

# **Bayesian Mixture Models and Expectation Maximization**

## **STATS 305C: Applied Statistics**

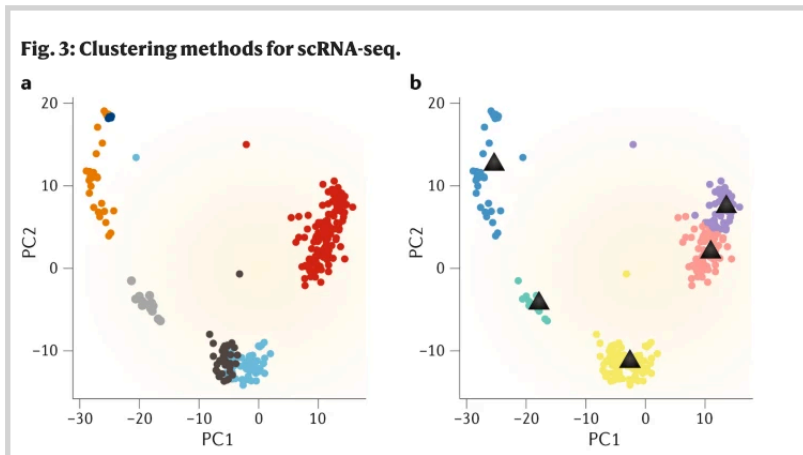
Scott Linderman

April 20, 2022

# Outline

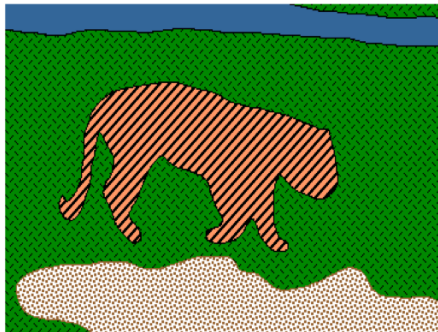
- ▶ Model: Bayesian mixture models
- ▶ Algorithm: MAP Estimation / K-Means

## Motivation: Clustering scRNA-seq data



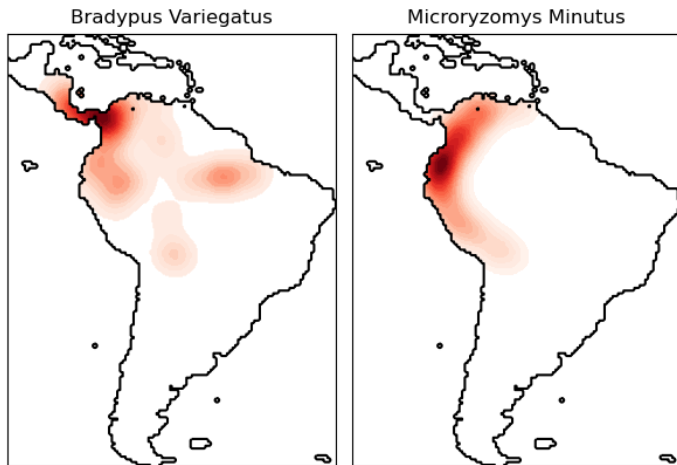
From Kiselev et al. [2019]

## Motivation: Foreground/background segmentation



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

## Motivation: Density estimation



# Notation

**Constants:** Let

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

**Data:** Let

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

**Latent Variables:** Let

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

# Notation II

## Parameters: Let

- ▶  $\theta_k$  denote the *natural parameters* of component  $k$
- ▶  $\pi \in \Delta_{K-1}$  denote the component *proportions* (i.e. probabilities).

## Hyperparameters: Let

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

# Generative Model

1. Sample the proportions from a Dirichlet prior:

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



# The beta distribution

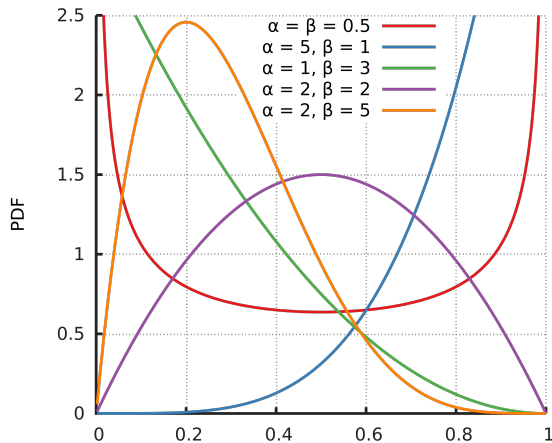


Figure: The beta distribution over  $\pi \in [0, 1]$  is a special case of the Dirichlet distribution.

[https://en.wikipedia.org/wiki/Beta\\_distribution](https://en.wikipedia.org/wiki/Beta_distribution)

# The Dirichlet distribution

If the beta distribution generates weighted coins, the Dirichlet generates weighted dice.

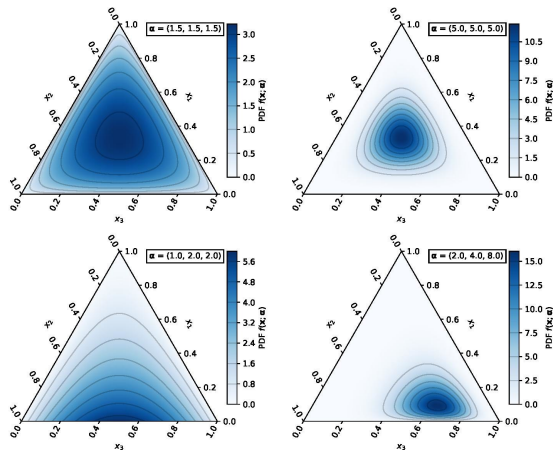


Figure: The Dirichlet distribution over  $\pi \in \Delta_2$ ; i.e. distributions over  $K = 3$  outcomes. From [https://en.wikipedia.org/wiki/Dirichlet\\_distribution](https://en.wikipedia.org/wiki/Dirichlet_distribution)

# Generative Model

1. Sample the proportions from a Dirichlet prior:

$$\boldsymbol{\pi} \sim \text{Dir}(\boldsymbol{\alpha}) \quad (2)$$

2. Sample the parameters for each component:

$$\boldsymbol{\theta}_k \stackrel{\text{iid}}{\sim} p(\boldsymbol{\theta} \mid \boldsymbol{\phi}, \nu) \quad \text{for } k = 1, \dots, K \quad (3)$$

3. Sample the assignment of each data point:

$$z_n \stackrel{\text{iid}}{\sim} \boldsymbol{\pi} \quad \text{for } n = 1, \dots, N \quad (4)$$

4. Sample data points given their assignments:

$$\mathbf{x}_n \sim p(\mathbf{x} \mid \boldsymbol{\theta}_{z_n}) \quad \text{for } n = 1, \dots, N \quad (5)$$

## Joint distribution

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

$$p(\pi, \{\theta_k\}_{k=1}^K, \{(z_n, \mathbf{x}_n)\}_{n=1}^N \mid \phi, \nu, \alpha) = p(\pi \mid \alpha) \prod_{k=1}^K p(\theta_k \mid \phi, \nu) \prod_{n=1}^N p(z_n \mid \pi) p(\mathbf{x}_n \mid z_n, \{\theta_k\}_{k=1}^K) \quad (6)$$

- Equivalently,

$$p(\pi, \{\theta_k\}_{k=1}^K, \{(z_n, \mathbf{x}_n)\}_{n=1}^N \mid \phi, \nu, \alpha) = p(\pi \mid \alpha) \prod_{k=1}^K p(\theta_k \mid \phi, \nu) \prod_{n=1}^N \prod_{k=1}^K [\Pr(z_n = k \mid \pi) p(\mathbf{x}_n \mid \theta_k)]^{\mathbb{I}[z_n=k]} \quad (7)$$

- Substituting in the assumed forms

$$p(\pi, \{\theta_k\}_{k=1}^K, \{(z_n, \mathbf{x}_n)\}_{n=1}^N \mid \phi, \nu, \alpha) = \text{Dir}(\pi \mid \alpha) \prod_{k=1}^K p(\theta_k \mid \phi, \nu) \prod_{n=1}^N \prod_{k=1}^K [\pi_k p(\mathbf{x}_n \mid \theta_k)]^{\mathbb{I}[z_n=k]}$$

# Joint distribution

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

$$p(\pi, \{\theta_k\}_{k=1}^K, \{(z_n, \mathbf{x}_n)\}_{n=1}^N \mid \phi, \nu, \alpha) = p(\pi \mid \alpha) \prod_{k=1}^K p(\theta_k \mid \phi, \nu) \prod_{n=1}^N p(z_n \mid \pi) p(\mathbf{x}_n \mid z_n, \{\theta_k\}_{k=1}^K) \quad (6)$$

- Equivalently,

$$p(\pi, \{\theta_k\}_{k=1}^K, \{(z_n, \mathbf{x}_n)\}_{n=1}^N \mid \phi, \nu, \alpha) = p(\pi \mid \alpha) \prod_{k=1}^K p(\theta_k \mid \phi, \nu) \prod_{n=1}^N \prod_{k=1}^K [\Pr(z_n = k \mid \pi) p(\mathbf{x}_n \mid \theta_k)]^{\mathbb{I}[z_n=k]} \quad (7)$$

- Substituting in the assumed forms

$$p(\pi, \{\theta_k\}_{k=1}^K, \{(z_n, \mathbf{x}_n)\}_{n=1}^N \mid \phi, \nu, \alpha) = \text{Dir}(\pi \mid \alpha) \prod_{k=1}^K p(\theta_k \mid \phi, \nu) \prod_{n=1}^N \prod_{k=1}^K [\pi_k p(\mathbf{x}_n \mid \theta_k)]^{\mathbb{I}[z_n=k]}$$

## Joint distribution

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

$$p(\pi, \{\theta_k\}_{k=1}^K, \{(z_n, \mathbf{x}_n)\}_{n=1}^N \mid \phi, \nu, \alpha) = p(\pi \mid \alpha) \prod_{k=1}^K p(\theta_k \mid \phi, \nu) \prod_{n=1}^N p(z_n \mid \pi) p(\mathbf{x}_n \mid z_n, \{\theta_k\}_{k=1}^K) \quad (6)$$

- Equivalently,

$$p(\pi, \{\theta_k\}_{k=1}^K, \{(z_n, \mathbf{x}_n)\}_{n=1}^N \mid \phi, \nu, \alpha) = p(\pi \mid \alpha) \prod_{k=1}^K p(\theta_k \mid \phi, \nu) \prod_{n=1}^N \prod_{k=1}^K [\Pr(z_n = k \mid \pi) p(\mathbf{x}_n \mid \theta_k)]^{\mathbb{I}[z_n=k]} \quad (7)$$

- Substituting in the assumed forms

$$p(\pi, \{\theta_k\}_{k=1}^K, \{(z_n, \mathbf{x}_n)\}_{n=1}^N \mid \phi, \nu, \alpha) = \text{Dir}(\pi \mid \alpha) \prod_{k=1}^K p(\theta_k \mid \phi, \nu) \prod_{n=1}^N \prod_{k=1}^K [\pi_k p(\mathbf{x}_n \mid \theta_k)]^{\mathbb{I}[z_n=k]}$$

# Exponential family mixture models

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

Let's assume an **exponential family** likelihood,

$$p(\mathbf{x} \mid \boldsymbol{\theta}_k) = h(\mathbf{x}_n) \exp \left\{ \langle t(\mathbf{x}_n), \boldsymbol{\theta}_k \rangle - A(\boldsymbol{\theta}_k) \right\}. \quad (9)$$

Then assume a **conjugate prior**,

$$p(\boldsymbol{\theta}_k \mid \boldsymbol{\phi}, \nu) \propto \exp \left\{ \langle \boldsymbol{\phi}, \boldsymbol{\theta}_k \rangle - \nu A(\boldsymbol{\theta}_k) \right\}. \quad (10)$$

The hyperparameters  $\boldsymbol{\phi}$  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  $\mathbf{x}_n$  is a Gaussian with mean  $\boldsymbol{\theta}_k \in \mathbb{R}^D$  and identity covariance,

$$p(\mathbf{x}_n | \boldsymbol{\theta}_k) = \mathcal{N}(\mathbf{x}_n | \boldsymbol{\theta}_k, I) \quad (11)$$

$$= (2\pi)^{-D/2} \exp \left\{ -\frac{1}{2} (\mathbf{x}_n - \boldsymbol{\theta}_k)^\top (\mathbf{x}_n - \boldsymbol{\theta}_k) \right\} \quad (12)$$

$$= (2\pi)^{-D/2} \exp \left\{ -\frac{1}{2} \mathbf{x}_n^\top \mathbf{x}_n + \mathbf{x}_n^\top \boldsymbol{\theta}_k - \frac{1}{2} \boldsymbol{\theta}_k^\top \boldsymbol{\theta}_k \right\}, \quad (13)$$

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

The conjugate prior is a Gaussian prior on the mean,

$$p(\boldsymbol{\theta}_k | \boldsymbol{\phi}, \nu) = \mathcal{N}(\nu^{-1} \boldsymbol{\phi}, \nu^{-1} I) \propto \exp \left\{ \boldsymbol{\phi}^\top \boldsymbol{\theta}_k - \frac{\nu}{2} \boldsymbol{\theta}_k^\top \boldsymbol{\theta}_k \right\} = \exp \left\{ \boldsymbol{\phi}^\top \boldsymbol{\theta}_k - \nu A(\boldsymbol{\theta}_k) \right\}. \quad (14)$$

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



# Outline

- ▶ Model: Bayesian mixture models
- ▶ **Algorithm: MAP Estimation / K-Means**

# MAP inference via coordinate ascent

Let's first consider **maximum a posteriori (MAP) inference**.

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

For now, set  $\phi = \mathbf{0}$ , 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  $n = 1, \dots, N$ , fix all variables but  $z_n$  and find  $z_n^\star$  that maximizes

$$p(\pi, \{\boldsymbol{\theta}_k\}_{k=1}^K, \{(z_n, \mathbf{x}_n)\}_{n=1}^N \mid \boldsymbol{\phi}, \nu, \boldsymbol{\alpha}) \propto p(\mathbf{x}_n \mid z_n, \{\boldsymbol{\theta}_k\}_{k=1}^K) = \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\theta}_{z_n}, I) \quad (15)$$

The cluster assignment that maximizes the likelihood is the one with the closest mean to  $\mathbf{x}_n$ , so set

$$z_n^\star = \arg \min_{k \in \{1, \dots, K\}} \|\mathbf{x}_n - \boldsymbol{\theta}_k\|_2. \quad (16)$$

## Coordinate ascent in the Gaussian mixture model II

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

$$p(\pi, \{\theta_k\}_{k=1}^K, \{(z_n, \mathbf{x}_n)\}_{n=1}^N \mid \phi, \nu, \alpha) \propto \prod_{n=1}^N p(\mathbf{x}_n \mid \theta_k)^{\mathbb{I}[z_n=k]} \quad (17)$$

$$\propto \exp \left\{ \sum_{n=1}^N \mathbb{I}[z_n = k] \left( \mathbf{x}_n^\top \theta_k - \frac{1}{2} \theta_k^\top \theta_k \right) \right\} \quad (18)$$

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

$$\theta_k^* = \frac{1}{N_k} \sum_{n=1}^K \mathbb{I}[z_n = k] \mathbf{x}_n, \quad (19)$$

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

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

## EM in the Gaussian mixture model

K-Means made **hard assignments** of data points to clusters in each iteration. What if we used **soft assignments** instead?

Instead of assigning  $z_n^*$  to the closest cluster, we compute *responsibilities* for each cluster:

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

$$\omega_{nk} = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\theta}_k, I)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\theta}_j, I)}. \quad (20)$$

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

$$\boldsymbol{\theta}_k^* = \frac{1}{N_k} \sum_{n=1}^K \omega_{nk} \mathbf{x}_n, \quad (21)$$

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

This is called the **expectation maximization (EM)** algorithm.

## 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.