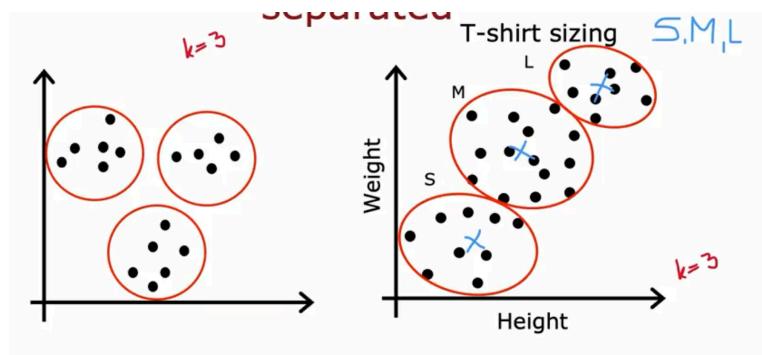
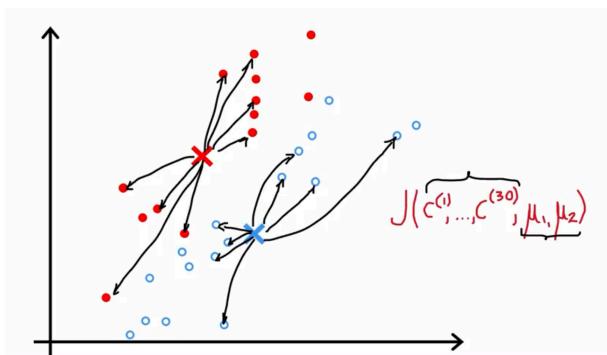


Course 3: Unsupervised Learning, Recommenders, Reinforcement Learning

Week 1:

Clustering

- K-Means Clustering:
 1. Initialize K amount of cluster centroid on random data points
 2. Find class membership (square distance) of n amount of object in relation to the nearest cluster centroid
 3. Relocate centroids to the average distance (mean) from all members of said centroid and reevaluate the objects until the membership (cost function) score no longer changes



- The cost function in k-means clustering is the aggregate of the distances computed between data points and their assigned centroids. It helps determine whether the overall clustering has improved in each iteration. A decreasing cost function indicates better-defined clusters, guiding convergence. Without it, there would be no clear way to measure or compare clustering quality across iterations.

Cost function for K-means

$$J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K) = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \mu_{c^{(i)}}\|^2$$

Repeat {

Assign points to cluster centroids

for $i = 1$ to m

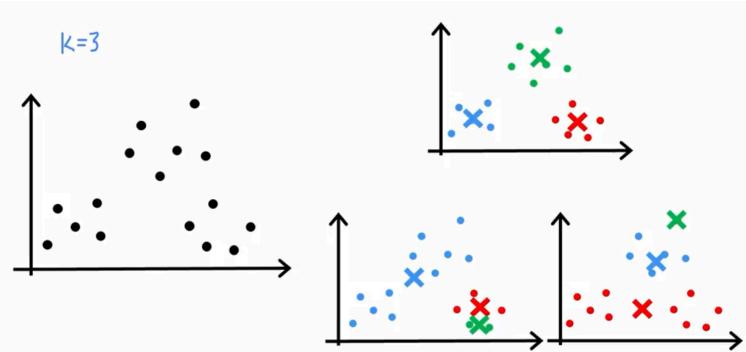
$c^{(i)} :=$ index of cluster centroid closest to $x^{(i)}$

Move cluster centroids

for $k = 1$ to K

$\mu_k :=$ average of points in cluster k

- Initializing K-means centroids randomly on data points helps overcome local minima and vastly different clusters.



Random initialization

```

For i = 1 to 100 {  

    Randomly initialize K-means.  

    Run K-means. Get  $c^{(1)}, \dots, c^{(m)}, \mu_1, \mu_2, \dots, \mu_k$   

    Computer cost function (distortion)  

     $J(c^{(1)}, \dots, c^{(m)}, \mu_1, \mu_2, \dots, \mu_k)$   

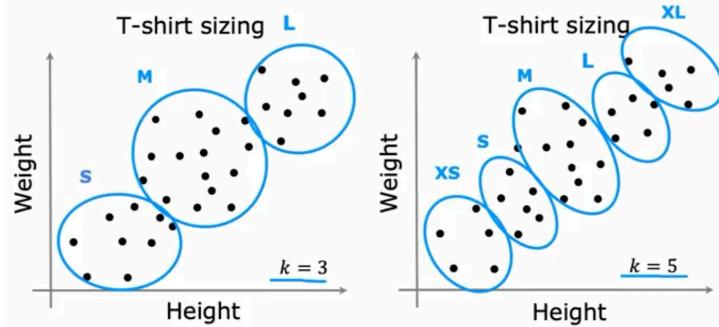
}
  
```

Pick set of clusters that gave lowest cost(J)

- Choosing amount of clusters:
 - Use square error to **1 to n** amount of clusters to determine the ideal amount of clusters (elbow method).
 - Evaluate based on overarching purpose of clustering (**Ideal**)

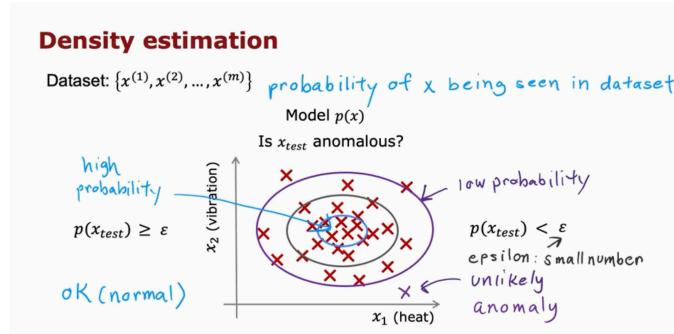
Choosing the value of K

Often, you want to get clusters for some later (downstream) purpose.
Evaluate K-means based on how well it performs on that later purpose.

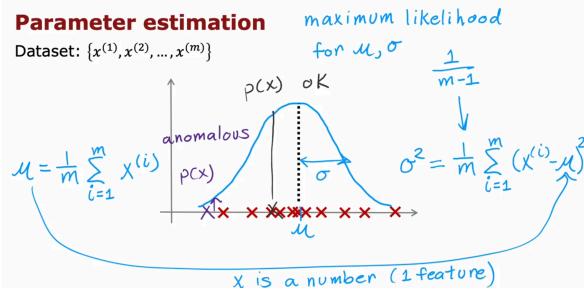


Anomaly Detection:

- Anomaly detection places data instances based off of density which is dependent on high probability to low probability. By then comparing test data with the density ranges we can determine if the test data is an anomaly or not.



- Normal Distribution leverages standard deviation and variance as parameters to plot data points. The farther away from the center you get, the higher the likelihood of the data point being *anomalous*.



- To determine if a new instance of x is an anomaly or not, we determine its probability ' $p(x)$ '. If $p(x) < \text{epsilon}$ (*a small number*) we consider the new instance x as anomalous.

Anomaly detection algorithm

1. Choose n features x_i that you think might be indicative of anomalous examples.

2. Fit parameters $\mu_1, \dots, \mu_n, \sigma_1^2, \dots, \sigma_n^2$

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$$

3. Given new example x , compute $p(x)$:

$$p(x) = \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2) = \prod_{j=1}^n \frac{1}{\sqrt{2\pi}\sigma_j} \exp\left(-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}\right)$$

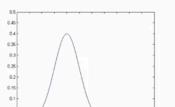
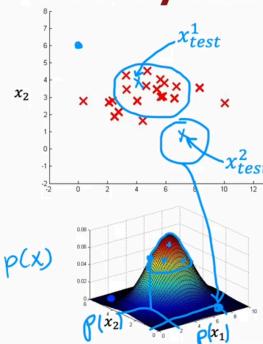
Anomaly if $p(x) < \epsilon$

Vectorized formula

$$\vec{\mu} = \frac{1}{m} \sum_{i=1}^m \vec{x}^{(i)}$$

$$\vec{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{bmatrix}$$

Anomaly detection example



$$\epsilon = 0.02$$

$$p(x_1^{(1)}) = 0.0426 \rightarrow \text{"OK"}$$

$$p(x_2^{(2)}) = 0.0021 \rightarrow \text{anomaly}$$

- Evaluating an Anomaly Algorithm by using CV and Test sets to tune *epsilon* and predict *anomaly* of data instance using $p(x)$:

- Use precision, recall and F1-scores as possible evaluation metrics for the data.

The importance of real-number evaluation

When developing a learning algorithm (choosing features, etc.), making decisions is much easier if we have a way of evaluating our learning algorithm.

Assume we have some labeled data, of anomalous and non-anomalous examples.

$$y=1 \quad y=0$$

Training set: $x^{(1)}, x^{(2)}, \dots, x^{(m)}$ (assume normal examples/not anomalous)

$y=0$ for all training examples

Cross validation set: $(x_{cv}^{(1)}, y_{cv}^{(1)}), \dots, (x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})})$

} include a few anomalous examples
mostly normal examples
 $y=1$
 $y=0$

Test set: $(x_{test}^{(1)}, y_{test}^{(1)}), \dots, (x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$

Aircraft engines monitoring example

10000 good (normal) engines
20 flawed engines (anomalous)
 $y=0$
 $y=1$

Training set: 6000 good engines

CV: 2000 good engines ($y = 0$)
use cross validation set

Test: 2000 good engines ($y = 0$), 10 anomalous ($y = 1$)
tune E tune x_j

2 to 50

$y=1$

train algorithm on training set

Alternative: No test set use if very few labeled anomalous examples

Training set: 6000 good engines 2 higher risk of overfitting

CV: 4000 good engines ($y = 0$), 20 anomalous ($y = 1$)

tune E tune x_j

- Anomaly Detection vs Supervised Learning

- Anomaly detection models are modeled based off of the normal examples in the dataset, as such anything that deviates from the normal is considered an anomaly allowing for brand new instances of deviations to be quickly flagged as anomalies even if the model has never seen it before. (Detecting previously unseen/new issues that may arise in a products (data centers, engines, phones) lifetime)
- Supervised models are trained to detect similar types of data or “anomalies”. As such supervised learning models wont perform well on new variations of data it has never seen before as it can't predict what it has never experienced during training. (Detecting known issues, instances (weather)).

Anomaly detection vs. Supervised learning

Very small number of positive examples ($y = 1$). (0-20 is common).
Large number of negative ($y = 0$) examples.
 $p(x)$
 $y=1$
Many different “types” of anomalies. Hard for any algorithm to learn from positive examples what the anomalies look like; future anomalies may look nothing like any of the anomalous examples we've seen so far.

Enam

Large number of positive and negative examples.
20 positive examples

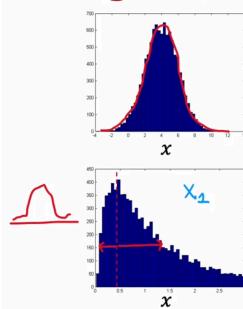
Enough positive examples for algorithm to get a sense of what positive examples are like, future positive examples likely to be similar to ones in training set.

Snam

- Features to choose:

- Turn non-gaussian features into gaussian using relevant mathematical transformations (transformation mustang be same for train/CV/test sets)
- Find or create a new feature (by combining features) to fit gaussian better (sometimes $p(x)$ can be comparable to normal and anomalous examples)

Non-gaussian features



$$\begin{aligned} x_1 &\leftarrow \log(x_1) \\ x_2 &\leftarrow \log(x_2 + 1) \quad \text{log}(x_2 + c) \\ x_3 &\leftarrow \sqrt{x_3} = x_3^{1/2} \\ x_4 &\leftarrow x_4^{1/3} \end{aligned}$$

np.log(x)

Week 2:

Collaborative filtering

Note: For the following notes, consider the example of users rating various movies. This way the notes will be easier to understand.

- To predict the rating of a movie:** W is the weight parameter determined using previous ratings of movies, x is the feature data of how much a certain movie is considered romance or action. (**User must have rated the movie, unrated movies don't count**)

Movie					What if we have features of the movies?		$n_u = 4$	$n_w = 5$	$n = 2$
	Alice(1)	Bob(2)	Carol(3)	Dave(4)	x_1 (romance)	x_2 (action)			
Love at last	5	5	0	0	0.9	0			
Romance forever	5	?	?	0	1.0	0.01			
Cute puppies of love	?	4	0	?	0.99	0			
Nonstop car chases	0	0	5	4	0.1	1.0			
Swords vs. karate	0	0	5	?	0	0.9			

For user 1: Predict rating for movie i as: $w^{(1)} \cdot x^{(i)} + b^{(1)}$ just linear regression
 $w^{(1)} = [5] \quad b^{(1)} = 0 \quad x^{(3)} = [0 \ 0]$ $w^{(1)} \cdot x^{(3)} + b^{(1)} = 4.95$

For user j : Predict user j 's rating for movie i as $w^{(j)} \cdot x^{(i)} + b^{(j)}$

- Cost function will only sum over movies that have been rated by user j . Note: n in the regularization term is the number of features.

Notation:

$r(i,j) = 1$ if user j has rated movie i (0 otherwise)
 $y^{(i,j)}$ = rating given by user j on movie i (if defined)
 $w^{(j)}$, $b^{(j)}$ = parameters for user j
 $x^{(i)}$ = feature vector for movie i

For user j and movie i , predict rating: $w^{(j)} \cdot x^{(i)} + b^{(j)}$
 $m^{(j)}$ = no. of movies rated by user j

To learn $w^{(j)}$, $b^{(j)}$

Cost function

To learn parameters $w^{(j)}, b^{(j)}$ for user j :

$$J(w^{(j)}, b^{(j)}) = \frac{1}{2} \sum_{i:r(i,j)=1} (w^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{k=1}^n (w_k^{(j)})^2$$

To learn parameters $w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}, \dots, w^{(n_u)}, b^{(n_u)}$ for all users :

$$J\left(\begin{array}{l} w^{(1)}, \dots, w^{(n_u)} \\ b^{(1)}, \dots, b^{(n_u)} \end{array}\right) = \frac{1}{2} \sum_{j=1}^{n_u} \sum_{i:r(i,j)=1} (w^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (w_k^{(j)})^2$$

- In a scenario where features are not available, we can use parameter w , x and b to learn what values of the features (**this is considering that the parameters have been learnt in advance**).
- In collaborative filtering, we can come up with the values of features thanks to having data from multiple sources (Ex: Ratings from various users), as such this methodology can't work in typical linear regression as we lack the necessary information to make such deductions.

Movie	Problem motivation			
	Alice (1)	Bob (2)	Carol (3)	Dave (4)
Love at last	5	5	0	0
Romance forever	5	?	?	0 ←
Cute puppies of love	?	4	0	? ←
Nonstop car chases	0	0	5	4 ←
Swords vs. karate	0	0	5	? ←

$w^{(1)} = [5] \quad w^{(2)} = [5] \quad w^{(3)} = [0] \quad w^{(4)} = [0]$ } using $w^{(j)}, x^{(i)} + b^{(j)}$
 $w^{(1)} \cdot x^{(1)} \approx 5 \quad w^{(2)} \cdot x^{(1)} \approx 5 \quad w^{(3)} \cdot x^{(1)} \approx 0 \quad w^{(4)} \cdot x^{(1)} \approx 0$ } $\rightarrow x^{(1)} = [1] \quad b^{(1)} = 0, b^{(2)} = 0, b^{(3)} = 0, b^{(4)} = 0$ ← Activate when learning to make deductions

Cost function

Given $w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}, \dots, w^{(n_u)}, b^{(n_u)}$

to learn $x^{(i)}$:

$$J(x^{(i)}) = \frac{1}{2} \sum_{j:r(i,j)=1} (w^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{k=1}^n (x_k^{(i)})^2$$

→ To learn $x^{(1)}, x^{(2)}, \dots, x^{(n_m)}$:

$$J(x^{(1)}, x^{(2)}, \dots, x^{(n_m)}) = \frac{1}{2} \sum_{i=1}^{n_m} \sum_{j:r(i,j)=1} (w^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2$$

- In order to get the **parameters** for different users, we put together the algorithm to learn parameters and the algorithm to learn features.
- In this aggregated algorithm (**Collaborative Filtering**) we sum over **each movie** for **each user** (not the other way around) for **only the movies that have been rated**.
- Helps us find w, b and x all at the same time (x also becomes a parameter).

Collaborative filtering

Cost function to learn $w^{(1)}, b^{(1)}, \dots, w^{(n_u)}, b^{(n_u)}$:

$$\min_{w^{(1)}, \dots, w^{(n_u)}, b^{(n_u)}} \frac{1}{2} \sum_{j=1}^{n_u} \sum_{i=1, r(i,j)=1}^{n_m} (w^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (w_k^{(j)})^2$$

Cost function to learn $x^{(1)}, \dots, x^{(n_m)}$:

$$\min_{x^{(1)}, \dots, x^{(n_m)}} \frac{1}{2} \sum_{i=1}^{n_m} \sum_{j=1, r(i,j)=1}^{n_u} (w^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2$$

Put them together:

$$\min_{\substack{w^{(1)}, \dots, w^{(n_u)}, \\ b^{(1)}, \dots, b^{(n_u)}, \\ x^{(1)}, \dots, x^{(n_m)}}} J(w, b, x) = \frac{1}{2} \sum_{\substack{(i,j): r(i,j)=1}} (w^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (w_k^{(j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2$$

Gradient Descent
collaborative filtering

Linear regression (course 1)

repeat {

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

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

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

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

}

parameters w, b, x x is also a parameter

- Using collaborative data (ratings) from multiple users for a single/multiple subject (movies), gives us the necessary information to appropriately guess the features of the subject (movie), which in turn helps us predict what potential user's (who have not rated the movie) may rate the movie in the future.
- Recommender systems are particularly efficient when working with **binary labels**
- In order to generalize a linear regression collaborative filtering model to work with binary labels we predict the probability of the **output label** being **1** (meaning the subject has been engaged, i.e being shown a movie), whereas before we'd predict whatever the **value of the output label is**.

Example applications

- Did user j purchase an item after being shown? $1, 0, ?$
- Did user j fav/like an item? $1, 0, ?$
- Did user j spend at least 30sec with an item? $1, 0, ?$
- Did user j click on an item? $1, 0, ?$

Meaning of ratings:
 \rightarrow 1 - engaged after being shown item
 \rightarrow 0 - did not engage after being shown item
 \rightarrow ? - item not yet shown

From regression to binary classification

- Previously:
- Predict $y^{(i,j)}$ as $w^{(j)} \cdot x^{(i)} + b^{(j)}$
- For binary labels:
Predict that the probability of $y^{(i,j)} = 1$ is given by $g(w^{(j)} \cdot x^{(i)} + b^{(j)})$
where $g(z) = \frac{1}{1+e^{-z}}$

Cost function for binary application

Previous cost function:

$$\frac{1}{2} \sum_{(i,j): r(i,j)=1} (w^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (w_k^{(j)})^2$$

Loss for binary labels $y^{(i,j)} : f_{(w,b,x)}(x) = g(w^{(j)} \cdot x^{(i)} + b^{(j)})$

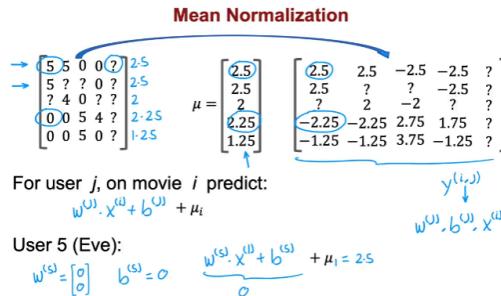
$$L(f_{(w,b,x)}(x), y^{(i,j)}) = -y^{(i,j)} \log(f_{(w,b,x)}(x)) - (1 - y^{(i,j)}) \log(1 - f_{(w,b,x)}(x))$$

Loss for single example

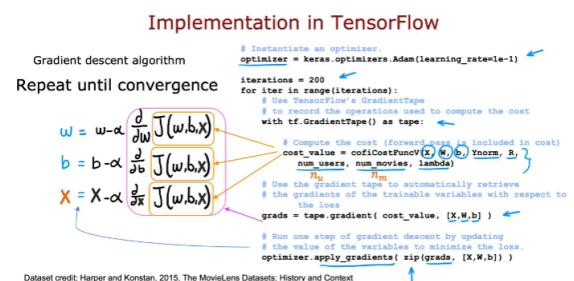
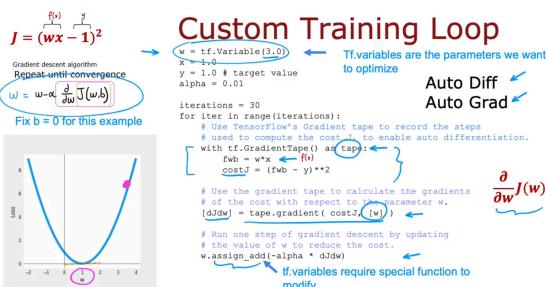
$$J(w, b, x) = \sum_{(i,j): r(i,j)=1} L(f_{(w,b,x)}(x), y^{(i,j)})$$

Recommender system implementation:

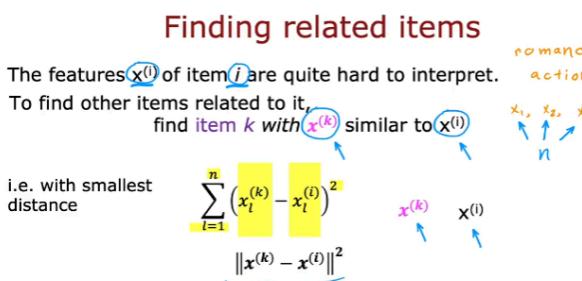
- Normalization also helps recommendation systems through the usage of mean normalization. This is important as without normalization, for a user without any ratings (considering movie example), the algorithm will predict '0', for all movies due to their parameters w and b being initialized as '0'
- In **mean normalization**, you convert all data points into a matrix and discover the mean of each class (in our case, each movie's mean rating). The mean is then subtracted from all the original data points (ratings) unless it is unknown/missing/empty.
- By implementing the mean into the algorithm's calculations, we gain a baseline prediction (for a movie's rating) when a user has no prior/reference data for said subject (rating)
- Normalizing works regardless of application on the column or the row, but its better to choose based on context (choosing row for users with no prior ratings, choosing column for movies with no prior rating).



- Tensorflow implementation of gradient descent and differentiation followed by **implementation of collaborative filtering**:



- Collaborative filtering algorithms can also help to find related items by finding features of item k similar to features of item i . This is achieved with the help of distance, as the shortest distances between the features of two items, the more related they are.



- Such an algorithm also suffers from **cold start** (How to rank new items that have very little data such as ratings? How to show new users with few rated items a reasonable set of items?) and **the inability to use side/additional information about items or users** (genre, actors, genders, location, age).

Content-based filtering:

- While **collaborative filtering** recommends items based off of similarity of ratings between users, **content-based filtering** recommends based on the features of the user and items to find a good match.

Examples of user and item features

User features:

- Age
- Gender (1 hot)
- Country (1 hot, 200)
- Movies watched (1000)
- Average rating per genre
- ...

Movie features:

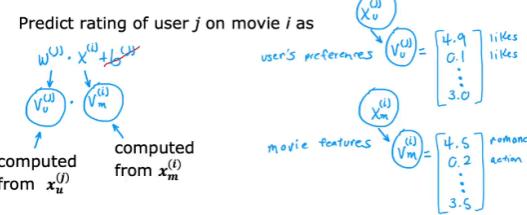
- Year
- Genre/Genres
- Reviews
- Average rating
- ...

$x_u^{(j)}$ for user j

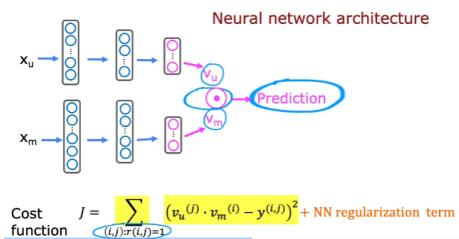
$x_m^{(i)}$ for movie i

- Try to determine if movie j will be a good match for user i . \mathbf{V}_m (Movie vector, parameter captures the likelihood of a set of genres) and \mathbf{V}_u (User vector, parameter may capture preferences, etc.) [must be the same size].

Content-based filtering: Learning to match



- To compute \mathbf{V}_u and \mathbf{V}_m using \mathbf{X}_u and \mathbf{X}_m , we make use of neural networks with \mathbf{X}_u and \mathbf{X}_m as inputs and \mathbf{V}_u and \mathbf{V}_m as outputs. The dot product of \mathbf{V}_u and \mathbf{V}_m is then used to predict the probability of the output (likelihood of a good match) being 1.

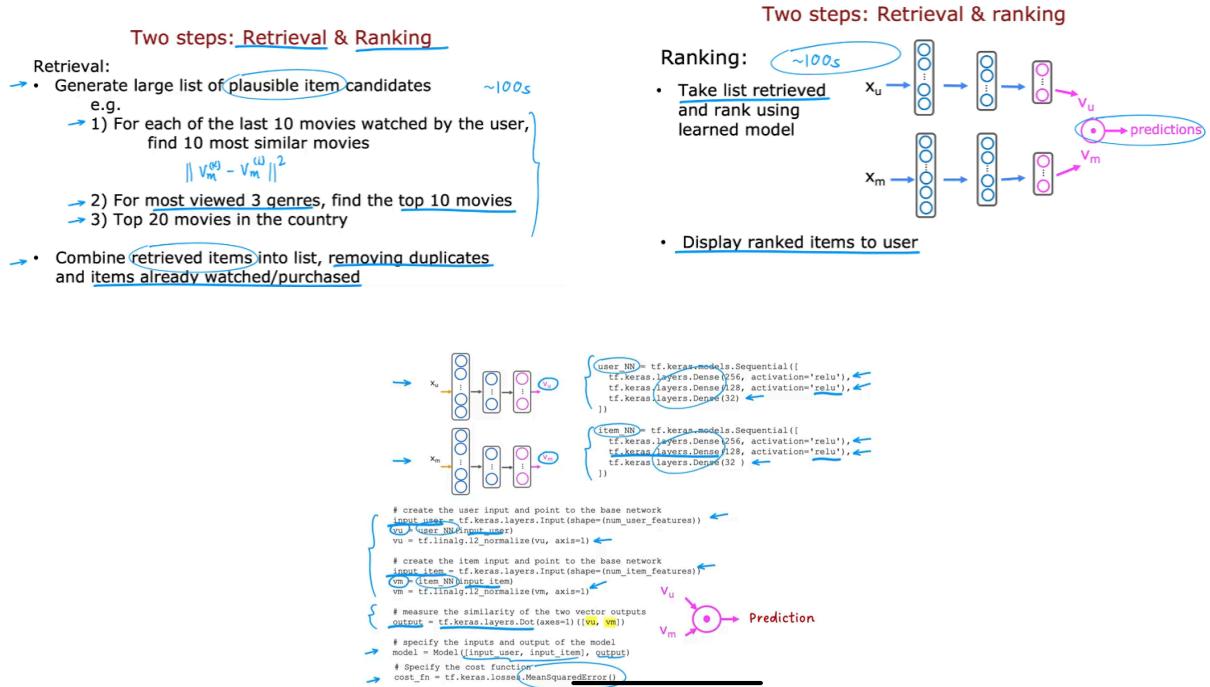


Learned user and item vectors:
 $v_u^{(j)}$ is a vector of length 32 that describes user j with features $x_u^{(j)}$
 $v_m^{(i)}$ is a vector of length 32 that describes movie i with features $x_m^{(i)}$

To find movies similar to movie i : $\|V_m^{(i)} - V_m^{(j)}\|^2$ small
 $\|x_m^{(i)} - x_m^{(j)}\|^2$

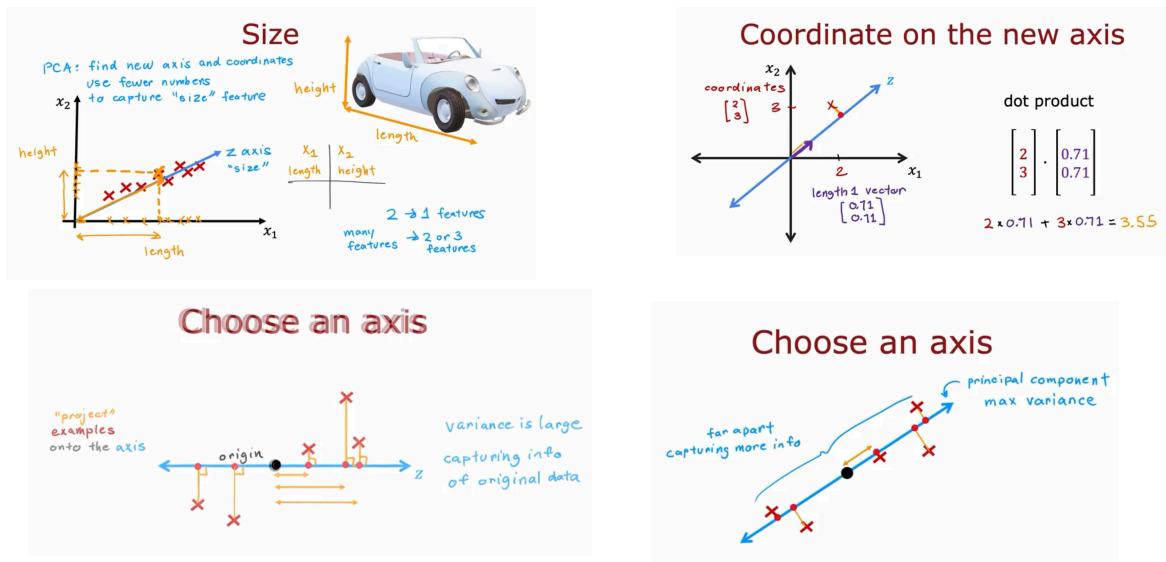
Note: This can be pre-computed ahead of time

- To efficiently recommend a handful of relevant items we employ **retrieval & ranking**. Retrieving more items yields better performance at slower recommendations. Optimize trade-off by checking if retrieving additional items yields better recommendations

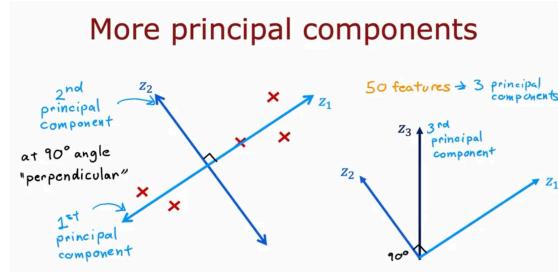


Principal Component Analysis:

- PCA is an unsupervised learning algorithm that can help convert 3d or multi dimensional plots to 2d for easier visualization
- PCA tends to compress features to decrease large amounts of features into a smaller amount
- Ensure to normalize and feature scale the features before finding the new axis
- The new axis is similar to a best fit line where variance is maximum, and we capture the most amount of information from the original feature data. This new axis is called the principal component

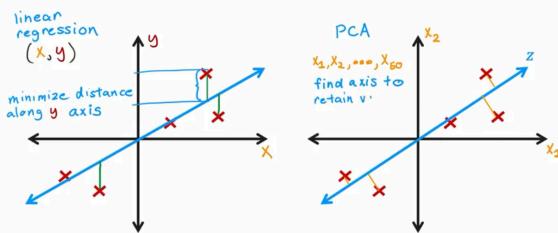


- A second principal component will be **90 degree to the first principal component** (similar for additional principal components)

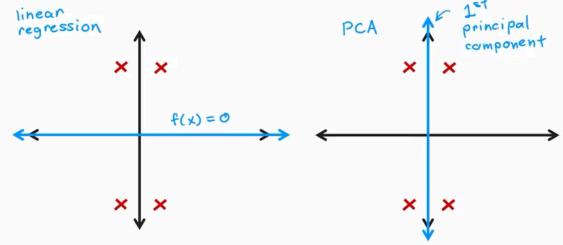


- Linear regression however focuses on finding the best fit line to predict output class y . In PCA there is no output class, as such the new axis is treating all features equally and trying to retain as much information from all features as possible. So even if the 2D lines may seem the same, they are in truth, not the same.

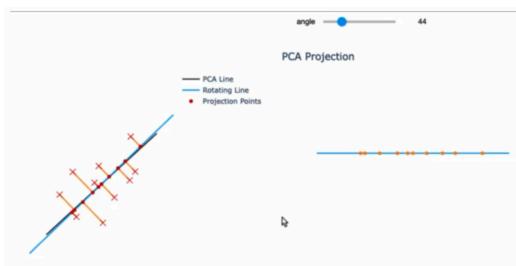
PCA is not linear regression



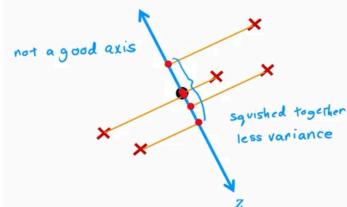
PCA is not linear regression



- Another way of thinking is that PCA will find an axis that distributes the features such that the features don't **overlap** or aren't **squished** together (if they are then we can't draw much information from that new axis)

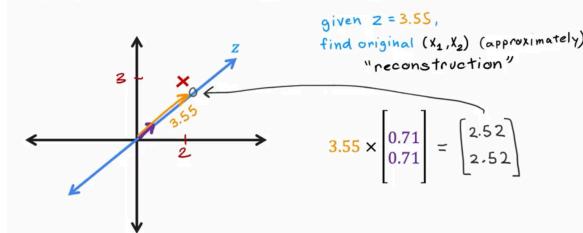


Choose an axis

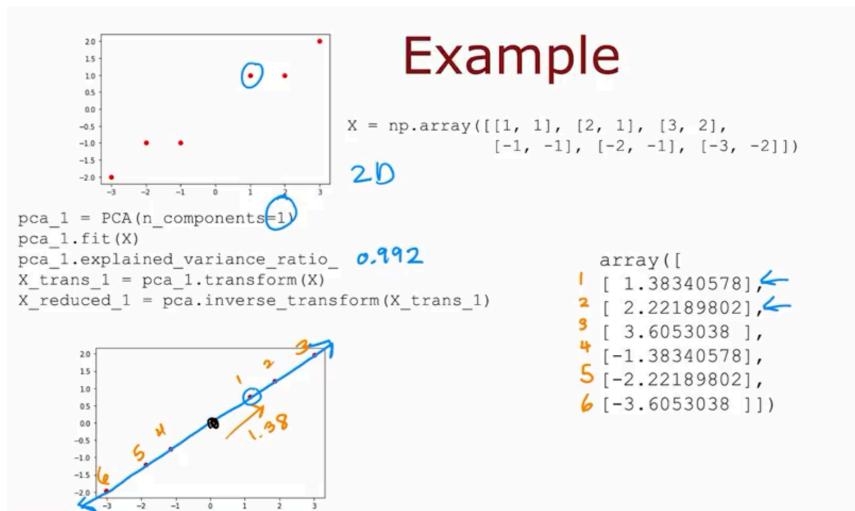


- We can reconstruct the original training example using the distance from origin to the new coordinate point on the principal component and multiplying with the vector of said principal component. This way we get an approximation of the original training examples from the original axis

Approximation to the original data



- Implementation:



Week 3:

Reinforcement learning:

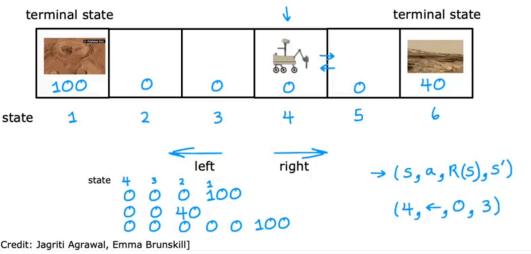
Applications

- • Controlling robots
- • Factory optimization
- • Financial (stock) trading
- Playing games (including video games)



- We know reinforcement learning uses reward systems to incentivize models to train themselves in a manner that progresses them towards a goal. This is achieved by creating **states s** , with each state having a set **reward $R(s)$** , and the model starts at a certain state and can then take an **action a** , to move to some new state s' at every **time step**. Combined we get $(s, a, R(s), s')$.

Mars Rover Example



- The **Return** is the aggregate of all rewards experienced in each state until reaching the terminal state, with each reward being multiplied by a weight factor called the **discount factor**.
- The **discount factor** makes the algorithm less patient by decreasing the **value** or **credit** of the rewards, with the value decreasing the more states/time the algorithm takes to reach the terminal. This urges the algorithm to pick greater rewards sooner for higher overall yield.
- If **negative rewards** are part of the system, the discount factor is actually incentivized to **push out the negative rewards as far into the future as possible**.

Return

Example of Return

100	50	25	12.5	6.25	40
100	0	0	0	0	40

← return
← reward

100	0	0	0	0	40	
state	1	2	3	4	5	6

The return depends on the actions you take.

100	2.5	5	10	20	40
100	0	0	0	0	40

$$\text{Return} = 0 + (0.9)^1 0 + (0.9)^2 100 = 0.729 \times 100 = 72.9$$

$$\text{Return} = R_1 + \gamma R_2 + \gamma^2 R_3 + \dots \text{ (until terminal state)}$$

$$\text{Discount Factor } \gamma = 0.9 \quad 0.99 \quad 0.999$$

$$R = 0.5$$

$$\text{Return} = 0 + (0.5)^0 0 + (0.5)^1 100$$

100	2.5	5	10	20	40
100	0	0	0	0	40

$$+ (0.5)^2 0 + (0.5)^3 40$$

100	2.5	5	12.5	20	40
100	0	0	0	0	40

$$+ (0.5)^4 0 + (0.5)^5 40$$

- We make use of a **policy π** , to map what action to take for a certain step. Thus: $\pi(s) = a$

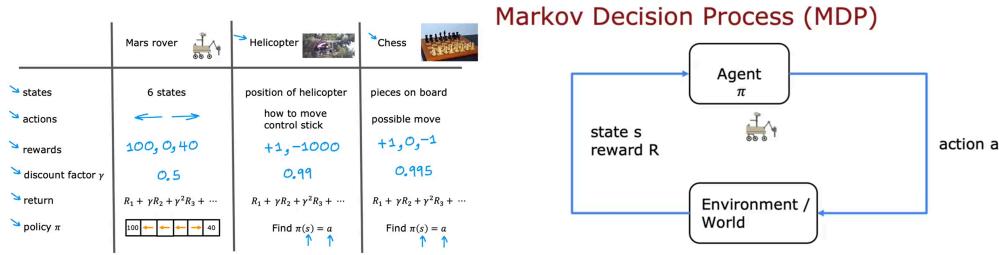
Policy

100	←	←	←	←	←	40
100	←	←	←	←	←	40
100	→	→	→	→	→	40
100	←	←	←	←	→	40

$$\begin{aligned}\pi(1) &= a \\ \pi(2) &= \leftarrow \\ \pi(3) &= \leftarrow \\ \pi(4) &= \leftarrow \\ \pi(5) &= \rightarrow\end{aligned}$$

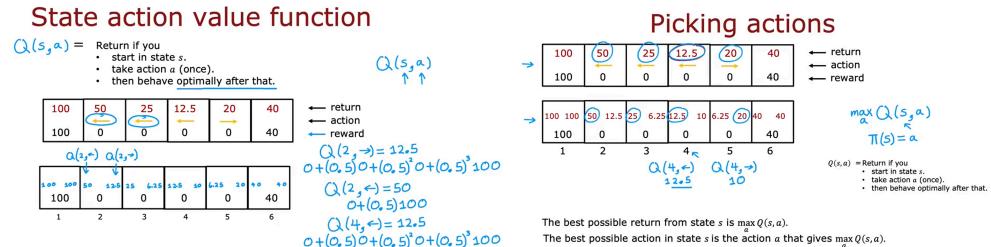
A policy is a function $\pi(s) = a$ mapping from states to actions, that tells you what action a to take in a given state s .

- The formalization of **states**, **rewards**, **actions**, **discount factor**, **return** and **policy** are formally known as **Markov Decision Process (MDP)** where the future is dependent only on the current state, the past matters not.



State-action value function

- State-action value function will provide a **return** based on the state and singular action specified, considering that the system behaves optimally after executing the specified action.
- This requires a system with defined policies and rewards to work.
- If, however, the **State-action value is defined $Q(s, a)$** , we can use it to determine the best action for a given state based on the maximum state-action value in a given state.



- Bellman equation** is the simplification of the **state-action value function** to calculate the state-action values for a given state with different initial actions as long as we know the max state-action value for the next state. (Ex: In the Mars rover's case, we calculate the state-action value for state 4 with action for going left and right. If we choose to go left, we need the MAX state-action value of state 3 and if we go right we will need that of state 5)

Explanation of Bellman Equation

$$\begin{aligned} Q(s, a) &= \text{Return if you} \\ &\quad \cdot \text{ start in state } s, \\ &\quad \cdot \text{ take action } a \text{ (once),} \\ &\quad \cdot \text{ then behave optimally after that.} \end{aligned}$$

$s \rightarrow s'$

\rightarrow The best possible return from state s' is $\max Q(s', a')$

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

Reward you get right away Return from behaving optimally starting from state s'

$$Q(s, a) = R_1 + \gamma [R_2 + \gamma R_3 + \gamma^2 R_4 + \dots]$$

Note: While bellman's equation and the state-action value function looks similar and are virtually the same, notice that bellman's requires us to know the maximum state-action value of proceeding steps or, at the very least, calculate it.

- Sometimes going in the best direction may not be the most reliable/feasible (going left might be the quickest way to get home, but I might slip and go right instead).
- To account for this, we may assign a probability to the potential states (a higher probability to the direction with higher state-action value and vice-versa)
- Over many iterations we will receive varying kinds of reward sequences/state action values

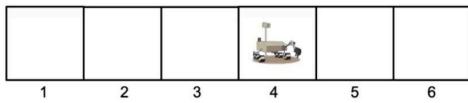
- In a **stochastic reinforcement learning** problem, what we're interested in is not maximizing the return because that's a random number. What we're interested in is **maximizing the average value of the sum of discounted rewards**. Say, if we were to take a policy and try it out a large amount of times, we'd get **lots of different reward sequences**. If we were to take the **average over all of these different sequences of the sum of discounted rewards**, then that's what we call the **expected return**.

Continuous state reinforcement learning algorithm

- In a continuous state reinforcement learning algorithm (**Continuous MTP**) we use vectors of numbers or **continuous state spaces** to quantify states.
- Instead of dealing with **discrete** set of states (rover being in one of six positions), we consider a large number of **continuous** value positions (rover could have a position anywhere between 0 km to 6 km, such as 4.2 km)
- Additionally, a state can consist of many numbers (vector of numbers) that may track multiple factors such as the velocity of the rover in **x** and **y** directions, its coordinates **x'** and **y'**, its angle **theta**.
- For any specific state, any of these numbers can take on any value in a **valid range**

Discrete vs Continuous State

Discrete State:



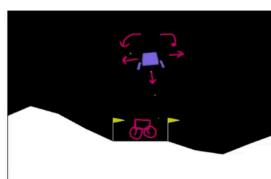
Continuous State:



$$s = \begin{bmatrix} x \\ y \\ \theta \\ \dot{x} \\ \dot{y} \\ \dot{\theta} \end{bmatrix}$$

- By creating a complex reward function we can incentivize correct behaviors and actions while demotivating incorrect ones, doing this is far less laborious than specifying the exact correct action for every state

Lunar Lander

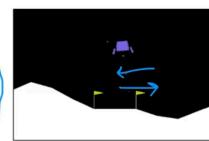


actions:
do nothing
left thruster
main thruster
right thruster

$$s = \begin{bmatrix} x \\ y \\ \dot{x} \\ \dot{y} \\ \theta \\ \dot{\theta} \\ l \\ r \end{bmatrix}$$

Reward Function

- Getting to landing pad: 100 – 140
- Additional reward for moving toward/away from pad.
- Crash: -100
- Soft landing: +100
- Leg grounded: +10
- Fire main engine: -0.3
- Fire side thruster: -0.03



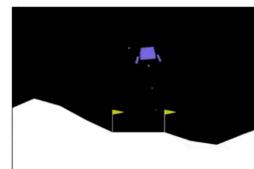
Note - In the lunar landing example, the state spaces comprise of:

- Position x** (Left/Right)
- Position y** (Up/Down)
- x'** (Horizontal velocity)
- y'** (Vertical velocity)
- Theta** (Angle: Left/Right Tilt)
- Theta'** (Angle Velocity)
- L** (Left leg grounded)
- R** (Right leg grounded)

Lunar Lander Problem

Learn a policy π that, given

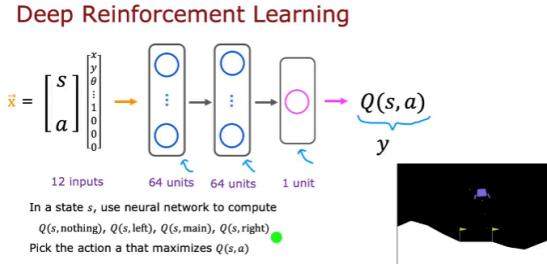
$$s = \begin{bmatrix} x \\ y \\ \dot{x} \\ \dot{y} \\ \theta \\ \dot{\theta} \\ l \\ r \end{bmatrix}$$



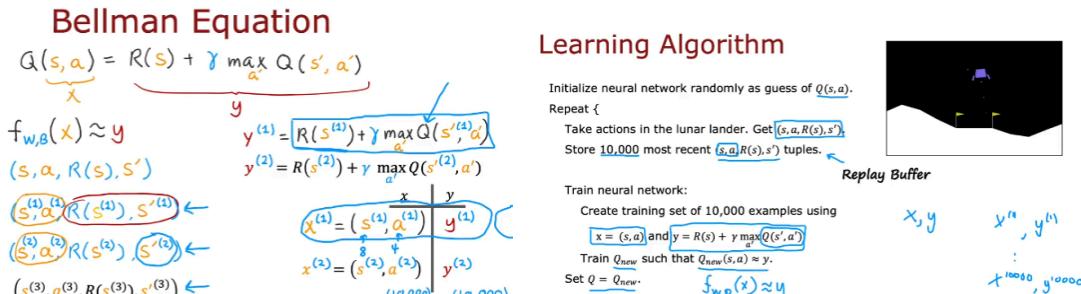
picks action $a = \pi(s)$ so as to maximize the return.

$$\gamma = 0.985$$

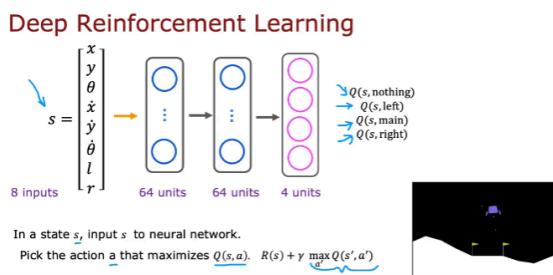
- By feeding a Deep Reinforcement Learning Model **a combination of vectorized state values (8 numbers as seen previously) and one-hot encoded actions (one-hot encoded 4 numbers corresponding to the action)**, we try to receive a **continuous state-action value, $Q(s,a)$** .
 - By computing all actions a for a constant state s (so the state remains the same, but we simply change the actions), we can determine which action yields the highest state-action value, helping us find the maximum $Q(s, a)$ or y .



- We can use the **bellman equation** to generate a **training set** of inputs x to corresponding outputs y by either taking the **policy based actions** or **random actions** if no policy is set. We also take a **random initial guess** for the **state-action value of the next state $Q(s', a')$** as we don't know it yet.
 - By taking different actions, we will get multiple sets of: state s , action a , reward for state $R(s)$ and next state s' .
 - Training the model will slowly give us a more **accurate Q** for us to replace during initialization instead of **randomly guessing it**



- Carrying out inference/computation for each action separately in the deep reinforcement learning model is **inefficient**. **The updated model can compute all actions for a singular state, simultaneously**



- Taking **initial actions** (does not need to be optimal) is one of the most **important** of the reinforcement algorithms, as without doing so, we cannot generate a training set to train the model on. **But what is the best way to choose these actions?**

How to choose actions while still learning?

In some state s

Option 1:
Pick the action a that maximizes $Q(s, a)$.

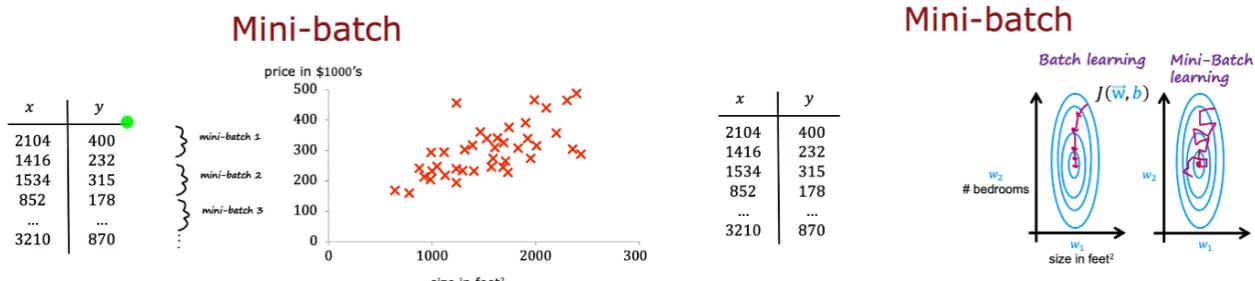
Option 2:

- With probability 0.95, pick the action a that maximizes $Q(s, a)$. **Greedy**, "Exploitation"
- With probability 0.05, pick an action a randomly. **Exploration**"

ϵ -greedy policy ($\epsilon = 0.05$)
0.95

Start ϵ high
 $1.0 \rightarrow 0.01$
Gradually decrease

- Always trying to maximize $Q(s, a)$ is not a good idea as it can trap us in a local minimum. **Ex:** Assume that parameters were initialized such that $Q(s, \text{Thruster})$ is always low, the model will never fire the Thruster while trying to maximize $Q(s, a)$ because it thinks its a bad decision, as such it will never discover that firing the Thrusters could possibly be a good action that would lead it to a better $Q(s, a)$.
- Thus, **Option 2 is better**, known as **ϵ -greedy policy**. It is good to start with a epsilon value that's initially high and gradually decrease it until the ideal performance is found
- We can use **mini-batches** to speed-up a reinforcement learning algorithm, where we set samples/batches from an entire set of training examples (For a training set of size 100,000,000, we may create batches for every 1000 examples).
- For every **one iteration** we choose one unique/different batch that has not been previously used



- Mini-batch gradient descent is useful when dataset is exceptionally large, allowing for faster training and computation, Batch gradient descent is useful when we have smaller datasets and it converges to the minima consistently
- Going back to the lunar landing, we may use 1000 examples per training iteration for our training set out of the 10,000
- Additionally updating Q to the newly computed Q is not always ideal as it may overwrite it with a worse one
- To mitigate use a **soft update** where we may set hyperparameters which determine how aggressively Q will move towards Q_{new} . This soft update helps to converge more reliably

Soft Update

Set $Q = Q_{\text{new}}$ ← $Q(s, a)$

w, b $w_{\text{new}}, b_{\text{new}}$

$$W = 0.01 W_{\text{new}} + 0.99 W$$

$$B = 0.01 B_{\text{new}} + 0.99 B$$