

Clustering

Clustering looks at a dataset and tries to find points that are similar to each other and group them together.

The data has no labels, only $x^{(i)}$, no $y^{(i)}$. The algorithm has to figure out the structure on its own.

Some applications:

- * Grouping similar news article
- * Market Segmentation
- * DNA microarray data
- * Astronomy

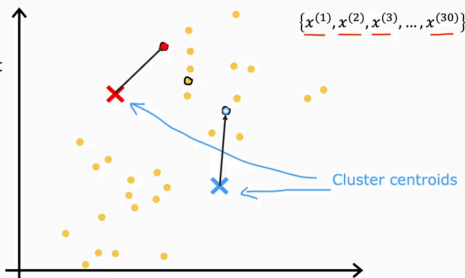
k-means intuition

k-means repeatedly does two things:

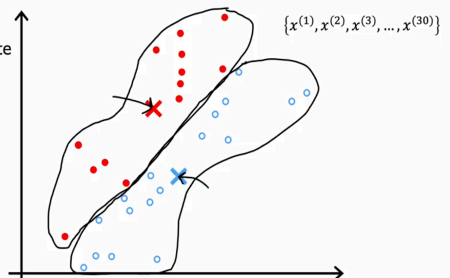
- * Assign each point to the closest cluster centroid
- * Move each centroid to the average of the points assigned to it.

We start with k randomly placed centroids. Color each point by which centroid is nearest. Then we move each centroid to the mean of its colored points. We then recolor and repeat. The centroids stop moving once assignments stabilize.

Step 1:
Assign each point to its closest centroid



Step 2:
Recompute the centroids



K-means Algorithm

We use the following notation:

- * k is the number of clusters
- * $x^{(i)} \in \mathbb{R}^n$
- * μ_k is the position of the k centroid

We first randomly initialize $\mu_1, \mu_2, \dots, \mu_k$

Repeat:

Assign points to cluster centroids

for $i = 1$ to m :

$c^{(i)} :=$ index k (from 1 to K) that minimizes $\|x^{(i)} - \mu_k\|^2$

Move cluster centroids

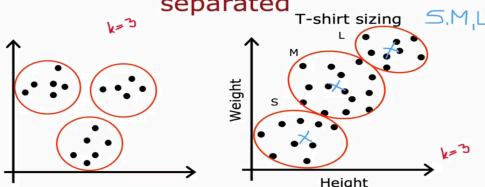
for $k = 1$ to K :

$\mu_k :=$ average (mean) of points assigned to cluster k

Note: In practice, the squared distance $\|x^{(i)} - \mu_k\|^2$ is minimized, not the distance because it is cheaper to work with the squared distance in terms of compute and they have the same minimum.

Corner case: If a centroid ends up with zero points assigned, the most common move is to delete it, ending with $k-1$ clusters. If we really need k , we reinitialized the centroid.

K-means for clusters that are not well separated



Optimization Objective.

Some notations

- * $c^{(i)}$ is the index of cluster $(1, 2, \dots, K)$ to which $x^{(i)}$ is currently assigned
- * μ_k is the cluster centroid k
- * $\mu_{c^{(i)}}$ is the cluster centroid of cluster to which example $x^{(i)}$ has been assigned.

K-means is minimizing a specific cost function, called the distortion function:

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

The average squared distance from each training example to the centroid it has been assigned to.

The two steps of k-means corresponds to minimizing J over different variables:

- * The assignment step minimizes J over $c^{(i)}$ with μ_k held fixed. Closest centroid is the best choice.
- * The move step minimizes J over μ_k with $c^{(i)}$ held fixed. The mean is the point that minimizes the sum of squared distance to a set of points.

So J goes down (or stays flat) every iteration. It should never go up. If implementation shows J increasing, there is a bug!

Convergence occurs when J stops decreasing or decreases very slowly.

Note: K-means converges to a local minimum, not necessarily the global one.

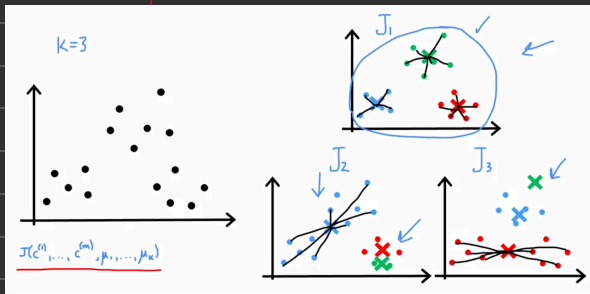
Initializing k-means

Random initialization!

- * Choose $k < m$
- * Randomly pick k distinct training examples.
- * Set $\mu_1, \mu_2, \dots, \mu_k$ equal to those k examples.

The clustering we end up with depends on initialization. Bad initial centroids give bad local minima of J . The fix is to run k-means multiple times with different random inits. Compute J at the end of each run. Keep the run with the lowest J .

50 to 1000 random restarts is typical. More than that gives diminishing returns.

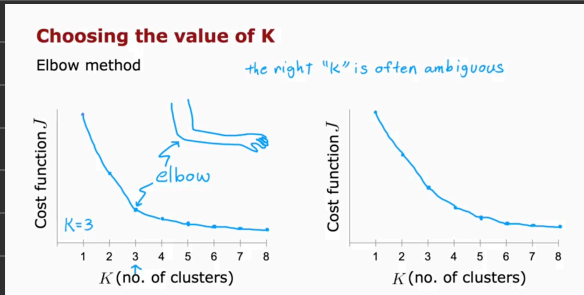


Random initialization

For $i = 1$ to 100 {
Randomly initialize K-means. $\leftarrow k$ random examples
Run K-means. Get $c^{(1)}, \dots, c^{(m)}, \mu_1, \mu_1, \dots, \mu_k \leftarrow$
Compute cost function (distortion)
 $J(c^{(1)}, \dots, c^{(m)}, \mu_1, \mu_1, \dots, \mu_k) \leftarrow$
}
Pick set of clusters that gave lowest cost (J)

Choosing the number of clusters k

Elbow method: plot J as a function of k . We look for a kink or elbow where J stops dropping sharply. The problem is that the curve is usually smooth with no clear elbow.



A wrong technique is choosing k to minimize J . J always decreases as k grows (more centroids means each point can sit closer to one). Following this rule always pick as large as possible which is bad.

The right approach is to pick k based on what clusters are for downstream.

Anomaly Detection.

Anomaly detection looks at unlabeled data of normal examples, learns what normal looks like, then flags new examples that look different.

The technique used to do anomaly detection is called density estimation.

Density estimation means building a model of the probability distribution that generated the training dataset. Given the dataset $\{x^{(1)}, \dots, x^{(m)}\}$, you fit a model $p(x)$ that tells us, for any point x , how likely it is to come from the same distribution as the training dataset.

After fitting $p(x)$, anomaly detection becomes simple:

- * Compute $p(x_{\text{test}})$ for a new example
- * If $p(x_{\text{test}}) < \epsilon$, flag it as anomalous
- * If $p(x_{\text{test}}) \geq \epsilon$, treat it as normal.

ϵ is a small threshold we have to choose.

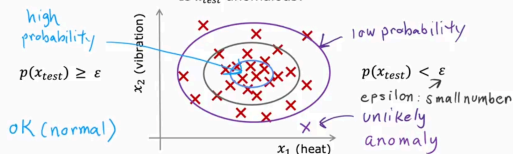
Visually: regions where the training data is dense have high $p(x)$. Regions far from any training example have low $p(x)$. The anomalies sit in the low-density regions.

Density estimation

Dataset: $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$ probability of x being seen in dataset

Model $p(x)$

Is x_{test} anomalous?



Gaussian (normal) distribution

To do density estimation, we need a model for $p(x)$. We use Gaussian here.

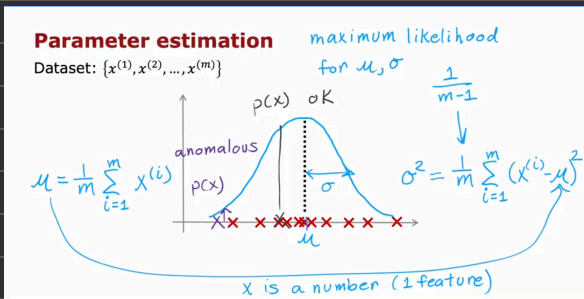
A random variable x follows a Gaussian distribution with mean μ and variance σ^2 if:

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\left(\frac{x-\mu}{\sigma}\right)^2}$$

- * Changing μ shifts the curve left or right.
- * Smaller σ makes it taller or narrower.
- * Larger σ makes it shorter and wider.
- * Total area under curve is always 1.

Estimating the parameters from data (maximum likelihood):

$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}, \quad \sigma^2 = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)^2$$



Anomaly Detection Algorithm:

Each training example has n features: $x^{(i)} \in \mathbb{R}^n$

Model $p(x)$ by assuming the features are statistically independent and that each feature is Gaussian:

$$p(x) = P(x_1; \mu_1, \sigma_1^2) \cdot P(x_2; \mu_2, \sigma_2^2) \cdots P(x_n; \mu_n, \sigma_n^2) \approx \prod_{j=1}^n P(x_j; \mu_j, \sigma_j^2)$$

Algorithm:

* Choose n features x_j we think indicate anomalous examples

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

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}, \quad \sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$$

* Given a new example x , compute:

$$p(x) = \prod_{j=1}^n \frac{1}{\sqrt{2\pi}\sigma_j} e^{-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}}$$

Flag as anomalous if $p(x) < \epsilon$

Intuition: $p(x)$ is large only when every x_j sits near its typical range. A single weird feature drags the whole product close to zero.

Anomaly detection algorithm

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

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

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}, \quad \sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$$

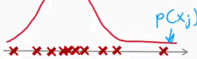
3. Given new example x , compute $p(x)$:

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

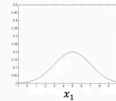
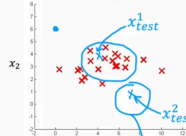
Anomaly if $p(x) < \epsilon$

Vectorized formula

$$\bar{\mu} = \frac{1}{m} \sum_{i=1}^m \bar{x}^{(i)}, \quad \bar{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \dots \\ \mu_n \end{bmatrix}$$

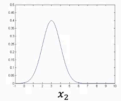


Anomaly detection example



$$\mu_1 = 5, \sigma_1 = 2$$

$$p(x_1; \mu_1, \sigma_1^2)$$



$$\mu_2 = 3, \sigma_2 = 1$$

$$p(x_2; \mu_2, \sigma_2^2)$$

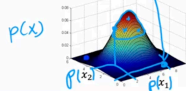
$\epsilon = 0.02$

$$p(x_{test}^{(1)}) = 0.0426$$

→ "ok"

$$p(x_{test}^{(2)}) = 0.0021$$

→ anomaly



Developing and evaluating an anomaly detection system

To tune the system, we need need a real-number evaluation metric. This requires some anomalous examples (even if very few) to validate against.

Setup: assume we have some normal examples ($y=0$) and a small number of known anomalies ($y=1$)

Training set: $x^{(1)}, x^{(2)}, \dots, x^{(m)}$ (assume normal examples)

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

test set: $(x_{test}^{(1)}, y_{test}^{(1)}), \dots, (x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$ } anomalous examples mostly normal.

Aircraft engines monitoring example

10000 good (normal) engines 2 to 50
20 flawed engines (anomalous) $y=1$

Training set: 6000 good engines $y=0$ train algorithm on training set
2 flawed engines $y=1$

CV: 2000 good engines ($y=0$) 10 anomalous ($y=1$)
use cross validation set tune ϵ tune x_j

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

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 ϵ tune x_j

Algorithm evaluation

course 2 weeks
skewed datasets

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

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

10
2000

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

If we have some labeled anomalies, why not just train a supervised classifier? it depends on the problem.

Anomaly detection is the right choice when:

- * Very small number of positive (anomalous) examples, say 0 to 20.
- * Large number of negative (normal) examples
- * Many different "types" of anomalies, hard for an algorithm to learn what they look like from the positive alone.
- * Future anomalies may look nothing like the ones we have seen so far.

Use cases: fraud detection, manufacturing defects (especially new types), monitoring machines in a data center.

Supervised learning is the right choice when:

- * Large numbers of positive and negative examples.
- * Enough positive for the algorithm to learn what they look like.
- * Future positives are likely to look similar to training positives.

Use cases: Spam classification, weather prediction, disease classification.

Choosing what features to use

Feature choice matters more in anomaly detection than in supervised learning because supervised learning can downweight bad features using the labels. Anomaly detection has no labels at training time, so it has to trust whatever features we give it.

We can plot a histogram of each feature. If it looks bell-shaped, we leave it. If it is skewed, we transform it until it looks more Gaussian. Common transformations!

* $\log(x)$

* $\log(x+c)$ for some small constant c (handles zeros)

* $x^{1/2}$

* $x^{1/3}$

It is important to apply whatever transformation we apply to training set to cv and test set.

Recommender Systems

Collaborative filtering

Making recommendations

A recommender system predicts how a user would rate or respond to items they have not seen yet, then surfacel with the highest predicted score.

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	?

Ratings				
★				
★	★			
★	★	★		
★	★	★	★	
★	★	★	★	★

→ n_u = no. of users

→ n_m = no. of movies

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

$n_u = 4$ $r(1,1) = 1$ $y^{(i,j)}$ = rating given by user j to movie i
 (defined only if $r(i,j)=1$)

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

Collaborative filtering algorithm

What if we don't have features for romance and action in the movie example? But if we already know users' parameters $w^{(i)}, b^{(i)}$ and their actual rating, we can work backwards to learn the features $x^{(i)}$

Given $w^{(1)}, \dots, w^{(n_u)}$ and $b^{(1)}, \dots, b^{(n_u)}$, we learn $x^{(i)}$ by minimizing:

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

For all movies at once

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

Now the trick: we don't have w 's and b 's either and we don't have x 's. But if we fix w, b we can solve for x and if we fix x we can solve for w, b . So we just learn both simultaneously by writing one combined cost function:

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

Minimizing with gradient descent over $w^{(j)}, b^{(j)}, x^{(i)}$ all together. The features $x^{(i)}$ are now learned parameters, not given inputs. We pick n (number of features) ourselves.

The name "collaborative filtering" comes from this, many users collaborating on rating movies effectively teach the algorithm what feature each movie has, which then helps predict ratings for other users.

$$w_k^{(i)} := w_k^{(i)} - \alpha \frac{\partial J}{\partial w_k^{(i)}} \quad b^{(i)} := b^{(i)} - \alpha \frac{\partial J}{\partial b^{(i)}} \quad x_k^{(i)} := x_k^{(i)} - \alpha \frac{\partial J}{\partial x_k^{(i)}}$$

Mean normalization

Introducing a new user who has rated nothing causes a problem. The cost function only has the regularization terms left for that user, which pushes all $w^{(i)}$ to zero. Predicted ratings become $0 + b^{(i)} = 0$ for every movie.

We fix this using mean normalization. We build the rating matrix Y where row i is movie i and column j is user j . For each movie i , we compute the average rating μ_i across users who actually rated it. Subtract μ_i from every observed rating in that row to get Y_{norm} .

We train collaboratively on Y_{norm} , learning $W^{(i)}$, $b^{(i)}$, $x^{(i)}$ as before. When predicting, we add the mean back:

$$\hat{y}^{(i,j)} = W^{(i)} x^{(j)} + b^{(i)} + \mu_i$$

For a brand new user with no ratings, $w^{(i)} = 0$ and $b^{(i)} = 0$, so the prediction collapses to μ_i .

Mean Normalization

5	0	0	2	2.5	
5	?	?	0	?	2.5
?	4	0	?	?	2
0	0	5	4	?	2.25
0	0	5	0	?	1.25

$\mu =$

2.5
2.5
2
2.25
1.25

2.5	2.5	-2.5	-2.5	?
2.5	?	?	-2.5	?
?	2	-2	?	?
-2.25	-2.25	2.75	1.75	?
-1.25	-1.25	3.75	-1.25	?

For user j , on movie i predict:
 $w^{(i)} \cdot x^{(j)} + b^{(i)} + \mu_i$

User 5 (Eve):
 $w^{(5)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$ $b^{(5)} = 0$ $\frac{w^{(5)} \cdot x^{(5)} + b^{(5)}}{0} + \mu_1 = 2.5$

Tensorflow Implementation of Collaborative filtering.

Collaborative filtering's gradients are non-trivial to derive by hand. Tensorflow's auto-diff (also called auto-grad) computes them automatically.

$J = \frac{1}{2} \sum_i (w^T x_i - y_i)^2$

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 are the parameters we want to optimize

```

w = tf.Variable(...)
x = tf.Variable(...)
b = tf.Variable(...)

optimizer =
tf.keras.optimizers.Adam(learning_rate=1e-1)

for iter in range(iterations):
    with tf.GradientTape() as tape:
        cost = cofi_cost_func(X, W, b, Y, R, lambda_)
        grads = tape.gradient(cost, [X, W, b])
        optimizer.apply_gradients(zip(grads, [X, W, b]))
                
```

Implementation in TensorFlow

Gradient descent algorithm
Repeat until convergence

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

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

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

```

# 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 = cofiCostFunc(X, W, b, Ynorm, R,
                                  num_users, num_movies, lambda)

        # Use the GradientTape 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 minimise 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

After training, the learned features $x^{(i)}$ characterizes each item. They are not human readable (the algorithm did not learn romance and action, it learned something), but items with similar x vectors play similar roles in user preferences.

To find similar items to i : compute the squared distance

$$\|x^{(k)} - x^{(i)}\|^2 = \sum_{l=1}^n (x_l^{(k)} - x_l^{(i)})^2$$

for every other item k we sort ascending, return the top few. That gives "more like this" recommendations.

Limitations of collaborative filtering

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 preferences, ...

Content-based filtering

Two ways to build recommender:

- * Collaborative filtering recommends item i to user j by looking at what other users with similar ratings to user j thought of item i . It uses the matrix only. It does not look at who the users are or what the items are.
- * Content-based filtering recommends item i to user j based on features of user j and features of item i , asking whether they are good match. It uses side information about users and items, not just past ratings.

Notation stays the same: $r(i,j) = 1$ if user j rated item i , $y^{(i,j)}$ is that rating.

Content-based filtering needs feature vectors for both users and items.

User features $x_u^{(j)}$ for user j might include:

- * Age (often binned or normalized)
- * Gender (one-hot, or three-way to allow unknown)
- * Country (one-hot)
- * Movies watched
- * Average rating given per genre

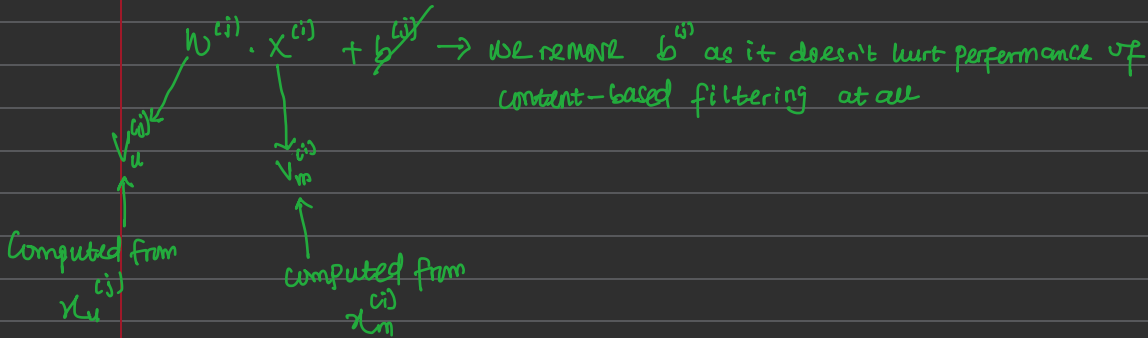
Movie features $x_m^{(i)}$ for movie i might include:

- * Year
- * Genre (one-hot or multi-hot, a movie can be multiple genres)
- * Reviews / Critic scores
- * Average rating across all users
- * Cast

The two vectors need not have the same dimension.

The Model

Predict rating of user j on movie i as:



So we have:

$$\hat{y}^{(ij)} = v_u^{(ij)} \cdot v_m^{(ij)}$$

where $v_u^{(ij)}$ is a vector computed from $x_u^{(ij)}$ and $v_m^{(ij)}$ is a vector computed from $x_m^{(ij)}$. Both v 's must have the same dimension (otherwise we cannot dot-product them)

Deep learning for content based filtering

How do we compute v_u from x_u and v_m from x_m ? we will use two neural networks. One for users, one for items.

User network: takes $x_u^{(ij)}$ (say 1500 numbers), passes it through hidden layers, outputs $v_u^{(ij)}$ (say 32 numbers)

Movie network: takes $x_m^{(ij)}$ (say 50 numbers), passes it through hidden layers, outputs $v_m^{(ij)}$ (say 32 numbers)

The two networks are independent (different weights, possibly different architectures), but their output layers must have the same number of units, because prediction is:

$$\hat{y}^{(ij)} = v_u^{(ij)} \cdot v_m^{(ij)}$$

For binary labels, we wrap top dot product in a sigmoid:

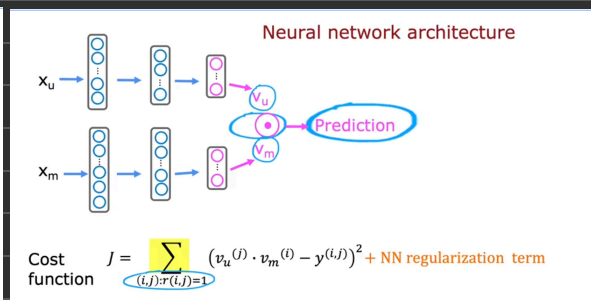
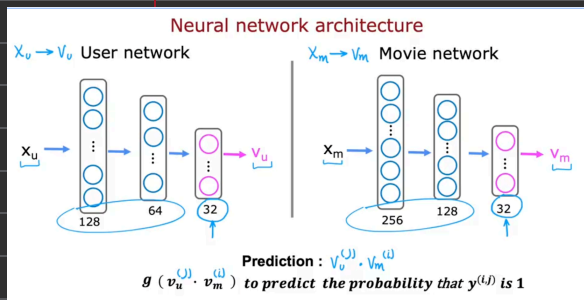
$$\hat{y}^{(i,j)} = g(v_u^{(i)} \cdot v_m^{(j)}), \quad g(z) = \frac{1}{1+e^{-z}}$$

Training the model

Both networks are trained together as a single combined model.

The cost function for the regression version is:

$$J = \sum_{(i,j): r(i,j)=1} (v_u^{(i)} \cdot v_m^{(j)} - y^{(i,j)})^2 + \text{NN regularization terms.}$$



Finding similar items

After training, $v_m^{(i)}$ describes the character of items i in the learned space. Items with similar v_m vectors are similar to each other.

To find items similar to item i :

$$\|v_m^{(k)} - v_m^{(i)}\|^2 = \sum_{L=1}^m (v_m^{(k),L} - v_m^{(i),L})^2$$

for every other items k , sort ascending and return the top few.

Note that fees can be precomputed. We do not need to wait until a user shows up to compute similar items. Build the list ahead of time and serve it from cache.

Recommending from a large catalogue.

For a service with millions of items (YouTube videos, products, news articles), running the full neural network on every item every time a user shows up is too slow. The standard fix is a two-stage architecture.

Retrieval: Generate a large list of plausible candidates, several hundred to a few thousand, using fast and approximate methods. For the movie example:

- * For each of the last 10 movies the user watched, fetch the 10 most similar movies (precomputed using V_m distances)
- * Top movies in the user's three most watched genres
- * Top movies in the user's country

We combine these lists, deduplicate, remove items the user has already seen. End up with maybe 100 to 1000 candidates.

Ranking: Take that shortlist, run the full content-based model on every candidate, sort by predicting rating $\hat{y}^{(i,j)}$, return the top k to the user.

The point of two stages is that retrieval is cheap and inexact, ranking is expensive and accurate. We only pay the ranking cost on a small set.

Trade-off in retrieval is that more candidates means better final recommendations but slower response. We tune the size on our held-out data, measuring offline metrics like recall@k or directly through A/B tests.

Tensorflow Implementation

```
user_NN = tf.keras.models.Sequential([
    tf.keras.layers.Dense(256, activation='relu'),
    tf.keras.layers.Dense(128, activation='relu'),
    tf.keras.layers.Dense(32)
])

item_NN = tf.keras.models.Sequential([
    tf.keras.layers.Dense(256, activation='relu'),
    tf.keras.layers.Dense(128, activation='relu'),
    tf.keras.layers.Dense(32)
])

# Inputs
input_user =
tf.keras.layers.Input(shape=(num_user_features,))
vu = user_NN(input_user)
vu = tf.linalg.l2_normalize(vu, axis=1)

input_item =
tf.keras.layers.Input(shape=(num_item_features,))
vm = item_NN(input_item)
vm = tf.linalg.l2_normalize(vm, axis=1)

# Dot product
output = tf.keras.layers.Dot(axes=1)([vu, vm])

model = tf.keras.Model([input_user, input_item],
output)
model.compile(optimizer=tf.keras.optimizers.Adam(0.01), loss=tf.keras.losses.MeanSquaredError())
```

Principal Component Analysis

PCA allows us to take data with a lot of features say hundreds or thousands of features and reduce it to very small number of features, say 2 or 5 features. It is commonly used to visualize data to figure out what might be going on in the data.

PCA Algorithm

Step 1: Normalize features. PCA is sensitive to feature scales because it works with variance. Standardize: Subtract the mean of each feature, divide by its standard deviation. Each feature now has mean 0 and variance 1.

Step 2: Pick the first principal component. The first PC is the direction (unit vector) along which the projected data has maximum variance.

Step 3: Pick subsequent components. The second principal component is the direction of maximum variance among direction perpendicular to the first. The third is perpendicular to the first two. And so on. The PCs form an orthogonal basis.

Step 4: Project data. To reduce n features to k features, keep the first k principal components and project each data point onto them. The projection of point x onto a unit vector u is:

$$z = x \cdot u$$

Step 5 (optional): Reconstruct. We can approximately recover the original point from its low dimensional representation.

$$x_{\text{approx}} = z \cdot u$$

PCA in code

```
from sklearn.decomposition import PCA

# Step 1: optionally normalize features first
# (StandardScaler if features have very different
scales)

# Step 2: fit
pca = PCA(n_components=2)
pca.fit(X)

# Step 3: how much variance does each component
explain?
pca.explained_variance_ratio_
# e.g. array([0.992, 0.008]) means PC1 explains
99.2% of variance

# Step 4: project to lower dimension
X_reduced = pca.transform(X)

# Step 5: reconstruct back to original space (lossy)
X_reconstructed = pca.inverse_transform(X_reduced)
```

PCA is mostly used for data visualization nowadays.
Used historically, less so now!

* Data compression

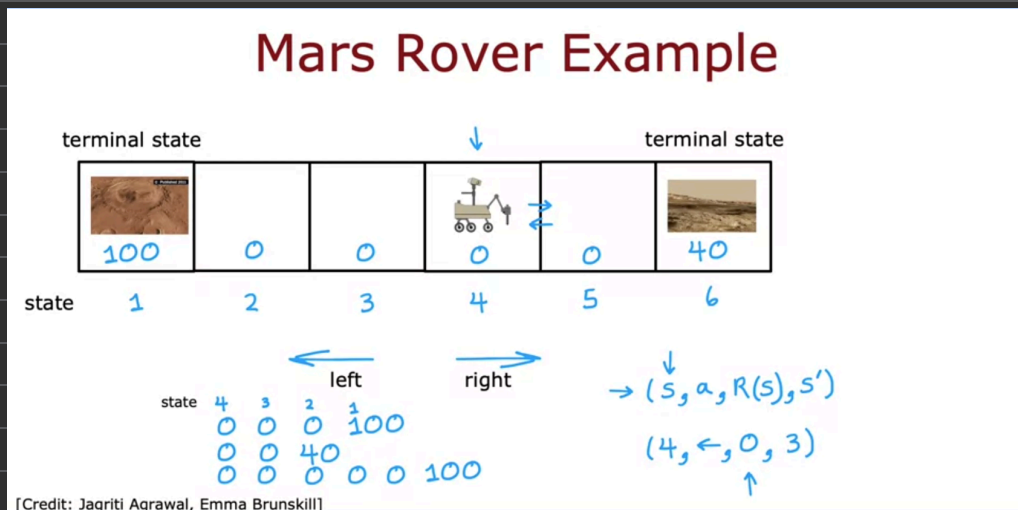
* Speeding up supervised learning: Reduce features before training a model. Largely superseded by modern algorithms that handle high-dimensional data fine! And by feature engineering, regularization, deep learning that learns its own representations

Reinforcement Learning.

Reinforcement learning (RL) trains an agent to take actions in an environment. The agent does not get told what action is correct. It gets a reward signal instead: a number that says how well things are going.

We can compare this with supervised learning which needs (x, y) pairs where y is the correct output. For tasks like flying helicopters or running a Mars rover, nobody can write down what the correct action at this state is for every state. But it is easy to score the outcome. Helicopter flying smoothly = good (positive reward like 10), helicopter crashed = bad (large negative reward). RL takes that score and figures actions on its own.

Mars Rover Example



The Return in RL

The return is the sum of rewards the agent collects from a given starting point, but with future rewards discounted.

$$G = R_1 + \gamma R_2 + \gamma^2 R_3 + \gamma^3 R_4 + \dots$$

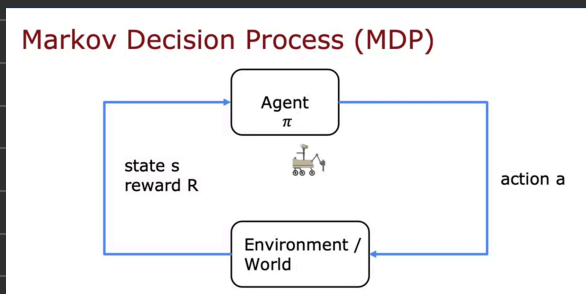
Where R_t is the reward received at step t and $\gamma \in [0, 1)$ is the discount factor. Reasons for the discount is that the infinite sum converges and a reward is worth more than the same reward in 100 steps. In financial application, discount factor is the interest rate.

Making decisions: Policies in RL

A policy π is a function that says, given a state s , what actions to take: $\pi(s) = a$.

So the goal of RL is to find a policy π that maximizes expected return.

This framework is called a Markov Decision Process (MDP). Markov because the next state depends only on the current state and the action, not the history.



State-action Value Function (Q-function)

The state-action value function, written $Q(s, a)$, is the return we get if we:

- * Start in state S
- * Take action a once
- * Then behave optimally from then on.

Sometimes called the Q -function or optimal Q -function (Q^*)

State action value function (Q-function)

$Q(s, a) =$ Return if you

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

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

← return

← action

← reward

$Q(s, a)$

↑ ↑

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

$Q(2, \leftarrow) = 50$
 $0 + (0.5)100$

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

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

If we know $Q(s, a)$ for every state and action, we would know the optimal policy. We just pick the highest Q

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

Large γ makes it patient and wants to take more steps to get the best reward while small γ makes it impatient and wants to take fewer steps to get to the terminal state even if it gets less reward.

Continuous State Spaces

Our Mars rover example has 6 discrete states, real problems have continuous states.

A self-driving car's state vector might be:

$$s = (x, y, \theta, \dot{x}, \dot{y}, \dot{\theta})$$

x and y position, heading θ , and their three velocities - six continuous numbers.

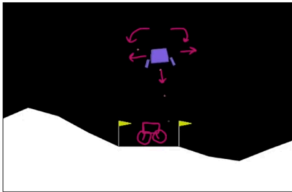
Lunar lander has 8-dimensional continuous state:

$$s = (x, y, \dot{x}, \dot{y}, \theta, \dot{\theta}, l, r)$$

x and y position, x and y velocity, angle, angular velocity, and two binary flags for whether left and right legs are touching the ground. We cannot store $Q(s, a)$ as a table of state-action pairs when s is continuous. There are infinitely many steps. We approximate $Q(s, a)$ with a neural network instead.

Lunar Lander

Lunar Lander



actions:

- do nothing
- left thruster
- main thruster
- right thruster

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

0 or 1

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



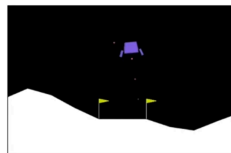
Lunar Lander Problem

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.

$$\gamma = 0.985$$



Deep Q-Network (DQN)

We want to train a neural network to compute or approximate the state-action value function $Q(s, a)$ and that in turn will let us pick good actions.

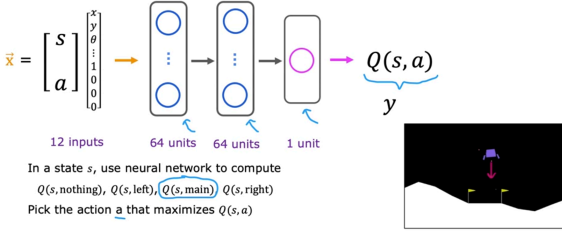
Input: state s (7 numbers for lunar lander) plus action a (one-hot encoded as 4 numbers), so input dimension is 12.

Output: scalar predicting $Q(s, a)$.

Architecture:

- * Input layer: 12 units
- * Hidden layer 1: 64 units, ReLU
- * Hidden layer 2: 64 units, ReLU
- * Output layer: 1 unit, linear (no activation).

Deep Reinforcement Learning



Bellman Equation

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

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

$y^{(1)} = R(s^{(1)}) + \gamma \max_{a'} Q(s'^{(1)}, a')$
 $y^{(2)} = R(s^{(2)}) + \gamma \max_{a'} Q(s'^{(2)}, a')$

$x^{(1)} = (s^{(1)}, a^{(1)})$
 $x^{(2)} = (s^{(2)}, a^{(2)})$
 $x^{(3)} = (s^{(3)}, a^{(3)})$

$y = (10,000)$, $y = (10,000)$

Training Q-network:

Training data is built by interacting with the environment.

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

Replay Buffer

Train neural network:

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

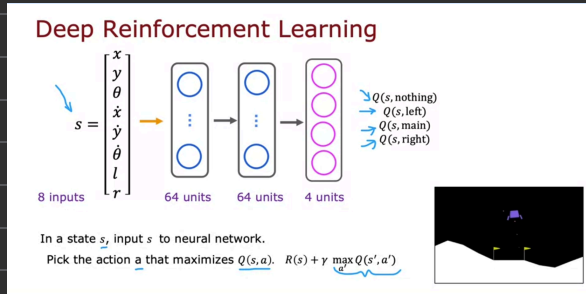
$$f_{w, \theta}(x) \approx y$$



$$x, y \quad x^{(1)}, y^{(1)} \\ \vdots \\ x^{(1000)}, y^{(1000)}$$

Algorithm refinement: Improved neural network architecture.

For the previous architecture, we have to carry out inference for each of the four actions, this is inefficient as we have to do inference four times for each actions, we can just train a single neural network to output all four actions at once.



ϵ -greedy policy

How does the agent decide which action to take while collecting experience?

- * Pure exploitation: always pick $\arg \max_a Q(s|a)$, the problem is that if the network is wrong about some action being bad, the agent never tries it and never learns.
- * Pure exploration: random actions. The agent never uses what it has learned.

ϵ -greedy balances the two:

- * With probability $1 - \epsilon$, pick $\arg \max_a Q(s|a)$ (greedy/exploit)
- * With probability ϵ , pick a random action (explore)

We typically start ϵ high (say 1.0) and decay it over training to a small value (say 0.01). Early on, we mostly explore. Later, mostly exploit.

RL algorithms are more finicky than supervised learning. Hyperparameter tuning matters more and a bad ϵ -schedule or learning rate can kill training.

Mini-batch and soft update.

We can use two more refinements:

* Mini-batch gradient descent: Instead of computing the gradient over all training data, we use a small random subset (say 64 examples) per step. Faster per step, noisier per step, usually converges faster overall. This is standard for any neural network training.

* Soft update: Naively, the target $y = R + \gamma \max_{a'} Q(s', a')$ use the same network that is being trained, which makes the targets move every step. This causes instability. The fix is to use two networks: A Q-network that gets trained and a target Q-network used to compute the targets y . The target network is updated slowly:

$$W_{\text{target}} := \tau W + (1 - \tau) W_{\text{target}}$$

With small τ (e.g. 0.001). The target network drifts slowly toward the Q-network, stabilizing the target. Without this, training often diverges.

Soft Update

Set $Q = Q_{\text{new}}$. $\leftarrow Q(s, a)$
 \uparrow \uparrow
 W, B $W_{\text{new}}, B_{\text{new}}$

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

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

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

State of RL

RL is great at simulated environments, games, robotics in controlled settings.

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.