**Module 3: Bayesian Nonparametrics**

# Infinite Mixture Models

1ˢᵗ finite MM recap

STAT/BIOSTAT 527, University of Washington

Emily Fox

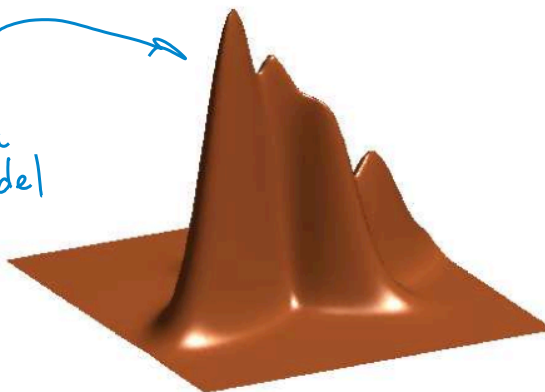May 7ᵗʰ, 2013

©Emily Fox 2013

---

# Density Estimation

- Estimate a density based on $x_1,…,x_N$

$x_1,…,x_N \sim P$

Let's consider a parametric model
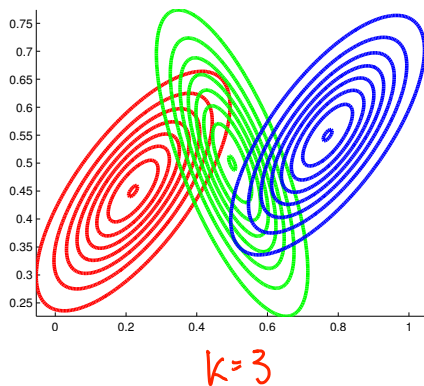
©Emily Fox 2013

1

# Density as Mixture of Gaussians

- Approximate density with a mixture of Gaussians
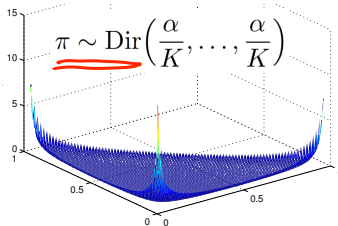
*Mixture of 3 Gaussians*

$$P = [\pi_1, \cdots, \pi_K]$$
$$\{_k \{\mu_k, \Sigma_k\}$$

$$p(x_i \mid \pi, \mu, \Sigma) =$$

Gauss. kernel, just like in KDE, but not centered at obs.

$$\sum_{k=1}^{K} \pi_k N(x_i \mid \mu_k, \Sigma_k)$$

# of mix comp.    mix. weights    shape params

k=3

In 1D:    P = target density

$$\sum \pi_k = 1$$

$\pi_2$   $\pi_1$   $\pi_3$
$\mu_2$   $\mu_1$   $\mu_3$

©Emily Fox 2013

---

# Model Summary

- Prior on model parameters
  - E.g., symmetric Dirichlet for $\pi$

$$\pi \sim \text{Dir}\left(\frac{\alpha}{K}, \ldots, \frac{\alpha}{K}\right)$$

  - Normal inverse Wishart prior for $\underline{\theta_k} = \{\mu_k, \Sigma_k\}$

$$\{1, \cdots, K\}$$

- Sample observations as

$$z_i \sim \pi \quad \text{choose a cluster}$$

$$x_i \mid z_i \sim N(\mu_{z_i}, \Sigma_{z_i}) \quad \text{sim obs. from selected Gauss.}$$
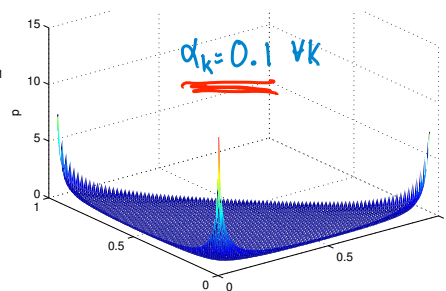
©Emily Fox 2013

# Dirichlet Distributions

- The Dirichlet distribution is defined on the simplex

$\pi \sim Dir(\alpha_1, \cdots, \alpha_K)$

$\Rightarrow \sum \pi_k = 1$

$$p(\pi \mid \alpha) = \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_{k=1}^{K} \pi_k^{\alpha_k - 1}$$

$\alpha_k = 10 \quad \forall k$

$\alpha_k = 0.1 \quad \forall k$

*Moments:* $\mathbb{E}_\alpha[\pi_k] = \dfrac{\alpha_k}{\alpha_0}$
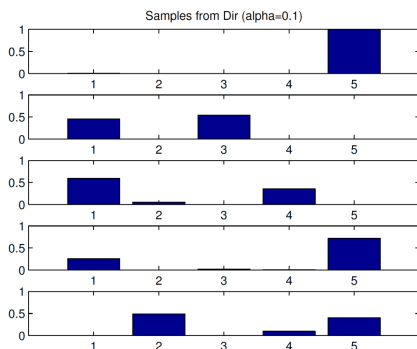
$\mathrm{Var}_\alpha[\pi_k] = \dfrac{K-1}{K^2(\alpha_0 + 1)}$
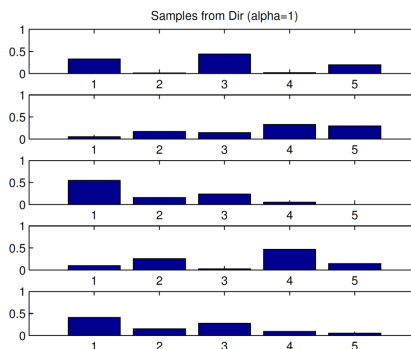
©Emily Fox 2013

---

# Dirichlet Samples

$\mathbb{E}_\alpha[\pi_k] = \dfrac{\alpha_k}{\alpha_0}$

- Samples are **sparse** for small values of $\alpha_i$

Samples from Dir (alpha=0.1)

Samples from Dir (alpha=1)

$\mathrm{Dir}(\pi \mid 0.1, 0.1, 0.1, 0.1, 0.1)$

puts mass @ corners

$\mathrm{Dir}(\pi \mid 1.0, 1.0, 1.0, 1.0, 1.0)$

uniform

©Emily Fox 2013

# Model In Pictures

■ Mixture weights

$\pi$

$z_2 = 2$

$z_1 = 3$

■ For each observation,

$$z_i \sim \pi$$
$$x_i \mid z_i \sim N(\mu_{z_i}, \Sigma_{z_i})$$

---

# GMM Sampler

■ Recall model

□ Observations: $x_1, \ldots, x_N$

want these { □ Cluster indicators: $z_1, \ldots, z_N$

□ Parameters: $\pi, \theta_k$ → $\pi = [\pi_1, \ldots, \pi_K]$
$$\theta_k = \{\mu_k, \Sigma_k\}$$

□ Generative model:

$$\pi \sim \mathrm{Dir}(\alpha_1, \ldots, \alpha_K) \qquad z_i \sim \pi$$
$$\{\mu_k, \Sigma_k\} \sim \mathrm{NIW}(\lambda) \qquad x_i \mid z_i, \{\theta_k\} \sim N(\mu_{z_i}, \Sigma_{z_i})$$

■ Iteratively sample

$$z_i \mid \pi, \{\theta_k\}, \{x_i\} \qquad i = 1, \ldots, N$$
$$\pi \mid \{z_i\}, \{\theta_k\}, \{x_i\}$$
$$\theta_k \mid \pi, \{z_i\}, \{x_i\} \qquad k = 1, \ldots K$$

# Standard Finite Mixture Sampler

Given mixture weights $\pi^{(t-1)}$ and cluster parameters $\{\theta_k^{(t-1)}\}_{k=1}^K$ from the previous iteration, sample a new set of mixture parameters as follows:

1. Independently assign each of the $N$ data points $x_i$ to one of the $K$ clusters by sampling the indicator variables $z = \{z_i\}_{i=1}^N$ from the following multinomial distributions:

$$z_i^{(t)} \sim \frac{1}{Z_i} \sum_{k=1}^K \pi_k^{(t-1)} f(x_i \mid \theta_k^{(t-1)}) \, \delta(z_i, k) \qquad Z_i = \sum_{k=1}^K \pi_k^{(t-1)} f(x_i \mid \theta_k^{(t-1)})$$

2. Sample new mixture weights according to the following Dirichlet distribution:

$$\pi^{(t)} \sim \text{Dir}(N_1 + \alpha/K, \ldots, N_K + \alpha/K) \qquad N_k = \sum_{i=1}^N \delta(z_i^{(t)}, k)$$

3. For each of the $K$ clusters, independently sample new parameters from the conditional distribution implied by those observations currently assigned to that cluster:

$$\theta_k^{(t)} \sim p(\theta_k \mid \{x_i \mid z_i^{(t)} = k\}, \lambda)$$
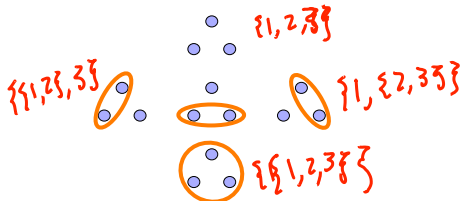
©Emily Fox 2013

# Mixtures Induce Partitions

- If our goal is clustering, the output grouping is defined by assignment *indicator variables:*
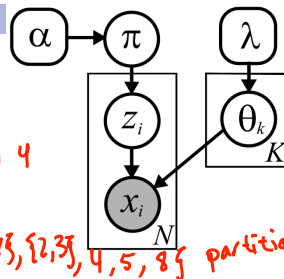
$$z_i \sim \pi$$

$z_1, \ldots, z_N$

1 5 5 3 2 1 1 4

$\{\{1,6,7\}, \{2,3\}, 4, 5, 8\}$ partition

- The number of ways of assigning $N$ data points to $K$ mixture components is $K^N$

- If $K \geq N$ this is much larger than the number of ways of partitioning that data:

$\{1, 2, 3\}$

$\{\{1,2\},3\}$

$\{1, \{2,3\}\}$

$\{\{1,3\}\}$

$\{\{1,2,3\}\}$

*N=3: 5 partitions versus* $3^3 = 27$

# Mixtures Induce Partitions

- If our goal is clustering, the output grouping is defined by assignment *indicator variables:*

$$z_i \sim \pi$$

- The number of ways of assigning *N* data points to *K* mixture components is $K^N$

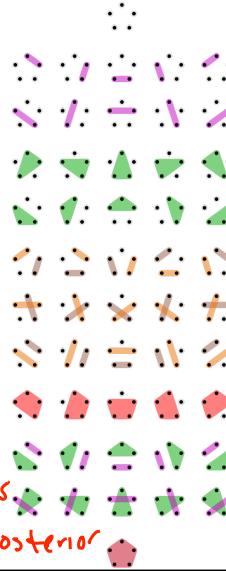- If $K \geq N$ this is much larger than the number of ways of partitioning that data:

*For any clustering, there is a unique partition, but many ways to label that partition's blocks.*

Note: sampler can switch between eq. partitions

*N=5: 52 partitions versus* $5^5 = 3125$

If prior sym., then eq. under posterior

*Courtesy Wikipedia*

---

## Module 3: Bayesian Nonparametrics

# Infinite Mixture Models
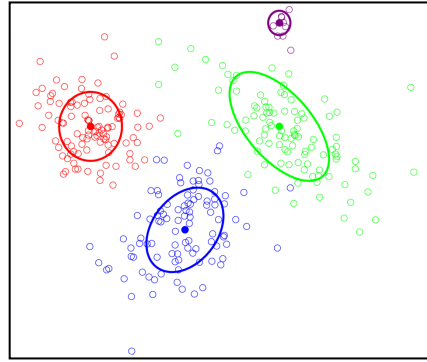
going infinite ...

STAT/BIOSTAT 527, University of Washington

Emily Fox

May 7th, 2013

# Motivating Nonparametric GMM

- What if current model doesn't fit new data?
- Bayesian nonparametric approach: $K \to \infty$
  - Allows infinite # clusters
  - Uses sparse subset
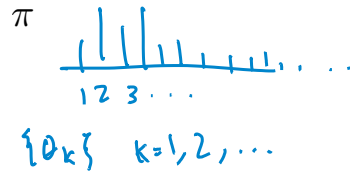  - Model **complexity adapts** to observations

Mixture of Gaussians

← allows us to add in new model comp.

$\theta_1$  $\theta_2$  $\theta_3$  $\theta_4$  $\theta_5$  $\theta_6$  $\theta_7$  $\cdots$

---

# Nonparam. Model In Pictures

- Mixture weights

$\pi$

1 2 3 $\cdots$
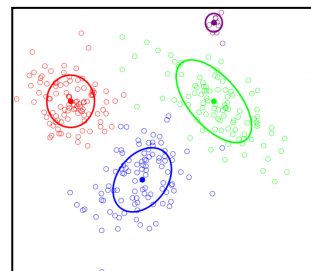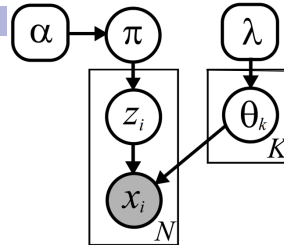
$\{\theta_k\}$  $k = 1, 2, \ldots$
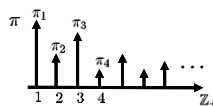
- For each observation, draw

$$z_i \sim \pi$$
$$x_i \mid z_i \sim N(\mu_{z_i}, \Sigma_{z_i})$$

How to define $\pi$?

7

# Stick-Breaking Process



- Start with stick of unit probability mass
- Repeatedly break portions of the remaining stick

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

$$\pi_1 = \beta_1$$

$$\pi_2 = \beta_2(1-\beta_1)$$

$$\vdots$$

$$\pi_k = \beta_k \prod_{\ell=1}^{k-1}(1-\beta_\ell)$$

$$\vdots$$

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

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

Stick of unit probability mass

---

# Stick-Breaking Process Summary



only a few w/ large weights

more w/ sig. weights

$\alpha = 1$    $\alpha = 5$

$$\pi_k = \beta_k \prod_{\ell=1}^{k-1}(1-\beta_\ell) = \beta_k\left(1-\sum_{\ell=1}^{k-1}\pi_\ell\right) \qquad \beta_k \sim \text{Beta}(1,\alpha)$$

$$\mathbb{E}[\beta_k] = \frac{1}{1+\alpha}$$

For large $\alpha$, breaking small portions

# Stick Breaks + Dirichlet Process



# Dirichlet Process Mixture Model

- Place Dirichlet process prior on weights and mixture parameters:

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

$$G = \sum_{k=1}^{\infty} \pi_k \delta_{\theta_k} \quad \begin{cases} \pi \sim \text{Stick}(\alpha) \\ \theta_k \overset{iid}{\sim} H \quad (\text{e.g. NIW}) \quad k=1,2,\ldots \end{cases}$$

- For each observation, draw

$$z_i \sim \pi$$

$$x_i \mid z_i \sim N(\mu_{z_i}, \Sigma_{z_i})$$

# Finite versus DP Mixtures

*Finite Mixture*    *DP Mixture*

$$\pi \sim \mathrm{Dir}\left(\frac{\alpha}{K}, \dots, \frac{\alpha}{K}\right) \qquad \pi \sim \mathrm{Stick}(\alpha)$$
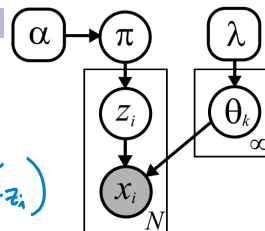
*sym.*

$$z_i \sim \pi$$
$$x_i \sim F(\theta_{z_i}) \qquad e.g. \ N(M_{z_i}, \Sigma_{z_i})$$

$\alpha \to \pi \qquad \lambda$

$z_i \qquad \theta_k$
$\infty$

$x_i$
$N$

**THEOREM:** For any measureable function *f*, as $K \to \infty$

$$\int_\Theta f(\theta)\, dG^K(\theta) \xrightarrow[K\to\infty]{\mathcal{D}} \int_\Theta f(\theta)\, dG(\theta)$$

$$G^K(\theta) = \sum_{k=1}^{K} \pi_k \delta_{\theta_k}(\theta) \qquad\qquad G \sim \mathrm{DP}(\alpha, H)$$

$$\pi^K \sim \mathrm{Dir}\left(\frac{\alpha}{K}, \dots, \frac{\alpha}{K}\right)$$

---

# Induced Partitions

- Recall that mixture models induce partitions of the data
$$z_i \sim \pi$$

- For a given prior on mixture weights, some partitions are more likely than others apriori

  □ Example 1: $\pi \sim \mathrm{Dir}(1, \dots, 1)$

  *uniform*

  $\pi$

  $1 \ 2 \ \cdots \ K$

  $\theta_1 \quad \theta_2$
  $\theta_5$

  *expect roughly the same # of obs. in each cluster*

  □ Example 2: $\pi \sim \mathrm{Dir}(0.01, \dots, 0.01)$

# Induced Partitions

- Recall that mixture models induce partitions of the data
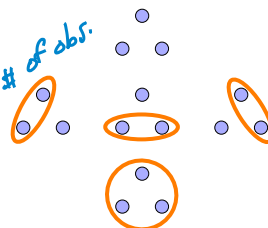$$z_i \sim \pi$$

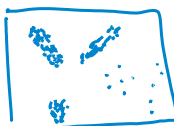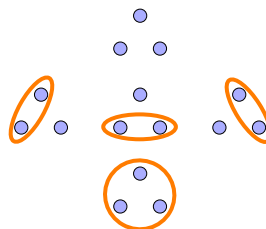- For a given prior on mixture weights, some partitions are more likely than others apriori
  - Example 3 (DP mix):  $\pi \sim \text{Stick}(\alpha)$

*[handwritten annotations: only a few w/ sig. mass, only many tiny weights]*

- What is the induced distribution on $z_1, \ldots, z_N$?
  - Do we expect many unique clusters?

*[handwritten: Answer: Chinese restaurant process]*

---

# Chinese Restaurant Process (CRP)

- Distribution on induced partitions described via the CRP
- Visualize clustering as a sequential process of customers sitting at tables in an (infinitely large) restaurant:

*customers*  ⟷  *observed data to be clustered*  $X_i$

*tables*  ⟷  *distinct clusters*  *[handwritten: each serving unique dish $\theta_k$]*

- The first customer sits at a table. Subsequent customers randomly select a table according to:  *[handwritten: # of unique clusters in $z_{1:N}$]*

$$p(z_{N+1} = z \mid z_1, \ldots, z_N, \alpha) = \frac{1}{\alpha + N}\left(\sum_{k=1}^{K} N_k \delta(z, k) + \alpha \delta(z, \bar{k})\right)$$

*[circles: ① $\frac{1}{\alpha+2}$  ② $\frac{1}{\alpha+2}$  ③  $\frac{\alpha}{\alpha+2}$ ]*

*[handwritten: # $z_i$ already at table k; weight on new cluster; new cluster index]*

# Chinese Restaurant Process (CRP)



Top row tables (After customer 7):
- Table with customers 1, 4, 5: $\frac{2}{7+\alpha}$
- Table with customers 2, 3, 7: $\frac{4}{7+\alpha}$
- Table with customer 6: $\frac{1}{7+\alpha}$
- Empty table: $\frac{\alpha}{7+\alpha}$   • • •   *After customer 7*

Middle row (cust. 8):
- Table with customers 1, 4, 5: $\frac{2}{8+\alpha}$
- Table with customers 2, 3, 5, 7, **8**: $\frac{5}{8+\alpha}$
- Table with customer 6: $\frac{1}{8+\alpha}$
- Empty table: $\frac{\alpha}{8+\alpha}$   • • •   *cust. 8*

Bottom row (cust. 9):
- Table with customers 1, 4, 5: $\frac{2}{9+\alpha}$
- Table with customers 2, 3, 5, 7, 8: $\frac{5}{9+\alpha}$
- Table with customer 6: $\frac{1}{9+\alpha}$
- Table with **9**: $\frac{1}{9+\alpha}$
- Empty table: $\frac{\alpha}{9+\alpha}$   • • •   *cust. 9*

*clustering /reinforcement induced by DP*
*"rich get richer"*

---

# CRPs & Exchangeable Partitions

$$p(z_{N+1} = z \mid z_1, \ldots, z_N, \alpha) = \frac{1}{\alpha + N}\left(\sum_{k=1}^{K} N_k \delta(z, k) + \alpha \delta(z, \bar{k})\right)$$
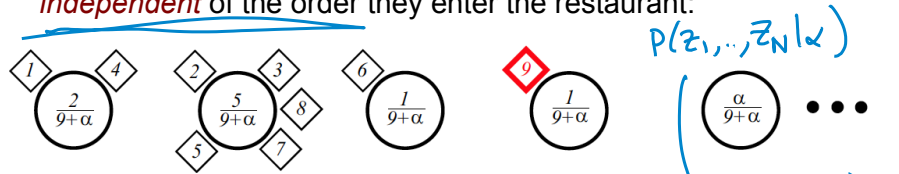
- The probability of a seating arrangement of *N* customers is *independent* of the order they enter the restaurant:



$P(z_1, \ldots, z_N \mid \alpha)$

  - Denominator terms:

$$1 \cdot \frac{1}{1+\alpha} \cdot \frac{1}{2+\alpha} \cdots \frac{1}{N-1+\alpha} = \frac{\Gamma(\alpha)}{\Gamma(N+\alpha)}$$

*2nd cust.   3rd cust.*

$P(z_1) P(z_2 \mid z_1) \cdots P(z_N \mid z_{1:N-1})$

# CRPs & Exchangeable Partitions

$$p(z_{N+1} = z \mid z_1, \ldots, z_N, \alpha) = \frac{1}{\alpha + N}\left(\sum_{k=1}^{K} N_k \delta(z, k) + \alpha \delta(z, \bar{k})\right)$$

- The probability of a seating arrangement of *N* customers is *independent* of the order they enter the restaurant:
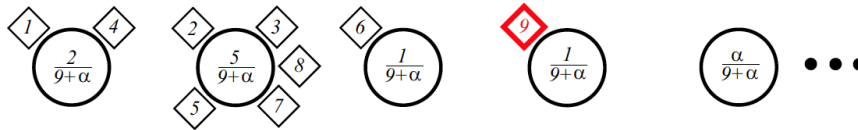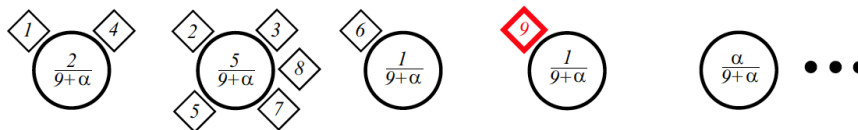


- Denominator terms: $\dfrac{1}{1+\alpha} \cdot \dfrac{1}{2+\alpha} \cdots \dfrac{1}{N-1+\alpha} = \dfrac{\Gamma(\alpha)}{\Gamma(N+\alpha)}$

- Number of new tables: $K$
  Numerator term for each new table: $\alpha$
  Combined: $\alpha^K$

---

# CRPs & Exchangeable Partitions

$$p(z_{N+1} = z \mid z_1, \ldots, z_N, \alpha) = \frac{1}{\alpha + N}\left(\sum_{k=1}^{K} N_k \delta(z, k) + \alpha \delta(z, \bar{k})\right)$$

- The probability of a seating arrangement of *N* customers is *independent* of the order they enter the restaurant:
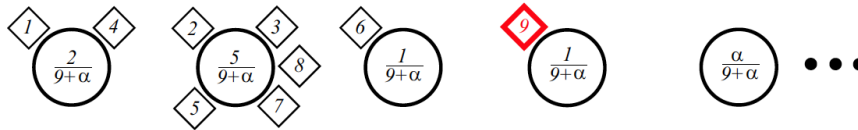


- Denominator terms: $\dfrac{1}{1+\alpha} \cdot \dfrac{1}{2+\alpha} \cdots \dfrac{1}{N-1+\alpha} = \dfrac{\Gamma(\alpha)}{\Gamma(N+\alpha)}$

- New table numerator terms: $\alpha^K$
- Customers joining $k$th occupied table: $1 \cdot 2 \cdots (N_k - 1) = (N_k - 1)! = \Gamma(N_k)$
  ↑ already 1 person sitting @ $k$th table

# CRPs & Exchangeable Partitions

$$p(z_{N+1} = z \mid z_1, \ldots, z_N, \alpha) = \frac{1}{\alpha + N}\left(\sum_{k=1}^{K} N_k \delta(z, k) + \alpha \delta(z, \bar{k})\right)$$

- The probability of a seating arrangement of *N* customers is *independent* of the order they enter the restaurant:
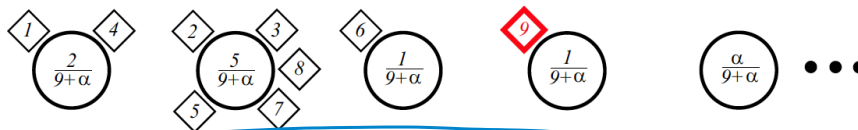


  □ Denominator terms: $\dfrac{1}{1+\alpha} \cdot \dfrac{1}{2+\alpha} \cdots \dfrac{1}{N-1+\alpha} = \dfrac{\Gamma(\alpha)}{\Gamma(N+\alpha)}$

  □ New table numerator terms: $\alpha^K$

  □ Customers joining $k^{\text{th}}$ occupied table:
    $$1 \cdot 2 \cdots (N_k - 1) = (N_k - 1)! = \Gamma(N_k)$$

---

# CRPs & Exchangeable Partitions

$$p(z_{N+1} = z \mid z_1, \ldots, z_N, \alpha) = \frac{1}{\alpha + N}\left(\sum_{k=1}^{K} N_k \delta(z, k) + \alpha \delta(z, \bar{k})\right)$$
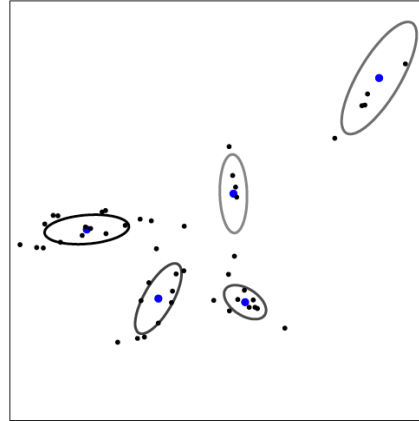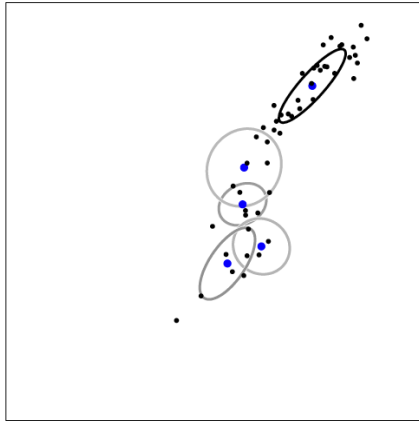
- The probability of a seating arrangement of *N* customers is *independent* of the order they enter the restaurant:



$$p(z_1, \ldots, z_N \mid \alpha) = \frac{\Gamma(\alpha)}{\Gamma(N+\alpha)} \alpha^K \prod_{k=1}^{K} \Gamma(N_k)$$
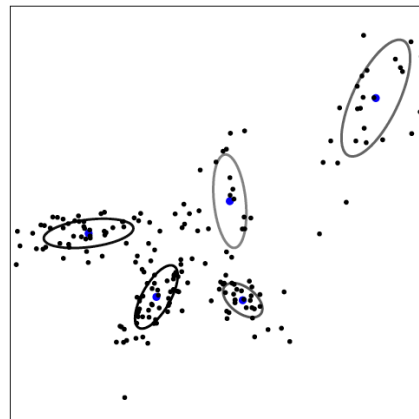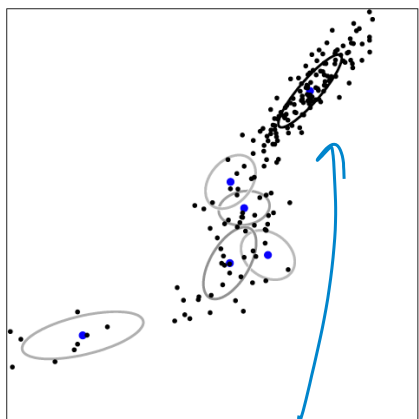
- Thus, the CRP is a prior on an *infinitely exchangeable* sequence
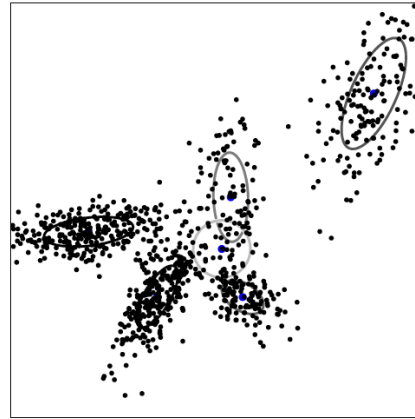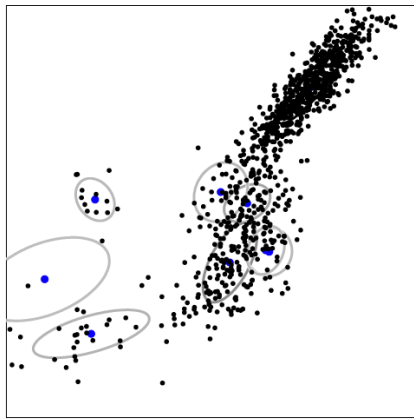
# Samples from DP Mixture Priors
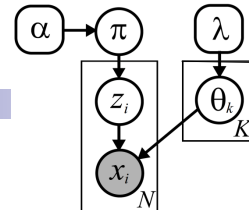


*N=50*

# Samples from DP Mixture Priors



*"rich get richer"*　　*N=200*

# Samples from DP Mixture Priors



*N=1000*

# Finite GMM Sampler



- Recall model
  - Observations: $x_1, \ldots, x_N$
  - Cluster indicators: $z_1, \ldots, z_N$  ← want these
  - Parameters: $\pi, \theta_k$ → $\pi = [\pi_1, \ldots, \pi_K]$
    $$\theta_k = \{\mu_k, \Sigma_k\}$$
  - Generative model:
    $$\pi \sim \text{Dir}(\alpha_1, \ldots, \alpha_K) \qquad z_i \sim \pi$$
    $$\{\mu_k, \Sigma_k\} \sim \text{NIW}(\lambda) \qquad x_i \mid z_i, \{\theta_k\} \sim N(\mu_{z_i}, \Sigma_{z_i})$$

- Iteratively sample

$$z_i \mid \pi, \{\theta_k\}, \{x_i\} \qquad i = 1, \ldots, N$$
$$\pi \mid \{z_i\}, \{\theta_k\}, \{x_i\}$$
$$\theta_k \mid \pi, \{z_i\}, \{x_i\} \qquad k = 1, \ldots K$$

# Collapsed DP Mixture Sampler

- Can't sample $\pi$ directly
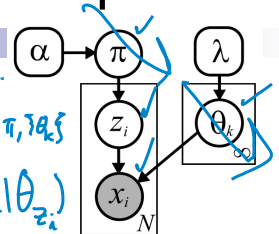- Integrate out all infinite-dimensional params $\pi, \{\theta_k\}$

$$P(z_{1:N}, x_{1:N}) = \int_\pi \int_{\theta_1, \theta_2, \ldots} p(\pi|\alpha) \prod_{k=1}^\infty p(\theta_k) \prod_{n=1}^N p(z_i|\pi) p(x_i|\theta_{z_i})$$

$$= P(z_1, \ldots, z_N | \alpha) \prod_{k=1}^\infty p(\{x_i : z_i = k\} | \lambda)$$

CRP      "likelihood"    "likelihood"

- Iteratively sample the cluster indicators

treat as last obs. in CRP

$$z_i^{(t)} \sim p\left(z_i = k \mid z_{\backslash i}^{(t-1)}, \alpha\right) p\left(x_i \mid \{x_j : z_j = k, i \neq j\}\right)$$

"prior"    all other indicators      all other obs. assigned to $k^{th}$ cluster

---

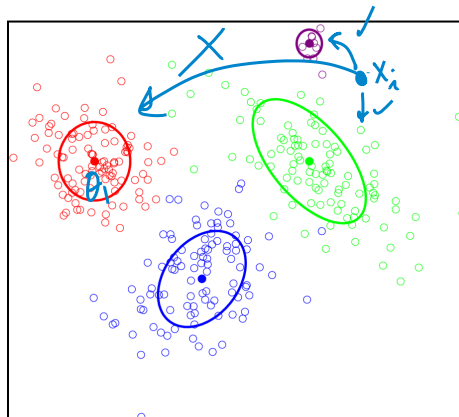# Collapsed Sampler Intuition

- Previously, $p(z_i = k \mid x_i, \pi, \theta) \propto \pi_k p(x_i \mid \theta_k)$

- If you're not told $\pi, \theta_k$

Approx $\pi$ by CRP
→ "prior" is based on cluster occupancy

Approx $\theta_k$
→ "likelihood" based on obs. already assigned to cluster

17

# Predictive Likelihood Term

- Recall NIW prior…Let's consider 1D example → N-IG

$$\mu_k \mid \sigma_k^2 \sim N(0, \gamma \sigma_k^2) \quad \sigma_k^2 \sim \text{IG}\left(\frac{\nu_0}{2}, \frac{\nu_0 S_0}{2}\right)$$

- Normal inverse gamma posterior
  → Student t predictive likelihood

*normal likelihood*
*NIG posterior*

$$p(x_i \mid \{x_j : z_j = k, j \neq i\}) = \int p(x_i \mid \theta_k) p(\theta_k \mid \{x_j : z_j = k, j \neq i\}) \, d\theta_k$$

$$p(x \mid \{x_j \mid z_j = k, j \neq i\}) = t_{\nu_0 + N_k^{-i}}\left(\frac{1}{\gamma + N_k^{-i}} \sum_{j: z_j = k, j \neq i} x_j,\right.$$

$$\left. \frac{N_k^{-i} + \gamma^{-1} + 1}{(N_k^{-i} + \gamma^{-1})(\nu_0 + N_k^{-i})}\left(\nu_0 S_0 + \sum_{j: z_j = k, j \neq i} x_j^2 - (N_k + \gamma^{-1})^{-1}(\sum_{j: z_j = k, j \neq i} x_j)^2\right)\right)$$

- ☐ Conjugacy: This integral is **tractable**

# Collapsed DP Mixture Sampler

1. Sample a random permutation $\tau(\cdot)$ of the integers $\{1, \ldots, N\}$.
2. Set $\alpha = \alpha^{(t-1)}$ and $z = z^{(t-1)}$. For each $i \in \{\tau(1), \ldots, \tau(N)\}$, resample $z_i$ as follows:

   (a) For each of the $K$ existing clusters, determine the predictive likelihood
   $$f_k(x_i) = p(x_i \mid \{x_j \mid z_j = k, j \neq i\}, \lambda)$$
   Also determine the likelihood $f_{\bar{k}}(x_i)$ of a potential new cluster $\bar{k}$
   $$p(x_i \mid \lambda) = \int_\Theta f(x_i \mid \theta) \, h(\theta \mid \lambda) \, d\theta$$

   (b) Sample a new cluster assignment $z_i$ from the following $(K+1)$–dim. multinomial:
   $$z_i \sim \frac{1}{Z_i}\left(\alpha f_{\bar{k}}(x_i)\delta(z_i, \bar{k}) + \sum_{k=1}^K N_k^{-i} f_k(x_i)\delta(z_i, k)\right) \qquad Z_i = \alpha f_{\bar{k}}(x_i) + \sum_{k=1}^K N_k^{-i} f_k(x_i)$$
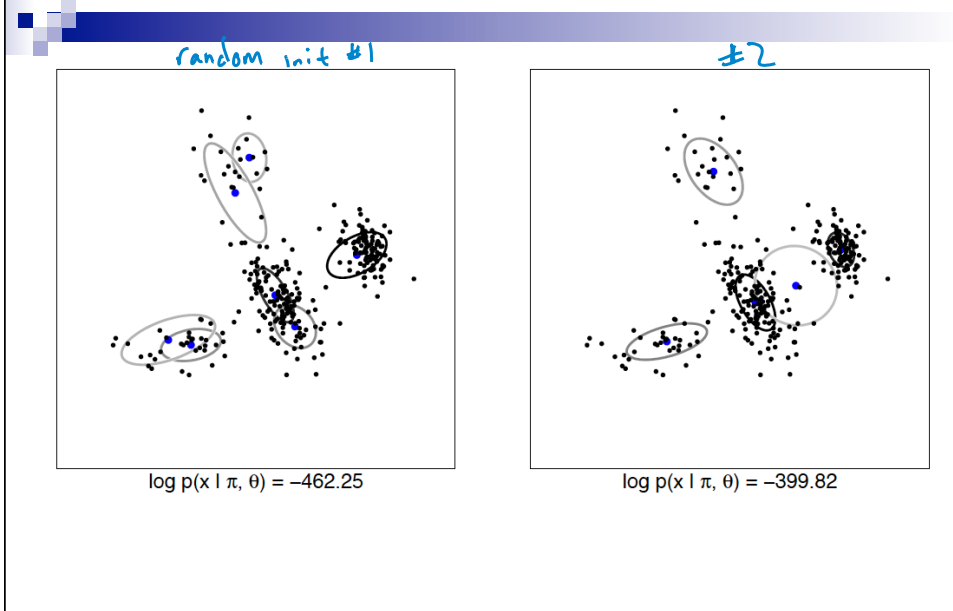   $N_k^{-i}$ is the number of other observations currently assigned to cluster $k$.

   (c) Update cached sufficient statistics to reflect the assignment of $x_i$ to cluster $z_i$. If $z_i = \bar{k}$, create a new cluster and increment $K$.
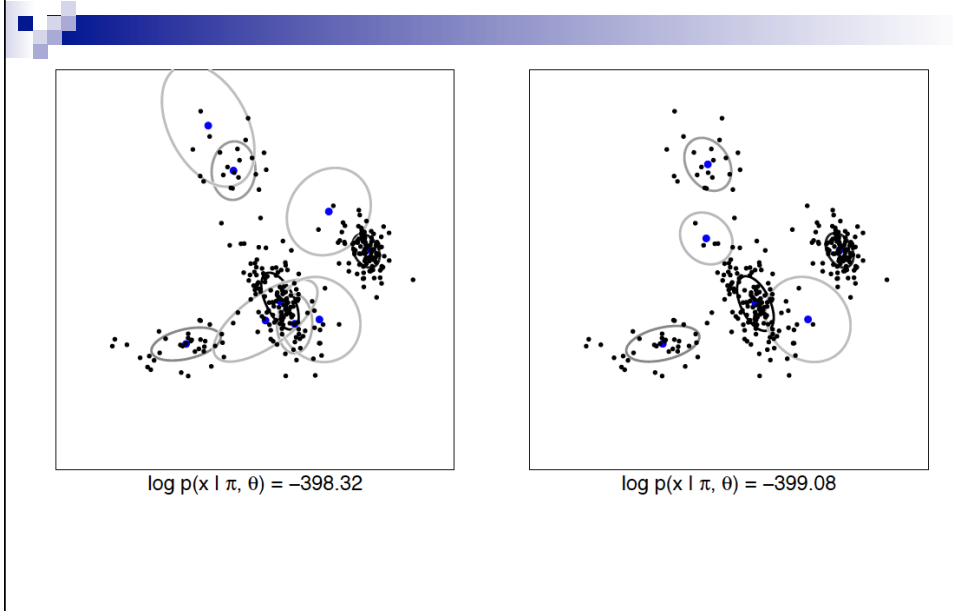
3. Set $z^{(t)} = z$.
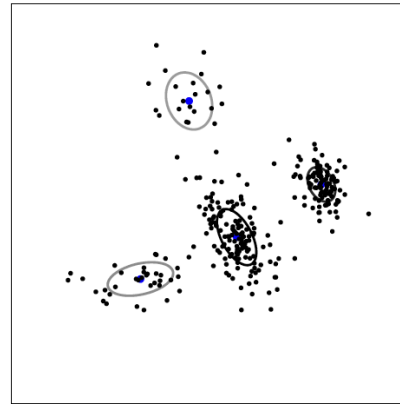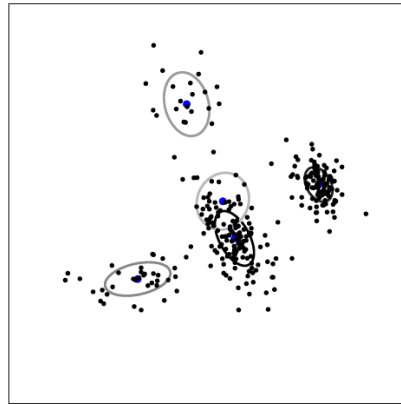4. If any current clusters are empty ($N_k = 0$), remove them and decrement $K$ accordingly.
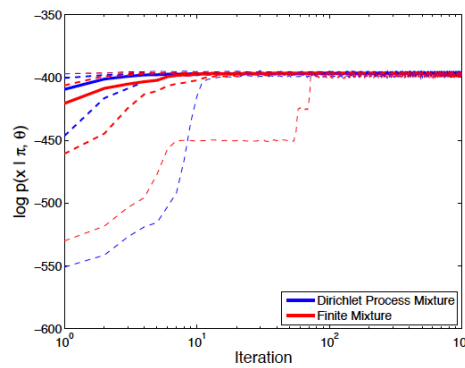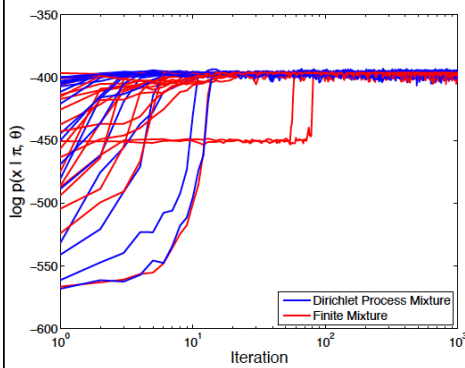
# Collapsed DP Sampler: 2 Iterations



random init #1

#2

log p(x | π, θ) = −462.25

log p(x | π, θ) = −399.82

# Collapsed DP Sampler: 10 Iterations



log p(x | π, θ) = −398.32

log p(x | π, θ) = −399.08

# Collapsed DP Sampler: 50 Iterations



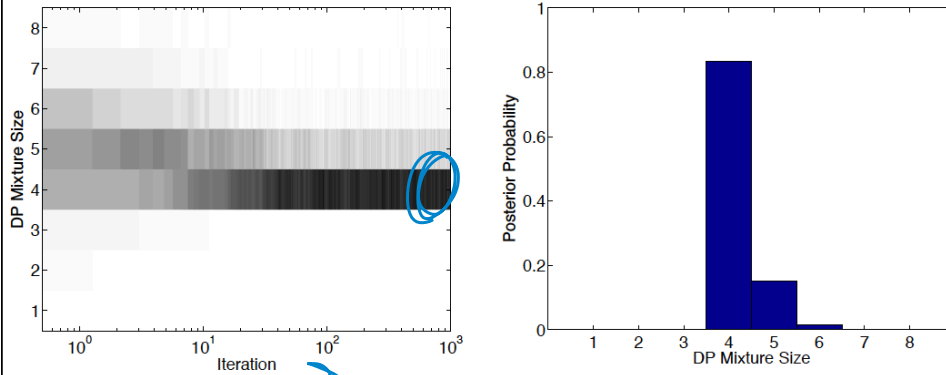log p(x | π, θ) = −397.67          log p(x | π, θ) = −396.71
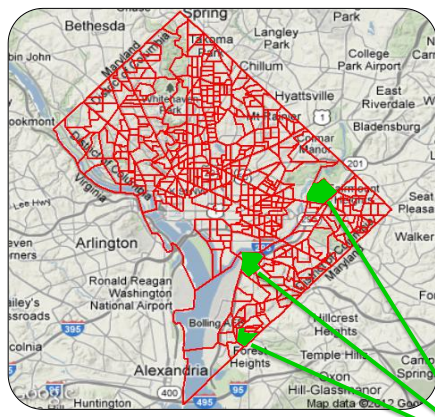
# DP vs. Finite Mixture Samplers



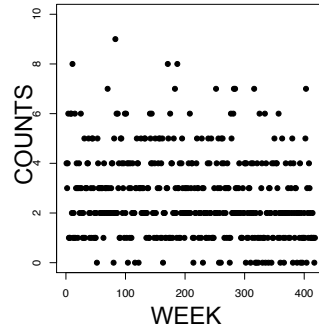pretty competitive

# DP Posterior Number of Clusters



Asy. consistency results for density est.
(assuming light tails on target density)
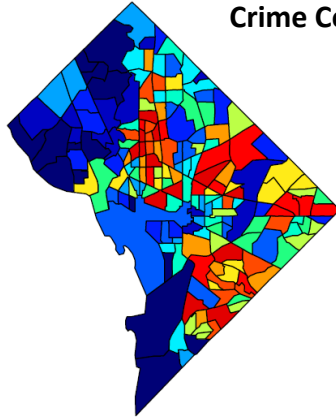Not asy. consistent on # of comp.

# DC Violent Crime Data



- 188 census tracts
- Weekly crime counts from 2001-2008
- Violent crime types:
  - ADW, arson, robbery, rape

Time series = crime counts

# DC Violent Crime Data
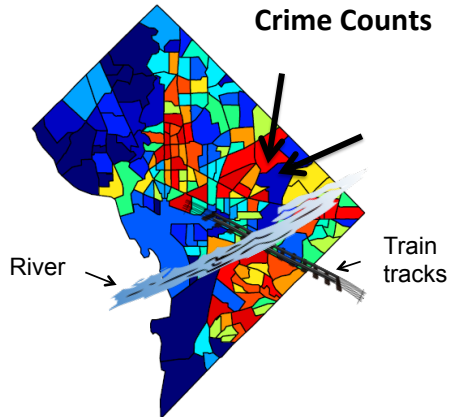


**Average Weekly Crime Counts**

Average Crime Count
- (1.665,2.641]
- (1.381,1.665]
- (1.155,1.381]
- (1.032,1.155]
- (0.9077,1.032]
- (0.8223,0.9077]
- (0.7368,0.8223]
- (0.6203,0.7368]
- (0.5085,0.6203]
- (0.4087,0.5085]
- (0.3307,0.4087]
- (0.229,0.3307]
- (0.07365,0.229]
- (0.009569,0.07365]

Goal: Forecast next week's map

---

# DC Violent Crime Data



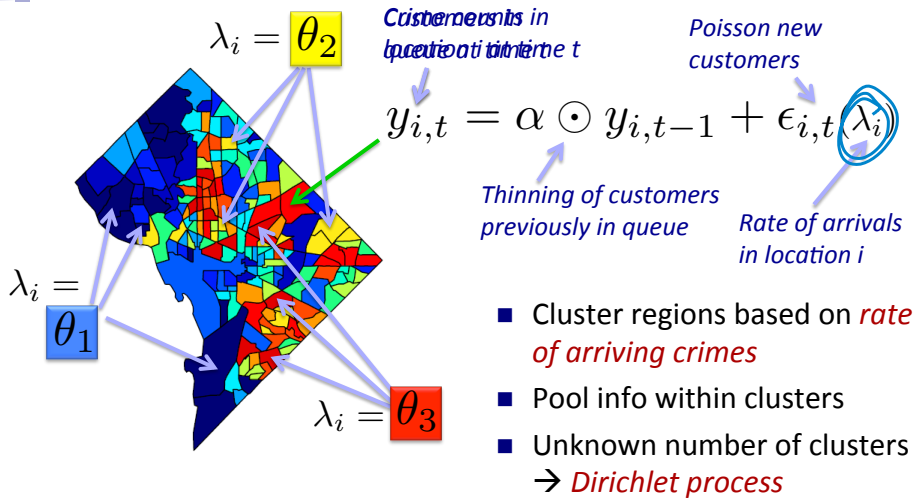**Average Weekly Crime Counts**

Average Crime Count
- (1.665,2.641]
- (1.381,1.665]
- (1.155,1.381]
- (1.032,1.155]
- (0.9077,1.032]
- (0.8223,0.9077]
- (0.7368,0.8223]
- (0.6203,0.7368]
- (0.5085,0.6203]
- (0.4087,0.5085]
- (0.3307,0.4087]
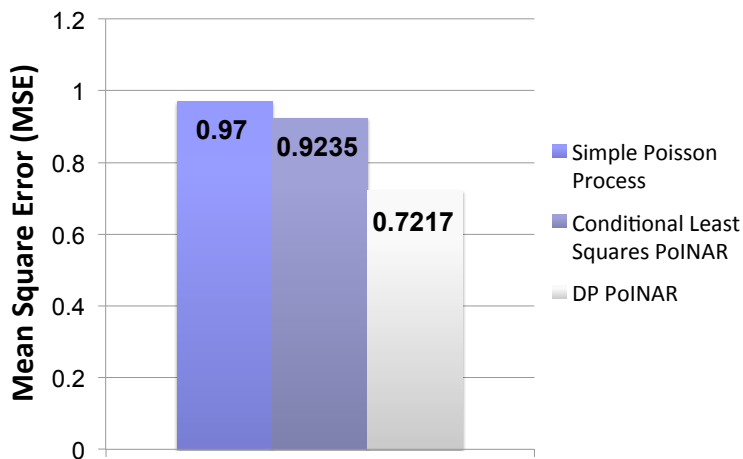- (0.229,0.3307]
- (0.07365,0.229]
- (0.009569,0.07365]

River

Train tracks

*Similar behavior in spatially disjoint tracts*
*→ Cluster census tracts*

22

# Poisson Integer-Valued Autoregressions

$\lambda_i = \theta_2$

*Customer counts in location i at time t*

*Poisson new customers*

$$y_{i,t} = \alpha \odot y_{i,t-1} + \epsilon_{i,t}(\lambda_i)$$

*Thinning of customers previously in queue*

*Rate of arrivals in location i*

$\lambda_i = \theta_1$

$\lambda_i = \theta_3$

- Cluster regions based on *rate of arriving crimes*
- Pool info within clusters
- Unknown number of clusters
  → *Dirichlet process*

Aldor-Noiman, Brown, Fox, and Stine, *arXiv:1304.5642*, April 2013

---

# Prediction Results



Bar chart — Mean Square Error (MSE):
- Simple Poisson Process: 0.97
- Conditional Least Squares PoINAR: 0.9235
- DP PoINAR: 0.7217

Aldor-Noiman, Brown, Fox, and Stine, *arXiv:1304.5642*, April 2013

# Acknowledgements

*Slides based on parts of the lecture notes of Erik Sudderth for "Applied Bayesian Nonparametrics" at Brown University*