

Expectation Maximization

CS771: Introduction to Machine Learning

Nisheeth

Building intuition for EM: a case study

Let events be “grades in a class”

| | |
|-------------------|-----------------------------|
| $w_1 =$ Gets an A | $P(A) = \frac{1}{2}$ |
| $w_2 =$ Gets a B | $P(B) = \mu$ |
| $w_3 =$ Gets a C | $P(C) = 2\mu$ |
| $w_4 =$ Gets a D | $P(D) = \frac{1}{2} - 3\mu$ |

(Note $0 \leq \mu \leq 1/6$)

Assume we want to estimate μ from data. In a given class there were

- a A's
- b B's
- c C's
- d D's

What's the maximum likelihood estimate of μ given a,b,c,d ?



Max likelihood solution

$$P(A) = \frac{1}{2} \quad P(B) = \mu \quad P(C) = 2\mu \quad P(D) = \frac{1}{2} - 3\mu$$

$$P(a, b, c, d \mid \mu) = K \left(\frac{1}{2}\right)^a (\mu)^b (2\mu)^c \left(\frac{1}{2} - 3\mu\right)^d$$

$$\log P(a, b, c, d \mid \mu) = \log K + a \log \frac{1}{2} + b \log \mu + c \log 2\mu + d \log \left(\frac{1}{2} - 3\mu\right)$$

$$\text{FOR MAX LIKE } \mu, \text{ SET } \frac{\partial \text{LogP}}{\partial \mu} = 0$$

$$\frac{\partial \text{LogP}}{\partial \mu} = \frac{b}{\mu} + \frac{2c}{2\mu} - \frac{3d}{1/2 - 3\mu} = 0$$

$$\text{Gives max like } \mu = \frac{b + c}{6(b + c + d)}$$

So if class got

| A | B | C | D |
|----|---|---|----|
| 14 | 6 | 9 | 10 |

$$\text{Max like } \mu = \frac{1}{10}$$



Same Problem with Hidden Information

Someone tells us that

Number of High grades (A's + B's) = h

Number of C's = c

Number of D's = d

What is the max. like estimate of μ now?

REMEMBER

$$P(A) = \frac{1}{2}$$

$$P(B) = \mu$$

$$P(C) = 2\mu$$

$$P(D) = \frac{1}{2} - 3\mu$$



Same Problem with Hidden Information

Someone tells us that

Number of High grades (A's + B's) = h

Number of C's = c

Number of D's = d

What is the max. like estimate of μ now?

We can answer this question circularly:

REMEMBER
 $P(A) = \frac{1}{2}$
 $P(B) = \mu$
 $P(C) = 2\mu$
 $P(D) = \frac{1}{2} - 3\mu$

EXPECTATION

If we know the value of μ we could compute the expected value of a and b

Since the ratio $a:b$ should be the same as the ratio $\frac{1}{2} : \mu$

$$a = \frac{\frac{1}{2}}{\frac{1}{2} + \mu} h \quad b = \frac{\mu}{\frac{1}{2} + \mu} h$$

MAXIMIZATION

If we know the expected values of a and b we could compute the maximum likelihood value of μ

$$\mu = \frac{b + c}{6(b + c + d)}$$



EM for our problem

We begin with a guess for μ

We iterate between EXPECTATION and MAXIMIZATION to improve our estimates of μ and a and b .

REMEMBER

$$P(A) = \frac{1}{2}$$

$$P(B) = \mu$$

$$P(C) = 2\mu$$

$$P(D) = \frac{1}{2} - 3\mu$$

Define $\mu(t)$ the estimate of μ on the t 'th iteration

$b(t)$ the estimate of b on t 'th iteration

$\mu(0)$ = initial guess

$$b(t) = \frac{\mu(t)h}{\frac{1}{2} + \mu(t)} = E[b \mid \mu(t)]$$

$$\begin{aligned} \mu(t+1) &= \frac{b(t) + c}{6(b(t) + c + d)} \\ &= \text{max like est of } \mu \text{ given } b(t) \end{aligned}$$

E -step

M -step



EM Convergence

- Convergence proof based on fact that $\text{Prob}(\text{data} \mid \mu)$ must increase or remain same between each iteration
[NOT OBVIOUS]
 - But it can never exceed 1 [OBVIOUS]
- So it must therefore converge [OBVIOUS]

In our example, suppose we had

$$h = 20$$

$$c = 10$$

$$d = 10$$

$$\mu(0) = 0$$



Convergence is generally linear: error decreases by a constant factor each time step.

| t | $\mu(t)$ | b(t) |
|---|----------|-------|
| 0 | 0 | 0 |
| 1 | 0.0833 | 2.857 |
| 2 | 0.0937 | 3.158 |
| 3 | 0.0947 | 3.185 |
| 4 | 0.0948 | 3.187 |
| 5 | 0.0948 | 3.187 |
| 6 | 0.0948 | 3.187 |



ALT-OPT/EM for Gaussian Mixture Model



MLE for Gaussian Discriminant Analysis

- Assume a K class generative classification model with Gaussian class-conditionals
- Assume class $k = 1, 2, \dots, K$ is modeled by a Gaussian with mean μ_k and cov matrix Σ_k
- Can assume label y_n to be one-hot and then $y_{nk} = 1$ if $y_n = k$, and $y_{nk} = 0$, o/w
- Assuming class prior as $p(y_n = k) = \pi_k$, the model has params $\Theta = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$
- Given training data $\{\mathbf{x}_n, y_n\}_{n=1}^N$, the MLE solution will be

$$\hat{\pi}_k = \frac{1}{N} \sum_{n=1}^N y_{nk}$$

Same as $\frac{N_k}{N}$ where N_k is # of training ex. for which $y_n = k$

$$\hat{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N y_{nk} \mathbf{x}_n$$

Same as $\frac{1}{N_k} \sum_{n:y_n=k}^N \mathbf{x}_n$

$$\hat{\Sigma}_k = \frac{1}{N_k} \sum_{n=1}^N y_{nk} (\mathbf{x}_n - \hat{\mu}_k)(\mathbf{x}_n - \hat{\mu}_k)^\top$$

Same as $\frac{1}{N_k} \sum_{n:y_n=k}^N (\mathbf{x}_n - \hat{\mu}_k)(\mathbf{x}_n - \hat{\mu}_k)^\top$

See [here](#) for a derivation of the MLE for GDA

Observations on the GDA objective function

- Here is a formal derivation of the MLE solution for $\Theta = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$

$$\begin{aligned}
 \hat{\Theta} &= \operatorname{argmax}_{\Theta} p(\mathbf{X}, \mathbf{y} | \Theta) = \operatorname{argmax}_{\Theta} \prod_{n=1}^N p(\mathbf{x}_n, y_n | \Theta) \begin{array}{l} \text{multinoulli} \\ \text{Gaussian} \end{array} \\
 &= \operatorname{argmax}_{\Theta} \prod_{n=1}^N p(y_n | \Theta) p(\mathbf{x}_n | y_n, \Theta) \\
 &= \operatorname{argmax}_{\Theta} \prod_{n=1}^N \prod_{k=1}^K \pi_k^{y_{nk}} \prod_{k=1}^K p(\mathbf{x}_n | y_n = k, \Theta)^{y_{nk}} \\
 &= \operatorname{argmax}_{\Theta} \prod_{n=1}^N \prod_{k=1}^K [\pi_k p(\mathbf{x}_n | y_n = k, \Theta)]^{y_{nk}} \\
 &= \operatorname{argmax}_{\Theta} \log \prod_{n=1}^N \prod_{k=1}^K [\pi_k p(\mathbf{x}_n | y_n = k, \Theta)]^{y_{nk}} \\
 &= \operatorname{argmax}_{\Theta} \sum_{n=1}^N \sum_{k=1}^K y_{nk} [\log \pi_k + \log \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)]
 \end{aligned}$$

In general, in models with probability distributions from the **exponential family**, the MLE problem will usually have a simple analytic form

Also, due to the form of the likelihood (Gaussian) and prior (multinoulli), the MLE problem had a nice separable structure after taking the log

Can see that, when estimating the parameters of the k^{th} Gaussian (π_k, μ_k, Σ_k) , we only will only need training examples from the k^{th} class, i.e., examples for which $y_{nk} = 1$

The form of this expression is important; will encounter this in GMM too



Need for EM/ALT-OPT: Two Equivalent Perspectives

1. Consider an LVM with latent variables and parameters. Trying to estimate parameters without also estimating the latent variables (by marginalizing them) is difficult

$$p(\mathbf{x}_n|\Theta) = \sum_{k=1}^K p(\mathbf{x}_n, \mathbf{z}_n = k|\Theta) = \sum_{k=1}^K p(\mathbf{z}_n = k|\phi)p(\mathbf{x}_n|\mathbf{z}_n = k, \theta) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n|\mu_k, \Sigma_k)$$

A Gaussian Mixture Model (GMM)

MLE for GMM with cluster ids marginalized/summed/integrated out

$$\Theta_{MLE} = \operatorname{argmax}_{\Theta} \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n|\mu_k, \Sigma_k)$$

Can't get closed form expressions for the π_k, μ_k, Σ_k due to "log of sum". Have to use gradient based methods

This issue not just for MLE for GMM but MLE for other LVMs too

EM/ALT-OPT will help us "simulate" this condition by making guesses about the values of \mathbf{z}_n 's

If we knew the \mathbf{z}_n 's, the problem will be much simpler; just like MLE for generative classification with Gaussian class-conditional

Since no marginalization of the \mathbf{z}_n 's required

2. Consider a complex prob. density (without any latent vars) for which MLE is hard

Directly defining a probability density as a mixture of Gaussians (\mathbf{x}_n is generated by the k^{th} Gaussian with probability π_k) without any reference to any latent variable whatsoever (we didn't define it as an LVM)

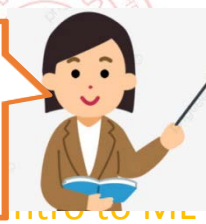
$$p(\mathbf{x}_n|\Theta) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n|\mu_k, \Sigma_k)$$

MLE for the params Θ of this distribution will again be hard (as we already saw above). However, we can artificially introduce a latent variable \mathbf{z}_n with each data point \mathbf{x}_n , denoting which Gaussian generated \mathbf{x}_n

Can now apply ALT-OPT/EM to estimate parameters Θ + we get the latent variables \mathbf{z}_n as a "by-product" (though we may not be interested in learning \mathbf{z}_n 's if our goal is just density estimation, not clustering)

Even though we didn't need the artificially introduced \mathbf{z}_n 's, their presence and doing ALT-OPT/EM made our job of estimating Θ easier!

Also in any LVM, given Θ , you can always estimate \mathbf{z}_n 's. Likewise, given \mathbf{z}_n , you can always estimate Θ



MLE for GMM

- Already saw that MLE is hard for GMM

$$\Theta_{MLE} = \operatorname{argmax}_{\Theta} \log p(\mathbf{X}|\Theta) = \operatorname{argmax}_{\Theta} \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)$$

- Two possible ways to solve this MLE problem

Will soon see how to get these guesses

- If someone gave us optimal “point” guesses $\hat{\mathbf{z}}_n$ ’s of cluster ids \mathbf{z}_n ’s, we could do MLE for the parameters just like we did for generative classification with Gaussian class-conditionals

$$\Theta_{MLE} = \operatorname{argmax}_{\Theta} \log p(\mathbf{X}, \hat{\mathbf{Z}} | \Theta) = \operatorname{argmax}_{\Theta} \sum_{n=1}^N \sum_{k=1}^K \hat{z}_{nk} [\log \pi_k + \log \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)]$$

In form of a probability distribution instead of a single “optimal” guess

- Alternatively, if someone gave a “probabilistic” guess of \mathbf{z}_n ’s, we can do MLE for Θ as follows

$$\Theta_{MLE} = \operatorname{argmax}_{\Theta} \mathbb{E}[\log p(\mathbf{X}, \mathbf{Z} | \Theta)] = \operatorname{argmax}_{\Theta} \sum_{n=1}^N \sum_{k=1}^K \mathbb{E}[z_{nk}] [\log \pi_k + \log \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)]$$

Similar to Approach 1 but maximizes an expectation

The expectation is w.r.t a distribution of \mathbf{Z} which we will see shortly

- Approach 1 is ALT-OPT and Approach 2 is Expectation Maximization (“soft” ALT-OPT).

Both require alternating between estimating \mathbf{Z} and Θ until convergence

ALT-OPT for GMM

Keep in mind: In LVMs, assuming i.i.d. data, the quantity $\log p(\mathbf{X}|\theta) = \sum_{n=1}^N \log p(\mathbf{x}_n|\theta)$ is called **incomplete data log-likelihood (ILL)** whereas $\log p(\mathbf{X}, \mathbf{Z}|\theta) = \sum_{n=1}^N \log p(\mathbf{x}_n, \mathbf{z}_n|\theta)$ is called **complete data log-likelihood (CLL)**. Goal is to maximize ILL but ALT-OPT maximizes CLL (EM too will maximize the expectation of CLL). The latent vars \mathbf{z}_n 's "complete" the data \mathbf{x}_n ☺



- We will assume we have a "hard" (most probable) guess of \mathbf{z}_n , say $\hat{\mathbf{z}}_n$

- Then ALT-OPT would look like this

- Initialize $\Theta = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$ as $\hat{\Theta}$
- Repeat the following until convergence

- For each n , compute most probable value (our best guess) of \mathbf{z}_n as

$$\hat{\mathbf{z}}_n = \operatorname{argmax}_{k=1,2,\dots,K} p(\mathbf{z}_n = k | \hat{\Theta}, \mathbf{x}_n)$$

Proportional to prior prob times likelihood, i.e.,
 $p(\mathbf{z}_n = k | \theta) p(\mathbf{x}_n | \mathbf{z}_n = k, \theta) = \pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)$

Posterior probability of point \mathbf{x}_n belonging to cluster k

- Solve MLE problem for Θ using most probable \mathbf{z}_n 's

$$\hat{\Theta} = \operatorname{argmax}_{\Theta} \sum_{n=1}^N \sum_{k=1}^K \hat{z}_{nk} [\log \pi_k + \log \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)]$$

Note: The objective function is $\sum_{n=1}^N \log p(\mathbf{x}_n, \hat{\mathbf{z}}_n | \theta) = \sum_{n=1}^N \log p(\hat{\mathbf{z}}_n | \theta) + \log p(\mathbf{x}_n | \hat{\mathbf{z}}_n, \theta)$

Same objective function as generative K -class classification with Gaussian class-conditionals

N_k : Effective number of points in cluster k

$$\hat{\pi}_k = \frac{1}{N} \sum_{n=1}^N \hat{z}_{nk}$$

$$\hat{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N \hat{z}_{nk} \mathbf{x}_n$$

$$\hat{\Sigma}_k = \frac{1}{N_k} \sum_{n=1}^N \hat{z}_{nk} (\mathbf{x}_n - \hat{\mu}_k)(\mathbf{x}_n - \hat{\mu}_k)^\top$$

Does that matter? Should we worry that we aren't solving the actual problem anymore?

Not really; will see the justification soon ☺

But wait! This is not the same as $\sum_{n=1}^N \log p(\mathbf{x}_n | \theta)$ which was the original MLE objective for this LVM ☹



Expectation-Maximization (EM) for GMM

.. which we maximized in ALT-OPT

Expectation of CLL

- EM finds Θ_{MLE} by maximizing $\mathbb{E}[\log p(\mathbf{X}, \mathbf{Z}|\Theta)]$ rather than $\log p(\mathbf{X}, \hat{\mathbf{Z}}|\Theta)$
- Note: Expectation will be w.r.t. the conditional posterior distribution of \mathbf{Z} , i.e., $p(\mathbf{Z}|\mathbf{X}, \Theta)$
- The EM algorithm for GMM operates as follows
 - Initialize $\Theta = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$ as $\hat{\Theta}$
 - Repeat until convergence
 - Compute conditional posterior $p(\mathbf{Z}|\mathbf{X}, \hat{\Theta})$. Since obs are i.i.d, compute separately for each n (and for $k = 1, 2, \dots, K$)

It is "conditional" posterior because it is also conditioned on Θ , not just data \mathbf{X}

Why w.r.t. this distribution? Will see justification in a bit

Needed to get the expected CLL

Requires knowing Θ

Same as $p(z_{nk} = 1 | \mathbf{x}_n, \hat{\Theta})$, just a different notation

$$p(\mathbf{z}_n = k | \mathbf{x}_n, \hat{\Theta}) \propto p(\mathbf{z}_n = k | \hat{\Theta}) p(\mathbf{x}_n | \mathbf{z}_n = k, \hat{\Theta}) = \hat{\pi}_k \mathcal{N}(\mathbf{x}_n | \hat{\mu}_k, \hat{\Sigma}_k)$$

- Update Θ by maximizing the expected complete data log-likelihood

Solution has a similar form as ALT-OPT (or gen. class.), except we now have the **expectation** of z_{nk} being used

N_k : Effective number of points in cluster k

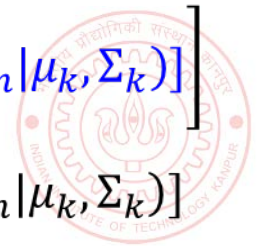
$$\hat{\pi}_k = \frac{1}{N} \sum_{n=1}^N \mathbb{E}[z_{nk}] \quad \hat{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N \mathbb{E}[z_{nk}] \mathbf{x}_n$$

$$\hat{\Sigma}_k = \frac{1}{N_k} \sum_{n=1}^N \mathbb{E}[z_{nk}] (\mathbf{x}_n - \hat{\mu}_k)(\mathbf{x}_n - \hat{\mu}_k)^\top$$

$$\hat{\Theta} = \operatorname{argmax}_{\Theta} \mathbb{E}_{p(\mathbf{Z}|\mathbf{X}, \hat{\Theta})} [\log p(\mathbf{X}, \mathbf{Z}|\Theta)] = \sum_{n=1}^N \mathbb{E}_{p(\mathbf{z}_n|\mathbf{x}_n, \hat{\Theta})} [\log p(\mathbf{x}_n, \mathbf{z}_n|\Theta)]$$

$$= \operatorname{argmax}_{\Theta} \mathbb{E} \left[\sum_{n=1}^N \sum_{k=1}^K z_{nk} [\log \pi_k + \log \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)] \right]$$

$$= \operatorname{argmax}_{\Theta} \sum_{n=1}^N \sum_{k=1}^K \mathbb{E}[z_{nk}] [\log \pi_k + \log \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)]$$



EM for GMM (Contd)

- The EM algo for GMM required $\mathbb{E}[z_{nk}]$. Note $z_{nk} \in \{0,1\}$

Reason: $\sum_{k=1}^K \gamma_{nk} = 1$

Need to normalize: $\mathbb{E}[z_{nk}] = \frac{\hat{\pi}_k \mathcal{N}(x_n | \hat{\mu}_k, \hat{\Sigma}_k)}{\sum_{\ell=1}^K \hat{\pi}_\ell \mathcal{N}(x_n | \hat{\mu}_\ell, \hat{\Sigma}_\ell)}$

$$\mathbb{E}[z_{nk}] = \gamma_{nk} = 0 \times p(z_{nk} = 0 | x_n, \hat{\Theta}) + 1 \times p(z_{nk} = 1 | x_n, \hat{\Theta}) = p(z_{nk} = 1 | x_n, \hat{\Theta}) \propto \hat{\pi}_k \mathcal{N}(x_n | \hat{\mu}_k, \hat{\Sigma}_k)$$

EM for Gaussian Mixture Model

1 Initialize $\Theta = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$ as $\Theta^{(0)}$, set $t = 1$

2 E step: compute the expectation of each z_n (we need it in M step)

Accounts for fraction of points in each cluster

Accounts for cluster shapes (since each cluster is a Gaussian)

Soft K-means, which are more of a heuristic to get soft-clustering, also gave us probabilities but didn't account for cluster shapes or fraction of points in each cluster

$$\mathbb{E}[z_{nk}^{(t)}] = \gamma_{nk}^{(t)} = \frac{\pi_k^{(t-1)} \mathcal{N}(x_n | \mu_k^{(t-1)}, \Sigma_k^{(t-1)})}{\sum_{\ell=1}^K \pi_\ell^{(t-1)} \mathcal{N}(x_n | \mu_\ell^{(t-1)}, \Sigma_\ell^{(t-1)})} \quad \forall n, k$$

3 Given "responsibilities" $\gamma_{nk} = \mathbb{E}[z_{nk}]$, and $N_k = \sum_{n=1}^N \gamma_{nk}$, re-estimate Θ via MLE

Effective number of points in the k^{th} cluster

$$\mu_k^{(t)} = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk}^{(t)} x_n$$

M-step:

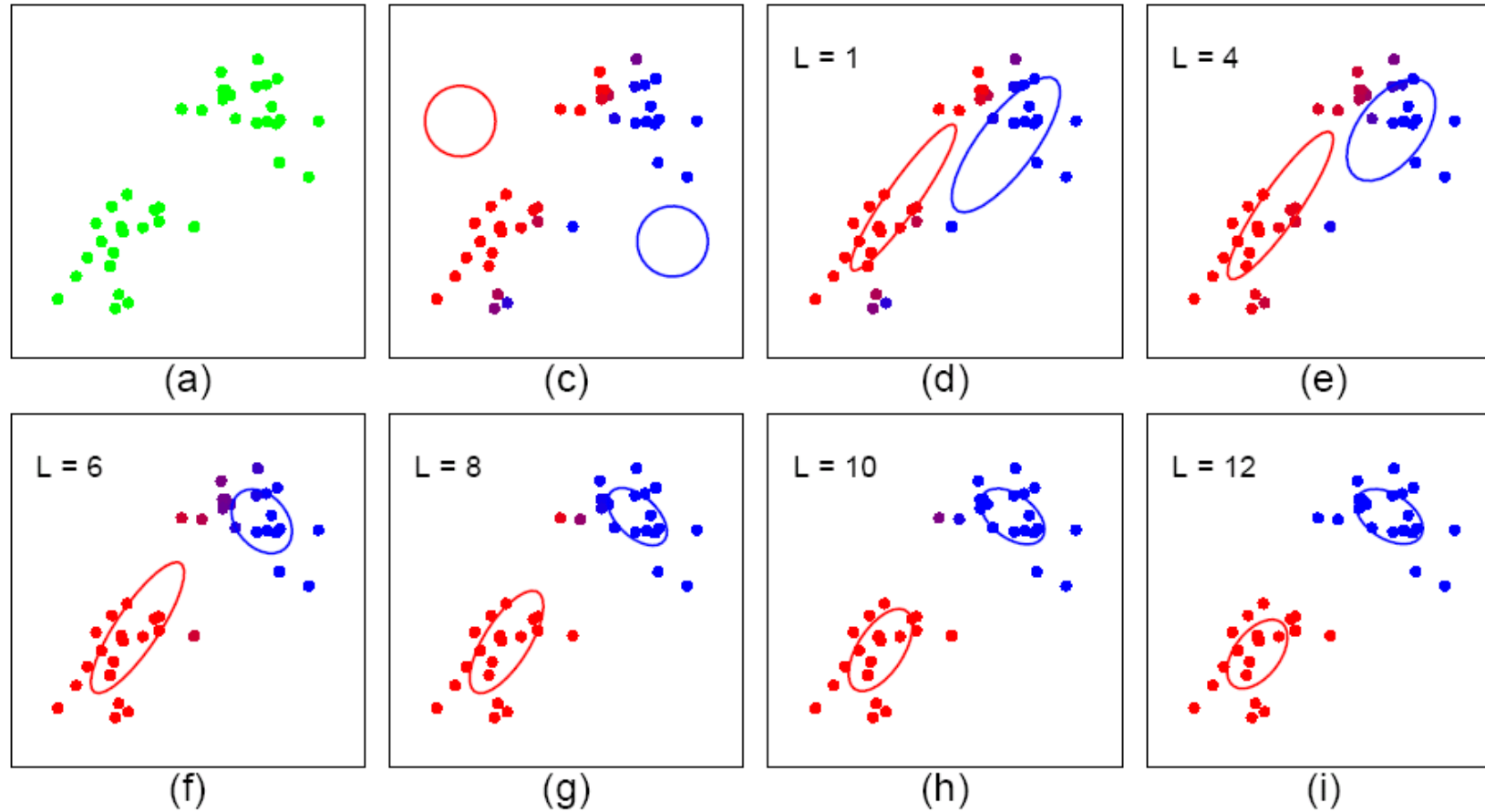
$$\Sigma_k^{(t)} = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk}^{(t)} (x_n - \mu_k^{(t)})(x_n - \mu_k^{(t)})^T$$

$$\pi_k^{(t)} = \frac{N_k}{N}$$

4 Set $t = t + 1$ and go to step 2 if not yet converged



EM for GMM in action



Note: Just like with k-means, cluster initialization matters. EM only finds local optima.

