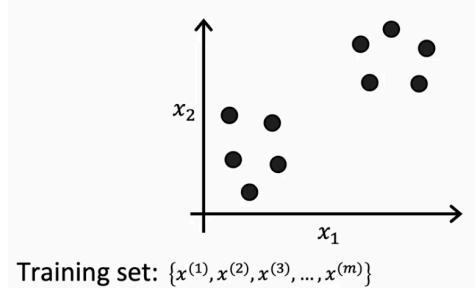
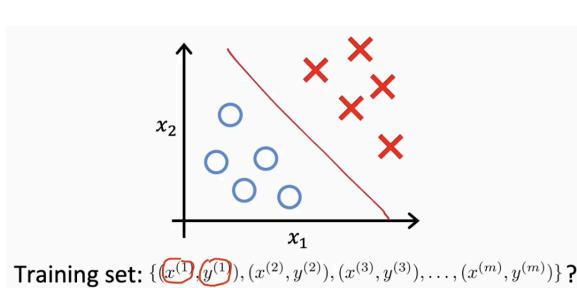


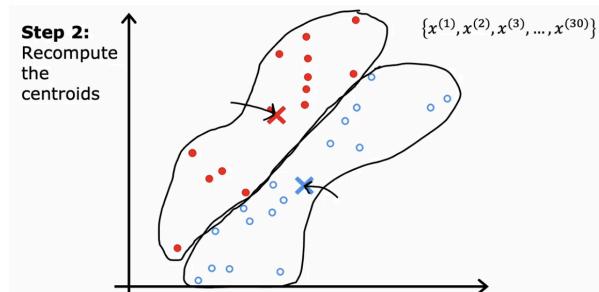
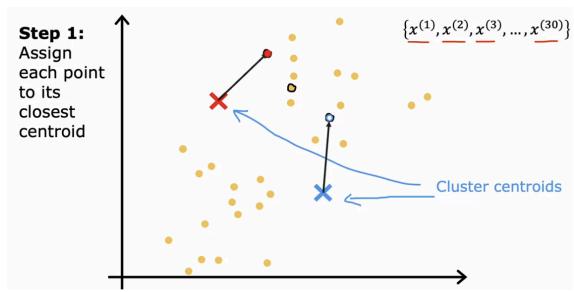
Week 1: Unsupervised Learning

Clustering algorithm

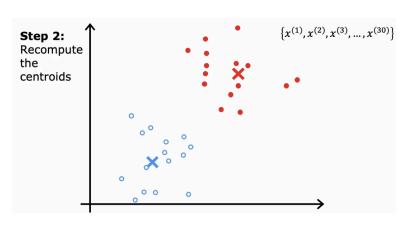
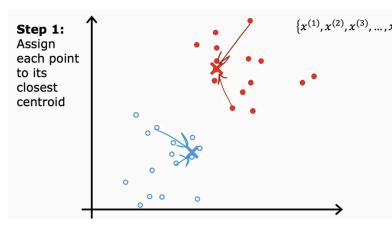
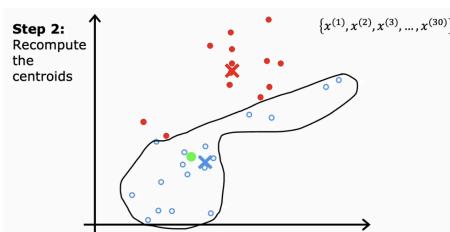
- Looks at a number of data points and evaluates which are similar
- In supervised learning, dataset includes input x and target output y



- Clustering algorithms look for clusters. Applications includes grouping similar news, market segmentation, DNA analysis, astronomical data analysis
 - Most common clustering algorithm is *K-means*:
 - Will take a random initial guess at centres of two clusters.
 - Will do two things: assign cluster centroids, and move cluster centroids
- 1) After making an initial guess, goes through all examples x , and for each will check whether it is closer to red or blue cluster centroid, and will assign each point to its closer centroid
 - 2) Will look at all red points and take an average, and will move red cross to average location of all the red dots



3) Then repeat, and check whether each point is closer to red or blue cluster centroid, and then compute average and move X again



4) Eventually K-means clustering algorithm converges i.e. there are no more changes

Details of K-means algorithm

Technical steps:

K-means algorithm

μ_1, μ_2 $x^{(1)}, x^{(2)}, \dots, x^{(n)}$
n = 2 K = 2

Randomly initialize K cluster centroids $\mu_1, \mu_2, \dots, \mu_K$

Repeat {

- # Assign points to cluster centroids
- for $i = 1$ to m
- $c^{(i)} :=$ index (from 1 to K) of cluster centroid closest to $x^{(i)}$
- # Move cluster centroids
- for $k = 1$ to K
- $\mu_k :=$ average (mean) of points assigned to cluster k

}

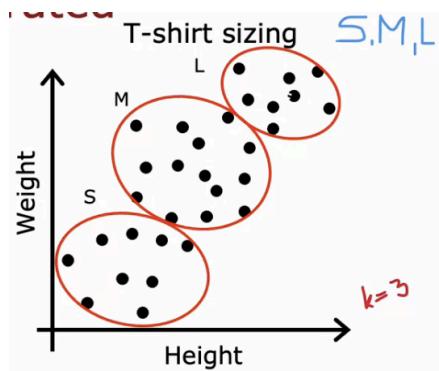
$$\mu_1 = \frac{1}{4} [x^{(1)} + x^{(5)} + x^{(6)} + x^{(10)}]$$

Notes;

- Mu1 and Mu2 have same dimensions as training data - here, they will be two-dimensional vectors
 - If a cluster has zero training examples assigned - second step would be trying to compute average of zero points (not well defined). Commonly, just eliminated cluster to give $K=k-1$ cluster; or randomly re-initialise initial guess.

K-means for clusters that are not well separated

E.g. T-shirt size:



Optimising objective (Cost function for K-means)

- Average square distance between every training examples $x(i)$ and the location of the cluster centroid to which $x(i)$ has been assigned

K-means optimization objective

$c^{(i)}$ = index of cluster ($1, 2, \dots, K$) to which example $x^{(i)}$ is currently assigned
 μ_k = cluster centroid k

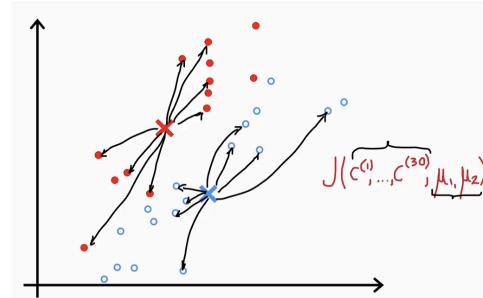
$\mu_{c^{(i)}}$ = cluster centroid of cluster to which example $x^{(i)}$ has been assigned

Cost function

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

$$\min_{c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K} J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$$

Distortion



- In more detail:
 - First part of K-means (assign points to cluster centroids) - trying to update $c(1)$ to $c(m)$ to minimise J as much as possible while holding $\mu(1)$ to $\mu(k)$ fixed
 - Second part (move cluster centroids): trying to update $\mu(1)$ to $\mu(k)$ while keeping $c(i)$ fixed

Cost function for K-means

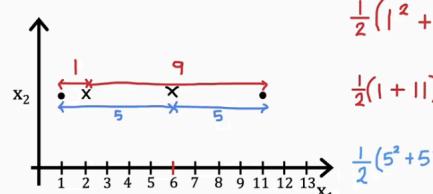
$J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K) = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \mu_{c^{(i)}}\|^2$

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
 - μ_k := average of points in cluster k

}

Moving the centroid



$$\frac{1}{2}(1^2 + 1^2) = \frac{1}{2}(1+1) = 1$$

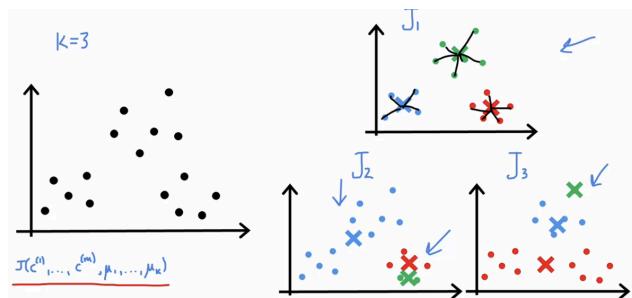
$$\frac{1}{2}(5^2 + 1^2) = 6$$

$$\frac{1}{2}(5^2 + 5^2) = 25$$

- Should always reduce (otherwise there will be a bug in the code)

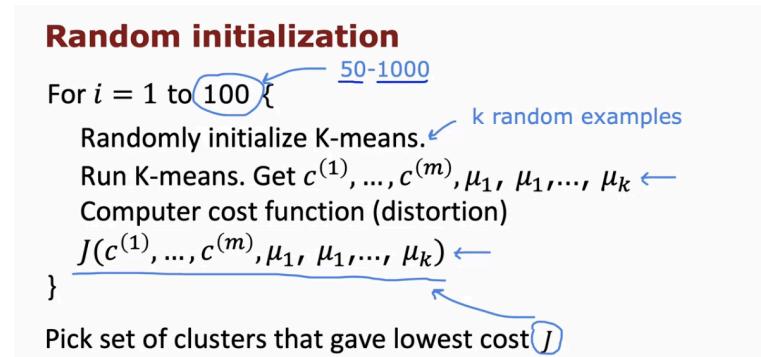
Initializing K-means: how to take the random guess

- Choose number of clusters to be less than training examples m (choose $K < m$)
- Randomly pick K training examples [not completely random]
- Set $\mu(1), \mu(2), \dots, \mu(K)$ equal to these K examples



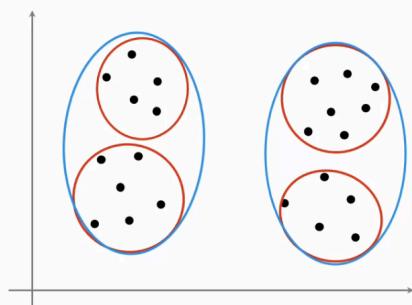
- Run initialisation multiple times to find best starting points - run cost function to evaluate these

The algorithm for this is:



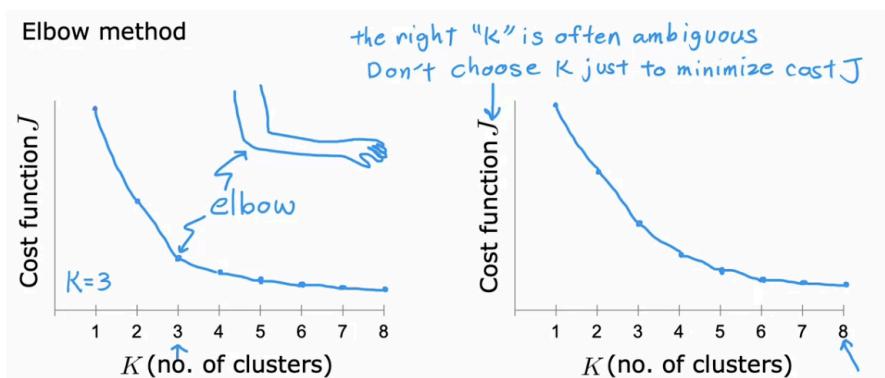
Choosing number of clusters (value of K)

What is the right value of K?

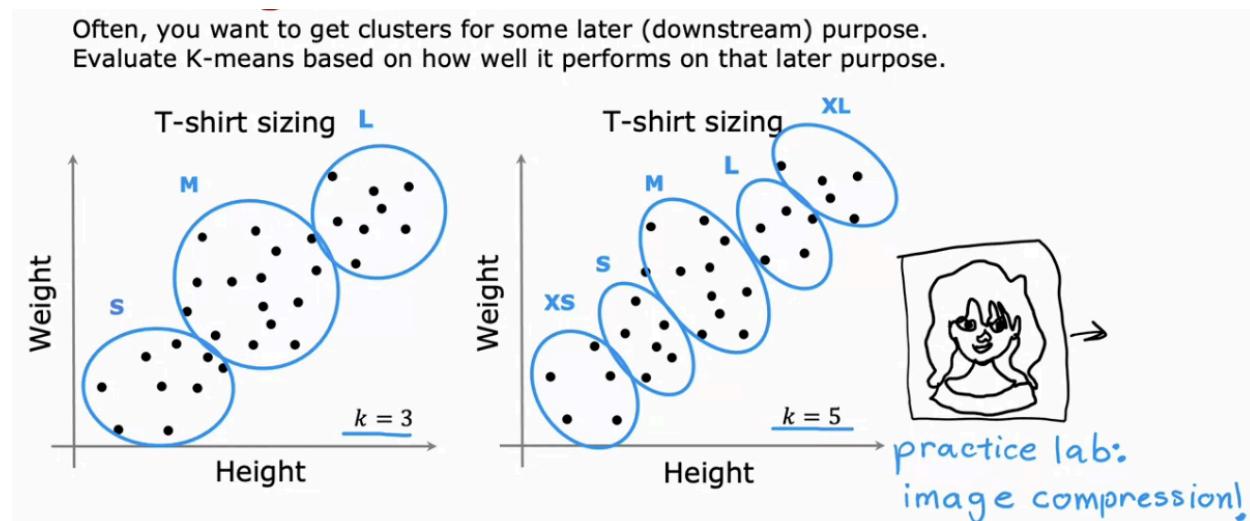


To automatically choose K:

- 1) Elbow method [but not v. good]
 - Plot cost function against value of K
 - But not always a clear elbow
 - Don't choose K just to minimise J (as will generally just be largest K)



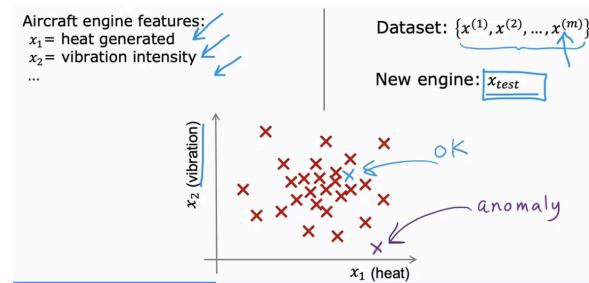
- 2) Evaluate K-means based on how well it performs on its later purpose
 E.g. whether you want to make three or five sizes for T-shirts - see trade-off between better fit vs making more types of t-shirts



- With image compression, there is a trade off between quality of compressed image vs how much you can compress image to save space

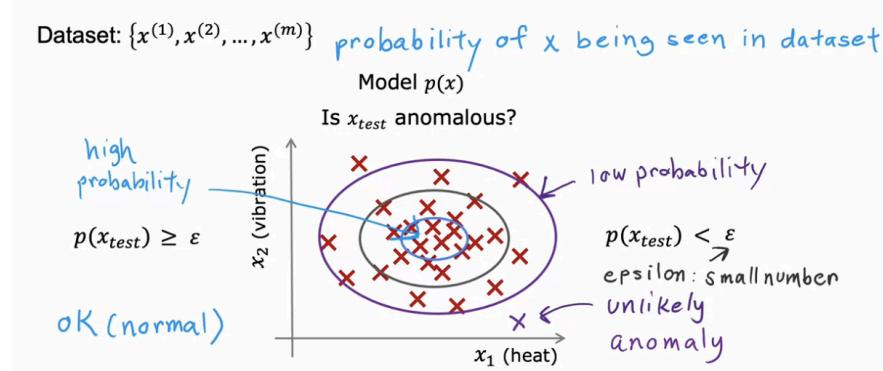
Anomaly detection

Example with aircraft engine manufacture:



The most common way to carry out anomaly detection is to use *density estimation*:

- Probability of x being seen in the dataset



Other examples:

1) Fraud detection:

- How often do they log in?
- How many web pages visited?
- Transactions?
- Posts?
- Typing speed?

Fraud detection:

- $x^{(i)}$ = features of user i 's activities
- Model $p(x)$ from data.
- Identify unusual users by checking which have $p(x) < \varepsilon$

transactions?
posts? typing speed?

perform additional checks to identify real fraud vs. false alarms

2) Manufacturing:

- To see if there is an issue that merits a closer look
- $X[i]$ = features of product i

3) Monitoring computers in a data centre:

$x^{(i)}$ = features of machine i

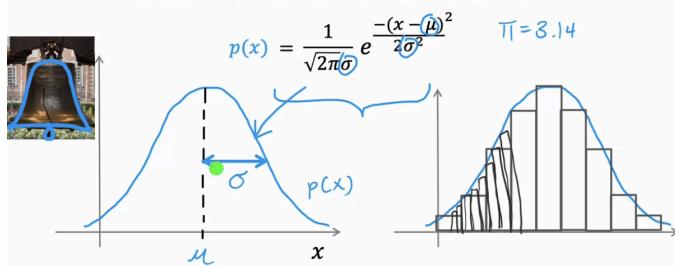
- x_1 = memory use,
- x_2 = number of disk accesses/sec,
- x_3 = CPU load,
- x_4 = CPU load/network traffic.

ratios

Gaussian (normal) distribution

Gaussian (Normal) distribution σ standard deviation σ^2 variance

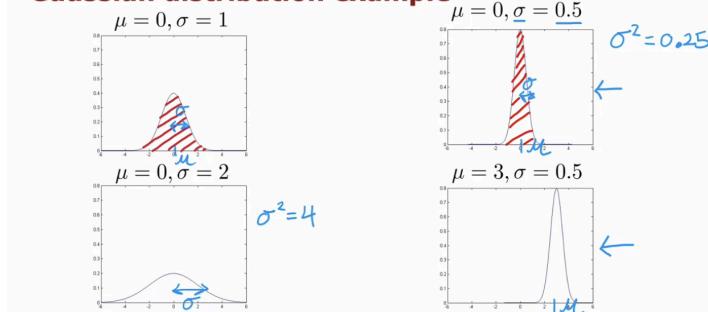
Say x is a number.
Probability of x is determined by a Gaussian with mean μ , variance σ^2 .



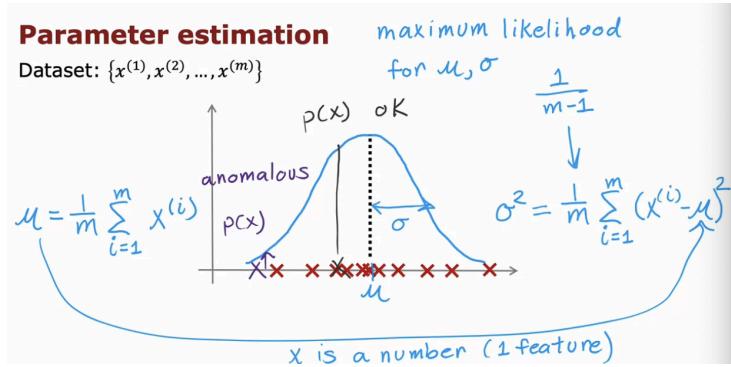
Changing mu and sigma:

- Area under curve stays the same

Gaussian distribution example



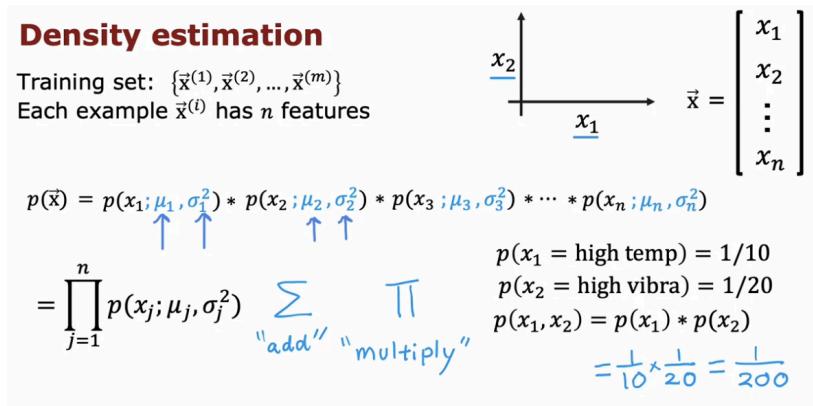
Applying this to anomaly detection: parameter estimation



- This is for if x is a single number (1 feature); but for practical applications you will have multiple features

Density estimation

- Build a model for $p(x)$
- Assumes features are statistically independent



Putting it together into an anomaly detection algorithm

- If just one of $p(x_j)$ is v. small, as $p(x)$ is a product, whole thing will be small

Anomaly detection algorithm

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

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

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

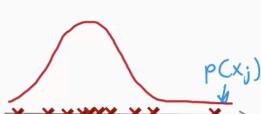
Vectorized formula

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

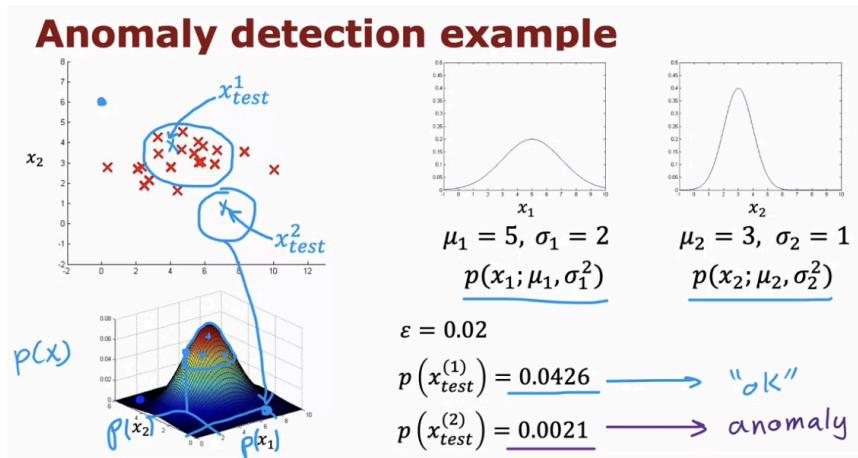
3. Given new example x , compute $p(x)$:

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

Anomaly if $p(x) < \varepsilon$



Anomaly detection example



Developing and evaluating an anomaly detection system

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 ($y=1$) and non-anomalous ($y=0$) examples
- Training set: $x(1), x(2), \dots, x(m)$ (assume normal examples/ not anomalous i.e. $y=0$ for all training examples)

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

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

E.g. aircraft engines monitoring example

10000	good (normal) engines	2 to 50
20	flawed engines (anomalous)	$y=1$
2		
Training set:	6000 good engines	train algorithm on training set
CV:	2000 good engines ($y=0$)	use cross validation set
		tune ϵ tune x_j
Test:	2000 good engines ($y=0$),	10 anomalous ($y=1$)
		10 anomalous ($y=1$)
Alternative:	No test set	use if very few labeled anomalous examples
Training set:	6000 good engines	higher risk of overfitting
CV:	4000 good engines ($y=0$), 20 anomalous ($y=1$)	tune ϵ tune x_j

Algorithm evaluation

Algorithm evaluation

Fit model $p(x)$ on training set $x^{(1)}, x^{(2)}, \dots, x^{(m)}$
On a cross validation/test example x , predict

$$y = \begin{cases} 1 & \text{if } p(x) < \varepsilon \text{ (anomaly)} \\ 0 & \text{if } p(x) \geq \varepsilon \text{ (normal)} \end{cases}$$
10
2000

course 2 week
skewed data

Possible evaluation metrics:

- True positive, false positive, false negative, true negative
- Precision/Recall
- F_1 -score

Use cross validation set to choose parameter ε

Anomaly detection vs supervised learning

Fraud techniques changes a lot (every few months get a new type of fraud), but spam has not changed much.

Anomaly detection vs. Supervised learning

Very small number of positive examples ($y = 1$). (0-20 is common).

Large number of negative ($y = 0$) examples.

 
 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.

Fraud

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.

Spam

Anomaly detection

- Fraud detection
 - Manufacturing - Finding new previously unseen defects in manufacturing.(e.g. aircraft engines)
 - Monitoring machines in a data center
- ⋮

vs. Supervised learning

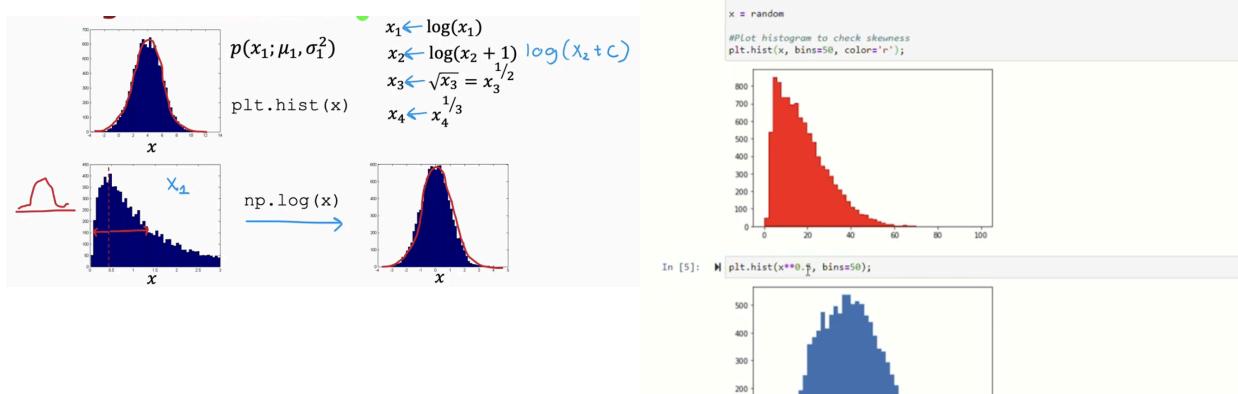
- Email spam classification
 - Manufacturing - Finding known, previously seen defects $y=1$ scratches
 - Weather prediction (sunny/rainy/etc.)
 - Diseases classification
- ⋮

Choosing what features to use

- Harder to find what features to ignore for an unsupervised algorithm

Non-gaussian features

- Try to make features as gaussian as possible
 - Look at them using plt.hist(x)



- Try things out. *Apply same transformation to training set and CV and test set*

Error analysis for anomaly detection

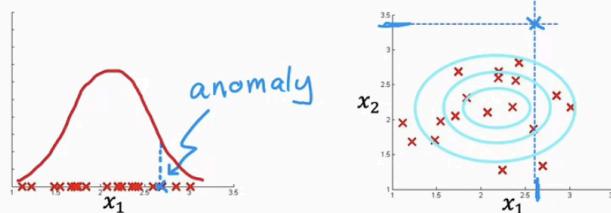
- May be useful to use a new feature
- Train a model, then see what anomalies in the CV set the algo has failed to detect. Then look at that examples to inspire new features that can flag these examples... (by giving unusually large or small values in the event of an anomaly)

Error analysis for anomaly detection

Want $p(x) \geq \epsilon$ large for normal examples x .
 $p(x) < \epsilon$ small for anomalous examples x .

Most common problem:

$p(x)$ is comparable for normal and anomalous examples.
 $(p(x)$ is large for both)



x_1 num transactions x_2 typing speed

Final example: monitoring computers in a data centre

Choose features that might take on unusually large or small values in the event of an anomaly.

$$\begin{aligned}x_1 &= \text{memory use of computer} \\x_2 &= \text{number of disk accesses/sec} \\x_3 &= \text{CPU load} \\x_4 &= \text{network traffic}\end{aligned}$$

high low

$$x_5 = \frac{\text{CPU load}}{\text{network traffic}}$$
$$x_6 = \frac{(\text{CPU load})^2}{\text{network traffic}}$$

not unusual

Deciding feature choice based on $p(x)$
Large for normal examples;
Becomes small for anomaly in the cross validation set

Recommender systems

Making Recommendations

Predicting movie ratings

User rates movies using one to five stars

Movie	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	?

$n_u = 4$ $r(1,1) = 1$
 $n_m = 5$ $r(3,1) = 0$ $y^{(3,2)} = 4$

Andrew Ng

Using per-item features

- Adding features of the movies

What if we have features of the movies?

Movie	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

$n_u = 4$
 $n_m = 5$
 $n = 2$

For user 1: Predict rating for movie i as: $w^{(1)} \cdot x^{(i)} + b^{(1)}$ just linear regression
 $w^{(1)} = \begin{bmatrix} 5 \\ 0 \end{bmatrix}$ $b^{(1)} = 0$ $x^{(3)} = \begin{bmatrix} 0.99 \\ 0 \end{bmatrix}$ $w^{(1)} \cdot x^{(3)} + b^{(1)} = 4.95$
For user j : Predict user j 's rating for movie i as $(w^{(j)}) \cdot x^{(i)} + b^{(j)}$

Andrew Ng

How can we make a cost function for this algorithm?

- Trying to minimise error between predicted rating and actual rating. Only summing over the movies that user j has actually rated...

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 :

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

Andrew Ng

- Training a different linear regression model for each of the $n(u)$ users

Collaborative filtering algorithm

- What if you don't have the values of features x_1 and x_2 ?
- Given the values of $w^T x$ that we have, what choice for $x(1)$ would cause these values to be right?...

Movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)	x_1 (romance)	x_2 (action)
	5	5	0	0	?	?
Love at last	5	5	0	0	?	?
Romance forever	5	?	?	0	?	$x^{(2)}$
Cute puppies of love	?	4	0	?	?	$x^{(3)}$
Nonstop car chases	0	0	5	4	?	?
Swords vs. karate	0	0	5	?	?	?

$w^{(1)} = \begin{bmatrix} 5 \\ 0 \end{bmatrix}$, $w^{(2)} = \begin{bmatrix} 5 \\ 0 \end{bmatrix}$, $w^{(3)} = \begin{bmatrix} 0 \\ 5 \end{bmatrix}$, $w^{(4)} = \begin{bmatrix} 0 \\ 5 \end{bmatrix}$ }
 $b^{(1)} = 0$, $b^{(2)} = 0$, $b^{(3)} = 0$, $b^{(4)} = 0$ ←
 using $w^{(j)} \cdot x^{(i)} + b^{(j)}$
 $w^{(1)} \cdot x^{(1)} \approx 5$
 $w^{(2)} \cdot x^{(1)} \approx 5$
 $w^{(3)} \cdot x^{(1)} \approx 0$
 $w^{(4)} \cdot x^{(1)} \approx 0$
 $\rightarrow x^{(1)} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$
 ↑

- Let's find a cost function for learning the values x_1 and x_2 (assuming you have the parameters for the different movies):

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

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

DeepLearning.AI Stanford ONLINE Andrew Ng

- Putting together to get collaborative filtering algorithm:

Collaborative filtering

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

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

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

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

Put them together:

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

DeepLearning.AI Stanford ONLINE Andrew Ng

So how do we minimise this cost function?

Option 1: Gradient descent

- Note that x is also a parameter here, so must also be minimised

collaborative filtering

Linear regression (course 1)

```

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

```

parameters w, b, x \times is also a parameter

Andrew Ng

Binary labels: favourites, likes and clicks

- Generalise the algorithm to this system (cf linear \rightarrow logistic regression)
- Predict how likely the people are to like the movies they have not yet rated \rightarrow Then can decide whether to recommend these movies or not

Movie	Binary labels				Example applications
	Alice(1)	Bob(2)	Carol(3)	Dave(4)	
Love at last	1	1	0	0	→ 1. Did user j purchase an item after being shown? 1, 0, ?
Romance forever	1	? ←	? ←	0	→ 2. Did user j fav/like an item? 1, 0, ?
Cute puppies of love	? ←	1	0	? ←	→ 3. Did user j spend at least 30sec with an item? 1, 0, ?
Nonstop car chases	0	0	1	1	→ 4. Did user j click on an item? 1, 0, ?
Swords vs. karate	0	0	1	? ←	

Meaning of ratings:
 → 1 - engaged after being shown item
 → 0 - did not engage after being shown item
 ? - 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 $g(z) = \frac{1}{1+e^{-z}}$

Adapting cost function for binary application:

Previous cost function:

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

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

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

cost for all examples

Recommender systems implementation detail

Mean Normalisation

Mean normalisation (i.e. normalise movie ratings to have a consistent average value) can speed up recommender system algorithms.

Movie	Alice(1)	Bob (2)	Carol (3)	Dave (4)	Eve (5)	
Love at last	5	5	0	0	?	0.5
Romance forever	5	?	?	0	?	2.5
Cute puppies of love	?	4	0	?	?	2.5
Nonstop car chases	0	0	5	4	?	2.5
Swords vs. karate	0	0	5	?	?	2.5

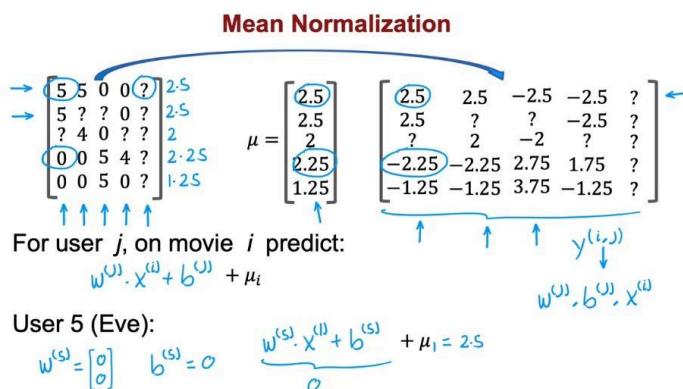
Users who have not rated any movies

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

$$w^{(s)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad b^{(s)} = 0 \quad w^{(s)} \cdot x^{(s)} + b^{(s)}$$

Computing average of rating of each movie:

- Subtracting mean from ratings. Can then year w,b,x for the modified dataset.
- Possible to also normalise columns; however, here rows make more sense. Normalising columns would work if it were a brand new movie no one had rated yet, rather than a new user



TensorFlow implementation of collaborative filtering

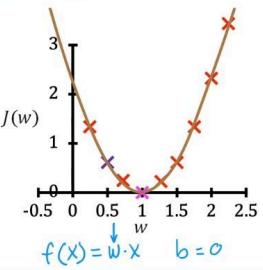
Derivatives in ML:

- Calculating update derivative for b can be complex - TensorFlow can help with this... automatically computing it
- Procedure below allows you how to implement gradient descent without having to know how to figure out derivative:

Gradient descent algorithm
Repeat until convergence

$$\underline{w} = w - \alpha \frac{\partial}{\partial w} J(w,b)$$

$$\underline{b} = b - \alpha \frac{\partial}{\partial b} J(w,b) \leftarrow b = 0$$



$J = (wx - 1)^2$

Gradient descent algorithm
Repeat until convergence

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

Fix $b = 0$ for this example

Custom Training Loop

```
w = tf.Variable(3.0)
x = 1.0
y = 1.0 # target value
alpha = 0.01

iterations = 30
for iter in range(iterations):
    # Use TensorFlow's GradientTape to record the steps
    # used to compute the cost J, to enable auto differentiation.
    with tf.GradientTape() as tape:
        fwb = w*x
        costJ = (fwb - y)**2

    # Use the gradient tape to calculate the gradients
    # of the cost with respect to the parameter w.
    [dJdw] = tape.gradient(costJ, [w])

    # Run one step of gradient descent by updating
    # the value of w to reduce the cost.
    w.assign_add(-alpha * dJdw)

# Tf.variables require special function to modify
```

Tf.variables are the parameters we want to optimize

Auto Diff

Auto Grad

$\frac{\partial}{\partial w} J(w)$

©DeepLearning.AI Stanford ONLINE Andrew Ng

Implementing collaborative filtering algorithm in TensorFlow (with adam):

- Using TensorFlow, you are not limited to gradient descent - can use a more powerful optimization algorithm such as adam

Gradient descent algorithm

Repeat until convergence

```
# Instantiate an optimizer.
optimizer = keras.optimizers.Adam(learning_rate=1e-1)

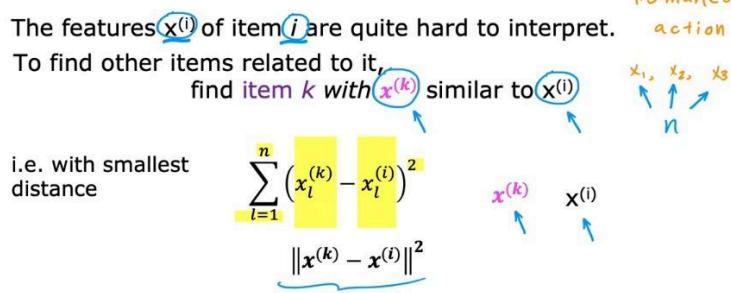
iterations = 200
for iter in range(iterations):
    # Use TensorFlow's GradientTape
    # to record the operations used to compute the cost
    with tf.GradientTape() as tape:
        # Compute the cost (forward pass is included in cost)
        cost_value = cofiCostFuncV(X, W, b, Ynorm, R,
                                    num_users, num_movies, lambda)
        n_u          n_m

    # Use the gradient tape to automatically retrieve
    # the gradients of the trainable variables with respect to
    # the loss
    grads = tape.gradient(cost_value, [X,W,b])

    # Run one step of gradient descent by updating
    # the value of the variables to minimize the loss.
    optimizer.apply_gradients(zip(grads, [X,W,b]))
```

Dataset credit: Harper and Konstan. 2015. The MovieLens Datasets: History and Context

Finding related items



Limitations of Collaborative Filtering

→ Cold start problem. How to

- • rank new items that few users have rated?
- • show something reasonable to new users who have rated few items?

→ Use side information about items or users:

- • Item: Genre, movie stars, studio,
- • User: Demographics (age, gender, location), expressed }

Collaborative filtering vs Content-based filtering

→ Collaborative filtering:

Recommend items to you based on ratings of users who gave similar ratings as you

→ Content-based filtering:

Recommend items to you based on features of user and item to find good match

$r(i,j) = 1$ if user j has rated item i

$y^{(i,j)}$ rating given by user j on item i (if defined)

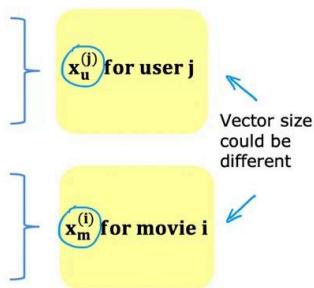
Examples of user and item features:

User features:

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

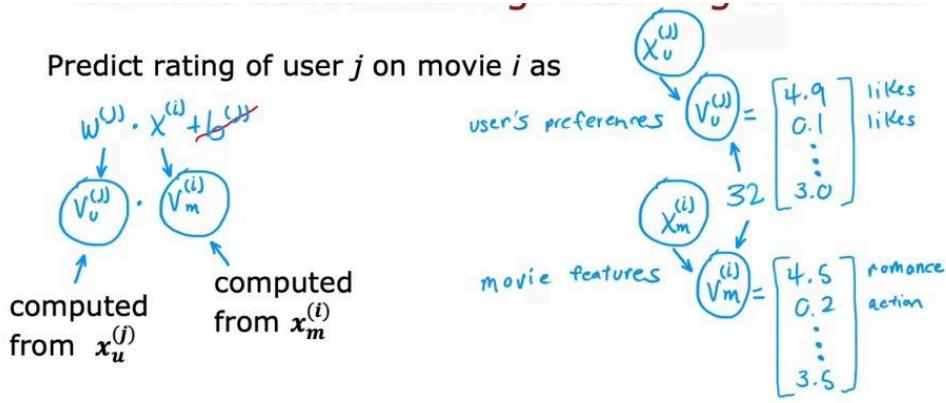
Movie features:

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



So how do we match user and movie in content-based filtering:

- Given features of a user $x(j)u$, how can we create the vector for the user's preferences vector and similarly, given features of a movie, how to compute movie features vector
- Xs can be different in size, vs must be same size so you can take a dot product

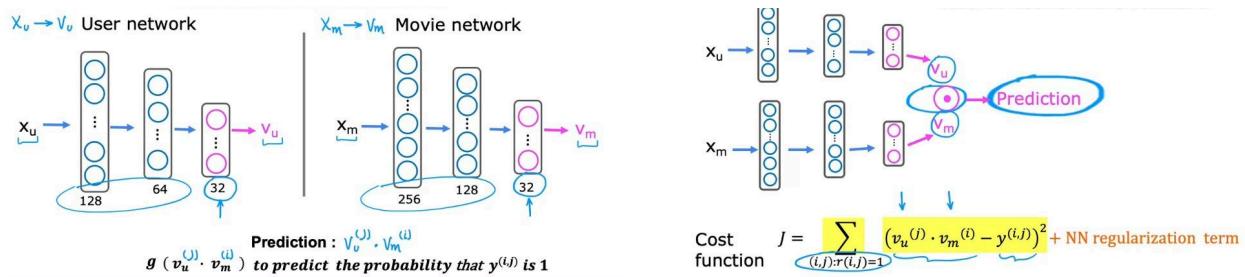


Summarising:

- In collaborative filtering, we have number of users give ratings of different items
- Vs in content-based filtering, we have features of users and features or items and want to find good matches by computing vectors for both and taking dot products
- But how to compute these vectors?...

Deep learning for content-based filtering

- Neural network architecture



Train parameters of neural network so result in vectors of users and movies to minimise square error.

- Regularisation encourages NN to keep parameters of NN small

After training this model, can also use it to find similar items... (can do this ahead of time)

- i.e. where the distance between two movie vectors is small

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 :

$$\|v_m^{(i)} - v_m^{(j)}\|^2 \text{ small}$$

$$\|x^{(i)} - x^{(j)}\|^2$$

Note: This can be pre-computed ahead of time

Note that this is reflective of NN benefits over decision trees - easier to take multiple neural networks and put them together to work in tandem.

We now consider practical issues of scaling up to very large item catalogues...

Recommending from a large catalogue

How can we efficiently find recommendation from a large set of items?

- Computationally unfeasible to run neural network influence huge numbers of times every time a user comes on website
- Instead, use two steps: retrieval and ranking...

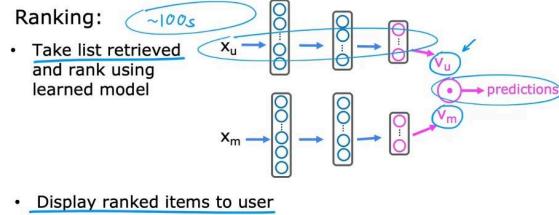
Two steps: Retrieval & Ranking

Retrieval:

- Generate large list of plausible item candidates e.g.
 ↗ 1) For each of the last 10 movies watched by the user, find 10 most similar movies
 $\|v_m^{(i)} - v_m^{(j)}\|^2$
 ↗ 2) For most viewed 3 genres, find the top 10 movies
 ↗ 3) Top 20 movies in the country
- Combine retrieved items into list, removing duplicates and items already watched/purchased

Two steps: Retrieval & ranking

Ranking:

- Take list retrieved and rank using learned model
 - Display ranked items to user
- 

- One additional optimization is that if you have computed VM for all the movies in advance, then all you need to do is to do inference on this part of the neural network a single time to compute VU, and then take that VU they just computed

for the user on your website right now and take the inner product between VU and VM for the movies that you have retrieved during the retrieval step.

Trade off:

- • Retrieving more items results in better performance, but slower recommendations.
- • To analyse/optimize the trade-off, carry out offline experiments to see if retrieving additional items results in more relevant recommendations (i.e., $p(y^{(i,j)}) = 1$ of items displayed to user are higher).

100 500

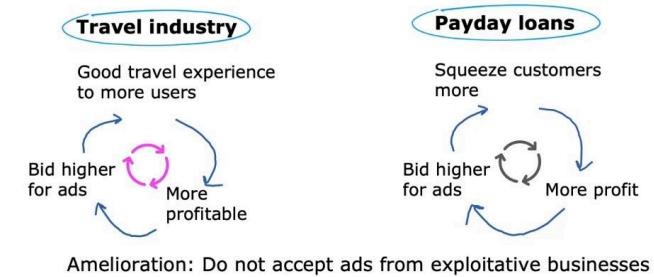
Ethical use of recommender systems

What is the goal of the recommender system?

Recommend:

- • Movies most likely to be rated 5 stars by user
- • Products most likely to be purchased
- • Ads most likely to be clicked on *+high bid*
- • Products generating the largest profit
- • Video leading to maximum watch time

Ethical considerations with recommender systems:

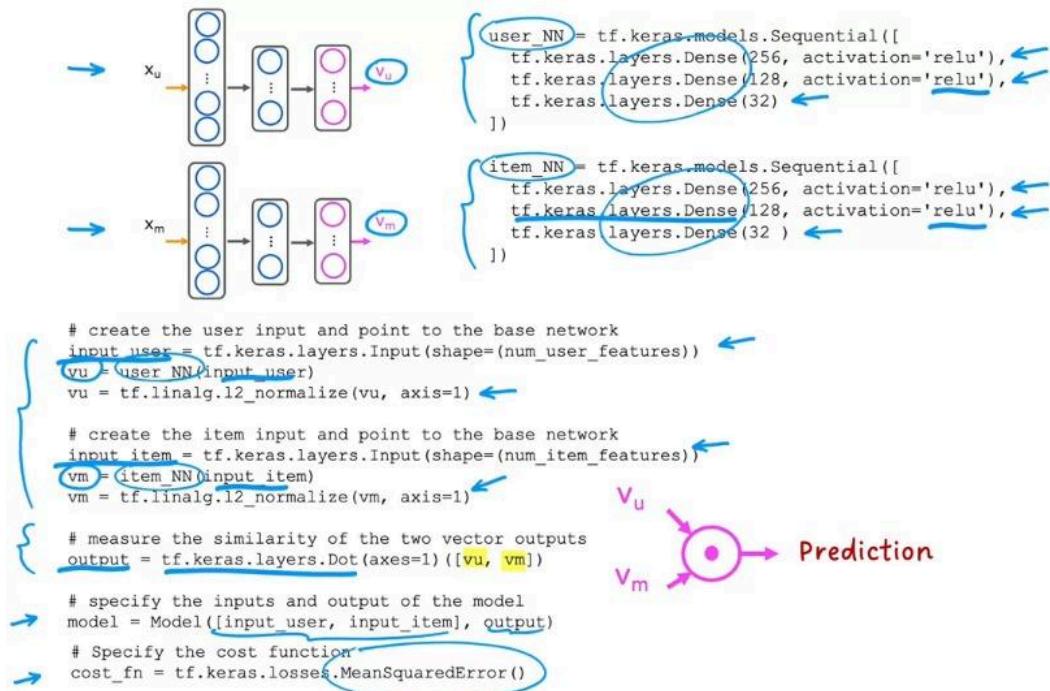


Other problematic cases:

- • Maximizing user engagement (e.g. watch time) has led to large social media/video sharing sites to amplify conspiracy theories and hate/toxicity
- • Amelioration : Filter out problematic content such as hate speech, fraud, scams and violent content
- • Can a ranking system maximize your profit rather than users' welfare be presented in a transparent way?
- • Amelioration : Be transparent with users

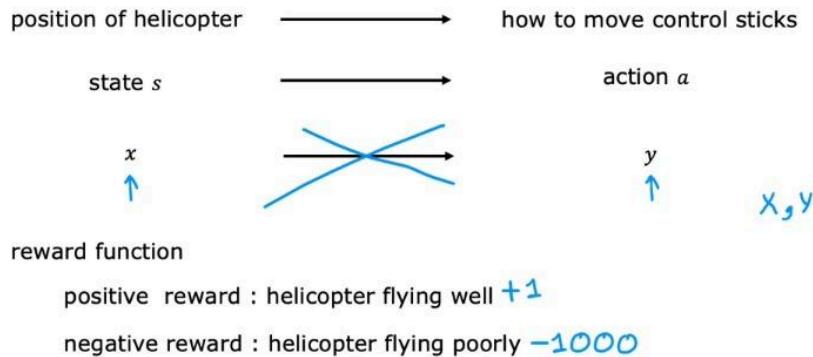
TensorFlow implementation of content-based filtering:

- axis=1 normalises vector to length 1
- keras.layers.Dot just takes dot product between vu and vm to give output
- Can also L2 normalise vector



Week 3: Reinforcement Learning

- Could use supervised learning to learn states and then output actions; however, hard to identify an exact state
 - Reinforcement learning - use a reward function

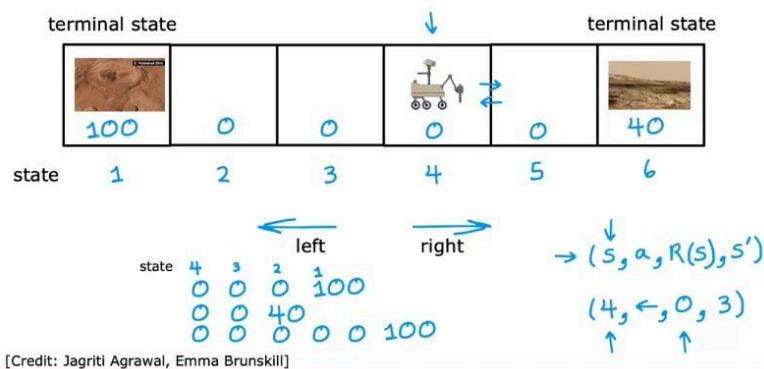


Applications:

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

The Mars rover example

- Rover can be any of six positions (states)
 - In this example, state 1 has a more interesting sample than 6 but is further away. To incentivise, give a higher reward at state 1.
 - At each step, rover can either go to left or go to the right.
 - 1 and 6 are terminal states - game ends/ nothing more happens after
 - $S = \text{state}$, $a = \text{action}$, $R(s) = \text{reward}$, $s' = \text{next state}$



The Return in reinforcement learning

How do you know if a particular set of rewards is better than another?

- i.e. rewards that you can get quicker may be better to take
- *Return* is the sum of rewards, with a discount factor applied (normally about 0.9, 0.99 etc - i.e. close to one)

	0	0		0	
state 1	2	3	4	5	6

$$\text{Return} = 0 + (0.9)0 + (0.9)^2 0 + (0.9)^3 100 = 0.729 \times 100 = 72.9$$

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

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

$$\gamma = 0.5$$

$$\text{Return} = 0 + (0.5)0 + (0.5)^2 0 + (0.5)^3 100$$

More concrete examples of returns:

- Return depends on the rewards, and the rewards depend on the actions you take
 - so the return depends on the actions you take
- Starting at different states and only going in one direction...

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

← return

$$\gamma = 0.5$$

← reward

The return depends on the actions you take.

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

$$0 + (0.5)0 + (0.5)^2 40 = 10$$

100	50	25	12.5	20	40
100	0	0	0	0	40

$$0 + (0.5)40 = 20$$

With negative rewards, this discount factor incentivises system to push out negative rewards as far into the future as possible.

Making decisions: Policies in reinforcement learning

Policies:

- Could decide to always go to a nearer/larger/smaller one etc.
- **A policy maps a state to an action**

Policy

$s \xrightarrow{\text{policy} \ \pi} a$	<table border="1" style="border-collapse: collapse; text-align: center;"> <tr><td>100</td><td><</td><td><</td><td>></td><td>></td><td>40</td></tr> <tr><td>100</td><td><</td><td><</td><td><</td><td><</td><td>40</td></tr> <tr><td>100</td><td>></td><td>></td><td>></td><td>></td><td>40</td></tr> <tr><td>100</td><td><</td><td><</td><td><</td><td>></td><td>40</td></tr> </table>	100	<	<	>	>	40	100	<	<	<	<	40	100	>	>	>	>	40	100	<	<	<	>	40
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 goal of reinforcement learning

100					40
-----	--	--	---	--	----

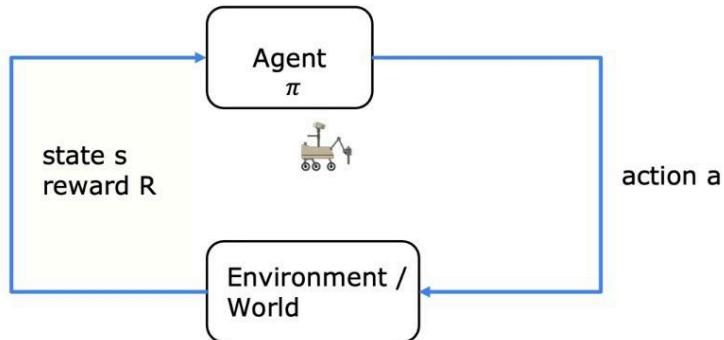
Find a policy π that tells you what action ($a = \pi(s)$) to take in every state (s) so as to maximize the return.

Summary

	Mars rover 	Helicopter 	Chess 						
states	6 states	position of helicopter	pieces on board						
actions	$\leftarrow \rightarrow$	how to move control stick	possible move						
rewards	$100, 0, 40$	$+1, -1000$	$+1, 0, -1$						
discount factor γ	0.5	0.99	0.995						
return	$R_1 + \gamma R_2 + \gamma^2 R_3 + \dots$	$R_1 + \gamma R_2 + \gamma^2 R_3 + \dots$	$R_1 + \gamma R_2 + \gamma^2 R_3 + \dots$						
policy π	<table border="1" style="border-collapse: collapse; text-align: center;"><tr><td>100</td><td><</td><td><</td><td><</td><td>></td><td>40</td></tr></table>	100	<	<	<	>	40	$\text{Find } \pi(s) = a$	$\text{Find } \pi(s) = a$
100	<	<	<	>	40				

Markov Decision Process (MDP):

- Markov - future only depends on current state, not on how you got here
- Can think of it as: you have an agent, choose action a , based on action a something will happen in world, we then see what state we in and what reward R we get



State-action value function, Q

$Q(s,a)$ - function of state and action

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

$Q(s,a)$
 $\uparrow \uparrow$

100	50	25	12.5	20	40
100	0	0	0	0	40

100	50	25	12.5	20	40
100	0	0	0	0	40

← return
 ← action
 ← reward

$$Q(2, \rightarrow) = 12.5 \\ 0 + (0.5)0 + (0.5)^2 0 + (0.5)^3 100 \\ Q(2, \leftarrow) = 50 \\ 0 + (0.5)100 \\ Q(4, \leftarrow) = 12.5 \\ 0 + (0.5)0 + (0.5)^2 0 + (0.5)^3 100$$

Q function gives you a way to pick actions too...

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

→	100	50	25	12.5	20	40
→	100	0	0	0	0	40

← return
 ← action
 ← reward

$\max_a Q(s,a)$
 $\pi(s) = a$

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

The best possible return from state s is $\max_a Q(s,a)$.

The best possible action in state s is the action a that gives $\max_a Q(s,a)$. Q^* Optimal Q function

Example of state-action value function:

- Discount factor is a measure of patience - how beneficial it is to wait
- See how values of $Q(s,a)$ change and how optimal policy changes

State Action Value Function Example

In this Jupyter notebook, you can modify the mars rover example to see how the values of $Q(s,a)$ will change depending on the rewards and discount factor changing.

```
In [1]: import numpy as np
from utils import *

In [2]: # Do not modify
num_states = 6
num_actions = 2

In [3]: terminal_left_reward = 100
terminal_right_reward = 40
each_step_reward = 0

# Discount factor
gamma = 0.5

# Probability of going in the wrong direction
misstep_prob = 0

In [4]: generate_visualization(terminal_left_reward, terminal_right_reward, each_step_reward, gamma, misstep_prob)
```

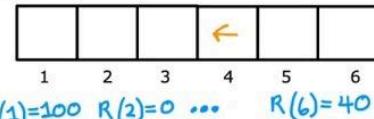


Bellman equation:

- Allows you to compute $Q(s,a)...$

$Q(s,a) = \text{Return if you}$

- start in state s .
- take action a (once).
- then behave optimally after that.



s : current state]

$R(s)$ = reward of current state

a : current action]

s' : state you get to after taking action a]

a' : action that you take in state s']

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

Looking at an example:

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

$$Q(s,a) = R(s)$$

100	100	50	12.5	25	6.25	12.5	10	6.25	20	40	40
100	0	0	0	0	0	0	0	0	0	40	40

$s=2$
 $a=\rightarrow$
 $s'=3$

$$Q(2, \rightarrow) = R(2) + 0.5 \max_{a'} Q(3, a')$$

$$= 0 + (0.5)25 = 12.5$$

$s=4$
 $a=\leftarrow$
 $s'=3$

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

$$= 0 + (0.5)25 = 12.5$$

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

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

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

$s \rightarrow s'$

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

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

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

$$Q(s, a) = R_1 + \gamma R_2 + \gamma^2 R_3 + \gamma^3 R_4 + \dots$$

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

To relate this back to the earlier example...

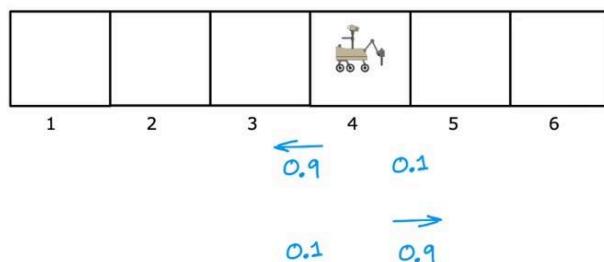
$$Q(s, a) = R(s) + \gamma \max_{a'} Q(s', a')$$

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

$Q(4, \leftarrow)$
 $= 0 + (0.5)0 + (0.5)^2 0 + (0.5)^3 100$
 $= R(4) + (0.5)(0 + (0.5)0 + (0.5)^2 100)$
 $= R(4) + (0.5) \max_{a'} Q(3, a')$

Random (stochastic) environment

- In some applications, when you take an action, outcome is not completely reliable (e.g. wind blowing or wheel slipping with robots etc.)...
- E.g. if you command it to go left, 10% of the time it might slip and go in the other direction, and vice versa



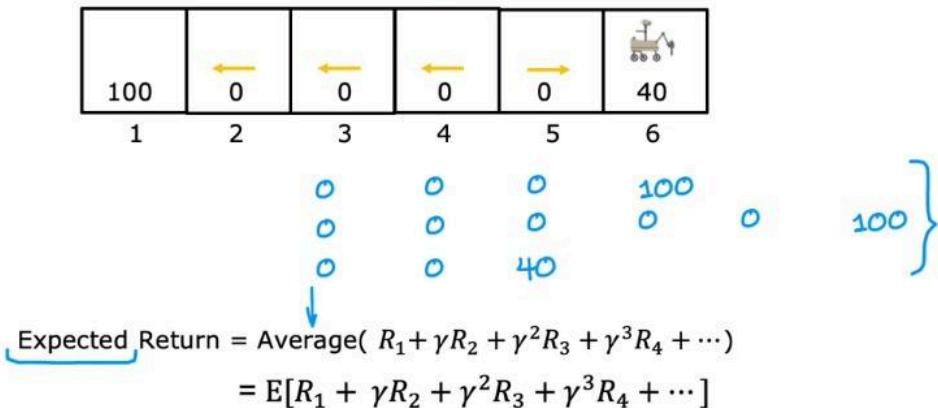
Expected Return:

- E.g. use policy where you go left in states 4,3,2 and right in state 5
- E.g. 1 states you visit will be a little random due to these percentages (0.9, 0.1) - might start in state 4, get a bit lucky and go to state 3, then to state 2, and to state 1 (giving 0 0 0 100)
- E.g. 2 but might get less lucky - tell it to go left, and works; but then tell it to go right, and it slips and go right, and then goes left for the rest (0 0 0 0 100)
- E.g. 3 start in four, say go left and it slips and goes right, then command it to go right (0 0 40)

We had previously written out return as sum of discounted rewards. But when the reinforcement learning problem is stochastic, there isn't one sequence of rewards you see for sure, but lots of different ones. Thus, in a stochastic learning problem, we don't maximise return (as this is a random number) - instead we want to maximise the average value of the sum of discounted rewards (i.e. if we take policy and try it out thousands of time, get lots of different sequences and then take average to get expected (average) return.)

(E = expected value)

Learning algorithm aims to choose a policy π to maximise the average of the expected sum of discounted rewards.



Here: s' is random, so we put an average E:

Goal of Reinforcement Learning:

Choose a policy $\pi(s) = a$ that will tell us what action a to take in state s so as to maximize the expected return.

Bellman
Equation:

$$Q(s, a) = \underbrace{R(s)}_{3 \leftarrow} + \gamma \underbrace{E[\max_{a'} Q(s', a')]}_{2 \text{ or } 4}$$

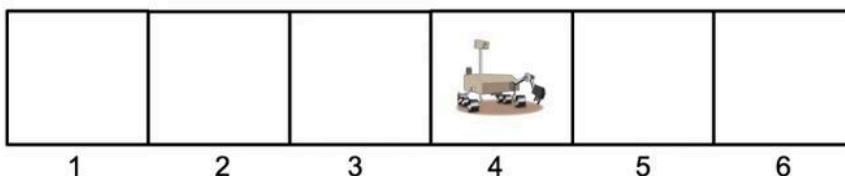
- In the lab, this probability of slipping etc. is called the “misstep_prob”

Example of continuous state space applications

- How to generalise to a continuous state instead of six discrete states (like the mars rover) e.g. if it could be anywhere along a line between 0 and 6km
- E.g. building a self-driving a car (xy position, orientation (angle theta), speed in a direction, and how quickly it is turning) - i.e. becomes a vector:

Discrete vs Continuous State

Discrete State:



Continuous State:



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

E.g.2 Position of autonomous helicopter

X (how far north or south), y (how far east or west), z (height above ground) + orientation (roll (rowing left to right), pitch (up or back), yaw (which compass direction it is facing))

Position (x,y,z); Roll (phi), pitch (theta), yaw (omega); speed (in x, y, z directions), rate of turning (how fast is roll, pitch, and yaw changing)

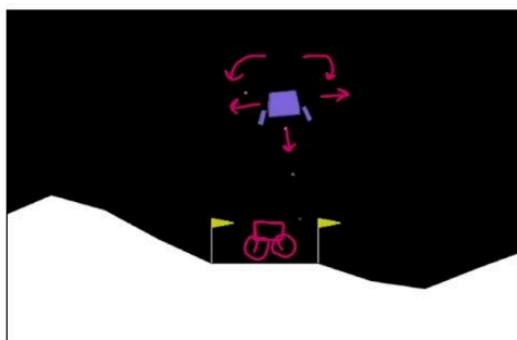
- This set of 12 numbers is input into a policy, which decides what the appropriate next action is

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

Thus, in a continuous state MDP, the state of the problem is a vector of continuous numbers.

Lunar Lander

Four possible actions:



actions:

do nothing	x
left thruster	y
main thruster	\dot{x}
right thruster	\dot{y}

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

0 or 1 {

- In vector, theta is tilt to left or right (angular velocity); l and r correspond to whether left or right leg is sitting on ground (binary values)

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

So overall problem is:

Learn a policy π that, given

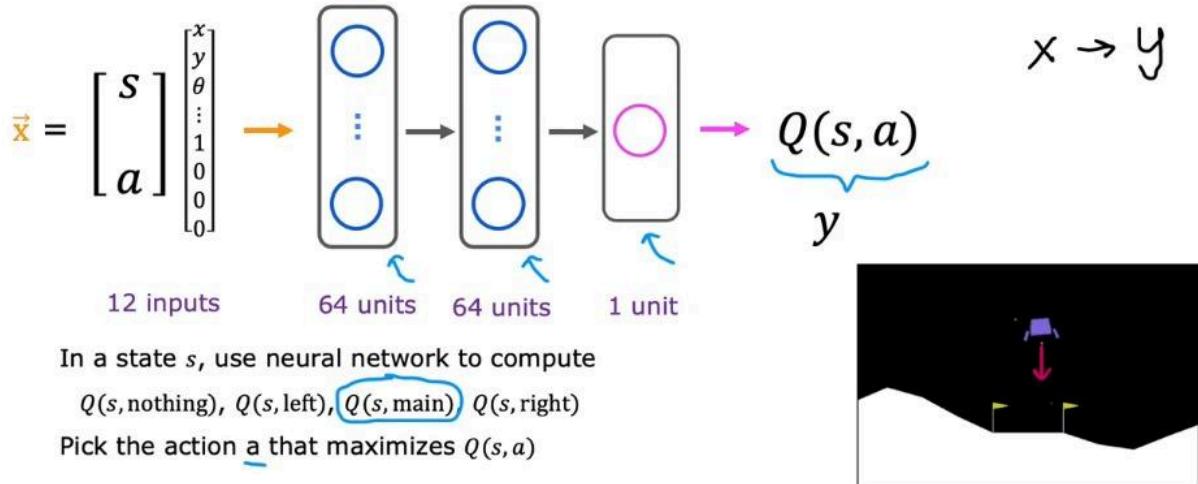
$$s = \begin{bmatrix} x \\ y \\ \dot{x} \\ \dot{y} \\ \theta \\ \dot{\theta} \\ l \\ r \end{bmatrix}$$

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

$$\gamma = 0.985$$

Deep Reinforcement Learning: Learning the state-value function

Train a neural network to approximate $Q(s,a)$, which will in turn allow us to pick good actions.



Encode any of four actions using a one-hot feature vector (e.g. for “do nothing”: encode with 1 0 0 0)

- Input a state-action pair and use it to output $Q(s,a)$ for each of the four actions, and pick the action a that maximises $Q(s,a)$

So, how do you train this neural network?

- Use Bellman's equation to create a training set with lots of examples x and y , and then use supervised learning to map from state-action pair to $Q(s,a)$ (i.e. accurately predict $x \rightarrow y$)

So how do you get the training set with values of x and y ?

- Take random actions to start with, to see what rewards you get. Each tuple can represent a single training example
- X s are inputs with 12 features, y s are just numbers, then train neural network with the mean squared error loss to try to predict y as a function of the input x

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

$f_{w, b}(x) \approx y$

$(s, a, R(s), s')$

$x^{(1)} = (s^{(1)}, a^{(1)})$

$y^{(1)} = R(s^{(1)}) + \gamma \max_{a'} Q(s'^{(1)}, a')$

$x^{(2)} = (s^{(2)}, a^{(2)})$

$y^{(2)} = R(s^{(2)}) + \gamma \max_{a'} Q(s'^{(2)}, a')$

$x^{(3)} = (s^{(3)}, a^{(3)})$

$y^{(3)} = R(s^{(3)}) + \gamma \max_{a'} Q(s'^{(3)}, a')$

$x^{(1)} = (s^{(1)}, a^{(1)})$

$y^{(1)} = R(s^{(1)}) + \gamma \max_{a'} Q(s'^{(1)}, a')$

$x^{(2)} = (s^{(2)}, a^{(2)})$

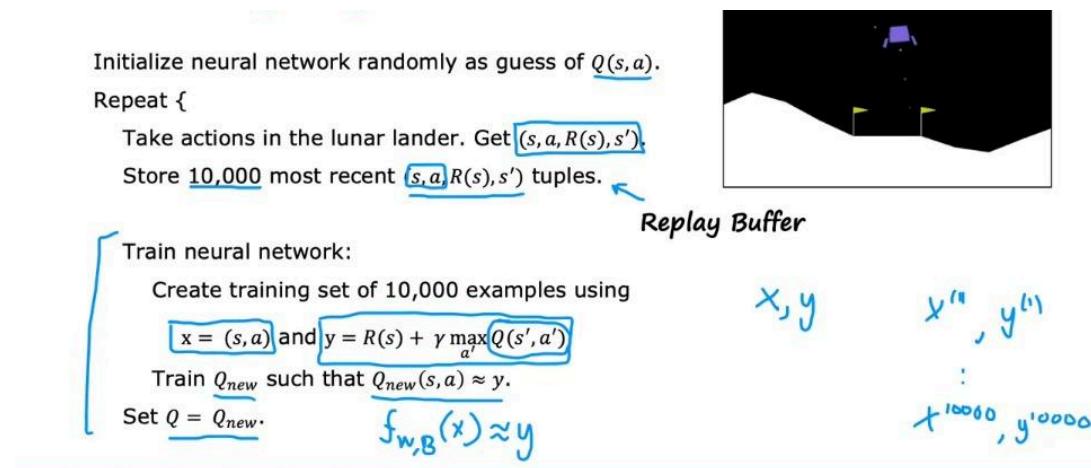
$y^{(2)} = R(s^{(2)}) + \gamma \max_{a'} Q(s'^{(2)}, a')$

$x^{(3)} = (s^{(3)}, a^{(3)})$

$y^{(3)} = R(s^{(3)}) + \gamma \max_{a'} Q(s'^{(3)}, a')$

Where does $Q(s', a')$ come from? Can take a random guess of the actual Q function, even though we do not know.

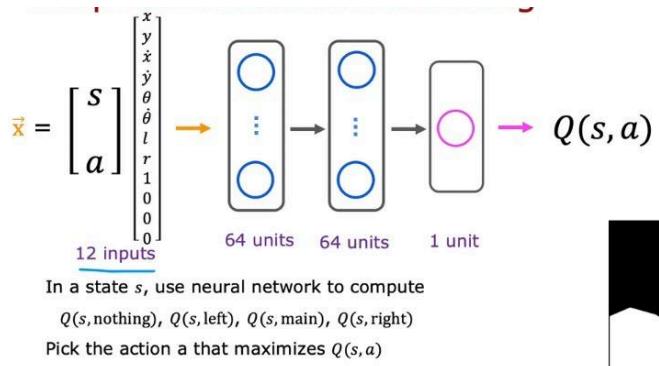
Full Learning Algorithm (DQN algorithm)



Starting a random guess of Q, using Bellman's equations to repeatedly improve guess of Q, slightly improves guess, so for next model you train will have a slightly better estimate of Q, which then updates Q.

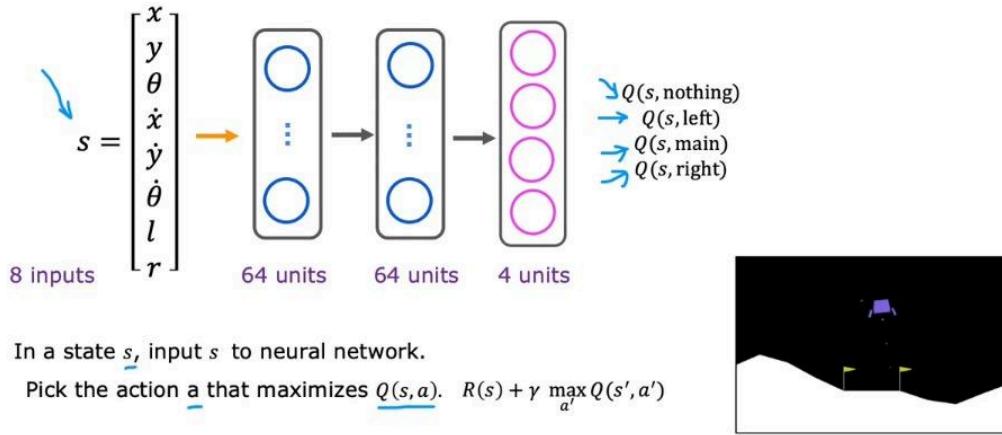
Algorithm refinement: Improved neural network architecture

- Possible change to architecture to make algo much more efficient...



Here, whenever we are in some state s , have to carry out inferences in the NN separately four times to compute four values of Q (nothing, right etc.), to pick action a that maximises $Q(s, a)$.

- More efficient to output all four simultaneously



This NN also makes it much more efficient to compute $R(s) + \dots$, because we are getting $Q(s', a')$ for all actions at the same time, so just pick the max to compute this value.

Algorithm refinement: epsilon-greedy policy

To pick actions while you're still learning (i.e. don't have a good estimate of $Q(s,a)$ yet), use the epsilon-greedy policy...

Initialize neural network randomly as guess of $Q(s, a)$.

Repeat {

Take actions in the lunar lander. Get $(s, a, R(s), s')$.

Store 10,000 most recent $(s, a, R(s), s')$ tuples.

Train model:

Create training set of 10,000 examples using

$$x = (s, a) \text{ and } y = R(s) + \gamma \max_a Q(s', a').$$

Train ϱ_{new} such that $\varrho_{new}(s, a) \approx y$, $f_{w,h}(x) \approx y$

Set $\theta = \theta_{new}$

How to choose actions while still learning? Some options... (Option 2 is best)

- Occasionally pick an action randomly. Suppose that $Q(s,a)$ was initialised randomly so that $Q(s,\text{main})$ is always low - then NN will never try to fire main thruster, and thus will never learn that firing main thruster is a good idea. Randomness allows it to try out new actions, so NN can overcome its own possible preconceptions - “exploration step”.
 - Exploration vs exploitation trade-off

- Trick: start with epsilon high, and then gradually decrease so that over time you are more likely to use improving estimate of Q to decide on actions instead of random actions

In some state s

Option 1:

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

$Q(s, \text{main})$ is low



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

Algorithm refinement: Mini-batch and soft updates

Can speed up supervised learning algorithms as well.

Mini-batches

- Problem with gradient descent (below) when dataset is very big is that every step of grad descent requires computing sum of this average over 10000000 examples. Then have to do it again after a tiny grad descent step.
- Mini-batch gradient descent does not use all training examples on every iteration, but instead use a smaller number m' (a subset) - much more efficient

How to choose actions while still learning?

x	y
2104	400
1416	232
1534	315
852	178
...	...
3210	870

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2$$

$m = 100,000,000$
 $m' = 1,000$

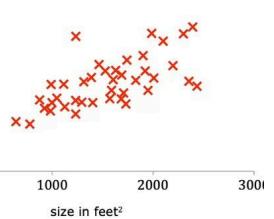
repeat {
 $w = w - \alpha \frac{\partial}{\partial w} \left[\frac{1}{2m'} \sum_{i=1}^{m'} (f_{w,b}(x^{(i)}) - y^{(i)})^2 \right]$
 $b = b - \alpha \frac{\partial}{\partial b} \left[\frac{1}{2m'} \sum_{i=1}^{m'} (f_{w,b}(x^{(i)}) - y^{(i)})^2 \right]$
}

x	y
2104	400
1416	232
1534	315
852	178
...	...
3210	870

Mini-batch

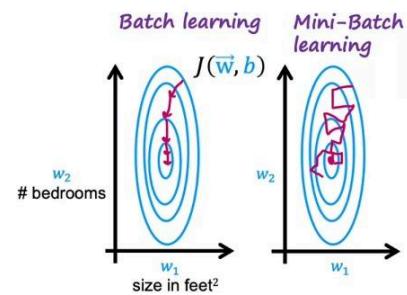
price in \$1000's

size in feet²



Recall *batch gradient descent*: every step causes parameters to reliably get closer to global minimum. In contrast, *mini-batch* will take a more convoluted route towards global minimum, but each iteration is much less computationally expensive.

- For a large dataset, mini batch (with other optimisation algos like adam, for example) is much more common than batch grad descent



Updated previous learning algorithm

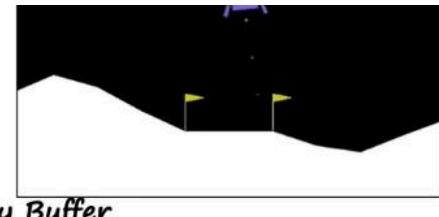
- Just take a subset of the examples of the tuples to train a few examples of the NN - accelerates it (even if it makes it a bit more noisy)
- Other refinement: Set $Q = Q_{\text{new}}$ can make an abrupt change to Q (as can be overwritten by a worse algorithm... use Soft Update to remedy this:)

Initialize neural network randomly as guess of $Q(s, a)$

Repeat {

Take actions in the lunar lander. Get $(s, a, R(s), s')$.

Store 10,000 most recent $(s, a, R(s), s')$ tuples.



Replay Buffer

Train model:

1,000

Create training set of 10,000 examples using

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

Train Q_{new} such that $Q_{\text{new}}(s, a) \approx y$.

Set $Q = Q_{\text{new}}$.

$$\begin{aligned} &x^{(1)}, y^{(1)} \\ &\vdots \\ &x^{(1000)}, y^{(1000)} \end{aligned}$$

Soft Update

- Soft update only adds a small amount of W_{new} (the new value) - controls how aggressively you move towards W_{new} and B_{new} from W and B (soft update allows a more gradual change)
- Allows more reliable convergence

$$\text{Set } Q = Q_{\text{new}} \leftarrow Q(s, a)$$

W, B $W_{\text{new}}, B_{\text{new}}$

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

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

$$W = 1 W_{\text{new}} + 0 W$$

- NB mini batching also applies well to supervised learning applications

Limitations of Reinforcement Learning

- Much easier to get to work in a simulation than a real robot!
- Far fewer applications than supervised and unsupervised learning.
- But ... exciting research direction with potential for future applications.