

# lecture 6

Thanks for sitting in the  
first 6 rows

Participation  
attendance

Poll method 1:

TB posted

HW 2 Monday  
TB posted

MMP OH: PNL  
Mon 4-5 pm B321

## Lecture II – Clustering – Part II: Non-parametric clustering

Marina Meilă  
[mmp@stat.washington.edu](mailto:mmp@stat.washington.edu)

Department of Statistics  
University of Washington

CSE 547/STAT 548  
Winter 2022

## 1 Paradigms for clustering



## 2 Methods based on non-parametric density estimation ✓

$x_1, x_2, \dots$  not iid

### 3 Model-based: Dirichlet process mixture models

$$f \in \{N(\mu, \Sigma), \dots\}$$

$$\{f(\cdot; \theta), \theta \in \Theta\}$$

$k = 1 : K$  finite mixture (parametric)

$$f(x) = \sum_{k=1}^K \pi_k f(x; \theta_k)$$

K is r.v. in  $\{1, 2, \dots\}$

$$\sum_{k=1}^K \pi_k = 1 \quad \text{mixture proportions}$$

$\pi_k \geq 0 \text{ for all } k$

# The Dirichlet distribution

- $Z \in \{1 : r\}$  a discrete random variable, let  $\theta_j = P_z(j)$ ,  $j = 1, \dots, r$ .
- Multinomial distribution Probability of i.i.d. sample of size  $N$  from  $P_z$

$$\Pr \text{ on } S = \{1, \dots, r\}$$

$$\theta_j = \Pr_{z \sim P_z} [z=j]$$

$$j = 1 : r$$

$z_1, \dots, z_n \sim \text{iid } P$

$$P(z^{1, \dots, n}) = \prod_{j=1}^r \theta_j^{n_j}$$

$$n_j = \#\{z_i = j\}$$

- where  $n_j = \#$  the value  $j$  is observed,  $j = 1, \dots, r$
- $n_{1:r}$  are the **sufficient statistics** of the data.
  - The **Dirichlet distribution** is defined over domain of  $\theta_{1:r}$ , with **real** parameters  $N'_{1:r} > 0$  by

Bayesian

$$\text{Prior: } D(\theta_{1:r}; N'_{1:r}) = \frac{\Gamma(\sum_j n'_j)}{\prod_j \Gamma(n'_j)} \prod_j \theta_j^{n'_j - 1}$$

$$\text{where } \Gamma(p) = \int_0^\infty t^{p-1} e^{-t} dt.$$

$$\text{Posterior: } \Pr(\theta_{1:r} | \mathcal{D}) \propto \ell(\mathcal{D}) \cdot \text{Prior}$$

$$\propto \prod_j \theta_j^{n_j} \theta_j^{N'_j - 1} = \prod_{j=1}^r \theta_j^{n_j + N'_j - 1}$$

$$r=3, n=5$$

$$n_1=3$$

$$n_2=n_3=1$$

$$\mathcal{D}: z_{1:n} = 11312 \rightarrow n_1=n_2=1$$

counts

$$D(\theta, n) = \frac{\Gamma(n)}{\prod_{j=1}^r \Gamma(n_j)} \prod_{j=1}^r \theta_j^{n_j-1}$$

$\leftarrow$  distribution over Probability simplex

$$\Delta_r = \{ \theta ; \theta_j \geq 0 \text{ for all } j, \sum_{j=1}^r \theta_j = 1 \}$$

Ex  $r=2$     $n_1=3$     $n_2=2$     $\frac{n_1!}{n!} = \frac{4!}{2!1!} \theta_1^{3-1} \theta_2^{2-1}$

$$\int_{\Delta_r} \prod_{j=1}^r \theta_j^{n_j-1} d\theta_1 d\theta_2 \cdots d\theta_r = \frac{\prod_{j=1}^r \Gamma(n_j)}{\Gamma(n)}$$

$$n = \sum_{j=1}^r n_j$$

$$\Gamma(p) = \int_0^\infty t^{p-1} e^{-t} dt$$

$$\Gamma(p+1) = p! \Rightarrow \Gamma(5) = 4!$$

$$p = 1, 2, 3, \dots$$

$$\theta_1, \theta_2 \geq 0$$

$$\theta_1 + \theta_2 = 1$$

$r=3$

$$\theta_{1,2,3} \geq 0$$

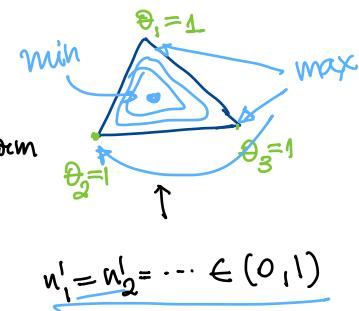
$$\theta_1 + \theta_2 + \theta_3 = 1$$

$$D(\cdot; n_1, n_2, n_3)$$

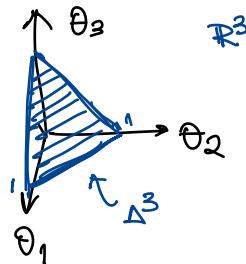


$$n_{1:r} = 1 \Rightarrow D = \text{uniform}$$

$$n_1 = n_2 = \dots > 1$$



$$n_1 = n_2 = \dots \in (0, 1)$$



$$\left\{ \begin{array}{l} \theta_1 + \theta_2 + \theta_3 = 1 \\ \theta_1, \theta_2, \theta_3 \geq 0 \end{array} \right.$$

# Dirichlet process mixtures

- Model-based
- generalization of mixture models to
  - infinite  $K$
  - Bayesian framework
- denote  $\theta_k$  = parameters for component  $f_k$
- assume  $f_k(x) \equiv f(x, \theta_k) \in \{f(x, \theta)\}$
- assume prior distributions for parameters  $g_0(\theta)$
- prior with hyperparameter  $\alpha > 0$  on the number of clusters
- very flexible model

# A sampling model for the data

## Dirichlet process

$\alpha = \text{hyperparam}$

- Example: Gaussian mixtures,  $d = 1$ ,  $\sigma_k = \sigma$  fixed
- $\theta = \mu$
- prior for  $\mu$  is  $\text{Normal}(0, \sigma_0^2 I_d)$
- Sampling process
  - for  $i = 1 : n$  sample  $x_i, k(i)$  as follows

denote  $\{1 : K\}$  the clusters after step  $i - 1$   
 define  $n_k$  the size of cluster  $k$  after step  $i - 1$

①

$$k(i) = \begin{cases} k & \text{w.p. } \frac{n_k}{i-1+\alpha} \\ K+1 & \text{w.p. } \frac{\alpha}{i-1+\alpha} \end{cases}, \quad k = 1 : K$$

$$\frac{\alpha}{i+\alpha-1}$$

- ② if  $k(i) = K + 1$  sample  $\mu_i \equiv \mu_{K+1}$  from  $\text{Normal}(0, \sigma_0^2)$   
 ③ sample  $x_i$  from  $\text{Normal}(\mu_{k(i)}, \sigma^2)$

- can be shown that the distribution of  $x_{1:n}$  is **interchangeable** (does not depend on data permutation)

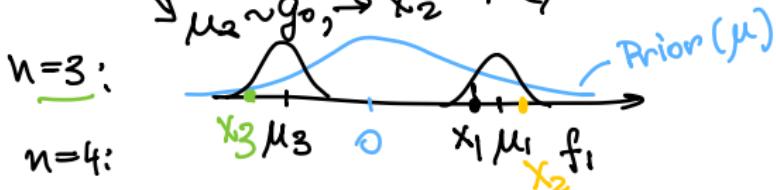
$$n=1: \mu_1 \sim g_0, x_1 \sim N(\mu_1, \sigma^2) \quad \mu_2 = \mu_1$$

$$n=2: \xrightarrow{x_2 \sim N(\mu_1, \sigma^2)} \mu_2 \sim g_0, \xrightarrow{x_2 \sim N(\mu_2, \sigma^2)}$$

$$f_e \in \{N(\mu, \sigma^2)\}$$

$$f_k = N(\mu_k, \sigma^2)$$

$$\mu \sim N(0, \sigma_0^2) = g_0$$



## A sampling model for the data

Chinese restaurant

at step  $i$ :

- Example: Gaussian mixtures,  $d = 1$ ,  $\sigma_k = \sigma$  fixed
- $\theta = \mu$
- prior for  $\mu$  is  $\text{Normal}(0, \sigma_0^2 I_d)$
- Sampling process

- for  $i = 1 : n$  sample  $x_i, k(i)$  as follows

denote  $\{1 : K\}$  the clusters after step  $i - 1$   
 define  $n_k$  the size of cluster  $k$  after step  $i - 1$

①

$$k(i) = \begin{cases} k & \text{w.p. } \frac{n_k}{i-1+\alpha} \\ K+1 & \text{w.p. } \frac{\alpha}{i-1+\alpha} \end{cases} \quad k = 1 : K \quad n(C_1) = n_1 = 2 \quad n(C_2) = 1 \quad (1)$$

② if  $k(i) = K + 1$  sample  $\mu_i \equiv \mu_{K+1}$  from  $\text{Normal}(0, \sigma_0^2)$ ③ sample  $x_i$  from  $\text{Normal}(\mu_{k(i)}, \sigma^2)$ 

- can be shown that the distribution of  $x_{1:n}$  is **interchangeable** (does not depend on data permutation)

•  $K \nearrow$  with  $n$  (slowly)

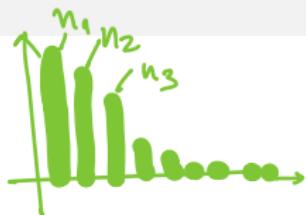
•  $\frac{\alpha}{n-1+\alpha} \rightarrow$  with  $n$

•  $K$  at  $n \nearrow$  with  $\alpha$

$$\sum n_k = i-1$$

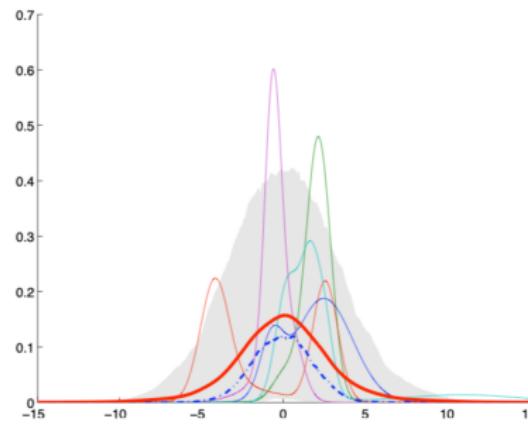
# The hyperparameters

- $\sigma_0$  controls spread of centers
  - should be large
- $\alpha$  controls number of cluster centers
  - $\alpha$  large  $\Rightarrow$  many clusters
- cluster sizes non-uniform (larger clusters attract more new points)
- many single point clusters possible



## General Dirichlet mixture model

- cluster densities  $\{f(x, \theta)\}$
- parameters  $\theta$  sampled from prior  $g_0(\theta, \beta)$
- cluster membership  $k(i)$  sampled as in (1)
- $x_i$  sampled from  $f(x, \theta_{k(i)})$
- **Model Hyperparameters  $\alpha, \beta$**



# Clustering with Dirichlet mixtures

← fitting model to data

## The clustering problem

- $\alpha, g_0, \beta, \{f\}$  given **model data**
- $\mathcal{D}$  given
- wanted  $\theta_{1:n}$  (not all distinct!)  $\mu_{1:k}$
- note:

- $\theta_{1:n}$  determines a hard clustering  $\Delta$   $l_k(i) = \text{cluster assignment for } x_i$
- the posterior of  $\theta_{1:n}$  given the data determines a soft clustering via  $P(x_i | k) \propto \int f(x_i | \theta_k) g_k(\theta_k) d\theta_k$

Estimating  $\theta_{1:n}$  cannot be solved in closed form

Usually solved by **MCMC (Markov Chain Monte Carlo)** sampling

Clustering with Dirichlet mixtures via MCMC (alternate  $k(i) \sim \theta_k \sim$ )

## MCMC estimation for Dirichlet mixture

**Input**  $\alpha, g_0, \beta, \{f\}, \mathcal{D}$ **State** cluster assignments  $k(i), i = 1 : n$ , ← random  
parameters  $\theta_k$  for all distinct  $k$  ← ( $\mu_i = x_i$ )**Iterate** ① for  $i = 1 : n$  (reassign data to clusters)

- ① remove  $i$  from its cluster (hence  $\sum_k n_k = n - 1$ )
- ② resample  $k(i)$  by

$$k(i) = \begin{cases} \text{existing } k & \text{w.p. } \propto \frac{n_k}{n-1+\alpha} f(x_i, \theta_k) \\ \text{new cluster} & \text{w.p. } \frac{\alpha}{n-1+\alpha} \int f(x_i, \theta) g_0(\theta) d\theta \end{cases} \quad (2)$$

- ③ if  $k(i)$  is new label, sample a new  $\theta_{k(i)} \propto g_0 f(x_i, \theta)$

- ② for  $k \in \{k(1 : n)\}$  (resample cluster parameters)

- ① sample  $\theta_k$  from posterior  $g_k(\theta) \propto g_0(\theta, \beta) \prod_{i \in C_k} f(x_i, \theta)$

$g_k$  can be computed in closed form if  $g_0$  is conjugate prior

**Output** a state with high posterior