



S&DS 365 / 665  
**Intermediate Machine Learning**

# **Graph Neural Nets** (continued) **Reinforcement Learning**

April 4

**Yale**

# Stuff

- Assignment 3 out; due week from today
- Quiz 3 on Wednesday
  - ▶ Variational inference and VAEs
  - ▶ Undirected graphs and glasso
  - ▶ Graph neural nets
- Some practice today :)

# Outline for today

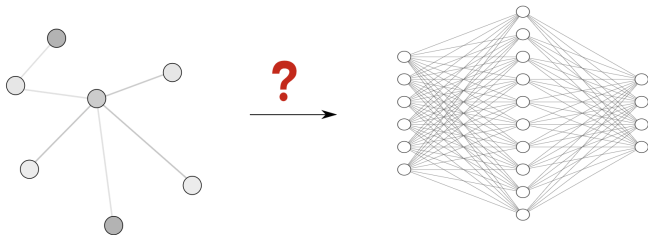
- Graph neural nets (continued)
- RL concepts
- Q-learning

# Graph neural networks

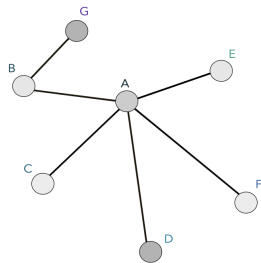
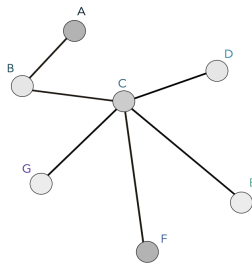
Let's pick up graph neural networks from last time

<https://distill.pub/2021/understanding-gnns/>

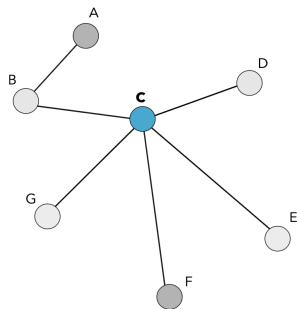
# Equivariance problem



# Equivariance problem



# Graph Laplacian



$$\begin{array}{c} \text{A} \\ \text{B} \\ \text{C} \\ \text{D} \\ \text{E} \\ \text{F} \\ \text{G} \end{array} \begin{bmatrix} 1 & -1 & & & & & \\ -1 & 2 & -1 & & & & \\ & -1 & 5 & -1 & -1 & -1 & -1 \\ & & -1 & 1 & & & \\ & & -1 & & 1 & & \\ & & -1 & & & 1 & \\ & & -1 & & & & 1 \end{bmatrix}$$

Laplacian  $L$  of  $G$

# Polynomials of the Laplacian

$$p_w(L) = w_0 I_n + w_1 L + w_2 L^2 + \cdots w_d L^d$$

If  $\text{dist}(u, v) > i$  then the  $(u, v)$  entry of  $L^i$  is zero

- This is analogous to a CNN filter (kernel)
- The weights  $w_i$  play role of filter coefficients
- Degree  $d$  of polynomial plays role of the size of the kernel



# The Laplacian is a Mercer kernel

- Symmetric  $L_{uv} = L_{vu}$
- Positive-definite:

$$f^T L f = \sum_{(u,v) \in E} (f_u - f_v)^2 \geq 0$$

# Whence equivariance

A transformation  $f : \mathbb{R}^n \longrightarrow \mathbb{R}^n$  is equivariant if

$$f(Px) = Pf(x)$$

for any permutation matrix  $P$ , where  $PP^T = I$ .

The transformed data and Laplacian are

$$x \longrightarrow Px$$

$$L \longrightarrow PLP^T$$

$$L^i \longrightarrow PL^iP^T$$

# Whence equivariance

A transformation  $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$  is equivariant if

$$f(Px) = Pf(x)$$

for any permutation matrix  $P$ , where  $PP^T = I$ .

The transformed polynomial kernels are

$$\begin{aligned} f(Px) &= \sum_{i=0}^d w_i (PL^i P^T) Px \\ &= \sum_{i=0}^d w_i PL^i x \\ &= P \sum_{i=0}^d w_i L^i x \\ &= Pf(x) \end{aligned}$$

# Building layers

Let  $h^{(k)}$  be the neurons at layer  $k$ .

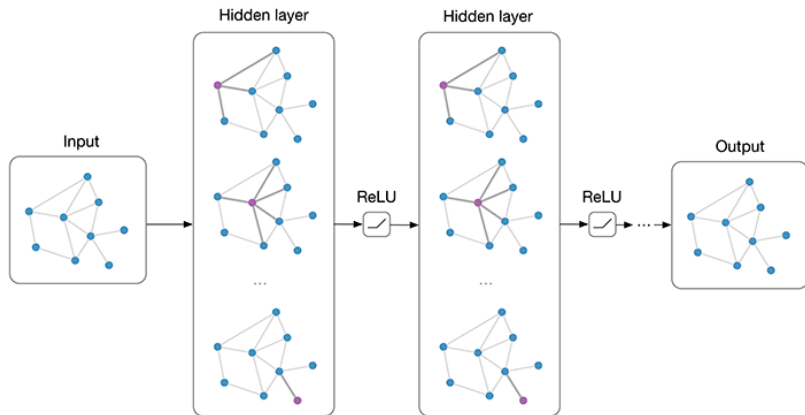
We start with  $h^{(0)} = x$ , a value  $x_j$  at each node  $j$

The next layer is

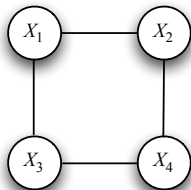
$$h^{(k+1)} = \varphi \left( p_w(L) h^{(k)} \right)$$

See tutorial for other ways of building layers

# Building layers

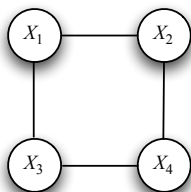


# Quiz practice



What is the Laplacian  $L$ ?

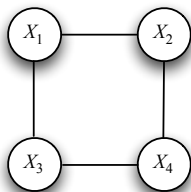
# Quiz practice



What is the Laplacian  $L$ ?

$$L = \begin{pmatrix} 2 & -1 & -1 & 0 \\ -1 & 2 & 0 & -1 \\ -1 & 0 & 2 & -1 \\ 0 & -1 & -1 & 2 \end{pmatrix}$$

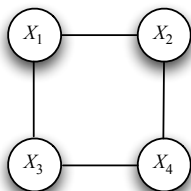
# Quiz practice



What is  $L^2$ ?



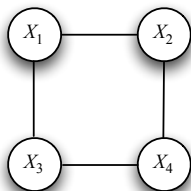
# Quiz practice



What is  $L^2$ ?

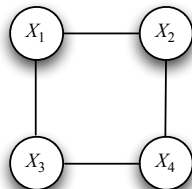
$$L^2 = \begin{pmatrix} 6 & -4 & -4 & 2 \\ -4 & 6 & 2 & -4 \\ -4 & 2 & 6 & -4 \\ 2 & -4 & -4 & 6 \end{pmatrix}$$

# Quiz practice



If  $x = (1, 2, 3, 4)^T$  what is  $h = \text{ReLU}(Lx)$ ?

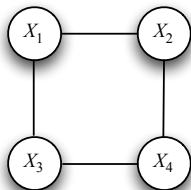
## Quiz practice



If  $x = (1, 2, 3, 4)^T$  what is  $h = \text{ReLU}(Lx)$ ?

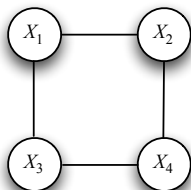
$$\text{ReLU}(Lx) = \text{ReLU}((-3, -1, 1, 3)^T) = (0, 0, 1, 3)^T$$

# Quiz practice



If  $x = (1, 2, 3, 4)^T$  what is  $x^T L x$ ?

## Quiz practice



If  $x = (1, 2, 3, 4)^T$  what is  $x^T L x$ ?

$$x^T L x = \sum_{(u,v) \in E} (x_u - x_v)^2 = 10$$

# Summary: Graph neural nets

- Certain data have natural graphical structure
- GNNs are analogues of CNNs for graphs
- Based on use of graph Laplacian
- Independent of ordering of nodes (equivariant)

**Next topic: Reinforcement learning**

# Reinforcement learning

- An agent interacts with an environment
- The actions the agent takes change the state of the environment
- The agent receives rewards for each action, and seeks to maximize the total cumulative reward

*Reinforcement learning is a framework for sequential decision making to achieve a long-term goal.*



# Reinforcement learning: Motivating examples

- Playing chess, backgammon, Atari games, etc.
- Learning to walk or ride a bike
- A robot vacuum cleaning up the house
- Drug discovery, personalized health, energy management

# Reinforcement learning: Formalization

- The environment is in state  $s$  at a given time
- The agent takes action  $a$
- The environment transitions to state  $s' = \text{next}(s, a)$
- The agent receives reward  $r = \text{reward}(s, a)$

# Reinforcement learning: Formalization

- The environment is in state  $s$  at a given time
- The agent takes action  $a$
- The environment transitions to state  $s' = \text{next}(s, a)$
- The agent receives reward  $r = \text{reward}(s, a)$

This is said to be a *Markov decision process*. It's “Markov” because the next state only depends on the current state and the action selected. It's a “decision process” because the agent is making choices of actions in a sequential manner.

# Characteristics

- RL is inherently sequential
- In between supervised and unsupervised learning
- Agent can't act too greedily; needs to be strategic

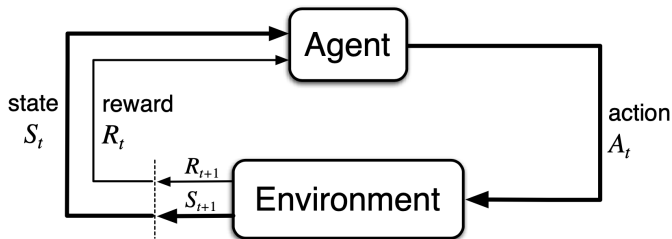
The aim of RL is to learn to make optimal decisions from experience

# Principles of RL

Some key RL concepts and principles:

Policy, reward signal, value function, model, Bellman equation

# Principles of RL



Rewards and state transitions are probabilistic, in general

# Principles of RL

*Policy*: A mapping from states to actions. An algorithm/rule to make decisions at each time step, designed to maximize the long term reward.

# Principles of RL

*Reward signal*: The sequence of rewards received at each time step. An abstraction of “pleasure” (positive reward) and “pain” (negative reward) in animal behavior.



# Principles of RL

*Value function*: A mapping from states to total reward. The total reward the agent can expect to accumulate in the future, starting from that state.

Rewards are short term. Values are predictions of future rewards.

# Principles of RL

*Model*: Used for planning to mimic the behavior of the environment, to predict rewards and next states.

A *model-free* approach directly estimates a value function, without modeling the environment.

Analogous to distinction between generative and discriminative classification models

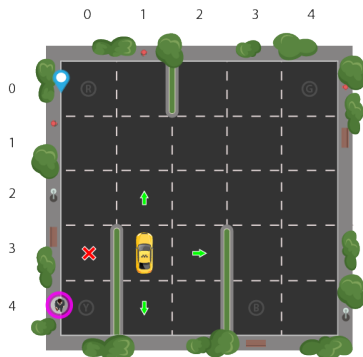
# Taxi problem

We'll introduce the important Q-learning algorithm with the toy “Taxi problem”

The code uses OpenAI gym.

# Taxi problem

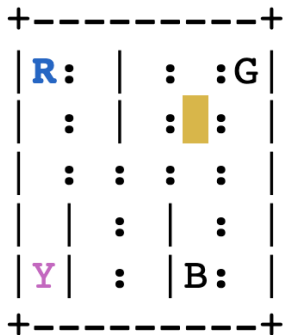
A taxicab drives around the environment, picking up and delivering a passenger at four locations



# Taxi problem

A taxicab drives around the environment, picking up and delivering a passenger at four locations

"Ascii art" rendition:



# Taxi problem: Description

- Four designated locations: R(ed), G(reen), Y(ellow), and B(lue)
- Taxi starts off at random square and passenger is at random location
- Taxi drives to passenger's location, picks up the passenger, drives to passenger's destination, drops off passenger
- Once the passenger is dropped off, the episode ends.

# Taxi problem: State space

- 25 taxi positions
- 5 possible locations of passenger: At waiting location or in taxi
- 4 possible destination locations
- Total number of states:  $25 \times 5 \times 4 = 500$

# Taxi problem: State space

Passenger location coded as integers:

- 0: R(ed)
- 1: G(reen)
- 2: Y(ellow)
- 3: B(lue)
- 4: in taxi



# Taxi problem: State space

Destinations coded as:

- 0: R(ed)
- 1: G(reen)
- 2: Y(ellow)
- 3: B(lue)

# Taxi problem: State space

Agent actions coded as:

- 0: move south
- 1: move north
- 2: move east
- 3: move west
- 4: pickup passenger
- 5: drop off passenger

# Taxi problem: State space

Rewards:

- Default reward per step: -1
- Reward for delivering passenger: +20
- Illegal “pickup” or “drop-off”: -10

# Taxi problem: State space

State space represented as a tuple:

state = (taxi row, taxi column, passenger location, destination)

# Q-learning

- Maintains a “quality” variable  $Q(s, a)$  for taking action  $a$  in state  $s$

# Q-learning

- Maintains a “quality” variable  $Q(s, a)$  for taking action  $a$  in state  $s$
- A measure of the cumulative rewards obtained by the algorithm when it takes action  $a$  in state  $s$

# Q-learning

- Maintains a “quality” variable  $Q(s, a)$  for taking action  $a$  in state  $s$
- A measure of the cumulative rewards obtained by the algorithm when it takes action  $a$  in state  $s$
- Quality should not be assessed purely based on the reward the action has in the current time step

# Q-learning

- Maintains a “quality” variable  $Q(s, a)$  for taking action  $a$  in state  $s$
- A measure of the cumulative rewards obtained by the algorithm when it takes action  $a$  in state  $s$
- Quality should not be assessed purely based on the reward the action has in the current time step
- Need to take into account the future rewards



# Q-learning update

## Update:

$$Q(s, a) \leftarrow Q(s, a) + \alpha \left( \text{reward}(s, a) + \gamma \max_{a'} Q(\text{next}(s, a), a') - Q(s, a) \right)$$

# Q-learning update

## Update:

$$Q(s, a) \leftarrow Q(s, a) + \alpha \left( \text{reward}(s, a) + \gamma \max_{a'} Q(\text{next}(s, a), a') - Q(s, a) \right)$$

- When action  $a$  is taken in state  $s$ , reward  $\text{reward}(s, a)$  is given
- Then, the algorithm moves to a new state  $\text{next}(s, a)$

# Q-learning update

## Update:

$$Q(s, a) \leftarrow Q(s, a) + \alpha \left( \text{reward}(s, a) + \gamma \max_{a'} Q(\text{next}(s, a), a') - Q(s, a) \right)$$

- For example, if the taxi is location (2, 2) and takes the “West” action ( $a = 3$ ), then there is a reward of -1, and the taxi moves to the new location (2, 1)
- If cab is empty, it remains empty, and if it contains the passenger, the passenger remains.

# Q-learning update

## Update:

$$Q(s, a) \leftarrow Q(s, a) + \alpha \left( \text{reward}(s, a) + \gamma \max_{a'} Q(\text{next}(s, a), a') - Q(s, a) \right)$$

- Cumulative future reward of this action is  $\max_{a'} Q(\text{next}(s, a), a')$
- Future rewards discounted by factor  $\gamma < 1$
- Trades off short-term against long-term rewards
- A gradient *ascent* algorithm, with step size  $\alpha$

Let's go to the notebook!

# Bellman equation



The optimal value function is the largest expected discounted long term reward starting from that state.

# Bellman equation: Deterministic case

The optimality condition for the value function  $v_*$  is

$$v_*(s) = \max_a \left\{ \text{reward}(s, a) + \gamma v_*(\text{next}(s, a)) \right\}$$

# Bellman equation: Deterministic case

The optimality condition for the Q-function is

$$Q_*(s, a) = \text{reward}(s, a) + \gamma \max_{a'} Q_*(\text{next}(s, a), a')$$

and then  $v_*(s) = \max_{a'} Q_*(s, a')$



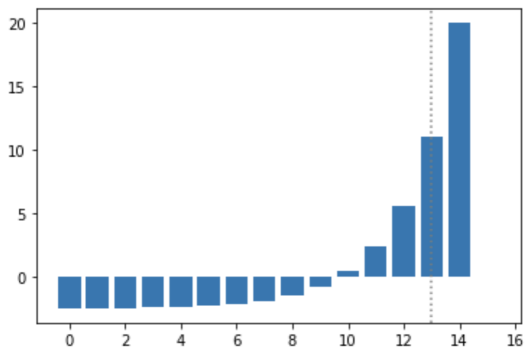
# Q-learning update

Note how this makes sense in terms of the update rule:

$$Q(s, a) \leftarrow Q(s, a) + \alpha \left( \text{reward}(s, a) + \gamma \max_{a'} Q(\text{next}(s, a), a') - Q(s, a) \right)$$

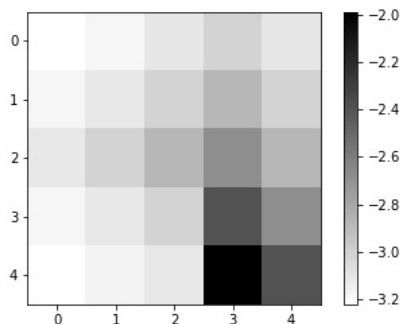
+-----+				
R:	:	:	G	
:	:	:		
:	:	:	:	
	:	:	:	
Y	:	B:	:	
+-----+				

(East)



# Question

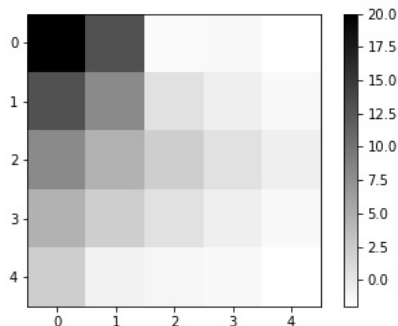
For a fixed passenger location and destination, the value function  $v(\text{row}, \text{col})$  assigns a value to each of the  $25 = 5 \times 5$  grid points.



Is the passenger waiting, or in the taxi?

# Question

For a fixed passenger location and destination, the value function  $v(\text{row}, \text{col})$  assigns a value to each of the  $25 = 5 \times 5$  grid points.



Is the passenger waiting, or in the taxi?

# Summary

- Reinforcement learning is a framework for sequential decision making to achieve a long-term goal
- The agent receives rewards for each action, and seeks to maximize the total cumulative reward
- The value of a state is the total reward the agent can expect to accumulate in the future, starting from that state
- Q-learning is an iterative algorithm that maximizes the “quality” of each state-action pair
- The Bellman equations are optimality conditions that characterize the value function