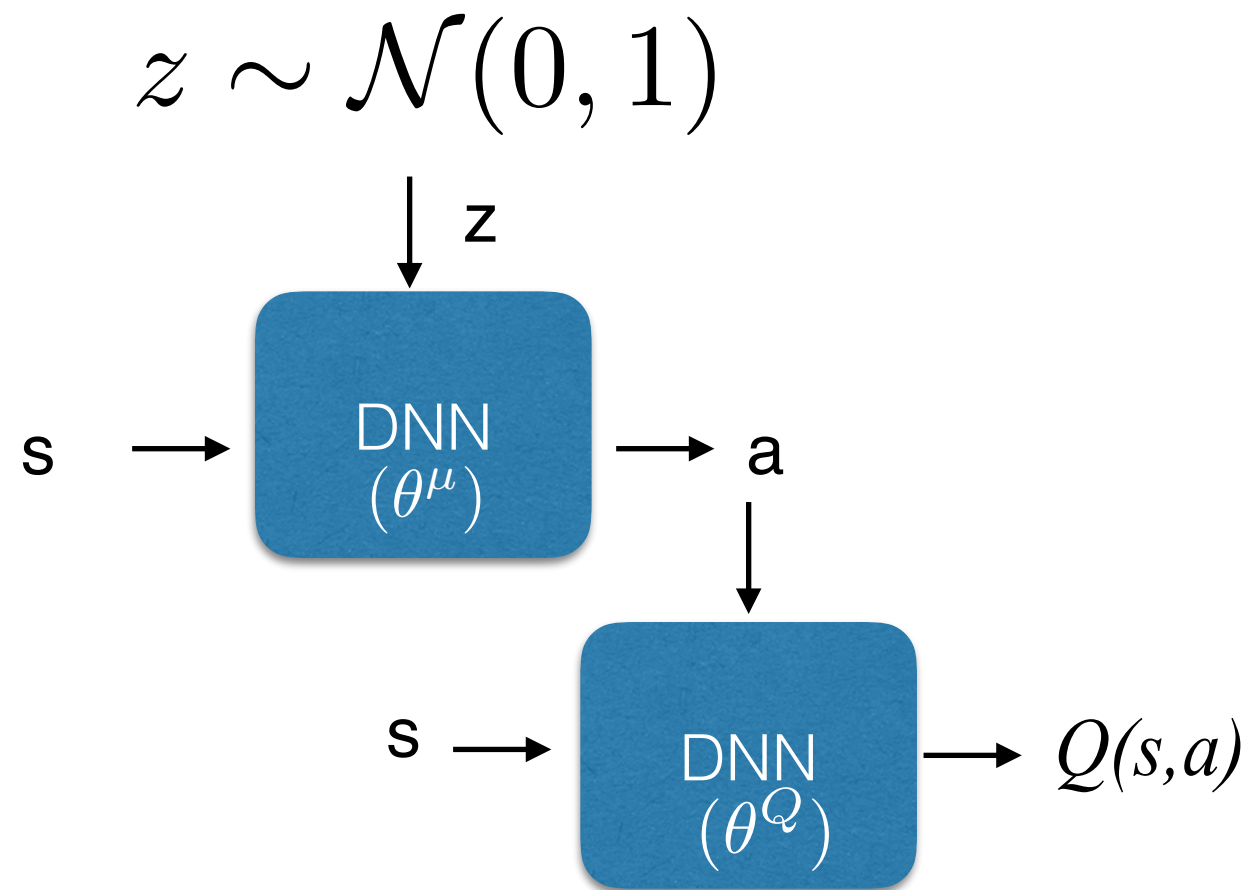
$$\begin{aligned} \frac{d}{d\theta} \mathbb{E}[R_T] &= \mathbb{E} \left[\sum_{t=1}^T \frac{dR_T}{da_t} \frac{da_t}{d\theta} \right] = \mathbb{E} \left[\sum_{t=1}^T \frac{d}{da_t} \mathbb{E}[R_T | a_t] \frac{da_t}{d\theta} \right] \\ &= \mathbb{E} \left[\sum_{t=1}^T \frac{dQ(s_t, a_t)}{da_t} \frac{da_t}{d\theta} \right] = \mathbb{E} \left[\sum_{t=1}^T \frac{d}{d\theta} Q(s_t, \pi(s_t, z_t; \theta)) \right] \end{aligned}$$

- Learn Q_ϕ to approximate $Q^{\pi, \gamma}$, and use it to compute gradient estimates.

Stochastic Value Gradients V0

- ▶ Learn Q_ϕ to approximate $Q^{\pi,\gamma}$, and use it to compute gradient estimates.
- ▶ Pseudocode:
 - for** iteration=1, 2, ... **do**
 - Execute policy π_θ to collect T timesteps of data
 - Update π_θ using $g \propto \nabla_\theta \sum_{t=1}^T Q(s_t, \pi(s_t, z_t; \theta))$
 - Update Q_ϕ using $g \propto \nabla_\phi \sum_{t=1}^T (Q_\phi(s_t, a_t) - \hat{Q}_t)^2$, e.g. with TD(λ)
 - end for**

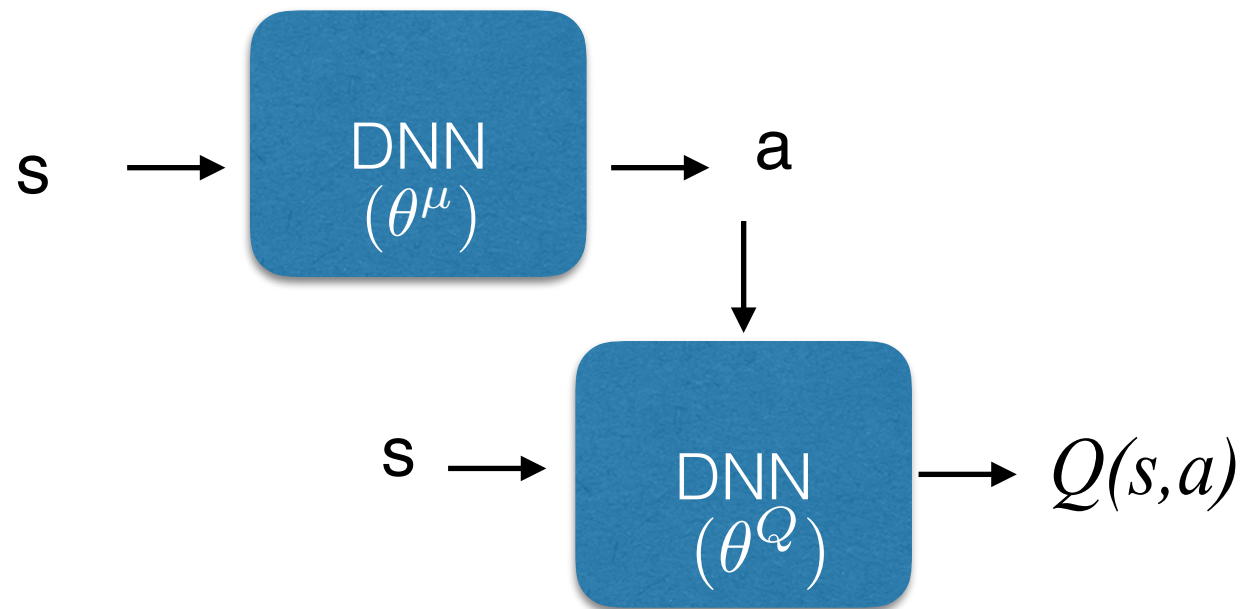
Stochastic Value Gradients V0



$$a = \mu(s; \theta) + z\sigma(s; \theta)$$

Compare with: Deep Deterministic Policy Gradients

$$a = \mu(\theta)$$



No z !