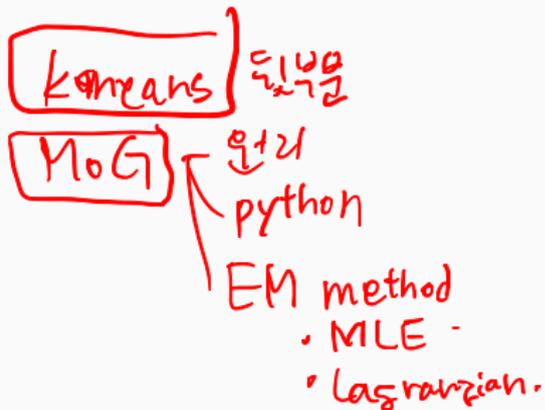


Lecture 07: Clustering II

[SCS4049-02] Machine Learning and Data Science

Seongsik Park (s.park@dgu.edu)

AI Department, Dongguk University



K-means algorithm: Advanced

K-means: centroid initialization method

K-means, MoG.

random \rightarrow K-means \rightarrow MoG
initial

If you happen to know approximately where the centroids should be (e.g., if you ran another clustering algorithm earlier), then you can set the initial hyperparameter containing the list of centroids.

Another solution is to run the algorithm multiple times with different random initializations and keep the best solution.

But, how exactly does it know which solution is the best? It uses a performance metric, called inertia, which is the mean squared distance between each instance and its closest centroid.

AIC
BIC
+
Information
criterion

K-means: inertia $\sum_{i=1}^n \|x_i - \mu_k\|^2$

MoG: log likelihood



mean squared distance
 $\text{mean. } \|x_n - \mu_k\|^2$

K-means: mini-batch k-means

Instead of using the full dataset at each iteration, the algorithm is capable of using mini-batches, moving the centroids just slightly at each iteration. This speeds up the algorithm typically by a factor of three or four and makes it possible to cluster huge datasets that do not fit in memory.

Although the mini-batch k-means algorithm is much faster than the regular k-means algorithm, its inertia is generally slightly worse, especially as the number of clusters increases.

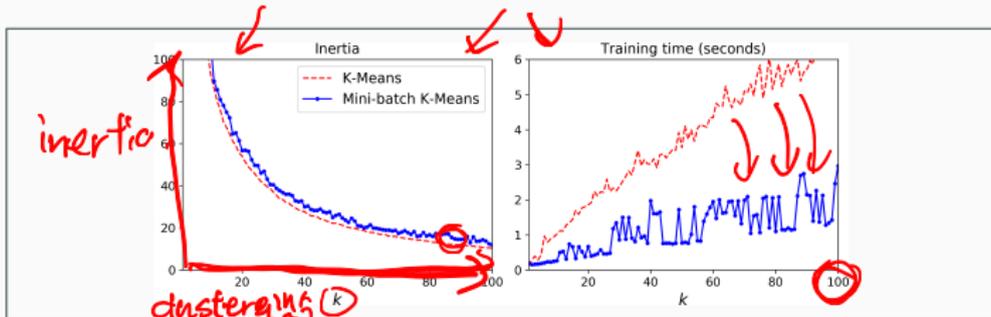


Figure 9-6. Mini-batch K-Means vs K-Means: worse inertia as k increases (left) but much faster (right)

K-means: finding the optimal number of clusters

Elbow rule: any lower value would be dramatic, while any higher value would not help much.

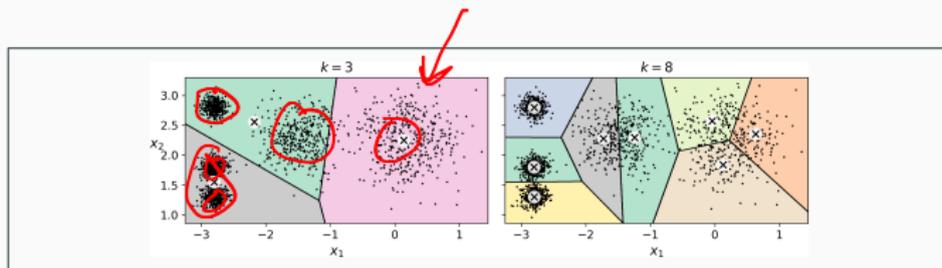


Figure 9-7. Bad choices for the number of clusters

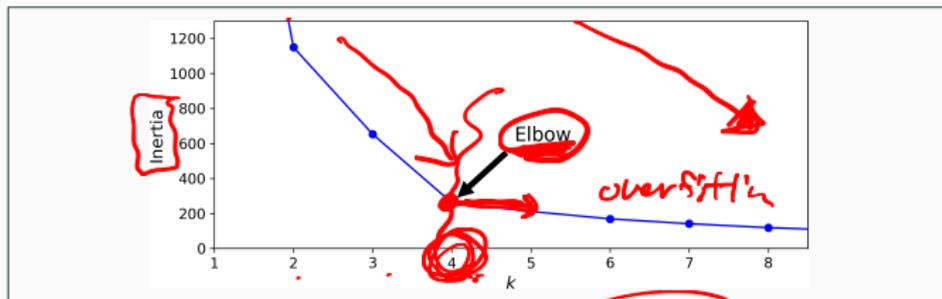


Figure 9-8. Selecting the number of clusters k using the "elbow rule"

K-means: finding the optimal number of clusters

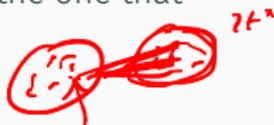
A more precise approach (but also more computationally expensive) is to use the silhouette score, which is the mean silhouette coefficient over all the instances.

- An instance's silhouette coefficient is equal to $(b - a) / \max(a, b)$
- where a is the mean distance to the other instances in the same cluster (i.e., the mean intra-cluster distance)
- and b is the mean nearest-cluster distance (i.e., the mean distance to the instances of the next closest cluster, defined as the one that minimizes b , excluding the instance's own cluster).

The silhouette coefficient can vary between -1 and +1.

- A coefficient close to +1 means that the instance is well inside its own cluster and far from other clusters,
- while a coefficient close to 0 means that it is close to a cluster boundary,
- and finally a coefficient close to -1 means that the instance may have been assigned to the wrong cluster.

$$b - a \geq \text{mean intra-cluster distance}$$



$$b \approx a$$

$$b < a \Rightarrow -\frac{a}{a} = -1$$



K-means: finding the optimal number of clusters

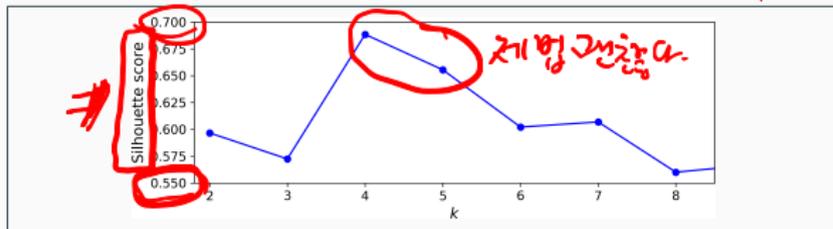


Figure 9-9. Selecting the number of clusters k using the silhouette score coefficient. This is called a silhouette diagram (see Figure 9-10):

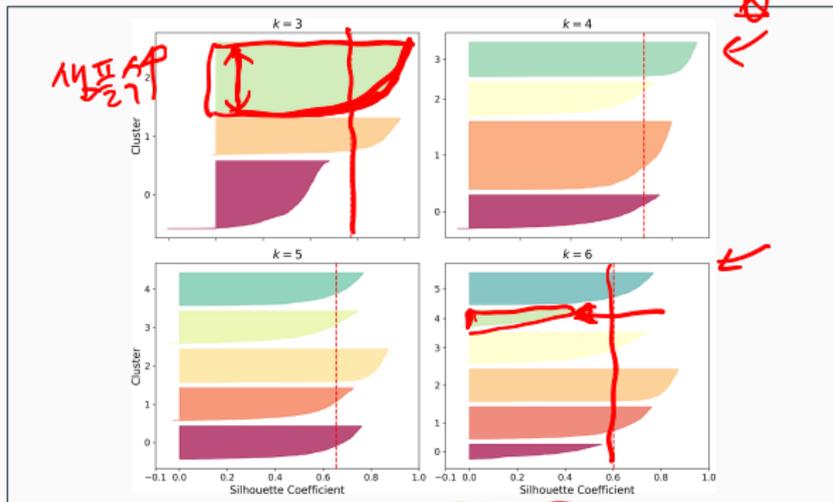


Figure 9-10. Silhouette analysis: comparing the silhouette diagram for various values of k

K-means: limitations

- It is necessary to run the algorithm several times to avoid suboptimal solution
- You need to specify the number of clusters
- K-means does not behave very well when the clusters have varying sizes, different densities, or nonspherical shapes

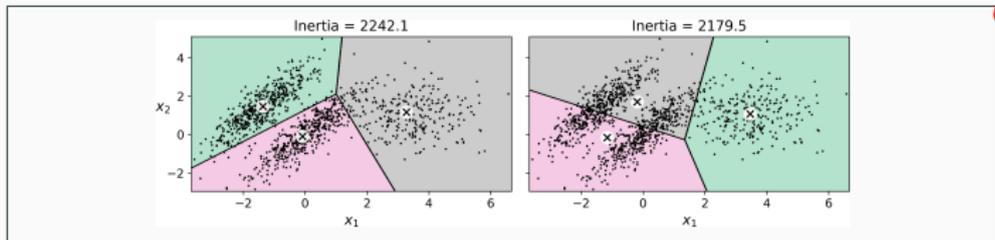


Figure 9-11. K-Means fails to cluster these elliptical blobs properly

k-means, MoG.



Mixtures of Gaussians

1. Initialize the means μ_k covariances Σ_k and mixing coefficients π_k
2. E-step Evaluate the responsibilities using the current parameter values

Expectation

$$\gamma(z_{nk}) = \frac{\pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(x_n | \mu_j, \Sigma_j)} \quad \begin{matrix} K=1,2,\dots,K \\ n=1,2,\dots,N \end{matrix} \quad (1)$$

3. M-step Re-estimate the parameters using the current responsibilities

μ, Σ, π
MLE

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) x_n \quad \begin{matrix} \text{K-means} \\ \text{이항분포의 경우와 비슷해?} \\ \text{0.0 1.0 0.0 0.0 0.0} \end{matrix} \quad (2)$$

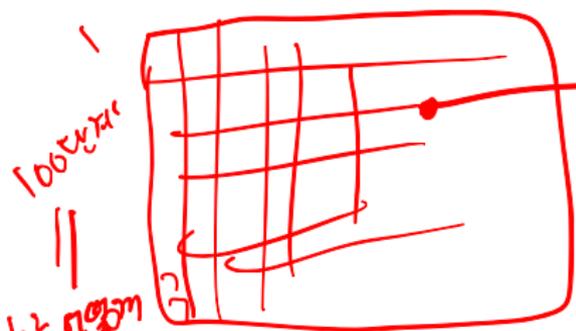
$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (x_n - \mu_k)(x_n - \mu_k)^T \quad (3)$$

$$\pi_k = \frac{N_k}{N} \quad \begin{matrix} N_k \leftarrow \sum \gamma(z_{nk}) \end{matrix} \quad (4)$$

4. Evaluate the log likelihood

$$\log p(\mathbf{X} | \mu, \Sigma, \pi) = \sum_{n=1}^N \log \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k) \right\} \quad (5)$$

and check for convergence of either the parameters or the log likelihood. If the convergence criterion is not satisfied, return to step 2.



$$\left[\frac{i}{100} \right] \text{ (circled)} = \frac{i}{100} \mathbf{e}_2$$

100x100
 1/2 resolution
 vector

100x100

$$1 \Rightarrow \frac{1}{100} \mathbf{e}_2 = \frac{1}{100} \mathbf{e}_2$$

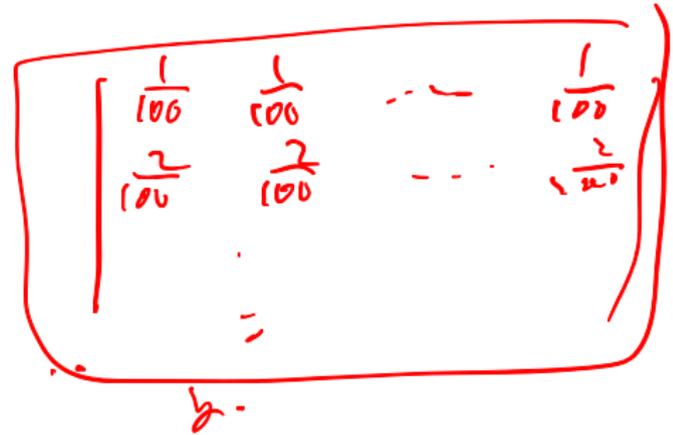
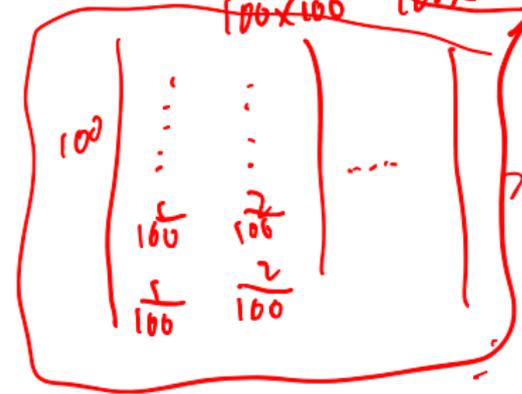
mesh grid (0,0)



100x100

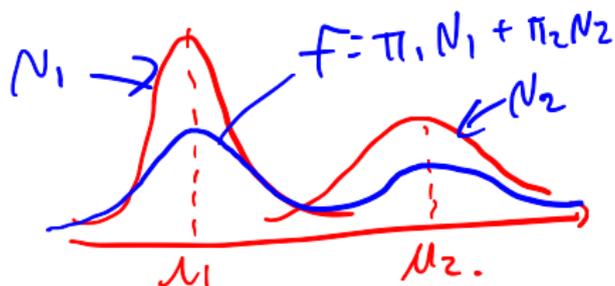
100x100

\mathbf{x}



mixture distribution.

1-D Gaussian $N(\mu_1, \sigma_1^2)$, $N(\mu_2, \sigma_2^2)$.



mixture distribution

$$f(x; \mu_1, \sigma_1^2, \mu_2, \sigma_2^2, \pi_1, \pi_2)$$

$$= \underbrace{\pi_1}_{0.5} \underbrace{N(x; \mu_1, \sigma_1^2)}_{0.5} + \underbrace{\pi_2}_{0.5} \underbrace{N(x; \mu_2, \sigma_2^2)}_{0.5}$$

$$N(x; \mu, \sigma^2)$$

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

mixture coefficient

$$\underline{\pi_1 + \pi_2 = 1}$$

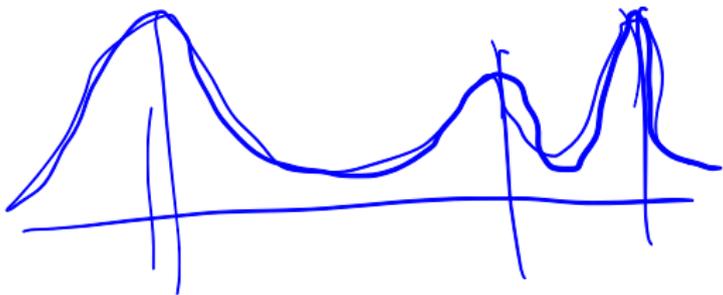


mode
distribution

$$N(\mu, \sigma^2)$$

mode = 가장 높은 값을 가진 μ .

mixture, 여러개의 mode가 가능해짐. 2, 3, 4, ...



$$\text{mode} = \mu_k$$

$$\sum_{k=1}^3 N(x; \mu_k, \sigma_k^2)$$

MoG: model

MoG, GMM ...

Mixture of Gaussian distributions can be written as a linear superposition of Gaussians.

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \mu_k, \Sigma_k) \quad (6)$$

K개의 분포를 각각 μ_k, Σ_k

Let us introduce K -dimensional binary random variable \mathbf{z} having a 1-of- K representation in which a particular element z_k is equal to 1 and all other elements are equal to 0. The values of z_k therefore satisfy $z_k \in \{0, 1\}$ and $\sum_k z_k = 1$, and we see that there are K possible states for the vector \mathbf{z} according to which element is nonzero.

one-hot 벡터.

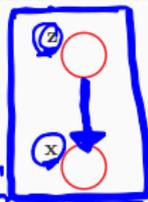
The marginal distribution over \mathbf{z} is specified in terms of the mixing coefficients π_k , such that

$$p(z_k = 1) = \pi_k$$
$$0 \leq \pi_k \leq 1 \text{ and } \sum_{k=1}^K \pi_k = 1$$

$\pi_k = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \\ \vdots & & & \end{bmatrix}$

MoG: 1-of-K representation

Figure 9.4 Graphical representation of a mixture model, in which the joint distribution is expressed in the form $p(\mathbf{x}, \mathbf{z}) = p(\mathbf{z})p(\mathbf{x}|\mathbf{z})$.



conditional probability

$z_k, \text{node} = \text{r.v.}$

$z_k \rightarrow x = \text{dependence}$

probabilistic graphical model, PGM.

$$p(\mathbf{x}, \mathbf{z}) = p(\mathbf{z}|\mathbf{x}) p(\mathbf{z})$$

Because \mathbf{z}_k uses a 1-of-K representation, we can also write this distribution in the form

$$p(\mathbf{z}) = p(z_1, \dots, z_K) = \prod_{k=1}^K \pi_k^{z_k}$$

$$\mathbf{z} = (z_1, \dots, z_K)$$

$$p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k} \quad (8)$$

Similarly, the conditional distribution of \mathbf{x} given a particular value for \mathbf{z} is a Gaussian

\leftarrow in cluster k

$$p(\mathbf{x}|z_k = 1) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (9)$$

which can also be written in the form

$z_k \in \{0, 1\}$

$$p(\mathbf{x}|\mathbf{z}) = \prod_{k=1}^K \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_k} \quad (10)$$

Bishop

MoG: model

The joint distribution is given by $p(\mathbf{x}|\mathbf{z})p(\mathbf{z})$, and the marginal distribution of \mathbf{x} is then obtained by summing the joint distribution over all possible states of \mathbf{z} to give

$$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{x}|\mathbf{z})p(\mathbf{z}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (11)$$

Another quantity that will play an important role is the condition probability of \mathbf{z} given \mathbf{x} . We shall use $\gamma(z_k)$ to denote $p(z_k = 1|\mathbf{x})$, whose value can be found using Bayes' theorem

$$\gamma(z_k) \equiv p(z_k = 1|\mathbf{x}) = \frac{p(\mathbf{x}|z_k = 1)p(z_k = 1)}{\sum_{j=1}^K p(\mathbf{x}|z_j = 1)p(z_j = 1)} = \frac{\pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \quad (12)$$

posterior

responsibility

normalize

$$= \frac{\pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \quad (13)$$

marginal posterior

$$P(x) = \sum_z \underbrace{P(x|z)}_{\text{likelihood}} \underbrace{P(z)}_{\text{prior}}$$

posterior

$$P(z|x) = \frac{P(z, x)}{P(x)} = \frac{P(x|z)P(z)}{\sum_z P(x|z)P(z)}$$

분리적인지

x 주어졌을 때

k=1	0.1
:	0.1
:	0.3
5	0.5
	0.0

x가 주어졌을 때, z의 값을
= k번째 cluster에 속할 확률

arg max
IC

$$\gamma(z_k) = P(z|x)$$

$$\begin{aligned}
 \underline{p(z_1 | x)} &= \gamma(z_1) = \underline{0.1} \\
 \underline{p(z_2 | x)} &= \gamma(z_2) = \underline{0.05} \\
 \underline{p(z_3 | x)} &= \gamma(z_3) = \underline{0.55} \\
 \underline{p(z_4 | x)} &= \gamma(z_4) = \underline{0.3}
 \end{aligned}$$

즉어진 x 에서
 각각의
 cluster에
 속할
 확률

$$\underset{k}{\operatorname{argmax}} \gamma(z_k)$$

$$p(z|x) = \frac{p(z,x)}{\sum_z p(z,x)} = \frac{p(x|z)p(z)}{\sum_z p(x|z)p(z)}$$

k-means.

center lb,
mean
 μ_k



각 데이터가
어떤 cluster에 속하는지

파괴
cov, π_k

$$\arg \min_k \|x_n - \mu_k\|^2$$

MoG.

max. $\gamma(z_k)$ 파괴

We shall view π_k as the prior probability of $z_k = 1$, and the quantity $\gamma(z_k)$ as the corresponding posterior probability once we have observed \mathbf{x} . As we shall see later, $\gamma(z_k)$ can also be viewed as the *responsibility* that component k takes for ‘explaining’ the observation \mathbf{x} .

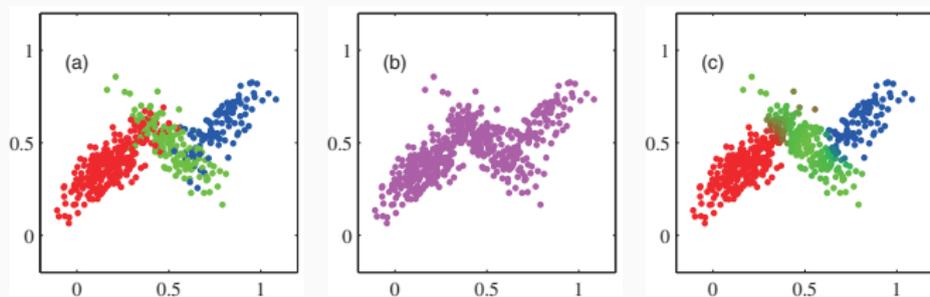
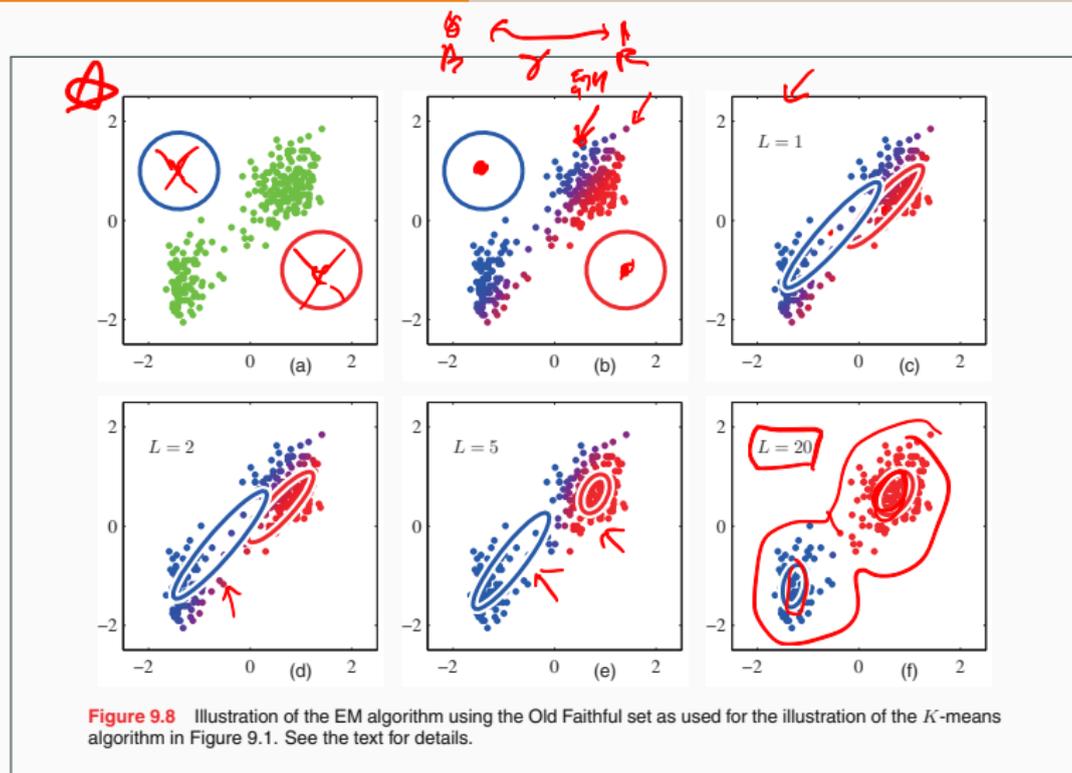
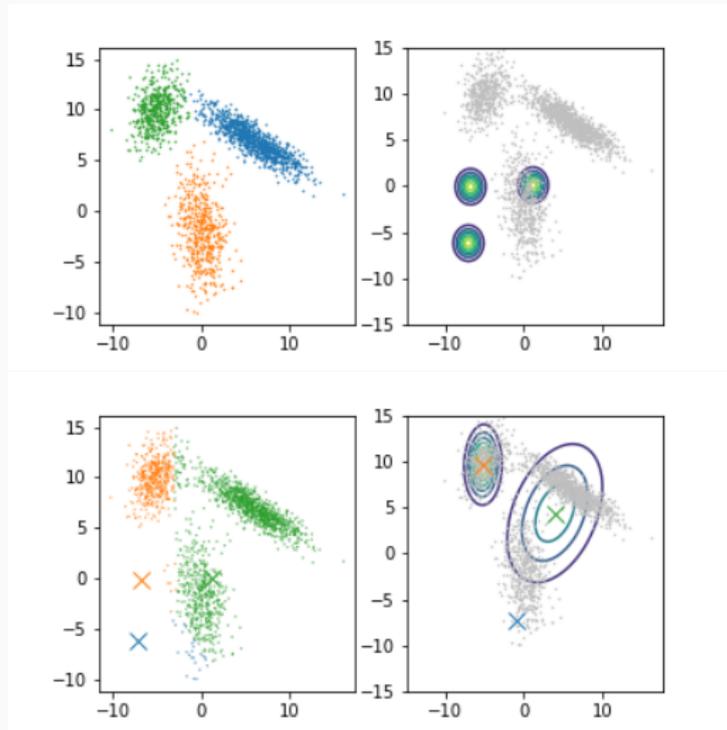


Figure 9.5 Example of 500 points drawn from the mixture of 3 Gaussians shown in Figure 2.23. (a) Samples from the joint distribution $p(\mathbf{z})p(\mathbf{x}|\mathbf{z})$ in which the three states of \mathbf{z} , corresponding to the three components of the mixture, are depicted in red, green, and blue, and (b) the corresponding samples from the marginal distribution $p(\mathbf{x})$, which is obtained by simply ignoring the values of \mathbf{z} and just plotting the \mathbf{x} values. The data set in (a) is said to be *complete*, whereas that in (b) is *incomplete*. (c) The same samples in which the colours represent the value of the responsibilities $\gamma(z_{nk})$ associated with data point \mathbf{x}_n , obtained by plotting the corresponding point using proportions of red, blue, and green ink given by $\gamma(z_{nk})$ for $k = 1, 2, 3$, respectively

MoG: example



MoG: python example



example

Figure 1: MoG python example: dataset and initialization (top) and the 1st iteration (bottom).

MoG: python example

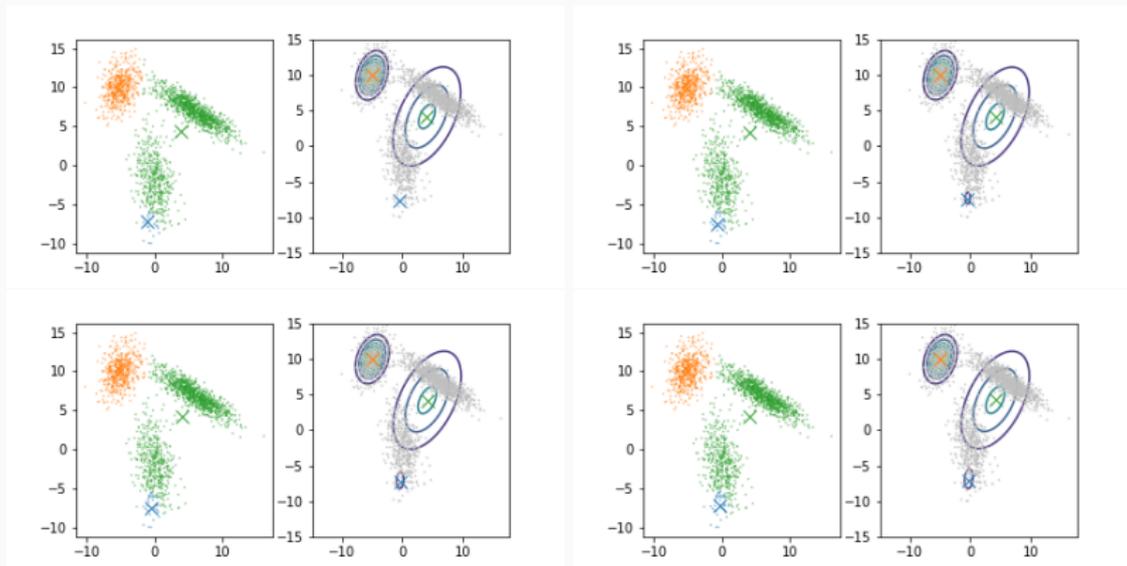


Figure 2: From the 2nd to 5th iterations.

MoG: python example

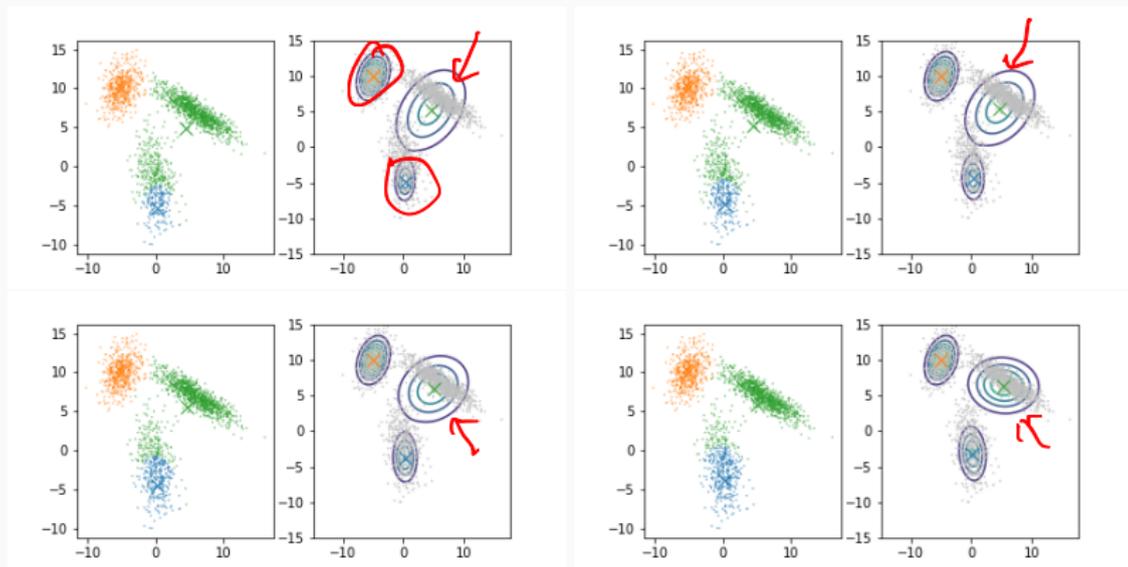


Figure 3: From the 12th to 15th iterations.

MoG: python example

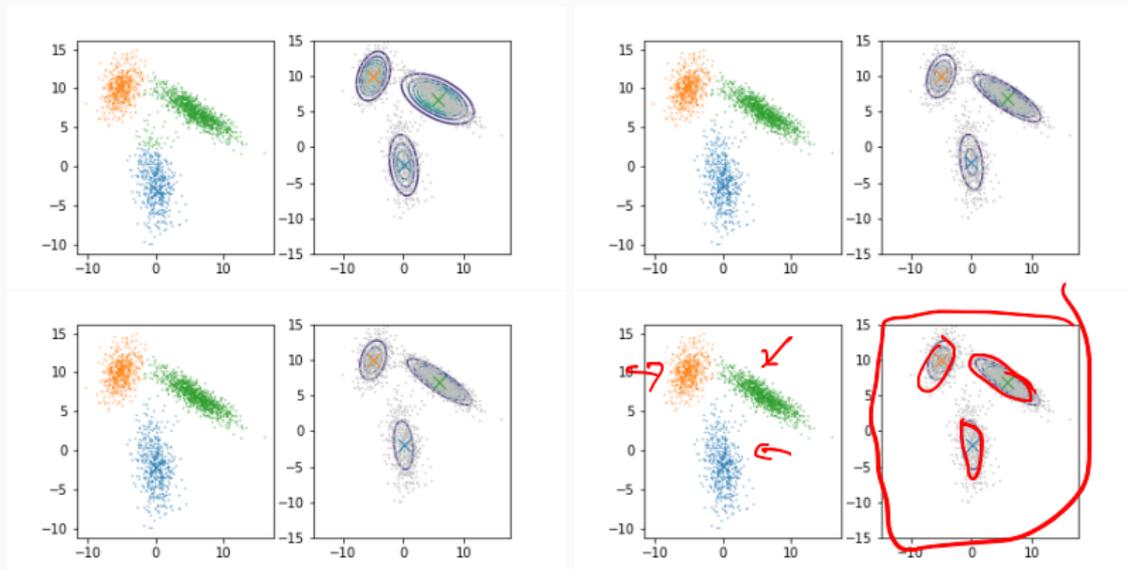


Figure 4: From the 16th to 19th iterations.

Mixtures of Gaussians: EM algorithm

- Maximum Likelihood Estimate (MLE)
- Lagrangian multiplier.

MoG: maximum likelihood estimate

Maximum likelihood estimation, or MLE, is on flavor of parameter estimation in machine learning. In order to perform parameter estimation, we need:

- some data \mathbf{x} ← \mathcal{M}
- some hypothesized generating function of the data $f(\mathbf{x}, \theta)$
- a set of parameters from that function θ μ
- some evaluation of the goodness of our parameters (an objective function)

$$\mathbf{x} \sim \text{분포.}$$

→ Likelihood.

In MLE, the objective function (evaluation) we chose is the likelihood of the data given our model. To find the best θ then, we need to find the θ which maximizes our evaluation function (the likelihood).

Therefore, in its general form the MLE is:

$$\hat{\theta}_{MLE} = \arg \max_{\theta} p(\mathbf{x}|\theta) \quad (14)$$

MLE

θ
MLE

→ likelihood function $p(\mathbf{x}|\theta)$

MoG: maximum likelihood estimate

$\sim N(\mu, \sigma^2)$ known n

Gaussian distribution을 따르는 i.i.d. 샘플 $\mathbf{x} = (x_1, x_2, \dots, x_N)$ 로부터
평균 $\theta = \mu$ 를 MLE로 추정하면,

unknown n

NTH.

↓ ↓ ↓ ↓ ↓

→ independent & identically distributed

\mathcal{L} : Likelihood function $\mathcal{L} = p(\mathbf{x}|\theta) = \prod_{n=1}^N \mathcal{N}(x_n|\mu)$ (15)

log-likelihood

$\log \mathcal{L} = \sum_{n=1}^N \log \mathcal{N}(x_n|\mu)$ (16)

$\frac{d}{d\mu} \log \mathcal{L} = -\text{const} \sum_{n=1}^N (x_n - \mu) = a$ (17)

$\frac{d}{d\mu} \log \mathcal{L} = 0 \iff \hat{\mu} = \frac{1}{N} \sum_{n=1}^N x_n$ (18)

arithmetic mean

$\arg \max_{\theta} \mathcal{L} = \arg \max_{\theta} \log \mathcal{L}$. $\hat{\mu}_{MLE} = \frac{1}{N} \sum_{n=1}^N x_n$

MoG: EM algorithm

An elegant and powerful method for finding maximum likelihood solutions for models with latent variables is called the expectation-maximization algorithm, or EM algorithm.

iterative algorithm

Let us begin by writing down the conditions that must be satisfied at a maximum of the likelihood function.

$p(\mathbf{x}_n | \mu, \Sigma, \pi)$

MoG

$$\log p(\mathbf{X} | \mu, \Sigma, \pi) = \sum_{n=1}^N \log \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k) \right\} \quad (19)$$

Setting the derivatives of log likelihood with respect to the means μ_k of the Gaussian components to zero,

$$\frac{\partial}{\partial \mu} \log p = - \sum_{n=1}^N \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \mu_j, \Sigma_j)} \Sigma_k (\mathbf{x}_n - \mu_k) \quad (20)$$

$\gamma(z_{nk})$

$\frac{\partial}{\partial \mu} \log p$

MoG: EM algorithm

Multiplying by Σ_k^{-1} and rearranging

$$0 = - \sum_{n=1}^N \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\underbrace{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}_{\gamma(Z_{nk})}} \boldsymbol{\Sigma}_k (\mathbf{x}_n - \boldsymbol{\mu}_k) \quad (21)$$

we obtain

$$\boxed{0} = \sum_{n=1}^N \gamma(Z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k) \quad (22)$$

$$\overset{1}{N_k} \boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(Z_{nk}) \mathbf{x}_n \quad (23)$$

$$N_k = \sum_{n=1}^N \gamma(Z_{nk}) \quad \leftarrow \text{1c-means samples} \quad (24)$$

$$\frac{\partial}{\partial \Sigma_k} \log P$$

If we set the derivative of log likelihood with respect to Σ_k to zero, and follow a similar line of reasoning, making use of the result for the maximum likelihood solution for the covariance matrix of a single Gaussian, we obtain

$$\Delta \Sigma_{k,MLE} \quad \star \quad \Sigma_k = \frac{1}{N_k} \sum_{n=1}^N \underbrace{\gamma(z_{nk})}_{\text{weight}} \underbrace{(\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T}_{\text{cov. defin}} \quad (25)$$

each data point weighted by the corresponding posterior probability and with the denominator given by the effective number of points associated with the corresponding component

MoG: EM algorithm

Finally, we maximize log likelihood with respect to the mixing coefficients π_k . Here we must take account of the constraint $\sum_k \pi_k = 1$. This can be achieved using a Lagrange multiplier and maximizing the following quantity

$$\log p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) + \lambda \left(\sum_{k=1}^K \pi_k - 1 \right) \quad (26)$$

which gives

$$0 = \sum_{n=1}^N \frac{\mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_j \pi_j \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} + \lambda \quad (27)$$

$$0 = \sum_{k=1}^K \sum_{n=1}^N \gamma(z_{nk}) + \lambda \sum_{k=1}^K \pi_k \quad (28)$$

Hence

$$\pi_k = \frac{N_k}{N} = \frac{1}{N} \sum \gamma(z_{nk}) \quad (29)$$

MoG: Lagrange multiplier

optimization + equality
(min, max) constraint.

The method can be summarized as follows: in order to find the maximum or minimum of a function $f(x)$ subjected to the equality constraint $g(x) = 0$, form the Lagrangian function constraint.

$$x, \lambda \quad \mathcal{L}(x, \lambda) = f(x) - \lambda g(x) \quad f(x) - \lambda g(x) \quad (30)$$

and find the stationary points of \mathcal{L} considered as a function of x and the Lagrange multiplier λ . The solution corresponding to the original constrained optimization is always a saddle point of the Lagrangian function, which can be identified among the stationary points from the definiteness of the bordered Hessian matrix.

MoG: Lagrange multiplier

Minimize $f(x, y) = x + y$ subject to the constraint $x^2 + y^2 = 1$, i.e.,



$$g(x, y) = x^2 + y^2 - 1 = 0 \quad (31)$$

$$1 - x^2 - y^2 = 0$$

Hence,

$$\mathcal{L}(x, y, \lambda) = f(x, y) + \lambda g(x, y) = x + y + \lambda(x^2 + y^2 - 1) \quad (32)$$

Gradient

$$\nabla_{x,y,\lambda} \mathcal{L}(x, y, \lambda) = \left(\frac{\partial}{\partial x}, \frac{\partial}{\partial y}, \frac{\partial}{\partial \lambda} \right) (1 + 2\lambda x, 1 + 2\lambda y, x^2 + y^2 - 1) \quad (33)$$

and therefore,

$$\nabla_{x,y,\lambda} \mathcal{L}(x, y, \lambda) = 0 \iff \left\{ \begin{array}{l} 1 + 2\lambda x = 0 \\ 1 + 2\lambda y = 0 \\ x^2 + y^2 - 1 = 0 \end{array} \right\} \quad (34)$$

MoG: Lagrange multiplier

$$\nabla_{x,y,\lambda} \mathcal{L}(x,y,\lambda) \stackrel{=0}{=} \iff \begin{cases} 1 + 2\lambda x = 0 \\ 1 + 2\lambda y = 0 \\ x^2 + y^2 - 1 = 0 \end{cases} \quad (35)$$

This yields

$$x = y = -\frac{1}{2\lambda}, \quad \lambda \neq 0 \quad (36)$$

$$\frac{1}{4\lambda^2} + \frac{1}{4\lambda^2} - 1 = 0 \quad (37)$$

So,

$$\lambda = \pm \frac{1}{\sqrt{2}} \quad (38)$$

which implies that the stationary points of \mathcal{L} are

$$\left(\frac{\sqrt{2}}{2}, \frac{\sqrt{2}}{2}, \frac{\sqrt{2}}{2} \right), \quad \left(-\frac{\sqrt{2}}{2}, -\frac{\sqrt{2}}{2}, \frac{\sqrt{2}}{2} \right) \quad (39)$$

Appendix

Reference and further reading

- “Chap 9” of A. Geron, Hands-On Machine Learning with Scikit-Learn, Keras & TensorFlow
- “Chap 9” of C. Bishop, Pattern Recognition and Machine Learning
- D. Arthur, S. Vassilvitskii, k-means++: The advantages of careful seeding (2006)
- Z. Ghahramani, M. J. Beal, Variational Inference for Bayesian mixtures of Factor Analysers (2000)

Due: 9월 22일 화요일, 23시 59분까지

- pdf로 업로드하세요.
- 손으로 작성한 파일을 스캔앱(Adobe scan, Office lens 등)을 써서 pdf로 저장해주세요.
- 컴퓨터로 작성(latex, word, ppt, 한글 등)한 파일도 가능합니다. pdf로 저장해주세요.
- 가독성이 떨어지는 파일도 불량처리 합니다.

1. Gaussian distribution $\mathcal{N}(x; \mu, \sigma)$ 를 따르는 N 개의 i.i.d. 샘플 $\mathbf{x} = (x_1, x_2, \dots, x_N)$ 에 대해 $\theta = (\mu, \sigma)$ 에 대한 MLE를 유도과정과 함께 구하세요. (20점)

2. python 코드로서 만능함수 부분 추가해서 올릴게요