ECS171: Machine Learning

L16: Unsupervised Learning III and reinforcement learning

Instructor: Prof. Maike Sonnewald TAs: Pu Sun & Devashree Kataria



Intended Learning Outcomes

- Describe the use of graphs in unsupervised learning
- Describe example of a k-NN graph
 - Computational cost implications
 - Concept of k
- Describe and appreciate the application of geometry to learning UMAP example
 - Role of weights in graph
 - Parameters n_neighbours and min_dist roles in simple examples
- Describe and compare/contrast reinforcement learning vrt supervised and unsupervised learning
 - Describe concepts of state, action, revard, policy and value in simple example like tic, tac, toe
 - Markov Decision Problem
 - Be familiar with Q-learning

Graphs

- Very powerful method to find patterns in data
- Can be based on a similarity or dissimilarity measure
- Non-parametric
- A vertex for each point, and an edge from *p* to *q*
 - Whenever q is a nearest neighbor of p: a point whose distance from p is minimum among all the given points other than p itself
- Can be used to understand the topology of a dataset

	nodes (or vertices)
~ /	
edges	
(or links)	\checkmark
_	

k-Nearest Neighbor graph

- A graph in which two vertices *p* and *q* are connected by an edge
 - Condition: distance between *p* and *q* is among the *k*-th smallest distances from *p* to other objects from *P*.
- Put differently: it is a graph where each node is connected to the *k* most similar entities, as an attempt to insert some bias about what a good graph structure should look like

Setting k allows looking at global vs local structure





Pairwise distance matrix between points



1		.62	.15	.15	.82	.25	.65	.33	.45	.61	.73	.84	.99	.41	.20	.24	.26	.34	.67	.15	.01	.35	.73	.51
2-	.77		.48	.50	.41	.52	.39	.29		.02	.12	.25	.50	.28	.43	.41	.43	.35	.25	.49	.62	.37	.25	.18
3-	.36	.45		.02	.69	.16	.52		.32	.48	.58	.71	.88	.29	.05	.09	.12	.19	.52	.03	.15		.58	.36
4-	.36	.48	.07		.70	.14	.52		.33	.49	.60	.72	.90	.30	.07	.10	.13	.20	.53	.01	.15		.58	.37
5-	.83	.21	.53	.58		.65	.22	.56	.47	.43	.35	.34	.67	.44	.65	.64	.66	.57	.21	.69	.83	.48	.24	.41
6-	.43	.41	.12	.15	.47		.45	.28	.35	.52	.61	.73	.96	.31	.17	.17	.19	.23	.49	.13	.25	.19	.55	.37
7-	.68	.20	.34	.35	.34	.30		.42	.35	.40	.38	.44	.78	.31	.48	.47	.48	.40	.18	.51	.65	.30	.25	.30
8-	.47	.34	.21	.26	.39	.16	.31		.13	.29	.40	.52	.72	.15	.15	.16	.20	.16	.38	.21	.33	.17	.43	.22
9-	.54	.26	.24	.28	.33		.23	.10		.19	.30	.40	.66	.14	.27	.25	.28	.20	.29	.33	.45	.21	.34	.16
10	.74	.05	.42	.45	.23	.38	.18	.31	.22		.13	.26	.50	.28	.43	.41	.43	.35	.26	.49	.62	.37	.27	.19
11-	.90	.15	.58	.61	.28	.55	.28	.49	.40	.18		17	.50	.37	.53	.51	.53	.45	.22	.59	.73	.45	.19	.24
12-	.20	.59	.25	.27	.65	.31	.54	.30	.36	.56	.73		.48	.46	.66	.65	.67	.58	.30	.72	.85	.55	.28	.39
13	.32	.45		.21	.52	.21	.40	.17	.22	.42	.60	.14		.69	.84	.83	.84	.77	.65	.90	1.0	.81	.61	.66
14	38	41	15	19	47	15	34	11	16	37	54	21	09		25	25	26	19	30	29	42	15	37	21
15	41	39	11		45	10	32	11	16	36	54	26	14	11	12.5	07	11	16	48	07	20	19	53	31
16	48	35	13	14	47	15	22	22	19	32	47	35	23	19	16		06	11	46	11	24		51	29
17.	40	40	06	10	19	08	30	15	19	37	54	27	16	12	08	10		10	48	14	25		52	30
10	46	33	12	16	43	13	24	14	12	30	47	30	17	12	10	08	08		30		34	18	43	22
10	03	20	50	61	30	55	27	57	43	22	10	76	63	57	55	47	55	18		57	67	33	10	21
19	35	48		.01	57	13	36	24	.45	45	61	25	1.00	17	15	14	.08	.40	61	.32	.07	21	58	37
20-		.40	25	22	0 /	42		.24	.Z /	.43	.01	.23	24	20	42	46	20	45	.01	22	.10	26		51
21-	.00	26	.33		04	.45	.07	.40	.54	.74	.90	-25	.54	.50	.42	.40	.59	.45	.92		5.0	.50	.15	.51
22-	.52	.26	.19	.23	.37	.18	.20	.14	.08	.23	.40	.35	.21	.10	.15	.12	.15	.08	.42	.23	.52	10	.40	.24
23-	1.0	.25	.66	.68	.33	.62	.35	.58	.50	.28	.12	.83	.69	.64	.62	.55	.62	.55	.08	.69	.99	.49		.23
24-	.87	.18	.53	.54	.30	.49	.20	.47	.38	.19	.13	.71	.58	.52	.50	.41	.49	.42	.07	.55	.86	.37	.15	
	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24

Nobel et al., 2021

Topology and data in ML



Learning more/dimensionality reduction with graphs

Going from high to low dimensions can be done with graphs

Avoids making assumptions of e.g. orthogonality as in PCA

Use distances to determine and preserve the 'latent space' in fewer dimensions





Graph layouts

We have a simple graph if we only care about the 0- and 1-simplices

FInding the topological representation is just a graph

Finding a low dimensional representation can be described as a graph layout problem.



More details on what follows: https://umap-learn.readthedocs.io/en/latest/how_umap_works.html#id6







Uniform Manifold Approximation Projection (UMAP)

FInding a low-dimensional representation using weighted graphs as an optimisation problem

Desirable properties:

- Lower dimensions for visualisation
- Weights assure that 'related' data-points stay together

Assume data is uniformly distributed

- Compute an approximation of local distance to each point
- Riemannian *geometry*



Open balls over uniformly distributed data



Open balls of radius one with a locally varying metric



Fuzzy open balls of radius one with a locally varying metric



Local connectivity and fuzzy open sets



Edges with incompatible weights



Graph with combined edge weights

Having determined the weights we can optimise

- Colour on previous slides proportional to weight
- Move from a high (h) to a low (l) dimensional representation by minimizing the cross entropy
- In algorithm, set 'min-dist' and 'n-neighbours'





Figure 5: UMAP projections of a 3D woolly mammoth skeleton (50k points, 10k shown) into 2 dimensions, with various settings for the n_neighbors and min_dist parameters.



Figure 5: UMAP projections of a 3D woolly mammoth skeleton (50k points, 10k shown) into 2 dimensions, with various settings for the n_neighbors and min_dist parameters.



Figure 5: UMAP projections of a 3D woolly mammoth skeleton (50k points, 10k shown) into 2 dimensions, with various settings for the n_neighbors and min_dist parameters.

Machine learning categories



Reinforcement learning

Learning algorithms differ in the information available to model we are developing:

- Supervised: correct outputs
- Unsupervised: no feedback, must construct measure of good output
- Reinforcement learning

More realistic learning scenario:

- Continuous stream of input information, and actions
- Effects of action depend on state of the world
- Obtain reward that depends on world state and actions
- Not "correct" answers... just some feedback

Reinforcement learning



Reinforcement learning: Tic Tac Toe example, Notation



environment

Reinforcement learning: Tic Tac Toe example, Notation



(current) state





reward (here: -1)

Formulating Reinforcement Learning

- World described by a discrete, finite set of states and actions
- At every time step t, we are in a state s_t, and we:
 - Take an action at (possibly null action)
 - Receive some reward **r**_{t+1}
 - Move into a new state **s**_{t+1}
- An RL agent may include one or more of these components:
 - **Policy** π : agent's behaviour function
 - Value function: how good is each state and/or action
 - **Model:** agent's representation of the environment

Reinforcement Learning: Policy

- A policy is the agent's behaviour.
- It's a selection of which action to take, based on the current state
- Deterministic policy: $a = \pi(s)$
- Stochastic policy: $\pi(a|s) = P[a_t = a|s_t = s]$

Reinforcement Learning: Value Function

- Value function is a prediction of future reward
- Used to evaluate the goodness/badness of states
- Our aim will be to maximize the value function (the total reward we receive over time):
 - find the policy with the highest expected reward
- By following a policy π , the value function is defined as:

 $V^{\pi}(s_t) = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \cdots$

- γ is called a discount rate, and it is always $0 \le \gamma \le 1$
- If γ is close to 1, rewards further in the future count more, and we say that the agent is "farsighted"
- If γ is less than 1 because there is usually a time limit to the sequence of actions needed to solve a task (we prefer rewards sooner rather than later)

Reinforcement Learning: Model

The model describes the environment by a distribution over rewards and state transitions:

$$P(s_{t+1} = s', r_{t+1} = r' | s_t = s, a_t = a)$$

We assume the Markov property: the future depends on the past only through the current state





• Rewards:



- Rewards: -1 per time-step
- Actions:



- Rewards: -1 per time-step
- Actions: N, E, S, W
- States:



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

[Slide credit: D. Silver]



 Arrows represent policy π(s) for each state s

[Slide credit: D. Silver]



 Numbers represent value V^π(s) of each state s

[Slide credit: D. Silver]

Example: Tic-Tac-Toe

- Consider the game tic-tac-toe:
 - Reward: win/lose/tie the game (+1/ 1/0) [only at final move in given game]
 - State: positions of X's and O's on the board
 - Policy: mapping from states to actions
 - Based on rules of game: choice of one open position
 - Value function: prediction of reward in future, based on current state
- In tic-tac-toe, since state space is tractable, can use a table to represent value function

Example: Tic-Tac-Toe

- Each board position (taking into account symmetry) has some probability
- Simple learning process:
 - start with all values = 0.5
 - policy: choose move with highest probability of winning given current legal moves from current state
 - update entries in table based on outcome of each game
 - After many games value function will represent true probability of winning from each state
- Can try alternative policy: sometimes select moves randomly (exploration!)

State	Probability of a win (Computer plays "o")
× 0 × 0	0.5
	0.5
× 0 × 0	1.0
×0 ×0	0.0
0 *	0.5
etc	

Basic Problems

- Markov Decision Problem (MDP): tuple (S, A, P, γ) where P is

$$P(s_{t+1} = s', r_{t+1} = r' | s_t = s, a_t = a)$$

- Standard MDP problems:
 - 1. Planning: given complete Markov decision problem as input, compute policy with optimal expected return



Basic Problems

- Markov Decision Problem (MDP): tuple (S, A, P, γ) where P is

$$P(s_{t+1} = s', r_{t+1} = r' | s_t = s, a_t = a)$$

- Standard MDP problems:
 - 1. Planning: given complete Markov decision problem as input, compute policy with optimal expected return
 - 2. Learning: We don't know which states are good or what the actions do.
 - We must try out the actions and states to learn what to do...

Example of Standard MDP Problem

- 1. Planning: given complete Markov decision problem as input, compute policy with optimal expected return
- 2. Learning: We don't know which states are good or what the actions do.
 - We must try out the actions and states to learn what to do...



r(s, a) (immediate reward)

Example of Standard MDP Problem

- 1. Planning: given complete Markov decision problem as input, compute policy with optimal expected return
- 2. Learning: We don't know which states are good or what the actions do.
 - We must try out the actions and states to learn what to do...



We will focus on learning, but discuss planning along the way...

Exploration vs. Exploitation

- If we knew how the world works (embodied in P), then the policy should be deterministic
 - Select optimal action in each state
- Reinforcement learning is like trial-and-error learning
- The agent should discover a good policy from its experiences of the environment
- Without losing too much reward along the way
- Since we do not have complete knowledge of the world, taking what appears to be the optimal action may prevent us from finding better states/actions
- Interesting trade-off:
 - Immediate reward (exploitation) vs. gaining knowledge that might enable higher future reward (exploration)

Exploitation/Exploration examples

- Restaurant Selection
 - Exploitation: Go to your favourite restaurant
 - Exploration: Try a new restaurant
- Online Banner Advertisements
 - Exploitation: Show the most successful advert
 - Exploration: Show a different advert
- Oil Drilling
 - Exploitation: Drill at the best known location
 - Exploration: Drill at a new location
- Game Playing
 - Exploitation: Play the move you believe is best
 - Exploration: Play an experimental move

Markov Decision Problem (MDP) formulation

- Goal: find policy π that maximizes expected accumulated future rewards V^{π}(s_t), obtained by following π from state s_t:

$$V^{\pi}(s_t) = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \cdots$$
$$= \sum_{i=0}^{\infty} \gamma^i r_{t+i}$$

- Game show example:
 - Assume series of questions, increasingly difficult, but increasing payoff
 - Choice: accept accumulated earnings and quit; or continue and risk losing everything
- Notice that:

$$V^{\pi}(s_t) = r_t + \gamma V^{\pi}(s_{t+1})$$

What to Learn?

- We might try to learn the function V (which we write as V*)

$$V^*(s) = \max_{a} \left[r(s, a) + \gamma V^*(\delta(s, a)) \right]$$

- Here $\delta(s, a)$ gives the next state, if we perform action a in current state s
- We could then do a lookahead search to choose best action from any state s:

$$\pi^*(s) = \arg\max_{a} \left[r(s, a) + \gamma V^*(\delta(s, a)) \right]$$

- But there's a problem:
 - This works well if we know $\delta()$ and r ()
 - But when we don't, we cannot choose actions this way

Q Learning

- Define a new function very similar to V*

$$Q(s,a) = r(s,a) + \gamma V^*(\delta(s,a))$$

- If we learn Q, we can choose the optimal action even without knowing δ

$$\pi^*(s) = \arg \max_{a} [r(s, a) + \gamma V^*(\delta(s, a))]$$

= $\arg \max_{a} Q(s, a)$

- Q is then the evaluation function we will learn

Q Learning



Training Rule to Learn Q

- Q and V* are closely related:

$$V^*(s) = \max_a Q(s,a)$$

- So we can write Q recursively:

$$Q(s_t, a_t) = r(s_t, a_t) + \gamma V^*(\delta(s_t, a_t))$$

= $r(s_t, a_t) + \gamma \max_{a'} Q(s_{t+1}, a')$

- Let \hat{Q} denote the learner's current approximation to Q
- Consider training rule, where s' is state resulting from applying action a in state s

$$\hat{Q}(s,a) \leftarrow r(s,a) + \gamma \max_{a'} \hat{Q}(s',a')$$

Q Learning for Deterministic World

- For each s, a initialize table entry $\hat{Q}(s, a) \leftarrow 0$
- Start in some initial state s
- Do forever:
 - Select an action a and execute it
 - Receive immediate reward r
 - Observe the new state s'
 - Update the table entry for $\hat{Q}(s, a)$ using Q learning rule:

$$\hat{Q}(s,a) \leftarrow r(s,a) + \gamma \max_{a'} \hat{Q}(s',a')$$

- $s \leftarrow s'$
- If we get to absorbing state, restart to initial state, and run thru "Do forever" loop until reach absorbing state

Updating Estimated Q

- Assume the robot is in state s1; some of its current estimates of Q are as shown; executes rightward move



- Important observation: at each time step (making an action a in state s only one entry of will \hat{Q} change (the entry $\hat{Q}(s, a)$)
- Notice that if rewards are non-negative, then \hat{Q} values only increase from 0, approach true Q

Q Learning: Summary

- Training set consists of series of intervals (episodes): sequence of (state, action, reward) triples, end at absorbing state
- Each executed action a results in transition from state s_i to s_j ; algorithm updates $\hat{Q}(s_i, a)$ using the learning rule
- Intuition for simple grid world, reward only upon entering goal state \rightarrow Q estimates improve from goal state back
 - 1. All $\hat{Q}(s, a)$ start at 0
 - 2. First episode only update $\hat{Q}(s, a)$ for transition leading to goal state
 - 3. Next episode if go thru this next-to-last transition, will update $\hat{Q}(s, a)$ another step back
 - 4. Eventually propagate information from transitions with non-zero reward throughout state-action space

Q Learning: Exploration/Exploitation

- Have not specified how actions chosen (during learning)
- Can choose actions to maximize $\hat{Q}(s, a)$
- Good idea?
- Can instead employ stochastic action selection (policy):

$$P(a_i|s) = rac{\exp(k\hat{Q}(s,a_i))}{\sum_j \exp(k\hat{Q}(s,a_j))}$$

- Can vary k during learning
 - More exploration early on, shift towards exploitation