# LVMs (Contd), Expectation Maximization (1)

CS771: Introduction to Machine Learning

Piyush Rai

# Plan

- ALT-OPT and EM
  - Example: Gaussian Mixture Model for data clustering
- A deeper look at ALT-OPT and EM
- General recipe for doing ALT-OPT and EM for any LVM

# 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

A Gaussian Mixture Model (GMM)

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

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

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

Can't get closed form expressions for the $\pi_k, \mu_k, \Sigma_k$ 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 $\boldsymbol{z}_n$'s

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

Since no marginalization of the $\boldsymbol{z}_n$'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 ($\boldsymbol{x}_n$ is generated by the $k^{th}$ Gaussian with probability $\pi_k$) without any reference to any latent variable whatsoever (we didn't define it as an LVM)

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

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 $\boldsymbol{z}_n$ with each data point $\boldsymbol{x}_n$, denoting which Gaussian generated $\boldsymbol{x}_n$

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

Now this prob. density estimation problem also becomes Problem 1 above - a clustering problem with latent variables

Even though we didn't need the artificially introduced $\boldsymbol{z}_n$'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 $\boldsymbol{z}_n$'s. Likewise, given $\boldsymbol{z}_n$, you can always estimate Θ

Remember that GMM is just like generative classification with Gaussian class-conditionals and training data labels unknown

# ALT-OPT/EM for Gaussian Mixture Model

# Detour: MLE for Generative Classification

- Assume a $K$ class generative classification model with Gaussian class-conditionals
- Assume class $k = 1, 2, \ldots, K$ is modeled by a Gaussian with mean $\mu_k$ and cov matrix $\Sigma_k$
- 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 $\{x_n, y_n\}_{n=1}^{N}$, 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 $y_n = k$

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

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

$$\hat{\Sigma}_k = \frac{1}{N_k} \sum_{n=1}^{N} y_{nk} (x_n - \hat{\mu}_k)(x_n - \hat{\mu}_k)^{\mathsf{T}}$$

Same as $\frac{1}{N_k} \sum_{n:y_n=k}^{N} (x_n - \hat{\mu}_k)(x_n - \hat{\mu}_k)^{\mathsf{T}}$

# Detour: MLE for Generative Classification

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

$$\widehat{\Theta} = \text{argmax}_\Theta\, p(\boldsymbol{X}, \boldsymbol{y}|\Theta) = \text{argmax}_\Theta \prod_{n=1}^N p(\boldsymbol{x}_n, y_n|\Theta)$$

multinoulli    Gaussian

$$= \text{argmax}_\Theta \prod_{n=1}^N p(y_n|\Theta)\, p(x_n|y_n, \Theta)$$

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

$$= \text{argmax}_\Theta \prod_{n=1}^N \prod_{k=1}^K \pi_k^{y_{nk}} \prod_{k=1}^K p(x_n|y_n = k, \Theta)^{y_{nk}}$$

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

$$= \text{argmax}_\Theta \prod_{n=1}^N \prod_{k=1}^K [\pi_k p(x_n|y_n = k, \Theta)]^{y_{nk}}$$

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

$$= \text{argmax}_\Theta \log \prod_{n=1}^N \prod_{k=1}^K [\pi_k p(x_n|y_n = k, \Theta)]^{y_{nk}}$$

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

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

# Detour: Exponential Family

Exp-fam dist also used for Genealized Linear Models (GLM) with $p(y|\boldsymbol{x}, \boldsymbol{w})$ modeled by an exp-fam distribution whose natural parameter is defined by $\boldsymbol{w}^\top \boldsymbol{x}$ (thus "linear"). Useful in problems where $y$ is not real/categorical but a count, or positive real, etc

- Exponential Family is a family of prob. distributions that have the form

Lin reg, logistic reg, softmax reg are also instances of GLMs

$$p(x|\theta) = h(x)\exp[\theta^\top T(x) - A(\theta)]$$

Even though their standard form may not look like this, they can be rewritten in this form after some algebra

- Many well-known distribution (Bernoulli, Binomial, multinoulli, Poisson, beta, gamma, Gaussian, etc.) are examples of exponential family distributions

- $\theta$ is called the natural parameter of the family

Natural params are a function of the distribution parameters in the standard form

- $h(x), T(x),$ and $A(\theta)$ are known functions (specific to the distribution)

- $T(x)$ is called the sufficient statistics: estimates of $\theta$ contain $x$ in form of suff-stats

- Every exp. family distribution also has a conjugate distribution (often also in exp. family)

- Also, MLE/MAP is usually quite simple since $\log p(x|\theta)$ will have a simple expression

- Also useful in fully Bayesian inference since they have conjugate priors

https://en.wikipedia.org/wiki/Exponential_family

# MLE for GMM

- Already saw that MLE is hard for GMM

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

- Two possible ways to solve this MLE problem

  > Will soon see how to get these guesses

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

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

  > In form of a probability distribution instead of a singe "optimal" guess

  2. Alternatively, if someone gave a "probabilistic" guess of $\boldsymbol{z}_n$'s, we can do MLE for $\Theta$ as follows

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

  > Similar to Approach 1 but maximizes an expectation

  > The expectation is w.r.t a distribution of 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 $\boldsymbol{Z}$ and $\Theta$ until convergence

# ALT-OPT for GMM

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

- Then ALT-OPT would look like this

  - Initialize $\Theta = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^{K}$ as $\widehat{\Theta}$

    Proportional to prior prob times likelihood, i.e.,
    $p(z_n = k|\Theta) \, p(x_n|z_n = k, \Theta) = \pi_k \, \mathcal{N}(x_n|\mu_k, \Sigma_k)$

  - Repeat the following until convergence

    Posterior probability of point $x_n$ belonging to cluster $k$

    - For each $n$, compute most probable value (our best guess) of $z_n$ as

$$\hat{z}_n = \text{argmax}_{k=1,2,\ldots,K} \; p(z_n = k|\widehat{\Theta}, x_n)$$

    - Solve MLE problem for $\Theta$ using most probable $z_n$'s

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

$$\widehat{\Theta} = \text{argmax}_{\Theta} \; \sum_{n=1}^{N} \sum_{k=1}^{K} \hat{z}_{nk}[\log \pi_k + \log \mathcal{N}(x_n|\mu_k, \Sigma_k)]$$

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

$N_k$ : Effective number of points in cluster k

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

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

$\hat{\Sigma}_k = \frac{1}{N_k} \sum_{n=1}^{N} \hat{z}_{nk}(x_n - \hat{\mu}_k)(x_n - \hat{\mu}_k)^{\mathsf{T}}$

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

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

Not really; will see the justification soon ☺

# Expectation-Maximization (EM) for GMM

.. which we maximized in ALT-OPT

Expectation of CLL

- EM finds $\Theta_{MLE}$ by maximizing $\mathbb{E}[\log p(X, Z|\Theta)]$ rather than $\log p(X, \widehat{Z}|\Theta)$

- Note: Expectation will be w.r.t. the <u>conditional</u> posterior distribution of $Z$, i.e., $p(Z|X, \Theta)$

It is "conditional" posterior because it is also conditioned on $\Theta$, not just data $X$

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

- The EM algorithm for GMM operates as follows
    - Initialize $\Theta = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^{K}$ as $\widehat{\Theta}$
    - Repeat until convergence

    Needed to get the expected CLL

    Requires knowing $\Theta$

        - Compute conditional posterior $p(Z|X, \widehat{\Theta})$. Since obs are i.i.d, compute separately for each $n$ (and for $k = 1,2,..K$)

Same as $p(z_{nk} = 1| x_n, \widehat{\Theta})$, just a different notation

$$p(z_n = k|x_n, \widehat{\Theta}) \propto p(z_n = k|\widehat{\Theta}) \, p(x_n|z_n = k, \widehat{\Theta}) = \hat{\pi}_k \mathcal{N}(x_n|\hat{\mu}_k, \hat{\Sigma}_k)$$

        - Update $\Theta$ 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

$$\widehat{\Theta} = \text{argmax}_{\Theta} \mathbb{E}_{p(Z|X, \widehat{\Theta})}[\log p(X, Z|\Theta)] = \sum_{n=1}^{N} \mathbb{E}_{p(z_n|x_n, \widehat{\Theta})}[\log p(x_n, z_n|\Theta)]$$

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

$N_k$ : Effective number of points in cluster k

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

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

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

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

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

## EM for Gaussian Mixture Model

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

❷ 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

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

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

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

M-step:

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

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

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

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