

Lecture 14: Unsupervised Learning

Introduction to Machine Learning [25737]

Sajjad Amini

Sharif University of Technology

1 Approach Definition

2 Principle Component Analysis

- Interpretation Via Maximum Projection Spread
- Interpretation Via Reconstruction

3 Clustering

4 Mixture Models

The material in the slides except cited are inspired from the following reference:

- Murphy, K. P. (2022). *Probabilistic machine learning: an introduction*. MIT press.

Section 1

Approach Definition

Principle Component Analysis

- Experience E : Set of N samples $\mathcal{D} = \{\mathbf{x}_n\}_{n=1}^N$
- Task T : Projecting data into low dimensional subspace which captures its main aspects
- Performance measure: Preserving data variations

Clustering

- Experience E : Set of N samples $\mathcal{D} = \{\mathbf{x}_n\}_{n=1}^N$
- Task T : Partition the input into regions that contains *similar* points.
- Performance measure in *Compression*: Compression loss

Section 2

Principle Component Analysis

Subsection 1

Interpretation Via Maximum Projection Spread

Data Matrix

- Assume $\mathbf{x} \in \mathbb{R}^D$ is a random variable and you have observed N copies of it as $\{\mathbf{x}_i\}_{i=1}^N$ (Equivalently the dataset).
- As before, we stack these copies into a Matrix \mathbf{X} as:

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \dots \\ \mathbf{x}_N^T \end{bmatrix} = \begin{bmatrix} x_1^1 & \dots & x_1^D \\ x_2^1 & \dots & x_2^D \\ \dots & & \dots \\ x_N^1 & \dots & x_N^D \end{bmatrix} \in \mathbb{R}^{N \times D}$$

- Each column is a feature (covariate or predictor)
- Each row is an observation

Characterizing Dataset

The dataset point create a point cloud in \mathbb{R}^D space.

- The expectation of this point cloud, calculated below, determines the center of point cloud.

$$\mathbb{E}[\mathbf{x}] = \begin{bmatrix} \mathbb{E}[x^1] \\ \vdots \\ \mathbb{E}[x^D] \end{bmatrix}$$

- The covariance matrix of this point cloud, calculated below, determines the spread of point cloud.

$$\begin{aligned} \text{Cov}[\mathbf{x}] &\triangleq \mathbb{E} [(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^T] = \mathbb{E}[\mathbf{x}\mathbf{x}^T] - \mathbb{E}[\mathbf{x}]\mathbb{E}[\mathbf{x}]^T = \mathbf{\Sigma} \\ &= \begin{bmatrix} \text{Cov}[X_1, X_1] & \text{Cov}[X_1, X_2] & \cdots & \text{Cov}[X_1, X_D] \\ \text{Cov}[X_2, X_1] & \text{Cov}[X_2, X_2] & \cdots & \text{Cov}[X_2, X_D] \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}[X_D, X_1] & \text{Cov}[X_D, X_2] & \cdots & \text{Cov}[X_D, X_D] \end{bmatrix} \end{aligned}$$

Utilizing Empirical Distribution

The empirical distribution for dataset $\{\mathbf{x}_i\}_{i=1}^N$ is defined as:

$$p_D(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^N \delta(\mathbf{x} - \mathbf{x}_n)$$

We can use it to compute the empirical (sample) mean and empirical (sample) covariance matrix as:

$$\mathbb{E}_D[\mathbf{x}] = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i = \bar{\mathbf{x}}$$
$$\text{Cov}_D[\mathbf{x}] = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T - \bar{\mathbf{x}} \bar{\mathbf{x}}^T = \mathbf{S}$$

Eliminating Summation Using Linear Algebra

For sample mean, we have:

$$\bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i = \frac{1}{N} \begin{bmatrix} | & & | \\ \mathbf{x}_1 & \dots & \mathbf{x}_N \\ | & & | \end{bmatrix}_{D \times N} \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}_{N \times 1} = \frac{1}{N} \mathbf{X}^T \mathbf{1}$$

For sample covariance matrix, we use forth method for matrix multiplication as:

$$\begin{aligned} \mathbf{S} &= \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T - \bar{\mathbf{x}} \bar{\mathbf{x}}^T = \frac{1}{N} \begin{bmatrix} | & & | \\ \mathbf{x}_1 & \dots & \mathbf{x}_N \\ | & & | \end{bmatrix} \begin{bmatrix} - & \mathbf{x}_1^T & - \\ & \vdots & \\ - & \mathbf{x}_N^T & - \end{bmatrix} - \bar{\mathbf{x}} \bar{\mathbf{x}}^T \\ &= \frac{1}{N} \mathbf{X}^T \mathbf{X} - \frac{1}{N^2} \mathbf{X}^T \mathbf{1} \mathbf{1}^T \mathbf{X} = \frac{1}{N} \mathbf{X}^T \underbrace{\left(\mathbf{I} - \frac{1}{N} \mathbf{1} \mathbf{1}^T \right)}_{\mathbf{H}} \mathbf{X} \end{aligned}$$

Characterizing Dataset [1]

Idempotent Matrix

Matrix \mathbf{P} is said to be idempotent matrix if $\mathbf{P}^2 = \mathbf{P}$

Projection Matrix

Matrix \mathbf{Z} is said to be projection matrix if it is symmetric and idempotent.

Working on \mathbf{H} Matrix

Matrix \mathbf{H} is a projection matrix because:

- $\mathbf{H}^T = (\mathbf{I} - \frac{1}{N}\mathbf{1}\mathbf{1}^T)^T = \mathbf{I} - \frac{1}{N}\mathbf{1}\mathbf{1}^T$
- Idempotent property:

$$\begin{aligned}\mathbf{H}^2 = \mathbf{H}\mathbf{H} &= (\mathbf{I} - \frac{1}{N}\mathbf{1}\mathbf{1}^T)(\mathbf{I} - \frac{1}{N}\mathbf{1}\mathbf{1}^T) = \mathbf{I} - \frac{2}{N}\mathbf{1}\mathbf{1}^T + \frac{1}{N^2}\mathbf{1}\overbrace{\mathbf{1}^T\mathbf{1}}^N\mathbf{1}^T \\ &= \mathbf{I} - \frac{2}{N}\mathbf{1}\mathbf{1}^T + \frac{1}{N}\mathbf{1}\mathbf{1}^T = \mathbf{H}\end{aligned}$$

Characterizing the Projection \mathbf{H}

Assume $\mathbf{v} \in \mathbb{R}^N$, then:

$$\mathbf{H}\mathbf{v} = \mathbf{v} - \frac{1}{N}\mathbf{1}\mathbf{1}^T\mathbf{v} = \mathbf{v} - \frac{\mathbf{1}^T\mathbf{v}}{N}\mathbf{1} = \mathbf{v} - \bar{v}\mathbf{1}$$

Thus \mathbf{H} removes the mean of the vector from each coordinate. Equivalently $\overline{\mathbf{H}\mathbf{v}} = \mathbf{0}$

Thus \mathbf{H} is the projection onto the subspace of vectors with zero mean (Projection onto hyperplane which is orthogonal to $\mathbf{1}$ vector).

Re-writing \mathbf{S}

Based on the projection matrix \mathbf{H} , we have:

$$\mathbf{S} = \frac{1}{N}\mathbf{X}^T\mathbf{H}\mathbf{X} = \frac{1}{N}\mathbf{X}^T\mathbf{H}^2\mathbf{X} = \frac{1}{N}\mathbf{X}^T\mathbf{H}^T\mathbf{H}\mathbf{X} = \frac{1}{N}(\mathbf{H}\mathbf{X})^T(\mathbf{H}\mathbf{X})$$

where $\mathbf{H}\mathbf{X}$ result in centered features.

Linear Combination of Features [1]

Original Formulation

Assume an arbitrary direction of $\mathbf{u} \in \mathbb{R}^D$, then consider the following value:

$$\begin{aligned}\mathbf{u}^T \Sigma \mathbf{u} &= \mathbf{u}^T [\mathbb{E}[\mathbf{x}\mathbf{x}^T] - \mathbb{E}[\mathbf{x}]\mathbb{E}[\mathbf{x}^T]] \mathbf{u} \stackrel{(a)}{=} \mathbb{E}[(\mathbf{u}^T \mathbf{x})(\mathbf{u}^T \mathbf{x})^T] - \mathbb{E}[\mathbf{u}^T \mathbf{x}]\mathbb{E}[(\mathbf{u}^T \mathbf{x})^T] \\ &= \mathbb{E}[(\mathbf{u}^T \mathbf{x})^2] - \mathbb{E}[\mathbf{u}^T \mathbf{x}]^2 = \text{var}(\mathbf{u}^T \mathbf{x})\end{aligned}$$

Switching to Empirical Distribution

Using empirical distribution, we have the empirical variance for $\{\mathbf{u}^T \mathbf{x}_i\}_{i=1}^N$:

$$\begin{aligned}\mathbf{u}^T \mathbf{S} \mathbf{u} &= \mathbf{u}^T \left[\frac{1}{N} \sum_{i=1}^N \mathbf{x}_i^T \mathbf{x}_i - \bar{\mathbf{x}} \bar{\mathbf{x}}^T \right] \mathbf{u} = \frac{1}{N} \sum_{i=1}^N (\mathbf{u}^T \mathbf{x}_i)(\mathbf{x}_i^T \mathbf{u}) - (\mathbf{u}^T \bar{\mathbf{x}})(\bar{\mathbf{x}}^T \mathbf{u}) \\ &= \frac{1}{N} \sum_{i=1}^N (\mathbf{u}^T \mathbf{x}_i)^2 - (\mathbf{u}^T \bar{\mathbf{x}})^2 = \frac{1}{N} \sum_{i=1}^N (\mathbf{u}^T \mathbf{x}_i)^2 - (\overline{\mathbf{u}^T \mathbf{x}})^2 = s^2\end{aligned}$$

Intuition Behind Principle Component Analysis [1]

Intuition

Finding the direction \mathbf{u} which result in the high projection value spread measured by project value variance

Extreme Cases

- Zero variance: The projection of all points onto \mathbf{u} is equal (The points are in the hyper-plane whose normal vector is \mathbf{u}).
- Large variance: The points are spread along the \mathbf{u} direction.

Objective

Fining the direction that maximize the projection variance (or equivalently projection spread)

Formulation

The problem for PCA can be formulated as:

$$\max_{\mathbf{u} \in \mathbb{R}^D} \mathbf{u}^T \mathbf{S} \mathbf{u}$$

The maximum value of objective function is infinity, thus we need to constrained \mathbf{u} as:

$$\max_{\mathbf{u} \in \mathbb{R}^D} \mathbf{u}^T \mathbf{S} \mathbf{u} \quad \text{subject to} \quad \|\mathbf{u}\|_2 = 1$$

Spectral Theorem

Eigenvalues and Eigenvectors of Symmetric Matrices

Based on *Spectral Theorem*, for symmetric matrix \mathbf{S} we have:

- All eigenvalues are real
- Eigenvectors are orthonormal (\mathbf{U} is orthogonal thus $\mathbf{P}^{-1} = \mathbf{P}^T$)

Then we have:

$$\begin{aligned}\mathbf{S} = \mathbf{P}\mathbf{\Lambda}\mathbf{P}^T &= \begin{bmatrix} | & | & | \\ \mathbf{p}_1 & \mathbf{p}_2 & \mathbf{p}_n \\ | & | & | \end{bmatrix} \begin{bmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \ddots \\ & & & \lambda_n \end{bmatrix} \begin{bmatrix} - & \mathbf{p}_1^T & - \\ - & \mathbf{p}_2^T & - \\ & \vdots & \\ - & \mathbf{p}_m^T & - \end{bmatrix} \\ &= \sum_{i=1}^n \lambda_i \mathbf{p}_i \mathbf{p}_i^T\end{aligned}$$

Covariance Matrices

Covariance matrices are positive semi-definite, equivalently, all their eigenvalues are non-negative ($\mathbf{u}^T \mathbf{\Sigma} \mathbf{u} \geq 0, \forall \mathbf{u}$ and $\mathbf{u}^T \mathbf{S} \mathbf{u} \geq 0, \forall \mathbf{u}$).

Characterizing S

Characterizing S

Using Spectral theorem, we can write S as:

$$S = P\Lambda P^T, \quad \left\{ \begin{array}{l} P = \begin{bmatrix} | & | & & | \\ \mathbf{p}_1 & \mathbf{p}_2 & & \mathbf{p}_D \\ | & | & & | \end{bmatrix} \\ \Lambda = \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \dots & \\ & & & \lambda_D \geq 0 \end{bmatrix}, \lambda_1 \geq \lambda_2 \dots \lambda_D \\ P^T P = I \end{array} \right.$$

Transforming Using Eigenvectors

Assume we define $\mathbf{y} = \mathbf{P}^T \mathbf{x} \in \mathbb{R}^D$ and $\bar{\mathbf{x}} = \mathbf{0}$, then:

$$\bar{\mathbf{y}} = \overline{\mathbf{P}^T \mathbf{x}} = \mathbf{P}^T \bar{\mathbf{x}} = \mathbf{0}$$

Thus the sample covariance matrix for \mathbf{y} is:

$$\mathbf{S}^y = \frac{1}{N} \sum_{i=1}^N (\mathbf{P}^T \mathbf{x}_i)(\mathbf{P}^T \mathbf{x}_i)^T = \mathbf{P}^T \left(\frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T \right) \mathbf{P} = \mathbf{P}^T \mathbf{S} \mathbf{P} = \mathbf{D}$$

Thus we take one step through whitening:

$$\text{cov}(Y^i, Y^j) = \begin{cases} 0 & i \neq j \\ \lambda_i & i = j \end{cases}$$

Finding Maximum Spread Direction [1]

Finding Maximum Spread Direction

Assume the maximum spread direction is \mathbf{u} and consider the following definition:

$$\mathbf{b} = \mathbf{P}^T \mathbf{u} \Rightarrow \mathbf{u} = \mathbf{P} \mathbf{b}$$

Now we measure the spread as:

$$\mathbf{u}^T \mathbf{S} \mathbf{u} = (\mathbf{P} \mathbf{b})^T (\mathbf{P} \mathbf{D} \mathbf{P}^T) (\mathbf{P} \mathbf{u}) = \mathbf{b}^T \overbrace{(\mathbf{P} \mathbf{P}^T)}^{\mathbf{I}} \mathbf{D} \overbrace{(\mathbf{P}^T \mathbf{P})}^{\mathbf{I}} \mathbf{b} = \sum_{j=1}^D \lambda_j b_j^2 \leq \lambda_1 \overbrace{\sum_{j=1}^D b_j^2}^{\|\mathbf{b}\|^2}$$

On the other hand, for $\|\mathbf{b}\|^2$, we have:

$$\|\mathbf{b}\|^2 = \|\mathbf{P}^T \mathbf{u}\|^2 = (\mathbf{P}^T \mathbf{u})^T (\mathbf{P}^T \mathbf{u}) = \mathbf{u}^T (\mathbf{P} \mathbf{P}^T) \mathbf{u} = \|\mathbf{u}\|^2 = 1$$

Thus:

$$\forall \mathbf{u} \in \mathbb{R}^D : \mathbf{u}^T \mathbf{S} \mathbf{u} \leq \lambda_1$$

Finding Maximum Spread Direction [1]

Finding Maximum Spread Direction

We see:

$$\forall \mathbf{u} \in \mathbb{R}^D : \mathbf{u}^T \mathbf{S} \mathbf{u} \leq \lambda_1$$

Now check the variance for $\mathbf{u} = \mathbf{p}_1$:

$$\mathbf{b} = \begin{bmatrix} - & \mathbf{p}_1^T & - \\ - & \mathbf{p}_2^T & - \\ & \vdots & \\ - & \mathbf{p}_m^T & - \end{bmatrix} \mathbf{p}_1 = \begin{bmatrix} \mathbf{p}_1^T \mathbf{p}_1 \\ \mathbf{p}_2^T \mathbf{p}_1 \\ \vdots \\ \mathbf{p}_D^T \mathbf{p}_1 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

Thus:

$$\mathbf{p}_1^T \mathbf{S} \mathbf{p}_1 = \mathbf{b}^T \mathbf{D} \mathbf{b} = \sum_{j=1}^D \lambda_j b_j = \lambda_1$$

And $\mathbf{u} = \mathbf{p}_1$ is the direction of maximum spread.

Finding Next Maximum Spread Directions [1]

Finding Next Maximum Spread Directions

Assume $\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_D$ to be the eigenvectors of \mathbf{S} matrix corresponding to eigenvalues sorted in the descending order. Then, we have seen:

$$\mathbf{p}_1 \in \underset{\|\mathbf{u}\|=1}{\operatorname{argmax}} \mathbf{u}^T \mathbf{S} \mathbf{u}$$

We can show the following in an almost similar way:

$$\mathbf{p}_2 \in \underset{\|\mathbf{u}\|=1, \mathbf{u} \perp \mathbf{p}_1}{\operatorname{argmax}} \mathbf{u}^T \mathbf{S} \mathbf{u}$$

$$\mathbf{p}_3 \in \underset{\|\mathbf{u}\|=1, \mathbf{u} \perp \mathbf{p}_i, i=1,2}{\operatorname{argmax}} \mathbf{u}^T \mathbf{S} \mathbf{u}$$

\vdots

$$\mathbf{p}_j \in \underset{\|\mathbf{u}\|=1, \mathbf{u} \perp \mathbf{p}_k, k=1, \dots, (j-1)}{\operatorname{argmax}} \mathbf{u}^T \mathbf{S} \mathbf{u}$$

Subsection 2

Interpretation Via Reconstruction

PCA Interpretation Using Reconstruction

Assume we have a high-dimensional data $\mathbf{x} \in \mathbb{R}^D$ and we want to project it to a low dimensional subspace $\mathbf{z} \in \mathbb{R}^L$ such that low dimensional representation is a *good representation*. To approach a mathematical formulation, we need:

- A projection (encoding) operator: $\mathbf{z} = \text{Encode}(\mathbf{x}; \boldsymbol{\theta})$
- An un-projection (decoding) operator: $\hat{\mathbf{x}} = \text{Decode}(\mathbf{z}; \boldsymbol{\theta})$
- A goodness measure: $\|\mathbf{x} - \hat{\mathbf{x}}\|^2$

Parameters

- Representation in the low dimensional space $\mathbf{z} \in \mathbb{R}^L$
- Basis functions for reconstruction $\hat{\mathbf{x}} = \sum_{i=1}^L z_i \mathbf{w}_i$ such that:

$$\mathbf{w}_i^T \mathbf{w}_j = \begin{cases} 1 & i = j \\ 0 & i \neq j \end{cases}$$

Or equivalently if $\mathbf{W} = [\mathbf{w}_1 \quad \mathbf{w}_2 \quad \dots \quad \mathbf{w}_L] \in \mathbb{R}^{D \times L}$ then:

$$\mathbf{W}^T \mathbf{W} = \begin{bmatrix} \mathbf{w}_1^T \\ \mathbf{w}_2^T \\ \vdots \\ \mathbf{w}_L^T \end{bmatrix} [\mathbf{w}_1 \quad \mathbf{w}_2 \quad \dots \quad \mathbf{w}_L] = \begin{bmatrix} \mathbf{w}_1^T \mathbf{w}_1 & \mathbf{w}_1^T \mathbf{w}_2 & \dots & \mathbf{w}_1^T \mathbf{w}_L \\ \mathbf{w}_2^T \mathbf{w}_1 & \mathbf{w}_2^T \mathbf{w}_2 & \dots & \mathbf{w}_2^T \mathbf{w}_L \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{w}_L^T \mathbf{w}_1 & \mathbf{w}_L^T \mathbf{w}_2 & \dots & \mathbf{w}_L^T \mathbf{w}_L \end{bmatrix} = \mathbf{I}$$

PCA Interpretation Using Reconstruction

You are given a dataset $\{\mathbf{x}_i\}_{i=1}^N$ in \mathbb{R}^D . You should design $\mathbf{W} \in \mathbb{R}^{D \times L}$ and $\{\mathbf{z}_i\}_{i=1}^N$ using the following problem:

$$\min_{\mathbf{W}, \{\mathbf{z}_k\}} \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - \mathbf{W} \mathbf{z}_i\|_2^2$$

Simplifying the Loss

In this case, the loss function is:

$$\begin{aligned}\mathcal{L}(\mathbf{w}_1, \{z_k^1\}) &= \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - z_i^1 \mathbf{w}_1\|^2 = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - z_i^1 \mathbf{w}_1)^T (\mathbf{x}_i - z_i^1 \mathbf{w}_1) \\ &= \frac{1}{N} \sum_{i=1}^N \left[\mathbf{x}_i^T \mathbf{x}_i - 2z_i^1 \mathbf{w}_1^T \mathbf{x}_i + (z_i^1)^2 \overbrace{\mathbf{w}_1^T \mathbf{w}_1}^{=1} \right] \\ &= \frac{1}{N} \sum_{i=1}^N [\mathbf{x}_i^T \mathbf{x}_i - 2z_i^1 \mathbf{w}_1^T \mathbf{x}_i + (z_i^1)^2]\end{aligned}$$

Basic Problem $L = 1$

Derivative w.r.t. Representation

$$\frac{\partial \mathcal{L}(\mathbf{w}_1, \{z_k^1\})}{\partial z_n^1} = \frac{1}{N} [-2\mathbf{w}_1^T \mathbf{x}_n + 2z_n^1] = 0 \Rightarrow z_n^1 = \mathbf{w}_1^T \mathbf{x}_n$$

Updating Loss Function

$$\mathcal{L}(\mathbf{w}_1) = \frac{1}{N} \sum_{i=1}^N [\mathbf{x}_i^T \mathbf{x}_i - (z_i^1)^2] = \text{const} - \frac{1}{N} \sum_{i=1}^N (z_i^1)^2$$

Dropping the constant term, we have:

$$\mathcal{L}(\mathbf{w}_1) = -\frac{1}{N} \sum_{i=1}^N (z_i^1)^2 = -\frac{1}{N} \sum_{i=1}^N \mathbf{w}_1^T \mathbf{x}_i \mathbf{x}_i^T \mathbf{w}_1 = -\mathbf{w}_1^T \mathbf{S} \mathbf{w}_1$$

Note that in the above, we assumed the empirical mean vector to be zero ($\bar{\mathbf{x}} = \mathbf{0}$)

Basic Problem $L = 1$

Solving for \mathbf{w}_1

We have the following optimization problem:

$$\min_{\mathbf{w}_1} \mathbf{w}_1^T \mathbf{S} \mathbf{w}_1 \quad \text{subject to } \mathbf{w}_1^T \mathbf{w}_1 = 1$$

Thus we form the Lagrangian as:

$$\tilde{\mathcal{L}}(\mathbf{w}_1) = \mathbf{w}_1^T \mathbf{S} \mathbf{w}_1 - \lambda_1 (\mathbf{w}_1^T \mathbf{w}_1 - 1)$$

The partial derivative for the Lagrangian is:

$$\frac{\partial}{\partial \mathbf{w}_1} \tilde{\mathcal{L}}(\mathbf{w}_1) = 2\mathbf{S} \mathbf{w}_1 - 2\lambda_1 \mathbf{w}_1 = 0 \Rightarrow \mathbf{S} \mathbf{w}_1 = \lambda_1 \mathbf{w}_1$$

Thus $(\lambda_1, \mathbf{w}_1)$ is a pair of (eigenvalue, eigenvector). But which of them?

$$\mathbf{w}_1^T \mathbf{S} \mathbf{w}_1 = \mathbf{w}_1^T \mathbf{w}_1 = \lambda_1$$

Thus \mathbf{w}_1 is the direction of eigenvector corresponding to largest eigenvalue.

General Case

Assume we want to find $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_L]$ and $\mathbf{z} = [z^1, \dots, z^L]$. Then we have the following problem:

$$\mathcal{L}(\mathbf{W}, \{\mathbf{z}_k\}) = \frac{1}{N} \sum_{i=1}^N \left\| \mathbf{x}_i - \sum_{j=1}^L z_i^j \mathbf{w}_j \right\|^2$$

And the solution is:

$$\begin{aligned} \mathbf{w}_i &= \mathbf{p}_i, \quad i = 1, \dots, L \\ z_i^j &= \mathbf{p}_j^T \mathbf{x}_i, \quad \begin{cases} i = 1, \dots, N \\ j = 1, \dots, L \end{cases} \end{aligned}$$

where $\{\mathbf{p}_i\}$ is the set of eigenvector for \mathbf{S} matrix corresponding to eigenvalues sorted in descending order.

Encoding

$$\mathbb{R}^L \ni \mathbf{z} = \text{Encode}(\mathbf{x}, \mathbf{W}) = \mathbf{W}^T \mathbf{x} = \begin{bmatrix} \mathbf{w}_1^T \mathbf{x} \\ \mathbf{w}_2^T \mathbf{x} \\ \vdots \\ \mathbf{w}_L^T \mathbf{x} \end{bmatrix} = \begin{bmatrix} z^1 \\ z^2 \\ \vdots \\ z^L \end{bmatrix}$$

Decoding

$$\mathbb{R}^D \ni \hat{\mathbf{x}} = \text{Decode}(\mathbf{z}, \mathbf{W}) = \mathbf{W} \mathbf{z} = \begin{bmatrix} \mathbf{w}_1 & \mathbf{w}_2 & \dots & \mathbf{w}_L \end{bmatrix} \begin{bmatrix} z^1 \\ z^2 \\ \vdots \\ z^L \end{bmatrix} = \sum_i z^i \mathbf{w}_i$$

Section 3

Clustering

Clustering Problem

Clustering

- Experience E : Set of N samples $\mathcal{D} = \{\mathbf{x}_n\}_{n=1}^N$
- Task T : Partition the input into regions that contains *similar* points.
- Performance measure in *Compression*: Compression loss

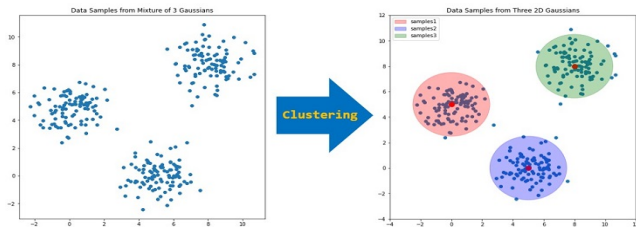


Figure: Sample GMM distribution

Section 4

Mixture Models

Mixture Models

One way to create more complex probability models is to take a convex combination of simple distributions. This is called a mixture model. This has the form $p(\mathbf{y}|\boldsymbol{\theta}) = \sum_{k=1}^K \pi_k p_c(\mathbf{y}|\boldsymbol{\theta}_k)$ where:

- $p_c(\cdot|\boldsymbol{\theta}_k)$ is the k -th mixture component
- $\{\pi_k\}_{k=1}^K$ are mixture weights with the following constraints:
 - $0 \leq \pi_k \leq 1, k = 1, \dots, K$
 - $\sum_{k=1}^K \pi_k = 1$

Mixture Models - Generative Story

Suppose latent variable z to be a categorical RV and distributed as $p(z|\boldsymbol{\theta}) = \text{Cat}(z|\boldsymbol{\pi})$ and conditional $p(\mathbf{y}|z = k, \boldsymbol{\theta}) = p_c(\mathbf{y}|\boldsymbol{\theta}_k)$. We can interpret mixture models as follows:

- We sample a specific component.
- We generate \mathbf{y} using sampled value of z .

Using the above procedure, we have:

$$p(\mathbf{y}|\boldsymbol{\theta}) = \sum_{k=1}^K p(z = k|\boldsymbol{\theta})p(\mathbf{y}|z = k, \boldsymbol{\theta}) = \sum_{k=1}^K \pi_k p(\mathbf{y}|\boldsymbol{\theta}_k)$$

Gaussian Mixture Model

Gaussian Mixture Model

Gaussian Mixture Model (GMM) or Mixture of Gaussian (MoG) is defined as:

$$p(\mathbf{y}|\boldsymbol{\theta}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{y}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

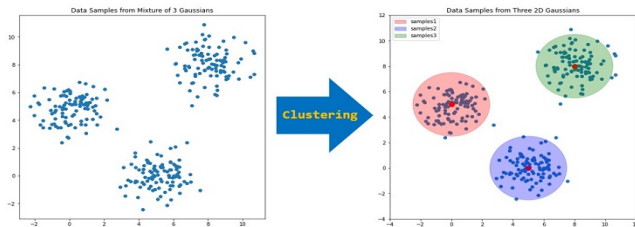


Figure: Sample GMM distribution

Problem Formulation

- Observed data samples $\{\mathbf{x}_i\}_{i=1}^n$
- Unobserved mixture element corresponding to each data sample $\{z_i\}_{i=1}^N$

Using the above two formulation, the complete dataset likelihood is:

$$p(\mathcal{D}|\boldsymbol{\theta}) = p(\{\mathbf{x}_i\}, \{z_i\}|\boldsymbol{\theta})$$

The marginal likelihood of dataset is:

$$p(\{\mathbf{x}_i\}|\boldsymbol{\theta}) = \sum_{\{z_i\}} p(\{\mathbf{x}_i\}, \{z_i\}|\boldsymbol{\theta})$$

and the maximum likelihood estimation for $\theta = \{\theta_1, \dots, \theta_K, \boldsymbol{\pi}\}$ can be calculated as:

$$\hat{\boldsymbol{\theta}}_{mle} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} p(\{\mathbf{x}_i\}|\boldsymbol{\theta})$$

Challenge and Solution

Challenge

As the scale of the problem increases (dimension of \mathbf{x} and number of dataset sample N), it becomes computationally intractable to exactly evaluate (or even optimize) the marginal likelihood.

Solution

One solution is to use expectation maximization algorithm as:

- Initialize θ randomly (or by using problem-specific heuristics) as $\theta^{(0)}$
- For $t = 1, 2, \dots, T$, repeat:
 - **E-step:** Compute posterior distribution of $\{z_i\}$ given $\{\mathbf{x}_i\}$ and $\theta^{(t-1)}$ as:

$$q^{(t)}(\{z_i\}) = p(\{z_i\}|\{\mathbf{x}_i\}, \theta^{(t-1)})$$

- **M-step:** Find $\theta^{(t)}$ as the maximizer of complete log-likelihood with respect to $q^{(t)}(\{z_i\})$ as:

$$\theta^{(t)} = \underset{\theta}{\operatorname{argmax}} \mathbb{E}_{q^{(t)}} [\log p(\{\mathbf{x}_i\}, \{z_i\}|\theta)] = \underset{\theta}{\operatorname{argmax}} \sum_{\{z_i\}} q^{(t)}(\{z_i\}) \log p(\{\mathbf{x}_i\}, \{z_i\}|\theta)$$

General Mixture Model

General Mixture Model

For a general mixture model, the samples are generated using the following distribution:

$$p(x|\boldsymbol{\theta}) = \sum_{k=1}^K \pi_k p_c(x|\boldsymbol{\theta}_k)$$

where we have:

$$\boldsymbol{\theta} = \left\{ \boldsymbol{\pi} = \begin{bmatrix} \pi_1 \\ \vdots \\ \pi_K \end{bmatrix}, \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_K \right\}$$

and $z \sim \text{Cat}(\boldsymbol{\pi})$

Complete log-Likelihood Formulation

$$\log p(\{\mathbf{x}_i\}, \{z_i\} | \boldsymbol{\theta}) = \log \prod_{i=1}^N p(\mathbf{x}_i, z_i | \boldsymbol{\theta}) = \log \prod_{i=1}^N p(\mathbf{x}_i | z_i, \boldsymbol{\theta}) p(z_i | \boldsymbol{\theta})$$

On the other hand, we have:

$$\begin{aligned} p(\mathbf{x}_i | z_i, \boldsymbol{\theta}) &= p_c(\mathbf{x}_i | \boldsymbol{\theta}_{z_i}) \\ p(z_i | \boldsymbol{\theta}) &= \pi_{z_i} \end{aligned}$$

Thus we have:

$$\begin{aligned} \log p(\{\mathbf{x}_i\}, \{z_i\} | \boldsymbol{\theta}) &= \sum_{i=1}^N (\log \pi_{z_i} + \log p_c(\mathbf{x}_i | \boldsymbol{\theta}_{z_i})) \\ &= \sum_{i=1}^N \sum_{k=1}^K \delta_{k, z_i} (\log \pi_k + \log p_c(\mathbf{x}_i | \boldsymbol{\theta}_k)) \end{aligned}$$

E-step

$$p(\{z_i\}|\{\mathbf{x}_i\}, \boldsymbol{\theta}) = \prod_{i=1}^N p(z_i|\mathbf{x}_i, \boldsymbol{\theta})$$

To compute $p(z_i|\mathbf{x}_i, \boldsymbol{\theta})$, we use Bayes rule as:

$$p(z_i = k|\mathbf{x}_i, \boldsymbol{\theta}) = \frac{p(\mathbf{x}_i|z_i = k, \boldsymbol{\theta})p(z_i = k|\boldsymbol{\theta})}{\sum_{l=1}^K p(\mathbf{x}_i|z_i = l, \boldsymbol{\theta})p(z_i = l|\boldsymbol{\theta})} = \frac{\pi_k p_c(\mathbf{x}_i|\boldsymbol{\theta}_k)}{\sum_{l=1}^K \pi_l p_c(\mathbf{x}_i|\boldsymbol{\theta}_l)}$$

Thus we have:

$$q^{(t)}(\{z_i\}) \prod_{i=1}^N q_i^{(t)}(z_i), \quad q_i^{(t)}(z_i) = p(z_i|\mathbf{x}_i, \boldsymbol{\theta}^{(t-1)})$$

M-step

$$\begin{aligned} & \mathbb{E}_{q^{(t)}} \left(\sum_{i=1}^N \sum_{k=1}^K \delta_{k,z_i} (\log \pi_k + \log p_c(\mathbf{x}_i | \boldsymbol{\theta}_k)) \right) \\ &= \sum_{i=1}^N \sum_{k=1}^K \mathbb{E}_{q^{(t)}} [\delta_{k,z_i} (\log \pi_k + \log p_c(\mathbf{x}_i | \boldsymbol{\theta}_k))] \\ &= \sum_{i=1}^N \sum_{k=1}^K \mathbb{E}_{q^{(t)}} [\delta_{k,z_i}] (\log \pi_k + \log p_c(\mathbf{x}_i | \boldsymbol{\theta}_k)) \\ &= \sum_{i=1}^N \sum_{k=1}^K q_i^{(t)}(k) (\log \pi_k + \log p_c(\mathbf{x}_i | \boldsymbol{\theta}_k)) \end{aligned}$$

Now we should maximize the above over all parameters $\boldsymbol{\theta}$.

M-step

The optimization problem for different parameters is:

$$\hat{\boldsymbol{\theta}}_k^{(t)} = \operatorname{argmax}_{\boldsymbol{\theta}_k} \sum_{i=1}^N q_i^{(t)}(k) \log p_c(\mathbf{x}_i | \boldsymbol{\theta}_k)$$

$$\hat{\boldsymbol{\pi}}^{(t)} = \operatorname{argmax}_{\boldsymbol{\pi}} \sum_{i=1}^N q_i^{(t)}(k) \log \pi_k, \text{ subject to } \sum_{k=1}^K \pi_k = 1, \pi_k \leq 0$$

The second optimization problem result in the following answer:

$$\hat{\pi}_k^{(t)} = \frac{1}{N} \sum_{i=1}^N q_i^{(t)}(k)$$

Algorithm

The algorithm is as follows:

- Initialize $\{\boldsymbol{\mu}_k^{(0)}, \boldsymbol{\Sigma}_k^{(0)}\}_{k=1}^K$ randomly and $\boldsymbol{\pi}^{(0)} = \frac{1}{K} \mathbf{1}$.
- For $t = 1, 2, \dots, T$, repeat:
 - **E-step:**

$$q_i^{(t)}(z_i = k) = \frac{\pi_k^{(t-1)} p_c(\mathbf{x}_i | \boldsymbol{\mu}_k^{(t-1)}, \boldsymbol{\Sigma}_k^{(t-1)})}{\sum_{l=1}^K \pi_l^{(t-1)} p_c(\mathbf{x}_i | \boldsymbol{\mu}_l^{(t-1)}, \boldsymbol{\Sigma}_l^{(t-1)})}, \quad \begin{cases} k = 1, \dots, K \\ i = 1, \dots, N \end{cases}$$

- **M-step:**

$$\pi_k^{(t)} = \frac{1}{N} \sum_{i=1}^N q_i^{(t)}(k)$$

$$\boldsymbol{\mu}_k^{(t)} = \frac{1}{N \pi_k^{(t)}} \sum_{i=1}^N q_i^{(t)}(k) \mathbf{x}_i$$

$$\boldsymbol{\Sigma}_k^{(t)} = \frac{1}{N \pi_k^{(t)}} \sum_{i=1}^N q_i^{(t)} \left(\mathbf{x}_i - \boldsymbol{\mu}_k^{(t)} \right) \left(\mathbf{x}_i - \boldsymbol{\mu}_k^{(t)} \right)^T$$



“Lecture 9: Principal component analysis (pca),” https://ocw.mit.edu/courses/18-650-statistics-for-applications-fall-2016/resources/mit18_650f16_pca/,
Accessed: 2022-09-24.