

THE UNIVERSITY OF BRITISH COLUMBIA

Topics in AI (CPSC 532L): **Multimodal Learning with Vision, Language and Sound**

Lecture 12: Deep Reinforcement Learning



Types of Learning

Supervised training

- Learning from the teacher
- Training data includes desired output

Unsupervised training

Training data does not include desired output

Reinforcement learning

Learning to act under evaluative feedback (rewards)

What is Reinforcement Learning

Agent-oriented learning — learning by interacting with an environment to achieve a goal

 More realizing and ambitious than other kinds of machine learning

Learning by trial and error, with only delayed evaluative feedback (reward)

- The kind go machine learning most like natural learning
- Learning that can tell for itself when it is right or wrong





Example: Hajime Kimura's RL Robot







After

Example: Hajime Kimura's RL Robot







After

Example: Hajime Kimura's RL Robot







After

Challenges of RL

- Evaluative feedback (reward)
- Sequentiality, delayed consequences
- Need for trial and error, to explore as well as exploit
- Non-stationarity
- The fleeting nature of time and online data

How does RL work?



At each step t the agent:

- Executes action a_t
- Receives observation o_t
- Receives scalar reward r_t
- The environment:
 - Receives action a_t
 - Emits observation o_{t+1}
 - Emits scalar reward r_{t+1}

Robot Locomotion



Objective: Make the robot move forward

State: Angle and position of the joints Action: Torques applied on joints **Reward**: 1 at each time step upright + forward movement



Atari Games



Objective: Complete the game with the highest score

State: Raw pixel inputs of the game stateAction: Game controls e.g. Left, Right, Up, DownReward: Score increase/decrease at each time step

Go Game (AlphaGo)



Objective: Win the game!

State: Position of all pieces Action: Where to put the next piece down **Reward**: 1 if win at the end of the game, 0 otherwise

Markov Decision Processes

- Mathematical formulation of the RL problem

Defined by:

- S: set of possible states
- \mathcal{A} : set of possible actions
- \mathcal{R} : distribution of reward given (state, action) pair
- \mathbb{P} : transition probability i.e. distribution over next state given (state, action) pair
- γ : discount factor

ir t state given (state, action) pair

Markov Decision Processes

— Mathematical **formulation** of the RL problem

Defined by:

- *S* : set of possible states
- \mathcal{A} : set of possible actions
- \mathcal{R} : distribution of reward given (state, action) pair
- ■: transition probability i.e. distribution over next state given (state, action) pair
- γ : discount factor

- Life is trajectory: $...S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1}, R_{t+2}, S_{t+2}, ...$

Markov Decision Processes

- Mathematical formulation of the RL problem

Defined by:

- S: set of possible states
- \mathcal{A} : set of possible actions
- \mathcal{R} : distribution of reward given (state, action) pair
- $oldsymbol{\gamma}$: discount factor

- Life is trajectory: $\ldots S_t, A_t, R_{t+1}, S_t$

Markov property: Current state completely characterizes the state of the world

$$p(r, s'|s, a) = Prob\Big[R_{t+1} = r, S_{t+1} = s' \mid S_t = s, A_t = a\Big]$$

ir kt state given (state, action) pair

$$A_{t+1}, A_{t+1}, R_{t+2}, S_{t+2}, \ldots$$

Components of the RL Agent

Policy

— How does the agent behave?

Value Function

— How good is each state and/or action pair?

Model

Agent's representation of the environment

Policy

- The policy is how the agent acts
- Formally, map from states to actions:
 - **Deterministic** policy: $a = \pi(s)$ **Stochastic** policy: $\pi(a|s) = \mathbb{P}[A_t = a|S_t = s]$

Policy

- The policy is how the agent acts
- Formally, map from states to actions:
 - **Deterministic** policy: $a = \pi(s)$ **Stochastic** policy: $\pi(a|s) = \mathbb{P}[A_t = a|S_t = s]$



What is a good policy?

What is a good policy?

Maximizes current reward? Sum of all future rewards?

What is a good policy?

Maximizes current reward? Sum of all future rewards?

Discounted future rewards!

What is a good policy?

Maximizes current reward? Sum of all future rewards?

Discounted future rewards!

Formally: $\pi^* = \arg \max_{\pi} \mathbb{E} \left[\sum_{t \ge 0} \gamma^t r_t | \pi \right]$

with $s_0 \sim p(s_0), a_t \sim \pi(\cdot | s_t), s_{t+1} \sim p(\cdot | s_t, a_t)$

Components of the RL Agent



— How does the agent behave?

Value Function

— How good is each state and/or action pair?

Model

Agent's representation of the environment

Value Function

A value function is a prediction of future reward

- "State Value Function" or simps "Value Function"
 - How good id a state?
 - Am I screwed? Am I winning this game?

"Action Value Function" or **Q-function**

- How good is a state action-pair?
- Should I do this now?

Value Function and Q-value Function

Following a policy produces sample trajectories (or paths) s_0 , a_0 , r_0 , s_1 , a_1 , r_1 , ...

— The value function (how good is the state) at state s, is the expected cumulative reward from state s (and following the policy thereafter):

$$V^{\pi}(s) = \mathbb{E}\left[\sum_{t\geq 0} \gamma^t r_t | s_0 = s, \pi
ight]$$

Value Function and Q-value Function

Following a policy produces sample trajectories (or paths) s_0 , a_0 , r_0 , s_1 , a_1 , r_1 , ...

- The value function (how good is the state) at state s, is the expected cumulative reward from state s (and following the policy thereafter):

— The **Q-value function** (how good is a state-action pair) at state s and action a, is the expected cumulative reward from taking action a in state s (and following) the policy thereafter):



$$V^{\pi}(s) = \mathbb{E}\left[\sum_{t\geq 0} \gamma^t r_t | s_0 = s, \pi
ight]$$

$$\sum_{t\geq 0}\gamma^t r_t |s_0=s,a_0=a,\pi$$

Components of the RL Agent

VPolicy

— How does the agent behave?

Value Function

— How good is each state and/or action pair?

Model

Agent's representation of the environment

Model

Model predicts what the world will do next



Model

Model predicts what the world will do next



Components of the RL Agent

VPolicy

— How does the agent behave?

Value Function

— How good is each state and/or action pair?



Agent's representation of the environment

Maze Example



Reward: -1 per time-step Actions: N, E, S, W States: Agent's location

Maze Example: Policy



Arrows represent a policy $\pi(s)$ for each state S

Maze Example: Value

									l
		-14	-13	-12	-11	-10	-9		
Start	-16	-15			-12		-8		
		-16	-17			-6	-7		
			-18	-19		-5			
		-24		-20		-4	-3		
		-23	-22	-21	-22		-2	-1	

Goal

Numbers represent value $\, \mathrm{v}_{\pi}(s) \,$ of each state S

Maze Example: Model



Goal

Grid layout represents transition model

Numbers represent the immediate reward for each state (same for all states)

Components of the RL Agent

Policy

— How does the agent behave?

Value Function

— How good is each state and/or action pair?

Model

Agent's representation of the environment

Approaches to RL: Taxonomy

Model-free RL

Value-based RL

- Estimate the optimal action-value function $Q^*(s, a)$
- No policy (implicit)

Policy-based RL

- Search directly for the optima policy π^*
- No value function

Model-based RL

- Build a model of the world
- Plan (e.g., by look-ahead) using model



Approaches to RL: Taxonomy

Model-free RL

Value-based RL

- Estimate the optimal action-value function $Q^*(s, a)$
- No policy (implicit)

Policy-based RL

- Search directly for the optima policy π^*
- No value function

Model-based RL

- Build a model of the world
- Plan (e.g., by look-ahead) using model

Actor-critic RL

- Value function

Policy function


Deep RL

Value-based RL

Policy-based RL

 $\pi_{ heta}$ Use neural nets to represent the policy

Model-based RL

- Use neural nets to represent and learn the model

- Use neural nets to represent Q function $Q(s,a;\theta)$ $Q(s,a;\theta^*) \approx Q^*(s,a)$

$\pi_{\theta^*} \approx \pi^*$

* slide from Dhruv Batra

Approaches to RL

Value-based RL

- Estimate the optimal action-value function $Q^*(s, a)$
- No policy (implicit)

* slide from Dhruv Batra

Optimal Q-function is the maximum achievable value

 $Q^*(s,a) = \max_{\pi} Q^{\pi}(s,a) = Q^{\pi^*}(s,a)$



Optimal Q-function is the maximum achievable value

Once we have it, we can act optimally

 $Q^*(s,a) = \max_{\pi} Q^{\pi}(s,a) = Q^{\pi^*}(s,a)$

 $\pi^*(s) = \operatorname{argmax} Q^*(s, a)$



- Optimal Q-function is the maximum achievable value
 - $Q^*(s, a) = \max_{\pi} Q^{\pi}(s, a) = Q^{\pi^*}(s, a)$
- Once we have it, we can act optimally

$$\pi^*(s) = rgmargmatrix _a$$

Optimal value maximizes over all future decisions

$$Q^*(s,a) = r_{t+1} + \gamma \max_{a_{t+1}} q_{t+1}$$

 $ax r_{t+2} + \gamma^2 \max_{a_{t+2}} r_{t+3} + \dots$ $= r_{t+1} + \gamma \max_{a_{t+1}} Q^*(s_{t+1}, a_{t+1})$

 $\max Q^*(s,a)$

* slide from David Silver

- Optimal Q-function is the maximum achievable value
 - $Q^*(s,a) = \max_{\pi} Q^{\pi}(s,a) = Q^{\pi^*}(s,a)$
- Once we have it, we can act optimally

Optimal value maximizes over all future decisions

$$Q^*(s,a) = r_{t+1} + \gamma \max_{a_{t+1}} q_{t+1}$$

 $ax r_{t+2} + \gamma^2 \max_{a_{t+2}} r_{t+3} + \dots$ $= r_{t+1} + \gamma \max_{a_{t+1}} Q^*(s_{t+1}, a_{t+1})$

Formally, Q* satisfied Bellman Equations $Q^*(s,a) = \mathbb{E}_{s'} \left[r + \gamma \max_{a'} Q^*(s',a') \mid s,a \right]$

 $\max Q^*(s,a)$

* slide from David Silver

 Q_i will converge to Q^* as $i \rightarrow$ infinity

 $Q_{i+1}(s,a) = \mathbb{E}\left[r + \gamma \max_{a'} Q_i(s',a')|s,a\right]$

$$Q_{i+1}(s,a) = \mathbb{E}\left[r\right]$$

 Q_i will converge to Q^* as $i \rightarrow$ infinity



 $r + \gamma \max_{a'} Q_i(s', a') |s, a|$

What's the problem with this?

$$Q_{i+1}(s,a) = \mathbb{E}\left[r\right]$$

 Q_i will converge to Q^* as $i \rightarrow$ infinity

What's the problem with this?

Not scalable. Must compute Q(s,a) for every state-action pair. If state is e.g. game pixels, computationally infeasible to compute for entire state space!

 $r + \gamma \max_{a'} Q_i(s', a') |s, a]$

$$Q_{i+1}(s,a) = \mathbb{E}\left[r\right]$$

 Q_i will converge to Q^* as $i \rightarrow$ infinity

What's the problem with this?

Not scalable. Must compute Q(s,a) for every state-action pair. If state is e.g. game pixels, computationally infeasible to compute for entire state space!

Solution: use a function approximator to estimate Q(s,a). E.g. a neural network!

 $r + \gamma \max_{a'} Q_i(s', a') |s, a|$

Q-Networks



$Q(s, a, \mathbf{w}) \approx Q^*(s, a)$



* slide from David Silver



Case Study: Playing Atari Games



Objective: Complete the game with the highest score

State: Raw pixel inputs of the game state Action: Game controls e.g. Left, Right, Up, Down **Reward**: Score increase/decrease at each time step

[Mnih *et al.*, 2013; Nature 2015]



$Q(s, a; \theta)$: neural network with weights θ







Current state st: 84x84x4 stack of last 4 frames (after RGB->grayscale conversion, downsampling, and cropping)

[Mnih *et al.*, 2013; Nature 2015]



$Q(s, a; \theta)$: neural network with weights θ







Current state st: 84x84x4 stack of last 4 frames (after RGB->grayscale conversion, downsampling, and cropping)

[Mnih *et al.*, 2013; Nature 2015]



$Q(s, a; \theta)$: neural network with weights θ







[Mnih *et al.*, 2013; Nature 2015]

Current state st: 84x84x4 stack of last 4 frames (after RGB->grayscale conversion, downsampling, and cropping)



$Q(s, a; \theta)$: neural network with weights θ







Current state Stat (after RGB->grayscale conversion, downsampling, and cropping)

[Mnih et al., 2013; Nature 2015]

Last FC layer has 4-d output (if 4 actions), corresponding to $Q(s_t, a_1)$, $Q(s_t, a_2), Q(s_t, a_3), Q(s_t, a_4)$



$Q(s, a; \theta)$: neural network with weights θ







[Mnih et al., 2013; Nature 2015]

Last FC layer has 4-d output (if 4 actions), corresponding to $Q(s_t, a_1)$, $Q(s_t, a_2), Q(s_t, a_3), Q(s_t, a_4)$

Number of actions between 4-18 depending on Atari game

Current state Stat (after RGB->grayscale conversion, downsampling, and cropping)



$Q(s, a; \theta)$: neural network with weights θ

A single feedforward pass to compute Q-values for all actions from the current state => efficient!



Mnih *et al.*, 2013; Nature 2015]

Last FC layer has 4-d output (if 4 actions), corresponding to $Q(s_t, a_1)$, $Q(s_t, a_2), Q(s_t, a_3), Q(s_t, a_4)$

Number of actions between 4-18 depending on Atari game

Current state st: 84x84x4 stack of last 4 frames (after RGB->grayscale conversion, downsampling, and cropping)



Remember: want to find a Q-function that satisfies the Bellman Equation: $Q^*(s, a) = \mathbb{E}[r + \gamma \max_{a'} Q^*(s', a') \mid s, a]$

Remember: want to find a Q-function that satisfies the Bellman Equation: $Q^*(s, a) = \mathbb{E}[r + \gamma \max_{a'} Q^*(s', a') \mid s, a]$

Forward Pass:

Loss function: $L_i(\theta_i) = \mathbb{E} \left| (y_i - Q(s, a; \theta_i)^2) \right|$ $y_i = \mathbb{E}[r + \gamma \max_{a'} Q^*(s', a') \mid s, a]$ where

 $Q^*(s,a) = \mathbb{E}[r + \gamma \max_{a'} Q^*(s',a') \mid s,a]$

Forward Pass:

Loss function: $L_i(\theta_i) = \mathbb{E} \left[(y_i - Q(s, a; \theta_i)^2) \right]$ $y_i = \mathbb{E}[r + \gamma \max_{a'} Q^*(s', a') \mid s, a]$ where

Backward Pass:

Gradient update (with respect to Q-function parameters θ):

$$\nabla_{\theta_i} L_i(\theta_i) = \mathbb{E}\left[r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) - Q(s, a; \theta_i)) \nabla_{\theta_i} Q(s, a; \theta_i)\right]$$

Remember: want to find a Q-function that satisfies the Bellman Equation:

Forward Pass:

Loss function: $L_i(\theta_i) = \mathbb{E} \left| (y_i - Q(s, a; \theta_i)^2) \right|$ $y_i = \mathbb{E}[r + \gamma \max_{a'} Q^*(s', a') \mid s, a]$ where

Backward Pass:

Gradient update (with respect to Q-function parameters θ):

$$\nabla_{\theta_i} L_i(\theta_i) = \mathbb{E}\left[r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) - Q(s, a; \theta_i)) \nabla_{\theta_i} Q(s, a; \theta_i)\right]$$



Training the Q-Network: **Experience Replay**

Learning from **batches of consecutive samples is problematic**:

 Samples are correlated => inefficient learning - Current Q-network parameters determines next training samples (e.g. if maximizing) => can lead to bad feedback loops

Address these problems using experience replay - Continually update a replay memory table of transitions (s_t , a_t , r_t , s_{t+1}) as game (experience) episodes are played of consecutive samples

- action is to move left, training samples will be dominated by samples from left-hand size)

- Train Q-network on random minibatches of transitions from the replay memory, instead

Experience Replay

Experience Replay

To remove correlations, build data-set from agent's own experience



Example: Atari Playing

The algorithm tries to hit the ball back, but it is yet too clumsy to manage.

Starting out - 10 minutes of training

Example: Atari Playing

The algorithm tries to hit the ball back, but it is yet too clumsy to manage.

Starting out - 10 minutes of training

Deep RL

Value-based RL

Use neural nets to represent Q function

Policy-based RL

– Use neural nets to represent the policy π_{θ}

 π

Model-based RL

- Use neural nets to represent and learn the model

$\begin{aligned} Q(s,a;\theta) \\ Q(s,a;\theta^*) &\approx Q^*(s,a) \end{aligned}$

$\pi_{ heta}$ $\pi_{ heta^*} pprox \pi^*$

* slide from Dhruv Batra

Deep RL

Value-based RL

- Use neural nets to represent Q function $Q(s, a; \theta)$

Policy-based RL

– Use neural nets to represent the policy π_{θ}

Model-based RL

- Use neural nets to represent and learn the model

$\begin{aligned} Q(s,a;\theta) \\ Q(s,a;\theta^*) &\approx Q^*(s,a) \end{aligned}$

$\pi_{\theta^*} \approx \pi^*$

* slide from Dhruv Batra

Formally, let's define a class of parameterized policies:

For each policy, define its value:

 $J(\theta) = \mathbb{E} \left| \sum_{t \ge 0} \gamma^t r_t | \pi_{\theta} \right|$



Formally, let's define a class of parameterized policies:

For each policy, define its value:

$$J(heta) = \mathbb{E}\left[\sum_{t\geq 0} \gamma^t
ight.$$

We want to find the optimal policy $\theta^* = \arg \max_{\alpha} J(\theta)$



Formally, let's define a class of parameterized policies:

For each policy, define its value:

$$J(heta) = \mathbb{E}\left[\sum_{t\geq 0} \gamma^t
ight]$$

We want to find the optimal policy $\theta^* = \arg \max_{\alpha} J(\theta)$

How can we do this?

 $r_t r_t |\pi_{\theta}|$

Formally, let's define a class of parameterized policies:

For each policy, define its value:

$$J(heta) = \mathbb{E}\left[\sum_{t\geq 0} \gamma^t
ight]$$

We want to find the optimal policy $\theta^* = rg \max_{\theta} J(\theta)$

How can we do this?

Gradient ascent on policy parameters!

 $|\pi_{ heta}| = |\pi_{ heta}|$

REINFORCE algorithm

Expected reward:

 $J(\theta) = \mathbb{E}_{\tau \sim p(\tau;\theta)} \left[r(\tau) \right]$ $= \int_{\tau} r(\tau) p(\tau; \theta) \mathrm{d}\tau$

Where $r(\tau)$ is the reward of a trajectory



$$\tau = (s_0, a_0, r_0, s_1, \ldots)$$

REINFORCE algorithm

Expected reward:

$$J(\theta) = \mathbb{E}_{\tau \sim p(\tau;\theta)} [r(\tau)]$$
$$= \int_{\tau} r(\tau) p(\tau;\theta) \phi$$

Where $r(\tau)$ is the reward of a trajectory Now let's differentiate this: $\nabla_{\theta} J(\theta) =$

Intractable! Expectation of gradient is problematic when p depends on θ

)]

 $\mathrm{d} au$

$$\tau = (s_0, a_0, r_0, s_1, \ldots)$$
$$\int_{\tau} r(\tau) \nabla_{\theta} p(\tau; \theta) \mathrm{d}\tau$$

REINFORCE algorithm

Expected reward:

$$J(\theta) = \mathbb{E}_{\tau \sim p(\tau;\theta)} [r(\tau)]$$
$$= \int_{\tau} r(\tau) p(\tau;\theta) d\tau$$

Where $r(\tau)$ is the reward of a trajectory Now let's differentiate this: $\nabla_{\theta} J(\theta) =$

However, we can use a nice trick: $\nabla_{\theta} p($

If we inject this back: $\nabla_{\theta} J(\theta) = \int_{\tau} (r(\tau) \nabla_{\theta} \log p(\tau; \theta)) p(\tau; \theta) d\tau$ $= \mathbb{E}_{\tau \sim p(\tau; \theta)} [r(\tau) \nabla_{\theta} \log p(\tau; \theta)]$ Can estimate with Monte Carlo sampling

$$\begin{aligned} \tau &= (s_0, a_0, r_0, s_1, \ldots) \\ \int_{\tau} r(\tau) \nabla_{\theta} p(\tau; \theta) d\tau \\ (\tau; \theta) &= p(\tau; \theta) \frac{\nabla_{\theta} p(\tau; \theta)}{p(\tau; \theta)} = p(\tau; \theta) \nabla_{\theta} \log p(\tau; \theta) \end{aligned}$$
Gradient estimator:

$\nabla_{\theta} J(\theta) \approx \sum r(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$ $t \ge 0$

Interpretation:

- If $r(\tau)$ is high, push up the probabilities of the actions seen - If $r(\tau)$ is low, push down the probabilities of the actions seen

Gradient estimator:

$\nabla_{\theta} J(\theta) \approx \sum r(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$ $t \ge 0$

Interpretation:

- If $r(\tau)$ is high, push up the probabilities of the actions seen - If $r(\tau)$ is low, push down the probabilities of the actions seen

Might seem simplistic to say that if a trajectory is good then all its actions were good. But in expectation, it averages out!



* slide from Dhruv Batra

Gradient estimator:

$\nabla_{\theta} J(\theta) \approx \sum r(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$ $t \ge 0$

Interpretation:

- If $r(\tau)$ is high, push up the probabilities of the actions seen - If $r(\tau)$ is low, push down the probabilities of the actions seen

Might seem simplistic to say that if a trajectory is good then all its actions were good. But in expectation, it averages out!

However, this also suffers from high variance because credit assignment is really hard. Can we help the estimator?

Objective: Image Classification

Take a sequence of "glimpses" selectively focusing on regions of the image, to predict class

- Inspiration from human perception and eye movements
- Saves computational resources => scalability
- Able to ignore clutter / irrelevant parts of image



glimpse

[Mnih et al., 2014]





Objective: Image Classification

Take a sequence of "glimpses" selectively focusing on regions of the image, to predict class

- Inspiration from human perception and eye movements
- Saves computational resources => scalability
- Able to ignore clutter / irrelevant parts of image

State: Glimpses seen so far **Action:** (x,y) coordinates (center of glimpse) of where to look next in image **Reward**: 1 at the final timestep if image correctly classified, 0 otherwise



glimpse

[Mnih et al., 2014]





Objective: Image Classification

Take a sequence of "glimpses" selectively focusing on regions of the image, to predict class

- Inspiration from human perception and eye movements
- Saves computational resources => scalability
- Able to ignore clutter / irrelevant parts of image

State: Glimpses seen so far **Action**: (x,y) coordinates (center of glimpse) of where to look next in image **Reward**: 1 at the final timestep if image correctly classified, 0 otherwise

Glimpsing is a **non-differentiable operation** => learn policy for how to take glimpse actions using REINFORCE Given state of glimpses seen so far, use RNN to model the state and output next action



glimpse

[Mnih et al., 2014]







[Mnih et al., 2014]







[Mnih et al., 2014]







[Mnih et al., 2014]







[Mnih et al., 2014]







Has also been used in many other tasks including fine-grained image recognition, image captioning, and visual question-answering!

[Mnih et al., 2014]







Has also been used in many other tasks including fine-grained image recognition, image captioning, and visual question-answering!

[Mnih et al., 2014]







Has also been used in many other tasks including fine-grained image recognition, image captioning, and visual question-answering!

[Mnih et al., 2014]







Has also been used in many other tasks including fine-grained image recognition, image captioning, and visual question-answering!

[Mnih et al., 2014]





Summary

lot of samples. *Challenge*: sample-efficiency

efficient. Challenge: exploration

Guarantees:

— Policy Gradients: Converges to a local minima of $J(\theta)$, often good enough! complicated function approximator

- **Policy gradients:** very general but suffer from high variance so requires a
- Q-learning: does not always work but when it works, usually more sample-

- Q-learning: Zero guarantees since you are approximating Bellman equation with a