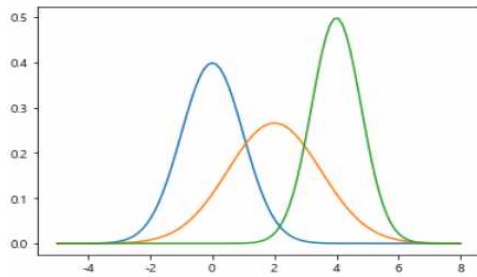


## GMM(가우시안 혼합 모델)

### <개요>

- 샘플이 파라미터가 알려지지 않은 여러 개의 혼합된 가우시안 분포에서 생성되었다고 가정하는 확률 모델
- 각 클러스터는 각각의 정규분포
- 각 샘플이 각 클러스터(각 파라미터를 가지는 정규분포)에서 생성될 확률을 구해서 가장 확률이 높은 클러스터에 할당



### ▶ GMM 변수

$x_i$  : i번째 샘플

$\pi_k$  : k번째 클러스터가 선택될 확률

$z_{ik}$  : i번째 샘플을 k번째 클러스터가 생성하였는지에 대한 여부(0과 1의 이진변수)

$\mu_k$  : k번째 클러스터의 평균

$\Sigma_k$  : k번째 클러스터의 분산

### ▶ 각 샘플이 발생할 확률

$$P(x_i) = \sum_j \pi_j N(x_i | \mu_j, \Sigma_j)$$
$$\sum_j \pi_j = 1$$
$$0 \leq \pi_j \leq 1$$

- 각 모든 클러스터에서 i번째 데이터가 생성될 확률의 합
- 각 클러스터가 전체망라, 상호배타이기 때문에 해당 식이 성립

=> GMM을 학습 : 적절한  $\pi, \mu, \Sigma$ 를 찾는 것

### <GMM 학습>

- EM 알고리즘을 적용하여 학습

#### ▶ E-step

- $\gamma(z_{ik}) = P(z_{ik} = 1 | x_i)$  를 바탕으로 함.
- 어떠한 데이터가 있는데, 이 데이터가 k번째 클러스터에서 생성되었을 확률을 구함
  - > 해당 확률을 가장 크게 만드는 클러스터가 해당 데이터의 클러스터

$$\begin{aligned}\gamma(z_{ik}) &= P(z_{ik} = 1 | x_i) \\ &= \frac{P(x_i | z_{ik}) P(z_{ik} = 1)}{\sum_j \pi_j N(x_i | \mu_j, \Sigma_j)} \\ &= \frac{\pi_k N(x_i | \mu_k, \Sigma_k)}{\sum_j \pi_j N(x_i | \mu_j, \Sigma_j)}\end{aligned}$$

#### ► M-step

- $\operatorname{argmax}_{\theta} L(X; \theta) \Rightarrow \operatorname{argmax}_{\theta} \ln L(X; \theta)$ 를 바탕으로 함
  - 각 클러스터를 구성하는 parameter가 존재할 때, 해당 데이터들이 나올 확률을 가장 크게 만드는 파라미터를 찾는 것(log가 증가함수이기 때문에 같은 의미)
  - $\ln L(X; \theta) = \ln P(X | \pi, \mu, \Sigma) = \ln \prod_{i=1}^n P(x_i | \pi, \mu, \Sigma) = \sum_{i=1}^n \ln P(x_i | \pi, \mu, \Sigma) = \sum_i \ln (\sum_j \pi_j N(x_i | \mu_j, \Sigma_j))$
  - > 목적식 :  $\operatorname{argmax}_{\pi, \Sigma, \mu} \sum_i \ln (\sum_j \pi_j N(x_i | \mu_j, \Sigma_j))$
  - > 이 때, 목적식이  $\mu_k, \Sigma_k$ 에 대한 음의 이차식이기 때문에 미분하여 0을 만족시키는 값을 최적값으로 사용
  - $\mu_k = \frac{1}{N_k} \sum_{i=1}^n \gamma(z_{ik}) x_i, \quad N_k = \sum_{i=1}^n \gamma(z_{ik}), \quad \Sigma_k = \frac{1}{N} \sum_{i=1}^n \gamma(z_{ik}) (x_i - \mu_k)(x_i - \mu_k)^T$
  - >  $\pi_k$ 의 경우 제약식( $\sum_j^k \pi_j = 1$ )이 존재하기 때문에 라그랑지안 승수 방식 사용
  - $\sum_i \ln (\sum_j \pi_j N(x_i | \mu_j, \Sigma_j)) + \lambda (1 - \sum_j^k \pi_j)$
  - >  $\lambda = N, \quad \pi_k = \frac{1}{N} \sum_{i=1}^n \gamma(z_{ik})$

#### <GMM 총정리>

- 초기화단계 :  $\mu_k, \Sigma_k, \pi_k$ 를 적당한 값으로 초기화
- E단계 :  $\gamma(z_{ik})$ 를 계산
- M단계 :  $\gamma(z_{ik})$ 를 이용하여 파라미터 업데이트
- 이를 반복

참고 : <https://untitledblog.tistory.com/133>, PRML