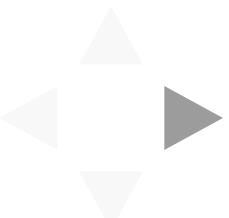


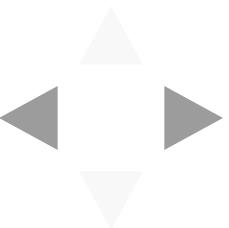
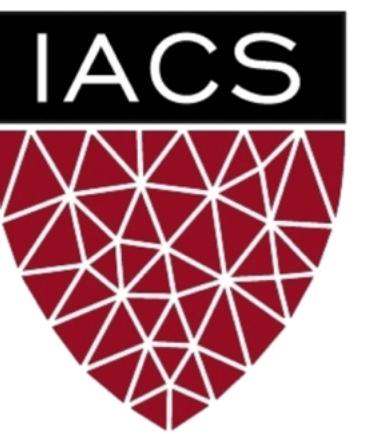
Lecture #9: Latent Variable Models and MLE

AM 207: Advanced Scientific Computing

Stochastic Methods for Data Analysis, Inference and Optimization

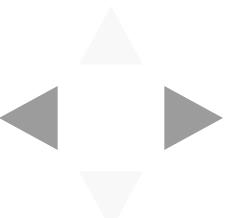
Fall, 2020



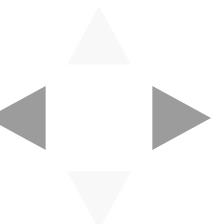


Outline

1. Motivation for Latent Variable Models
2. Common Latent Variable Models
3. Maximum Likelihood Estimation for Latent Variable Models: Expectation Maximization
4. Mixture of Gaussians



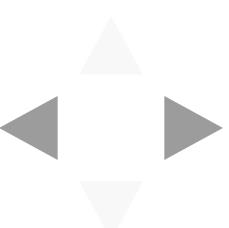
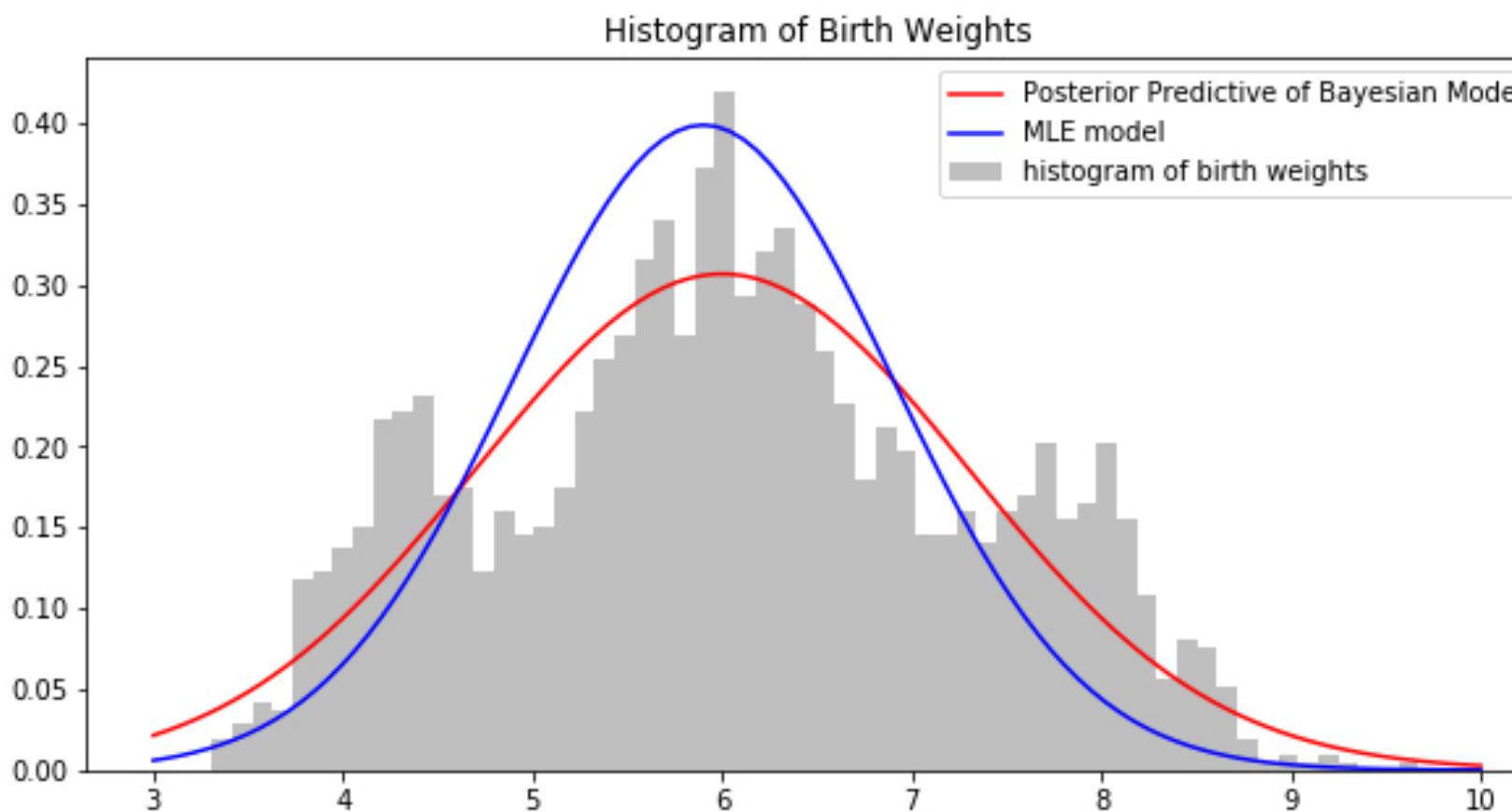
Motivation for Latent Variable Models



A Model for Birth Weights

Recall our model for birth weights, Y_1, \dots, Y_N . We posited that the birth weights are iid normally distributed with known σ^2 , $Y_n \sim \mathcal{N}(\mu, 1)$.

Compare the maximum likelihood model and the Bayesian model for birth weight. Which model would you use to make clinical decisions? What's hard about this comparison?



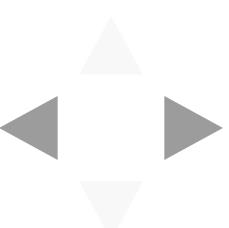
A Similarity Measure for Distributions: Kullback–Leibler Divergence

Visually comparing models to the *empirical distribution* of the data is impractical. Fortunately, there are a large number of quantitative measures for comparing two distributions, these are called *divergence measures*. For example, the *Kullback–Leibler (KL) Divergence* is defined for two distributions $p(\theta)$ and $q(\theta)$ supported on Θ as:

$$D_{\text{KL}}[q \parallel p] = \int_{\Theta} \log \left[\frac{q(\theta)}{p(\theta)} \right] q(\theta) d\theta$$

The KL-divergence $D_{\text{KL}}[q \parallel p]$ is bounded below by 0, which happens if and only if $q = p$. The KL-divergence has information theoretic interpretations that we will explore later in the course.

Note: The KL-divergence is defined in terms of the pdf's of p and q . If p is a distribution from which we only have samples and not the pdf (like the empirical distribution), we can nonetheless estimate $D_{\text{KL}}[q \parallel p]$. Techniques that estimate the KL-divergence from samples are called *non-parametric*. We will use them later in the course.



Class Membership as a Latent Variable

We observe that there are three *clusters* in the data. We posit that there are three *classes* of infants in the study: infants with low birth weights, infants with normal birth weights and those with high birth weights. The numbers of infants in the classes are not equal.

For each observation Y_n , we model its class membership Z_n as a categorical variable,

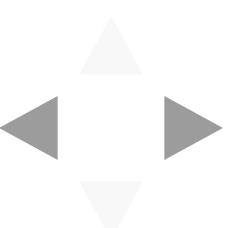
$$Z_n \sim Cat(\pi),$$

where π_i in $\pi = [\pi_1, \pi_2, \pi_3]$ is the class proportion. Note that we don't have the class membership Z_n in the data! So Z_n is called a *latent variable*.

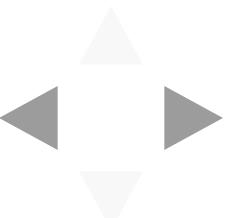
Depending on the class, the n -th birth weight Y_n will have a different normal distribution,

$$Y_n | Z_n \sim \mathcal{N} (\mu_{Z_n}, \sigma_{Z_n}^2)$$

where μ_{Z_n} is one of the three class means $[\mu_1, \mu_2, \mu_3]$ and $\sigma_{Z_n}^2$ is one of the three class variances $[\sigma_1^2, \sigma_2^2, \sigma_3^2]$.

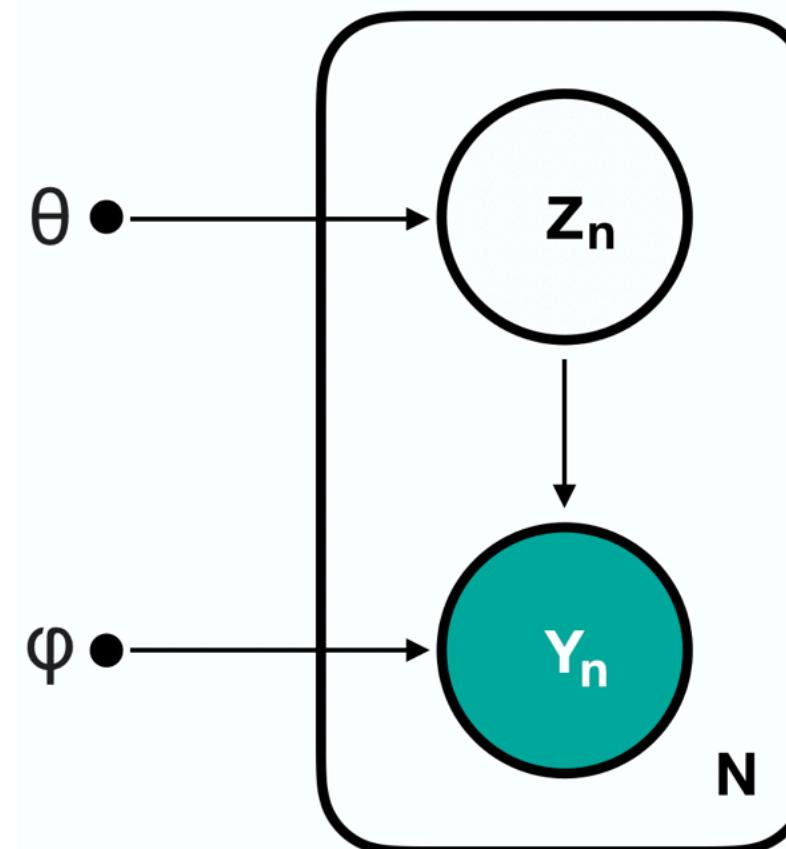


Common Latent Variable Models

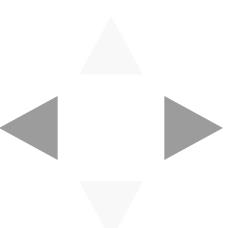


Latent Variable Models

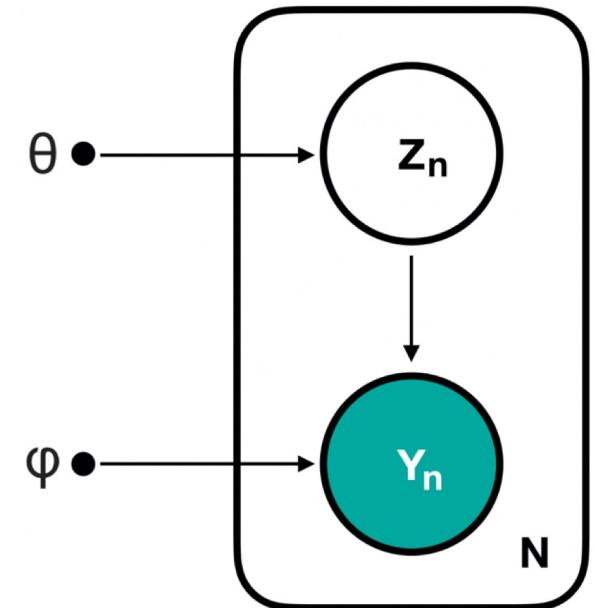
Models that include an observed variable Y and at least one unobserved variable Z are called *latent variable models*. In general, our model can allow Y and Z to interact in many different ways. Today, we will study models with one type of interaction:



$$\begin{aligned} Z_n &\sim p(Z|\theta) \\ Y_n|Z_n &\sim p(Y|Z, \phi) \\ n &= 1, \dots, N \end{aligned}$$



Item-Response Models



In *item-response models*, we measure an real-valued unobserved trait Z of a subject by performing a series of experiments with binary observable outcomes, Y :

$$Z_n \sim \mathcal{N}(\mu, \sigma^2),$$

$$\theta_n = g(Z_n)$$

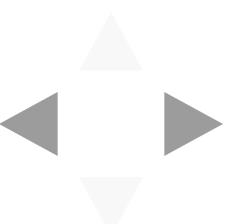
$$Y_n | Z_n \sim Ber(\theta_n),$$

where $n = 1, \dots, N$ and g is some fixed function of Z_n .

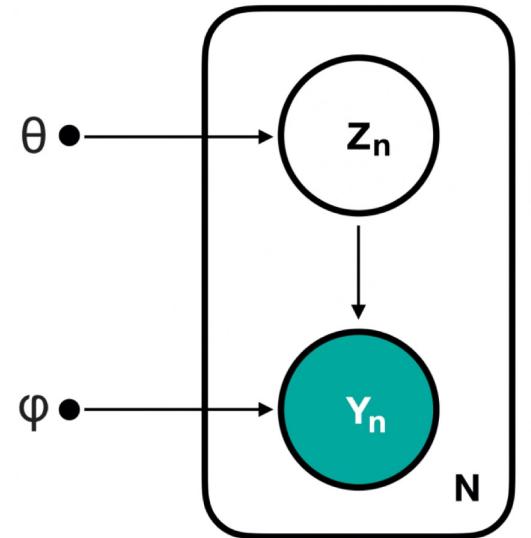
Applications

Item response models are used to model the way "underlying intelligence" Z relates to scores Y on IQ tests.

Item response models can also be used to model the way "suicidality" Z relates to answers on mental health surveys. Building a good model may help to infer when a patient is at psychiatric risk based on in-take surveys at points of care through out the health-care system.



Factor Analysis Models



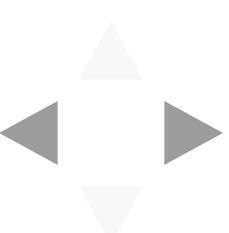
In **factor analysis models**, we posit that the observed data Y with many measurements is generated by a small set of unobserved factors Z :

$$Z_n \sim \mathcal{N}(0, I),$$
$$Y_n | Z_n \sim \mathcal{N}(\mu + \Lambda Z_n, \Phi),$$

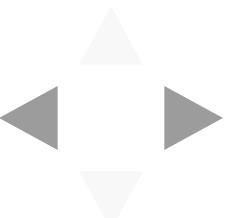
where $n = 1, \dots, N$, $Z_n \in \mathbb{R}^{D'}$ and $Y_n \in \mathbb{R}^D$. We typically assume that D' is much smaller than D .

Applications

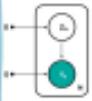
Factor analysis models are useful for biomedical data, where we typically measure a large number of characteristics of a patient (e.g. blood pressure, heart rate, etc), but these characteristics are all generated by a small list of health factors (e.g. diabetes, cancer, hypertension etc). Building a good model means we may be able to infer the list of health factors of a patient from their observed measurements.



Maximum Likelihood Estimation for Latent Variable Models: Expectation Maximization



MODEL



$$\begin{aligned} z_n &\sim p(z|n) \\ y_n &\sim p(y_n|z_n, \theta) \end{aligned}$$

parameters: θ, β Calculus plots we need:

1. $E_{x \sim p(x)}[f(x)] = \int f(x)p(x)dx$

2. properties of E :

A. $E_{x \sim p(x)}[af(x)] = aE_{x \sim p(x)}[f(x)]$

(Jensen's Inequality) \rightarrow B. $\log E_{x \sim p(x)}[f(x)] \geq E_{x \sim p(x)}[\log f(x)]$

C. $E_{x \sim p(x)}[f(x)] + E_{x \sim p(x)}[g(x)] = E_{x \sim p(x)}[f(x) + g(x)]$

3. $D_KL[q||p](x) = \int \log \frac{q(x)}{p(x)} q(x) dx$
 $= E_{x \sim q(x)}[\log \frac{q(x)}{p(x)}]$

Likelihood over $y|z$:

$$\log \prod_{n=1}^N p(y_n|z_n, \theta) = \sum_{n=1}^N \log p(y_n|z_n, \theta) + \log p(z_n|\theta)$$

can we evaluate this?

Likelihood over observed data:

$$\begin{aligned} \log \prod_{n=1}^N p(y_n|z_n, \theta) &= \sum_{n=1}^N \log p(y_n|z_n, \theta) \\ &= \sum_{n=1}^N \log \int p(y_n|z_n, \theta) p(z_n|z) dz_n \\ &= \sum_{n=1}^N \log E_{z \sim p(z|z_n)}[p(y_n|z_n, \theta)] \end{aligned}$$

$\downarrow (\theta, \beta)$

The maximum likelihood objective:

$$\hat{\theta}_{MLE}, \hat{\beta}_{MLE} = \arg \max_{\theta, \beta} \ell(\theta, \beta) = \arg \max_{\theta, \beta} \sum_{n=1}^N \log E_{z \sim p(z|z_n)}[p(y_n|z_n, \theta)]$$

is this hard?

Trying out the optimization:

$$\begin{aligned} \nabla_{\theta, \beta} \ell(\theta, \beta) &= \nabla_{\theta, \beta} \sum_{n=1}^N \log E_{z \sim p(z|z_n)}[p(y_n|z_n, \theta)] \\ &= \sum_{n=1}^N \nabla_{\theta, \beta} \log E_{z \sim p(z|z_n)}[p(y_n|z_n, \theta)] \\ &= \sum_{n=1}^N \nabla_{\theta, \beta} \frac{E_{z \sim p(z|z_n)}[p(y_n|z_n, \theta)]}{E_{z \sim p(z|z_n)}[1]} \end{aligned}$$

\downarrow

can we MC estimate this?

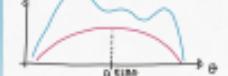


An objective we can actually work with:

$$\begin{aligned} \max_{\theta, \beta} \ell(\theta, \beta) &= \max_{\theta, \beta} \sum_{n=1}^N \log E_{z \sim p(z|z_n)}[p(y_n|z_n, \theta)] \\ &= \max_{\theta, \beta} \sum_{n=1}^N \log \int p(y_n|z_n, \theta) p(z_n|z) dz_n \\ &= \max_{\theta, \beta} \sum_{n=1}^N \log \int \frac{p(y_n|z_n, \theta) p(z_n|z)}{q(z_n)} \frac{q(z_n)}{q(z_n)} dz_n \\ &\quad \text{Introduce auxiliary variable } q(z_n)! \\ &q \text{ is same dist'n of your choice.} \\ &= \max_{\theta, \beta} \sum_{n=1}^N \log \int_{q(z_n)} \frac{p(y_n|z_n, \theta) p(z_n|z)}{q(z_n)} dz_n \\ &\quad \text{is the log helping? does } \nabla_{\theta, \beta} \text{ commute with } \int_{q(z_n)}? \\ &\geq \max_{\theta, \beta, q} \sum_{n=1}^N \mathbb{E}_{z \sim q(z_n)} \left[\log \left(\frac{p(y_n|z_n, \theta) p(z_n|z)}{q(z_n)} \right) \right] \end{aligned}$$

The Evidence Lower Bound
ELBO(θ, β, q)

Idea: Instead of maximizing the log-likelihood, we maximize the lower bound ELBO.

Why you should
When ELBO is maximized $\ell(\theta, \beta)$ is guaranteed to be as big.Why you shouldn't
If ELBO is maximized then $\ell(\theta, \beta)$ can be far from optimized

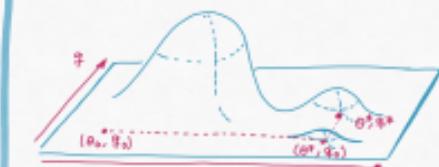
How to maximize ELBO: coordinate ascent

I. Maximize θ, β , fixing q^*

$$\begin{aligned} \theta^*, \beta^* &= \arg \max_{\theta, \beta} \text{ELBO}(\theta, \beta, \beta^*) \\ &= \arg \max_{\theta, \beta} \sum_{n=1}^N \mathbb{E}_{z \sim q(z_n)} \left[\log \left(\frac{p(y_n|z_n, \theta) p(z_n|z)}{q(z_n)} \right) \right] \\ &= \arg \max_{\theta, \beta} \sum_{n=1}^N \left(\mathbb{E}_{z \sim q(z_n)} \left[\log [p(y_n|z_n, \theta) p(z_n|z)] \right] \right. \\ &\quad \left. - \mathbb{E}_{z \sim q(z_n)} \left[\log q(z_n) \right] \right) \\ &\quad \text{irrelevant for } \max_{\theta, \beta} \\ &\equiv \arg \max_{\theta, \beta} \sum_{n=1}^N \mathbb{E}_{z \sim q(z_n)} \left[\log [p(y_n|z_n, \theta) p(z_n|z)] \right] \end{aligned}$$

is this problem easier?

Are we done?



Maximizing each coordinate once is not sufficient! we need to iterate!

II. Maximize q , fixing θ^*, β^*

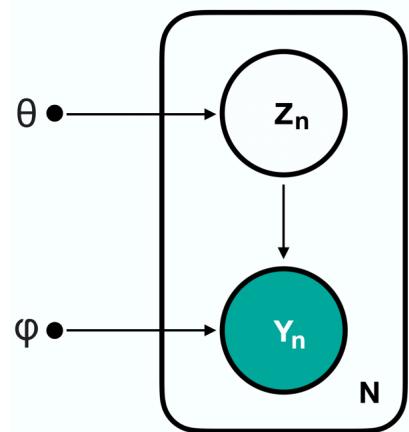
$$\begin{aligned} q^* &= \arg \max_q \text{ELBO}(\theta^*, \beta^*, q) \\ \text{Note: } \ell(\theta^*, \beta^*) - \text{ELBO}(\theta^*, \beta^*, q) &= \sum_{n=1}^N \left[\log p(y_n|z_n, \theta^*) - \mathbb{E}_{z \sim q(z_n)} \left[\log \frac{p(y_n|z_n, \theta^*)}{q(z_n)} \right] \right] \\ &= \sum_{n=1}^N \left[\mathbb{E}_{z \sim q(z_n)} \left[\log p(y_n|z_n, \theta^*) \right] - \mathbb{E}_{z \sim q(z_n)} \left[\log \frac{p(y_n|z_n, \theta^*)}{q(z_n)} \right] \right] \\ &= \sum_{n=1}^N \left[\mathbb{E}_{z \sim q(z_n)} \left[\log p(y_n|z_n, \theta^*) - \log \frac{p(y_n|z_n, \theta^*)}{q(z_n)} \right] \right] \\ &= \sum_{n=1}^N \left[\mathbb{E}_{z \sim q(z_n)} \log \left(\frac{p(y_n|z_n, \theta^*)}{p(y_n|z_n, \theta^*, \beta^*)} \right) \right] \\ &= \sum_{n=1}^N \left[\mathbb{E}_{z \sim q(z_n)} \log \left(\frac{q(z_n)}{p(z_n|y_n, \theta^*, \beta^*)} \right) \right] \\ &= \sum_{n=1}^N D_{KL}[q(z_n) || p(z_n|y_n, \theta^*, \beta^*)] \end{aligned}$$

$$q^* = \arg \max_q \text{ELBO}(\theta^*, \beta^*, q) = \arg \min_q D_{KL}[q(z_n) || p(z_n|y_n, \theta^*, \beta^*)]$$

$$q^* = p(z_n|y_n, \theta^*, \beta^*)$$

The Expectation Maximization Algorithm

The *expectation maximization (EM) algorithm* maximize the ELBO of the model,



$$\begin{aligned} Z_n &\sim p(Z|\theta) \\ Y_n|Z_n &\sim p(Y|Z, \phi) \\ n &= 1, \dots, N \end{aligned}$$

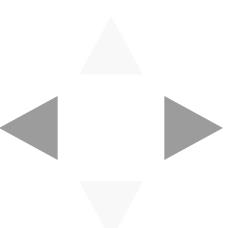
1. **Initialization:** Pick θ_0, ϕ_0 .
2. Repeat $i = 1, \dots, I$ times:

E-Step:

$$q_{\text{new}}(Z_n) = \underset{q}{\operatorname{argmax}} \text{ } ELBO(\theta_{\text{old}}, \phi_{\text{old}}, q) = p(Z_n|Y_n, \theta_{\text{old}}, \phi_{\text{old}})$$

M-Step:

$$\begin{aligned} \theta_{\text{new}}, \phi_{\text{new}} &= \underset{\theta, \phi}{\operatorname{argmax}} \text{ } ELBO(\theta, \phi, q_{\text{new}}) \\ &= \underset{\theta, \phi}{\operatorname{argmax}} \sum_{n=1}^N \mathbb{E}_{Z_n \sim p(Z_n|Y_n, \theta_{\text{old}}, \phi_{\text{old}})} [\log(p(y_n, Z_n|\phi, \theta)] . \end{aligned}$$



Example: EM for the Gaussian Mixture Model of Birth Weight

Solving the M-Step

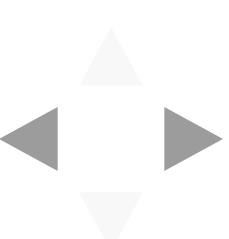
We see that the optimization problem in the M-step:

$\mu_{\text{new}}, \sigma^2_{\text{new}}, \pi_{\text{new}} = \underset{\mu, \sigma^2, \pi}{\operatorname{argmax}} \text{ELBO}(\mu, \sigma^2, \pi, q_{\text{new}})$ is equivalent to two problems

1. $\underset{\mu, \sigma^2}{\operatorname{argmax}} \sum_{n=1}^N \sum_{k=1}^K r_{n,k} \log \mathcal{N}(y_n; \mu_k, \sigma_k^2)$
2. $\underset{\pi}{\operatorname{argmax}} \sum_{n=1}^N \sum_{k=1}^K r_{n,k} \pi_k$

We can solve each optimization problem analytically by finding stationary points of the gradient (or the Lagrangian):

- $\mu_{\text{new}} = \frac{1}{\sum_{n=1}^N r_{n,k}} \sum_{n=1}^N r_{n,k} y_n$
- $\sigma^2_{\text{new}} = \frac{1}{\sum_{n=1}^N r_{n,k}} \sum_{n=1}^N r_{n,k} (y_n - \mu_{\text{new}})^2$
- $\pi_{\text{new}} = \frac{\sum_{n=1}^N r_{n,k}}{N}$



Example: EM for the Gaussian Mixture Model of Birth Weight

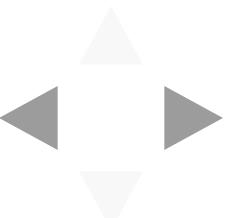
All Together

Initialization: Pick any π, μ, σ^2

E-Step: Compute $r_{n,k} = \frac{\pi_{k,\text{old}} \mathcal{N}(y_n; \mu_{k,\text{old}}, \sigma_{k,\text{old}}^2)}{\mathcal{Z}}$, where
 $\mathcal{Z} = \sum_{k=1}^K \pi_{k,\text{old}} \mathcal{N}(y_n; \mu_{k,\text{old}}, \sigma_{k,\text{old}}^2)$.

M-Step: Compute model parameters:

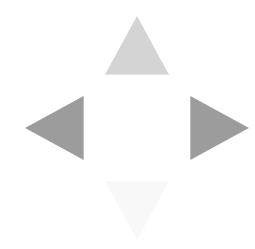
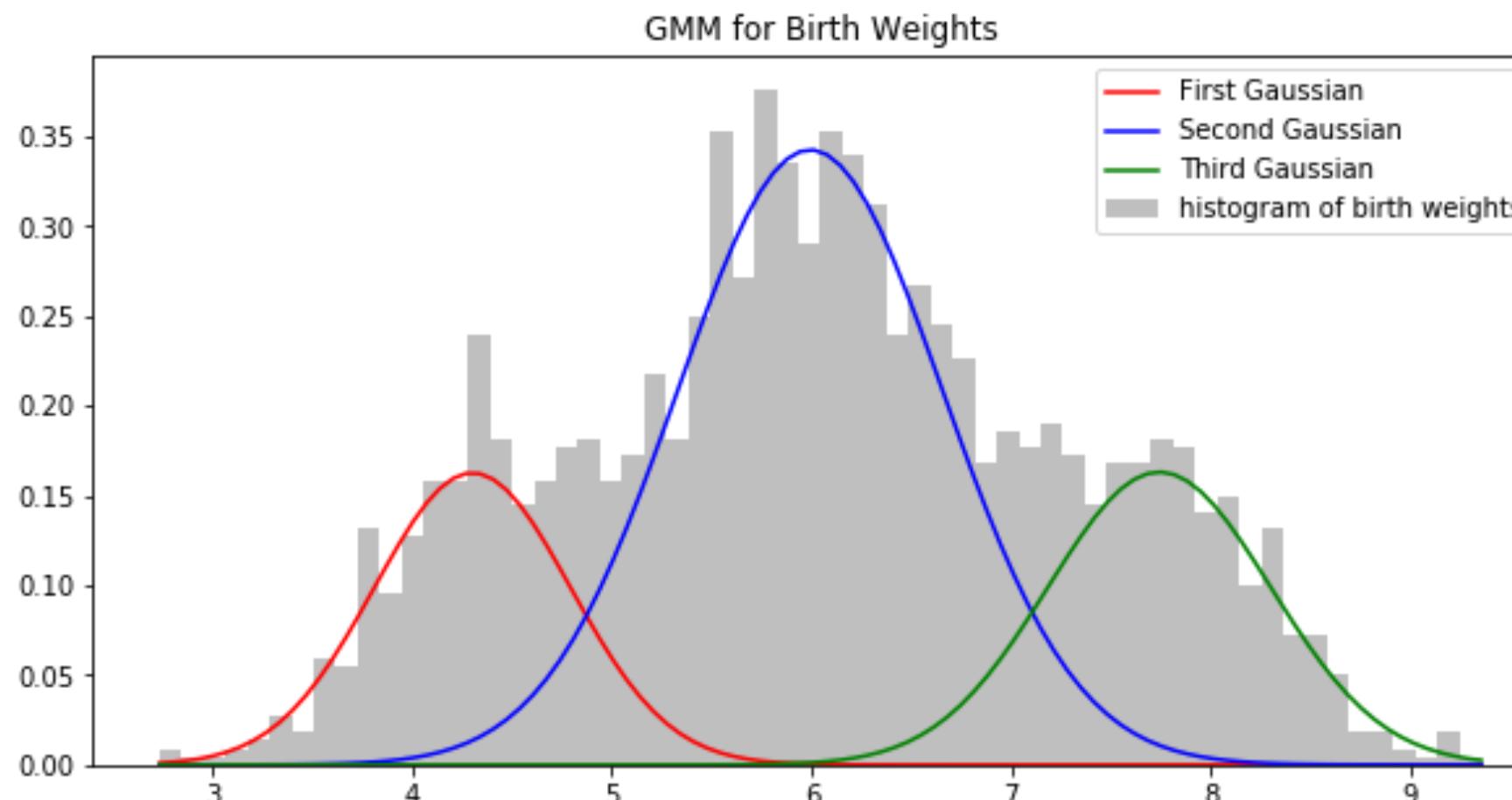
- $\mu_{\text{new}} = \frac{1}{\sum_{n=1}^N r_{n,k}} \sum_{n=1}^N r_{n,k} y_n$
- $\sigma_{\text{new}}^2 = \frac{1}{\sum_{n=1}^N r_{n,k}} \sum_{n=1}^N r_{n,k} (y_n - \mu_{\text{new}})^2$
- $\pi_{\text{new}} = \frac{\sum_{n=1}^N r_{n,k}}{N}$



Implementing EM for the Gaussian Mixture Model of Birth Weight

In [3]:

```
fig, ax = plt.subplots(1, 1, figsize=(10, 5))
ax.hist(y, bins=60, density=True, color='gray', alpha=0.5, label='histogram of birth weights')
ax.plot(x, pi_current[0] * sp.stats.norm(mu_current[0], sigma_current[0]**0.5).pdf(x), color='red', label='First Gaussian')
ax.plot(x, pi_current[1] * sp.stats.norm(mu_current[1], sigma_current[1]**0.5).pdf(x), color='blue', label='Second Gaussian')
ax.plot(x, pi_current[2] * sp.stats.norm(mu_current[2], sigma_current[2]**0.5).pdf(x), color='green', label='Third Gaussian')
ax.set_title('GMM for Birth Weights')
ax.legend(loc='best')
plt.show()
```



Sanity Check: Log-Likelihood During Training

Remember that plotting the MLE model against actual data is not always an option (e.g. high-dimensional data).

A sanity check for that your EM algorithm has been implemented correctly is to plot the observed data log-likelihood over the iterations of the algorithm:

$$\ell_y(\mu, \sigma^2, \pi) = \sum_{n=1}^N \log \sum_{k=1}^K \mathcal{N}(y_n; \mu_k, \sigma_k^2) \pi_k$$

```
In [4]: fig, ax = plt.subplots(1, 1, figsize=(10, 3))
ax.plot(range(len(log_lkhd)), log_lkhd, color='red', alpha=0.5)
ax.set_title('observed data log-likelihood over iterations of EM')
plt.show()
```

