# STATS271/371: Applied Bayesian Statistics
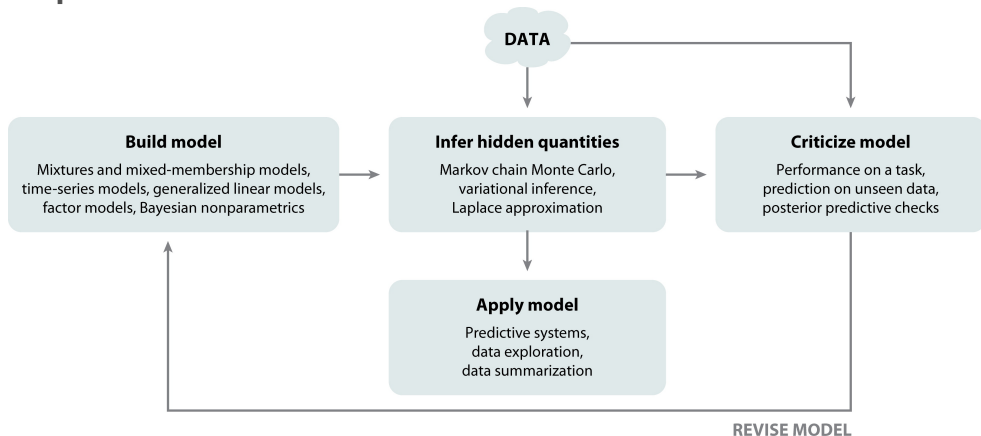
**Bayesian Mixture Models and (Collapsed) Gibbs Sampling**

Scott Linderman

April 21, 2021

# Box's Loop



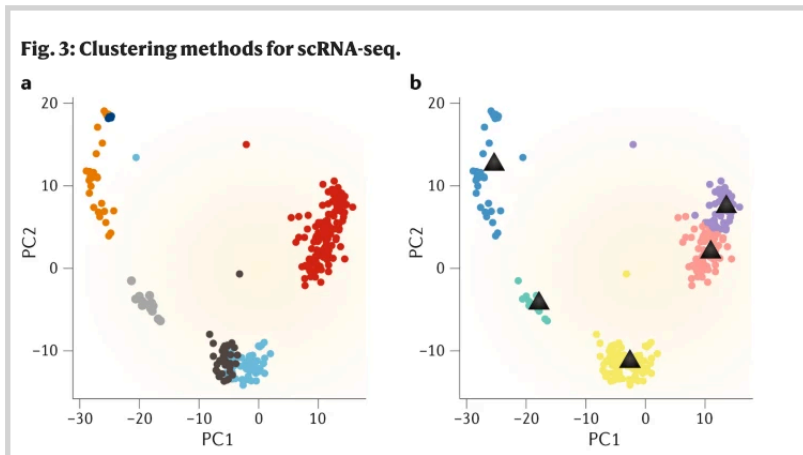**DATA**

**Build model**
Mixtures and mixed-membership models, time-series models, generalized linear models, factor models, Bayesian nonparametrics

**Infer hidden quantities**
Markov chain Monte Carlo, variational inference, Laplace approximation

**Criticize model**
Performance on a task, prediction on unseen data, posterior predictive checks

**Apply model**
Predictive systems, data exploration, data summarization

**REVISE MODEL**

Blei DM. 2014.
Annu. Rev. Stat. Appl. 1:203–32

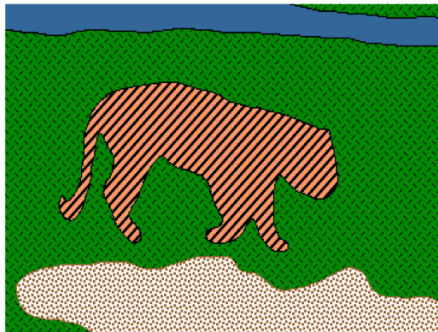# Lap 4: Bayesian Mixture Models and (Collapsed) Gibbs Sampling

- ▶ **Model:** Bayesian mixture models
- ▶ **Algorithm:** Gibbs sampling
- ▶ **Criticism:** Posterior predictive checks
- ▶ **Algorithm II:** Collapsed Gibbs sampling
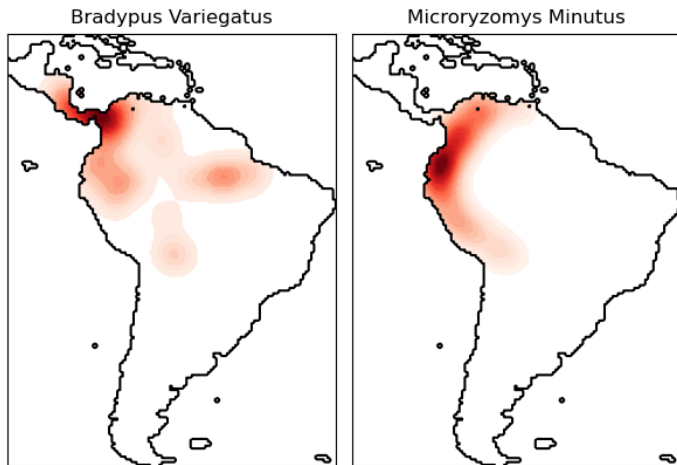
# Motivation: Clustering scRNA-seq data



Fig. 3: Clustering methods for scRNA-seq.

From Kiselev et al. [2019]

# Motivation: Foreground/background segmentation



https://ai.stanford.edu/~syyeung/cvweb/tutorial3.html

## Motivation: Density estimation



Bradypus Variegatus      Microryzomys Minutus

https://scikit-learn.org/stable/auto_examples/neighbors/plot_species_kde.html

# Notation

**Constants:** Let

- ▶ $N$ denote the number of data points.

- ▶ $K$ denote the number of mixture components (i.e. clusters)

**Data:** Let

- ▶ $x_n \in \mathbb{R}^D$ denote the $n$-th data point.

**Latent Variables:** Let

- ▶ $z_n \in \{1, \ldots, K\}$ denote the *assignment* of the $n$-th data point.

# Notation II

**Parameters:** Let

- ▶ $\eta_k$ denote the *natural parameters* of component $k$

- ▶ $\pi \in \Delta_K$ denote the component *proportions* (i.e. probabilities).

**Hyperparameters:** Let

- ▶ $\phi$, $\nu$ denote hyperparameters of the prior on $\eta$

- ▶ $\alpha \in \mathbb{R}_+^K$ denote the concentration of the prior on proportions.

## Generative Model

1. Sample the proportions from a Dirichlet prior:

$$\pi \sim \mathrm{Dir}(\boldsymbol{\alpha}) \tag{1}$$

2. Sample the parameters for each component:

$$\eta_k \overset{\mathrm{iid}}{\sim} p(\boldsymbol{\eta} \mid \boldsymbol{\phi}, \nu) \qquad \text{for } k = 1, \dots, K \tag{2}$$

3. Sample the assignment of each data point:

$$z_n \overset{\mathrm{iid}}{\sim} \pi \qquad \qquad \text{for } n = 1, \dots, N \tag{3}$$

4. Sample data points given their assignments:

$$\boldsymbol{x}_n \sim p(\boldsymbol{x} \mid \eta_{z_n}) \qquad \text{for } n = 1, \dots, N \tag{4}$$

## Joint distribution

This generative model corresponds to the following factorization of the joint distribution,

$$p(\pi, \{\eta_k\}_{k=1}^K, \{(z_n, \boldsymbol{x}_n)\}_{n=1}^N \mid \boldsymbol{\phi}, \nu, \boldsymbol{\alpha}) = p(\pi \mid \boldsymbol{\alpha}) \prod_{k=1}^K p(\eta_k \mid \boldsymbol{\phi}, \nu) \prod_{n=1}^N p(z_n \mid \pi) p(\boldsymbol{x}_n \mid \boldsymbol{z}_n, \{\eta_k\}_{k=1}^K) \tag{5}$$

Equivalently,

$$p(\pi, \{\eta_k\}_{k=1}^K, \{(z_n, \boldsymbol{x}_n)\}_{n=1}^N \mid \boldsymbol{\phi}, \nu, \boldsymbol{\alpha}) =$$
$$p(\pi \mid \boldsymbol{\alpha}) \prod_{k=1}^K p(\eta_k \mid \boldsymbol{\phi}, \nu) \prod_{n=1}^N \prod_{k=1}^K [\Pr(z_n = k \mid \pi) p(\boldsymbol{x}_n \mid \eta_k)]^{\mathbb{I}[z_n = k]} \tag{6}$$

Substituting in the assumed forms

$$p(\pi, \{\eta_k\}_{k=1}^K, \{(z_n, \boldsymbol{x}_n)\}_{n=1}^N \mid \boldsymbol{\phi}, \boldsymbol{\alpha}) = \mathrm{Dir}(\pi \mid \boldsymbol{\alpha}) \prod_{k=1}^K p(\eta_k \mid \boldsymbol{\phi}, \nu) \prod_{n=1}^N \prod_{k=1}^K [\pi_k p(\boldsymbol{x}_n \mid \eta_k)]^{\mathbb{I}[z_n = k]} \tag{7}$$

## Exponential family mixture models

What about $p(\boldsymbol{x} \mid \boldsymbol{\eta}_k)$ and $p(\boldsymbol{\eta}_k \mid \boldsymbol{\phi}, \nu)$?

Recall the *exponential family* distributions from Lap 2. Let's assume an exponential family likelihood,

$$p(\boldsymbol{x} \mid \boldsymbol{\eta}_k) = h(\boldsymbol{x}_n) \exp \left\{ \langle t(\boldsymbol{x}_n), \boldsymbol{\eta}_k \rangle - A(\boldsymbol{\eta}_k) \right\}. \tag{8}$$

Then assume a *conjugate prior*,

$$p(\boldsymbol{\eta}_k \mid \boldsymbol{\phi}, \nu) \propto \exp \left\{ \langle \boldsymbol{\phi}, \boldsymbol{\eta}_k \rangle - \nu A(\boldsymbol{\eta}_k) \right\}. \tag{9}$$

The hyperparmeters $\boldsymbol{\phi}$ are *pseudo-observations* of the sufficient statistics (like statistics from fake data points) and $\nu$ is a *pseudo-count* (like the number of fake data points).

Note that the product of prior and likelihood remains in the same family as the prior. That's why we call it conjugate.

## Example: Gaussian mixture model

Assume the conditional distribution of $\boldsymbol{x}_n$ is a Gaussian with mean $\boldsymbol{\eta}_{z_n} \in \mathbb{R}^D$ and identity covariance,

$$p(\boldsymbol{x}_n \mid \boldsymbol{\eta}_k) = \mathcal{N}(\boldsymbol{x}_n \mid \boldsymbol{\eta}_k, \boldsymbol{I}) \tag{10}$$

$$= (2\pi)^{-D/2} \exp\left\{-\tfrac{1}{2}(\boldsymbol{x}_n - \boldsymbol{\eta}_k)^\top (\boldsymbol{x}_n - \boldsymbol{\eta}_k)\right\} \tag{11}$$

$$= (2\pi)^{-D/2} \exp\left\{-\tfrac{1}{2}\boldsymbol{x}_n^\top \boldsymbol{x}_n + \boldsymbol{x}_n^\top \boldsymbol{\eta}_k - \tfrac{1}{2}\boldsymbol{\eta}_k^\top \boldsymbol{\eta}_k\right\}, \tag{12}$$

which is an exponential family distribution with base measure $h(\boldsymbol{x}_n) = (2\pi)^{-D/2}e^{-\frac{1}{2}\boldsymbol{x}_n^\top \boldsymbol{x}_n}$, sufficient statistics $t(\boldsymbol{x}_n) = \boldsymbol{x}_n$, and log normalizer $A(\boldsymbol{\eta}_k) = \tfrac{1}{2}\boldsymbol{\eta}_k^\top \boldsymbol{\eta}_k$.

Then assume a Gaussian prior on the component parameters. It's conjugate,

$$p(\boldsymbol{\eta}_k \mid \boldsymbol{\phi}, \nu) = \mathcal{N}(\nu^{-1}\boldsymbol{\phi}, \nu^{-1}\boldsymbol{I}) \propto \exp\left\{\boldsymbol{\phi}^\top \boldsymbol{\eta}_k - \tfrac{\nu}{2}\boldsymbol{\eta}_k^\top \boldsymbol{\eta}_k\right\} = \exp\left\{\boldsymbol{\phi}^\top \boldsymbol{\eta}_k - \nu A(\boldsymbol{\eta}_k)\right\}. \tag{13}$$

Note that $\boldsymbol{\phi}$ sets the location and $\nu$ sets the precision (i.e. inverse variance).

## MAP inference via coordinate ascent

Before diving into fully Bayesian inference algorithms, let's first consider **MAP inference**.

**Idea:** find the mode of $p(\pi, \{\eta_k\}_{k=1}^K, \{z_n\}_{n=1}^N \mid \{x_n\}_{n=1}^N, \phi, \nu, \alpha)$ by **coordinate ascent**.

For now, set $\phi = 0$, and $\nu = 0$ so that the prior is an (improper) uniform distribution. Then maximizing the posterior is equivalent to maximizing the likelihood.

While we're simplifying, let's even fix $\pi = \frac{1}{K}\mathbf{1}_K$.

## Coordinate ascent in the Gaussian mixture model

For the Gaussian mixture model (with uniform prior and $\pi = \frac{1}{K}\mathbf{1}_K$), coordinate ascent amounts to:

**1.** For each $n = 1, \ldots, N$, fix all variables but $z_n$ and find $z_n^\star$ that maximizes

$$p(\pi, \{\eta_k\}_{k=1}^K, \{(z_n, x_n)\}_{n=1}^N \mid \phi, \nu, \alpha) \propto p(x_n \mid z_n, \{\eta_k\}_{k=1}^K) = \mathcal{N}(x_n \mid \eta_{z_n}, I) \qquad (14)$$

The cluster assignment that maximizes the likelihood is the one with the closest mean to $x_n$, so set

$$z_n^\star = \underset{k \in \{1, \ldots, K\}}{\arg\min} \|x_n - \eta_k\|_2. \qquad (15)$$

## Coordinate ascent in the Gaussian mixture model II

**2** For each $k = 1, \ldots, K$, fix all variables but $\eta_k$ and find $\eta_k^\star$ that maximizes,

$$p(\pi, \{\eta_k\}_{k=1}^K, \{(z_n, \mathbf{x}_n)\}_{n=1}^N \mid \phi, \nu, \alpha) \propto \prod_{n=1}^N p(\mathbf{x}_n \mid \eta_k)^{\mathbb{I}[z_n = k]} \tag{16}$$

$$\propto \exp\left\{ \sum_{n=1}^N \mathbb{I}[z_n = k] \left( \mathbf{x}_n^\top \eta_k - \tfrac{1}{2} \eta_k^\top \eta_k \right) \right\} \tag{17}$$

Taking the derivative of the log and setting to zero yields,

$$\eta_k^\star = \frac{1}{N_k} \sum_{n=1}^K \mathbb{I}[z_n = k] \mathbf{x}_n, \tag{18}$$

where $N_k = \sum_{n=1}^N \mathbb{I}[z_n = k]$.

This is the **k-means algorithm**!

## Aside: EM in the Gaussian mixture model

We'll talk more about *coordinate ascent variational inference* (CAVI) and *expectation-maximization* (EM) next week. Not to spoil the surprise, but we'll see that they have a similar flavor. Instead of assigning $z_n^\star$ to the closest cluster, we compute *responsibilities* for each cluster:

**1.** For each data point $n$ and component $k$, set the *responsibility* to,

$$\omega_{nk} = \frac{\pi_k \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\eta}_k, \boldsymbol{I})}{\sum_{j=1}^{K} \pi_j \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\eta}_j, \boldsymbol{I})}. \tag{19}$$

**2.** For each component $k$, set the mean to

$$\boldsymbol{\eta}_k^\star = \frac{1}{N_k} \sum_{n=1}^{K} \omega_{nk} \mathbf{x}_n, \tag{20}$$

where $N_k = \sum_{n=1}^{N} \omega_{nk}$.

Note that EM allows for arbitrary proportions $\pi$. Those can be updated as well: for each component $k$, set $\pi_k = \frac{N_k}{N}$.

# Lap 4: Bayesian Mixture Models and (Collapsed) Gibbs Sampling

▶ Model: Bayesian mixture models

▶ **Algorithm: Gibbs sampling**

▶ Criticism: Posterior predictive checks

▶ Algorithm II: Collapsed Gibbs sampling

## Gibbs sampling in Bayesian mixture models

**Idea:** just like in coordinate ascent, update one variable at a time. *But rather than setting it to its conditional mode, sample from its conditional distribution.*

## Gibbs sampling in Bayesian mixture models II

1. For each data point $n$, sample a new assignment from the complete conditional distribution

$$z_n \sim p(z_n \mid \{\mathbf{x}_n\}_{n=1}^N, \{z_{n'}\}_{n' \neq n}, \{\boldsymbol{\eta}_k\}_{k=1}^K, \boldsymbol{\pi}, \boldsymbol{\phi}, \nu, \boldsymbol{\alpha}). \tag{21}$$

Thanks to the factorization of the joint distribution,

$$\Pr(z_n = k \mid -) \propto \Pr(z_n = k \mid \boldsymbol{\pi}) \, p(x_n \mid \boldsymbol{\eta}_k) \tag{22}$$

In the Gaussian mixture model, this is,

$$\Pr(z_n = k \mid -) = \frac{\pi_k \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\eta}_k, \boldsymbol{I})}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\eta}_j, \boldsymbol{I})} \tag{23}$$

$$\equiv \omega_{nk}. \tag{24}$$

I.e., Gibbs sampling generates *random* assignments by sampling according to the responsibilities.

## Gibbs sampling in Bayesian mixture models III

**2** For each component *k*, sample new parameters from their complete conditional,

$$\boldsymbol{\eta}_k \sim p(\boldsymbol{\eta}_k \mid \{(\mathbf{x}_n, z_n)\}_{n=1}^N, \{\boldsymbol{\eta}_{k'}\}_{k' \neq k}, \boldsymbol{\pi}, \boldsymbol{\phi}, \nu, \boldsymbol{\alpha}). \tag{25}$$

Thanks to the factorization of the joint distribution,

$$p(\boldsymbol{\eta}_k \mid -) \propto p(\boldsymbol{\eta}_k \mid \boldsymbol{\phi}, \nu) \prod_{n: z_n = k} p(x_n \mid \boldsymbol{\eta}_k). \tag{26}$$

In an Gaussian mixture model,

$$p(\boldsymbol{\eta}_k \mid -) \propto \exp\left\{ \left( \boldsymbol{\phi} + \sum_{n=1}^N \mathbb{I}[z_n = k]\mathbf{x}_n \right)^\top \boldsymbol{\eta}_k - \frac{\nu + N_k}{2} \boldsymbol{\eta}_k^\top \boldsymbol{\eta}_k \right\} \tag{27}$$

$$\propto \mathcal{N}\left( \boldsymbol{\eta}_k \mid (\nu + N_k)^{-1}\left( \boldsymbol{\phi} + \sum_{n=1}^N \mathbb{I}[z_n = k]\mathbf{x}_n \right), (\nu + N_k)^{-1}\boldsymbol{I} \right) \tag{28}$$

where $N_k = \sum_{n=1}^N \mathbb{I}[z_n = k]$. What happens when $N_k \to \infty$?

## Gibbs sampling in Bayesian mixture models IV

3 Finally, sample new component proportions from their complete conditional,

$$\pi \sim p(\pi \mid \{(x_n, z_n)\}_{n=1}^N, \{\eta_k\}_{k=1}^K, \phi, \nu, \alpha). \tag{29}$$

Thanks to the factorization of the joint distribution,

$$p(\pi \mid -) \propto \text{Dir}(\pi \mid \alpha) \prod_{n=1}^N p(z_n \mid \pi) \tag{30}$$

$$\propto \prod_{k=1}^K \pi_k^{\alpha_k - 1} \times \prod_{n=1}^N \prod_{k=1}^K \pi_k^{\mathbb{I}[z_n = k]} \tag{31}$$

$$\propto \text{Dir}\left(\pi \mid [\alpha_1 + N_1, \ldots, \alpha_K + N_K]\right) \tag{32}$$

where $N_k = \sum_{n=1}^N \mathbb{I}[z_n = k]$. What happens when $N_k \to \infty$?

## Gibbs sampling in Bayesian exponential family mixture models

What happens in general exponential family models? Step 2 becomes,

**2** For each component $k$, sample new parameters from their complete conditional,

$$p(\boldsymbol{\eta}_k \mid -) \propto \exp\left\{ \left\langle \boldsymbol{\phi} + \sum_{n=1}^{N} \mathbb{I}[z_n = k] t(\boldsymbol{x}_n), \boldsymbol{\eta}_k \right\rangle - (\nu + N_k) A(\boldsymbol{\eta}_k) \right\} \tag{33}$$

$$= p\left( \boldsymbol{\eta}_k \,\middle|\, \boldsymbol{\phi} + \sum_{n=1}^{N} \mathbb{I}[z_n = k] \, t(\boldsymbol{x}_n), \, \nu + N_k \right) \tag{34}$$

where $N_k = \sum_{n=1}^{N} \mathbb{I}[z_n = k]$.

## Opportunities for parallelism

▶ In all three algorithms above, the updates of $z_n$ are independent of one another (once you fix $\{\eta_k\}_{k=1}^K$) and hence can be performed in parallel.

▶ Likewise, the updates of $\eta_k$ are independent of one another (once you fix $\{z_n\}_{n=1}^N$) and hence can be performed in parallel.

▶ In fact, we can write these as simple map-reduce algorithms and take advantage of parallel hardware if it's available.

▶ In the Gibbs sampling case, updating many variables at once from their combined conditional distribution is called **blocked Gibbs sampling**, and it's particularly easy when the variables are conditionally independent, as in the Bayesian mixture model.

# Next time

► Show that **Gibbs is a special case of MH** and, as such, asymptotically generates samples from the posterior distribution.

► Talk about **posterior predictive checks** and ways of choosing $K$.

► Introduce **collapsed Gibbs sampling**, which will allow us to generalize to **nonparametric Bayesian mixture models**.

# References I

Vladimir Yu Kiselev, Tallulah S Andrews, and Martin Hemberg. Challenges in unsupervised clustering of single-cell RNA-seq data. *Nat. Rev. Genet.*, 20(5):273–282, May 2019.