

STA 360/602L: MODULE 8.3

FINITE MIXTURE MODELS: UNIVARIATE CONTINUOUS DATA

DR. OLANREWAJU MICHAEL AKANDE

CONTINUOUS DATA – UNIVARIATE CASE

- Suppose we have univariate continuous data $y_i \stackrel{iid}{\sim} f$, for i, \dots, n , where f is an unknown density.
- Turns out that we can approximate "almost" any f with a **mixture of normals**. Usual choices are

1. **Location mixture** (multimodal):

$$f(y) = \sum_{k=1}^K \lambda_k \mathcal{N}(\mu_k, \sigma^2)$$

2. **Scale mixture** (unimodal and symmetric about the mean, but fatter tails than a regular normal distribution):

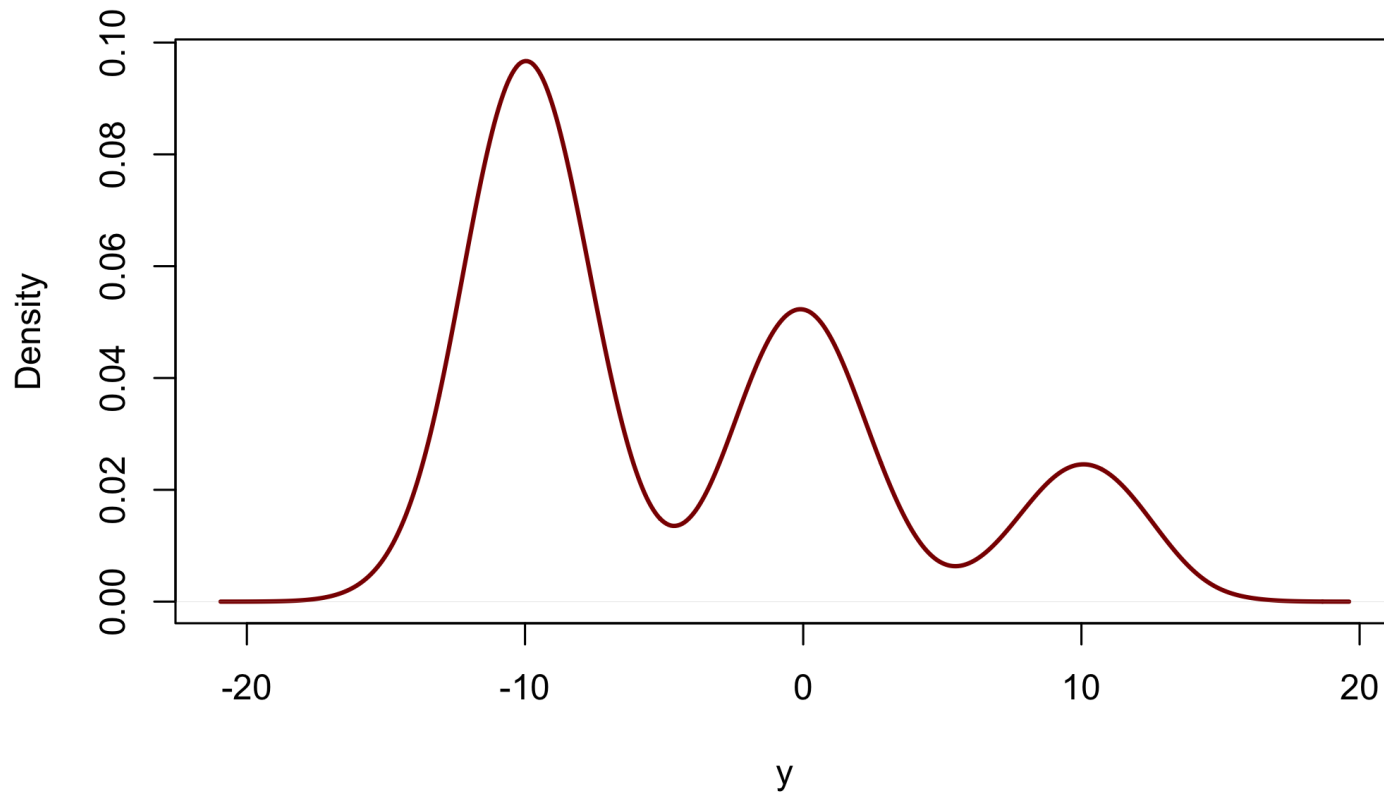
$$f(y) = \sum_{k=1}^K \lambda_k \mathcal{N}(\mu, \sigma_k^2)$$

3. **Location-scale mixture** (multimodal with potentially fat tails):

$$f(y) = \sum_{k=1}^K \lambda_k \mathcal{N}(\mu_k, \sigma_k^2)$$

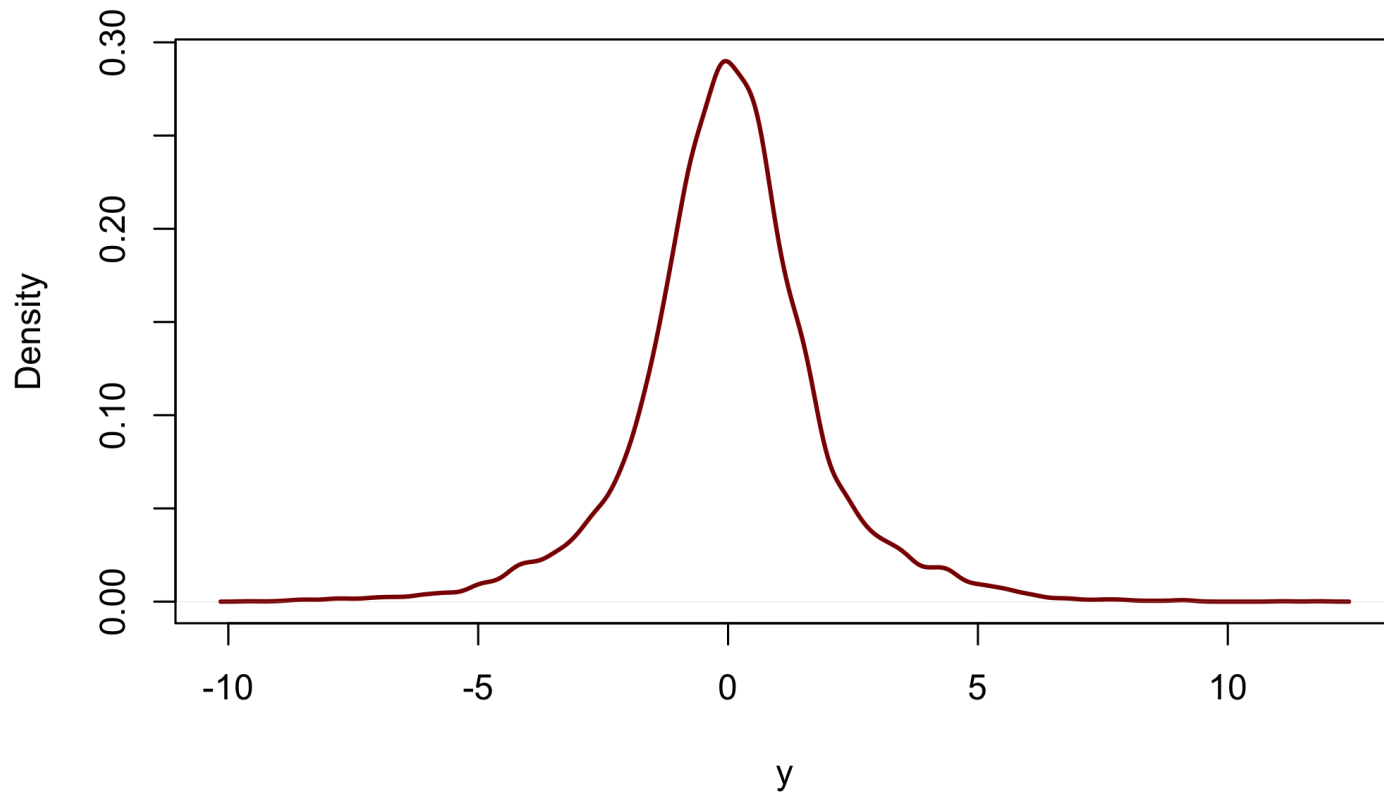
LOCATION MIXTURE EXAMPLE

$$f(y) = 0.55\mathcal{N}(-10, 4) + 0.30\mathcal{N}(0, 4) + 0.15\mathcal{N}(10, 4)$$



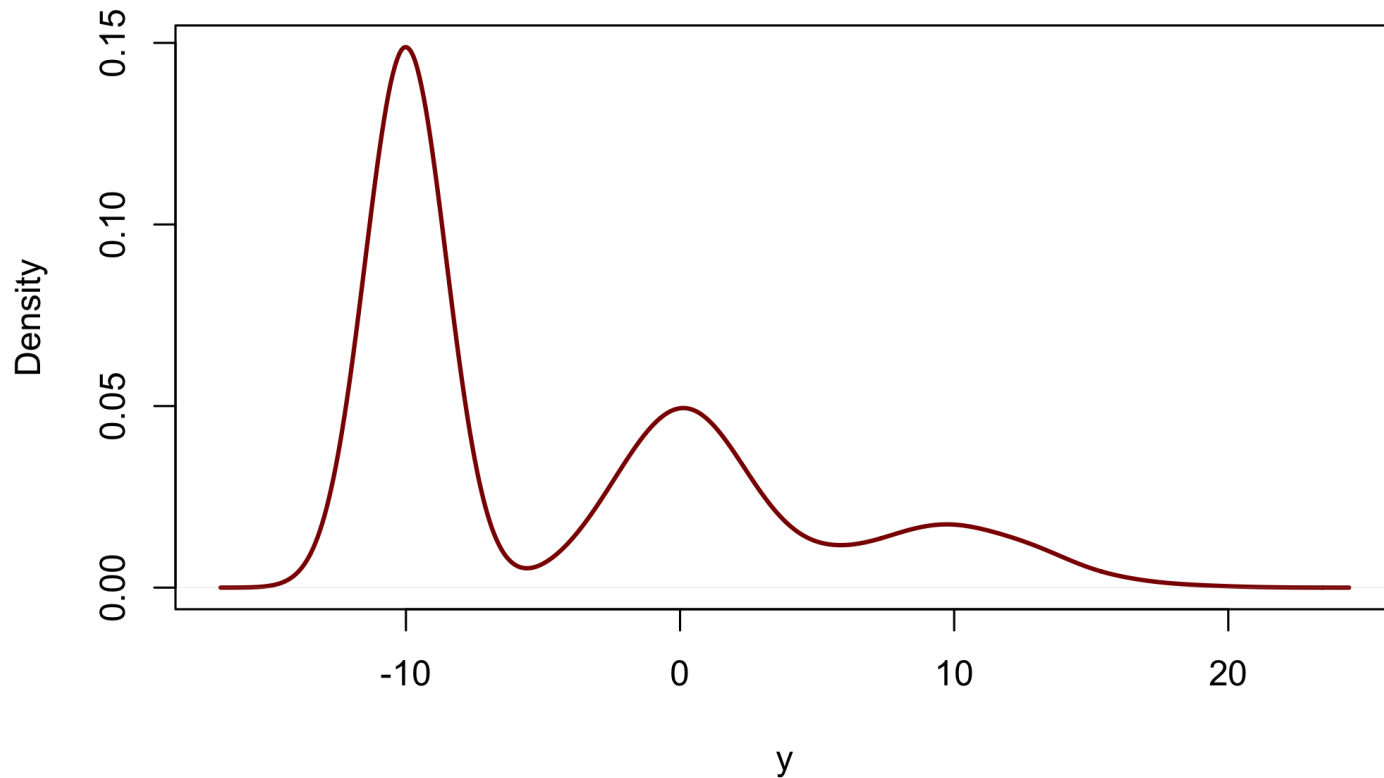
SCALE MIXTURE EXAMPLE

$$f(y) = 0.55\mathcal{N}(0, 1) + 0.30\mathcal{N}(0, 5) + 0.15\mathcal{N}(0, 10)$$



LOCATION-SCALE MIXTURE EXAMPLE

$$f(y) = 0.55\mathcal{N}(-10, 1) + 0.30\mathcal{N}(0, 5) + 0.15\mathcal{N}(10, 10)$$



LOCATION MIXTURE OF NORMALS

- Consider the location mixture $f(y) = \sum_{k=1}^K \lambda_k \mathcal{N}(\mu_k, \sigma^2)$. How can we do inference?
- Right now, we only have three unknowns: $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_K)$, $\boldsymbol{\mu} = (\mu_1, \dots, \mu_K)$, and σ^2 .
- For priors, the most obvious choices are
 - $\pi[\boldsymbol{\lambda}] = \text{Dirichlet}(\alpha_1, \dots, \alpha_K)$,
 - $\mu_k \sim \mathcal{N}(\mu_0, \gamma_0^2)$, for each $k = 1, \dots, K$, and
 - $\sigma^2 \sim \mathcal{IG}\left(\frac{\nu_0}{2}, \frac{\nu_0 \sigma_0^2}{2}\right)$.
- However, we do not want to use the likelihood with the sum in the mixture. We prefer products!

DATA AUGMENTATION

- This brings us to the concept of **data augmentation**, which we actually already used in the mixture of multinomials.
- Data augmentation is a commonly-used technique for designing MCMC samplers using **auxiliary/latent/hidden variables**. Again, we have already seen this.
- **Idea**: introduce variable Z that depends on the distribution of the existing variables in such a way that the resulting conditional distributions, with Z included, are easier to sample from and/or result in better mixing.
- Z 's are just latent/hidden variables that are introduced for the purpose of simplifying/improving the sampler.

DATA AUGMENTATION

- For example, suppose we want to sample from $p(x, y)$, but $p(x|y)$ and/or $p(y|x)$ are complicated.
- Choose $p(z|x, y)$ such that $p(x|y, z)$, $p(y|x, z)$, and $p(z|x, y)$ are easy to sample from. Note that we have $p(x, y, z) = p(z|x, y)p(x, y)$.
- Alternatively, rewrite the model as $p(x, y|z)$ and specify $p(z)$ such that

$$p(x, y) = \int p(x, y|z)p(z)dz,$$

where the resulting $p(x|y, z)$, $p(y|x, z)$, and $p(z|x, y)$ from the joint $p(x, y, z)$ are again easy to sample from.

- Next, construct a Gibbs sampler to sample all three variables (X, Y, Z) from $p(x, y, z)$.
- Finally, throw away the sampled Z 's and from what we know about Gibbs sampling, the samples (X, Y) are from the desired $p(x, y)$.

LOCATION MIXTURE OF NORMALS

- Back to location mixture $f(y) = \sum_{k=1}^K \lambda_k \mathcal{N}(\mu_k, \sigma^2)$.
- Introduce latent variable $z_i \in \{1, \dots, K\}$.
- Then, we have
 - $y_i | z_i \sim \mathcal{N}(\mu_{z_i}, \sigma^2)$, and
 - $\Pr(z_i = k) = \lambda_k \equiv \prod_{k=1}^K \lambda_k^{1_{[z_i=k]}}$.
- How does that help? Well, the observed data likelihood is now

$$\begin{aligned} p[Y = (y_1, \dots, y_n) | Z = (z_1, \dots, z_n), \boldsymbol{\lambda}, \boldsymbol{\mu}, \sigma^2] &= \prod_{i=1}^n p(y_i | z_i, \mu_{z_i}, \sigma^2) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2\sigma^2} (y_i - \mu_{z_i})^2 \right\} \end{aligned}$$

which is much easier to work with.

POSTERIOR INFERENCE

- The joint posterior is

$$\begin{aligned}\pi(Z, \boldsymbol{\mu}, \sigma^2, \boldsymbol{\lambda} | Y) &\propto \left[\prod_{i=1}^n p(y_i | z_i, \mu_{z_i}, \sigma^2) \right] \cdot \Pr(Z | \boldsymbol{\mu}, \sigma^2, \boldsymbol{\lambda}) \cdot \pi(\boldsymbol{\mu}, \sigma^2, \boldsymbol{\lambda}) \\ &\propto \left[\prod_{i=1}^n p(y_i | z_i, \mu_{z_i}, \sigma^2) \right] \cdot \Pr(Z | \boldsymbol{\lambda}) \cdot \pi(\boldsymbol{\lambda}) \cdot \pi(\boldsymbol{\mu}) \cdot \pi(\sigma^2) \\ &\propto \left[\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2\sigma^2} (y_i - \mu_{z_i})^2 \right\} \right] \\ &\quad \times \left[\prod_{i=1}^n \prod_{k=1}^K \lambda_k^{1[z_i=k]} \right] \\ &\quad \times \left[\prod_{k=1}^K \lambda_k^{\alpha_k - 1} \right] \cdot \\ &\quad \times \left[\prod_{k=1}^K \mathcal{N}(\mu_k; \mu_0, \gamma_0^2) \right] \\ &\quad \times \left[\mathcal{IG} \left(\sigma^2; \frac{\nu_0}{2}, \frac{\nu_0 \sigma_0^2}{2} \right) \right] \cdot\end{aligned}$$

FULL CONDITIONALS

- For $i = 1, \dots, n$, sample $z_i \in \{1, \dots, K\}$ from a categorical distribution (multinomial distribution with sample size one) with probabilities

$$\begin{aligned}\Pr[z_i = k | \dots] &= \frac{\Pr[y_i, z_i = k | \mu_k, \sigma^2, \lambda_k]}{\sum_{l=1}^K \Pr[y_i, z_i = l | \mu_l, \sigma^2, \lambda_l]} \\ &= \frac{\Pr[y_i | z_i = k, \mu_k, \sigma^2] \cdot \Pr[z_i = k | \lambda_k]}{\sum_{l=1}^K \Pr[y_i | z_i = l, \mu_l, \sigma^2] \cdot \Pr[z_i = l | \lambda_l]} \\ &= \frac{\lambda_k \cdot \mathcal{N}(y_i; \mu_k, \sigma^2)}{\sum_{l=1}^K \lambda_l \cdot \mathcal{N}(y_i; \mu_l, \sigma^2)}.\end{aligned}$$

- Note that $\mathcal{N}(y_i; \mu_k, \sigma^2)$ just means evaluating the density $\mathcal{N}(\mu_k, \sigma^2)$ at the value y_i .

FULL CONDITIONALS

- Next, sample $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_K)$ from

$$\pi[\boldsymbol{\lambda} | \dots] \equiv \text{Dirichlet}(a_1 + n_1, \dots, a_d + n_d),$$

where $n_k = \sum_{i=1}^n 1[z_i = k]$, the number of individuals assigned to cluster k .

- Sample the mean μ_k for each cluster from

$$\pi[\mu_k | \dots] \equiv \mathcal{N}(\mu_{k,n}, \gamma_{k,n}^2);$$
$$\gamma_{k,n}^2 = \frac{1}{\frac{n_k}{\sigma^2} + \frac{1}{\gamma_0^2}}; \quad \mu_{k,n} = \gamma_{k,n}^2 \left[\frac{n_k}{\sigma^2} \bar{y}_k + \frac{1}{\gamma_0^2} \mu_0 \right],$$

- Finally, sample σ^2 from

$$\pi(\sigma^2 | \dots) = \mathcal{IG}\left(\frac{\nu_n}{2}, \frac{\nu_n \sigma_n^2}{2}\right).$$
$$\nu_n = \nu_0 + n; \quad \sigma_n^2 = \frac{1}{\nu_n} \left[\nu_0 \sigma_0^2 + \sum_{i=1}^n (y_i - \mu_{z_i})^2 \right].$$

WHAT'S NEXT?

MOVE ON TO THE READINGS FOR THE NEXT MODULE!