

# STA 360/602L: MODULE 8.4

## FINITE MIXTURE MODELS: UNIVARIATE CONTINUOUS DATA (ILLUSTRATION)

DR. OLANREWaju MICHAEL AKANDE

# LOCATION MIXTURE OF NORMALS RECAP

- Sampling model with latent variable:

- $y_i | z_i \sim \mathcal{N}(\mu_{z_i}, \sigma^2)$ , and

- $\Pr(z_i = k) = \lambda_k \equiv \prod_{k=1}^K \lambda_k^{1_{[z_i=k]}}$ .

- Priors:

- $\pi[\boldsymbol{\lambda}] = \text{Dirichlet}(\alpha_1, \dots, \alpha_K)$ ,

- $\mu_k \sim \mathcal{N}(\mu_0, \gamma_0^2)$ , for each  $k = 1, \dots, K$ , and

- $\sigma^2 \sim \text{IG}\left(\frac{\nu_0}{2}, \frac{\nu_0 \sigma_0^2}{2}\right)$ .

# FULL CONDITIONALS RECAP

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

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

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

$$\pi[\boldsymbol{\lambda} | \dots] \equiv \text{Dirichlet}(\alpha_1 + n_1, \dots, \alpha_K + n_K),$$

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

# FULL CONDITIONALS RECAP

- Sample the mean  $\mu_k$  for each cluster from

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

- Finally, sample  $\sigma^2$  from

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

# PRACTICAL CONSIDERATIONS

- As we will see in the illustration very soon, the sampler for this model can suffer from label switching.
- For example, suppose our groups are men and women. Then, if we run the sampler multiple times (starting from the same initial values), sometimes it will settle on females as the first group, and sometimes on females are the second group.
- Specifically, MCMC on mixture models in general can suffer from label switching.
- Fortunately, results are still valid if we interpret them correctly.
- Specifically, we should focus on quantities and estimands that are invariant to permutations of the clusters. For example, look at marginal quantities, instead of conditional ones.

# OTHER PRACTICAL CONSIDERATIONS

- So far we have assumed that the number of clusters  $K$  is known.
- What if we don't know  $K$ ?
  - Compare marginal likelihood for different choices of  $K$  and select  $K$  with best performance.
  - Can also use other metrics, such as MSE, and so on.
  - Maybe a prior on  $K$ ?
  - Go Bayesian non-parametric: **Dirichlet processes!**

MOVE TO THE R SCRIPT HERE.

# WHAT'S NEXT?

MOVE ON TO THE READINGS FOR THE NEXT MODULE!