

INFO 7390 – Advances in Data Sciences and Architecture Quiz Solutions

Student Name: _____
Professor: Nik Bear Brown

Rules:

1. Actual exam will be on HackerRank
2. Use Slack to ask questions of TAs

Q1 (5 Points) What is a policy in reinforcement learning? What is the difference between a deterministic policy and a stochastic policy?

Solution:

A policy, denoted as π (or sometimes $\pi(a|s)$), is a mapping from some state s to the probabilities of selecting each possible action given that state. For example, a greedy policy outputs for every state the action with the highest expected Q-Value.

- A **deterministic policy** is a mapping $\pi : \mathcal{S} \rightarrow \mathcal{A}$. For each state $s \in \mathcal{S}$, it yields the action $a \in \mathcal{A}$ that the agent will choose while in state s .
- A **stochastic policy** is a mapping $\pi : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$. For each state $s \in \mathcal{S}$ and action $a \in \mathcal{A}$, it yields the probability $\pi(a|s)$ that the agent chooses action a while in state s .

Q2 (5 Points) What is the difference between VALID and SAME padding?

Solution:

'Valid' padding means no padding. The output size of the convolutional layer shrinks depending on the input size & kernel size. On the contrary, 'same' padding means using padding.

Q3 (5 Points) Which of the following describes the policy that is greedy with respect to the Q-table? What action would a greedy policy take in states 1 and 2?

	action 1	action 2
state 1	1	2
state 2	4	3

Solution:

For state 1, action 2 has the highest estimated return ($2 > 1$). For state 2, action 1 has the highest estimated return ($4 > 3$).

Q4 (5 Points) What is the difference between a value function V and a Q function in reinforcement learning?

Solution:

Anything that says that one is a state-value function and the other is an action-value function is accepted for full credit.

$V_{\pi}(s)$ is the state-value function. The V value function states what the expected overall value of a state s under the policy π .

$Q\pi(s,a)$ is the action-value function. The Q function states what the value of a state s and an action a under the policy π .

The relationship between $Q\pi$ and $V\pi$ (the value of being in that state) is

$$V\pi(s) = \sum_{a \in A} \pi(a|s) * Q\pi(a,s)$$

Q5 (5 Points) Which of the values for epsilon yields an epsilon(ϵ)-greedy policy that is guaranteed to always select the greedy action? What is the range of values that epsilon(ϵ) can take? Which of the values for epsilon yields an epsilon(ϵ)-greedy policy that is guaranteed to always select the non-greedy action?

Solution:

You can think of the agent who follows an ϵ -greedy policy as always having a (potentially unfair) coin at its disposal, with probability ϵ of landing heads.

If the coin is greater or equal to epsilon(ϵ) it will take the max (greedy) q-value so if the ϵ is 0 the agent will always take the max (greedy) q-value.

There is NO VALUE of epsilon yields an epsilon(ϵ)-greedy policy that is guaranteed to **always** select the non-greedy action as there is always the possibility that the coin returns a 1.

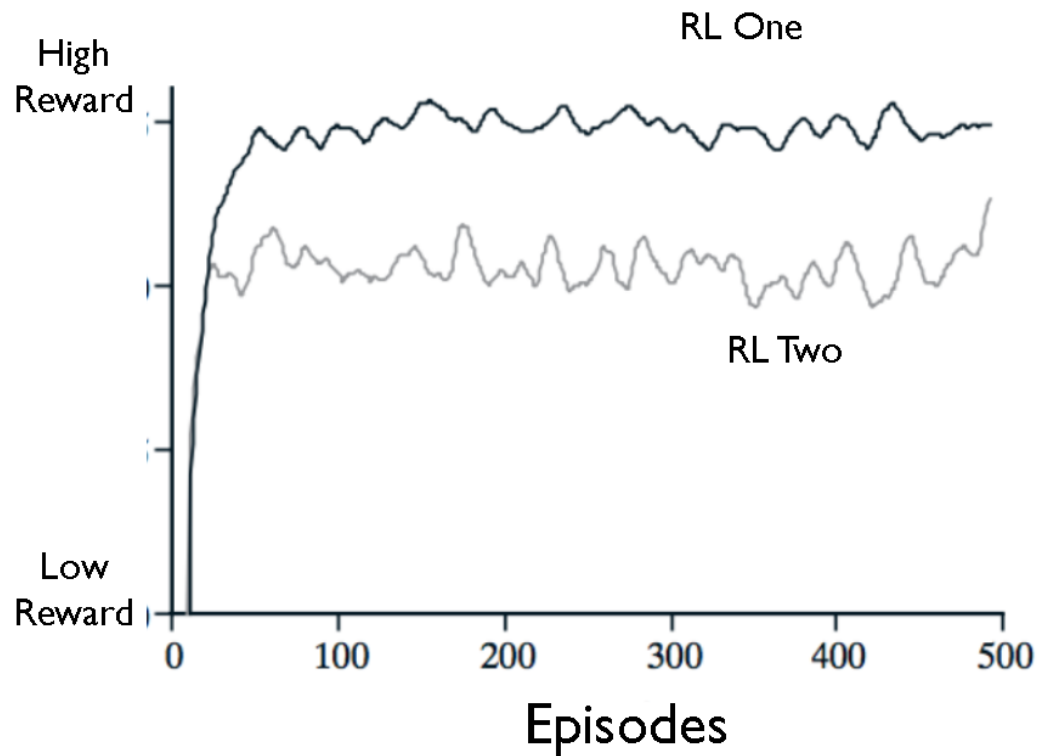
Note that ϵ must always be a value between 0 and 1, inclusive (that is, $\epsilon \in [0,1]$).

Q6 (5 Points) Given a uniform selection the actions in exploration, which of the values for epsilon yields an epsilon-greedy policy that is equivalent an approximately equiprobable random policy (where, from each state, each action is equally likely to be selected)?

Solution:

An $\epsilon = 1$ would mean that nearly all of the time the agent is exploring.

Q7 (5 Points) Which reinforcement learning algorithm is doing better? RL One (top) or RL Two (bottom)? Why? Do the algorithms plateau? If so, where?



Solution:

RL One (top) is clearly better as its plateau for average reward is well above RL Two (bottom) once the plateau starts around 50 episodes or so.

Q8 (5 Points) What steps can we take to prevent overfitting in a neural network?

Solution:

Note this was a practice quiz question to check if one reads the practice quiz.

Data Augmentation
Weight Sharing
Early Stopping
Dropout

Q9 (5 Points) In OpenAI- Gym's Taxi-v2 environment there are 4 locations (labeled by different letters) and your job is to pick up the passenger at one location and drop him off in another. You receive +20 points for a successful drop-off and lose 1 point for every timestep it takes. There is also a 10-point penalty for illegal pick-up and drop-off actions. Assume the environment is on a 5x5 grid and the taxi can move up, down, right, left, wait, pick up passenger or drop passenger. How many states are there? How many actions are there? What is the size of the Q-table?

Solution:

How many states are there?

There are 5x5 or 25 states.

How many actions are there?

There are 7 actions, move up, down, right, left, wait, pick up passenger or drop passenger.

The Q-table is 25x7 or 175.

Q10 (5 Points) Which of the following state spaces are continuous and which are discrete?

- A. Playing cards with a standard 52-card deck.
- B. GPS coordinates for self-driving cars.
- C. Force applied to the arm of an industrial robot.
- D. Board positions in a standard tic-tack-toe game.
- E. Keys to play on a keyboard.

Solution:

Discrete (A, D, E)

Playing cards with a standard 52-card deck.

Board positions in a standard tic-tack-toe game.

Keys to play on a keyboard.

Continuous (B, C)

GPS coordinates for self-driving cars.

Force applied to the arm of an industrial robot.

Q11 (10 Points) Describe the structure of a Convolutional neural network (CNN). Is it supervised or unsupervised? How is subsampling done in CNNs?

Solution:

A Convolutional neural network (CNN) is supervised (the output is labeled). Subsampling is usually done by pooling.

Input->Conv->RELU (or another non-linear activation function) -> Pool (usually Max-pool)
the above sequence can be repeated for a while until it goes to a MLP structure

fc->RELU (or another non-linear activation function) for each hidden layer then possibly a normalizing activation like softmax then output layer.

Note some consider the activation function to be part of fc, so something like fc-> fc-> fc->Softmax (maybe)->output is fine for full credit.

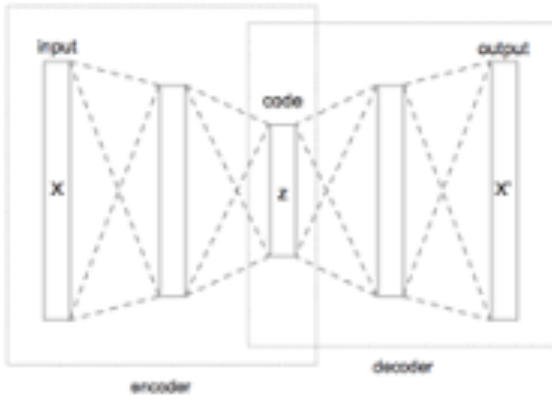
Q12 (10 Points) Describe the structure of an Autoencoder (AE). Is it supervised or unsupervised?

Solution:

An Autoencoder (AE) is unsupervised.

An **autoencoder** consists of two parts, the encoder and the decoder, which can be defined as transitions and such that: In the simplest case, given one hidden layer, the encoder stage of an **autoencoder** takes the input and maps it to : This image is usually referred to as code, latent variables, or latent representation. The latent layer is the basis of the latent representation.

No picture is required. Just words.



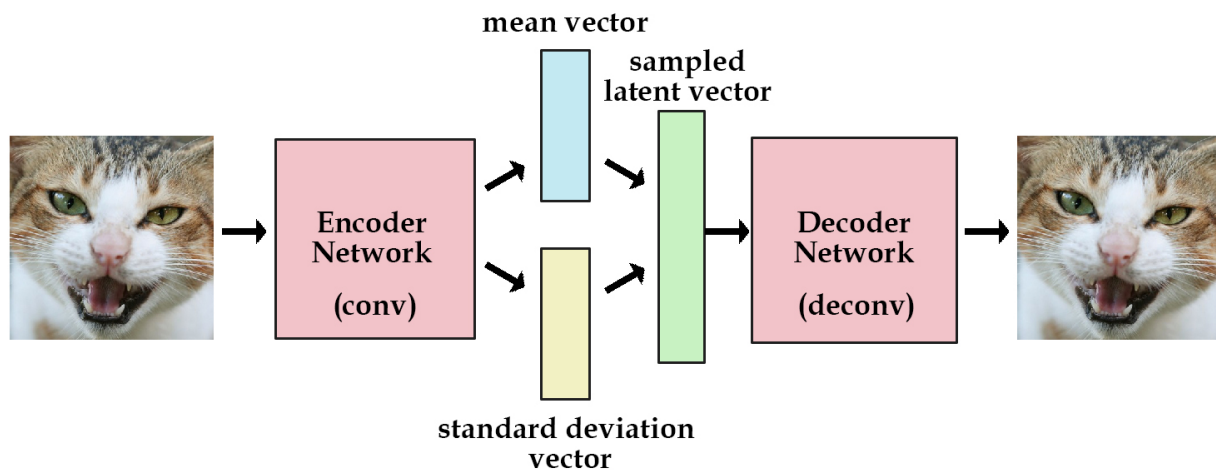
Q13 (10 Points) Describe the structure of a Variational Autoencoder (VAE). Is it supervised or unsupervised? How are the parameters of any probabilistic distributions represented in its structure?

Solution:

A Variational Autoencoder (VAE) is unsupervised the input layer is the same data as the output layer.

From the input layer there is an encoder network of one or more layers and before the output layer there is a corresponding decoder network. The latent layer is the basis of the latent representation. Unlike a deterministic autoencoder, additional layers which allow for the estimation of probability distributions are added. For example, if a Gaussian distribution is used layers that represent the parameters of a Gaussian distribution, the mean and standard deviation would be used.

No picture is required. Just words.



Q14 (10 Points) Describe the structure of a Generative Adversarial Network (GAN). Is it supervised or unsupervised?

Solution:

GANs are unsupervised learning algorithms that use a supervised loss as part of the training.

A generative adversarial network (GAN) has two parts:

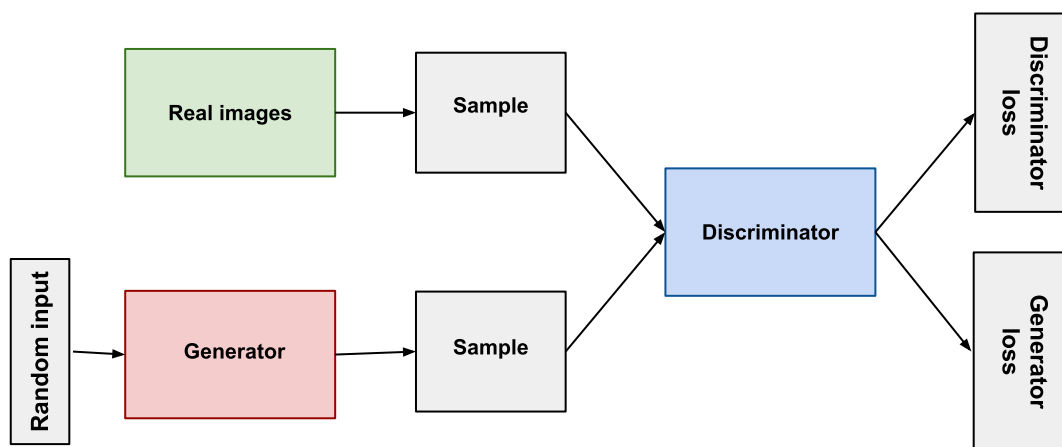
- The **generator** learns to generate plausible data. The generated instances become negative training examples for the discriminator.
- The **discriminator** learns to distinguish the generator's fake data from real data. The discriminator penalizes the generator for producing implausible results.

When training begins, the generator produces obviously fake data, and the discriminator quickly learns to tell that it's fake. As training progresses, the generator gets closer to producing output that can fool the discriminator: Finally, if generator training goes well, the discriminator gets worse at telling the difference between real and fake. It starts to classify fake data as real, and its accuracy decreases.

Both the generator and the discriminator are neural networks. The generator output is connected directly to the discriminator input.

Through backpropagation, the discriminator's classification provides a signal that the generator uses to update its weights.

No picture is required. Just words.:



See https://developers.google.com/machine-learning/gan/gan_structure

Q15 (10 Points) Assume that in building a deep reinforcement learning for games like pong or space invaders the pixels on a single screen capture are being used as input and an MLP for the deep learning. How could one input states that captured the motion? How could one input states identified the things in the screen rather than individual pixels?

Solution:

To capture motion, one would have to input several frames rather than a single frame. Say the last few frames.

To identify the things on the screen one would need to use more sophisticated deep learning that detected and localized the elements in an image such as deep learning based semantic segmentation.