



香港科技大學
THE HONG KONG
UNIVERSITY OF SCIENCE
AND TECHNOLOGY

COMP 5212
Machine Learning
Lecture 22

Generative Adversarial Networks, Reinforcement Learning

Junxian He
Nov 26, 2024

Announcement

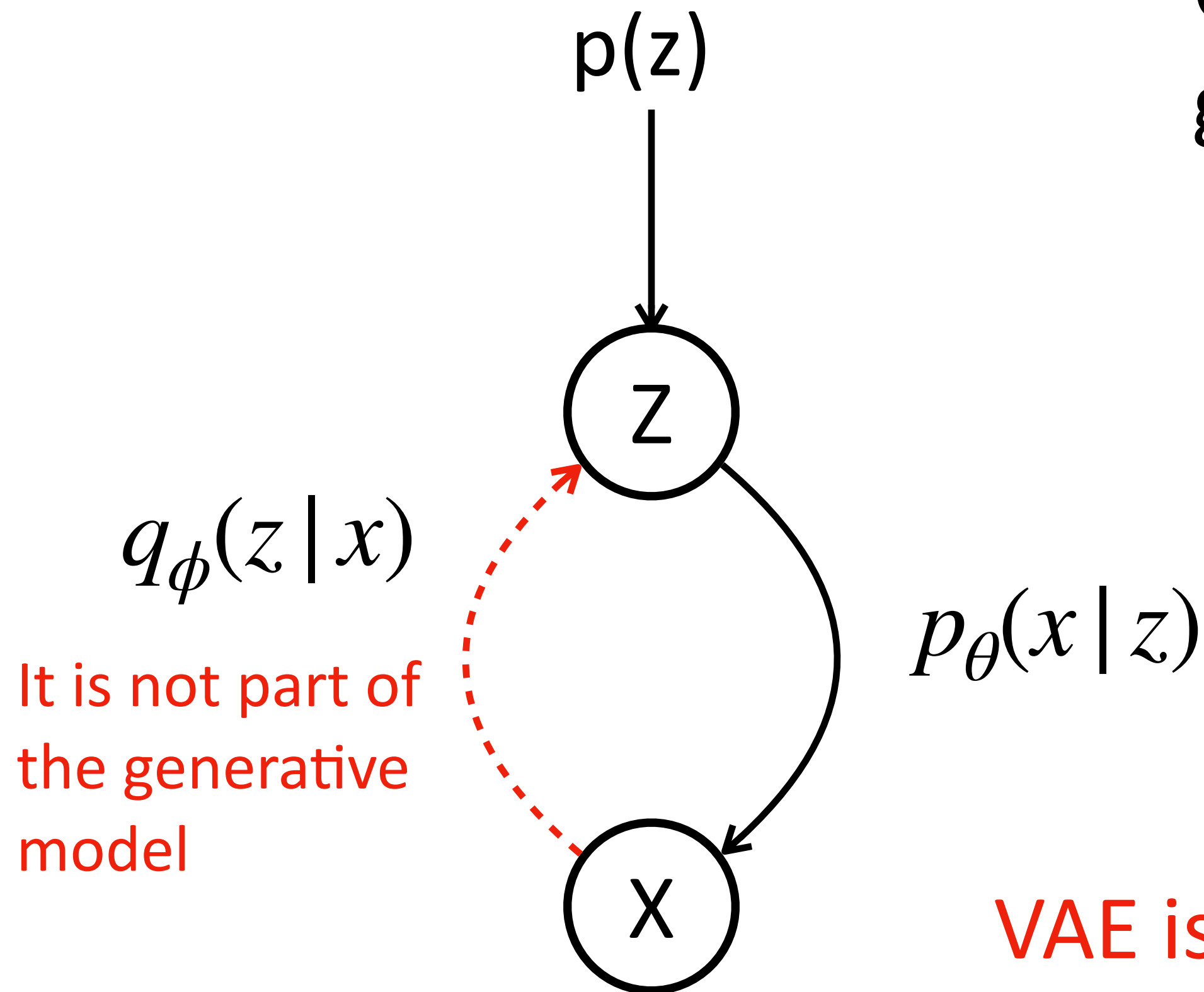
- HW4 is out, it is fairly easy, mainly a reflection of all the COMP5212 contents with only multi-choice questions
- The first round of Kaggle private leaderboard was released last night — do not overoptimize the public leaderboard too much

Recap: VAEs

Only the right (black) part defines the generative model, and the distribution

$p_{\theta}(x | z)$: generative network/decoder

$q_{\phi}(z | x)$: inference network/encoder



VAE is a name to represent both the model $p(x)$ and the inference network that is used to train the model, but do not confuse them together

Training VAEs

E-Step:

$$\operatorname{argmax}_{\phi} \underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction Loss}} - \underbrace{D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))}_{\text{KL Regularizer}}$$

M-Step:

$$\operatorname{argmax}_{\theta} \underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction Loss}} - \underbrace{D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))}_{\text{KL Regularizer}}$$

Intuitively we hope to approximate $p(\mathbf{z}|\mathbf{x})$ with $q(\mathbf{z}|\mathbf{x})$ accurately in the E-step, to approximate the true EM algorithm

Training VAEs

Algorithm 1 Minibatch version of the Auto-Encoding VB (AEVB) algorithm. Either of the two SGVB estimators in section 2.3 can be used. We use settings $M = 100$ and $L = 1$ in experiments.

$\theta, \phi \leftarrow$ Initialize parameters

repeat

$\mathbf{X}^M \leftarrow$ Random minibatch of M datapoints (drawn from full dataset)

$\epsilon \leftarrow$ Random samples from noise distribution $p(\epsilon)$

$\mathbf{g} \leftarrow \nabla_{\theta, \phi} \tilde{\mathcal{L}}^M(\theta, \phi; \mathbf{X}^M, \epsilon)$ (Gradients of minibatch estimator (8))

$\theta, \phi \leftarrow$ Update parameters using gradients \mathbf{g} (e.g. SGD or Adagrad [DHS10])

until convergence of parameters (θ, ϕ)

return θ, ϕ

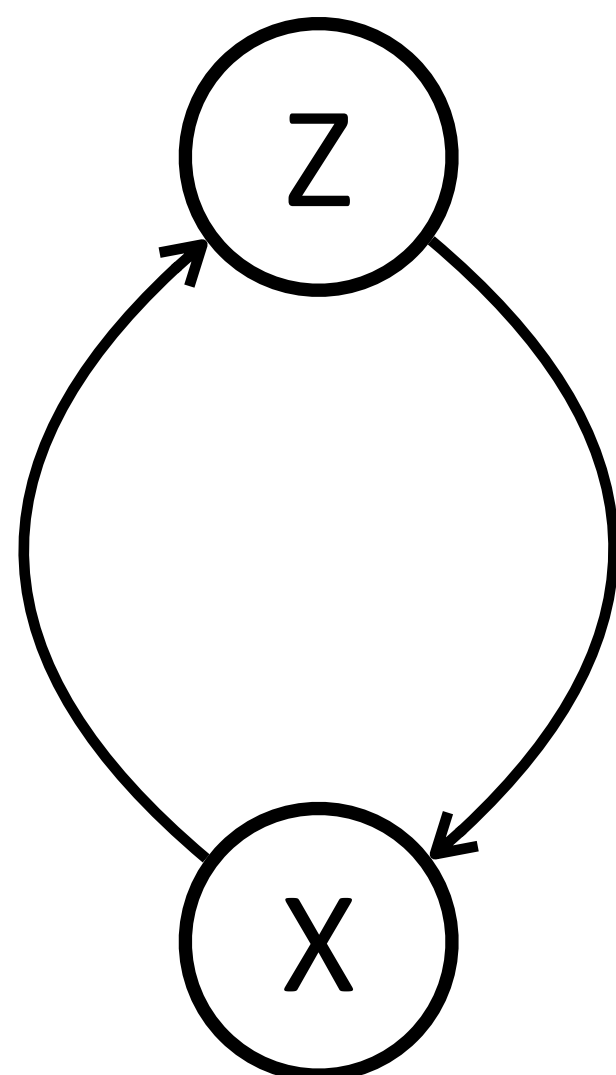
End-to-end, because the objectives are the same (ELBO)

VAE training is optimizing ELBO with gradient descent

AutoEncoders

$$\text{VAE: } \underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction Loss}} - \underbrace{D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))}_{\text{KL Regularizer}}$$

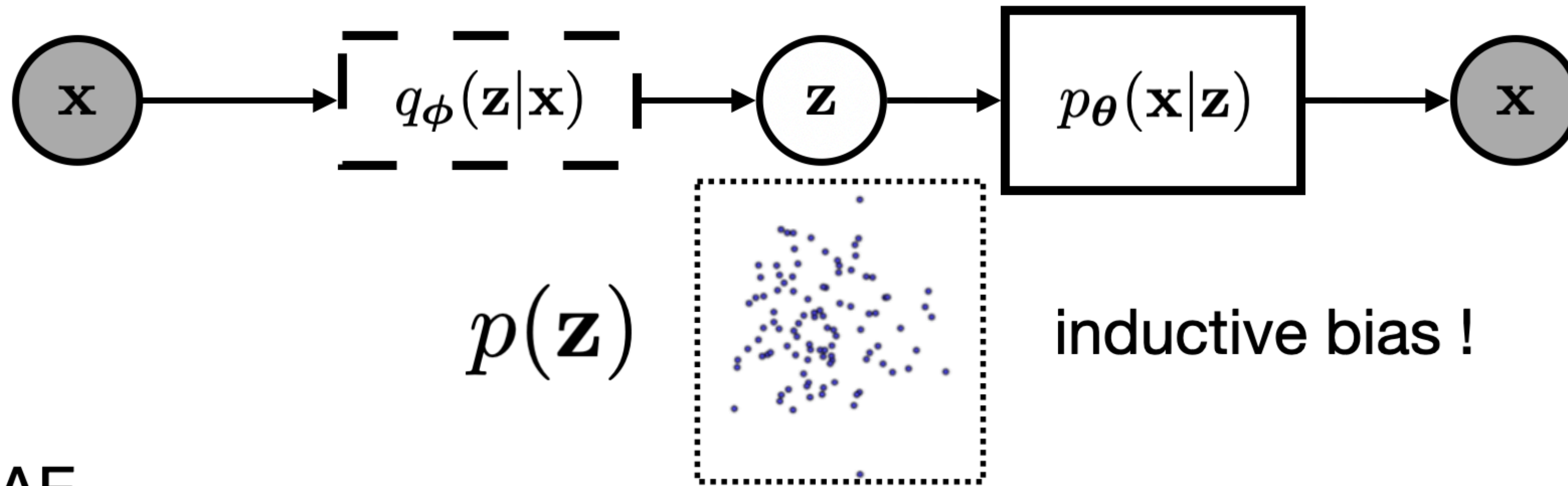
$$\text{AE: } \log p_{\theta}(x | q(x))$$



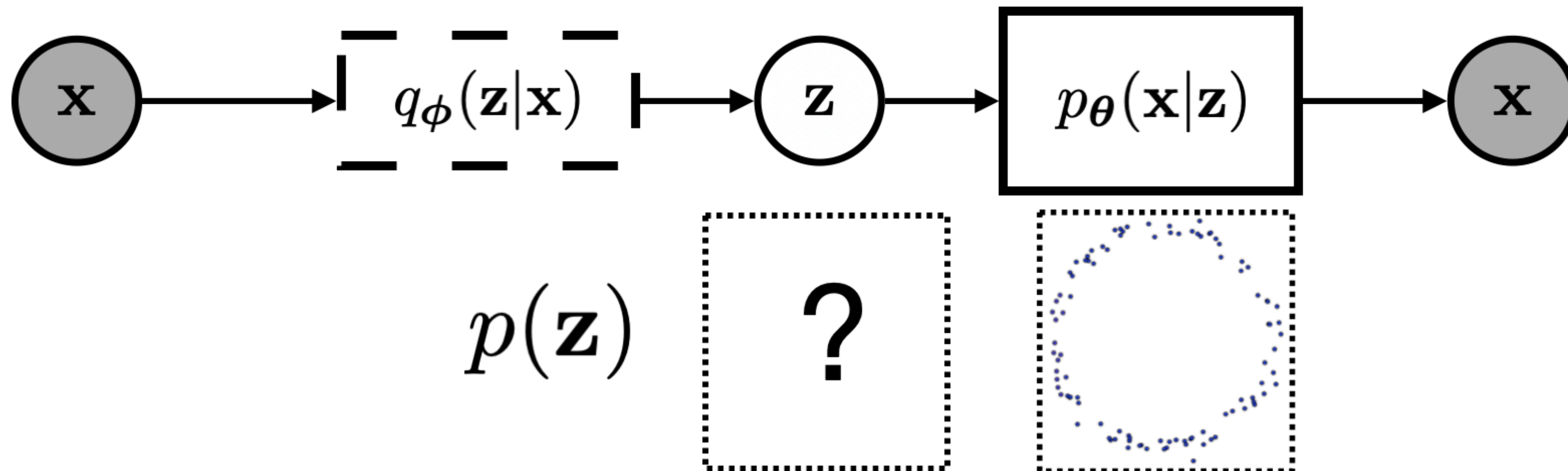
1. Can we generate X samples from an autoencoder?
2. Can we approximate $p(x)$ given x with an autoencoder?
3. What is the difference between the representation space from AE and VAE?

VAE v.s. AE

VAE



AE

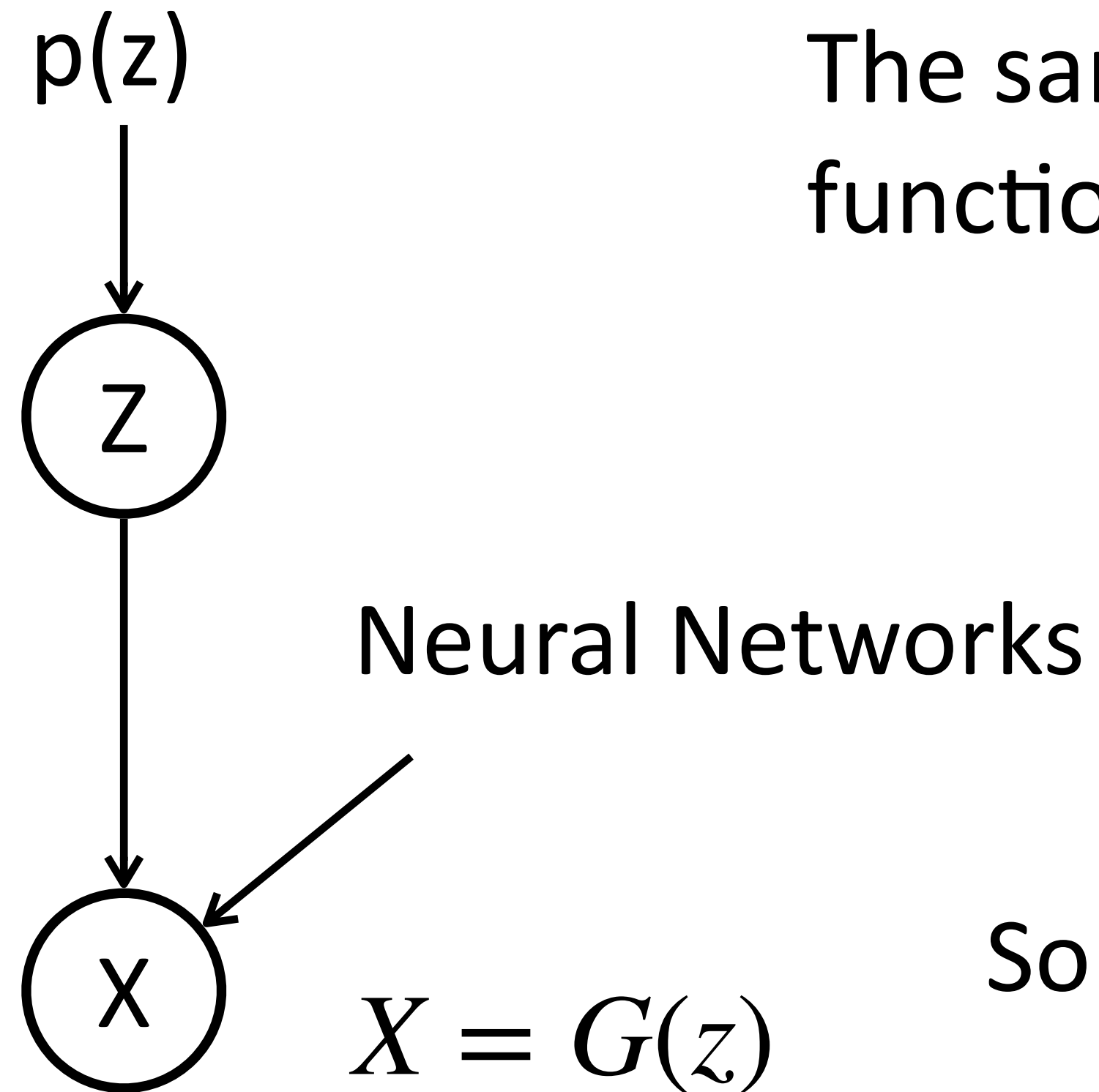


Generative Adversarial Nets

**Ian J. Goodfellow, Jean Pouget-Abadie*, Mehdi Mirza, Bing Xu, David Warde-Farley,
Sherjil Ozair†, Aaron Courville, Yoshua Bengio‡**
Département d'informatique et de recherche opérationnelle
Université de Montréal
Montréal, QC H3C 3J7

Generative Adversarial Networks

The GAN Model



The same as the VAE model, except that x is a deterministic function of z , but it can be a distribution as well

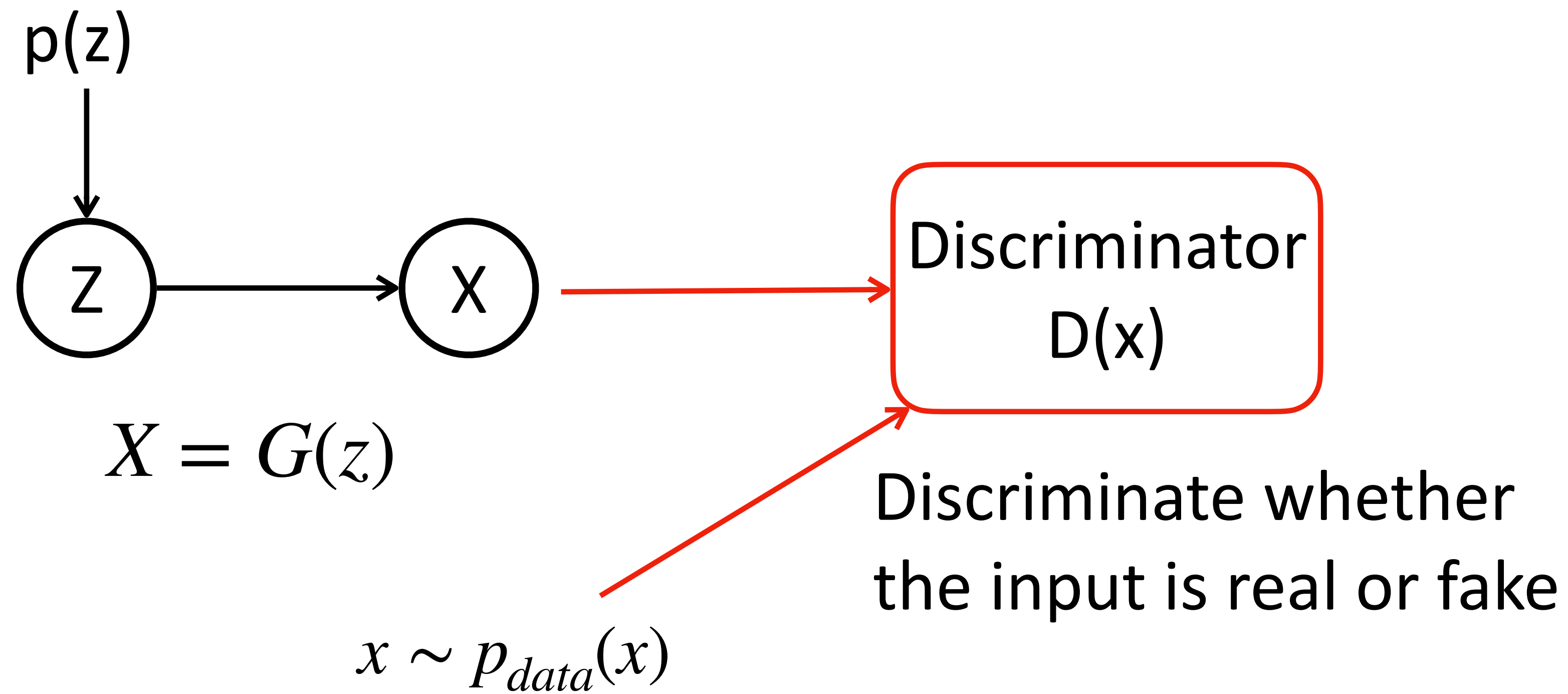
Can VAE use a deterministic $x = G(z)$?

Sometimes we call GANs *implicit* generative models

You can draw samples, but hard to evaluate $p(x)$

Training GANs

Computation Graph



1. Generator is trained to produce realistic examples to fool the discriminator
2. Discriminator is trained to discriminate real and fake examples

Training GANs

1. Generator is trained to produce realistic examples to fool the discriminator
2. Discriminator is trained to discriminate real and fake examples

The two objectives are against each other

Adversarial Game

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))].$$

Classification loss

$G(\mathbf{z})$ is trained to minimize the probability of $G(\mathbf{z})$ recognized as “fake” by D

$D(\mathbf{x})$ is trained with a standard classification loss

Training GANs

1. GAN is a new algorithm to train a common generative model (VAE as well)
2. GAN training is not MLE

Theory of GANs

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))] .$$

Proposition 1. For G fixed, the optimal discriminator D is

$$D_G^*(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_g(\mathbf{x})}$$

$$\begin{aligned} C(G) &= \max_D V(G, D) \\ &= \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log D_G^*(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}} [\log(1 - D_G^*(G(\mathbf{z})))] \\ &= \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log D_G^*(\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim p_g} [\log(1 - D_G^*(\mathbf{x}))] \\ &= \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \left[\log \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_g(\mathbf{x})} \right] + \mathbb{E}_{\mathbf{x} \sim p_g} \left[\log \frac{p_g(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_g(\mathbf{x})} \right] \end{aligned}$$

Theory of GANs

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))] .$$

Theorem 1. *The global minimum of the virtual training criterion $C(G)$ is achieved if and only if $p_g = p_{\text{data}}$. At that point, $C(G)$ achieves the value $-\log 4$.*

$$C(G) = -\log(4) + KL \left(p_{\text{data}} \left\| \frac{p_{\text{data}} + p_g}{2} \right. \right) + KL \left(p_g \left\| \frac{p_{\text{data}} + p_g}{2} \right. \right)$$

$$C(G) = -\log(4) + 2 \cdot JSD(p_{\text{data}} \| p_g)$$

Proposition 2. *If G and D have enough capacity, and at each step of Algorithm 1, the discriminator is allowed to reach its optimum given G , and p_g is updated so as to improve the criterion*

$$\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log D_G^*(\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim p_g} [\log(1 - D_G^*(\mathbf{x}))]$$

then p_g converges to p_{data}

Training GANs

Algorithm 1 Minibatch stochastic gradient descent training of generative adversarial nets. The number of steps to apply to the discriminator, k , is a hyperparameter. We used $k = 1$, the least expensive option, in our experiments.

for number of training iterations **do**

for k steps **do**

- Sample minibatch of m noise samples $\{z^{(1)}, \dots, z^{(m)}\}$ from noise prior $p_g(z)$.
- Sample minibatch of m examples $\{x^{(1)}, \dots, x^{(m)}\}$ from data generating distribution $p_{\text{data}}(x)$.
- Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[\log D(x^{(i)}) + \log \left(1 - D(G(z^{(i)})) \right) \right].$$

end for

- Sample minibatch of m noise samples $\{z^{(1)}, \dots, z^{(m)}\}$ from noise prior $p_g(z)$.
- Update the generator by descending its stochastic gradient:

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log \left(1 - D(G(z^{(i)})) \right).$$

end for

The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

Inner loop to update discriminator first

Training GANs

1. GAN is a new algorithm to train a common generative model (like VAE)
2. GAN training is not MLE **What is it then?**

Suppose the generator $G(x)$ is parameterized by θ , then what is the gradient when updating $G(x)$?

$$C(G) = -\log(4) + KL(p_{data} \parallel \frac{p_{data} + p_g^*}{2}) + KL(p_g \parallel \frac{p_{data} + p_g^*}{2})$$

p_g^* is from the solution of the discriminator, which is fixed when optimizing θ

$$\nabla_{\theta} C(G) = \nabla_{\theta} KL(p_g(x; \theta) \parallel \frac{p_{data} + p_g^*(x)}{2})$$

Training GANs

Recall that MLE is equivalent to minimizing $KL(p_{data}(x) || p_g(x))$

For GANs, the generator is to minimize $KL(p_g(x; \theta) || \frac{p_{data} + p_g^*(x)}{2})$

$$KL(p || q) \neq KL(q || p)$$

KL divergence is asymmetric, and GANs' KL divergence is in the opposite direction with respect to MLE

GANs v.s. VAEs

GANs are widely demonstrated to show superiority to VAEs on generating realistic, vivid images. In contrast, VAEs' generation is more blurred



GANs' generated images

GANs' generation can “miss mode” of the data distribution, where the generated images are not diverse to cover all the data distributions (VAEs do not have this issue)

Brock et al. LARGE SCALE GAN TRAINING FOR HIGH FIDELITY NATURAL IMAGE SYNTHESIS. ICLR 2019.

Implication of the KL divergence

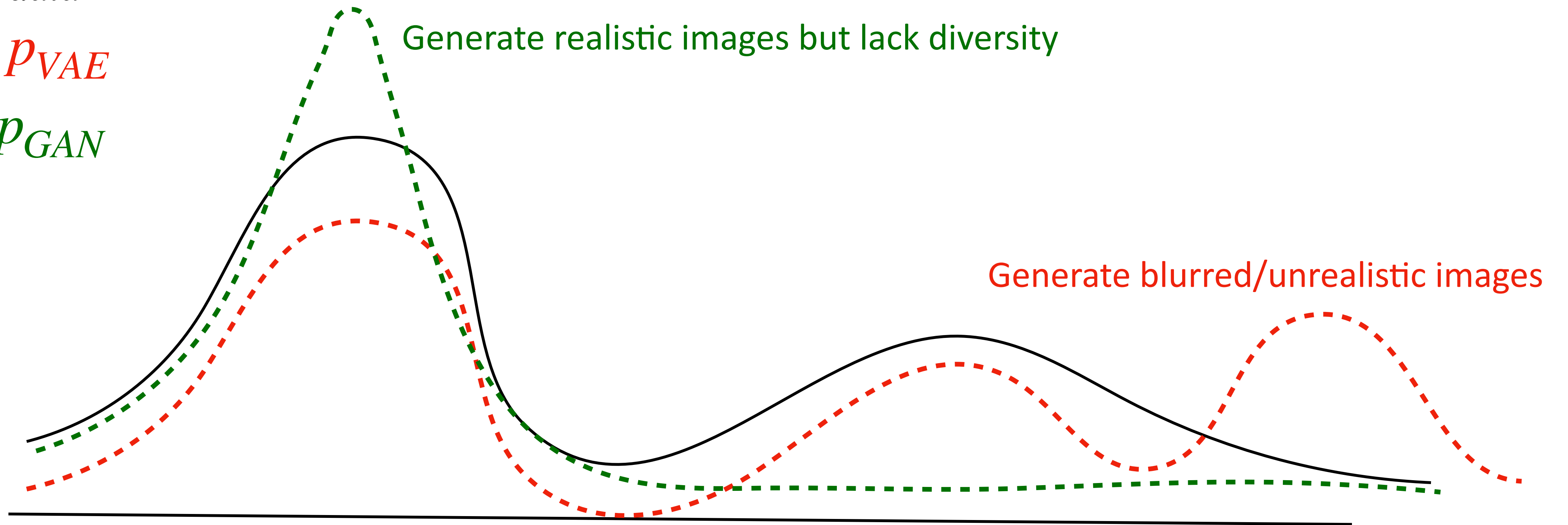
$$KL(p_{data}(x) || p_g(x)) \text{ v.s. } KL(p_g(x) || p_{data}(x))$$

VAEs GANs (approximately)

— P_{data}

- - - P_{VAE}

- - - P_{GAN}



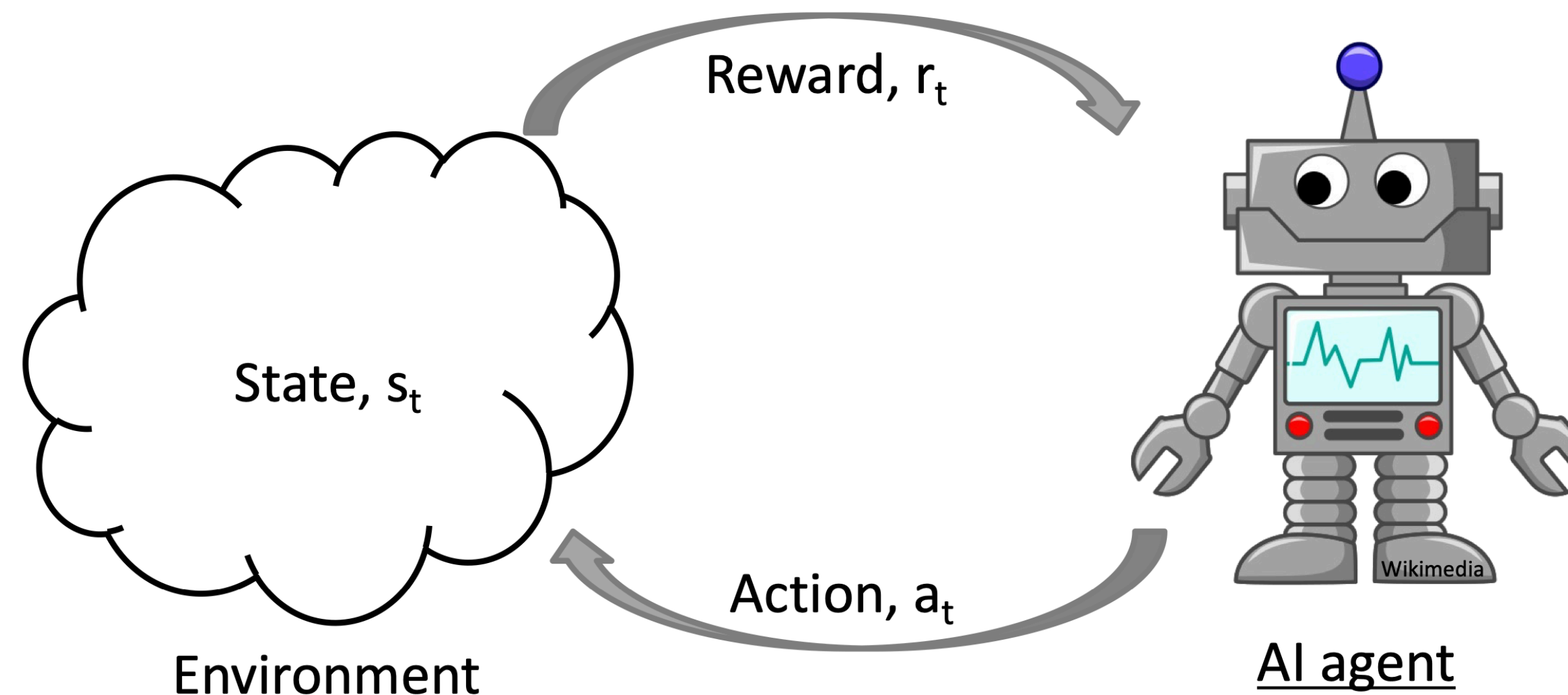
Reinforcement Learning

Learning Tasks

- Supervised learning - $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$
 - Regression - $y^{(i)} \in \mathbb{R}$
 - Classification - $y^{(i)} \in \{1, \dots, C\}$
- Unsupervised learning - $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^N$
 - Clustering
 - Dimensionality reduction
- Reinforcement learning - $\mathcal{D} = \{\mathbf{s}^{(t)}, \mathbf{a}^{(t)}, r^{(t)}\}_{t=1}^T$

RL Setup

In many cases, we cannot precisely define what the correct output is (think of we want to train a robot to walk)



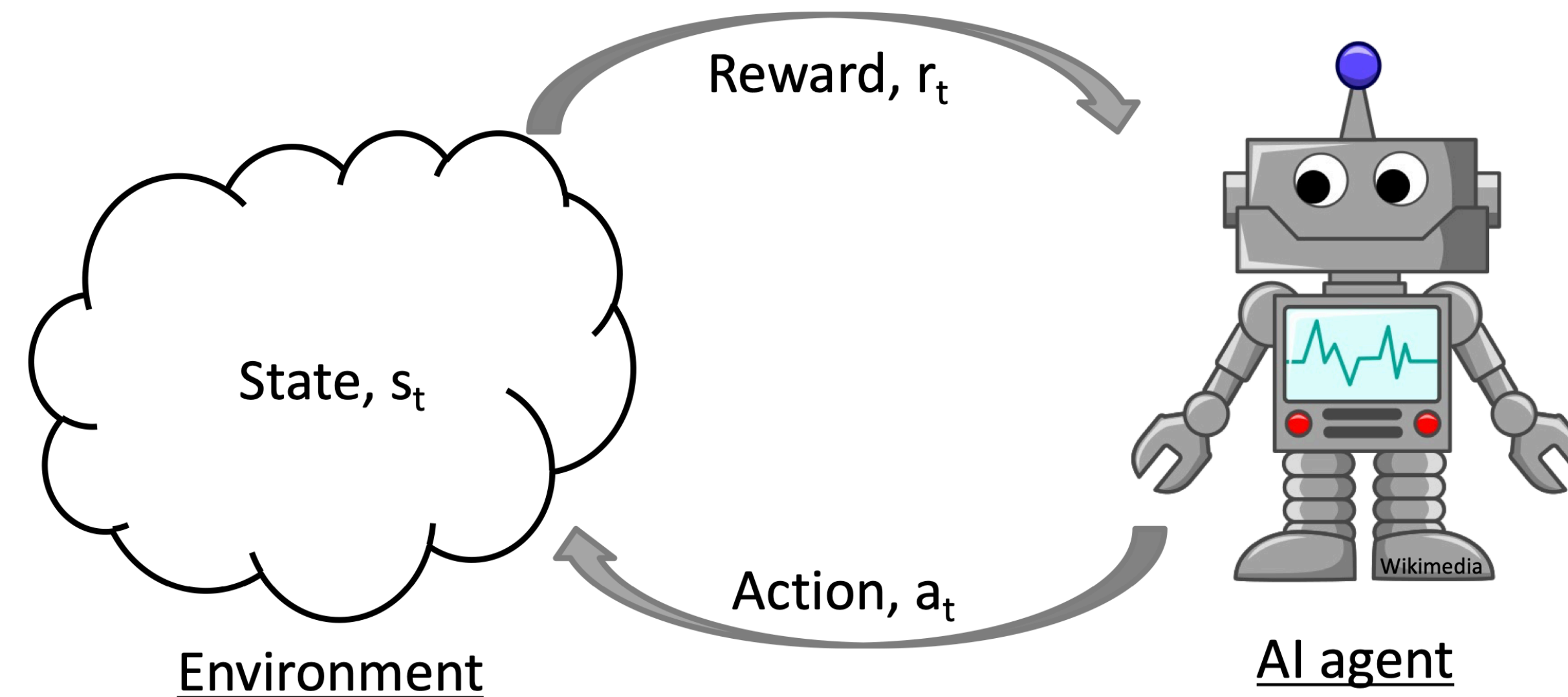
Agent chooses **actions** which can depend on past

Environment can change **state** with each action

Reward (Output) depends on (Inputs) action and state of environment

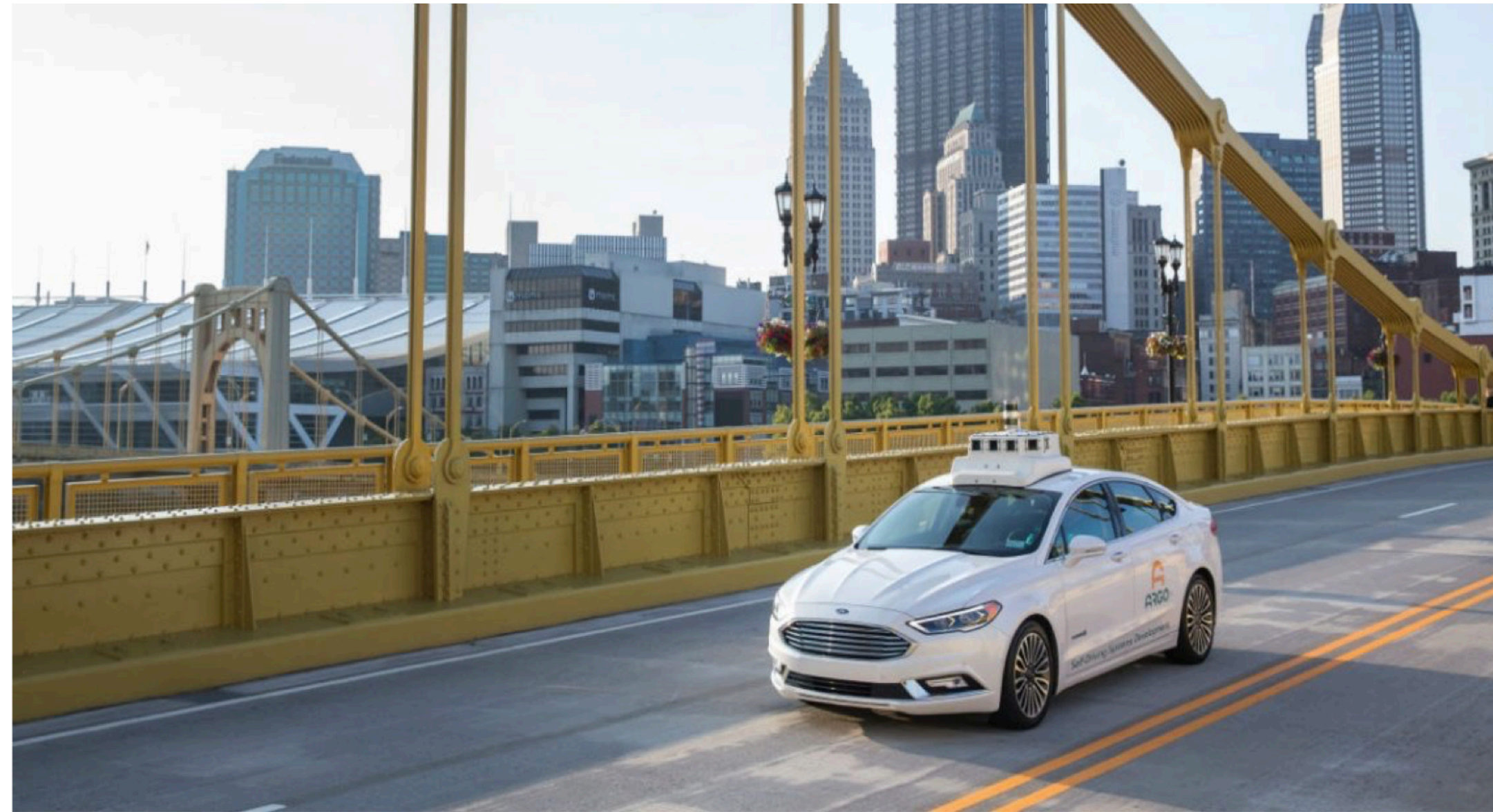
Goal: maximize the total reward

Differences from Supervised Learning



- Maximize reward (rather than learn reward) **Supervised training is like imitation**
- Inputs are not iid – state & action depends on past

RL Examples



RL Setup

- State space, \mathcal{S}
- Action space, \mathcal{A}
- Reward function
 - Stochastic, $p(r | s, a)$
 - Deterministic, $R: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$
- Transition function
 - Stochastic, $p(s' | s, a)$
 - Deterministic, $\delta: \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$

In this lecture, we assume they are known

- Reward and transition functions can be known or unknown

RL Setup

- Policy, $\pi : \mathcal{S} \rightarrow \mathcal{A}$
 - Specifies an action to take in *every* state
- Value function, $V^\pi : \mathcal{S} \rightarrow \mathbb{R}$
 - Measures the expected total reward of starting in some state s and executing policy π , i.e., in every state, taking the action that π returns

RL Example - gridworld

\mathcal{S} = all empty squares in the grid

\mathcal{A} = {up, down, left, right}

Deterministic transitions

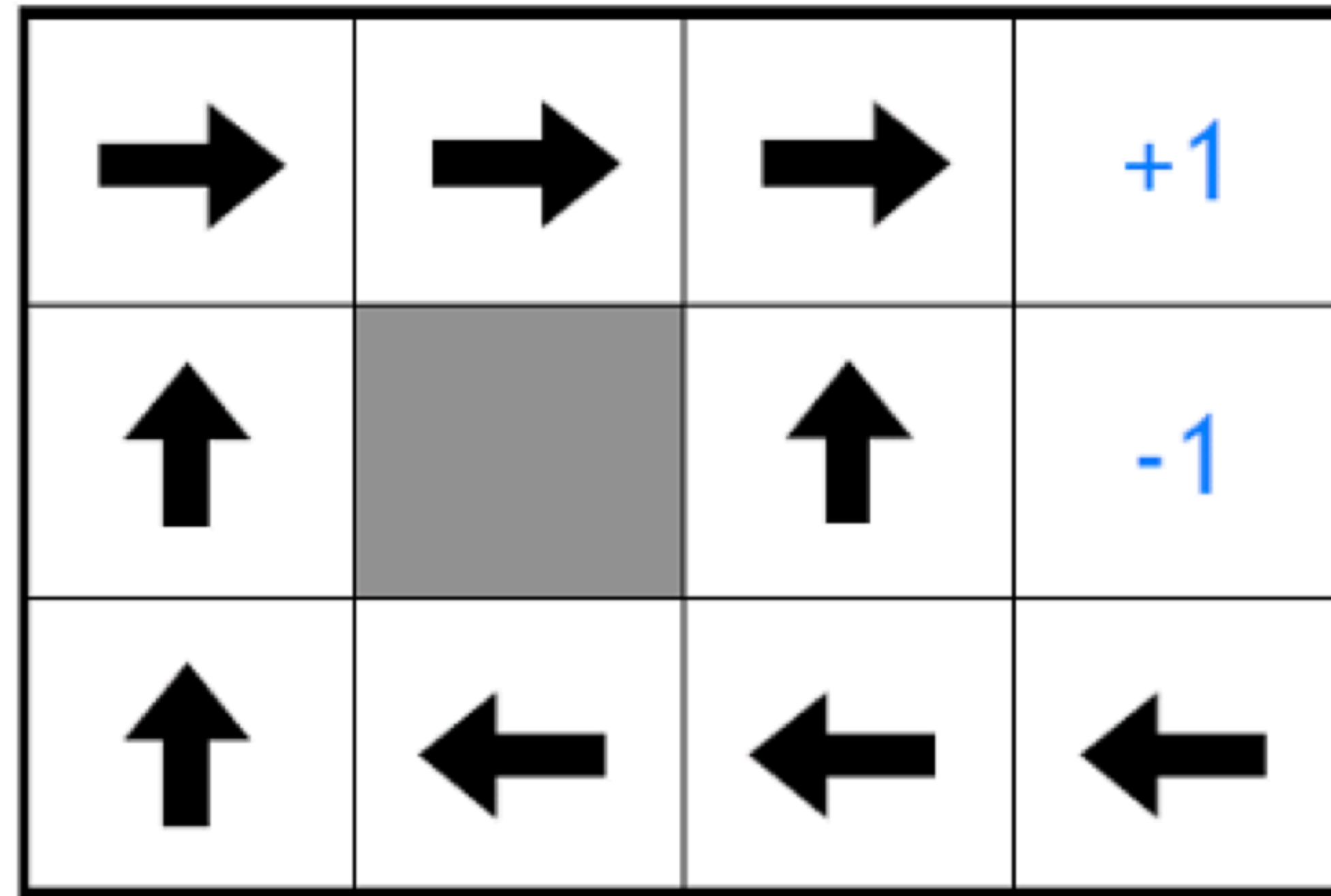
Rewards of +1 and -1 for entering the labelled squares

Terminate after receiving either reward

			+1
			-1

o

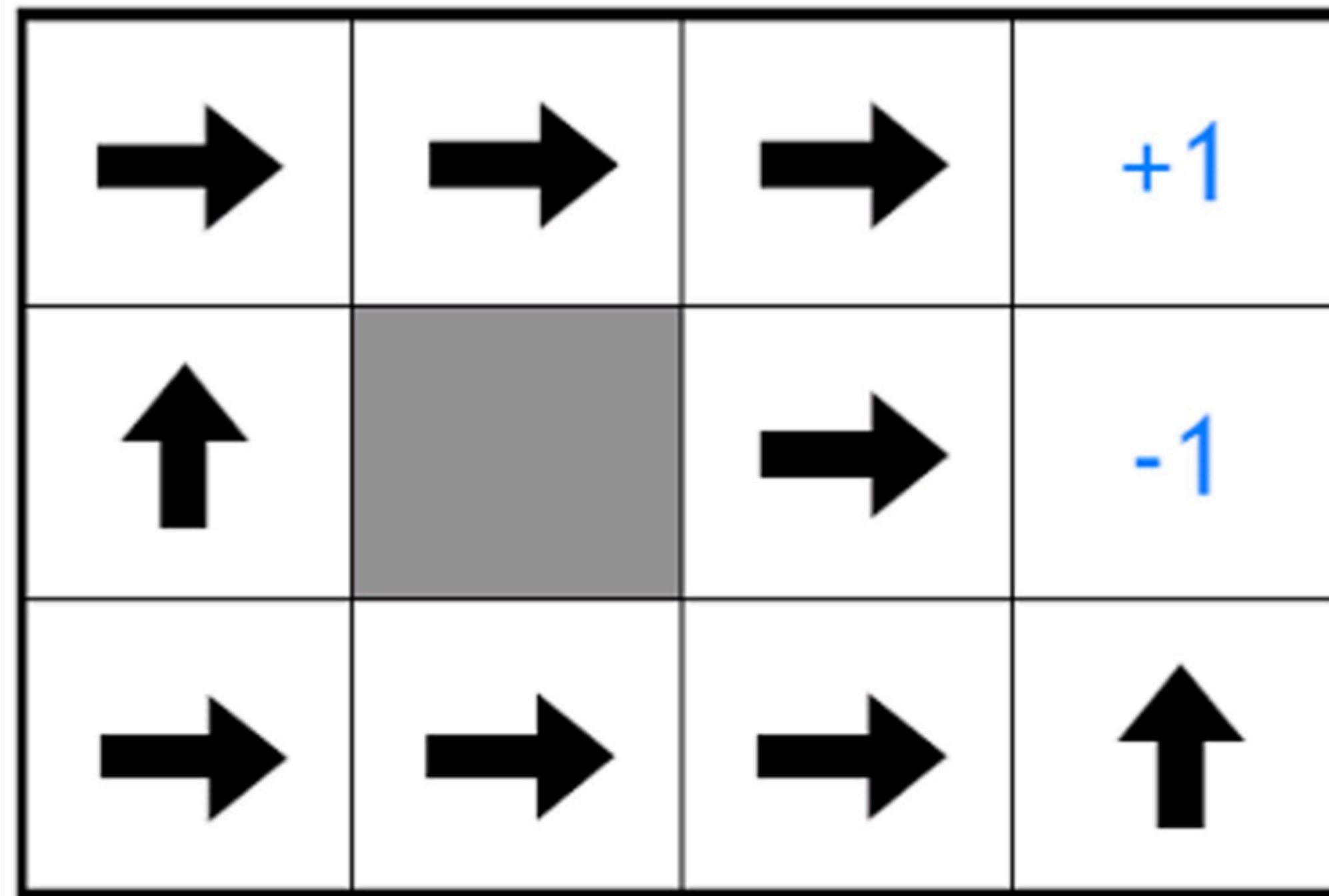
RL Example - gridworld



Is this policy optimal?

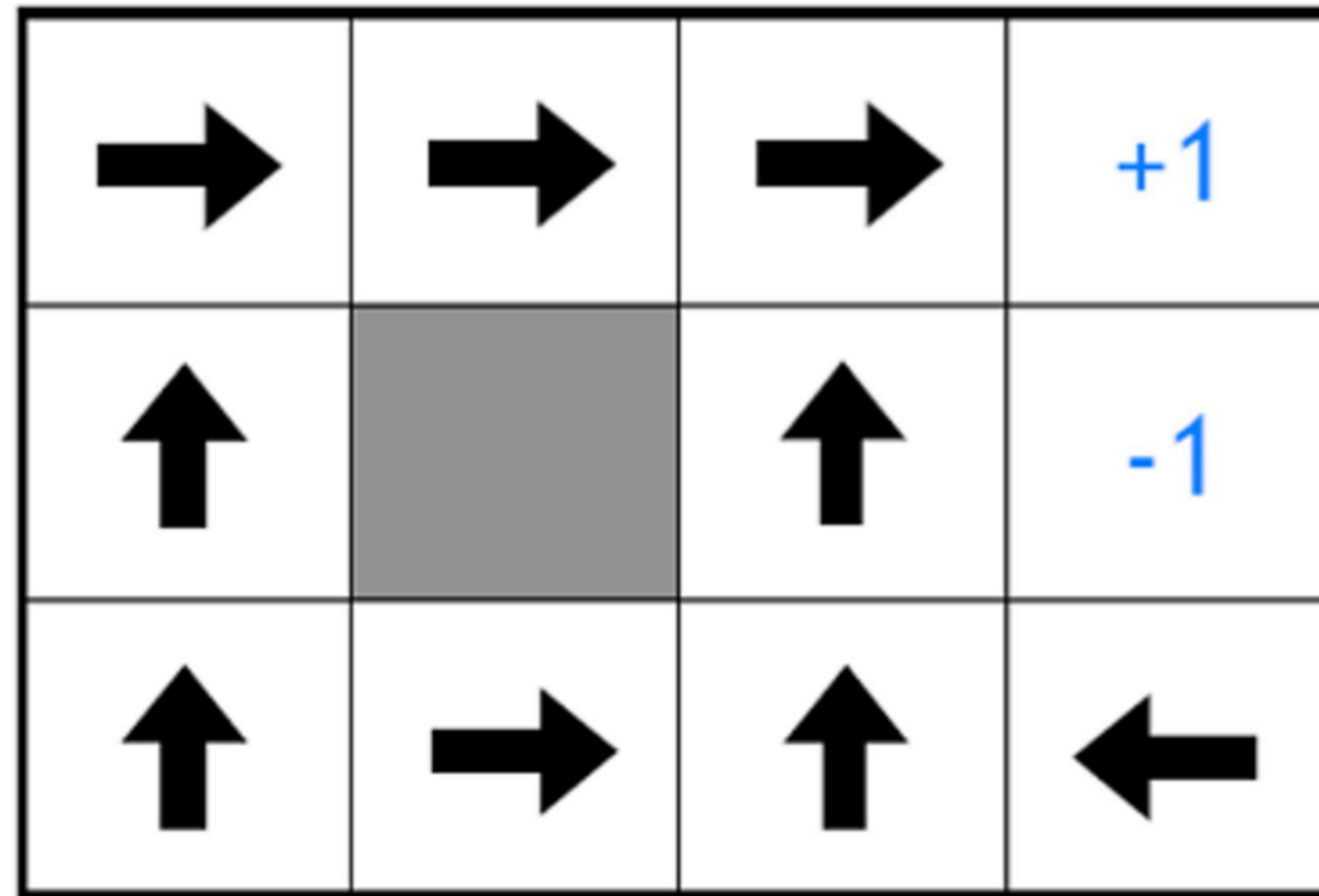
RL Example - gridworld

Optimal policy given a reward of -2 per step



RL Example - gridworld

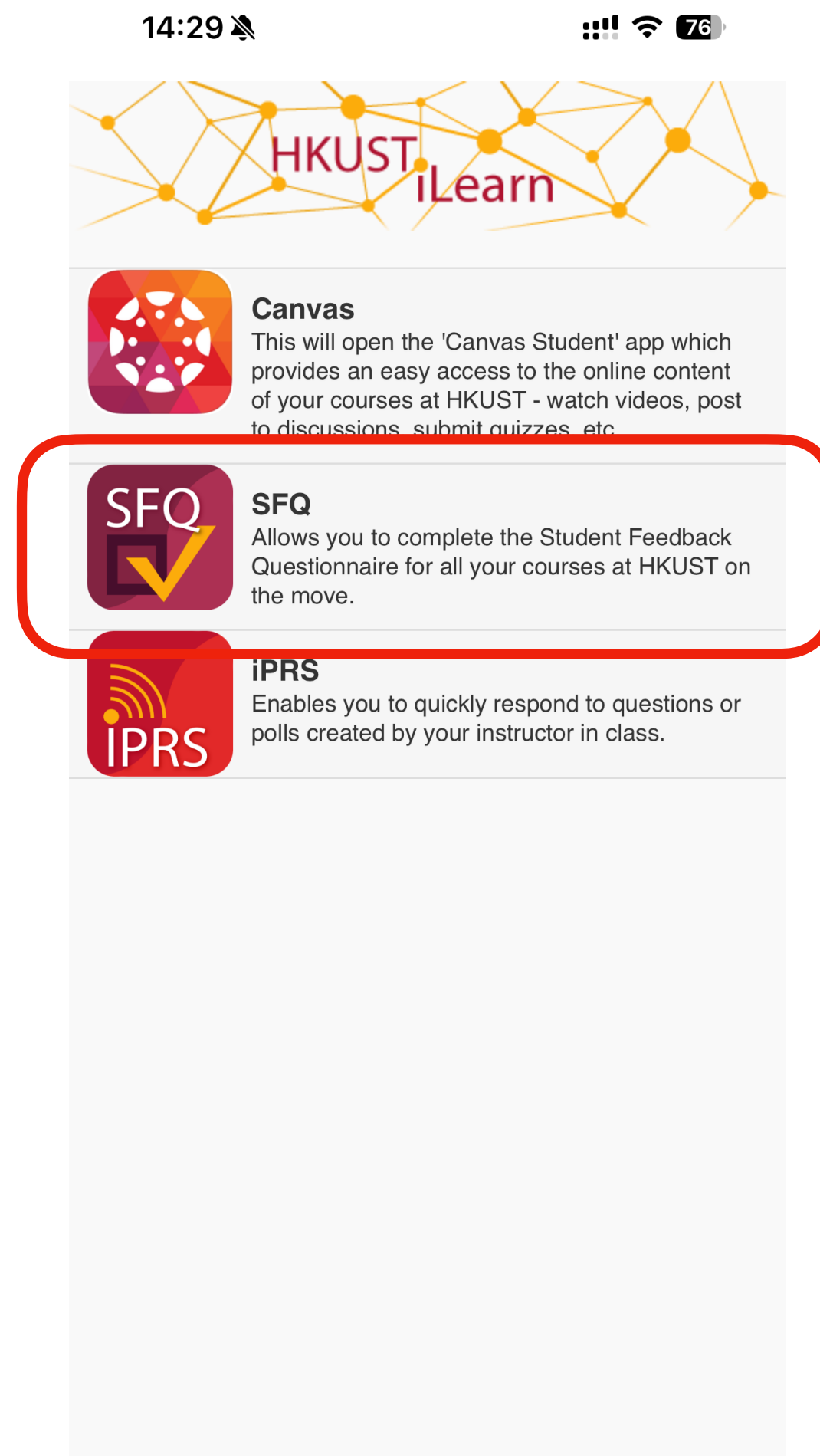
Optimal policy given a reward of -0.5 per step



What would be the algorithm to find the optimal policy automatically?

Course Evaluation

Anonymous to instructors



Or



<https://survey.ust.hk/hkust/>