

PGM - Week 2

Latent Variable Models and EM

the clustering problem

Consider a dataset $\{x_1, \dots, x_n\} \subseteq \mathbb{R}^d$. Our task is to partition this set into K subsets (K known for now). Intuitively, each subset (from now on referred to a cluster) of points should share some common patterns.

A natural soln for this is to define K prototypes $\{\mu_1, \mu_2, \dots, \mu_K\} \subseteq \mathbb{R}^d$ and determine the assignment of each datapoint x_n to each prototype μ_k according to a given criterion.

To solve this optimisation prob. we can define a set of binary variables r_{nk} , where $r_{nk} = 1$ iff x_n assigned to μ_k , and $r_{nk} = 0$ if not

Then, the objective can be

$$J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|^2$$

Goal: find r_{nk}, μ_k optimiser J . This is

$$r_{nk} = \begin{cases} 1 & \text{if } k = \underset{j}{\operatorname{argmin}} \|x_n - \mu_j\|^2 \\ 0 & \text{if not} \end{cases}$$

$$\mu_k = \frac{\sum_i r_{ik} x_i}{\sum_i r_{ik}}$$

\Rightarrow We found K-means: the rules above are applied iteratively - do they ensure convergence?

- Discuss: discrete / local minima / # of steps.

Drawbacks

- slow (to compute r_{nk} costs $O(nk)$)
- other dissimilarity.
- hard assignments

Uses: - decision making
- compression
-

Gaussian mixture model

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

We also introduce $z \in \{0, 1\}^K$, $\sum_k [z]_k = 1$ and/or

define

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

and also consider

$$p(z_k = 1) = \pi_k \quad (0 \leq \pi_k \leq 1, \sum \pi_k = 1)$$

the marginal p.m.f. is

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

likewise

$$p(x|z) = \prod_{k=1}^K N(\mu_k, \Sigma_k)^{z_k}$$

$$p(x) = \sum_z p(z) p(x|z) = \sum \pi_k N(\mu_k, \Sigma_k)$$

★) We have found an equivalent formulation of GMM with explicit latent variable

Another important quantity is $p(z|x)$

$$\begin{aligned} p(z_k) &\equiv p(z_k=1|x) = \frac{p(x|z_k=1) p(z_k=1)}{\sum_j p(x|z_j=1) p(z_j=1)} \\ &= \frac{\pi_k N(x|\mu_k, \Sigma_k)}{\sum_j \pi_j N(x|\mu_j, \Sigma_j)} \end{aligned}$$

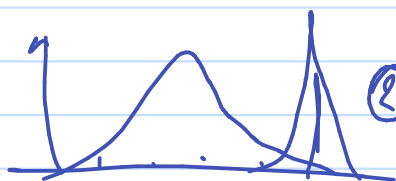
We will refer to this as the 'responsibility' of component to explain the obs x .

- Ancestral sampling to generate $x \sim \text{GMM}$

The loglike of GMM is $\theta = \{\pi, \mu, \Sigma\}$

$$\ell(\theta) = \log p(x|\mu, \pi, \Sigma) \quad (\star)$$

$$\textcircled{1} \quad = \sum_n \log \left(\sum_k \pi_k N(x_n | \mu_k, \Sigma_k) \right)$$



② Symmetries \rightarrow ill posed

log applies to Σ , not to $\log \rightarrow$ hard.

EM for GMMs. (we'll see this for the particular case of GMM, then integral)

Set the first order opt cond for μ_k

$$\mu_k = \frac{1}{N_k} \sum_n \underbrace{p_n}_{\gamma(z_{nk})} x_n, \quad N_k = \sum_{n=1}^n \gamma(z_{nk})$$

$$\Sigma_k = \frac{1}{N_k} \sum_n \gamma(z_{nk}) (x_n - \mu_k)(x_n - \mu_k)^T$$

$$\pi_k = N_k / N - \text{for this Laplace mlt}$$

Exercises: derive these

Remark: this is not a closed-form soln. since since we don't have $\gamma(z_{nk})$. Workaround: 2 steps

E: compute $\gamma(z_{nk}) = p(z_{nk} = z | x)$

M: use $\gamma(z_{nk})$ to compute μ, Σ, π (ML)

to sample: initialize with K-means.

Alternative view of EM

Let us now leave the GMM aside. In more general, perhaps abstract, terms, the goal of EM is to find ML solns for models involving latent variables.

Recall our notation X, Z

$$\ln p(X|\theta) = \ln \int p(X|\theta, z) p(z|\theta) dz$$

When can we optimise this? or even before that. in which cases can this be calculated? In very limited cases. Mixtures mix well with log, even in the exp. In other cases this is not. Mention Linear Gaussian

Let us consider the hypothetical scenario where we know the values of z (along with θ of X)

Def: $\{x_i, z_i\}$ will be the complete dataset, while $\{x_i\}$ will be the observed dataset. (imcomplete)

The related complete-data loglike would be

$$\begin{aligned}\log(p(X, z | \theta)) &= \log \prod p(x_i, z_i | \theta) \\ &= \log \prod p(x_i | z_i, \theta) p(z_i) \\ &= \log \prod \pi(z_i).\end{aligned}$$

easy to compute

This is impractical as z unknown. However one interesting interpretation of this is, to see the likelihood as a random fn, as it depends on z which is random.

A good approx of this would be to take its expectation and optimise that. (E)

Then with an explicit expression, we

maximise (μ) . So $\theta^{old} + E + \mu = \theta^{new}$ for θ .

E step: Use θ^{old} to compute $p(z|x, \theta^{old})$,
 compute expectation of the complete-data likelihood

$$Q(\theta, \theta^{old}) = \sum_z \ln p(x, z | \theta) \cdot p(z | x, \theta^{old})$$

M step: $\theta^{new} = \arg \max_{\theta} Q(\theta, \theta^{old})$

This can also be use for MAP (instead of ML)
 in which case $Q(\theta, \theta^{old}) + \log p(\theta) \leftarrow \text{prior}$.

Back to GMM

complete-data likelihood

$$p(X, Z | \mu, \Sigma, \pi) = \prod_{n=1}^N \prod_{k=1}^K \pi_k^{z_{nk}} p(x_n | \mu_k, \Sigma_k)^{z_{nk}}$$

$$\rightarrow l(\theta) = \sum_n \sum_k z_{nk} (\log \pi_k + \log N(x_n | \mu_k, \Sigma_k))$$

\rightarrow Tractable MLE

This is much easier: only one term in the

k -sum is non-zero $\rightarrow \mu, \Sigma$ are simple ests.

furthermore $\pi_k = \frac{1}{N} \sum_n z_{nk}$ (portion of data assumed to be that cluster)

Now let's see 2: $\propto p(z_{nk} | x_n, \theta)$

$$p(z | X, \theta) \propto p(X, z | \theta) = \prod_k \prod_n \overbrace{\pi_k N(x_n | \mu_k, \Sigma_k)}^{z_{nk}}$$

\rightarrow it factorizes over $n \rightarrow z_n$ all indep

so it looks like

$$p(z | X, \theta) = \prod_k \prod_n p(z_{nk} | x_n, \theta)$$

Now the expectation

binary var

$$\mathbb{E}[z_{nk}] = 1 \cdot p(z_{nk}=1 | x_n, \theta) + 0 \cdot p(z_{nk}=0 | x_n, \theta)$$

$$= p(z_{nk}=1 | x_n, \theta)$$

$$= \gamma(z_{nk}) \stackrel{\uparrow}{=} \frac{\pi_k N(x_n | \mu_k, \Sigma_k)}{\sum_j \pi_j N(x_n | \mu_j, \Sigma_j)}$$

Finally

$$\rightarrow \mathbb{E}_z \log p(X, z | \theta) = \sum_k \sum_n \gamma(z_n) (\log \pi_k + \log N(x_n | \mu_k, \Sigma_k))$$

Intuition says that EM recovers K-mean

Exercise: show that K-mean is recovered

1) choosing

$$p(X|\mu_k, \Sigma_k) = N(X|\mu_k, I\epsilon)$$

$\hookrightarrow \epsilon > 0$ fixed.

to show a soft-assignment version of K-mean.

take $\epsilon \rightarrow 0$ and see what happens.

Recommendation: See the examples for EM for Bayesian linear reg an mixture of Bernoulli for completion..

EM in general

$$p(X|\theta) = \sum_z p(X, z|\theta)$$

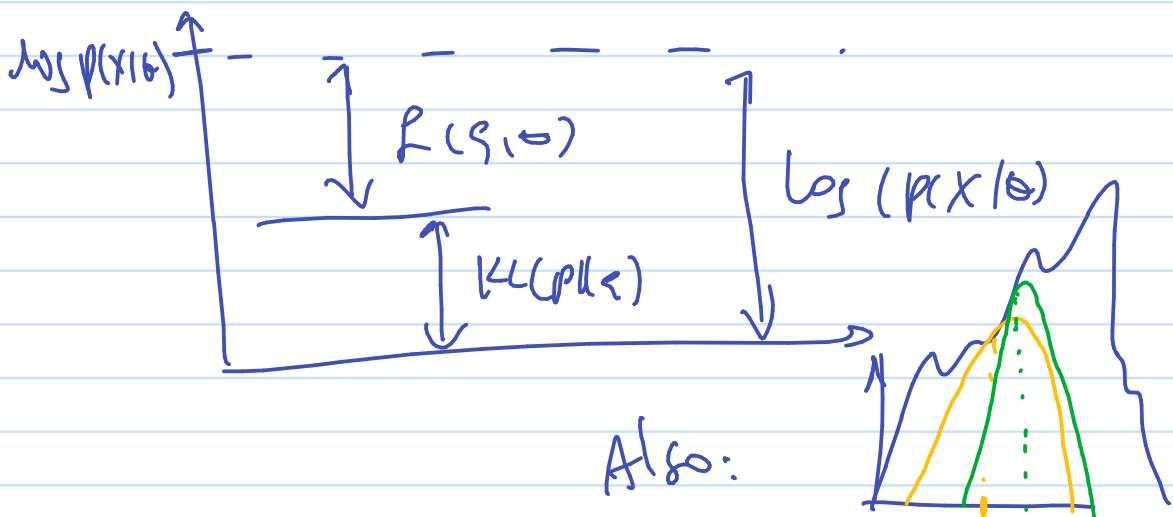
difficult \nearrow \searrow easy

Introduce a distribution $q(z)$ over z . For any choice of $q(\cdot)$ the following holds

$$\log p(X|\theta) = \sum_z q(z) \cdot \log p(X|\theta)$$

\longrightarrow

$$\begin{aligned}
\log p(X|\theta) &= \sum_k q(z_k) \log p(X|\theta) \\
&= \sum_k q(z_k) \log \left(\frac{\overbrace{p(X|\theta) p(z|X, \theta) \cdot q(z)}^{p(X, z|\theta)}}{p(z|X, \theta) \cdot q(z)} \right) \\
&= \sum_k q(z_k) \log \left(\frac{p(X, z|\theta)}{q(z)} \cdot \left(\frac{p(z|X, \theta)}{q(z)} \right) \right) \\
&= \sum_k q(z_k) \log \left(\frac{p(X, z|\theta)}{q(z)} \right) \\
&\quad - \sum_k q(z_k) \log \left(\frac{q(z)}{p(z|X, \theta)} \right) \\
&= \mathcal{L}(q, \theta) + \text{KL}(p \| q)
\end{aligned}$$



EM Breaks the prob of ML into two simpler problems

- compute an approx to $\log p(X|\theta)$
- maximise this approximation.

in this problem, this actually achieves ML (locally)

Next: $q(z) = p(z|X, \theta)$ is unattainable, so we'll try to get closer.