

Unsupervised Learning, Recommenders, Reinforcement Learning

Clustering is an **unsupervised machine learning** technique used to group similar data points together **without using labeled data**.

Definition:

Clustering involves organizing a dataset into **clusters** (groups) such that:

- Data points in the **same cluster** are **more similar** to each other.
 - Data points in **different clusters** are **more different**.
-

Goal:

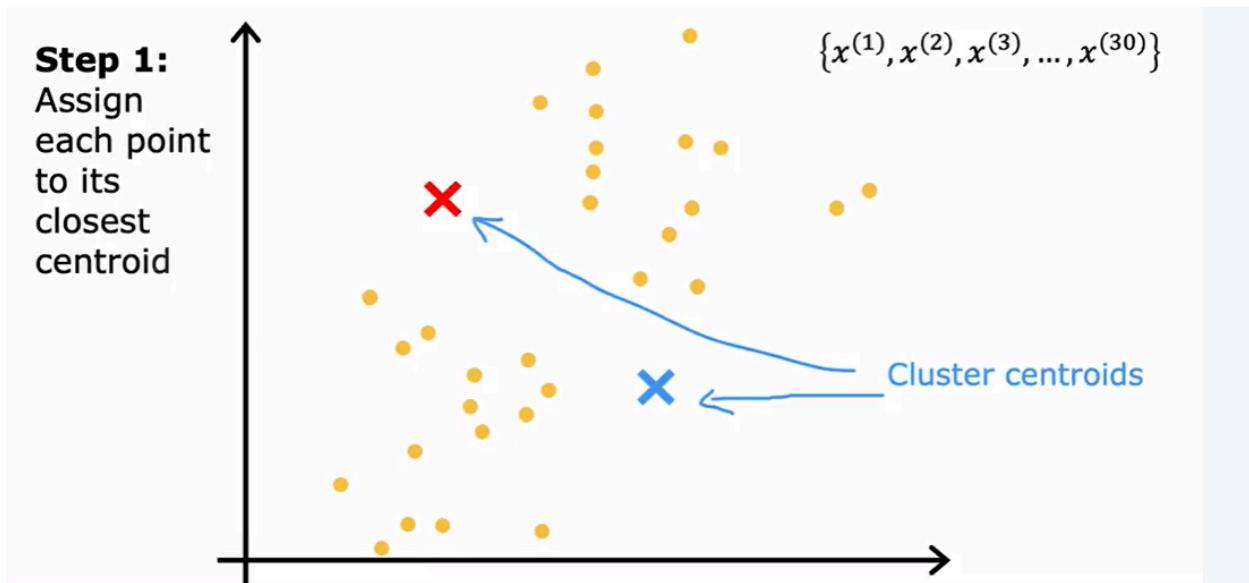
Find **natural groupings** or structure in the data — for example:

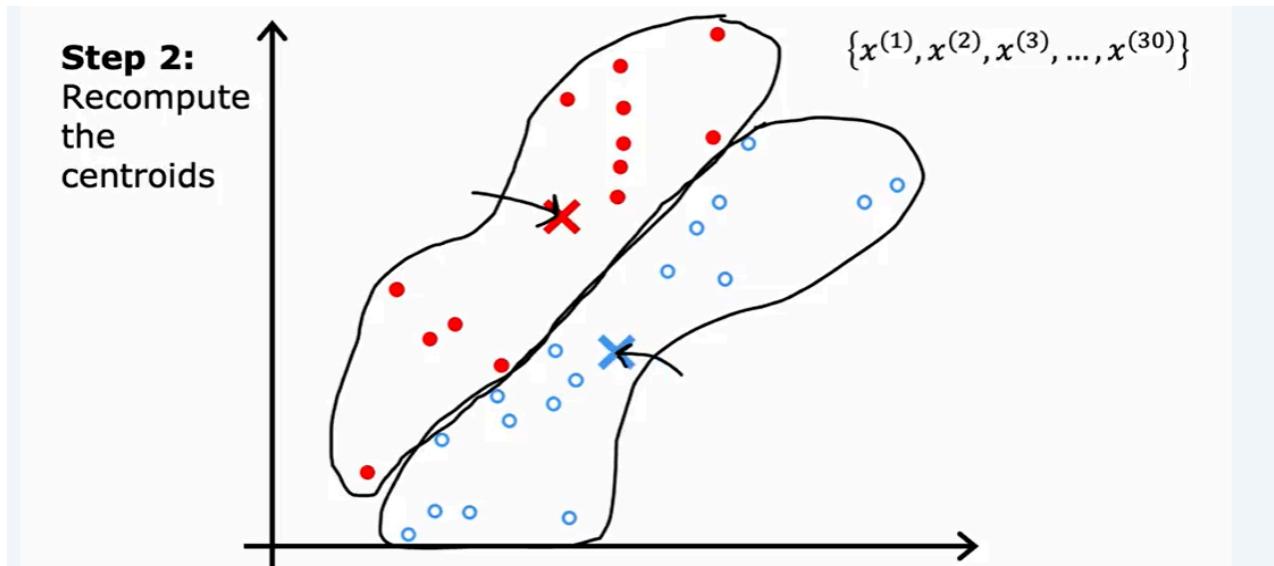
- Grouping customers by buying behavior.
- Organizing documents by topic.
- Segmenting images based on color or texture.

Popular Clustering Algorithm:

Algorithm	Description
K-Means	Partitions data into k clusters by minimizing the distance to cluster centroids.
Hierarchical Clustering	Builds nested clusters using a tree-like structure (dendrogram).
DBSCAN	Groups based on density, great for finding clusters of arbitrary shapes and removing outliers.
Gaussian Mixture Models (GMM)	Uses probabilistic approach assuming data comes from a mixture of Gaussians.

K-means Clustering :





K means clustering algorithm:

K-means algorithm

Randomly initialize K cluster centroids

$$\mu_1, \mu_2, \dots, \mu_K$$

$n=2$ ~~$K=2$~~ $K=K-1$

Repeat {

Assign points to cluster centroids

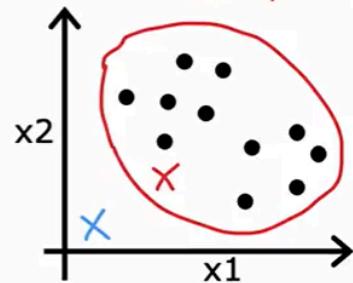
for $i = 1$ to m
 $c^{(i)} :=$ index (from 1 to K) of cluster
 centroid closest to $x^{(i)}$

Move cluster centroids

for $k = 1$ to K

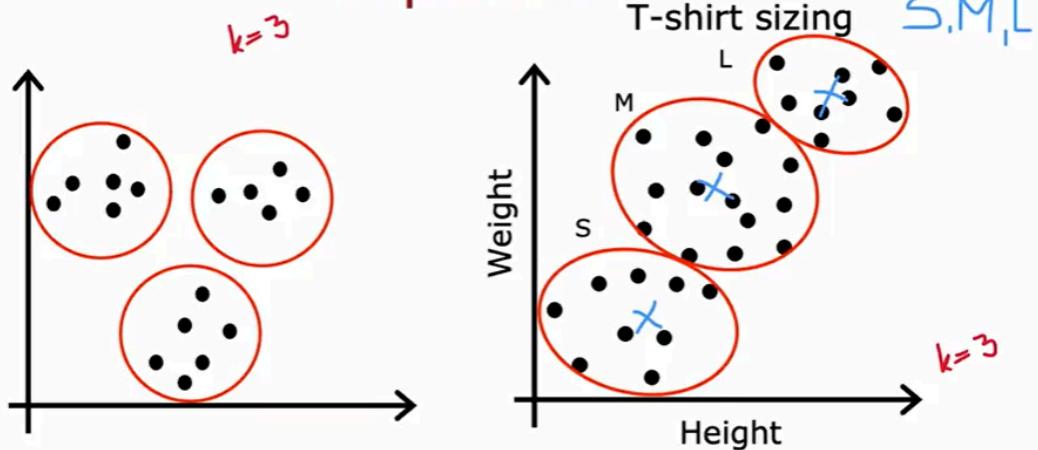
$\mu_k :=$ average (mean) of points assigned to cluster k

}



Here , if any cluster is not associated with any point then just remove that cluster. To find the new cluster position just find average coordinate point of both x and y for the point that cluster associated with. That will be new point of that cluster.

K-means for clusters that are not well separated



Even if the data are not lied in well separated , the k-means clustering can do clustering on them.

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

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

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

The objective function here is used to check for a classification problem how many number of cluster is needed and in this case number of cluster vs objective function value is taken where number of cluster will be in x -axis and objective function value will be in y-axis and an elbow shape curved will be created and from that graph the number of cluster of that elbow point is appropriate.

Cost Function in K-Means (a.k.a. Objective Function)

K-Means aims to minimize the **total intra-cluster variance**, also called the **Within-Cluster Sum of Squares (WCSS)**.

Mathematical Form:

$$J = \sum_{i=1}^K \sum_{x_j \in C_i} \|x_j - \mu_i\|^2$$

Where:

- K : number of clusters
- C_i : set of points in the **i-th** cluster
- μ_i : centroid (mean) of cluster C_i
- x_j : data point in cluster C_i
- $\|x_j - \mu_i\|^2$: squared Euclidean distance from the point to the cluster center

What It Means:

- For each cluster, compute the squared distance between each point and its centroid.
- Sum up all those distances across all clusters.
- The goal is to minimize this total cost, so that data points are as close as possible to their respective centroids.

Intuition:

- The better the clustering, the smaller the total distance between points and their cluster centers.
 - During training, K-Means **reassigns points and recomputes centroids** to keep reducing this cost.
-

Elbow Method & Cost:

You can use the cost function values (WCSS) at different values of **K** to visualize the **elbow method**, which helps decide the optimal number of clusters.

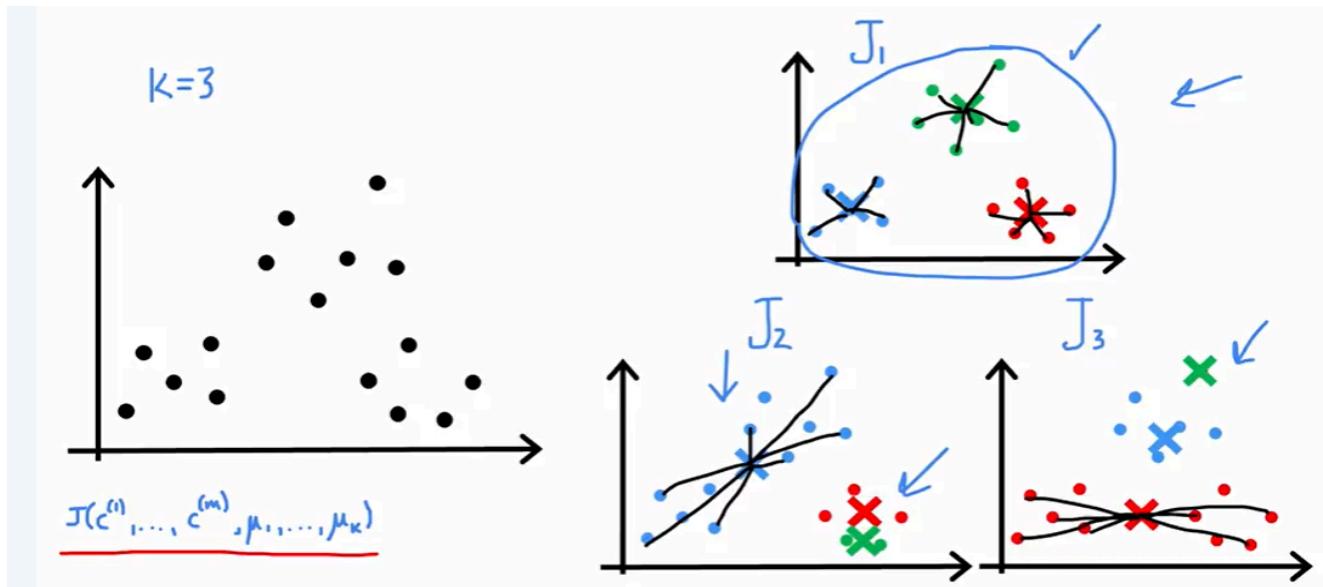
Optimization Strategy:

K-Means uses an **iterative greedy approach** (Lloyd's Algorithm):

1. Assign points to the nearest cluster.
2. Recompute the centroids.
3. Re-evaluate the cost.
4. Repeat until convergence (no major change in cost or assignments).

After every single iteration of k-means the WCSS function or distortion function value should go down otherwise there is a bug in the algorithm .

Initializing K-means:



K-means algorithm:

K-means algorithm

Step 0: Randomly initialize K cluster centroids $\mu_1, \mu_2, \dots, \mu_k$

Repeat {

Step 1: Assign points to cluster centroids

Step 2: Move cluster centroids

}

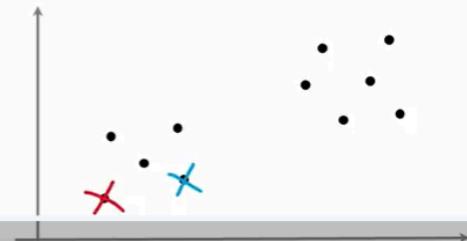
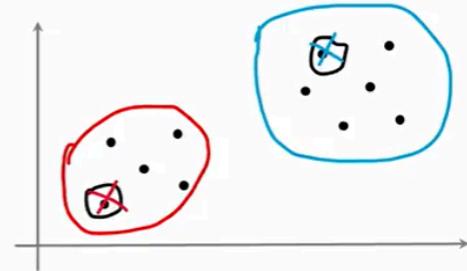
Step-0:

Random initialization

Choose $K < m$

Randomly pick K training examples.

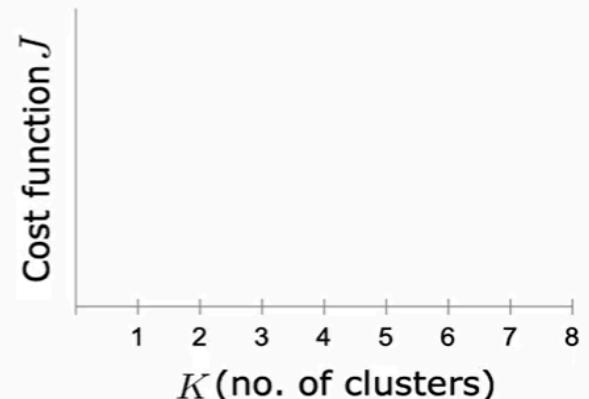
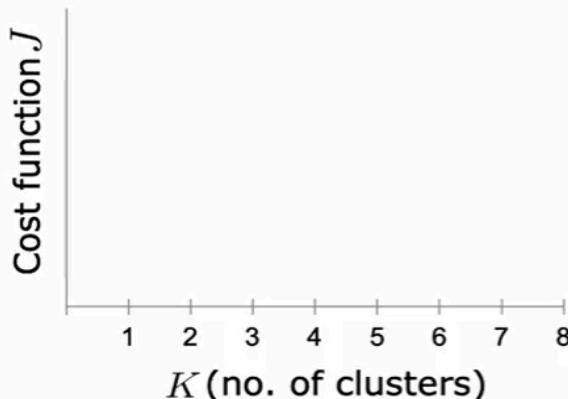
Set $\mu_1, \mu_2, \dots, \mu_k$ equal to these K examples.



Choosing the right value of K:

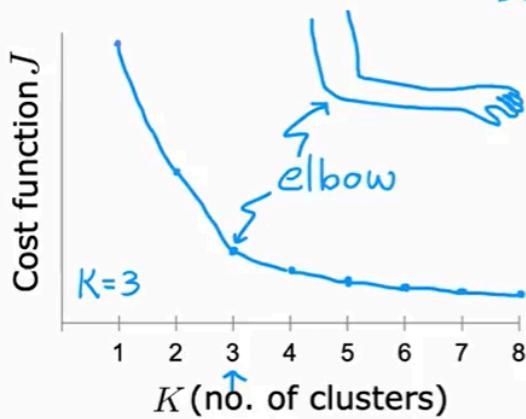
Choosing the value of K

Elbow method

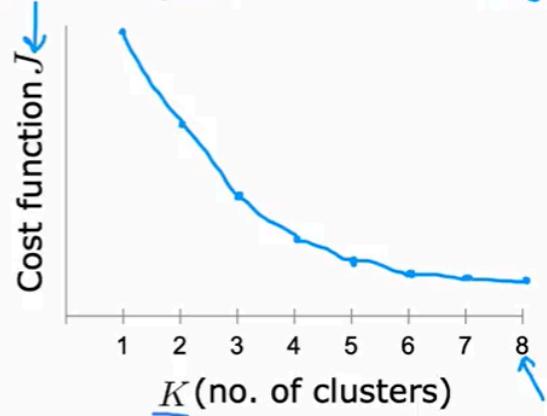


Choosing the value of K

Elbow method

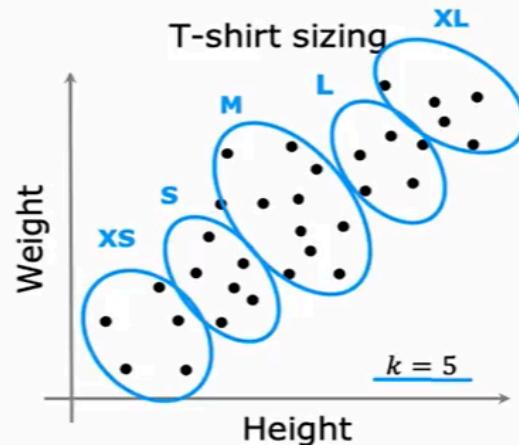
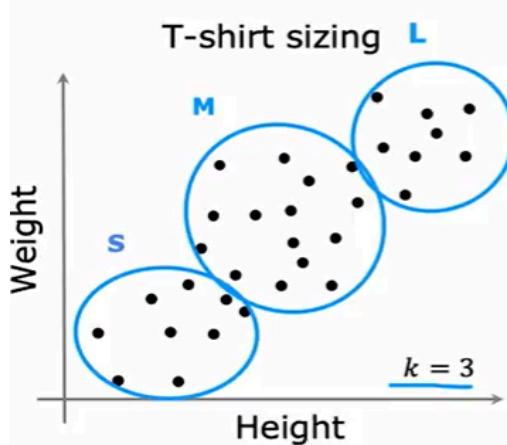


the right "K" is often ambiguous
Don't choose K just to minimize cost J

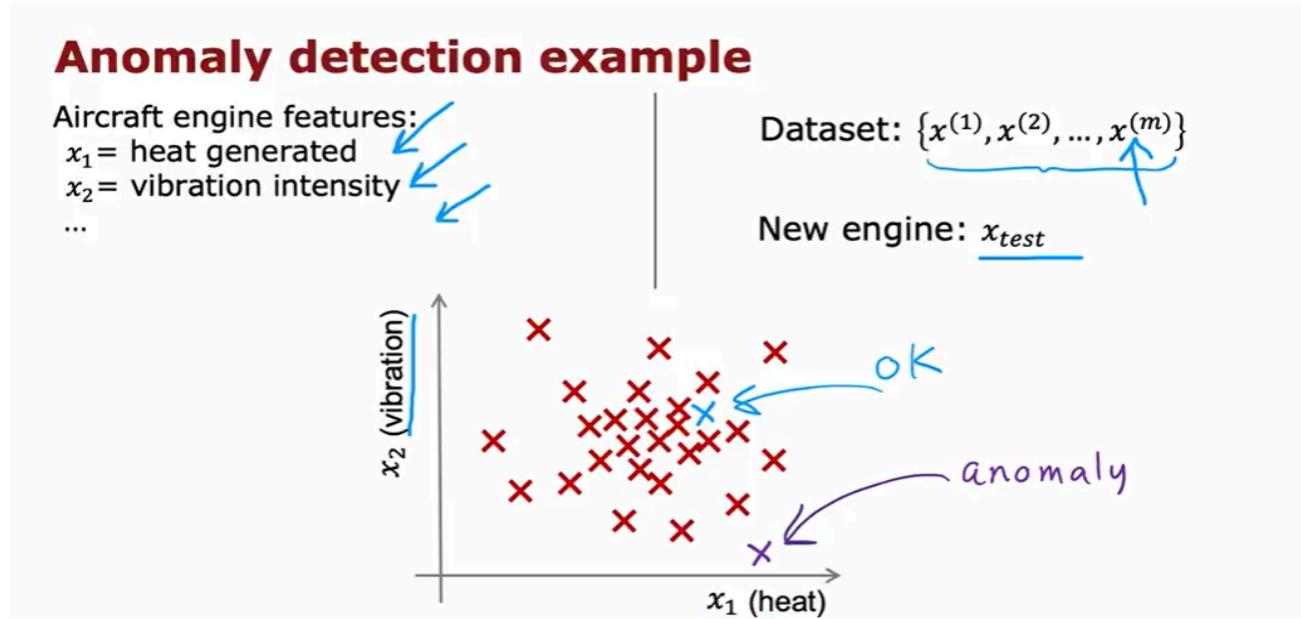


Choosing the value of K

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



Anomaly detection :



📊 The Data (Training Set)

- You have a **dataset** of normal engine behavior:
$$\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$$
- Each **x** is a point with features like:
 - x_1 : heat generated
 - x_2 : vibration intensity
- These are shown as **red crosses** clustered together on the plot.



New Data Point

- You get a **new engine** data point: x_{test}
- You want to check if it's behaving **normally** or it's an **anomaly**.



How to Detect Anomaly

You typically do this in one of two ways:

1. Distance-based Approach:

- If x_{test} is **too far** from the cluster of normal data → It's an anomaly.
- In the graph:
 - The blue point is **OK** (close to the cluster).
 - The purple point (bottom-left corner) is **far away**, so it's marked as **anomaly**.

2. Probability-based (Gaussian):

- Fit a probability distribution (like Gaussian) on normal data.
- If the **probability** of x_{test} under that model is **very low**, it's likely an anomaly.



Final Goal:

Identify whether the new data point is **consistent** with what you've seen before or **suspiciously different**.

*For a single number gaussian or normal distribution works to detect anomaly , if the probability of the single feature output is very low then it is an anomaly. But in case of multi featured case there are algorithms for anomaly detection.

Anomaly detection for multi featured dataset:

Density estimation

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

Each example $\vec{x}^{(i)}$ has n features

$$\vec{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$p(\vec{x}) = p(x_1; \mu_1, \sigma_1^2) * p(x_2; \mu_2, \sigma_2^2) * p(x_3; \mu_3, \sigma_3^2) * \dots * p(x_n; \mu_n, \sigma_n^2)$$

$$= \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2) \quad \sum \text{"add"} \quad \prod$$

$$\begin{aligned} p(x_1 = \text{high temp}) &= 1/10 \\ p(x_2 = \text{high vibra}) &= 1/20 \\ p(x_1, x_2) &= p(x_1) * p(x_2) \\ &= \frac{1}{10} \times \frac{1}{20} = \frac{1}{200} \end{aligned}$$

Anomaly detection example

Fraud detection:

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

how often log in?
how many web pages visited?
transactions?
posts? typing speed?

perform additional checks to identify real fraud vs. false alarms

Manufacturing:

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

airplane engine

circuit board

smartphone

ratios

Monitoring computers in a data center:

$x^{(i)}$ = 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.

Gaussian (Normal) Distributions:

The **Gaussian distribution** is the same as the **Normal distribution** — they are two names for the **same concept** in statistics.

Quick Explanation:

- "Normal distribution" is the more common name in general statistics.
 - "Gaussian distribution" is named after **Carl Friedrich Gauss**, who described it mathematically.
-

Characteristics (as shown in your image):

- Bell-shaped curve (like the **Liberty Bell** shown for analogy).
- Centered around **mean (μ)**.
- Spread is measured by **standard deviation (σ)**.
- The **variance** is σ^2 .
- The peak of the curve is at μ , and the spread is determined by σ .

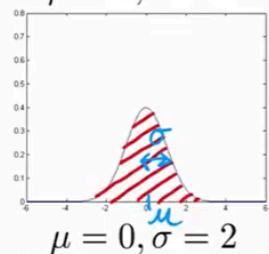
In math form:

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

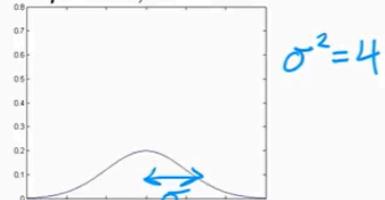
Different choices of mu and sigma affects the gaussian distribution.

Gaussian distribution example

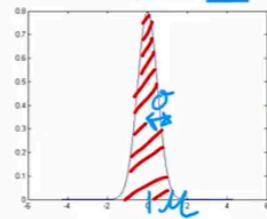
$$\mu = 0, \sigma = 1$$



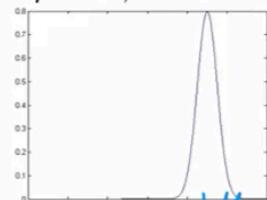
$$\mu = 0, \sigma = 2$$



$$\mu = 0, \sigma = 0.5$$

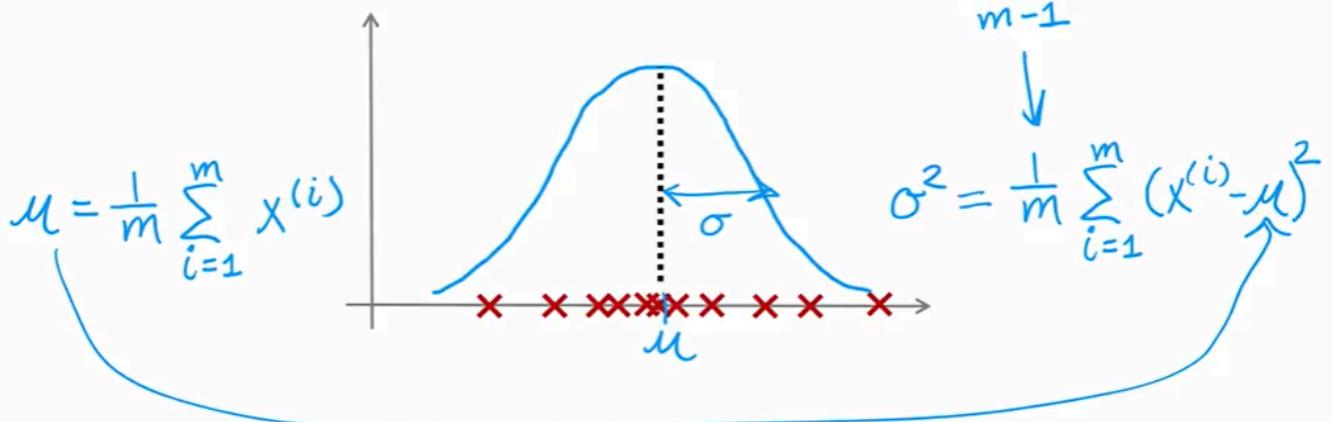


$$\mu = 3, \sigma = 0.5$$



Parameter estimation

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



maximum likelihood
for μ, σ

Algorithm for anomaly detection:

Anomaly detection algorithm

1. Choose n features x_i that you think might be indicative of anomalous examples.
2. Fit parameters $\mu_1, \dots, \mu_n, \sigma_1^2, \dots, \sigma_n^2$

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

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

Vectorized formula

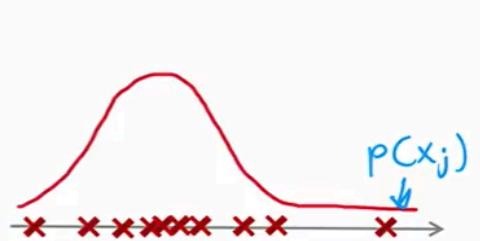
$$\vec{\mu} = \frac{1}{m} \sum_{i=1}^m \vec{x}^{(i)}$$

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

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

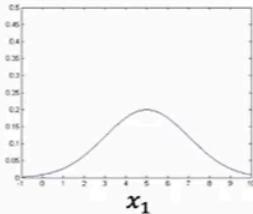
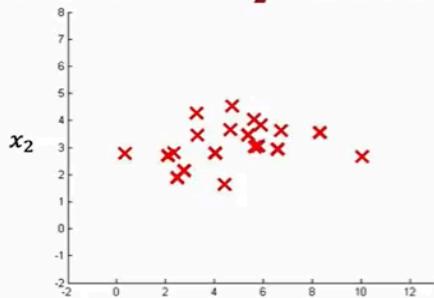
$$p(x) = \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2) = \prod_{j=1}^n \frac{1}{\sqrt{2\pi}\sigma_j} \exp\left(-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}\right)$$

Anomaly if $p(x) < \varepsilon$



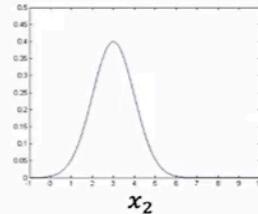
When multiplied $p(x_1)$ and $p(x_2)$ then the graph looks like a 3d shape.

Anomaly detection example



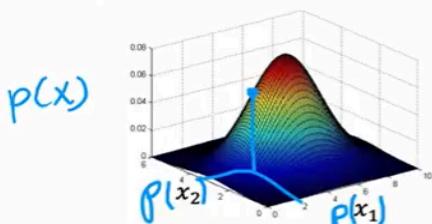
$$\mu_1 = 5, \sigma_1 = 2$$

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

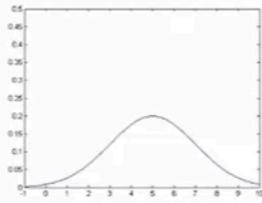
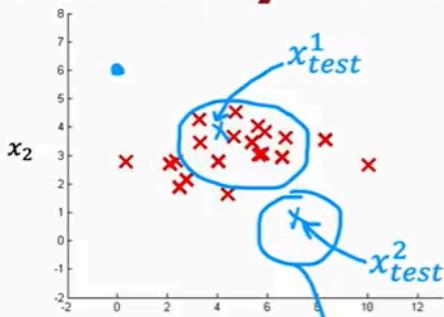


$$\mu_2 = 3, \sigma_2 = 1$$

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

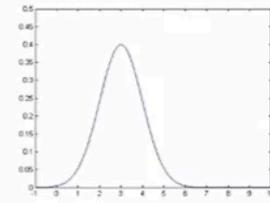


Anomaly detection example



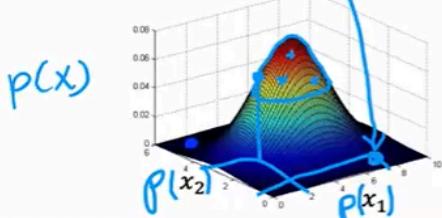
$$\mu_1 = 5, \sigma_1 = 2$$

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



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

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



$$\varepsilon = 0.02$$

$$p(x_{test}^{(1)}) = 0.0426 \longrightarrow "ok"$$

$$p(x_{test}^{(2)}) = 0.0021 \longrightarrow \text{anomaly}$$

Algorithm evaluation

Fit model $p(x)$ on training set $x^{(1)}, x^{(2)}, \dots, x^{(m)}$

On a cross validation/test example x , predict

course 2 week 3
skewed datasets

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

10
2000

Possible evaluation metrics:

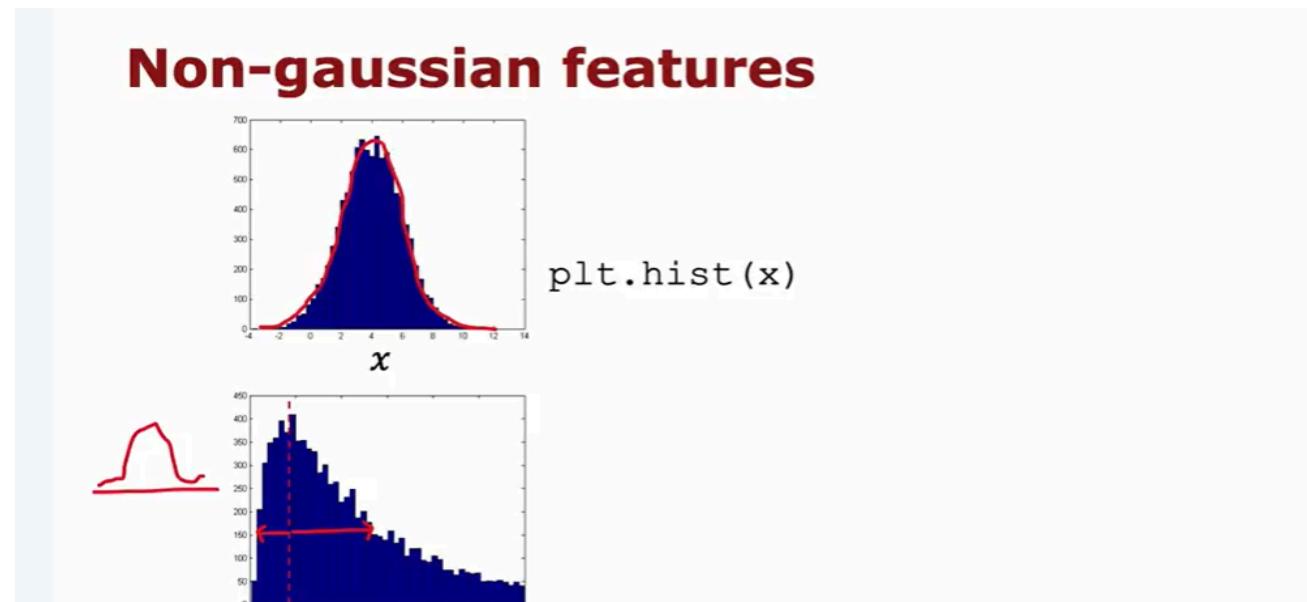
- True positive, false positive, false negative, true negative
- Precision/Recall
- F₁-score

Use cross validation set

Anomaly detection vs supervised learning :

Anomaly detection	vs.	Supervised learning
<ul style="list-style-type: none">→ Fraud detection• Manufacturing - Finding new previously unseen defects in manufacturing.(e.g. aircraft engines)• Monitoring machines in a data center<li style="text-align: center;">⋮		<ul style="list-style-type: none">• Email spam classification• Manufacturing - Finding known, previously seen defects• Weather prediction (sunny/rainy/etc.)• Diseases classification<li style="text-align: center;">⋮

Choosing what features to use:

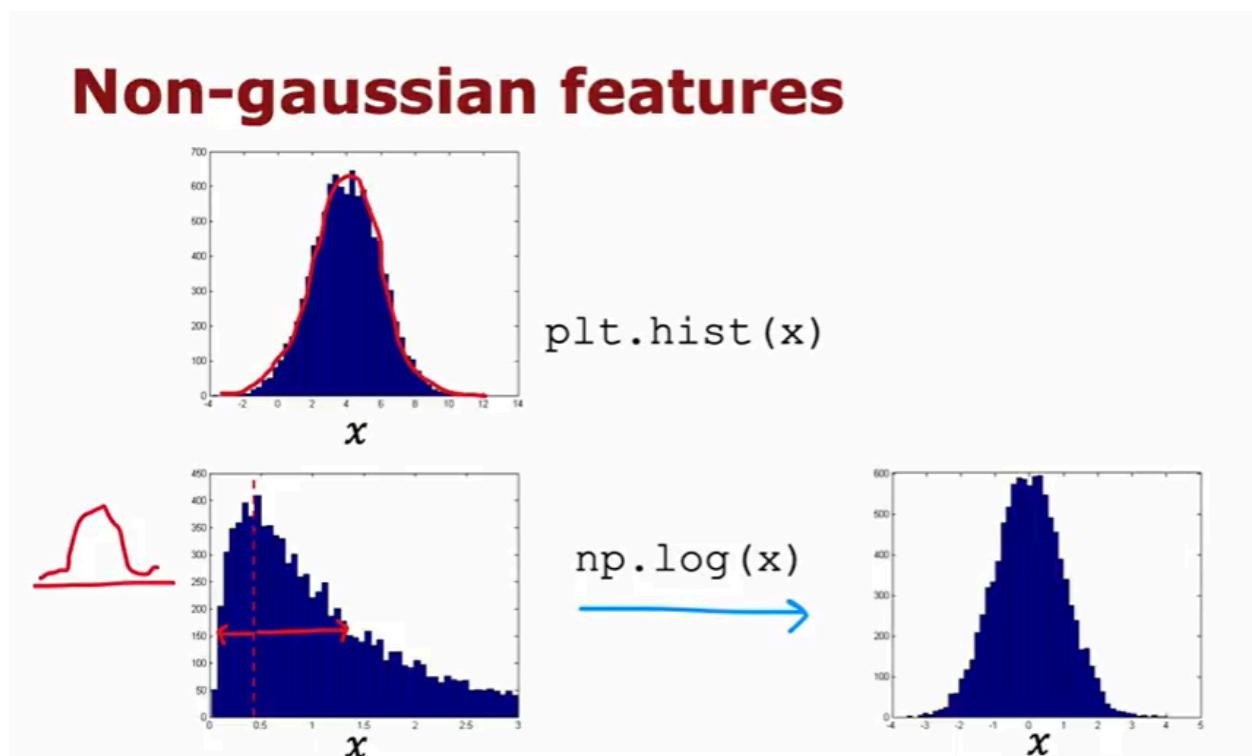


The first plot is looking like gaussian but actually it is not a gaussian features.

🔍 About the First Plot

- It looks like a bell (symmetric and unimodal), but:
 - Just looking bell-shaped ≠ truly Gaussian.
 - It might not follow the exact mathematical formula of a Gaussian.
 - Even small deviations (e.g., **heavy tails, skew, kurtosis**) make it *non-Gaussian*.

So, it's visually similar to a Gaussian, but the underlying distribution might not be a perfect Gaussian.



From non-gaussian features to turn into gaussian like features.

Common ways to make a feature non-gaussian:

Transformation	Use Case	Code Example
Log Transform	Skewed right (long tail on the right)	<code>np.log(x + 1)</code>
Square Root Transform	Moderate right-skew	<code>np.sqrt(x)</code>
Box-Cox Transform	Best for positive data, tries different powers	<code>scipy.stats.boxcox(x)</code>
Yeo-Johnson Transform	Handles both positive and negative values	<code>PowerTransformer(method='yeo-johnson')</code>
Z-score Standardization	Not Gaussianizing, but helps shape & scale	<code>(x - mean) / std</code>

Recommender System :

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	?

*one to five stars
zero*

Ratings

★				
★	★			
★	★	★		
★	★	★	★	
★	★	★	★	★

$n_u = \text{no. of users}$

$n_m = \text{no. of movies}$

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

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

$n_u = 4$ $r(1,1) = 1$ $y^{(1,1)} = 5$

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

PCA[Principal Component Analysis] :

Principal Component Analysis (PCA) is a statistical technique used for dimensionality reduction. It transforms a set of possibly correlated variables (features) into a smaller number of uncorrelated variables called principal components. These components capture the most significant variations in the data, with the first component accounting for the most variance, the second for the second most, and so on.

PCA is widely used for:

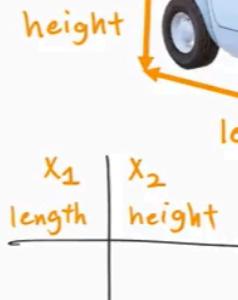
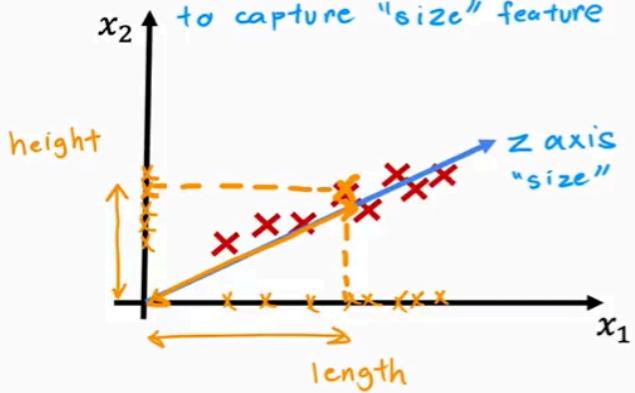
1. **Reducing the number of features:** This is particularly useful when you have high-dimensional data (many features) and want to reduce computational complexity without losing much information.
2. **Data visualization:** By reducing dimensions to 2 or 3 components, you can plot data points and observe patterns.
3. **Noise reduction:** By keeping only the most significant components, PCA helps remove noise and irrelevant variations.

The steps involved in PCA:

1. **Standardize the data** (subtract the mean and divide by the standard deviation) to ensure that all features are on the same scale.
2. **Calculate the covariance matrix** to understand how the features are correlated with each other.
3. **Compute the eigenvalues and eigenvectors** of the covariance matrix.
4. **Select the principal components** (the eigenvectors corresponding to the largest eigenvalues).
5. **Transform the data** into the new coordinate system formed by the principal components.

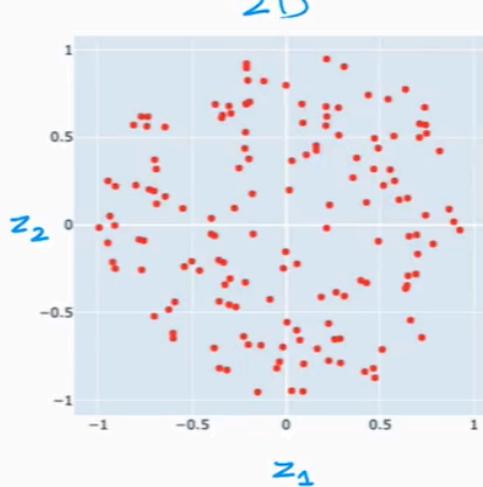
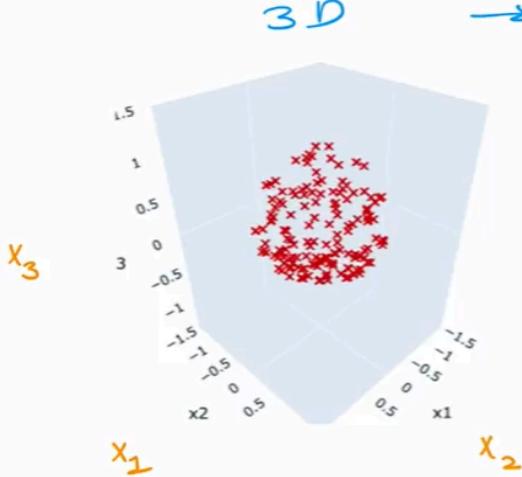
Size

PCA: find new axis and coordinates
use fewer numbers
to capture "size" feature



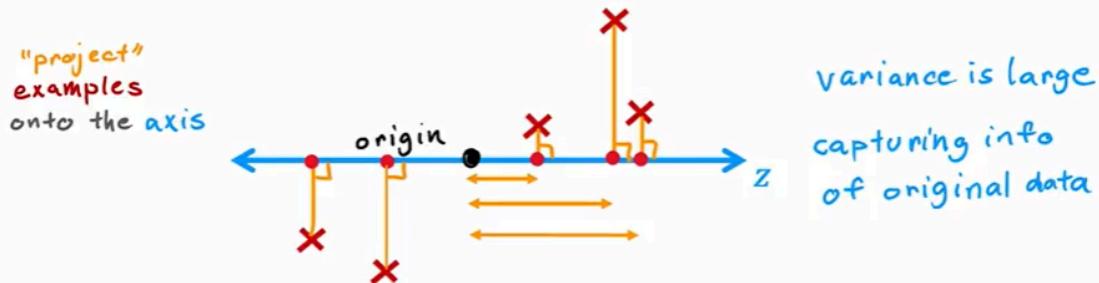
$2 \rightarrow 1$ features
many features $\rightarrow 2$ or 3 features

From 3D to 2D



Instead of having 3 features x_1, x_2, x_3 reduces it into z_1 and z_2 axis and it is convenient way to visualize.

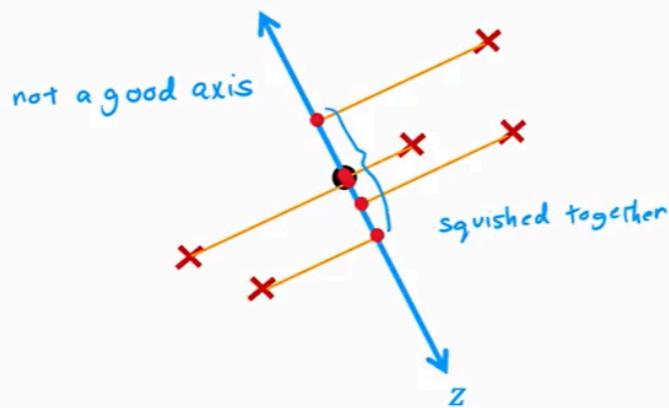
Choose an axis



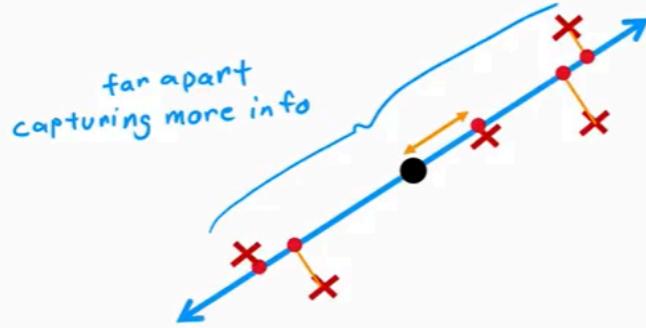
The projection from the real data point to the z axis should be in such a way that the variance among the data point projected onto the z axis should be maximum . This means the z axis should capture as much information as it can from the real data point .

This type of choice is squished and the variance is much less here. This type of distribution is discouraged to apply as principle component.

Choose an axis



Choose an axis

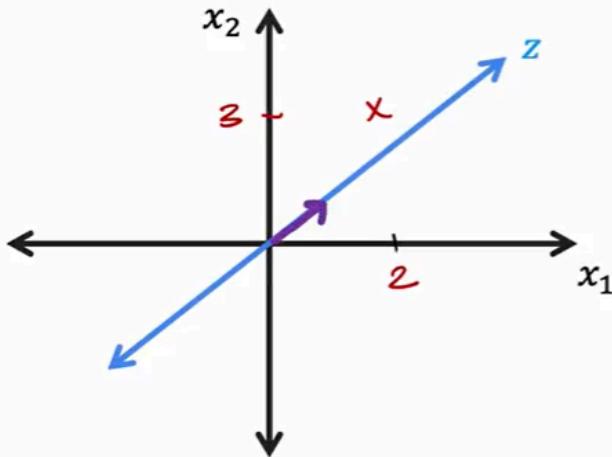


This is capturing a large amount of information and it is called principle component.

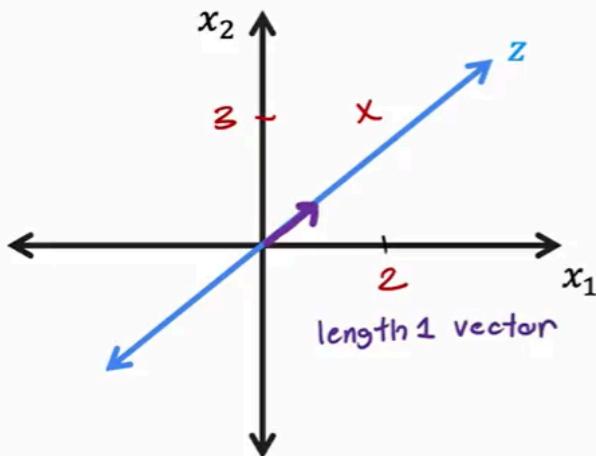
Example:

Let's say for 2 data point 2 and 3 i want to draw or find the PCA it can be done by doing so

Coordinate on the new axis

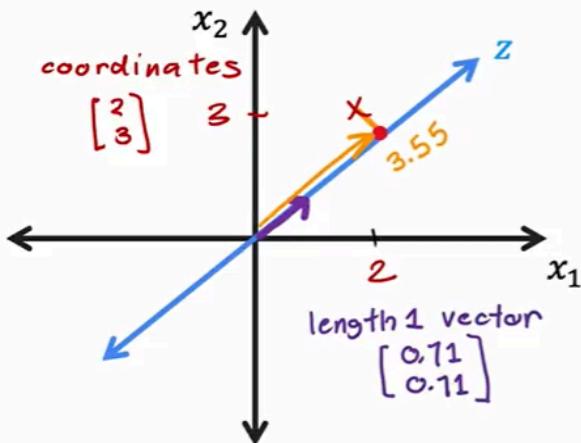


Coordinate on the new axis



The PCA axis is z axis. Length 1 vector .

Coordinate on the new axis



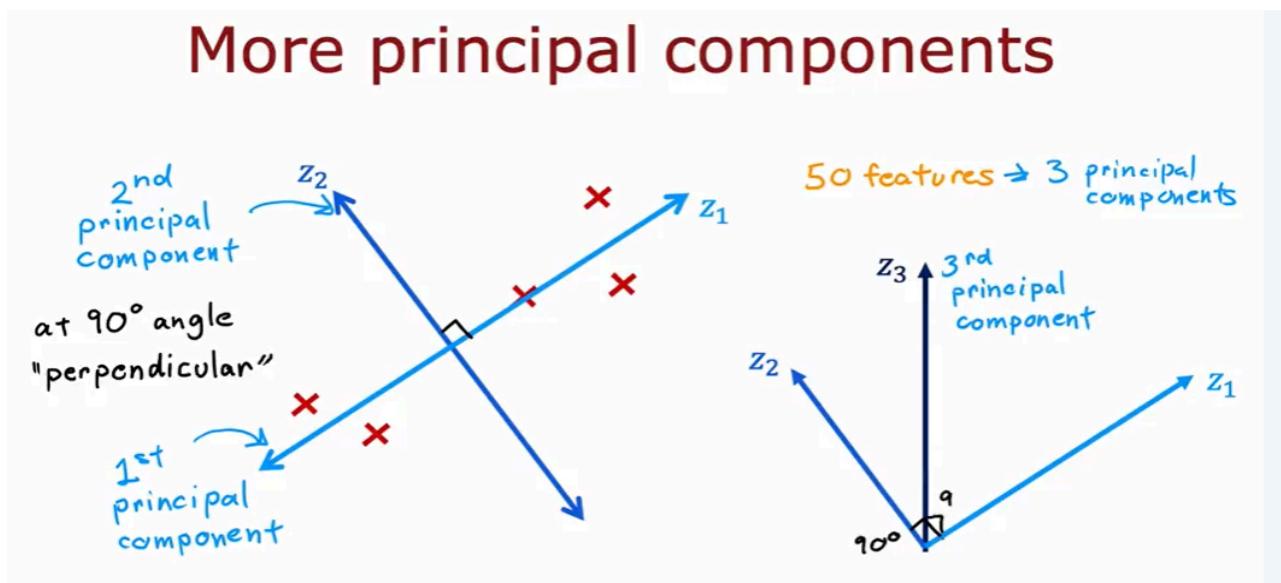
dot product

$$\begin{bmatrix} 2 \\ 3 \end{bmatrix} \cdot \begin{bmatrix} 0.71 \\ 0.71 \end{bmatrix}$$

$$2 \times 0.71 + 3 \times 0.71 = 3.55$$

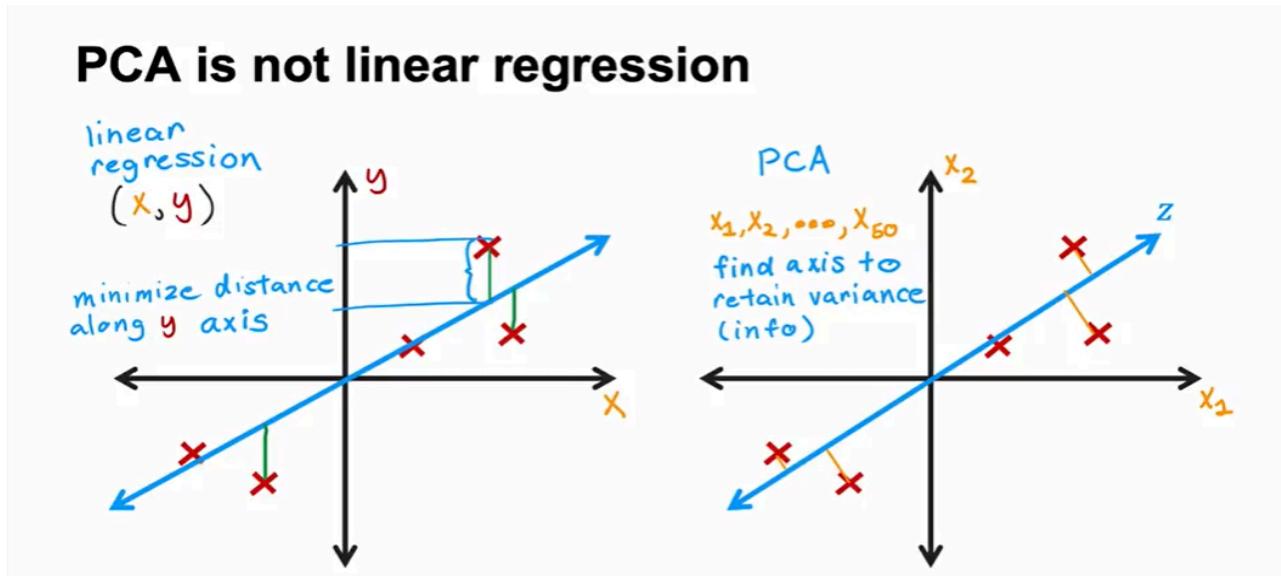
3.55 is the distance from the origin through the pca axis or z axis.

More PCA :



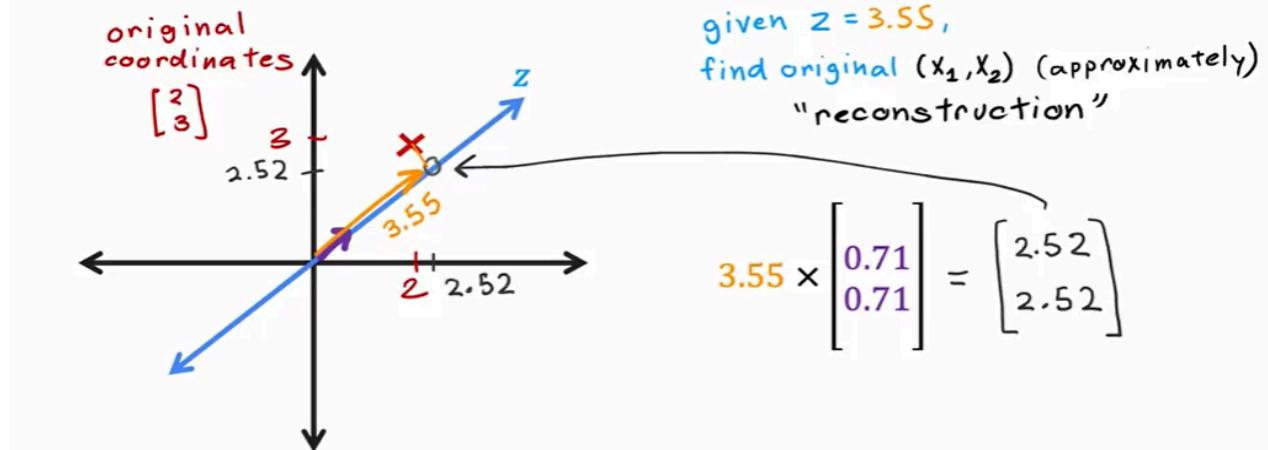
There can be more than one PCA for data point.

PCA vs Linear regression:



Approximate to the original datapoint:

Approximation to the original data



Reconstruction step of PCA.

PCA step in scikit learn:

PCA in scikit-learn

Optional pre-processing: Perform feature scaling

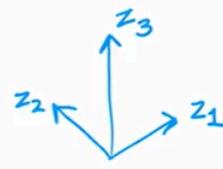


for visualization

1. "fit" the data to obtain 2 (or 3) new axes (principal components)
fit includes mean normalization



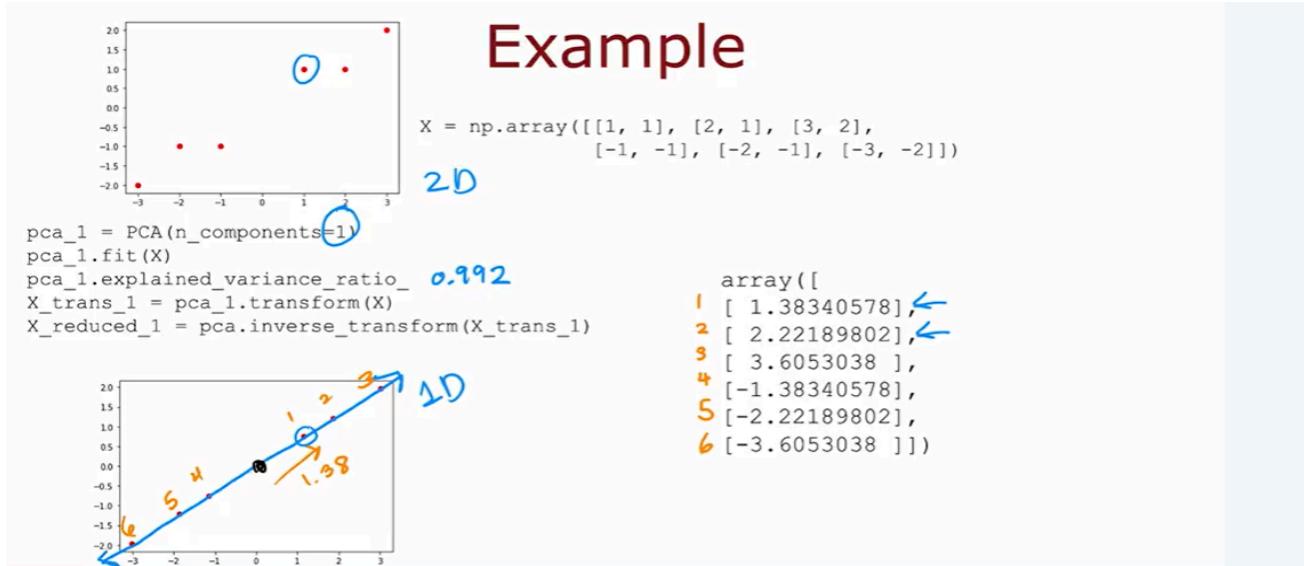
2. Optionally examine how much variance is explained by each principal component.
explained_variance_ratio



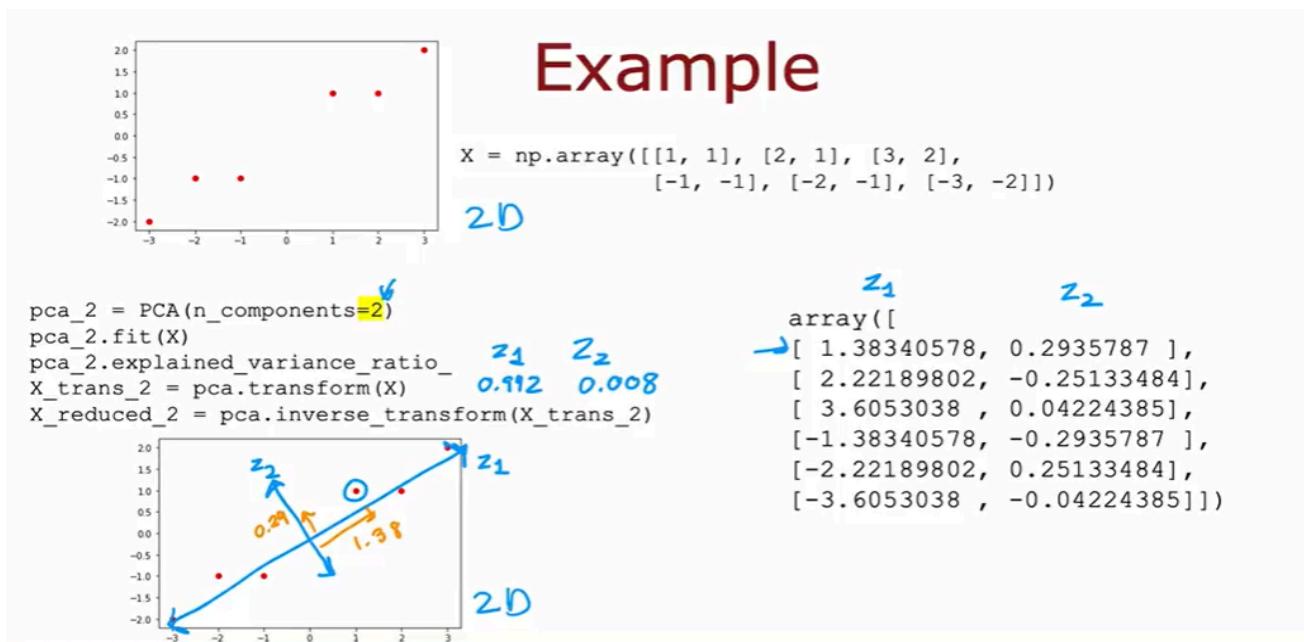
3. Transform (project) the data onto the new axes
transform



Example:



For 2 PCA



Applications of PCA

❖ Visualization *reduce to 2 or 3 features*

Less frequently used for:

- Data compression
(to reduce storage or transmission costs) $50 \rightarrow 10$
- Speeding up training of a supervised learning model

$$n = 1000 \rightarrow 100$$

Reinforcement Learning:

Reinforcement Learning (RL) is a type of machine learning where an **agent** learns how to behave in an environment by **interacting with it** and receiving **rewards or penalties**. The goal of the agent is to **maximize the total cumulative reward** over time.

Core concept:

1. **Agent**: The learner or decision-maker.
2. **Environment**: The external system the agent interacts with.
3. **State (S)**: The current situation of the environment.
4. **Action (A)**: A choice the agent makes.
5. **Reward (R)**: Feedback from the environment after an action.
6. **Policy (π)**: Strategy the agent uses to decide actions based on states.
7. **Value Function (V)**: Expected long-term reward of a state.
8. **Q-Function (Q)**: Expected long-term reward of a state-action pair.
9. **Episode**: One run of interaction from start to finish.

Return

	0	0		0	
state 1	2	3	4	5	6

$$\text{Return} = 0 + (0.9)0 + (0.9)^20 + (0.9)^3100 = 0.729 \times 100 = 72.9$$

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

Discount Factor $\gamma = 0.9 \quad 0.99 \quad 0.999$
 $\gamma = 0.5$

$$\text{Return} = 0 + (0.5)0 + (0.5)^20 + (0.5)^3100 = 1$$

Return in reinforcement learning is sum of the weighted reward

Example of Return

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

\leftarrow return $\gamma = 0.5$
 \leftarrow reward

The return depends on the actions you take.

100	2.5	5	10	20	40
100	0	0	0	0	40
1	2	3	4	5	6

$0 + (0.5)0 + (0.5)^240 = 10$

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

$0 + (0.5)40 = 20$

Making decision : policies in reinforcement learning:

Policy (π) defines the agent's way of behaving at a given time — that is, what action to take given a state.

Formally:

$$\pi(a|s) = P(A_t = a | S_t = s)$$

- It maps states (s) to a distribution over actions (a).
- A policy can be deterministic or stochastic:
 - Deterministic Policy: Always picks a fixed action for a state.

$$a = \pi(s)$$

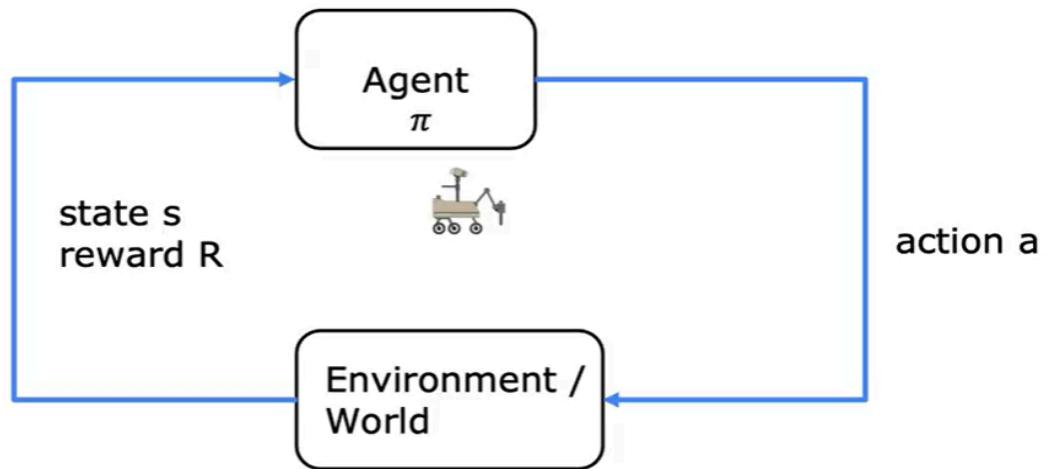
- Stochastic Policy: Picks actions based on a probability distribution.

$$\pi(a|s) = \text{Probability of taking action } a \text{ when in state } s$$

	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	 100 ← ← ← → 40	Find $\pi(s) = a$	Find $\pi(s) = a$

This formalism of reinforcement learning algorithm has a name , it is called MDP [Markov Decision Process]

Markov Decision Process (MDP)



State action value:

The **State-Action Value function**, often denoted as $Q(s, a)$, tells us:

“What is the expected cumulative reward if the agent starts in **state s**, takes **action a**, and follows a policy π afterward?”

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.

$Q(s, a)$

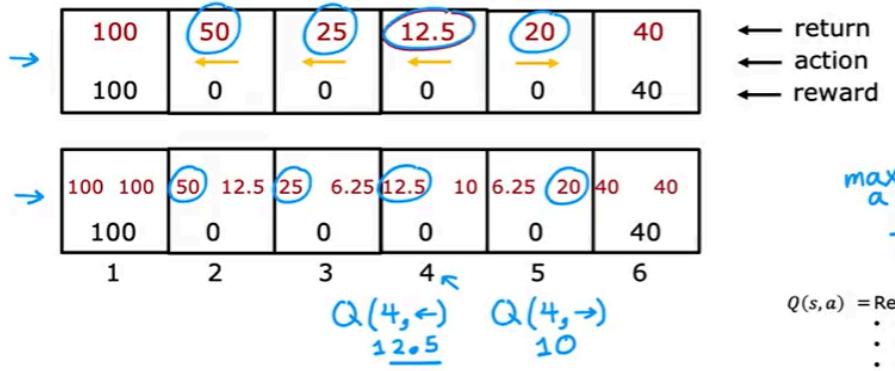
100	50	25	12.5	20	40
100	0	0	0	0	40

← return
← action
← reward

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

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

Picking actions



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

$$\pi(s) = a$$

$Q(s, a)$ = 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

★ What is $Q^*(s, a)$?

$Q^*(s, a)$ is the **maximum expected return** (cumulative discounted reward) **achievable by any policy**, starting from state s , taking action a , and **then following the optimal policy**.

$$Q^*(s, a) = \max_{\pi} Q^{\pi}(s, a)$$

It represents the **best possible value** for each (state, action) pair.

Bellman Equation:

🎯 What is the Bellman Equation **used for**?

The Bellman Equation helps us:

1. Break down the value of a policy or action into immediate reward + future value.
2. Formulate recursive relationships to solve for value functions or Q-values.
3. Enable dynamic programming methods like value iteration and policy iteration.
4. Serve as the foundation for learning algorithms like Q-learning, SARSA, and Deep Q-Networks (DQN).

Explanation of Bellman Equation

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

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

$$Q(s, a) = \underbrace{R(s)}_{\text{Reward you get right away}} + \gamma \max_{a'} Q(s', a')$$

Return from behaving optimally starting from state s' .

$$\begin{aligned} & R_1 + \gamma R_2 + \gamma^2 R_3 + \gamma^3 R_4 + \dots \\ & = R_1 + \gamma [R_2 + \gamma R_3 + \gamma^2 R_4 + \dots] \end{aligned}$$

Explanation of Bellman Equation

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

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

$$\begin{aligned} & Q(4, \leftarrow) \\ & = 0 + (0.5)0 + (0.5)^2 0 + (0.5)^3 100 \\ & = R(4) + (0.5)[0 + (0.5)0 + (0.5)^2 100] \\ & = R(4) + (0.5) \max_{a'} Q(3, a') \end{aligned}$$

It is nothing but breaking down the reward what u get right away and what u will get right into the future.

Bellman Equation

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

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

$$\begin{aligned} s &= 2 \\ a &= \rightarrow \\ s' &= 3 \end{aligned}$$

$$\begin{aligned} Q(2, \rightarrow) &= R(2) + 0.5 \max_{a'} Q(3, a') \\ &= 0 + (0.5)25 = 12.5 \end{aligned}$$

$$\begin{aligned} Q(4, \leftarrow) &= R(4) + 0.5 \max_{a'} Q(3, a') \\ &= 0 + (0.5)25 = 12.5 \end{aligned}$$

$$\begin{aligned} s &= 4 \\ a &= \leftarrow \\ s' &= 3 \end{aligned}$$

*** It is like dynamic programming .

Random(Stochastic) Environment:

A **stochastic (random) environment** is one where the **next state and reward are not fully predictable**, even if the agent takes the same action in the same state.

In other words: Taking the same action in the same situation **can lead to different results** — due to randomness in the environment.

Example:

Let's say an agent is in a grid world (like a maze), and it decides to move "right" from a certain state.

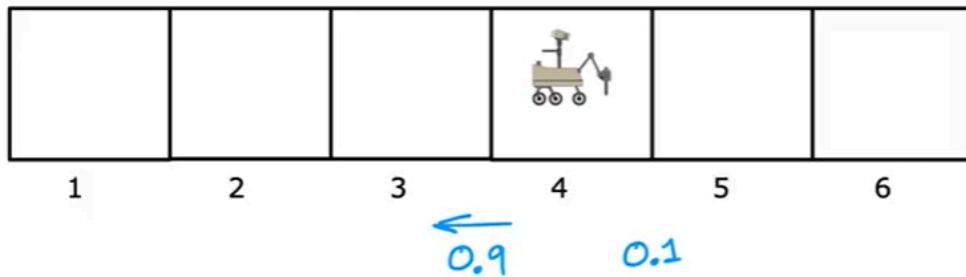
In a **deterministic environment**, it always goes right.

In a **stochastic environment**, it might:

- go right 80% of the time,
- go up 10% of the time,
- go down 10% of the time.

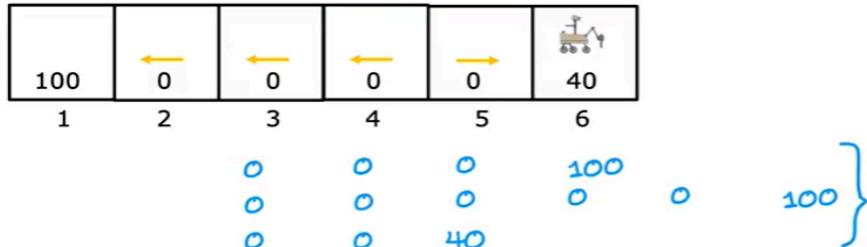
This randomness is what makes the environment **stochastic**.

Stochastic Environment



Here the mars rover is in 4 and here 0.9 percent chance that it will end up in index 3 and 0.1 percent chance that it will end up in index 5 so this kind of situation is called stochastic environment .

Expected Return



$$\begin{aligned}\text{Expected Return} &= \text{Average}(R_1 + \gamma R_2 + \gamma^2 R_3 + \gamma^3 R_4 + \dots) \\ &= E[R_1 + \gamma R_2 + \gamma^2 R_3 + \gamma^3 R_4 + \dots]\end{aligned}$$

Here, from a state randomly various sequence can be achieved by the agent and for each sequence the reward is averaged and this is called the expected value. The goal of this RL algorithm is to find the policy (π) to maximize the sum of discounted rewards.

The previous bellman equation will be changed slightly here as expected issues will be merged.

Expected Return

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:

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

↑ ↑ ↓
3 2 or 4

Here from state 3, agent can move to 2 or 4 randomly that means a' can be left or right and for this the average will be taken.

So the total return taking action a from state s is, once in a behaving optimally is equal to the reward r getting right away + what you expect to get on average of the future returns.

Optimal policy					
100.0 100	46.06 0	21.25 0	10.49 0	18.52 0	40.0 40
Q(s,a)					
100.0 100	100.0 0	46.06 0	14.56 0	21.25 0	7.02 0

Here the value of optimal policy is given and Q value is also given , with increase of mis_step the reward value gets reduced and the above is for mis_step 0.1 ,but when mis_step is 0.4 the value of optimal policy gets reduced more as mis_step 0.4 means 40% time the agent misses the step that is probability of going in wrong direction

Optimal policy					
100.0 100	32.18 0	10.88 0	6.15 0	13.23 0	40.0 40
Q(s,a)					
100.0 100	100.0 0	32.18 0	23.26 0	10.88 0	8.28 0

Continuous State Space:

In a **continuous state space**, the **states are not countable** — they can take **any real value within a range**.

Instead of a few discrete states (like s1, s2, s3), you might have **infinite possible states** like:

$$s=(x=3.52, y=-1.75)$$

📌 Examples of Continuous States

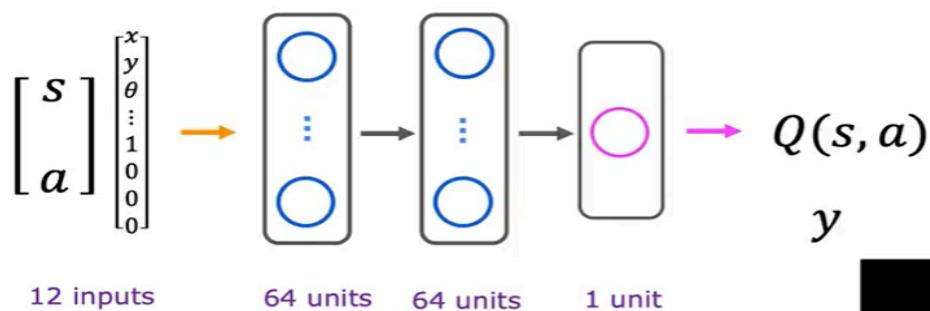
Scenario	State Variables
Self-driving car	Position, velocity, steering angle, acceleration
Robot arm	Joint angles, torque, velocity
Stock market	Price, volume, moving average

🧠 Why is it Challenging?

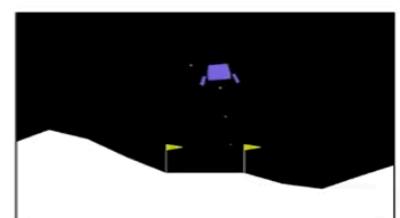
1. You can't store a value for every state, because there are infinitely many.
2. You can't use classic Q-tables like in small discrete environments.
3. You need function approximators (e.g., neural networks) to represent the value or policy.

Deep RL :

Deep Reinforcement Learning



In a state s , use neural network to compute
 $Q(s, \text{nothing}), Q(s, \text{left}), Q(s, \text{main}), Q(s, \text{right})$
 Pick the action a that maximizes $Q(s, a)$



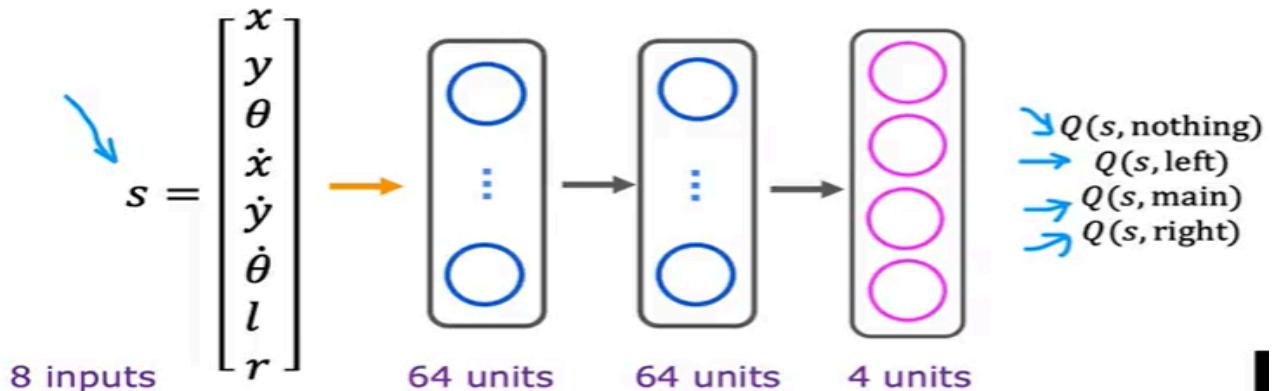
Deep Reinforcement Learning (Deep RL) combines **Reinforcement Learning (RL)** and **Deep Learning** to solve complex decision-making problems where traditional RL methods (like Q-tables) fall short.

What Is Deep RL?

Deep RL = Reinforcement Learning (learn by interacting with environment)
+ Deep Learning (learn from high-dimensional data using neural networks)

It allows agents to learn optimal behavior **from pixels, raw sensors, or complex inputs** — instead of hand-coded features or discrete states.

Deep Reinforcement Learning



In a state s , input s to neural network.

Pick the action a that maximizes $Q(s, a)$. $R(s) + \gamma \max_{a'} Q(s', a')$

Training Deep RL:

Bellman Equation

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

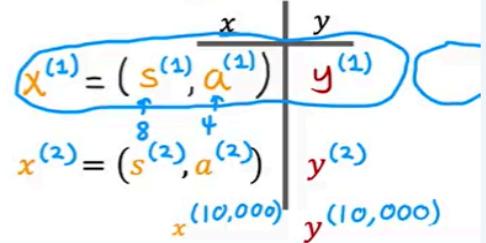
$$f_{w, B}(x) \approx y$$

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

$(s^{(1)}, a^{(1)}, R(s^{(1)}), s'^{(1)}) \leftarrow$
 $(s^{(2)}, a^{(2)}, R(s^{(2)}), s'^{(2)}) \leftarrow$
 $(s^{(3)}, a^{(3)}, R(s^{(3)}), s'^{(3)}) \leftarrow$

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

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



Here, This diagram illustrates how the **Bellman Equation** is used in **Deep Q-Learning** to update Q-values using neural networks.

Core Bellman Equation

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

- $Q(s, a)$: The value of taking action a in state s
- $R(s)$: Immediate reward
- γ : Discount factor
- s' : Next state
- a' : All possible next actions

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 $Q_{new}(s, a) \approx y$.

Set $Q = Q_{new}$.

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

$$x, y$$

$$x'', y''$$

:

$$x^{10,000}, y^{10,000}$$

