

Spring '21

Mixture Models

CS181

Data $D = \{x_1, \dots, x_n\}$ No target labels

Last lecture

1) K-means clustering

2) HAC

arguably a bit ad hoc, HAC is flexible but less useful in high dimensions

Mixture models

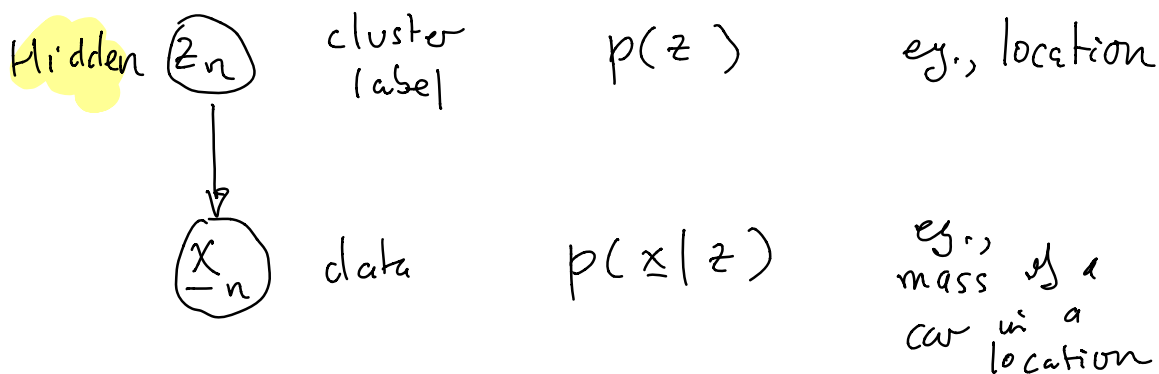
Probabilistic, generative view. Data $\{x_n\}$ as coming from a mixture of components of a distribution.



Use D to estimate parameters of a mixture model

cluster \equiv component

for x , cluster according to most likely component



$$p(x, z) = p(z) p(x|z)$$

Predict $p(z|x) \propto p(z) p(x|z)$

⊠ Gaussian Mixture Model (GMM)

Classes C_1, \dots, C_K $x_n \in \mathbb{R}^D$

$$p(z = C_k) = \theta_k \text{ categorical}$$

$$p(x|z = C_k) = N(x; \mu_k, \Sigma_k)$$



Parameters ω : $\{\theta_k\}_{k=1}^K, \{\mu_k, \Sigma_k\}_{k=1}^K$

\uparrow mean \uparrow covar

Max. Lik. Est.

$$\begin{aligned} \text{log likelihood } \mathcal{L}(x; \omega) &= \sum_n \ln p(x_n) \\ &= \sum_n \ln \sum_k p(x_n, C_k) \\ &= \sum_n \ln \sum_k p(C_k) p(x_n | C_k) \end{aligned}$$

$$= \sum_n \ln \sum_k \theta_k N(\underline{x}_n; \underline{\mu}_k, \underline{\Sigma}_k)$$

No analytical solution!

Log-sum structure prevents decomposing over the parameters.

ⓐ Consider a simpler, "complete data" problem. Suppose we observe class labels z_1, \dots, z_N .

Complete-data log likelihood (z_n 1-hot)

$$\begin{aligned} l(\underline{X}, \underline{z}; \underline{w}) &= \sum_n \ln p(\underline{x}_n, z_n) \\ &= \sum_n \ln p(z_n) p(\underline{x}_n | z_n) \\ &= \sum_n \ln p(z_n) + \sum_n \ln p(\underline{x}_n | z_n) \end{aligned}$$

$$p(z_n) = \prod_k \theta_k^{z_{nk}} \quad p(\underline{x}_n | z_n) = \prod_k N(\underline{x}_n; \underline{\mu}_k, \underline{\Sigma}_k)^{z_{nk}}$$

$$l(\underline{X}, \underline{z}; \underline{w}) = \sum_n \sum_k z_{nk} \ln \theta_k + \sum_n \sum_k z_{nk} \ln N(\underline{x}_n; \underline{\mu}_k, \underline{\Sigma}_k)$$

Solve analytically

$$\hat{\theta}_k = \frac{N_k}{N}$$

$$\hat{\underline{\mu}}_k = \text{mean}(\underline{x}_n)_{n \text{ s.t. } z_n = C_k}$$

$$\hat{\underline{\Sigma}}_k = \text{var}(\underline{x}_n)_{n \text{ s.t. } z_n = C_k}$$

$$N_k = \# \text{ examples in } C_k$$

Expectation-Maximization Algorithm

Iterative approach!

Initialize parameters $\{\theta_k\}$ $\{\mu_k, \Sigma_k\}$, repeat

- ① E-step. Use parameters to predict distribution on assignments z_n for each example
- ② M-step. Update parameters, maximizing expected complete data log likelihood given predicted assignments.

E-M algorithm. Very powerful, general method to maximize likelihood for models with latent variables (eg., assignments z).

M-step often has an analytical solution.

E-step Soft class assignment

$$p(z_n = c_k | \underline{x}_n) = q_{nk} \propto \theta_k N(\underline{x}_n; \underline{\mu}_k, \underline{\Sigma}_k)$$

M-step Work with the expected complete-data log likelihood:

$$\max_{\underline{\omega}} \mathbb{E}_{z \sim q} \left[\ln(\underline{x}, z; \underline{\omega}) \right]$$

$$\begin{aligned} \theta_1 &= \theta_2 = \frac{1}{2} \\ \underline{\mu}_1 &= \underline{\mu}_2 \\ \underline{\Sigma}_1 &= \underline{\Sigma}_2 \\ \text{Suppose} \end{aligned}$$

Analytical solution:

$$\hat{\theta}_k \leftarrow \frac{N_k}{N} \quad \hat{\underline{\mu}}_k = \frac{1}{N_k} \sum_n q_{nk} \underline{x}_n$$

$$N_k = \sum_n q_{nk} \quad \hat{\underline{\Sigma}}_k = \frac{1}{N_k} \sum_n q_{nk} (\underline{x}_n - \underline{\mu}_k)(\underline{x}_n - \underline{\mu}_k)^T$$

Properties

At each step, E-M algorithm adjusts parameters to improve likelihood of observed data \mathcal{D}

$$p(\mathcal{D}; \underline{\omega}^{(0)}) < p(\mathcal{D}; \underline{\omega}^{(1)}) < \dots$$

$\underline{\omega}^{(t)}$: parameters in step t

→ Converges to a local optimum (Restart can help)

Understanding M-step

Observed log likelihood

$$\ln p(\underline{x}_n) = \ln \sum_k \theta_k N(\underline{x}_n; \underline{\mu}_k, \underline{\Sigma}_k)$$

$$\begin{aligned} p(\underline{x}_n) &= \sum_k p(\underline{x}_n, c_k) \\ &= \sum_k p(z_n = c_k) p(\underline{x}_n | z_n = c_k) \end{aligned}$$

Complete-data log lik:

$$\begin{aligned} \ln p(\underline{x}_n, z_n) &= \ln p(z_n) p(\underline{x}_n | z_n) \\ &= \ln p(z_n) + \ln p(\underline{x}_n | z_n) \\ &= \ln \prod_k \theta_k^{z_{nk}} + \ln \prod_k N(\underline{x}_n; \underline{\mu}_k, \underline{\Sigma}_k)^{z_{nk}} \\ &= \sum_k z_{nk} \ln \theta_k + \sum_k z_{nk} \ln N(\underline{x}_n; \underline{\mu}_k, \underline{\Sigma}_k) \end{aligned}$$

Expected complete data log lik:

$$\begin{aligned} \mathbb{E}_{z_n} \ln p(\underline{x}_n, z_n) &= \sum_k q_{nk} \ln p(c_k) p(\underline{x}_n | c_k) \\ &= \sum_k q_{nk} \ln \theta_k + \sum_k q_{nk} \ln N(\underline{x}_n; \underline{\mu}_k, \underline{\Sigma}_k) \end{aligned}$$

M-step works with

$$\mathbb{E}_{z_n} \ln \prod_n [\cdot] = \mathbb{E}_{z_n} \sum_n \ln [\cdot] = \sum_n \mathbb{E}_{z_n} \ln [\cdot]$$

Notes

① Initialization matters!

eg, what happens if

(a) $\hat{\theta}_k = \frac{1}{K}$ all k , all $\hat{\Sigma}_k, \hat{\mu}_k$ same each k

(b) $\hat{\theta}_1 > \{\hat{\theta}_2, \dots, \hat{\theta}_K\}$,
all $\hat{\Sigma}_k, \hat{\mu}_k$ same each k

② Choose a variation, eg.

general $N(\underline{x}; \underline{\mu}_k, \underline{\Sigma}_k)$ θ # Param. μ Σ
K-1 KD KD²

diagonal $N(\underline{x}; \underline{\mu}_k, \underline{\Sigma}_k)$ K-1 KD KD
 \uparrow
diag.

spherical
(isotropic) $N(\underline{x}; \underline{\mu}_k, \sigma^2 \underline{I})$ K-1 KD K
(linear decision boundaries)

③ Connects to K-means!

K-means Repeats

1. Hard assign each example to nearest prototype μ_k
2. Update prototype to centroid of assigned points (mean)

E-M Fix class prob at $1/K$ with

Consider spherical Gaussian, ~~some~~
constant $\epsilon \in \mathbb{I}$ for $\epsilon > 0$, small.

E-step: Hard assign

M-step: Update μ_k to new
centroids of assigned
data.