

Lecture 11: Unsupervised Learning

ADDO AI

2018

**Instructor: Dr. Syed Waqar ul
Qounain Jaffry
Assistant Professor PUCIT**



Unsupervised Learning

Agenda

- Introduction
- Prerequisite Self Check
- Un Supervised Learning
- Principal Component Analysis (PCA)
- K-Mean Clustering
- Gaussian Mixture Models
- Expectation Maximization Algorithm

Introduction

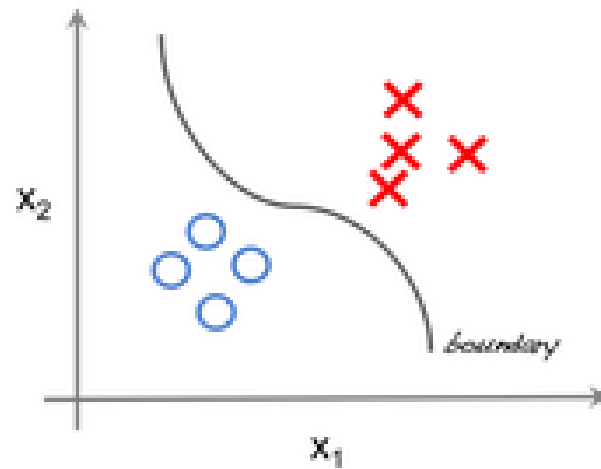
- Supervised Learning
 - This type of learning requires ground truth
 - Using the ground truth, machine learning models train to differentiate
 - Already have prior knowledge of what the output values should be
 - E.g. Classification and Regression
- If ground truth is not available
 - Data can still carry patterns
 - Then how to process data?
- Solution
 - Unsupervised Learning

Pre-requisite Self Check

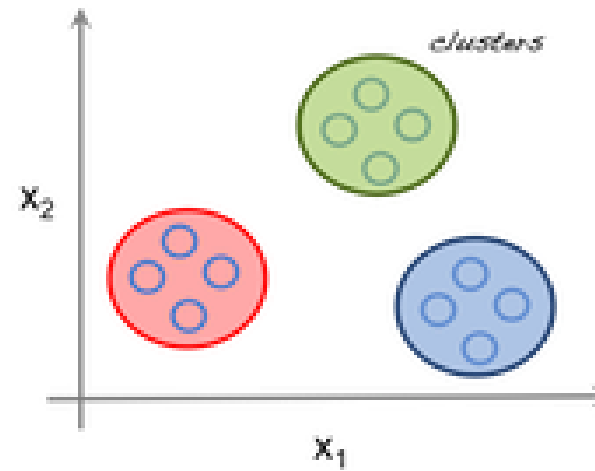
- Variance-Covariance-Matrix
- Langrange Multipliers
- Matrix Calculus $\frac{d x^T A x}{d x} =$
- Eigen Values of the Matrix
- Eigen Vectors of the Matrix
- Joint probability Multiplication rule: $P(AB)$

Supervised vs Unsupervised Learning

Supervised learning



Unsupervised learning



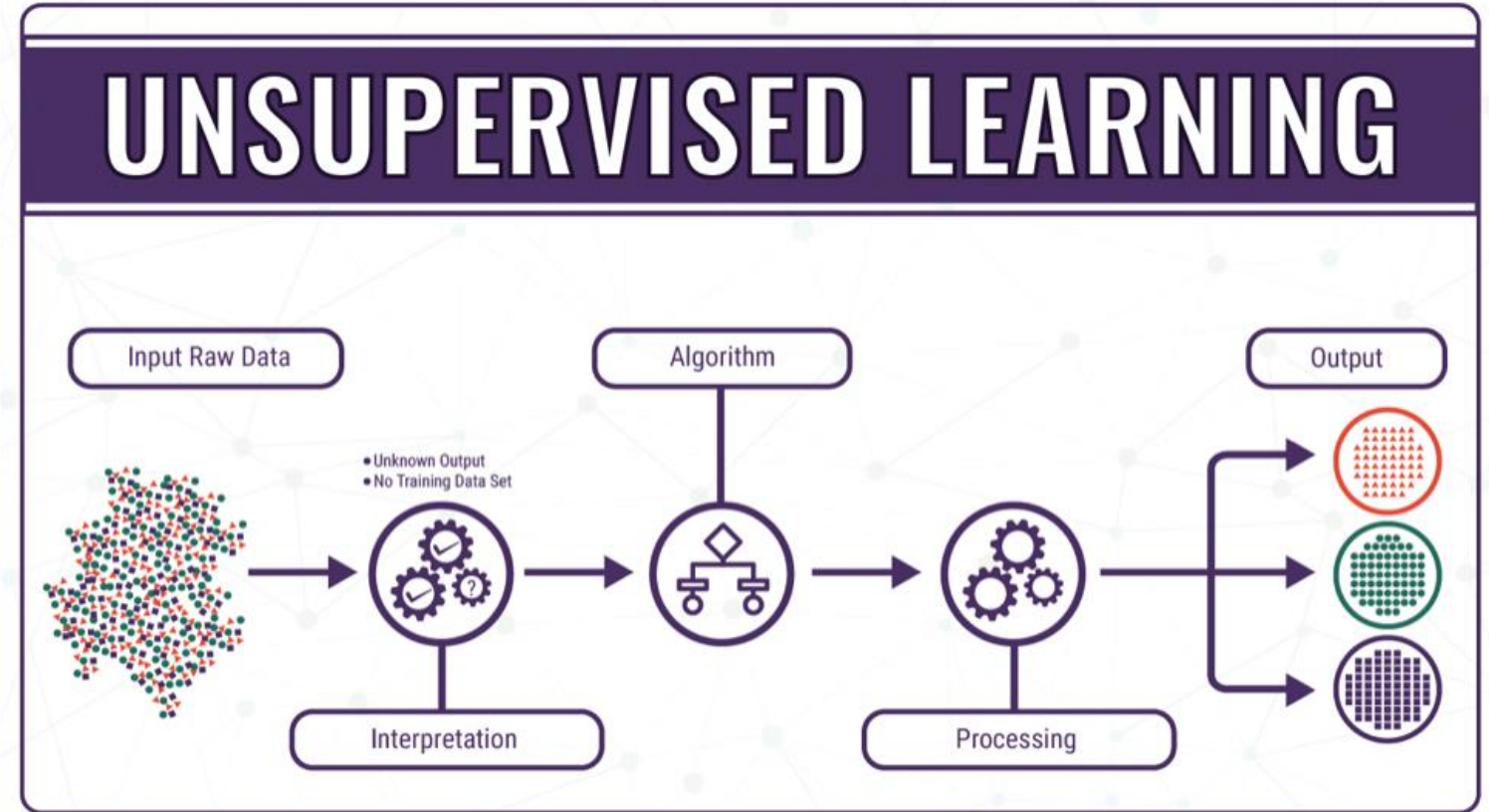
Unsupervised Learning

- It is focused on finding patterns in data
- Infer the natural structure present within data
- It is used for data exploration

	<i>Supervised Learning</i>	<i>Unsupervised Learning</i>
<i>Discrete</i>	classification or categorization	clustering
<i>Continuous</i>	regression	dimensionality reduction

Unsupervised Learning

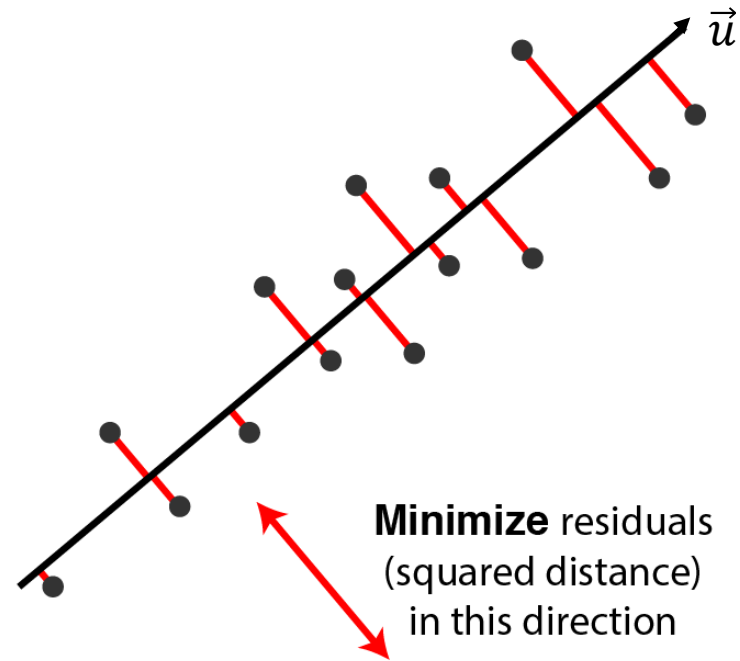
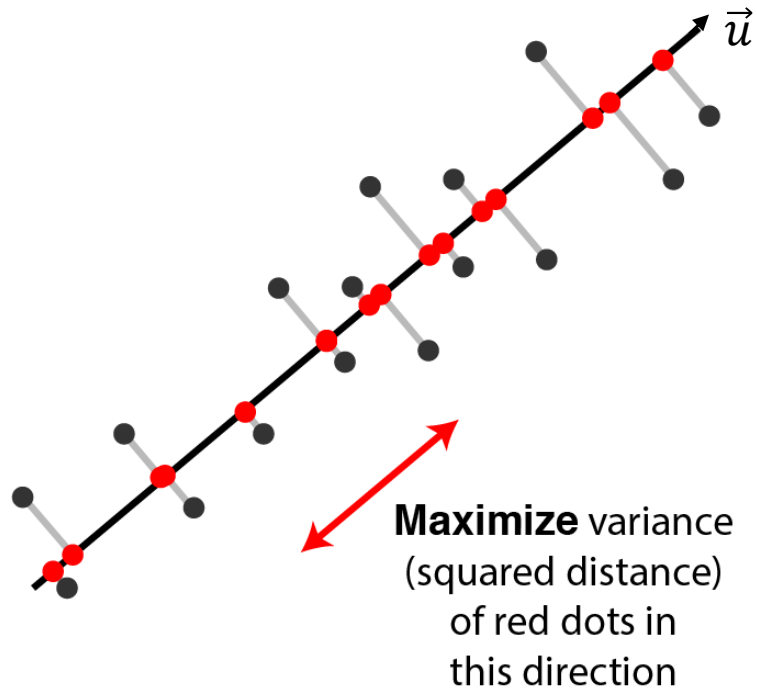
- Major applications
 - Dimensionality Reduction
 - PCA
 - Clustering
 - K-means
 - Density Estimation
 - Gaussian Mixture Models



Principal Component Analysis (PCA)

- Principal Component Analysis (PCA) is the primary example of linear unsupervised learning
- PCA is used for:
 - Dimensionality reduction
 - Feature extraction
 - Lossy data compression
- It is also called Karhunen-Loeve transform

Principal Component Analysis (PCA)



Maximum Variance Formulation

- Consider a set of data points
 - $X = [x_1, x_2, x_3, \dots, x_n]$ where each $x_i \in R_D$.
- Goal
 - To find a vector $\mathbf{u} \in R_D$ such that the variance of the projected data onto \mathbf{u} is maximum.
 - Projections of a data points x_i onto \mathbf{u} are obtained via dot-products $\mathbf{u}^T \mathbf{x}_i$ for $i = 1, 2, \dots, N$
 - Mean of projected data is computed as $\mathbf{u}^T \bar{\mathbf{x}}$ where $\bar{\mathbf{x}} = \sum_{i=1}^n \mathbf{x}_i$
 - Hence, variance of projected data along the direction \mathbf{u} is computed as
 - $Var(u) = \frac{1}{N} \sum_{i=1}^n (u^T x_i - u^T \bar{x})^2$

Maximum Variance Formulation

- $Var(u) = \frac{1}{N} \sum_{i=1}^n (u^T x_i - u^T \bar{x}) (u^T x_i - u^T \bar{x})^T$
- $Var(u) = \frac{1}{N} \sum_{i=1}^n (u^T x_i - u^T \bar{x}) (x_i^T u - \bar{x}^T u)$
- $Var(u) = \frac{1}{N} \sum_{i=1}^n u^T (x_i - \bar{x}) (x_i^T - \bar{x}^T) u$
- $Var(u) = u^T \frac{1}{N} \sum_{i=1}^n (x_i - \bar{x}) (x_i^T - \bar{x}^T) u$
- Say $S = \frac{1}{N} \sum_{i=1}^n (x_i - \bar{x}) (x_i^T - \bar{x}^T)$, that is a data-covariance matrix
- $Var(u) = u^T S u$

Maximum Variance Formulation

- Goal
 - To find a direction vector u that maximizes $u^T S u$
 - As we are interested in direction vector hence, constraint is applied
 - $u^T u = 1$ i.e. u is a unit vector
- To solve the optimization problem with constraint
 - Lagrang multipliers are used
 - Lagrang multiplier states that for a function $f(x)$ with given constraint $g(x)$
 - $L(x, \lambda) = f(x) - \lambda g(x)$ where $\lambda \neq 0$ is called lagrang multiplier
 - In order to find maximum value, this lagrangian function is derivated and equated to zero

Principal Component Analysis

- Currently
 - $f(\mathbf{u}) = \mathbf{u}^T S \mathbf{u}$
- Constraint
 - $\mathbf{u}^T \mathbf{u} = 1 \Rightarrow g(\mathbf{u}) = \mathbf{u}^T \mathbf{u} - 1$
- Hence, Lagrangian function becomes
 - $L(\mathbf{u}, \lambda) = f(\mathbf{u}) - \lambda g(\mathbf{u}) \Rightarrow \mathbf{u}^T S \mathbf{u} - \lambda(\mathbf{u}^T \mathbf{u} - 1)$
- To find optimal values, take gradient w.r.t \mathbf{u} :
 - $\frac{\partial L(\mathbf{u}, \lambda)}{\partial \mathbf{u}} = \frac{\partial}{\partial \mathbf{u}} (\mathbf{u}^T S \mathbf{u} - \lambda(\mathbf{u}^T \mathbf{u} - 1)) = 2S\mathbf{u} - 2\lambda\mathbf{u}$

Principal Component Analysis

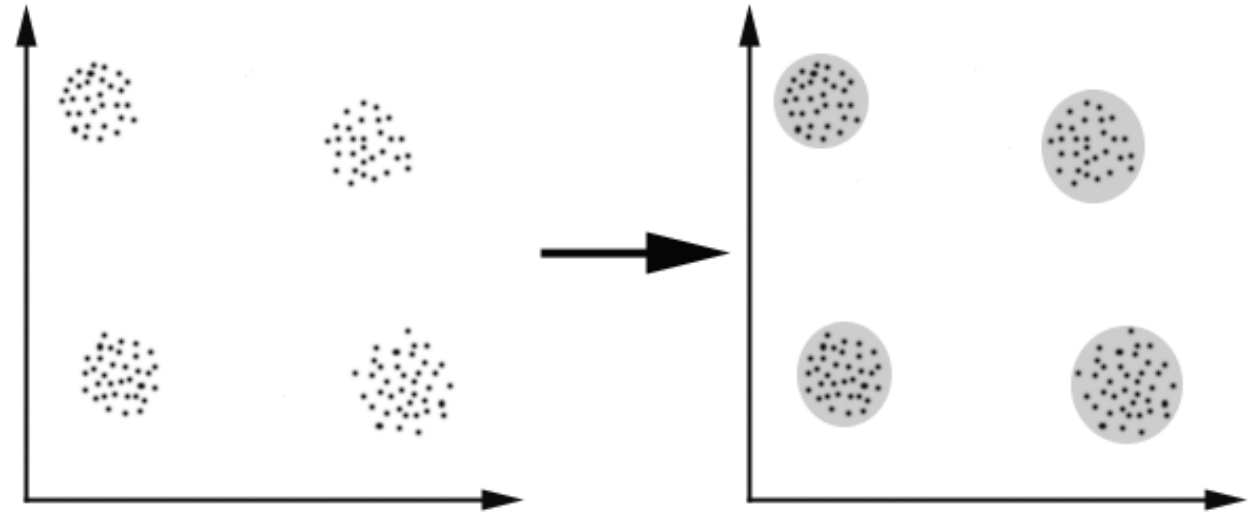
- Equation gradient to zero:

$$\Rightarrow 2S\mathbf{u} - 2\lambda\mathbf{u} = 0 \Rightarrow S\mathbf{u} = \lambda\mathbf{u}$$

- This should be recognizable as an eigenvector equation where \mathbf{u} is an eigenvector of S and λ is the associated eigenvalue.
- The eigenvector of S corresponding to the largest eigenvalue is called the first principal component.
- Additional principal components can be defined incrementally by choosing each new projection direction as the one with maximum projected variance among all directions orthogonal to those already considered .
- First M principal components correspond to the eigenvectors $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_M$ of S corresponding to the M largest eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_M$.

Clustering

- Unsupervised learning algorithm to identify groups or clusters of similar data points
- Come up with clusters such that:
 - Within cluster, points have minimum distance
 - Small intra-cluster distances
 - Between clusters, distance is large
 - Large inter-cluster distances



K-means clustering

- Can be seen as an instance of the more powerful framework of Expectation Maximization (to be covered later)
- Given data points $x_1, x_2, x_3, \dots, x_n$ where each $x_n \in R_D$, and an integer $K > 1$, the goal is to partition the data into K clusters .
- It performs mutually exclusive clustering
 - No data point can exist in multiple clusters
- Given K , randomly select K data points as clusters with each cluster mean denoted by μ_k
- At every step, calculate distance between remaining points with that of selected cluster.
 - For a data point x_i and any mean cluster point μ_k
 - the distance becomes $\|x_n - \mu_k\|^2$

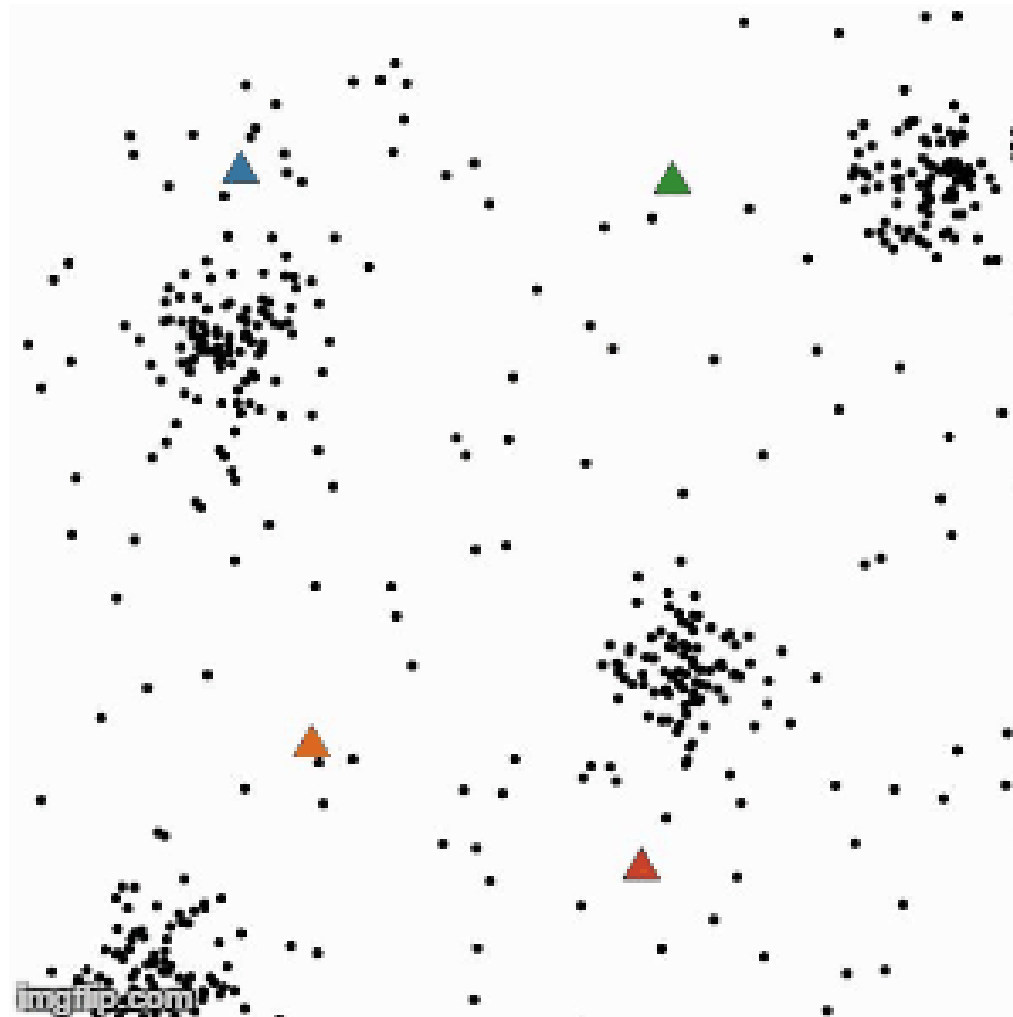
K-means clustering

- We also need a variable to denote assignment of any point x_n to the proper cluster.
- Let r_{nk} denote the assignment of data point x_n
- Then, as clusters to be formed are mutually exclusive:
 - Hence, a given data-point, can belong to only one cluster
 - This can be defined using hot-vector encoding or 1-of- K
 - Here $x_{nk} = 1$ if x_n belongs to cluster k and 0 otherwise
 - E.g. if a data point x_n belongs to cluster 2, whereas $K=3$:
 - respective vector would be [0 1 0]
 - Here $r_{n1} = r_{n3}=0$; as point x_n neither belongs to cluster 1 or 3, whereas $r_{n2}= 1$

K-means clustering

- Then for a particular set of clusters $\{\mu_1, \dots, \mu_k\}$ and cluster assignments $\{\mathbf{r}_1, \dots, \mathbf{r}_N\}$, where \mathbf{r}_k is a 1-of-K vector
- We can compute the sum-of-squared distances between data points and their assigned clusters as
 - $Distance(\mu_k, r_{ik}) = \sum_{k=1}^K \sum_{n=1}^N r_{nk} \|x_n - \mu_k\|^2$
- Goal:
 - $argminDistance(\mu_k, r_{nk})$
- Achieved via iterative, alternating optimization between assignments $\{r_{nk}\}$ and clusters $\{\mu_k\}$.

K-means clustering



K-means clustering

- Why alternating optimization?
- Finding cluster centers and cluster memberships simultaneously is a chicken-and-egg problem.
- However, individually these problems are much simpler.
 - Given memberships, computing cluster centers is trivial.
 - Given cluster centers, computing memberships is trivial.
- Alternating optimization gives us a powerful framework of solving complex problems by decomposing them into simpler ones .
- Notice that we appended the observed data x_i with some unobserved variables r_{ik} and then solved easy individual problems.
- These unobserved variables are called hidden or latent variables.

K-means clustering

- Data points $\{x_1, \dots, x_N\}$, integer $K > 1$
- Result : Cluster representatives $\{\mu_k\}$, assignments $\{r_{nk}\}$

Choose some initial μ_k ;

while not converged do

 Fix clusters and update assignments ($\{r_{nk}\} = \arg \min \{r_{nk}\} \text{ Distance}(\mu_k, r_{nk})$);

 Fix assignments and update clusters ($\{\mu_k\} = \arg \min \{\mu_k\} \text{ Distance}(\mu_k, r_{nk})$);

end

K-means clustering

- $\{r_{nk}\} = \arg \min \{r_{nk}\} \text{ Distance}(\mu_k, r_{nk})$
 - To obtain this, we have to ensure that $\|x_n - \mu_k\|^2$ is minimum
 - In other words, point x_n is placed in cluster that is closest.
- $r_{nk} = \begin{cases} 1 & \text{if } k = \operatorname{argmin}_j \|x_n - \mu_j\|^2 \\ 0 & \text{otherwise} \end{cases}$
- $\{\mu_k\} = \arg \min \{\mu_k\} \text{ Distance}(\mu_k, r_{nk})$
 - As distance function is quadratic in terms of μ_k
 - $\frac{\partial}{\partial \mu_k} \text{Distance}(\mu_k, r_{nk}) = \frac{\partial}{\partial \mu_k} \sum_{k=1}^K \sum_{i=1}^N r_{nk} \|x_n - \mu_k\|^2 = 0$
 - $\Rightarrow -2 \sum_{n=1}^N r_{nk} (x_n - \mu_k) = 0 \Rightarrow \mu_k = \frac{\sum_{n=1}^N r_{nk} x_n}{\sum_{n=1}^N r_{nk}}$

K-means clustering

- Requires value of K
- To determine value of K:
 - For various values of K , determine which value results in least error
- Hard decision
 - When points lie at cluster boundaries, hard decision is not also a good strategy
 - Use probability based soft assignments
 - Lead to mixture models

Gaussian Mixture Models

- Simple Gaussian Model

- Two parameters

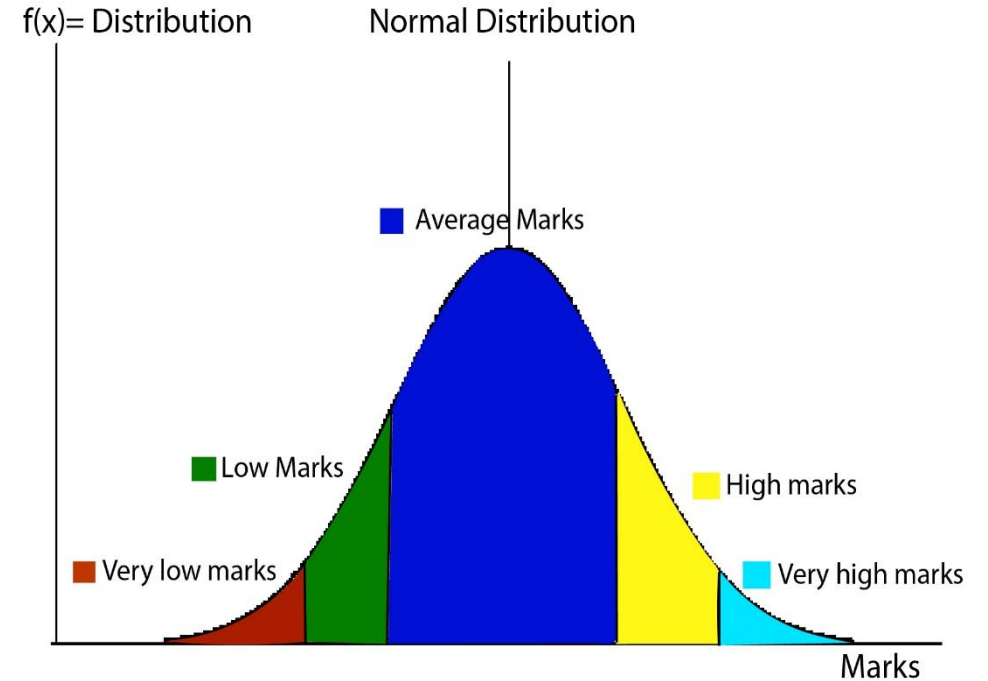
- Mean
 - Variance

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

- For a D-dimensional input vector \mathbf{x} :

- $N(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{\sqrt{|\boldsymbol{\Sigma}|}} e^{\frac{-1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})}$

- Where $\boldsymbol{\Sigma}$ is called the covariance matrix that has D*D dimensions
 - $|\boldsymbol{\Sigma}|$ is determinant and $\boldsymbol{\mu}$ is the D-dimensional mean vector

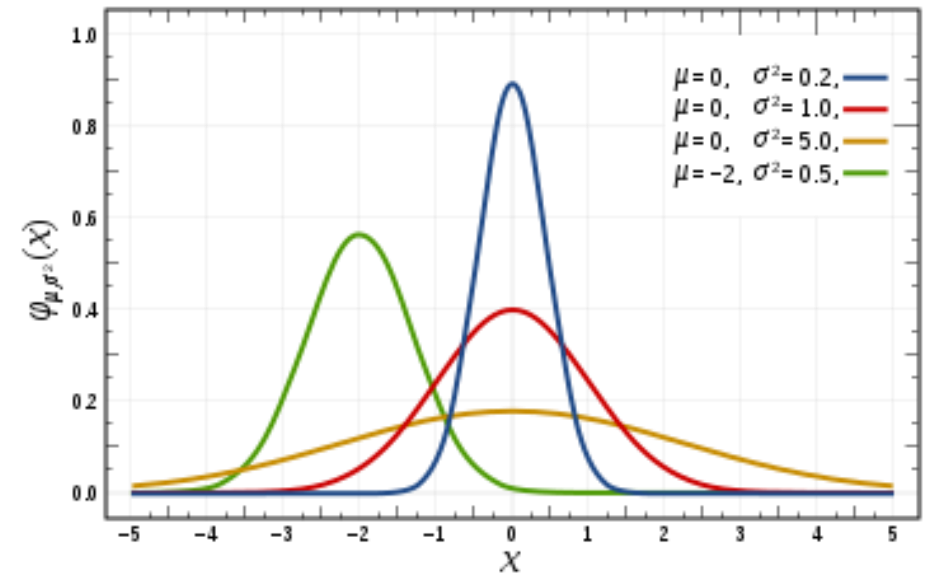


Gaussian Mixture Models (GMM)

- Multi model data
 - Can't be modeled using Simple Normal distribution/Uni-Modal Gaussian
 - Use mixture of models, that is linear superpositions of Uni-modal Gaussians

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k N(\mathbf{x} | \boldsymbol{\mu}_k, \Sigma_k)$$

- Here, π_k refers to mixing coefficients and they must satisfy
 - $0 < \pi_k < 1$
 - $\sum_{k=1}^K \pi_k = 1$



Gaussian Mixture Models (GMM)

Latent Variables

- Now, let's derive the $p(\mathbf{x})$ using latent variables
- Like K-means, let's assume a latent variable \mathbf{z} with 1-of-K encoding
- $p(\mathbf{x}) = \sum_{k=1}^K p(\mathbf{x}, z_k) = \sum_{k=1}^K p(z_k)p(\mathbf{x}|z_k)$
- As, now, soft-assignments are being made in form of probabilities:
 - $p(z_k = 1) = \pi_k$
- Due to 1-of-K encoding of latent variable i.e. vector \mathbf{z} :
 - $p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k}$ # equivalent to $p(z_k = 1)$
- Similarly, the conditional distribution of \mathbf{x} given particular value of \mathbf{z} is:
 - $p(\mathbf{x}|z_k = 1) = N(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$

Gaussian Mixture Models (GMM)

Latent Variables

- Using whole vector \mathbf{z} notation, it can be written as:
 - $p(\mathbf{x}|\mathbf{z}) = \prod_{k=1}^K N(\mathbf{x}|\boldsymbol{\mu}_k, \Sigma_k)^{z_k}$ # equivalent to $p(\mathbf{x}|z_k = 1)$
- Joint distribution then becomes:
- $p(\mathbf{x}) = \sum_{k=1}^K p(z_k)p(\mathbf{x}|z_k) = \sum_{k=1}^K \pi_k N(\mathbf{x}|\boldsymbol{\mu}_k, \Sigma_k)$
- Here, using a latent variable, we have derived the similar expression for GMM

Gaussian Mixture Models (GMM)

Responsibilities

- $p(\mathbf{x})$ is the marginal density that we are looking to model.
- $p(\mathbf{x}|z_k = 1)$ is the component conditional density . That is, probability density of \mathbf{x} according to component k .
- $p(z_k = 1)$ is the prior probability of component k .
- $p(z_k = 1 | \mathbf{x})$ is the posterior probability of component k .
 - It can be represented as responsibility r_k
 - Can be viewed as the responsibility that component k takes for explaining observation \mathbf{x} .
 - Can be computed via Bayes' theorem

$$r_k = p(z_k = 1 | \mathbf{x}) = \frac{p(z_k=1) p(\mathbf{x}|z_k=1)}{p(\mathbf{x})} = \frac{p(z_k=1) p(\mathbf{x}|z_k=1)}{\sum_{j=1}^K p(z_j)p(\mathbf{x}|z_j)} = \frac{\pi_k N(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j N(\mathbf{x} | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}$$

Gaussian Mixture Models (GMM)

Parameter Estimation

- Assuming independent and identically data distribution (IID)
 - It enables to express joint probability as product of marginals
 - $p(x,y) = p(x) p(y)$ (with iid) $p(x,y) = p(y|x) p(x) = p(x|y)p(y)$ (w/o iid)
- Likelihood of any model M with data is given by:
 - $L(D|\theta) = \prod_{n=1}^N p(x_n|\theta)$

where θ represents parameters of model and D represents input data $[x_1, x_2, \dots x_n]$

- As likelihood and its log would present the same pattern, hence, for simplicity, log-likelihood is used
 - $\text{Log}_L(D|\theta) = \text{Log}(\prod_{n=1}^N p(x_n|\theta)) = \sum_{n=1}^N \log(p(x_n|\theta))$

Gaussian Mixture Models (GMM)

Parameter Estimation

- As $p(\mathbf{x}) = \sum_{k=1}^K \pi_k N(\mathbf{x}|\boldsymbol{\mu}_k, \Sigma_k)$
- By taking likelihood of GMM:
 - $L(D|\theta) = \prod_{n=1}^N (\sum_{k=1}^K \pi_k N(\mathbf{x}_n|\boldsymbol{\mu}_k, \Sigma_k))$
- Then log-likelihood becomes:
 - $\text{Log}L(D|\theta) = \sum_{n=1}^N \log((\sum_{k=1}^K \pi_k N(\mathbf{x}_n|\boldsymbol{\mu}_k, \Sigma_k)))$
- Derivate w.r.t mean $\boldsymbol{\mu}_k$
- $\frac{\partial \text{Log}L}{\partial \boldsymbol{\mu}_k} = \frac{\partial}{\partial \boldsymbol{\mu}_k} (\sum_{n=1}^N \log((\sum_{k=1}^K \pi_k N(\mathbf{x}_n|\boldsymbol{\mu}_k, \Sigma_k)))) = \mathbf{0}$
- $\frac{\partial \text{Log}L}{\partial \boldsymbol{\mu}_k} = -2 \sum_{n=1}^N \left(\frac{\pi_k N(\mathbf{x}_n|\boldsymbol{\mu}_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(\mathbf{x}_n|\boldsymbol{\mu}_j, \Sigma_j)} \right) \Sigma_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) = \mathbf{0}$
- $\frac{\partial \text{Log}L}{\partial \boldsymbol{\mu}_k} = \sum_{n=1}^N r_{nk} \Sigma_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) = \mathbf{0}$

Gaussian Mixture Models (GMM)

Parameter Estimation

- $\Rightarrow \sum_{n=1}^N r_{nk} \Sigma_k^{-1} \mathbf{x}_n = \sum_{n=1}^N r_{nk} \Sigma_k^{-1} \boldsymbol{\mu}_k$
- By multiplying both sides by Σ_k , equation yields :
 - $\boldsymbol{\mu}_k = \frac{\sum_{n=1}^N r_{nk} \mathbf{x}_n}{\sum_{n=1}^N r_{nk}} = \frac{1}{N_k} \sum_{n=1}^N r_{nk} \mathbf{x}_n$ where $N_k = \sum_{n=1}^N r_{nk}$
- This expression is the same as derived using K-means distance.
- Only difference is that in k-means, r_{nk} could either be 1 or zero, whereas here, it is probability based soft-assignments.

Gaussian Mixture Models (GMM)

Parameter Estimation

- Similarly, log-likelihood function can be differentiated with respect to remaining parameters that include π_k and Σ_k
- $$\frac{\partial \text{Log} L}{\partial \Sigma_k} = \frac{\partial}{\partial \Sigma_k} \left(\sum_{n=1}^N \log \left(\sum_{k=1}^K \pi_k N(\mathbf{x}_n | \boldsymbol{\mu}_k, \Sigma_k) \right) \right) = 0$$
- It yields:
 - $$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T$$
- For π_k value, as it must follow a constraint
 - Lagrangian function will be used
 - $f(\mathbf{x}) = \text{log-likelihood function}$ $g(\mathbf{x}) : \sum_{k=1}^K \pi_k = 1 \Rightarrow \sum_{k=1}^K \pi_k - 1 = 0$

Gaussian Mixture Models (GMM)

Parameter Estimation

- Expression becomes:

- $L(x, \lambda) = \text{Log}L(D|\boldsymbol{\pi}, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) + \lambda(\sum_{k=1}^K \pi_k - 1) = 0$

- $\frac{\partial}{\partial \pi_k} L(x, \lambda) = \frac{\partial}{\partial \pi_k} (\sum_{n=1}^N \log(\sum_{k=1}^K \pi_k N(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)))$
 $+ \frac{\partial}{\partial \pi_k} [\lambda(\sum_{k=1}^K \pi_k - 1)] = 0$

- $\Rightarrow \sum_{n=1}^N \left(\frac{N(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j N(\mathbf{x}_n|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \right) + \frac{\partial}{\partial \pi_k} [\lambda(\pi_0 + \dots + \pi_k + \dots + \pi_K)]$

- $\Rightarrow \sum_{n=1}^N \left(\frac{r_{nk}}{\pi_k} \right) + \lambda = 0$ $\#r_{nk} = \frac{\pi_k N(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j N(\mathbf{x}_n|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \Rightarrow \frac{N(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j N(\mathbf{x}_n|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} = \frac{r_{nk}}{\pi_k}$

Gaussian Mixture Models (GMM)

Parameter Estimation

- Multiply both sides by π_k
 - $\sum_{n=1}^N (r_{nk}) + \lambda \pi_k = 0 \dots \dots \dots (1)$
 - $\Rightarrow \sum_{n=1}^N (r_{nk}) = -\lambda \pi_k$
- Taking summation over k on both sides:
 - $\Rightarrow \sum_{k=1}^K \sum_{n=1}^N (r_{nk}) = -\lambda \sum_{k=1}^K \pi_k$
 - $\Rightarrow \sum_{k=1}^K \sum_{n=1}^N (r_{nk}) = -\lambda$
 - $\Rightarrow \lambda = -N$
- By putting value of λ in equation 1:
 - $\sum_{n=1}^N (r_{nk}) - N \pi_k = 0 \dots \dots \dots (1)$
 - $N_k = N \pi_k \Rightarrow \pi_k = \frac{N_k}{N}$
- In words, mixing coefficient for component k is given by the average responsibility that it takes for explaining the training data points .

Gaussian Mixture Models (GMM)

Parameter Estimation

- Notice that solutions for $\pi_k, \boldsymbol{\mu}_k, \Sigma_k$ are dependent on the responsibilities r_{nk} .
- However, the responsibilities depend on $\pi_k, \boldsymbol{\mu}_k, \Sigma_k$
- We can now present the alternating optimization algorithm for GMMs

Gaussian Mixture Models (GMM)

Alternative optimization algorithms

- Data points $\{x_1, \dots, x_N\}$, integer $K > 1$
- Result : Component Parameters $\{\mu_k, \Sigma_k\}$, mixing coefficients $\{\pi_k\}$

Choose some initial values for μ_k, Σ_k, π_k ;

1. Fix parameters and find $r_{nk} = \frac{\pi_k N(x_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(x_n | \mu_j, \Sigma_j)}$
2. Fix responsibilities, and update parameters
 - $\mu_k' = \frac{1}{N_k} \sum_{n=1}^N r_{nk} x_n$ **where** $N_k = \sum_{n=1}^N r_{nk}$
 - $\Sigma_k' = \frac{1}{N_k} \sum_{n=1}^N r_{nk} (x_n - \mu_k')(x_n - \mu_k')^T$
 - $\pi_k = \frac{N_k}{N}$
3. Evaluate Log-likelihood:
 - $\text{Log}L(D | \pi, \mu_k, \Sigma_k) = \sum_{n=1}^N \log(\sum_{k=1}^K \pi_k N(x_n | \mu_k, \Sigma_k))$
4. Check for convergence, if not converged, go back to step 1 else stop.

Gaussian Mixture Models (GMM)

- Each iteration increases (or retains) the value of the log-likelihood.
 - Therefore, convergence to (local) maximum is guaranteed.
- This algorithm has a name – Expectation Maximization (EM)
- GMM converges slower than K-means and performs more computations per-iteration.
- Singularity issue:
 - If a component gets mapped to a training data point then:
 - $$\frac{1}{(2\pi)^{D/2}} \frac{1}{\sqrt{|\Sigma|}} e^{\frac{-1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)} = \frac{1}{(2\pi)^{D/2}} \frac{1}{\sqrt{|\Sigma|}}$$
 - Here if $|\Sigma|$ approaches to zero: this term can get infinite

Gaussian Mixture Models (GMM)

- Solution
 - Care must be taken to check if that has happened or is close to happening.
 - If so, the collapsing component should be reset to some other randomly chosen μ_k and large Σ_k
 - Other part of the optimization algorithm should proceed as before
- GMM is used for performing density estimation as well.

Expectation Maximization Algorithm

- K-means and GMMs are examples of latent variable models.
- Specifically for GMMs, we have seen an incremental algorithm for learning the parameters using Machine Learning.
- That algorithm is actually an instance of a powerful framework called Expectation-Maximization (EM) .
- EM is used for solving latent variable problems using Machine Learning.
- We will now present a more general explanation of the EM algorithm.

Expectation Maximization Algorithm

- Maximum likelihood is equivalent to maximizing the log-likelihood $\text{Log}L(D|\boldsymbol{\theta})$ where $\boldsymbol{\theta}$ represents all parameters of a model.
- Using the sum-rule $\text{Log}L(D|\boldsymbol{\theta}) = \sum_{n=1}^N \log(\sum_{k=1}^K p(x_n|\theta_k))$
- Maximization is no longer straight-forward since log is 'blocked' by the summation.
- So we take another approach

Expectation Maximization Algorithm

- We will denote $\{x, z\}$ as the complete dataset.
 - Z represents latent variables
- We will denote $\{x\}$ as the incomplete dataset.
- The goal now is to maximize the complete-data log-likelihood function $p((x, z) | \theta)$.
- But for that we need to know the values of Z which are unobserved.
 - What can be computed about Z , however, is the posterior $p(z | x, \theta)$.
 - So instead of the incomputable log-likelihood, the next best computable number would be its expected-value under the posterior $p(z | x, \theta)$

Expectation Maximization Algorithm

- This yields the E-step of the EM algorithm.
 - $E_{z|x, \theta_{old}} \log p(x, z|\theta) = Q(\theta, \theta_{old}) = \sum_{k=1}^K p(z|x, \theta_{old}) \log p(x, z|\theta)$
 - $= \sum_{k=1}^K p(z|x, \theta_{old}) \log p(z|\theta) p(x|z, \theta)$
- Since we are eventually interested in optimal parameters θ_{new} we treat this expectation as a function of θ and denote it by $Q(\theta, \theta_{old})$.
- The M-step corresponds to maximizing this expectation
 - $\theta_{new} = \operatorname{argmax}_{\theta} Q(\theta, \theta_{old})$
- In short, EM replaces the log-likelihood by the expected log-likelihood and maximizes it.

Expectation Maximization Algorithm

- Goal is to maximize likelihood $p(z|x, \theta)$ with respect to θ by introducing joint distribution $p(x, z | \theta)$ involving latent variables Z .
- Steps:
 1. Choose initial θ_{old}
 2. E-step: Evaluate $p(z|x, \theta_{old})$ i.e. responsibilities by using θ_{old}
 3. M-step : Obtain new estimate θ_{new} by maximizing the expectation $Q(\theta, \theta_{old})$
 - i. $\theta_{new} = \operatorname{argmax}_{\theta} Q(\theta, \theta_{old})$
 4. Check for convergence of either log-likelihood or parameters. If not converged, then
 - i. $\theta_{old} \leftarrow \theta_{new}(1)$
 - ii. return to step 2.