

Gaussian Mixture Model (GMM)

Assumptions:

① hot-encoded discrete latent variable $z_k \in \{0,1\}$

for the K clusters, with prior

$$P(z_k=1) = \pi_k \quad \pi_k \in [0,1] \quad \boxed{\sum_{k=1}^K \pi_k = 1}$$

② the clusters are Gaussian with different parameters

$$P(x|z_k=1) = N(x|\mu_k, \Sigma_k)$$

③ the joint distribution

$$\begin{aligned} P(x, z_k=1) &= P(x|z_k=1) P(z_k=1) \\ &= \pi_k N(x|\mu_k, \Sigma_k) \end{aligned}$$

④ full generative model

$$P(x) = P(x, z) = \sum_{k=1}^K \pi_k N(x|\mu_k, \Sigma_k)$$

Given the data $X = \{x_1, x_2, \dots, x_N\}$

$$\ln p(x|\pi, \mu, \Sigma) = \ln \prod_{n=1}^N p(x_n|\pi, \mu, \Sigma)$$

$$L = \sum_{n=1}^N \ln p(x_n | \Pi, M, \Sigma)$$

$$L = \sum_{n=1}^N \ln \left[\sum_{k=1}^K \pi_k N(x_n | \mu_k, \Sigma_k) \right]$$

because of this summation term
It's difficult to get a closed form
of the likelihood.

One alternative is to go for EM
(Expectation Maximization algorithm)

Multivariate Gaussian:

$$N(x | \mu_k, \Sigma_k) = \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma|^{\frac{1}{2}}} e^{-\frac{1}{2} (x - \mu_k)^T \Sigma^{-1} (x - \mu_k)}$$

$$\frac{\partial N(x | \mu_k, \Sigma_k)}{\partial \mu_k} = N(x | \mu_k, \Sigma_k) (x - \mu_k)^T \Sigma^{-1}$$

$$\frac{\partial L}{\partial \mu_k} = \sum_{n=1}^N \frac{\partial}{\partial \mu_k} \left[\ln \sum_{k=1}^K \pi_k N(x_n | \mu_k, \Sigma_k) \right]$$

$$= \sum_{n=1}^N \frac{1}{\sum_{k=1}^K \pi_k N(x_n | \mu_k, \Sigma_k)} \frac{\partial}{\partial \mu_k} \left[\sum_{k=1}^K \pi_k N(x_n | \mu_k, \Sigma_k) \right]$$

$$= \sum_{n=1}^N \left[\frac{\pi_k N(x_n | \mu_k, \Sigma_k)}{\sum_{k=1}^K \pi_k N(x_n | \mu_k, \Sigma_k)} \right] (x_n - \mu_k)^T \Sigma^{-1} \geq 0$$

↳ r_{nk} this is posterior probability or responsibilities.

$$\Rightarrow \sum_{n=1}^N r_{nk} (x_n - \mu_k)^T \Sigma^{-1} = 0$$

$$\Rightarrow \sum_{n=1}^N r_{nk} x_n = \sum_{n=1}^N r_{nk} \mu_k$$

$$\Rightarrow \mu_k = \frac{\sum_{n=1}^N r_{nk} x_n}{\sum_{n=1}^N r_{nk}}$$

μ_k = the weighted average over the points x_n for which cluster k takes the responsibility.

Maximize w.r.t π_k

but we need to take care of constraint that $\sum_{k=1}^K \pi_k = 1$, this can be achieve by lagrange multiplier.

$$L = \sum_{n=1}^N \ln p(x_n | \pi, \mu, \Sigma) + \lambda \left(\sum_{k=1}^K \pi_k - 1 \right)$$

$$\frac{\partial L}{\partial \pi_k} = \sum_{n=1}^N \frac{\partial}{\partial \pi_k} \left(\ln \left(\sum_{k=1}^K \pi_k N(x_n | \mu_k, \Sigma_k) \right) \right) + \lambda$$

$$= \sum_{n=1}^N \frac{\pi_k N(x_n | \mu_k, \Sigma_k)}{\sum_{k=1}^K \pi_k N(x_n | \mu_k, \Sigma_k)} + \lambda = 0$$

$$\Rightarrow \sum_{k=1}^N \pi_k + \lambda = 0$$

$$\Rightarrow \pi_k = -\frac{1}{\lambda} \sum_{n=1}^N r_{nk}$$

So, π_k also depends on λ .

$$\text{So, } \frac{\partial}{\partial \lambda} L(\pi_k, \lambda) = \frac{\partial}{\partial \lambda} \lambda \left(\sum_{k=1}^K \pi_k - 1 \right) = 0$$

$$\Rightarrow \sum_{k=1}^K \pi_k = 1$$

$$\Rightarrow \sum_{k=1}^K \frac{1}{\lambda} \sum_{n=1}^N r_{nk} = 1$$

$$\Rightarrow \lambda = -\sum_{k=1}^K \sum_{n=1}^N r_{nk}$$

$$\Rightarrow \boxed{\lambda = -N}$$

$$\boxed{\sum_{k=1}^K r_{nk} = 1}$$

$$\text{So, } \boxed{\pi_k = \frac{1}{N} \sum_{n=1}^N r_{nk}}$$

fraction of points
for which cluster
 k carries responsibility

so, for \sum_{K}

$$S_K = \frac{1}{N_K} \sum_{n=1}^N r_{nk} (x_n - \bar{x}_n) (\bar{x} - \bar{x}_n)^T \quad \checkmark$$

\bar{x}

$$P(A) = \gamma_2$$

$$P(B) = u \quad 0 \leq u \leq \gamma_6$$

$$P(C) = 2u$$

$$P(D) = \gamma_2 - 3u$$

estimate u from the data given.

$$P(D|\theta) = \underset{\theta}{\operatorname{argmax}} P(\theta|D)$$

$$L = P(a, b, c, d | u) = \left(\frac{1}{2}\right)^a (u)^b (2u)^c \left(\frac{1}{2} - 3u\right)^d$$

by likelihood

$$\begin{aligned} \frac{\partial L}{\partial u} &= \frac{\partial}{\partial u} \left[\log \left(\frac{1}{2} \right)^a + \log(u)^b + \log(2u)^c + \log \left(\frac{1}{2} - 3u \right)^d \right] = 0 \\ &\Rightarrow \frac{\partial}{\partial u} [a \log \frac{1}{2} + b \log u + c \log 2u + d \log \left(\frac{1}{2} - 3u \right)] = 0 \end{aligned}$$

$$\Rightarrow \frac{b}{u} + \frac{c}{2u} + \frac{d}{\frac{1}{2} - 3u} = 0$$

$$\Rightarrow \frac{1}{u} \left(b + \frac{c}{2} \right) + \frac{2d}{1 - 6u} = 0$$

$$\Rightarrow \frac{1-6u}{u} \left(b + \frac{c}{2} \right) + 2d = 0$$

$$\Rightarrow \left(\frac{1}{u} - 6 \right) \left(b + \frac{c}{2} \right) + 2d = 0$$

$$\Rightarrow \frac{b}{u} + \frac{c}{2u} - 6b - 3c + 2d = 0$$

$$\Rightarrow \frac{1}{u} \left(b + \frac{c}{2} \right) = 6b + 3c - 2d$$

$$\Rightarrow u = \frac{2b+c}{6b+3c-2d}$$

SVM

$$w^*, b^* = \underset{w, b}{\operatorname{argmin}} \frac{\|w\|}{2} + C \sum_{i=1}^n \epsilon_i \text{ such that}$$

for correctly
classified point
 $\epsilon_i = 0.$

$$y_i (w^T x_i + b) \geq 1 - \epsilon_i \quad \epsilon_i \geq 0$$

using lagrangian formulation for this optimization problem.

Introducing two Lagrangian multiplier for each inequality constraint.

$$L(w^*, b^*) = \underset{w, b}{\operatorname{argmin}} \frac{\|w\|}{2} + C \sum_{i=1}^n \epsilon_i - \sum_{i=1}^n \alpha_i \{ y_i (w^T x_i + b) - 1 + \epsilon_i \} - \sum_{i=1}^n \beta_i \epsilon_i$$

Take derivatives of L w.r.t w, b, ϵ_i

$$\frac{\partial L}{\partial w} = w - \sum_{i=1}^n \alpha_i y_i x_i = 0$$

$$\Rightarrow w = \sum_{i=1}^n \alpha_i y_i x_i$$

$$\frac{\partial L}{\partial b} = \boxed{\sum_{i=1}^n \alpha_i y_i = 0}$$

$$\frac{\partial L}{\partial \alpha_i} = c - \alpha_i - \beta_i = 0 \Rightarrow \boxed{c = \alpha_i + \beta_i}$$

Now Substitute the values.

$$L(w^*, b^*) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \boxed{x_i^T x_j}$$

S.t. $0 \leq \alpha_i \leq c$

dual form

Observations

- ① for every x_i there is an α_i
- ② x_i 's only occur in the form of $x_i^T x_j$
- ③ $f(x_q) = \sum_{i=1}^n \alpha_i y_i x_i^T x_q + b$
- ④ $\alpha_i > 0$ only for support vectors
 $= 0$ for all non-support vectors.

$x_i^T x_j$ can be written as $K(x_i, x_j)$

at runtime

$$f(x_q) = \sum_{i=1}^n \alpha_i y_i \boxed{x_i^T x_q} + b$$

$\rightarrow K(x_i, x_q)$

this is similarity function, can be any similarity.

↓

Replacing this from (dot product) to
 $K(x_i, x_j)$ is called kernel trick & such
SVM is called Kernel SVM.

Linear SVM = $x_i^T x_j$

Kernel SVM = $K(x_i, x_j)$