

C3 Unsupervised Learning, Recommender System and Reinforcement Learning

▼ Welcome!

- ▼ Unsupervised learning
 - Clustering
 - Anomaly detection
- Recommender systems
- Reinforcement learning

▼ W1 Unsupervised Learning

- ▼ Clustering
 - ▼ What is clustering
 - Training set $\{x(1), \dots, x(m)\}$ without the target label y
 - ▼ Applications of clustering
 - Grouping similar news
 - DNA analysis
 - Market segmentation
 - Astronomical data analysis
 - ▼ K-means algorithm
 - 1. Randomly initialize K cluster centroids $\mu_1, \mu_2, \dots, \mu_K$
 - 2.
Repeat {

#assign points to cluster centroids
for $i = 1$ to m :
 $c(i) := \text{index (from 1 to } K) \text{ of cluster centroid closest to } x(i)$

#move cluster centroids
for $k = 1$ to K :
 $\mu_k := \text{average(mean) of points assigned to cluster } k$
}
}

- ▼ Optimization objective
 - $c(i)$ = index of cluster(1,2..K) to which example $x(i)$ is currently assigned
 - μ_k = cluster centroid k
 - $\mu_{c(i)}$ = cluster centroid of cluster to which example $x(i)$ has been assigned
 - Cost function/distortion function(L2 norm):

$$J(c(1), \dots, c(m), \mu_1, \dots, \mu_k) = \sum \|x(i) - \mu_{c(i)}\|^2 / m$$
- ▼ Initializing K-means
 - ▼ Random initialization
 - Choose $K < m$
 - Randomly pick K training examples
 - Set μ_1, \dots, μ_k equal to these K examples
 - #better than running K-means once
 For $i = 1$ to 100 {
 randomly initialize K-means
 run K-means, get $c(1), \dots, c(m), \mu_1, \dots, \mu_k$
 compute $J(c(1), \dots, c(m), \mu_1, \dots, \mu_k)$
 pick set of clusters that gave lowest cost J
- ▼ Choosing the number of clusters
 - Elbow method: calculate the cost function with K clusters
 - Evaluate k-means based on how well it performs on that later purpose
- ▼ Anomaly detection
 - ▼ Finding unusual events
 - ▼ Density estimation
 - $p(x)$: probability of x being seen in dataset
 - $p(x) < \epsilon$, anomaly point
 - $p(x) \geq \epsilon$, normal
 - ▼ Anomaly detection example
 - (Fraud detection)
 $x(i)$ = features of user i 's activities
 model $p(x)$ from data
 identify unusual users by checking which have $p(x) < \epsilon$
 - Gaussian distribution
- ▼ Algorithm

- 1. Choose n features $x(i)$ that you think might be indicative of anomalous examples
- 2. Fit parameters $\mu_1, \mu_2, \dots, \mu_n, \delta_1^2, \dots, \delta_n^2$
- 3. Given new example x , compute $p(x)$:

$$p(x) = p(x_1, \mu_1, \delta_1^2) * p(x_2, \mu_2, \delta_2^2) * \dots * (x_n, \mu_n, \delta_n^2) = N(\mu_1, \delta_1^2) * N(\mu_2, \delta_2^2) * \dots * N(\mu_n, \delta_n^2)$$
 anomaly if $p(x) < \epsilon$
- If $p(x_i)$ is very small then the overall $p(x)$ will be small
- ▼ Developing and evaluating an anomaly detection system
 - Fit model $p(x)$ on training set
 - On a cross validation example x , predict:
 $y = 1$ if $p(x) < \epsilon$ else 0
 - Possible evaluation metrics:
 true positive, false positive, true negative
 precision/recall
 F1-score
 - Can also use cross validation set to choose parameter ϵ
- ▼ Anomaly detection vs supervised learning
 - The number of positive examples ≤ 20 then use anomaly detection
 - If the future anomalies look nothing like the examples before then use anomaly detection
- ▼ Choosing what features to use
 - Change the features to more-Gaussian features
 - Error analysis for anomaly detection
- ▼ Monitoring computers in a data center
 - choose features that might take on unusually large or small values in the event of an anomaly
 - Deciding feature choice based on $p(x)$:
 large for normal examples
 becomes small for anomaly in the cross validation set

▼ W2 Recommender System

- ▼ Collaborative filtering
 - ▼ Using per-item features

- nu: number of users
nm: number of movies
n: number of features
- For user j: predict rating for movie i as:
 $w(j) \cdot x(i) + b(j)$
- $r(i,j) = 1$ if user j has rated movie i
 $y(i,j)$: rating given by user j on movie i
 $x(i)$: feature vector for movie i
 $w(j), b(j)$: parameters for user j
 $m(j)$ = number of movies rated by user j
- Cost function to learn $w(j), b(j)$ for user j:
 $\min J(w(j), b(j))$
 $= \sum (w(j) \cdot x(i) + b(j) - y(i,j))^2 / (2m(j)) + \lambda \sum (w_k(j))^2 / (2m(j)), k=1,2..n$
- Cost function to learn $w(1), \dots, w(nu), b(1), \dots, b(nu)$ for all users:
 $\min J(\text{all parameters } w, b)$
 $= \sum \sum (w(j) \cdot x(i) + b(j) - y(i,j))^2 / 2 + \lambda \sum \sum (w_k(j))^2 / 2, k=1,2..n$

▼ Collaborative filtering algorithm

- Given $w(1), \dots, w(nu), b(1), \dots, b(nu)$ to learn $x(i)$: #single feature
 $J(x(i))$
 $= \sum (w(j) \cdot x(i) + b(j) - y(i,j))^2 / 2 + \lambda \sum (x(k))^2 / 2$
- To learn $x(1), \dots, x(nm)$: # all features
 $J(x(1), \dots, x(nm))$
 $= \sum \sum (w(j) \cdot x(i) + b(j) - y(i,j))^2 / 2 + \lambda \sum \sum (x(k))^2 / 2$
- To learn w, b and x :
 $J(w, b, x) = \sum (w(j) \cdot x(i) + b(j) - y(i,j))^2 / 2 + \lambda \sum \sum (x(k))^2 / 2 + \lambda \sum \sum (w_k(j))^2 / 2$
- Gradient descent:
repeat

$$w_i^{(j)} = w_i^{(j)} - \alpha \frac{\partial}{\partial w_i^{(j)}} J(w, b, x)$$

$$b^{(j)} = b^{(j)} - \alpha \frac{\partial}{\partial b^{(j)}} J(w, b, x)$$

$$x_k^{(i)} = x_k^{(i)} - \alpha \frac{\partial}{\partial x_k^{(i)}} J(w, b, x)$$

▼ Binary labels: favs, likes and clicks

- For binary labels: predict that the probability of $y(i,j) = 1$ is given $g(w(j) \cdot x(i) + b(j))$ where g is a sigmoid function

- Cost function for single example:

$$L(g(w(j) \cdot x(i) + b(j)), y(i, j)) = -y(i, j) \log(g(w(j) \cdot x(i) + b(j))) - (1 - y(i, j)) \log(1 - g(w(j) \cdot x(i) + b(j)))$$
- Cost function for all example:

$$J(w, b, x) = \sum L(g(w(j) \cdot x(i) + b(j)), y(i, j))$$

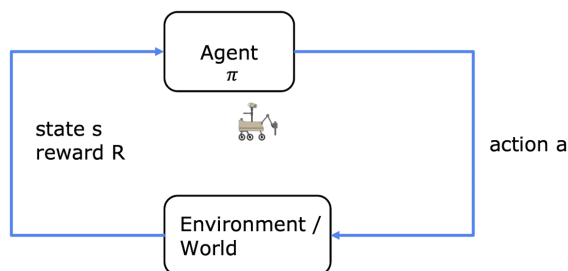
▼ W3 Reinforcement Learning

▼ key concepts

- policy π : what action ($a = \pi(s)$) to take in every state so as to maximize the return
- states \ actions \ rewards \ discount factor γ \ return \ policy

▼ markov decision process (MDP)

- the future only depends on the current state



▼ state-action value function

- $Q(s, a)$ = return if you
 .start in state s
 .take action a (once)
 .then behave optimally after that

- Bellman Equation:

$$Q(s, a) = R(s) + \gamma \max_{a'} Q(s', a')$$

$$Q(s, a) = R(s) + \gamma \max_{a'} Q(s', a')$$

100	100	50	12.5	25	6.25	12.5	10	6.25	20	40	40
100	0	0	0	0	0	0	40				
1	2	3	4	5	6						

$Q(4, \leftarrow)$
 $= 0 + (0.5)0 + (0.5)0 + (0.5)100$
 $= R(4) + (0.5)(0 + (0.5)0 + (0.5)100)$
 $= R(4) + (0.5) \max_{a'} Q(3, a')$

▼ random(stochastic) environment

- expected return

$$= \text{average}(R_1 + \gamma R_2 + \gamma^2 R_3 + \gamma^3 R_4 + \dots)$$

$$= E(R_1 + \gamma R_2 + \gamma^2 R_3 + \gamma^3 R_4 + \dots)$$

- maximum the average value
- Bellman Equation:

$$Q(s, a) = R(s) + \gamma E[\max_{a'} Q(s', a')]$$

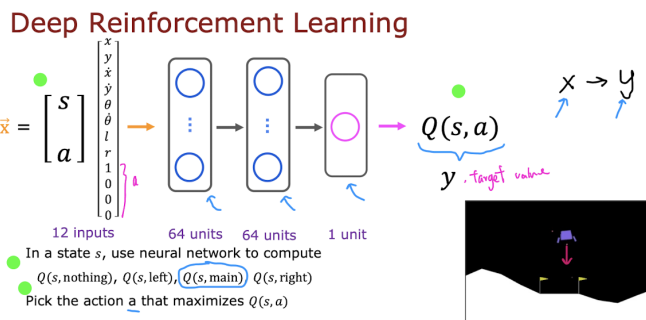
▼ continuous state spaces

▼ Lunar Lander Problem

- Learn a policy π
- given $s = [x, y, x', y', \theta, \theta', l, r]$
- picks action $a = \pi(s)$ so as to maximum the return
- $\gamma = 0.985$

▼ deep reinforcement learning

▪



- Bellman Equation:

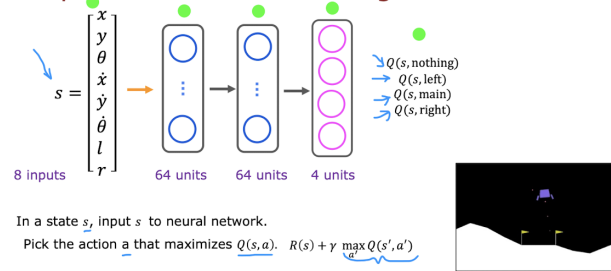
$$Q(s, a) = R(s) + \gamma \max_{a'} Q(s', a')$$

$$x = [s, a]$$

$$y = R(s) + \gamma \max_{a'} Q(s', a')$$
- learning algorithm
 Initialize neural network randomly as guess of $Q(s, a)$
 Repeat{
 Take actions in the lunar lander. Get($s, a, R(s), s'$)
 Store 10000 most recent ($s, a, R(s), s'$) tuples *replay buffer
 Train neural network:
 create training set of 10000 examples using
 $x = (s, a)$ and $y = R(s) + \gamma \max_{a'} Q(s', a')$
 Train Q_{new} such that $Q_{\text{new}} \sim y$
 set $Q = Q_{\text{new}}$
 }

▼ improved neural network architecture

Deep Reinforcement Learning



algorithm refinement: ϵ -greedy policy

- Take actions in the lunar lander

>>>

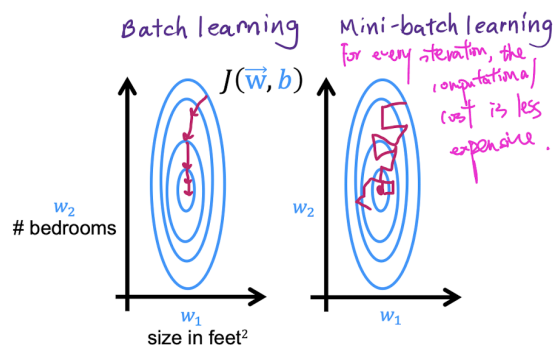
- with probability $1-\epsilon$, pick the action a that maximizes $Q(s, a)$
- with probability ϵ , pick an action a randomly

- 1 for greedy exploitation
2 for exploration
- greedy $1-\epsilon$ of the time and exploring ϵ of the time
- start ϵ very high and gradually decrease

mini-batch and soft update

- if the sample size is too much and when repeating gradient descent, it only minus a little
- so for each iteration, only choose the subset of training set
create training set of 10000 examples
>>> create training set of 1000 examples

Mini-batch



- (soft update) set $Q = Q_{\text{new}}$
>>>
if $Q(W, B)$ then $Q(W_{\text{new}}, B_{\text{new}})$
 $w = 0.01W_{\text{new}} + 0.99W$
 $B = 0.01B_{\text{new}} + 0.99B$
*0.01 and 0.99 can be modified

▼ further reading

- Mnih, V., Kavukcuoglu, K., Silver, D. et al. Human-level control through deep reinforcement learning. Nature 518, 529–533 (2015).
Lillicrap, T. P., Hunt, J. J., Pritzel, A., et al. Continuous Control with Deep Reinforcement Learning. ICLR (2016).
Mnih, V., Kavukcuoglu, K., Silver, D. et al. Playing Atari with Deep Reinforcement Learning. arXiv e-prints. arXiv:1312.5602 (2013).