Department of Computer Science
UNIVERSITY OF COLORADO BOULDER

Bayesian Nonparametrics and DPMM

Machine Learning: Jordan Boyd-Graber
University of Colorado Boulder
LECTURE 17

**Clustering as Probabilistic Inference**

- GMM is a probabilistic model (unlike *K*-means)
- There are several latent variables:
  - Means
  - Assignments
  - (Variances)

- GMM is a probabilistic model (unlike *K*-means)
- There are several latent variables:
  - Means
  - Assignments
  - (Variances)
- Before, we were doing EM

**Clustering as Probabilistic Inference**

- GMM is a probabilistic model (unlike *K*-means)
- There are several latent variables:
  - Means
  - Assignments
  - (Variances)
- Before, we were doing EM
- Today, new models and new methods

**Nonparametric Clustering**

- What if the number of clusters is not fixed?
- Nonparametric: can grow if data need it
- Probabilistic distribution over number of clusters

**Dirichlet Process**

- Distribution over distributions
- Parameterized by: $\alpha$, $G$

**Dirichlet Process**

- Distribution over distributions
- Parameterized by: $\alpha$, $G$
- Concentration parameter

**Dirichlet Process**

- Distribution over distributions
- Parameterized by: $\alpha$, *G*
- Concentration parameter
- Base distribution

- Distribution over distributions
- Parameterized by: $\alpha$, $G$
- Concentration parameter
- Base distribution
- You can then draw observations from $x \sim \text{DP}(\alpha, G)$.

## Defining a DP

- Break off sticks

$$V_1, V_2, \ldots \sim_{\text{iid}} \text{Beta}(1, \alpha) \qquad \text{and} \qquad C_k := V_k \prod_{j=1}^{k-1} (1 - V_k)$$

## Defining a DP

- Break off sticks

$$V_1, V_2, \ldots \sim_{\mathrm{iid}} \mathrm{Beta}(1, \alpha) \qquad \text{and} \qquad C_k := V_k \prod_{j=1}^{k-1} (1 - V_k)$$

- Draw atoms

$$\Phi_1, \Phi_2, \ldots \sim_{\mathrm{iid}} G$$

**Defining a DP**

- Break off sticks

$$V_1, V_2, \ldots \sim_{\text{iid}} \text{Beta}(1, \alpha) \qquad \text{and} \qquad C_k := V_k \prod_{j=1}^{k-1}(1 - V_k)$$

- Draw atoms

$$\Phi_1, \Phi_2, \ldots \sim_{\text{iid}} G$$

- Merge into complete distribution

$$\Theta = \sum_{k \in \mathbb{N}} C_k \delta_{\Phi_k}$$

**Properties of a DPMM**

- Expected value is the same as base distribution

$$\mathbb{E}_{\mathsf{DP}(\alpha,G)}[x] = \mathbb{E}_G[x] \tag{1}$$

- As $\alpha \to \infty$, $\mathsf{DP}(\alpha, G) = G$
- Number of components unbounded
- Impossible to represent fully on computer (truncation)
- You can nest DPs

# DP as mixture Model

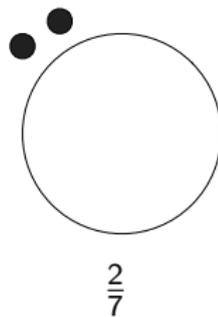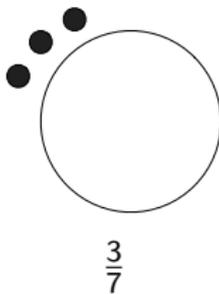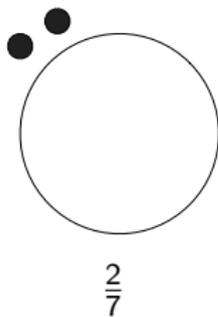To generate an observation, you first sit down at a table. You sit down at a table proportional to the number of people sitting at the table.

To generate an observation, you first sit down at a table. You sit down at a table proportional to the number of people sitting at the table.
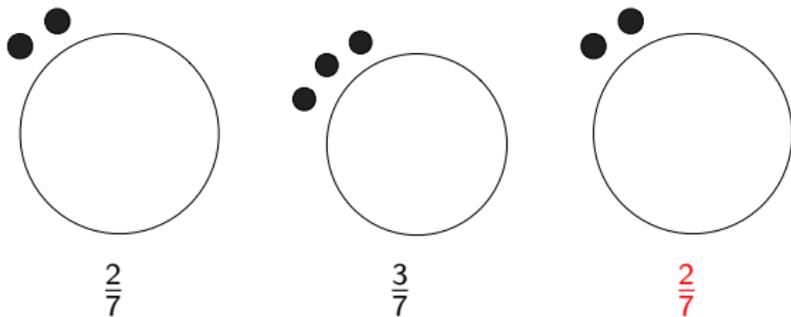


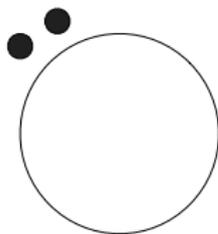$$\frac{2}{7} \qquad \frac{3}{7} \qquad \frac{2}{7}$$

To generate an observation, you first sit down at a table. You sit down at a table proportional to the number of people sitting at the table.



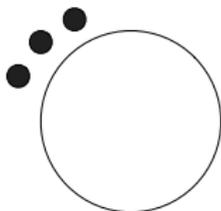$$\frac{2}{7} \qquad \frac{3}{7} \qquad \frac{2}{7}$$

To generate an observation, you first sit down at a table. You sit down at a table proportional to the number of people sitting at the table.
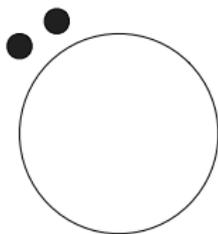


$$\frac{2}{7}$$
$$x \sim \mu_1$$
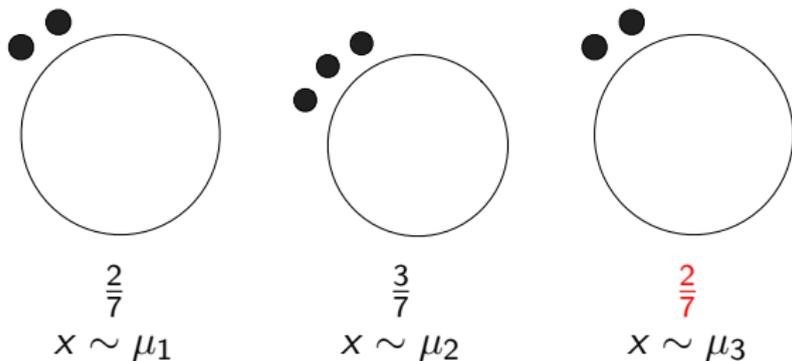
$$\frac{3}{7}$$
$$x \sim \mu_2$$

$$\frac{2}{7}$$
$$x \sim \mu_3$$

The Chinese Restaurant as a Distribution

To generate an observation, you first sit down at a table. You sit down at a table proportional to the number of people sitting at the table.



$$\frac{2}{7}$$
$$x \sim \mu_1$$

$$\frac{3}{7}$$
$$x \sim \mu_2$$

$$\frac{2}{7}$$
$$x \sim \mu_3$$

But this is just Maximum Likelihood

Why are we talking about Chinese Restaurants?

- The *posterior* of a DP is CRP
- A new observation has a new table / cluster with probability proportional to $\alpha$
- But this must be balanced against the probability of an observation *given a cluster*

$$\Theta = \sum_{k \in \mathbb{N}} C_k \delta_{\Phi_k}$$

**Gibbs Sampling**

- We want to know the cluster assignment of each observation
- Take a random guess initially

**Gibbs Sampling**

- We want to know the cluster assignment of each observation
- Take a random guess initially
- This provides a mean for each cluster

## Gibbs Sampling

- We want to know the cluster assignment of each observation
- Take a random guess initially
- This provides a mean for each cluster
- Let the number of clusters grow

## Gibbs Sampling

- We want to know the cluster assignment of each observation (tables)
- Take a random guess initially
- This provides a mean for each cluster
- Let the number of clusters grow

## Gibbs Sampling

- We want to know $\vec{z}$
- Compute $p(z_i \mid z_1 \ldots z_{i-1}, z_{i+1}, \ldots z_m, x, \alpha, G)$
- Update $z_i$ by sampling from that distribution
- Keep going . . .

**Gibbs Sampling**

- We want to know $\vec{z}$
- Compute $p(z_i \mid z_1 \ldots z_{i-1}, z_{i+1}, \ldots z_m, x, \alpha, G)$
- Update $z_i$ by sampling from that distribution
- Keep going . . .

Notation

$$p(z_i = k \mid z_{-i}) \equiv p(z_i \mid z_1 \ldots z_{i-1}, z_{i+1}, \ldots z_m) \tag{2}$$

$$p(z_i = k \,|\, \vec{z}_{-i}, \vec{x}, \{\theta_k\}, \alpha) \tag{3}$$

$$\tag{4}$$

$$p(z_i = k \mid \vec{z}_{-i}, \vec{x}, \{\theta_k\}, \alpha) \qquad (3)$$
$$= p(z_i = k \mid \vec{z}_{-i}, x_i, \vec{x}, \theta_k, \alpha) \qquad (4)$$
$$(5)$$

Dropping irrelevant terms

**Gibbs Sampling for DPMM**

$$p(z_i = k \,|\, \vec{z}_{-i}, \vec{x}, \{\theta_k\}, \alpha) \tag{3}$$

$$= p(z_i = k \,|\, \vec{z}_{-i}, x_i, \vec{x}, \theta_k, \alpha) \tag{4}$$

$$= p(z_i = k \,|\, \vec{z}_{-i}, \alpha) p(x_i \,|\, \theta_k, \vec{x}) \tag{5}$$

$$\tag{6}$$

Chain rule

**Gibbs Sampling for DPMM**

$$p(z_i = k \mid \vec{z}_{-i}, \vec{x}, \{\theta_k\}, \alpha) \tag{3}$$
$$= p(z_i = k \mid \vec{z}_{-i}, x_i, \vec{x}, \theta_k, \alpha) \tag{4}$$
$$= p(z_i = k \mid \vec{z}_{-i}, \alpha) p(x_i \mid \theta_k, \vec{x}) \tag{5}$$
$$= \begin{cases} \left(\frac{n_k}{n. + \alpha}\right) \int_\theta p(x_i \mid \theta) p(\theta \mid G, \vec{x}) & \text{existing} \\ \frac{\alpha}{n. + \alpha} \int_\theta p(x_i \mid \theta) p(\theta \mid G) & \text{new} \end{cases} \tag{6}$$
$$\tag{7}$$

Applying CRP

$$p(z_i = k \mid \vec{z}_{-i}, \vec{x}, \{\theta_k\}, \alpha) \tag{3}$$

$$= p(z_i = k \mid \vec{z}_{-i}, x_i, \vec{x}, \theta_k, \alpha) \tag{4}$$

$$= p(z_i = k \mid \vec{z}_{-i}, \alpha) p(x_i \mid \theta_k, \vec{x}) \tag{5}$$

$$= \begin{cases} \left(\frac{n_k}{n. + \alpha}\right) \int_\theta p(x_i \mid \theta) p(\theta \mid G, \vec{x}) & \text{existing} \\ \frac{\alpha}{n. + \alpha} \int_\theta p(x_i \mid \theta) p(\theta \mid G) & \text{new} \end{cases} \tag{6}$$

$$= \begin{cases} \left(\frac{n_k}{n. + \alpha}\right) \mathcal{N}\left(x, \frac{n\bar{x}}{n+1}, \mathbb{1}\right) & \text{existing} \\ \frac{\alpha}{n. + \alpha} \mathcal{N}\left(x, 0, \mathbb{1}\right) & \text{new} \end{cases} \tag{7}$$

Scary integrals assuming $G$ is normal distribution with mean zero and unit variance. (Derived in optional reading.)

1. Random initial assignment to clusters
2. For iteration $i$:
   2.1 "Unassign" observation $n$
   2.2 Choose new cluster for that observation
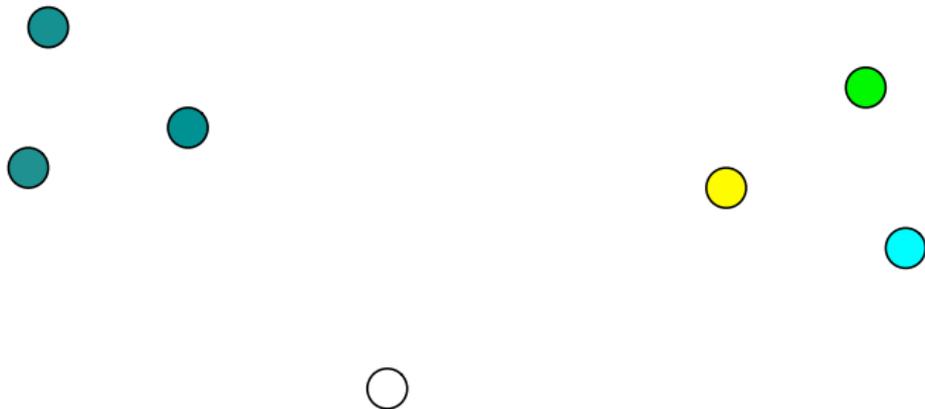
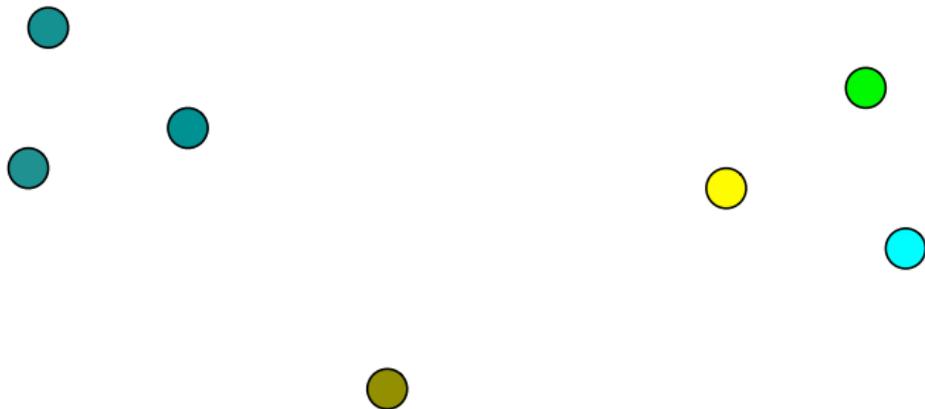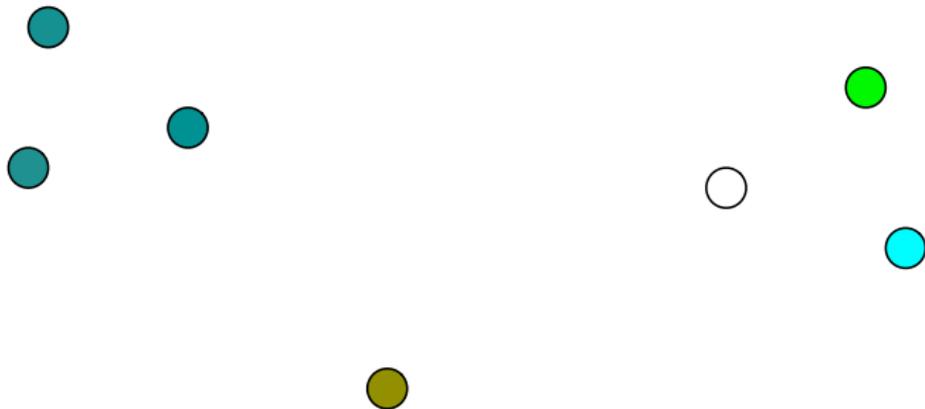# Toy Example

# Toy Example

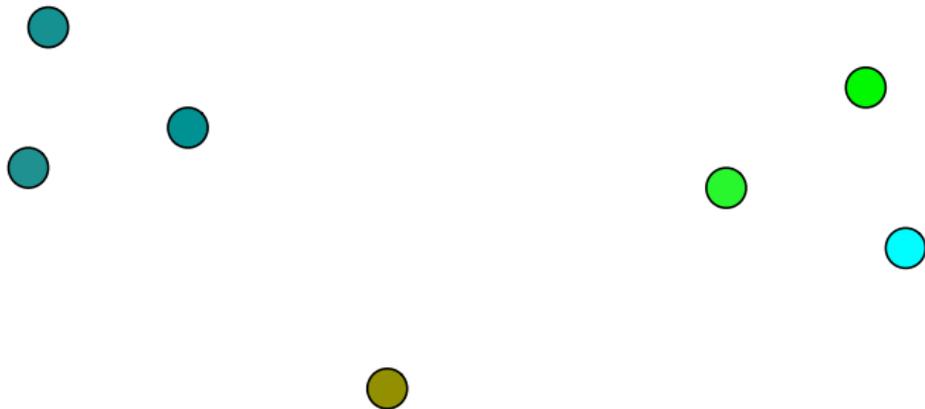# Toy Example
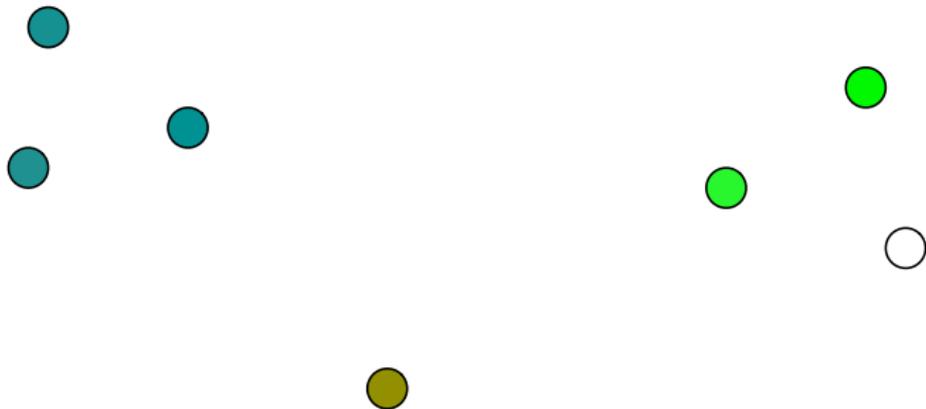
# Toy Example

# Toy Example
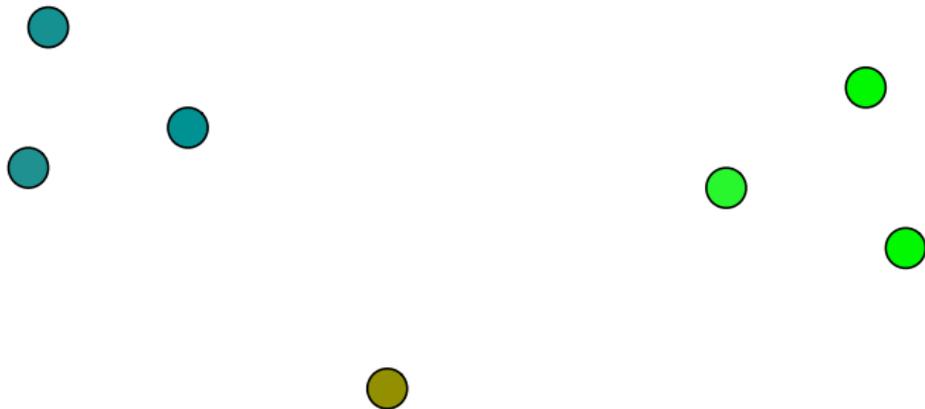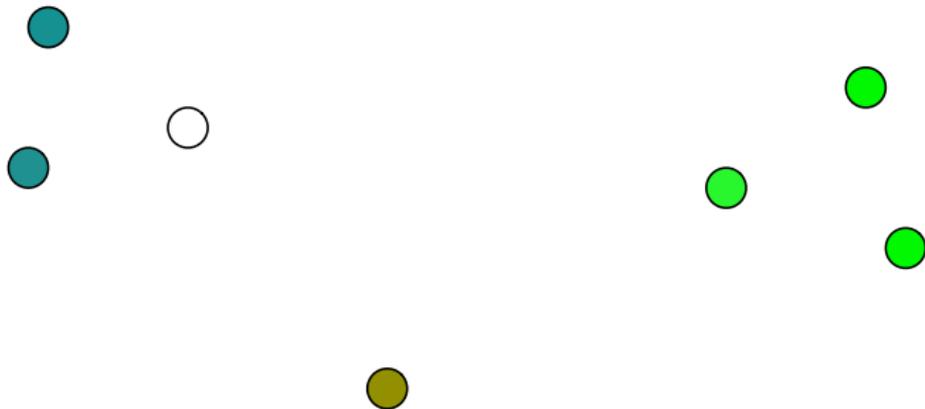
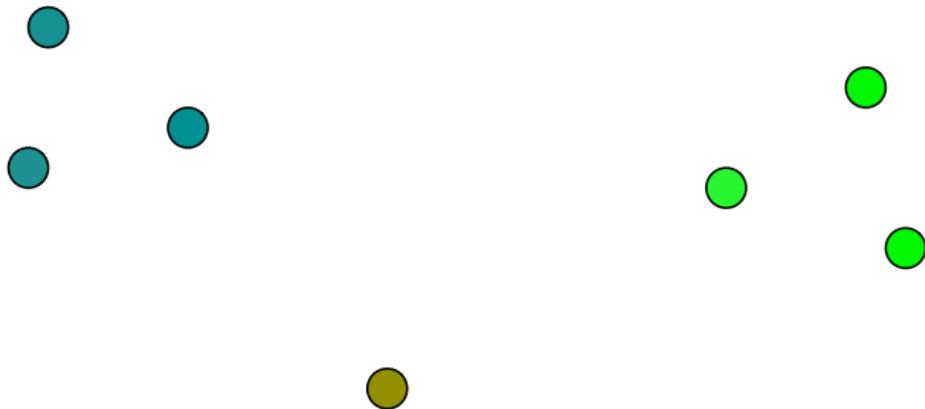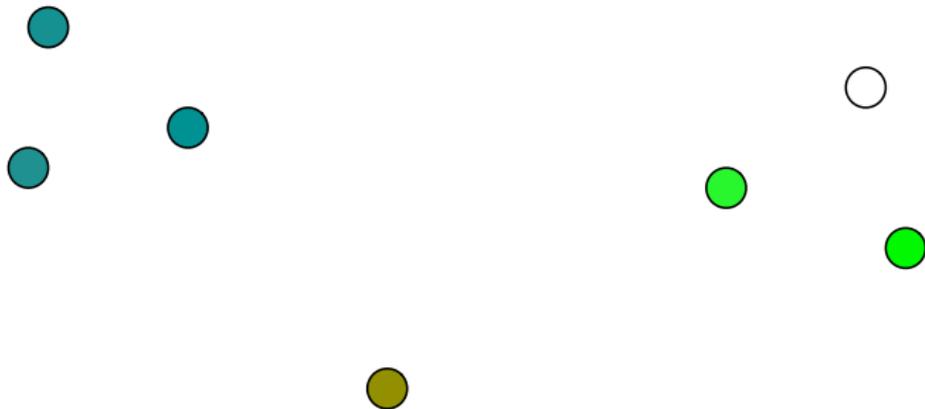**Toy Example**


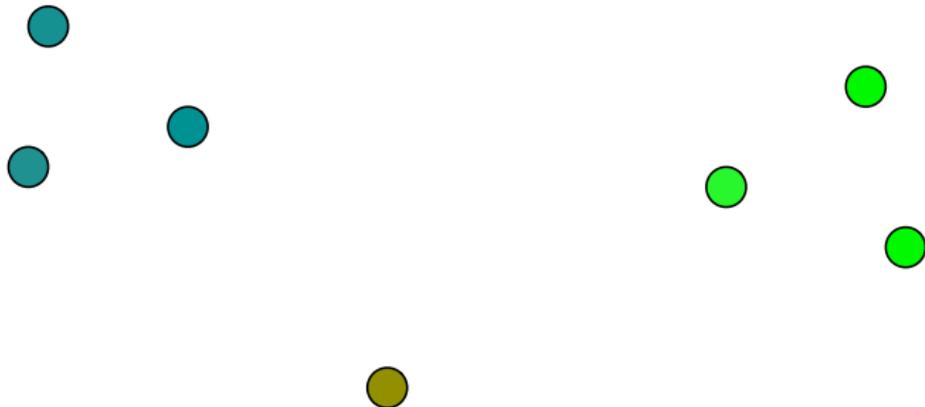
New cluster created!

# Toy Example

# Toy Example

# Toy Example

# Toy Example

And repeat . . .

- Gibbs often faster to implement
- EM easier to diagnose convergence
- EM can be parallelized
- Gibbs is more widely applicable

- Walking through DPMM clustering
- Clustering discrete data with more than one cluster per observation