

Clustering

Jiali Lin

Virginia Tech

January 15, 2017

Outline

Introduction

Finite mixture models

The Dirichlet process

Fitting Dirichlet processes to mixture modeling

Spectral clustering

Introduction

Similarity-based clustering

- ▶ Input: an $N \times N$ **dissimilarity matrix**.
- ▶ Output: **flat clustering**, where we partition the objects into disjoint sets.
- ▶ Sensitive to the initial conditions and requires some model selection method for K .

Feature-based clustering

- ▶ Input: an $N \times D$ feature matrix.
- ▶ Output: **hierarchical clustering**, where we create a nested tree of partitions.
- ▶ Most are deterministic and do not require the specification of K .

Clustering Evaluation: Rand Index

The validation of clustering structures is the most difficult.

- ▶ The number of assumed clusters may be different.
- ▶ No true cluster labels.

Rand index computes following for all pairs of data points

$$R = \frac{TP + TN}{TP + FP + FN + TN}$$

- ▶ **False positive (FP)**: target splits but algorithm clusters.
- ▶ **False negative (FN)**: target clusters but algorithm splits.
- ▶ **True positive (TP)**: algorithm and target both cluster together.
- ▶ **True negative (TN)**: algorithm and target both split apart.

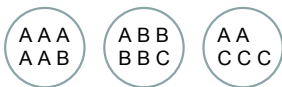


Figure: Circles are proposed clusters, letters are true cluster labels. Invariant to label choices, and takes time linear in N .

Outline

Introduction

Finite mixture models

The Dirichlet process

Fitting Dirichlet processes to mixture modeling

Spectral clustering

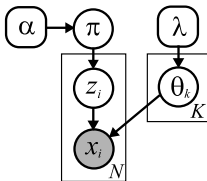
Finite mixture models

Traditional representation of a finite mixture model

$$p(\mathbf{x}_i | z_i = k, \boldsymbol{\theta}) = p(\mathbf{x}_i | \boldsymbol{\theta}_k)$$

$$p(z_i = k | \boldsymbol{\pi}) = \pi_k$$

$$p(\boldsymbol{\pi} | \alpha) = \text{Dir}(\boldsymbol{\pi} | (\alpha/K) \mathbf{1}_K)$$



- ▶ The form of $p(\boldsymbol{\theta}_k | \lambda)$ is chosen to be conjugate to $p(\mathbf{x}_i | \boldsymbol{\theta}_k)$.
- ▶ We can write $p(\mathbf{x}_i | \boldsymbol{\theta}_k)$ as $\mathbf{x}_i \sim F(\boldsymbol{\theta}_{z_i})$, where F is the observation distribution.
- ▶ We can write $\boldsymbol{\theta}_k \sim H(\lambda)$, where H is the prior.

Another representation

Consider

$$G(\boldsymbol{\theta}) = \sum_{k=1}^K \pi_k \delta_{\boldsymbol{\theta}_k}(\boldsymbol{\theta})$$

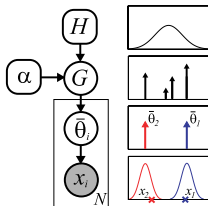


Figure: Here $\bar{\theta}_i$ is the parameter used to generate observation x_i ; these parameters are sampled from distribution G , which has the form.

- ▶ The discrete measure, G is a finite mixture of delta functions, K centered on the cluster parameters $\boldsymbol{\theta}_k$.
- ▶ The probability that $\bar{\theta}_i$ is equal to $\boldsymbol{\theta}_k$ is exactly π_k , the prior probability for that cluster.

Generative model

- Finite Gaussian mixture model ($K = 2$ clusters)

$$z_n \stackrel{iid}{\sim} \text{Categorical}(\rho_1, \rho_2)$$

$$\mathbf{x}_n \stackrel{indep}{\sim} N(\boldsymbol{\mu}_{z_n}, \boldsymbol{\Sigma})$$

- Don't know μ_1, μ_2

$$\boldsymbol{\mu}_k \stackrel{iid}{\sim} N(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$$

- Don't know ρ_1, ρ_2

$$\rho_1 \sim \text{Beta}(a_1, a_2)$$

$$\rho_2 = 1 - \rho_1$$

- **Inference goal:** assignments of data points to clusters, cluster parameters

Generative model

- Finite Gaussian mixture model (K clusters)

$$\boldsymbol{\rho}_{1:K} \sim \text{Dirichlet}(a_{1:K})$$

$$\boldsymbol{\mu}_k \stackrel{iid}{\sim} N(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$$

$$z_n \stackrel{iid}{\sim} \text{Categorical}(\boldsymbol{\rho}_{1:K})$$

$$\mathbf{x}_n \stackrel{indep}{\sim} N(\boldsymbol{\mu}_{z_n}, \boldsymbol{\Sigma})$$

What if $K \rightarrow \infty$

- ▶ Now, we will always (with probability one) get exactly K clusters.
- ▶ Want more flexible model: generate a variable number of clusters.
- ▶ The more data we generate, the more likely to see a new cluster.
- ▶ **Solutions:** replace the discrete distribution G with a **random probability measure**.
- ▶ Recall

$$G(\boldsymbol{\theta}) = \sum_{k=1}^K \pi_k \delta_{\boldsymbol{\theta}_k}(\boldsymbol{\theta})$$

$\boldsymbol{\theta}_i$ can take on the same value $\boldsymbol{\theta}_k$ for some k .

Outline

Introduction

Finite mixture models

The Dirichlet process

Fitting Dirichlet processes to mixture modeling

Spectral clustering

Dirichlet Distribution and Multinomial

- Consider

$$\begin{aligned}(\pi_1, \dots, \pi_K) &\sim \text{Discrete}(\alpha_1, \dots, \alpha_K) \\ z | (\pi_1, \dots, \pi_K) &\sim \text{Discrete}(\pi_1, \dots, \pi_K)\end{aligned}$$

- Then

$$\begin{aligned}z &\sim \text{Discrete}\left(\frac{\alpha_1}{\sum_i \alpha_i}, \dots, \frac{\alpha_K}{\sum_i \alpha_i}\right) \\ (\pi_1, \dots, \pi_K) | z &\sim \text{Discrete}(\alpha_1 + \delta_1(z), \dots, \alpha_K + \delta_K(z))\end{aligned}$$

where $\delta_i(z) = 1$ if z takes on value i , 0 otherwise.

Dirichlet process

- ▶ **Dirichlet process** is a distribution over probability measures $G : \Theta \rightarrow \mathbb{R}^+$.
- ▶ Require $G(\theta) \geq 0$ and $\int_{\theta} G(\theta) d\theta = 1$.
- ▶ For any finite partition (T_1, \dots, T_K) of Θ

$$(G(T_1), \dots, G(T_k)) \sim \text{Dir}(\alpha H(T_1), \dots, \alpha H(T_K))$$

- ▶ Define: $G \sim \text{DP}(\alpha, H)$, where α is called the **concentration parameter** and H is called the **base measure**.
- ▶ Intuitively, G needs to resemble with the basic distribution H .
- ▶ α determines how closely the histogram of spikes represents H .

Some properties of DP

- We are interested in

$$p(\theta) = \int p(\theta|G)p(G)dG$$
$$p(G|\theta) = \frac{p(\theta|G)p(G)}{p(\theta)}$$

- Recall **Dirichlet-multinomial conjugacy**.
- If $G \sim \text{DP}(\alpha, H)$, then $p(\theta_i \in T_i) = H(T_i)$ and the posterior is

$$(G(T_1), \dots, G(T_k)) | \theta \sim \text{Dir}(\alpha H(T_1) + \mathbb{I}(\theta \in T_1), \dots, \alpha H(T_K) + \mathbb{I}(\theta \in T_K))$$

- If we observe multiple samples $\theta_i \sim G$

$$G | \theta_{1:n}, \alpha, H \sim \text{DP}\left(\alpha + n, \frac{\alpha}{\alpha + n} H + \frac{1}{\alpha + n} \sum_{i=1}^n \delta_{\theta_i}\right)$$

Stick-breaking Construction

- ▶ Consider a partition $(\boldsymbol{\theta}, \mathbf{X} \setminus \boldsymbol{\theta})$ of \mathbf{X} .
- ▶ We know the the posterior process

$$\begin{aligned}(G(\boldsymbol{\theta}), G(\mathbf{X} \setminus \boldsymbol{\theta})) &\sim \text{Dir}((\alpha + 1) \frac{\alpha H + \delta_{\boldsymbol{\theta}}}{\alpha + 1}(\boldsymbol{\theta}), (\alpha + 1) \frac{\alpha H + \delta_{\boldsymbol{\theta}}}{\alpha + 1}(\mathbf{X} \setminus \boldsymbol{\theta})) \\ &= \text{Dir}(1, \alpha)\end{aligned}$$

- ▶ G has a point mass located at $\boldsymbol{\theta}$:

$$G = \beta \delta_{\boldsymbol{\theta}} + (1 - \beta) G' \quad \text{with} \quad \beta \sim \text{Beta}(1, \alpha)$$

- ▶ G' is the (renormalized) probability measure with the point mass removed.

Stick-breaking Construction (cont'd)

- Currently, we have

$$G \sim \text{DP}(\alpha + 1, \frac{\alpha H + \delta_{\theta}}{\alpha + 1})$$

$$G = \beta \delta_{\theta} + (1 - \beta) G'$$

$$\theta \sim H$$

$$\beta \sim \text{Beta}(1, \alpha)$$

- Consider a further partition $(\theta, T_1, \dots, T_K)$ of \mathbf{X}

$$\begin{aligned}(G(\theta), G(T_1), \dots, G(T_K)) &= (\beta, (1 - \beta)G'(T_1), \dots, (1 - \beta)G'(T_K)) \\ &\sim \text{Dir}(1, \alpha H(T_1), \dots, \alpha H(T_K))\end{aligned}$$

- The agglomerative/decimative property of Dirichlet implies

$$\begin{aligned}(G(T_1), \dots, G(T_k)) &\sim \text{Dir}(\alpha H(T_1), \dots, \alpha H(T_K)) \\ G' &\sim \text{DP}(\alpha, H)\end{aligned}$$

Stick-breaking Construction (cont'd)

- We have

$$G \sim \text{DP}(\alpha, H)$$

$$G = \beta_1 \delta_{\theta_1} + (1 - \beta_1) G_1$$

$$G = \beta_1 \delta_{\theta_1} + (1 - \beta_1)(\beta_2 \delta_{\theta_2} + (1 - \beta_2) G_2)$$

$$\vdots$$

$$G = \sum_{k=1}^{\infty} \pi_k \delta_{\theta_k}(\theta)$$

where

$$\beta_k \sim \text{Beta}(1, \alpha)$$

$$\pi_k = \beta_k \prod_{l=1}^{k-1} (1 - \beta_l) = \beta_k (1 - \sum_{l=1}^{k-1} \pi_l)$$

Stick-breaking Construction (cont'd)

- ▶ This is often denoted by

$$\pi \sim \text{GEM}(\alpha)$$

where $\pi \sim \text{GEM}(\alpha)$ and $\theta_k \sim H$.

- ▶ Samples from a DP are discrete with probability one.
- ▶ In other words, if you keep sampling it, you will get more and more repetitions of previously generated values.

Blackwell-MacQueen Urn Scheme

- ▶ Working with infinite dimensional sticks is problematic.
- ▶ We can exploit the clustering property to draw samples from a GP.
- ▶ **Key:** if $\theta_i \sim G$ are N observations from $G \sim \text{DP}(\alpha, H)$, taking on K distinct values θ_k , then the predictive distribution.

$$p(\bar{\theta}_{N+1} = \theta | \bar{\theta}_{1:N}, \alpha, H) = \frac{1}{\alpha + N} (\alpha H(\theta) + \sum_{k=1}^K N_k \delta_{\bar{\theta}_k}(\theta))$$

where N_k is the number of previous observations equal to θ_k .

- ▶ This is the **Polya urn** or **Blackwell-MacQueen** sampling scheme.
- ▶ The urn model can be equivalently expressed as

$$x_i | \theta_i \sim F(\theta_i)$$

$$\theta_i | G \sim G$$

$$G \sim \text{DP}(\alpha, H)$$

$$\theta_i | \theta_{1:i-1} \sim \frac{1}{i-1 + \alpha} \left(\sum_{j=1}^{i-1} N_j \delta_{\bar{\theta}_j}(\theta) + \alpha H \right)$$

The Chinese restaurant process (CRP)

- ▶ Let discrete variables z_i specify which value of θ_k to use.
- ▶ That is, we define $\bar{\theta}_i = \theta_{z_i}$.

$$p(z_{N+1} = z | z_{1:N}, \alpha) = \frac{1}{\alpha + N} (\alpha \mathbb{I}(z = k^*) + \sum_{k=1}^K N_k \mathbb{I}(z = k))$$

where k^* represents a new cluster index that has not yet been used.

- ▶ This is the **Chinese restaurant process** or **CRP**.
- ▶ The tables are like clusters, and the customers are like observations.
 - A person joins an existing table with probability $\frac{N_k}{\alpha + N}$.
 - He may choose to sit at a new table k^* with probability $1/(\alpha + N)$.
- ▶ The difference between the CRP and the Polya Urn Model is that the CRP specifies only a distribution over partitions (i.e., table assignments), but doesn't assign parameters to each group, whereas the Polya Urn Model does both.

Outline

Introduction

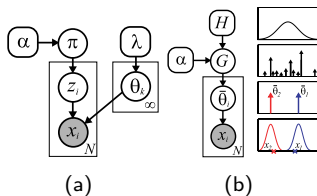
Finite mixture models

The Dirichlet process

Fitting Dirichlet processes to mixture modeling

Spectral clustering

Dirichlet Process Mixture



- The DP is not particularly useful as a model for data directly, since data vectors rarely repeat exactly.
- Useful as a prior for the parameters to generate data.
- Define $G \sim \text{DP}(\alpha, H)$. Equivalently, we can write the model

$$\pi \sim \text{GEM}(\alpha), \quad z_i \sim \pi$$

$$\theta_k \sim H(\lambda)$$

$$x_i \sim F(\theta_{z_i})$$

Gaussian mixture model

$$\boldsymbol{\pi} \sim \text{GEM}(\alpha), \quad z_i \sim \boldsymbol{\pi}$$

$$G = \sum_{k=1}^K \pi_k \delta_{\boldsymbol{\theta}_k}(\boldsymbol{\theta}) = \text{DP}(\alpha, N(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0))$$

$$\boldsymbol{\mu}_i \sim G$$

$$\boldsymbol{x}_i \sim N(\boldsymbol{\mu}_i, \boldsymbol{\Sigma})$$

Fitting a DP mixture modeling

- Fit a DPMM by modifying the collapsed Gibbs sampler.

$$p(z_i = k | \mathbf{z}_{-i}, \mathbf{x}, \alpha, \boldsymbol{\lambda}) \propto p(z_i = k | \mathbf{z}_{-i}, \alpha) p(\mathbf{x}_i | \mathbf{x}_{-i}, z = k, \mathbf{z}_{-i}, \boldsymbol{\lambda})$$

- The first term is given by

$$p(z_i | \mathbf{z}_{-i}, \alpha) = \frac{\mathbb{I}}{\alpha + N - 1} (\alpha \mathbb{I}(z = k^*) + \sum_{k=1}^K N_{k,-i} \mathbb{I}(z = k))$$

- If $z_i = k$, then \mathbf{x}_i is conditionally independent of all the data points except those assigned to cluster k . Hence

$$p(\mathbf{x}_i | \mathbf{x}_{-i}, \mathbf{z}_{-i}, z_i = k, \boldsymbol{\lambda}) = p(\mathbf{x}_i | \mathbf{x}_{-i,k}, \boldsymbol{\lambda}) = \frac{p(\mathbf{x}_i, \mathbf{x}_{-i,k} | \boldsymbol{\lambda})}{p(\mathbf{x}_{-i,k} | \boldsymbol{\lambda})}$$

where

$$p(\mathbf{x}_i, \mathbf{x}_{-i,k} | \boldsymbol{\lambda}) = \int p(\mathbf{x}_i | \boldsymbol{\theta}_k) \left[\prod_{j \neq i, z_j = k} p(\mathbf{x}_j | \boldsymbol{\theta}_k) \right] H(\boldsymbol{\theta}_k | \boldsymbol{\lambda}) d\boldsymbol{\theta}_k$$

Fitting a DP mixture modeling(cont'd)

- If $z_i = k^*$, corresponding to a new cluster, we have

$$p(\mathbf{x}_i | \mathbf{x}_{-i}, \mathbf{z}_{-i}, z_i = k^*, \boldsymbol{\lambda}) = p(\mathbf{x}_i | \boldsymbol{\lambda}) = \int p(\mathbf{x}_i | \boldsymbol{\theta}) H(\boldsymbol{\theta} | \boldsymbol{\lambda}) d\boldsymbol{\theta}$$

Initialize

for each $i = 1 : N$ in random order **do**

Remove \mathbf{x}_i 's sufficient statistics from old cluster z_i .

for each $k = 1 : K$ **do**

 Compute $p_k(\mathbf{x}_i) = p(\mathbf{x}_i | \mathbf{x}_{-i}(k))$.

 Set $N_{k,-i} = \dim(\mathbf{x}_{-i}(k))$.

 Compute $p(z_i = k | \mathbf{z}_{-i}, D) = \frac{N_{k,-i}}{\alpha + N - 1}$.

Compute $p_*(\mathbf{x}_i) = p(\mathbf{x}_i | \boldsymbol{\lambda})$.

Compute $p(z_i = | \mathbf{z}_{-i}, D) = \frac{\alpha}{\alpha + N - 1}$.

Normalize $p(z_i | \cdot)$.

Sample $z_i \sim p(z_i | \cdot)$.

Add \mathbf{x}_i 's sufficient statistics to new cluster z_i .

If any cluster is empty, remove it and decrease K .

Outline

Introduction

Finite mixture models

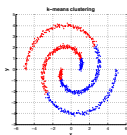
The Dirichlet process

Fitting Dirichlet processes to mixture modeling

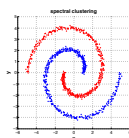
Spectral clustering

Data Clustering

- ▶ Two different criteria
 - Compactness, e.g., k-means, mixture models
 - Connectivity, e.g., spectral clustering



(c)



(d)

Figure: Clustering data consisting of 2 spirals. (a) K-means. (b) Spectral clustering. Figure generated by SpectralClusteringDemo, written by Wei-Lwun Lu.

Weighted Graph Partitioning

- Some graph terminology
 - Objects (e.g., pixels, data points) $i \in I$ = vertices of graph G .
 - Edges (ij) = pixel pairs with $W_{ij} > 0$.
 - Similarity matrix $\mathbf{W} = [W_{ij}]$.
 - Degree

$$d_i = \sum_{j \in G} W_{ij}$$

$$d_A = \sum_{i \in A} d_i \quad \text{degree of } A \subseteq G$$

- $\text{Assoc}(A, B) = \sum_{i \in A} \sum_{j \in B} W_{ij}$

Cuts in a Graph

- ▶ Edge cut = set of edges whose removal makes a graph disconnected.
- ▶ Weight of a cut:

$$\text{cut}(A, B) = \sum_{i \in A} \sum_{j \in B} W_{ij} = \text{Assoc}(A, B)$$

- ▶ Normalized Cut criteria: minimum $\text{cut}(A, \bar{A})$

$$\text{Ncut}(A, \bar{A}) = \frac{\text{cut}(A, \bar{A})}{d_A} + \frac{\text{cut}(A, \bar{A})}{d_{\bar{A}}}$$

Graph-based Clustering

- ▶ Affinity matrix: $\mathbf{W} = [W_{ij}]$.
- ▶ Degree matrix: $\mathbf{D} = \text{diag}(d_i)$.
- ▶ Laplacian matrix: $\mathbf{L} = \mathbf{D} - \mathbf{W}$.
- ▶ (Bipartite) partition vector:

$$\mathbf{x} = [x_1, \dots, x_n] = [1, 1, \dots, -1, \dots, -1]$$

Clustering via Optimizing Normalized Cut

- The normalized cut:

$$\text{Ncut}(A, B) = \frac{\text{cut}(A, B)}{d_A} + \frac{\text{cut}(A, B)}{d_B}$$

- Transform Ncut equation to a matrix form (Shi & Malik 2000):

$$\min_x \text{Ncut}(x) = \min_y \frac{y^T (\mathbf{D} - \mathbf{W}) y}{y^T \mathbf{D} y}$$

$$\text{subject to: } y \in \{1, -b\}^n$$

$$y^T \mathbf{D} \mathbf{1} = 0$$

- Relax the continuous domain by solving generalized eigenvalue system:

$$\min_y y^T (\mathbf{D} - \mathbf{W})y, \quad \text{s.t.} \quad y^T \mathbf{D}y = 1$$

- Which gives:

$$(\mathbf{D} - \mathbf{W})y = \lambda \mathbf{D}y$$

- Note that $(\mathbf{D} - \mathbf{W})\mathbf{1} = 0$ so, the first eigenvector is $y_0 = \mathbf{1}$ with eigenvalue 0.
- The second smallest eigenvector is the real valued solution to this problem.

Algorithm

- Define a similarity function between 2 nodes. i.e.

$$W_{i,j} = \exp\left(-\frac{\|X_{(i)} - X_{(j)}\|_2^2}{\sigma_X^2}\right)$$

- Compute affinity matrix (W) and degree matrix (D).
- Solve

$$(D - W)y = \lambda Dy$$

- Do singular value decomposition (SVD) of the graph Laplacian
 $L = D - W$.

$$L = V^T \Lambda V \Rightarrow y^*$$

- Use the eigenvector with the second smallest eigenvalue, bipartition the graph
 - For each threshold k ,

$$A_k = \{i | y_i \text{ among } k \text{ largest element of } y^*\}$$

$$B_k = \{i | y_i \text{ among } n - k \text{ smallest element of } y^*\}$$

- Compute $\text{Ncut}(A_k, B_k)$
- Output: $k^* = \operatorname{argmax} \text{Ncut}(A_k, B_k), A_{k^*}, B_{k^*}$