GP (contd), Latent Variable Models and EM Algorithm

CS772A: Probabilistic Machine Learning

Piyush Rai

GP Prediction with Gaussian Likelihood

In general, the PPD when using GP is defined as

 $p(y_*|\mathbf{y}) = \int p(y_*|f_*) p(f_*|\mathbf{f}) p(\mathbf{f}|\mathbf{y}) d\mathbf{f} df_*$

And don't even have to compute/use the posterior $p(\mathbf{f}|\mathbf{y})$ (which in this case is a Gaussian by the way O) to get the PPD

For Gaussian likelihood (and fixed hyperparams), we don't need to do above integral

Assuming zero Reason: The marginal likelihood is Gaussian GP prior mean function $\mathbf{p}(\mathbf{f}) = \mathcal{N}(\mathbf{f} \mid \mathbf{0}, \mathbf{K})$ Gaussian likelihood $p(\mathbf{y}|\mathbf{f}) = \mathcal{N}(\mathbf{y}|\mathbf{f}, \beta^{-1}\mathbf{I}_N)$ (assuming β is fixed) $p(\mathbf{y}) = \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f})d\mathbf{f} = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{K} + \beta^{-1}\mathbf{I}_N) = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{C}_N)$ Marginal likelihood of training outputs Marginal likelihood of training and test outputs $p\left(\begin{bmatrix} y \\ y_* \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} y \\ y_* \end{bmatrix} \mid \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} C_N & K_* \\ K_* & \kappa(x_*, x_*) + \beta^{-1} \end{bmatrix}\right)$ Marginal likelihood of PPD obtained using joint to conditional $p(y_*|\mathbf{y}) = \mathcal{N}(y_*|\mathbf{k}_*^{\mathsf{T}}\mathbf{C}_N^{-1}\mathbf{y}, \kappa(x_*, x_*) - \mathbf{k}_*^{\mathsf{T}}\mathbf{C}_N^{-1}\mathbf{k}_* + \beta^{-1})$ results of Gaussians

• $p(y_*|y)$ is almost identical to $p(f_*|\mathbf{f})$ with **K** replaced by $\mathbf{C}_{\mathbf{N}}$ + extra β^{-1} noise variance

Learning Hyperparameters in GP based Models

- Can learn the hyperparameters of the GP prior as well as of the likelihood model
- Assuming $\mu = 0$, the hyperparams of GP are cov/kernel function hyperparams

$$\kappa(\mathbf{x}_{n}, \mathbf{x}_{m}) = \exp\left(-\frac{||\mathbf{x}_{n} - \mathbf{x}_{m}||^{2}}{\gamma}\right) \qquad (RBF \text{ kernel})$$

$$\kappa(\mathbf{x}_{n}, \mathbf{x}_{m}) = \exp\left(-\sum_{d=1}^{D} \frac{(\mathbf{x}_{nd} - \mathbf{x}_{md})^{2}}{\gamma_{d}}\right) \qquad (ARD \text{ kernel})$$

$$\kappa(\mathbf{x}_{n}, \mathbf{x}_{m}) = \kappa_{\theta_{1}}(\mathbf{x}_{n}, \mathbf{x}_{m}) + \kappa_{\theta_{2}}(\mathbf{x}_{n}, \mathbf{x}_{m}) + \ldots + \kappa_{\theta_{M}}(\mathbf{x}_{n}, \mathbf{x}_{m}) \qquad (flexible composition of multiple kernels)$$

- MLE-II is a popular choice for learning these hyperparams (otherwise MCMC, VI, etc)
- Denoting the covariance/kernel matrix as \mathbf{K}_{θ} , for Gaussian likelihood case, the marg-lik $p(\mathbf{y}|\theta, \beta^{-1}) = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{K}_{\theta} + \beta^{-1}\mathbf{I}_{N})$
- This can be maximized to learn heta and eta
- For non-Gaussian likelihoods, the marg-lik itself will need to be approximated



Weight Space View vs Function Space View

- GPs are defined w.r.t. a function space that models input-output relationship
- In contrast, we have seen models that are defined w.r.t. a weight space, e.g.,

$$p(\mathbf{y}|\mathbf{X}, \mathbf{w}) = \mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{w}, \beta^{-1}\mathbf{I}_{N}) \xrightarrow{\text{Likelihood}} \text{Marginal likelihood}$$

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\boldsymbol{\mu}_{0}, \boldsymbol{\Sigma}_{0}) \xrightarrow{\text{Prior over weight vector}} \text{Marginal likelihood}$$

$$p(\mathbf{y}|\mathbf{X}) = \int p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})d\mathbf{w} = \mathcal{N}(\mathbf{y}|\mathbf{X}\boldsymbol{\mu}_{0}, \beta^{-1}\mathbf{I}_{N} + \mathbf{X}\boldsymbol{\Sigma}_{0}\mathbf{X}^{\mathsf{T}}) \xrightarrow{\text{out the weights}} \text{out the weights}$$

$$p(\mathbf{y}|\mathbf{X}) = \mathcal{N}(\mathbf{y}|\mathbf{0}, \beta^{-1}\mathbf{I}_{N} + \mathbf{X}\mathbf{X}^{\mathsf{T}}) \xrightarrow{\text{Marginal likelihood}} \text{assuming } \boldsymbol{\mu}_{0} = \mathbf{0} \text{ and } \boldsymbol{\Sigma}_{0} = \mathbf{I}$$

$$p(\mathbf{y}|\mathbf{X}) = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{X}\mathbf{X}^{\mathsf{T}}) \xrightarrow{\text{Assuming noise-free likelihood}}$$

• Thus the joint marginal of the N responses y_1, y_2, \dots, y_N is a multivariate Gaussian

This equivalence also shows that Bayesian linear regression is a special case of GP with linear kernel $p\left(\begin{bmatrix}y_1\\y_2\\\vdots\\y_N\end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix}0\\0\\\vdots\\0\end{bmatrix}, \begin{bmatrix}x_1^{\mathsf{T}}x_1 & \dots & x_1^{\mathsf{T}}x_N\\x_2^{\mathsf{T}}x_1 & \dots & x_2^{\mathsf{T}}x_N\\\vdots & \ddots & \vdots\\x_1^{\mathsf{T}}x_1 & \dots & x_1^{\mathsf{T}}x_N\end{bmatrix}\right)$

Same as a GP $f(x_i) = y_i$, $\mu(x) = 0$ and linear covariance/kernel function $\kappa(x_i, x_j) = x_i^{\mathsf{T}} x_j$

Thus GPs can be seen as bypassing the weight space and directly defining the model using a marginal likelihood via a function space defined by the GP

Scalability of GPs

- Computational costs in some steps of GP models scale in the size of training data
- For example, prediction cost is O(N) $p(y_*|y) = \mathcal{N}(y_*|\hat{\mu}, \hat{\sigma}^2)$ $\hat{\mu} = \mathbf{k}_*^{\mathsf{T}} \mathbf{C}_N^{-1} y$ $\hat{\sigma}^2 = \kappa(x_*, x_*) - \mathbf{k}_*^{\mathsf{T}} \mathbf{C}_N^{-1} \mathbf{k}_* + \beta^{-1}$
- GP models often require matrix inversions (e.g., in marg-lik computation when estimating hyperparameters) takes $O(N^3)$
- Storage also requires $O(N^2)$ since need to store the covariance matrix
- A lot of work on speeding up GPs^1 . Some prominent approaches include $M \ll N$ pseudo-inputs and pseudo-outputs
 - Inducing Point Methods (condition predictions only on a small set of "learnable" points)
 - Divide-and-Conquer (learn GP on small subsets of data and aggregate predictions)
 - Kernel approximations
- Note that nearest neighbor methods and kernel methods also face similar issues
 - Many tricks to speed up kernel methods can be used for speeding up GPs too

Neural Networks and Gaussian Process

- An infinitely-wide single hidden layer NN with i.i.d. priors on weights = GP
- Shown formally by (Neal², 1994). Based on applying the central limit theorem



This equivalence is useful for several reasons

- Can use a GP instead of an infinitely wide Bayesian NN (which is impractical anyway)
- With GPs, inference is easy (at least for regression and with known hyperparams)
- A proof that GPs can also learn any function (just like infinitely wide neural nets Hornik's theorem)

Connection generalized to infinitely wide multiple hidden layer NN (Lee et al³, 2018)
²Priors for infinite networks, Tech Report, 1994
³Deep Neural Networks as Gaussian Processes (ICLR 2018)

GP: Some other comments

- GPs can be thought of as Bayesian analogues of kernel methods
- Can get estimate of the uncertainty in the function and its predictions
- Can learn the kernel (by learning the hyperparameters of the kernels)
- In some ways, GPs and (Bayesian/ensembles of) deep neural nets have same goals
 - These methods are also very related (though appear different based on their formulation)
 - Several recent papers have investigated these connections
- GP can be a nice alternative to (Bayesian/ensembles of) deep neural networks
 - GP may be preferable if we don't have that much training data (deep networks requires lots of data to train well)
 - When we have lots of training data, training and test speed may be an issue for GP (but faster versions exist)
- Not limited to supervised learning problems
 - f could even define a mapping of low-dim latent variable z_n to an observation x_n

$$x_n = f(\mathbf{z}_n) + "noise"$$

GP latent variable model for dimensionality reduction (like a kernel version of probabilistic PCA)



GP: A Visualization

Assumed zero mean function and a squared exponential kernel



GP packages

- Many mature implementations of GP exist. You may check out
 - GPyTorch (PyTorch), GPFlow (Tensorflow)
 - sklearn (Python with some basic GP implementations)
 - GPML (MATLAB), GPsuff (MATLAB/Octave)
 - Many others such as Stan, GPJax
- A comparison of the various packages: <u>https://en.wikipedia.org/wiki/Comparison_of_Gaussian_process_software</u>



Conditional Posterior

• Consider a model with K unknown params/hyperparams $\Theta = (\theta_1, \theta_2, ..., \theta_K)$

Joint posterior $p(\Theta|X) = \frac{p(\Theta)p(X|\Theta)}{p(X)} = \frac{p(\Theta)p(X|\Theta)}{\int p(\Theta)p(X|\Theta)d\theta_1d\theta_2 \dots d\theta_K}$ Usually intractable integral so the full posterior can't be computed exactly

• We can however compute conditional posteriors (CP) which for each θ_i looks like $p(\theta_i | \text{whatever } \theta_i \text{ depends on})$ Can be data and/or other params/hyperparams given their fixed values (or current estimates)

• To compute each CP, look at the joint distribution $p(X, \Theta)$

 $p(\boldsymbol{X},\boldsymbol{\Theta}) = p(\boldsymbol{X},\boldsymbol{\theta}_1,\boldsymbol{\theta}_2,\ldots,\boldsymbol{\theta}_K) = p(\boldsymbol{X}|\boldsymbol{\theta}_1,\boldsymbol{\theta}_2,\ldots,\boldsymbol{\theta}_K)p(\boldsymbol{\theta}_1|\boldsymbol{\theta}_2,\ldots,\boldsymbol{\theta}_K)p(\boldsymbol{\theta}_2|\boldsymbol{\theta}_3,\ldots,\boldsymbol{\theta}_K) \dots p(\boldsymbol{\theta}_K)$

- CP of θ_i will be proportional to the product of all the terms involving θ_i
 - If those terms are conjugate to each other, it is called local conjugacy. CP is then easy to compute
- Many algorithms for computing point estimate/full posterior use the CPs
 - Expectation Maximization, Variational Inference , MCMC (especially Gibbs sampling)

Latent Variable Models

Application 1: Can use latent variables to learn latent properties/features of data, e.g.,

- Cluster assignment of each observation (in mixture models)
- Low-dim rep. or "code" of each observation (e.g., prob. PCA, variational autoencoders, etc)

Plate notation of a generic LVM



• In such apps, latent variables (z_n 's) are called "local variables" (specific to individual obs.) and other unknown parameters/hyperparams (θ, ϕ above) are called "global var"

Latent Variable Models

- Application 2: Sometimes, <u>augmenting</u> a model by latent variables simplifies inference
 - These latent variables aren't part of the original model definition
- Some of the popular examples of such augmentation include
 - In Probit regression for binary classification, we can model each label $y_n \in \{0,1\}$ as

$$y_n = \mathbb{I}[z_n > 0]$$
 where $z_n \sim \mathcal{N}(w^T x_n, 1)$ is an auxiliary latent variable

.. and use EM etc, to infer the unknowns \boldsymbol{w} and $\boldsymbol{z_n}$'s (PML-2, Sec 15.4)

Many sparse priors on weights can be thought of as Gaussian "scale-mixtures"

$$\mathsf{Laplace}(w_d|0, 1/\gamma) = \frac{\gamma}{2} \exp(-\gamma|w_d|) = \int \mathcal{N}(w_d|0, \tau_d^2) \mathsf{Gamma}(\tau_d^2|1, \gamma^2/2) d\tau_d^2$$

.. where au_d 's are latent vars. Can use EM to infer w, au (MLAPP 13.4.4 - EM for LASSO)

Such augmentations can often make a non-conjugate model a locally conjugate one

Conditional posteriors of the unknowns often have closed form in such cases

Nomenclature/Notation Alert

- Why call some unknowns as parameters and others as latent variables?
- Well, no specific reason. Sort of a convention adopted by some algorithms
 - EM: Unknowns estimated in E step referred to as latent vars; those in M step as params
 - Usually: Latent vars (Conditional) posterior computed; parameters point estimation
- Some algos won't make such distinction and will infer posterior over all unknowns
- Sometimes the "global" or "local" unknown distinction makes it clear
 - Local variables = latent variables, global variables = parameters
- But remember that this nomenclature isn't really cast in stone, no need to be confused so long as you are clear as to what the role of each unknown is, and how we want to estimate it (posterior or point estimate) and using what type of inference algorithm

12

Hybrid Inference (posterior infer. + point est.)

In many models, we infer posterior on some unknowns and do point est. for others

CP of $w: p(w|X, y, \hat{\lambda}, \hat{\beta})$

 $\{\hat{\lambda}, \hat{\beta}\} = \operatorname{argmax}_{\lambda,\beta} p(\boldsymbol{y}|\boldsymbol{X}, \lambda, \beta)$

- We have already seen MLE-II for lin reg. which alternates between
 - Inferring CP over the main parameter given the point estimates of hyperparams
 - Maximizing the marginal lik. to do point estimation for hyperparams
- The Expectation-Maximization algorithm (will see today) also does something similar
 - In E step, the CP of latent variables is inferred, given <u>current</u> point-est of params
 - M step maximizes expected complete data log-lik. to get point estimates of params
- If we can't (due to computational or other reasons) infer posterior over all unknowns, how to decide which variables to infer posterior on, and for which to do point-est?
- Usual approach: Infer posterior over local vars and point estimates for global vars
 - Reason: We typically have plenty of data to reliably estimate the global variables so it is okay even if we just do point estimation for those

Inference/Parameter Estimation in Latent Variable Models using Expectation-Maximization (EM)



15

Parameter Estimation in Latent Variable Models

• Assume each observation x_n to be associated with a "local" latent variable z_n



- Although we can do fully Bayesian inference for all the unknowns, suppose we only want a point estimate of the "global" parameters $\Theta = (\theta, \phi)$ via MLE/MAP
- Such MLE/MAP problems in LVMs are difficult to solve in a "clean" way
 - Would typically require gradient based methods with no closed form updates for Θ
 - However, EM gives a clean way to obtain closed form updates for Θ



16

Why MLE/MAP of Params is Hard for LVMs?

- Suppose we want to estimate $\Theta = (\theta, \phi)$ via MLE. If we knew \mathbf{z}_n , we could solve
- Easy to solve $\Theta_{MLE} = \arg \max_{\Theta} \sum_{n=1}^{N} \log p(\mathbf{x}_n, \mathbf{z}_n | \Theta) = \arg \max_{\Theta} \sum_{n=1}^{N} [\log p(\mathbf{z}_n | \phi) + \log p(\mathbf{x}_n | \mathbf{z}_n, \theta)]$ = Easy. Usually closed form if $p(\mathbf{z}_n | \phi)$ and $p(\mathbf{x}_n | \mathbf{z}_n, \theta)$ have simple forms exp-f In particular, if they are exp-fam distributions
- However, since in LVMs, \mathbf{z}_n is hidden, the MLE problem for Θ will be the following Basically, the marginal (elihood after tegrating out z_n $\Theta_{MLE} = \arg \max_{\Theta} \sum_{n=1}^n \log p(\mathbf{x}_n | \Theta) = \arg \max_{\Theta} \log p(\mathbf{X} | \Theta)$ $\bullet \log p(\mathbf{x}_n | \Theta)$ will not have a simple expression since $p(\mathbf{x}_n | \Theta)$ requires sum/integral likelihood after integrating out z_n

$$p(\mathbf{x}_n|\Theta) = \sum_{\mathbf{z}_n} p(\mathbf{x}_n, \mathbf{z}_n|\Theta)$$
 ... or if \mathbf{z}_n is continuous: $p(\mathbf{x}_n|\Theta) = \int p(\mathbf{x}_n, \mathbf{z}_n|\Theta) d\mathbf{z}_n$

CS772A: PML

• MLE now becomes difficult (basically MLE-II now), no closed form expression for Θ .

• Can we maximize some other quantity instead of $\log p(x_n | \Theta)$ for this MLE?

An Important Identity

• Assume $p_z = p(Z|X, \Theta)$ and q(Z) to be some prob distribution over Z, then

• In the above
$$\mathcal{L}(q, \Theta) = \sum_{Z} q(Z) \log \left\{ \frac{p(X, Z | \Theta)}{q(Z)} \right\}$$

•
$$KL(q||p_z) = -\sum_Z q(\mathbf{Z}) \log\left\{\frac{p(\mathbf{Z}|\mathbf{X},\Theta)}{q(\mathbf{Z})}\right\}$$

Assume **Z** discrete

- KL is always non-negative, so $\log p(X|\Theta) \ge \mathcal{L}(q,\Theta)$
- Thus $\mathcal{L}(q, \Theta)$ is a lower-bound on $\log p(X|\Theta)$
- Thus if we maximize $\mathcal{L}(q, \Theta)$, it will also improve $\log p(X|\Theta)$
- Also, as we'll see, it's easier to maximize $\mathcal{L}(q, \Theta)$



Verify the identity



18



The Expectation-Maximization (EM) Algorithm

• ALT-OPT of $\mathcal{L}(q, \Theta)$ w.r.t. q and Θ gives the EM algorithm (Dempster, Laird, Rubin, 1977)

20



The Expected CLL

Expected CLL in EM is given by (assume observations are i.i.d.)

$$\begin{aligned} \mathcal{Q}(\Theta, \Theta^{old}) &= \sum_{n=1}^{N} \mathbb{E}_{p(\boldsymbol{z}_n | \boldsymbol{x}_n, \Theta^{old})} [\log p(\boldsymbol{x}_n, \boldsymbol{z}_n | \Theta)] \\ &= \sum_{n=1}^{N} \mathbb{E}_{p(\boldsymbol{z}_n | \boldsymbol{x}_n, \Theta^{old})} [\log p(\boldsymbol{x}_n | \boldsymbol{z}_n, \Theta) + \log p(\boldsymbol{z}_n | \Theta)] \end{aligned}$$

- If $p(\mathbf{z}_n|\Theta)$ and $p(\mathbf{x}_n|\mathbf{z}_n,\Theta)$ are exp-family distributions, $Q(\Theta,\Theta^{\text{old}})$ has a very simple form
- In resulting expressions, replace terms containing z_n 's by their respective expectations, e.g.,
 - \boldsymbol{z}_n replaced by $\mathbb{E}_{p(\boldsymbol{z}_n | \boldsymbol{x}_n, \widehat{\Theta})}[\boldsymbol{z}_n]$
 - $\mathbf{z}_n \mathbf{z}_n^{\mathsf{T}}$ replaced by $\mathbb{E}_{p(\mathbf{z}_n | \mathbf{x}_n, \widehat{\Theta})}[\mathbf{z}_n \mathbf{z}_n^{\mathsf{T}}]$
- However, in some LVMs, these expectations are intractable to compute and need to be approximated (will see some examples later)



What's Going On?

Alternating between them until convergence to some local optima

KL becomes zero and $\mathcal{L}(q, \Theta)$ becomes equal to $\log p(X|\Theta)$; thus their curves touch at current Θ

- As we saw, the maximization of lower bound $\mathcal{L}(q,\Theta)$ had two steps
- Step 1 finds the optimal q (call it \hat{q}) by setting it as the posterior of Z given current Θ
- Step 2 maximizes $\mathcal{L}(\hat{q}, \Theta)$ w.r.t. Θ which gives a new Θ .



EM vs Gradient-based Methods

- Can also estimate params using gradient-based optimization instead of EM
 - We can usually explicitly sum over or integrate out the latent variables Z, e.g.,

$$\mathcal{L}(\Theta) = \log p(\mathbf{X}|\Theta) = \log \sum_{n} p(\mathbf{X}, \mathbf{Z}|\Theta)$$

- Now we can optimize $\mathcal{L}(\Theta)$ using first/second order optimization to find the optimal Θ
- EM is usually preferred over this approach because
 - $\hfill\blacksquare$ The M step has often simple closed-form updates for the parameters Θ
 - Often constraints (e.g., PSD matrices) are automatically satisfied due to form of updates
 - In some cases[†], EM usually converges faster (and often like second-order methods)
 - E.g., Example: Mixture of Gaussians with when the data is reasonably well-clustered
 - EM applies even when the explicit summing over/integrating out is expensive/intractable
 - EM also provides the conditional posterior over the latent variables Z (from E step)

Some Applications of EM

- Mixture Models and Dimensionality Reduction/Representation Learning
 - Mixture Models: Mixture of Gaussians, Mixture of Experts, etc
 - Dim. Reduction/Representation Learning: Probabilistic PCA, Variational Autoencoders
- Problems with missing features or missing labels (which are treated as latent variables)
 - $\widehat{\Theta} = \operatorname{argmax}_{\Theta} \log p(\mathbf{x}^{obs} | \Theta) = \operatorname{argmax}_{\Theta} \log \int p([\mathbf{x}^{obs}, \mathbf{x}^{miss}] | \Theta) d\mathbf{x}^{miss}$
 - $\widehat{\Theta} = \operatorname{argmax}_{\Theta} \sum_{n=1}^{N} \log p(x_n, y_n | \Theta) + \sum_{n=N+1}^{N+M} \log \sum_{c=1}^{K} p(x_n, y_n = c | \Theta)$

Hyperparameter estimation in probabilistic models (an alternative to MLE-II)

MLE-II estimates hyperparams by maximizing the marginal likelihood, e.g.,

 $\{\hat{\lambda}, \hat{\beta}\} = \operatorname{argmax}_{\lambda,\beta} p(\boldsymbol{y}|\boldsymbol{X}, \lambda, \beta) = \operatorname{argmax}_{\lambda,\beta} \int p(\boldsymbol{y}|\boldsymbol{w}, \boldsymbol{X}, \beta) p(\boldsymbol{w}|\lambda) d\boldsymbol{w}$

- With EM, can treat w as latent var, and λ , β as "parameters"
 - E step will estimate the CP of w given current estimates of λ, β
 - M step will re-estimate λ, β by maximizing the expected CLL

 $\mathbb{E}[\log p(\mathbf{y}, \mathbf{w} | \mathbf{X}, \beta, \lambda)] = \mathbb{E}[\log p(\mathbf{y} | \mathbf{w}, \mathbf{X}, \beta) + \log p(\mathbf{w} | \lambda)]$



For a Bayesian linear

regression model

Expectations w.r.t.

the CP of **w**



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 μ_k and cov matrix Σ_k
- The labels \mathbf{z}_n (known) are one-hot vecs. Also, $z_{nk}=1$ if $\mathbf{z}_n=k$, and $\mathbf{z}_{nk}=0$, o/w
- Assuming class prior as $p(\mathbf{z}_n = k) = \pi_k$, the model has params $\Theta = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$
- Given training data $\{x_n, z_n\}_{n=1}^N$, the MLE solution will be

$$\hat{\pi}_{k} = \frac{1}{N} \sum_{n=1}^{N} z_{nk} \qquad \text{Same as } \frac{N_{k}}{N} \text{ where } N_{k} \text{ is } \# \text{ of training ex. for which } y_{n} = k$$

$$\hat{\mu}_{k} = \frac{1}{N_{k}} \sum_{n=1}^{N} z_{nk} \boldsymbol{x}_{n} \qquad \text{Same as } \frac{1}{N_{k}} \sum_{n:\boldsymbol{z}_{n}=k}^{N} \boldsymbol{x}_{n}$$

$$\hat{\Sigma}_{k} = \frac{1}{N_{k}} \sum_{n=1}^{N} z_{nk} (\boldsymbol{x}_{n} - \hat{\mu}_{k}) (\boldsymbol{x}_{n} - \hat{\mu}_{k})^{\mathsf{T}} \qquad \text{Same as } \frac{1}{N_{k}} \sum_{n:\boldsymbol{z}_{n}=k}^{N} (\boldsymbol{x}_{n} - \hat{\mu}_{k}) (\boldsymbol{x}_{n} - \hat{\mu}_{k})^{\mathsf{T}}$$

Detour: MLE for Generative Classification

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

$$\widehat{\Theta} = \operatorname{argmax}_{\Theta} p(X, Z|\Theta) = \operatorname{argmax}_{\Theta} \prod_{n=1}^{N} p(x_n, z_n|\Theta)_{\text{multinoulli}} \qquad \text{Gaussian}$$

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

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

$$= \operatorname{argmax}_{\Theta} \prod_{n=1}^{N} \prod_{k=1}^{K} \pi_k^{Z_{nk}} \prod_{k=1}^{K} p(x_n|z_n = k, \Theta)^{Z_{nk}}$$

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

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

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

$$= \operatorname{argmax}_{\Theta} \log \prod_{n=1}^{N} \sum_{k=1}^{K} [\pi_k p(x_n|z_n = k, \Theta)]^{Z_{nk}}$$

$$= \operatorname{argmax}_{\Theta} \sum_{n=1}^{N} \sum_{k=1}^{K} [\pi_k p(x_n|z_n = k, \Theta)]^{Z_{nk}}$$
The form of this expression is important; will encounter this in GMM too $\sum_{n=1}^{N} \sum_{k=1}^{K} [\pi_k p(x_n|z_n = k, \Theta)]^{Z_{nk}}$

EM for Mixture Models

• So how do we estimate the parameters of a GMM where z_n 's are <u>unknown</u>?



- The guess about \boldsymbol{z}_n can be in one of the two forms
 - A "hard" guess a single best value $\hat{\boldsymbol{z}}_n$ (some "optimal" value of the random variable \boldsymbol{z}_n)
 - The "expected" value $\mathbb{E}[\boldsymbol{z}_n]$ of the random variable \boldsymbol{z}_n
- Using the hard guess \hat{z}_n of z_n will result in an ALT-OPT like algorithm of the latent variables
- Using the expected value of z_n will give the so-called Expectation-Maximization (EM) alo

EM is pretty much like ALT-OPT

EM for Gaussian Mixture Model (GMM)

- EM finds Θ_{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., $p(Z|X, \Theta)$
- 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
- 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)

Expectation of CLL

It is "conditional" posterior

on Θ , not just data X

because it is also conditioned

Same as $p(z_{nk} = 1 | x_{nk})$ different notation

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

- Update Θ by maximizing the <u>expected</u> 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

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

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

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

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

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

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

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

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

Requires knowing Θ

EM for GMM (Contd)

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

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

EM for Gaussian Mixture Model

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

points in each cluster

E step: compute the expectation of each z_n (we need it in M step) 2 Accounts for fraction of Accounts for cluster shapes (since

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

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

 $\mathbb{E}[\boldsymbol{z}_{nk}^{(t)}] = \gamma_{nk}^{(t)} = \frac{\pi_k^{(t-1)}\mathcal{N}(\boldsymbol{x}_n | \boldsymbol{\mu}_k^{(t-1)}, \boldsymbol{\Sigma}_k^{(t-1)})}{\sum_{k=1}^{K} \pi_k^{(t-1)}\mathcal{N}(\boldsymbol{x}_n | \boldsymbol{\mu}_k^{(t-1)}, \boldsymbol{\Sigma}_k^{(t-1)})}$

$$\begin{pmatrix} t \\ k \end{pmatrix} = \frac{1}{N_k} \sum_{n=1}^{N} \gamma_{nk}^{(t)} \mathbf{x}_n$$
 Effective number of points in the k^{th} cluster

M-step: $\boldsymbol{\Sigma}_{k}^{(t)} = \frac{1}{N_{k}} \sum_{n=1}^{N} \gamma_{nk}^{(t)} (\boldsymbol{x}_{n} - \boldsymbol{\mu}_{k}^{(t)}) (\boldsymbol{x}_{n} - \boldsymbol{\mu}_{k}^{(t)})^{\top}$

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

 $\boldsymbol{\mu}$





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

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

each cluster is a Gaussian

 $\forall n, k$

EM: Some Final Comments

- The E and M steps may not always be possible to perform exactly. Some reasons
 - The conditional posterior of latent variables p(Z|X, O) may not be easy to compute
 Will need to approximate p(Z|X, O) using methods such as MCMC or variational inference Results in
 - Even if $p(Z|X,\Theta)$ is easy, the expected CLL may not be easy to compute $\mathbb{E}[\log p(X, Z|\Theta)] = \int \log p(X, Z|\Theta) p(Z|X, \Theta) dZ$ Can often be approximated by Monte-Carlo using sample from the CP of Z
 - Maximization of the expected CLL may not be possible in closed form
- EM works even if the M step is only solved approximately (Generalized EM)
- If M step has multiple parameters whose updates depend on each other, they are updated in an alternating fashion - called Expectation Conditional Maximization (ECM)
- Other advanced probabilistic inference algos are based on ideas similar to EM
 - E.g., Variational EM, Variational Bayes (VB) inference, a.k.a. Variational Inference (VI)

Monte-Carlo EM