



# Nonparametric Bayesian Methods

Tamara Broderick

ITT Career Development Assistant Professor  
EECS  
MIT

Michael I. Jordan

Pehong Chen Distinguished Professor  
EECS, Statistics  
UC Berkeley

# Nonparametric Bayesian Methods: Part I

Tamara Broderick  
ITT Career Development Assistant Professor  
EECS  
MIT

# Nonparametric Bayes

# Nonparametric Bayes

- Bayesian

# Nonparametric Bayes

- Bayesian

$$\mathbb{P}(\text{parameters}|\text{data}) \propto \mathbb{P}(\text{data}|\text{parameters})\mathbb{P}(\text{parameters})$$

# Nonparametric Bayes

- Bayesian

$$\mathbb{P}(\text{parameters}|\text{data}) \propto \mathbb{P}(\text{data}|\text{parameters})\mathbb{P}(\text{parameters})$$

- Not parametric (i.e. not finite parameter, unbounded/growing/infinite number of parameters)

# Nonparametric Bayes

- Bayesian

$$\mathbb{P}(\text{parameters}|\text{data}) \propto \mathbb{P}(\text{data}|\text{parameters})\mathbb{P}(\text{parameters})$$

- Not parametric (i.e. not finite parameter, unbounded/growing/infinite number of parameters)



[wikipedia.org]

# Nonparametric Bayes

- Bayesian

$$\mathbb{P}(\text{parameters}|\text{data}) \propto \mathbb{P}(\text{data}|\text{parameters})\mathbb{P}(\text{parameters})$$

- Not parametric (i.e. not finite parameter, unbounded/growing/infinite number of parameters)



“Wikipedia phenomenon”

[wikipedia.org]



# Nonparametric Bayes

- Bayesian

$$\mathbb{P}(\text{parameters}|\text{data}) \propto \mathbb{P}(\text{data}|\text{parameters})\mathbb{P}(\text{parameters})$$

- Not parametric (i.e. not finite parameter, unbounded/growing/infinite number of parameters)



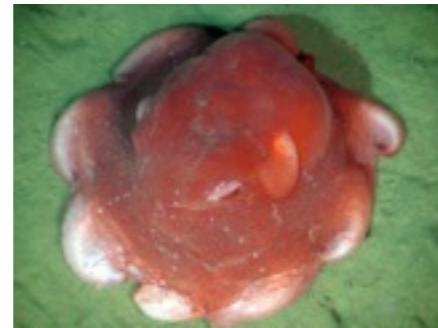
[wikipedia.org]

# Nonparametric Bayes

- Bayesian

$$\mathbb{P}(\text{parameters}|\text{data}) \propto \mathbb{P}(\text{data}|\text{parameters})\mathbb{P}(\text{parameters})$$

- Not parametric (i.e. not finite parameter, unbounded/growing/infinite number of parameters)



[Ed Bowlby, NOAA]

[wikipedia.org]

# Nonparametric Bayes

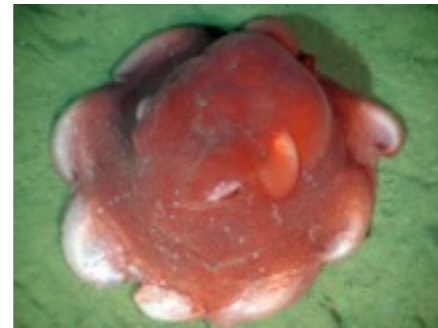
- Bayesian

$$\mathbb{P}(\text{parameters}|\text{data}) \propto \mathbb{P}(\text{data}|\text{parameters})\mathbb{P}(\text{parameters})$$

- Not parametric (i.e. not finite parameter, unbounded/growing/infinite number of parameters)



[wikipedia.org]



[Ed Bowlby, NOAA]



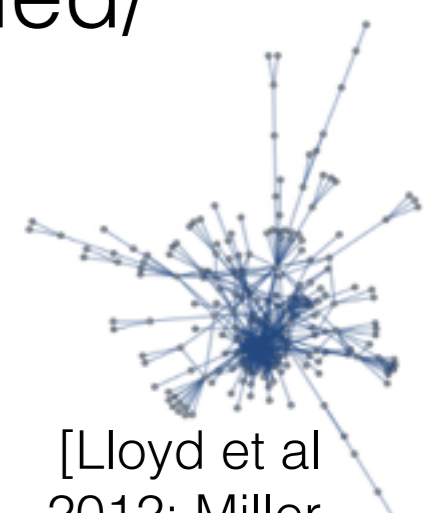
[Fox et al 2014]

# Nonparametric Bayes

- Bayesian

$$\mathbb{P}(\text{parameters}|\text{data}) \propto \mathbb{P}(\text{data}|\text{parameters})\mathbb{P}(\text{parameters})$$

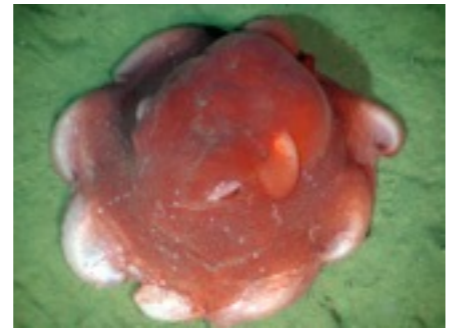
- Not parametric (i.e. not finite parameter, unbounded/growing/infinite number of parameters)



[Lloyd et al 2012; Miller et al 2010]



[wikipedia.org]



[Ed Bowlby, NOAA]



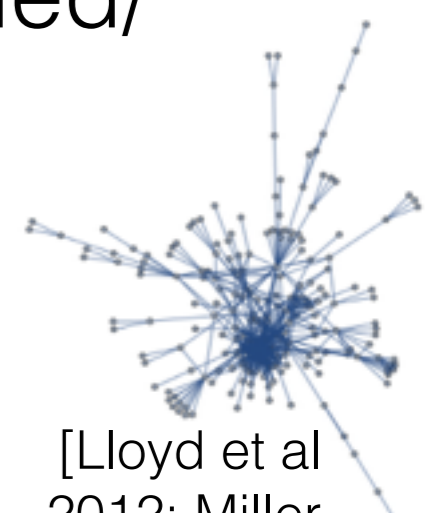
[Fox et al 2014]

# Nonparametric Bayes

- Bayesian

$$\mathbb{P}(\text{parameters}|\text{data}) \propto \mathbb{P}(\text{data}|\text{parameters})\mathbb{P}(\text{parameters})$$

- Not parametric (i.e. not finite parameter, unbounded/growing/infinite number of parameters)



[Lloyd et al 2012; Miller et al 2010]



[wikipedia.org]



[Ed Bowlby, NOAA]



[Fox et al 2014]



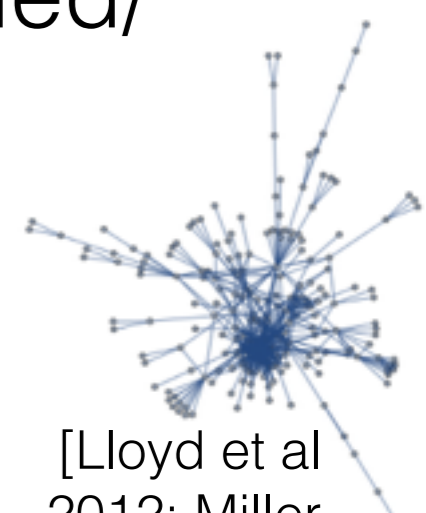
[Sudderth, Jordan 2009]

# Nonparametric Bayes

- Bayesian

$$\mathbb{P}(\text{parameters}|\text{data}) \propto \mathbb{P}(\text{data}|\text{parameters})\mathbb{P}(\text{parameters})$$

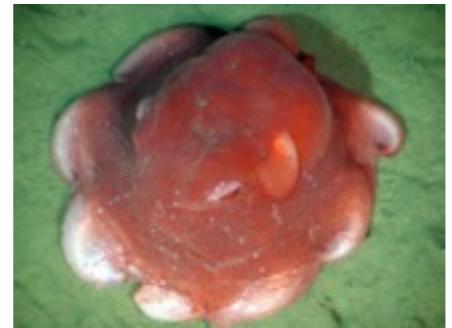
- Not parametric (i.e. not finite parameter, unbounded/growing/infinite number of parameters)



[Lloyd et al 2012; Miller et al 2010]



[wikipedia.org]



[Ed Bowlby, NOAA]



[Fox et al 2014]



[Ewens 1972; Hartl, Clark 2003]



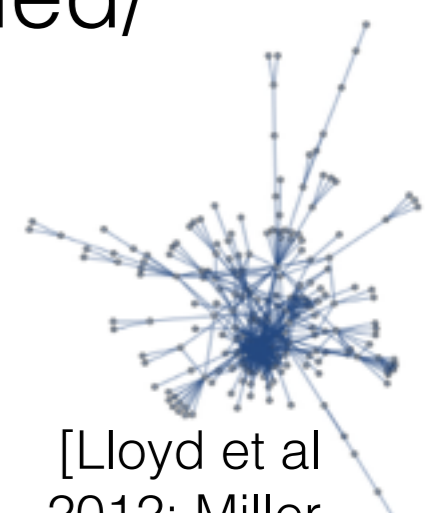
[Sudderth, Jordan 2009]

# Nonparametric Bayes

- Bayesian

$$\mathbb{P}(\text{parameters}|\text{data}) \propto \mathbb{P}(\text{data}|\text{parameters})\mathbb{P}(\text{parameters})$$

- Not parametric (i.e. not finite parameter, unbounded/growing/infinite number of parameters)



[Lloyd et al 2012; Miller et al 2010]



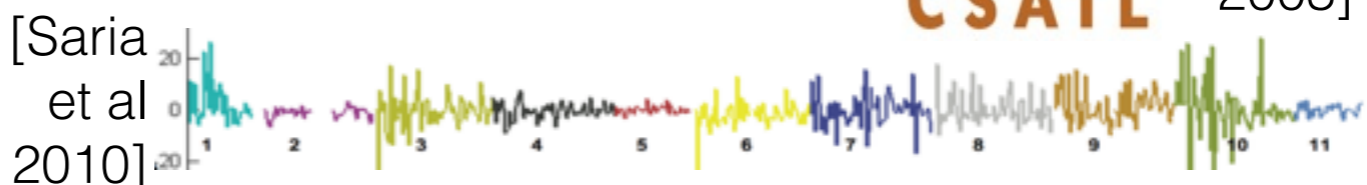
[Ed Bowlby, NOAA]



[Fox et al 2014]



[Ewens 1972; Hartl, Clark 2003]



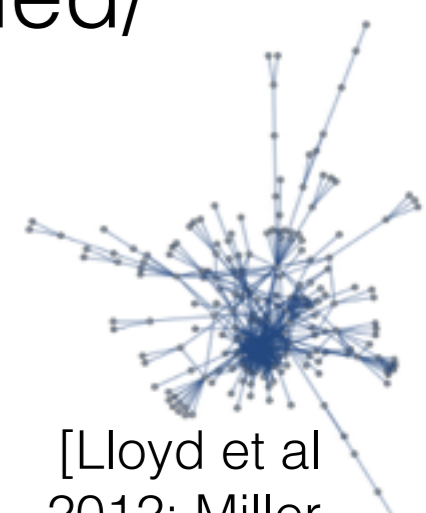
[Sudderth, Jordan 2009]

# Nonparametric Bayes

- Bayesian

$$\mathbb{P}(\text{parameters}|\text{data}) \propto \mathbb{P}(\text{data}|\text{parameters})\mathbb{P}(\text{parameters})$$

- Not parametric (i.e. not finite parameter, unbounded/growing/infinite number of parameters)



[Lloyd et al 2012; Miller et al 2010]



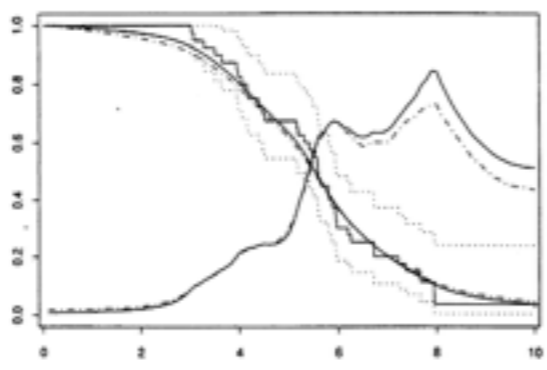
[wikipedia.org]



[Ed Bowlby, NOAA]



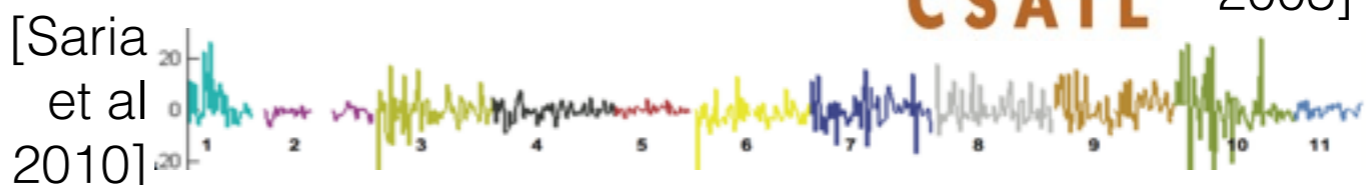
[Fox et al 2014]



[Arjas, Gasbarra 1994]



[Ewens 1972; Hartl, Clark 2003]



[Saria et al 2010]



[Sudderth, Jordan 2009]

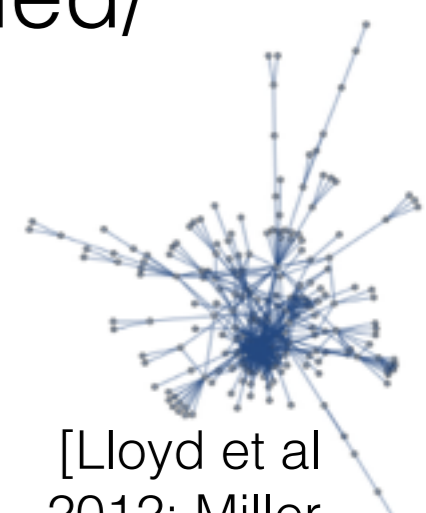


# Nonparametric Bayes

- Bayesian

$$\mathbb{P}(\text{parameters}|\text{data}) \propto \mathbb{P}(\text{data}|\text{parameters})\mathbb{P}(\text{parameters})$$

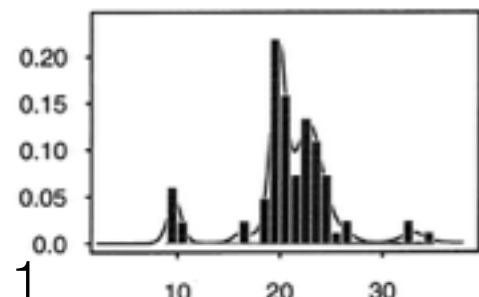
- Not parametric (i.e. not finite parameter, unbounded/growing/infinite number of parameters)



[Lloyd et al 2012; Miller et al 2010]



[wikipedia.org]



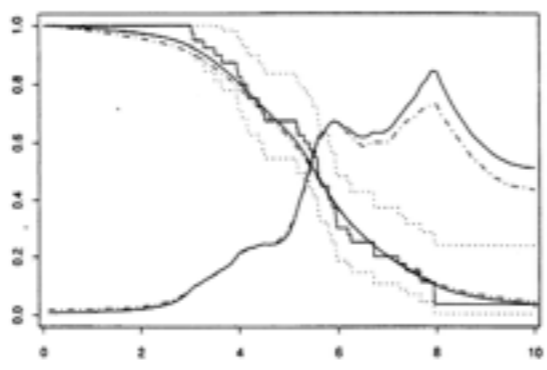
[Escobar, West 1995; Ghosal et al 1999]



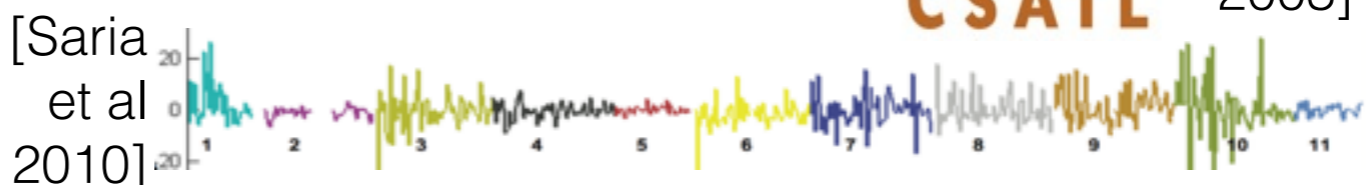
[Ed Bowlby, NOAA]



[Fox et al 2014]



[Arjas, Gasbarra 1994]



[Saria et al 2010]



[Ewens 1972; Hartl, Clark 2003]

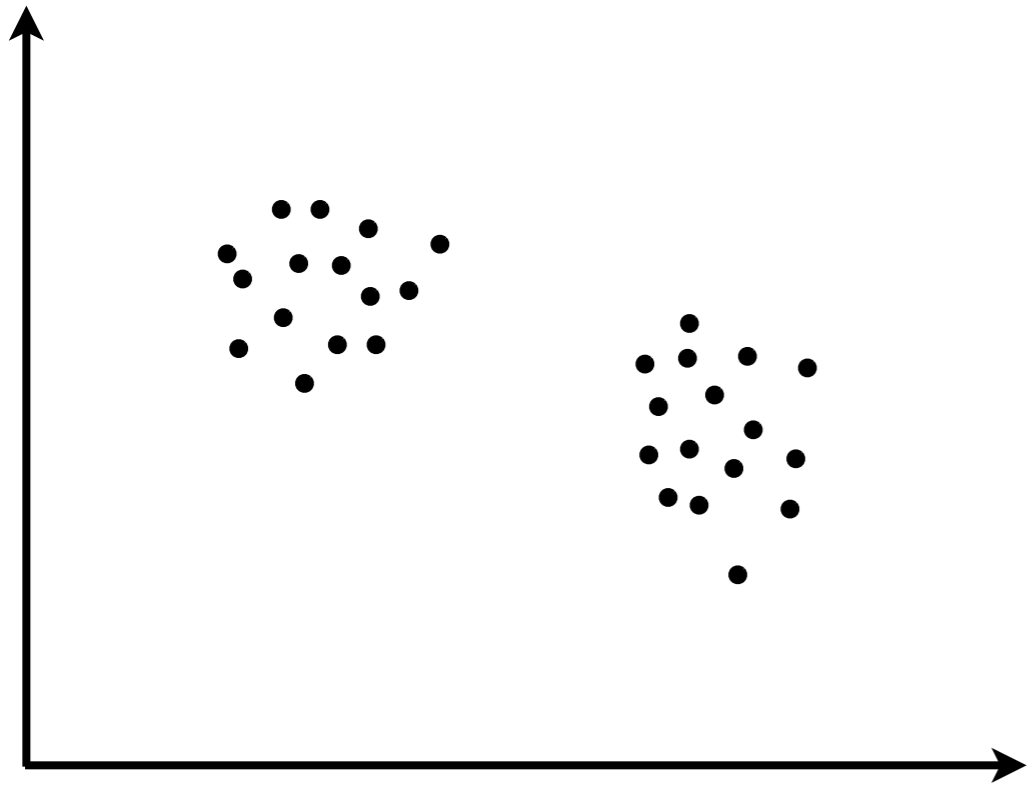


[Sudderth, Jordan 2009]

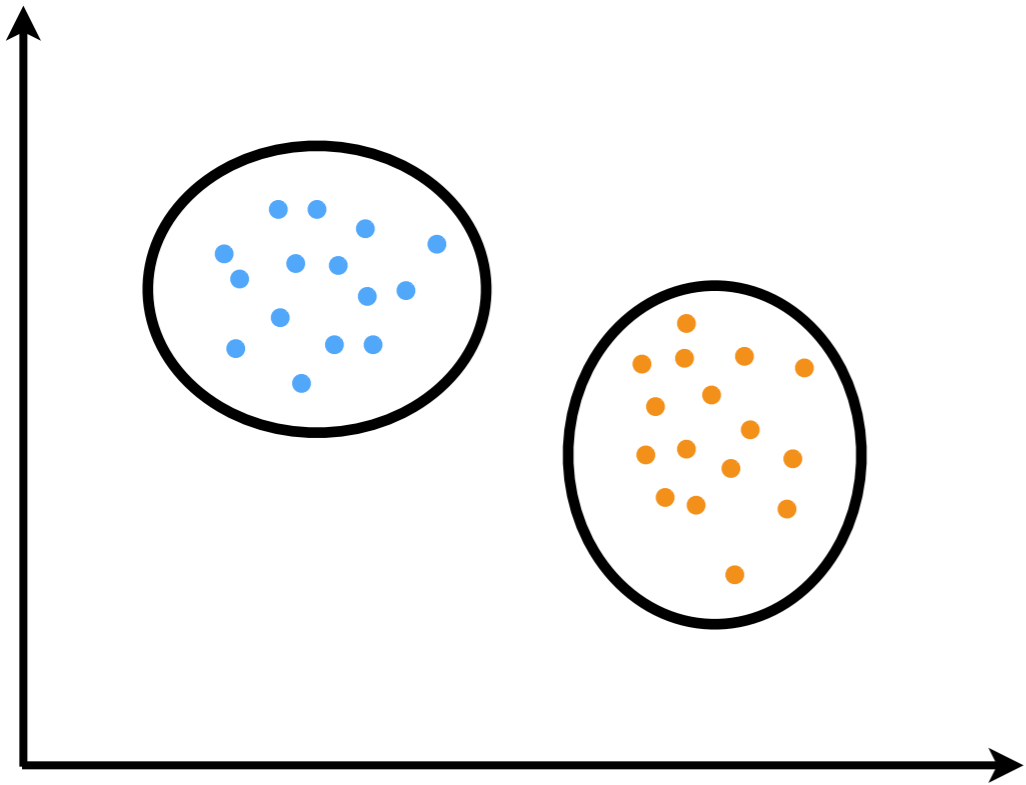
# Roadmap

- Example problem: clustering
- Example NPBayes model: Dirichlet process
- Big questions
  - Why NPBayes?
  - What does a growing/infinite number of parameters really mean (in NPBayes)?
  - Why is NPBayes challenging but practical?

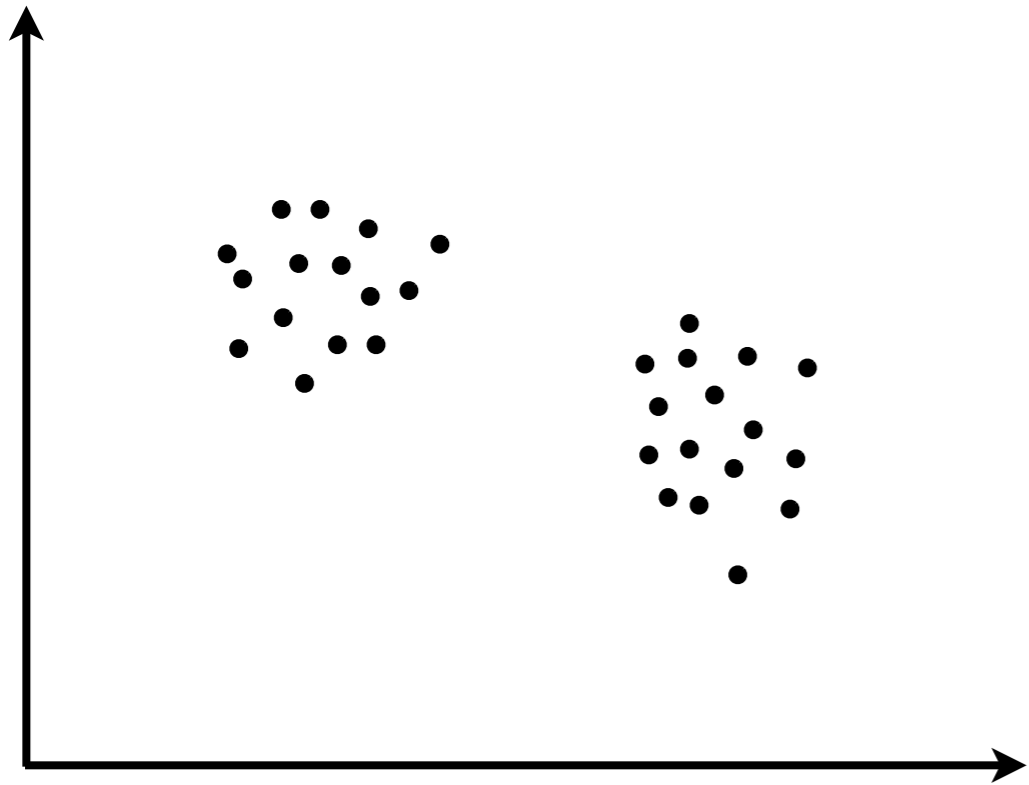
# Clustering



# Clustering

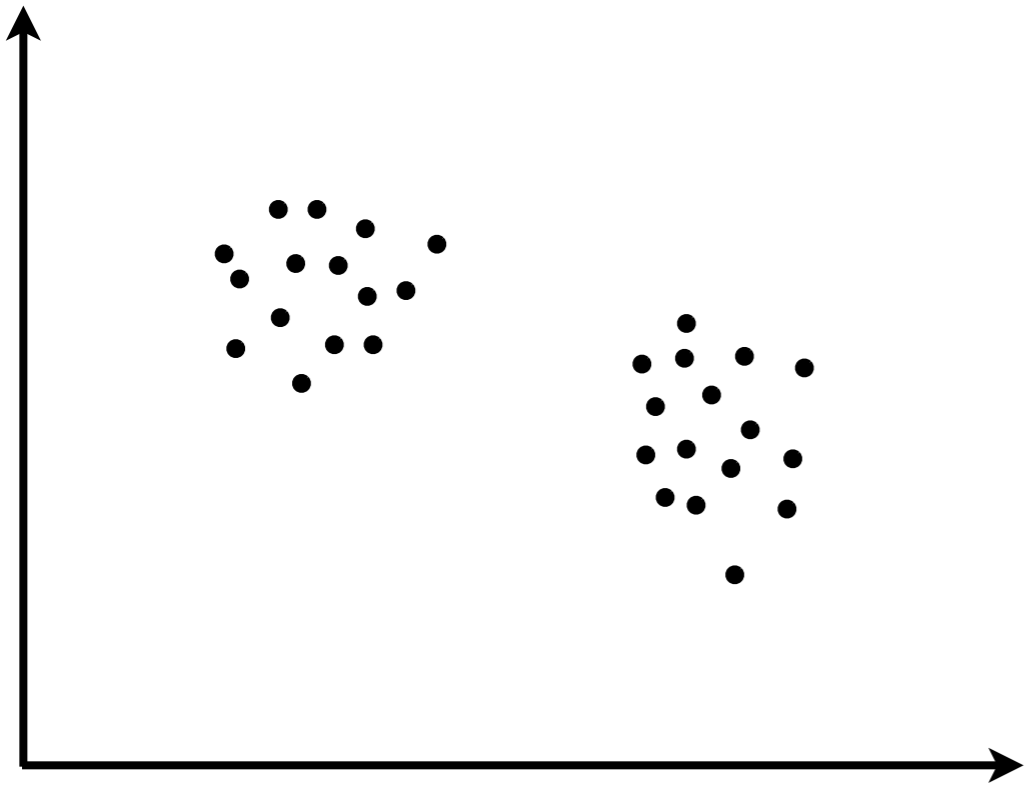


# Clustering



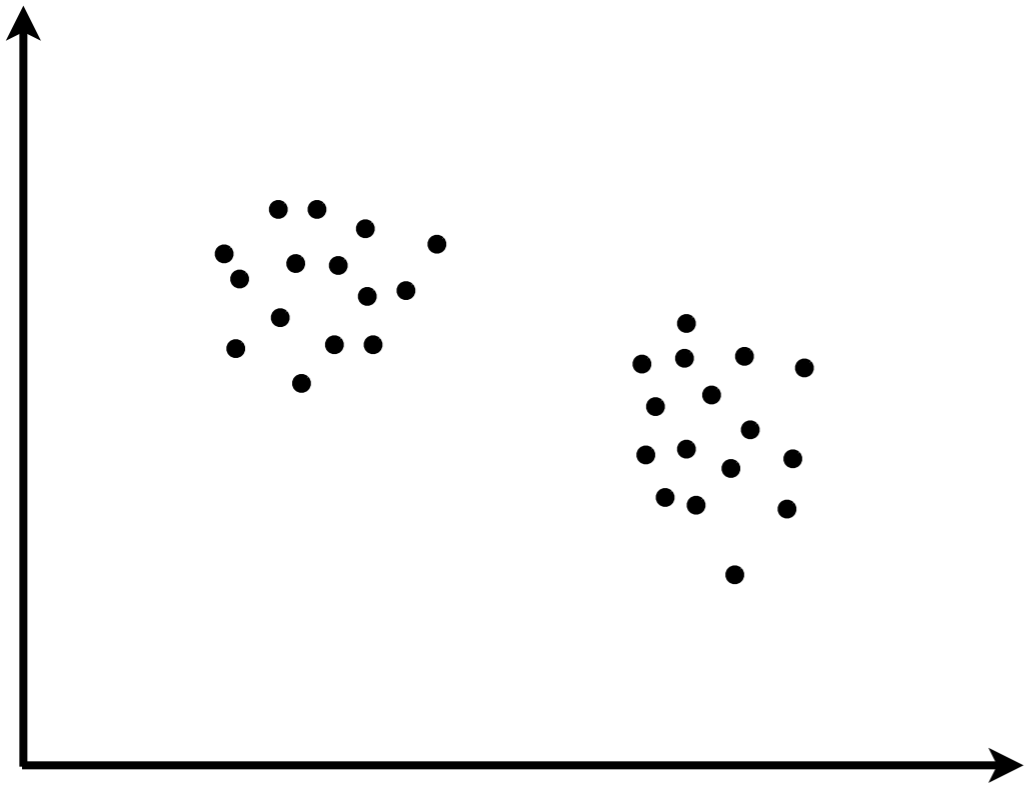
# Clustering

$$\mathbb{P}(\text{parameters}|\text{data}) \propto \mathbb{P}(\text{data}|\text{parameters})\mathbb{P}(\text{parameters})$$



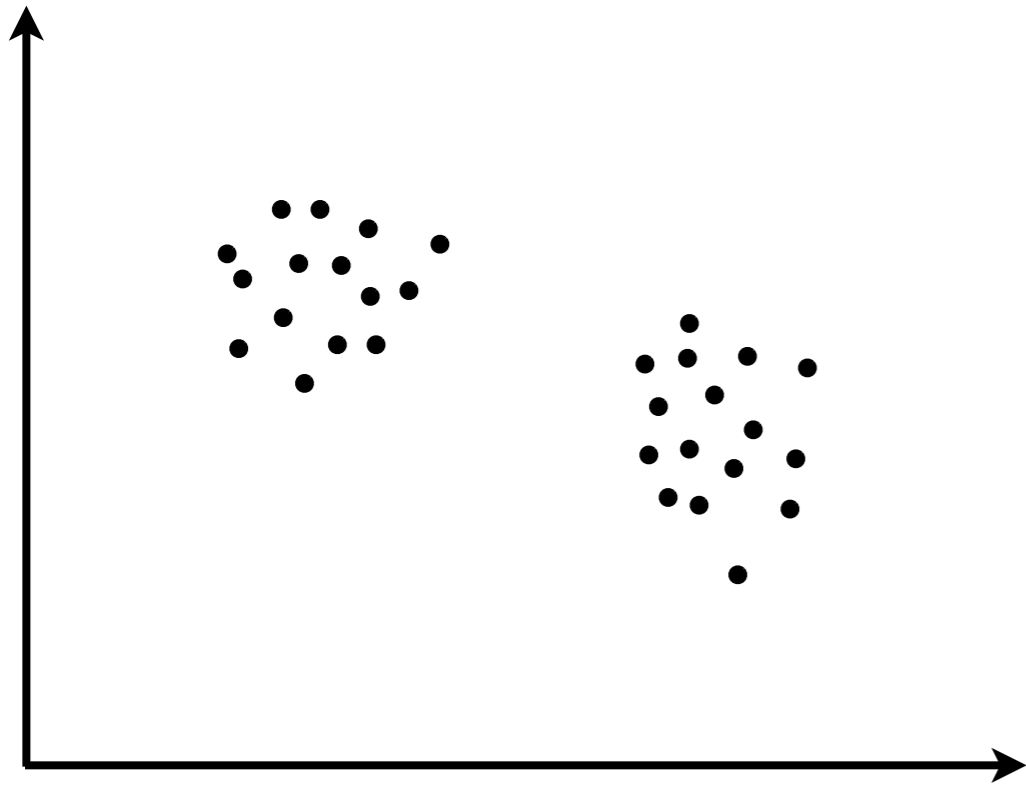
# Generative model

$$\mathbb{P}(\text{parameters}|\text{data}) \propto \mathbb{P}(\text{data}|\text{parameters})\mathbb{P}(\text{parameters})$$



# Generative model

$$\mathbb{P}(\text{parameters}|\text{data}) \propto \mathbb{P}(\text{data}|\text{parameters})\mathbb{P}(\text{parameters})$$

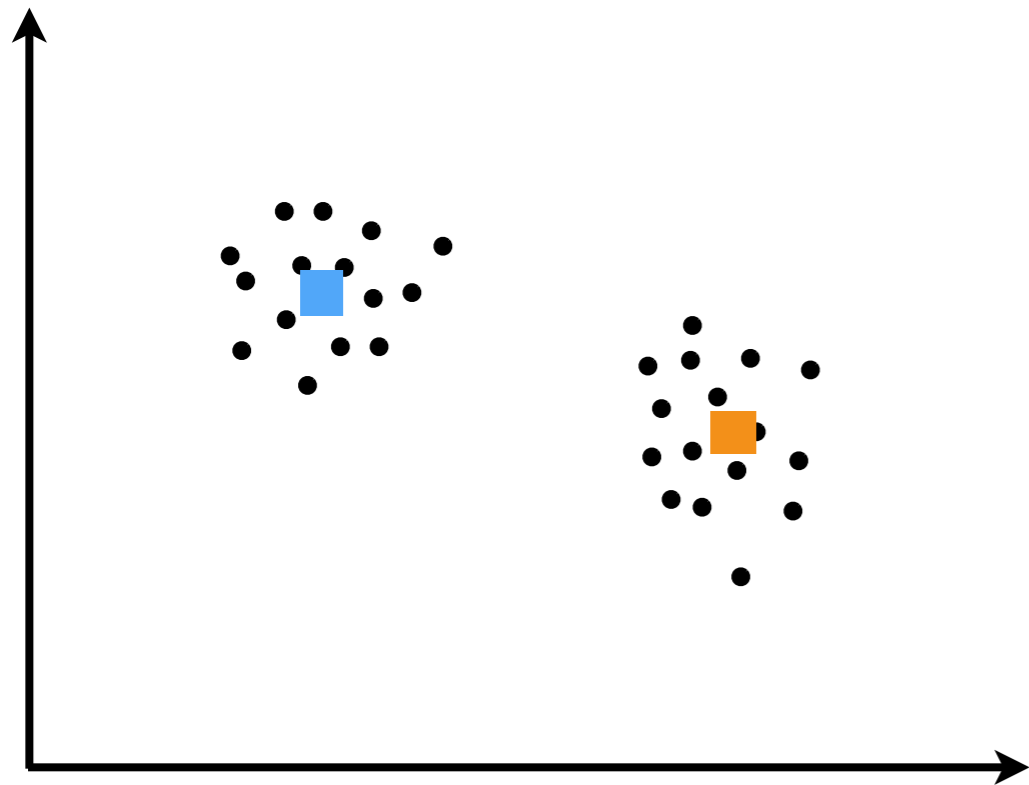


- Finite Gaussian mixture model ( $K=2$  clusters)



# Generative model

$$\mathbb{P}(\text{parameters}|\text{data}) \propto \mathbb{P}(\text{data}|\text{parameters})\mathbb{P}(\text{parameters})$$

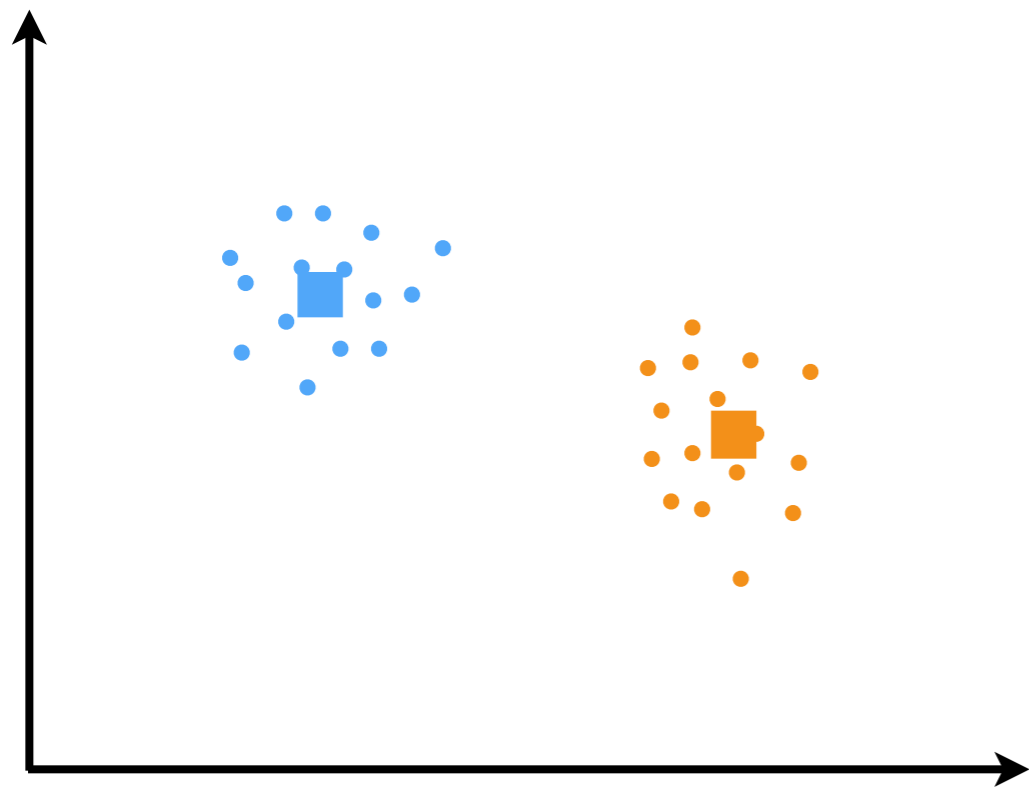


- Finite Gaussian mixture model ( $K=2$  clusters)

$\mu_k$

# Generative model

$$\mathbb{P}(\text{parameters}|\text{data}) \propto \mathbb{P}(\text{data}|\text{parameters})\mathbb{P}(\text{parameters})$$



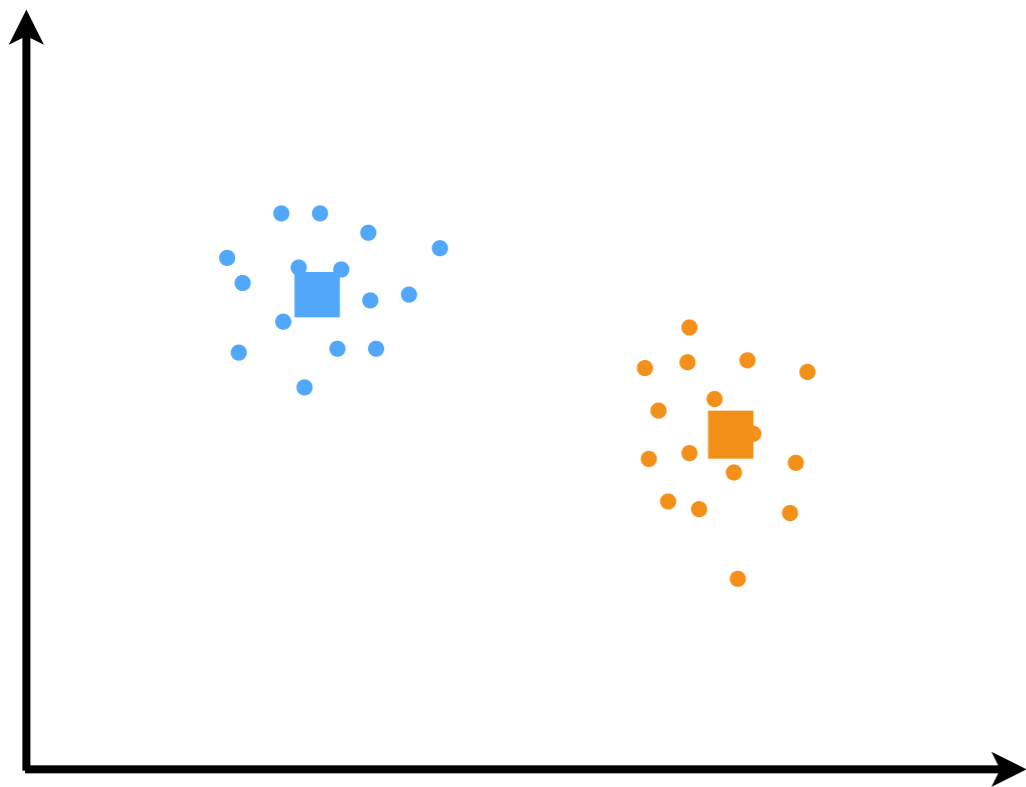
- Finite Gaussian mixture model ( $K=2$  clusters)

$\mu_k$

$z_n$

# Generative model

$$\mathbb{P}(\text{parameters}|\text{data}) \propto \mathbb{P}(\text{data}|\text{parameters})\mathbb{P}(\text{parameters})$$



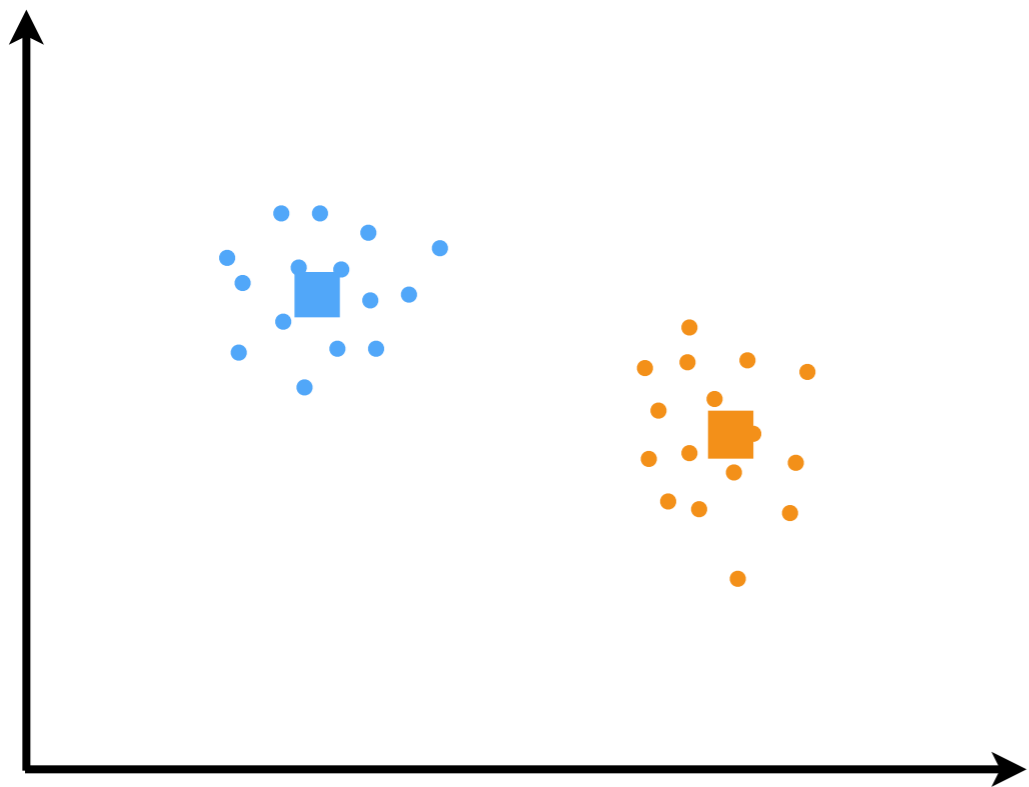
- Finite Gaussian mixture model ( $K=2$  clusters)

$\mu_k$

$$z_n$$
$$x_n \stackrel{\text{indep}}{\sim} \mathcal{N}(\mu_{z_n}, \Sigma)$$

# Generative model

$$\mathbb{P}(\text{parameters}|\text{data}) \propto \mathbb{P}(\text{data}|\text{parameters})\mathbb{P}(\text{parameters})$$



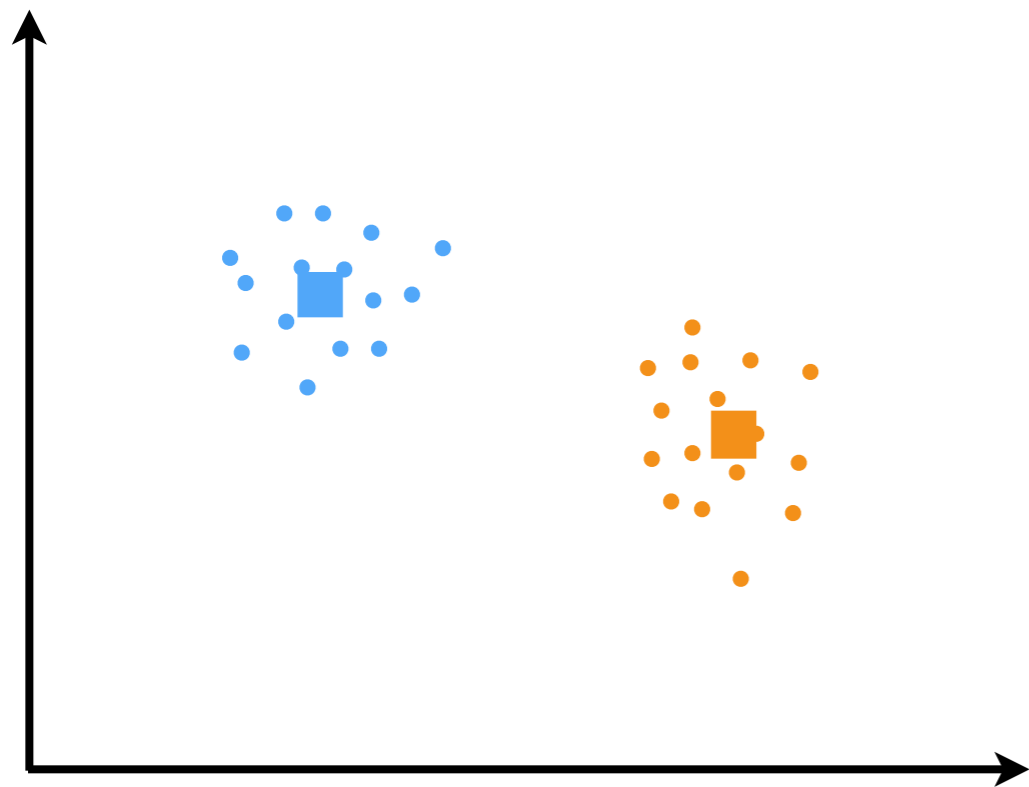
- Finite Gaussian mixture model ( $K=2$  clusters)

$$\mu_k \stackrel{iid}{\sim} \mathcal{N}(\mu_0, \Sigma_0)$$

$$z_n \stackrel{indep}{\sim} \mathcal{N}(\mu_{z_n}, \Sigma)$$

# Generative model

$$\mathbb{P}(\text{parameters}|\text{data}) \propto \mathbb{P}(\text{data}|\text{parameters})\mathbb{P}(\text{parameters})$$



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

$$\mu_k \stackrel{iid}{\sim} \mathcal{N}(\mu_0, \Sigma_0)$$

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

$$\rho_2 = 1 - \rho_1$$

$$z_n$$

$$x_n \stackrel{indep}{\sim} \mathcal{N}(\mu_{z_n}, \Sigma)$$

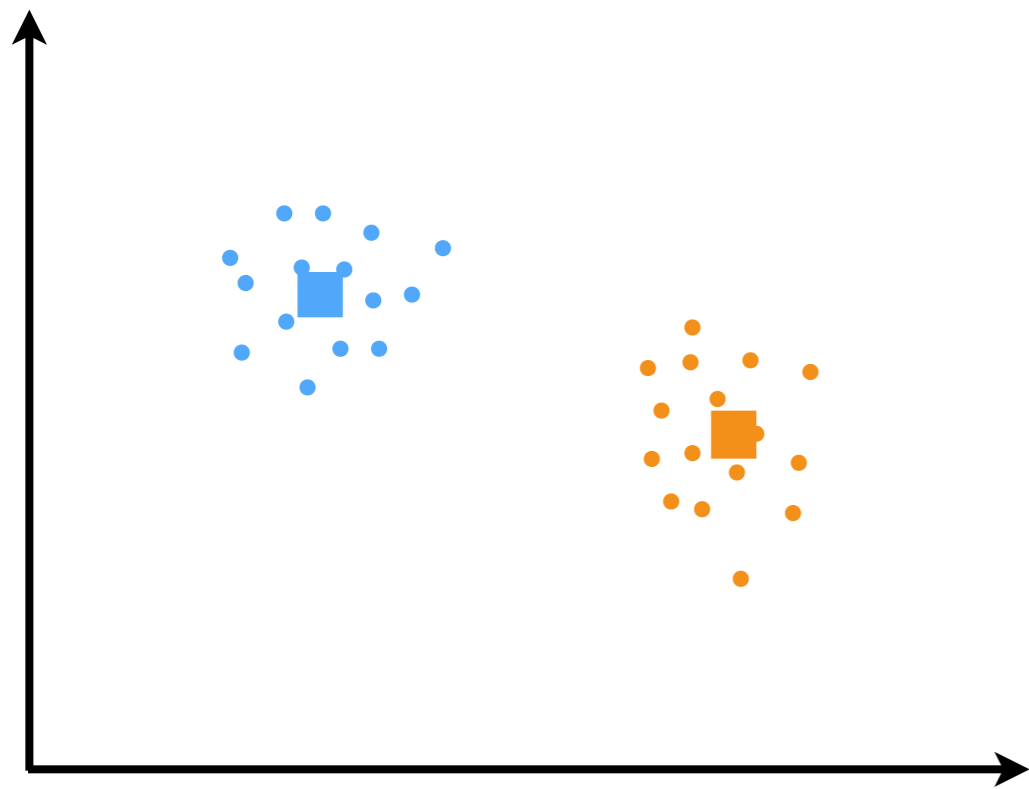


$\rho_1$

$\rho_2$

# Generative model

$$\mathbb{P}(\text{parameters}|\text{data}) \propto \mathbb{P}(\text{data}|\text{parameters})\mathbb{P}(\text{parameters})$$



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

$$\mu_k \stackrel{iid}{\sim} \mathcal{N}(\mu_0, \Sigma_0)$$

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

$$\rho_2 = 1 - \rho_1$$

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

$$x_n \stackrel{indep}{\sim} \mathcal{N}(\mu_{z_n}, \Sigma)$$

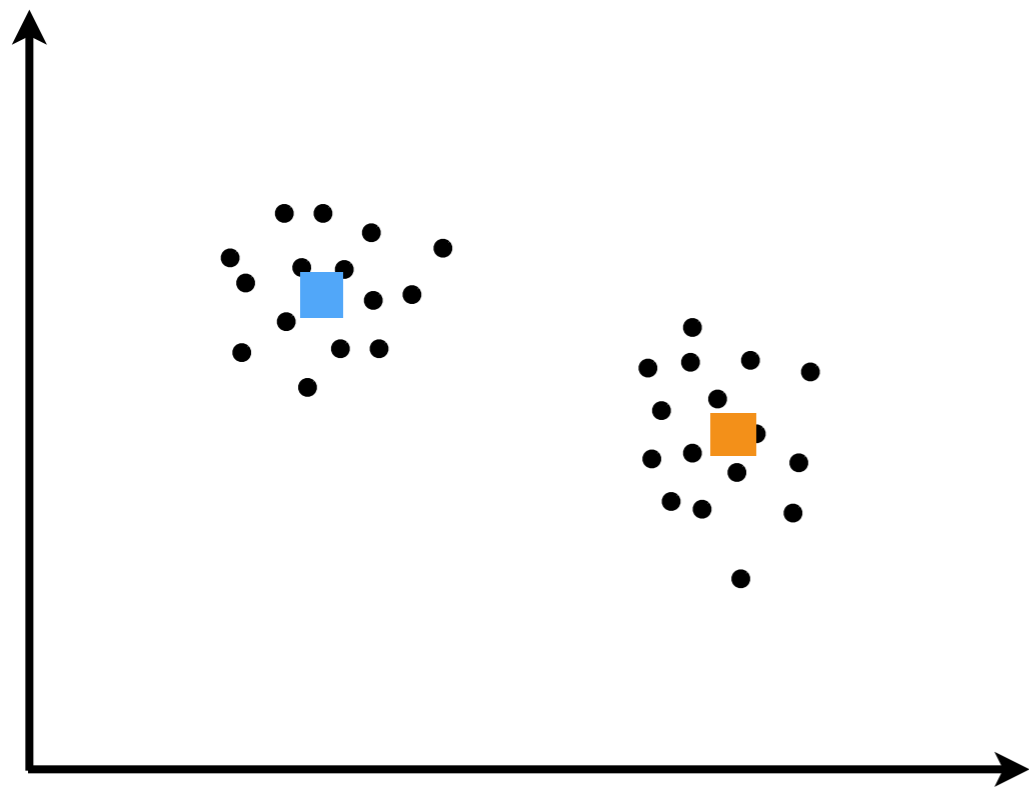


$\rho_1$

$\rho_2$

# Generative model

$$\mathbb{P}(\text{parameters}|\text{data}) \propto \mathbb{P}(\text{data}|\text{parameters})\mathbb{P}(\text{parameters})$$



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

$$\mu_k \stackrel{iid}{\sim} \mathcal{N}(\mu_0, \Sigma_0)$$

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

$$\rho_2 = 1 - \rho_1$$

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

$$x_n \stackrel{indep}{\sim} \mathcal{N}(\mu_{z_n}, \Sigma)$$



$\rho_1$

$\rho_2$

# Beta distribution review

$$\text{Beta}(\rho_1 | a_1, a_2) = \frac{\Gamma(a_1 + a_2)}{\Gamma(a_1)\Gamma(a_2)} \rho_1^{a_1-1} (1 - \rho_1)^{a_2-1}$$

$$\rho_1 \in (0, 1)$$

$$a_1, a_2 > 0$$

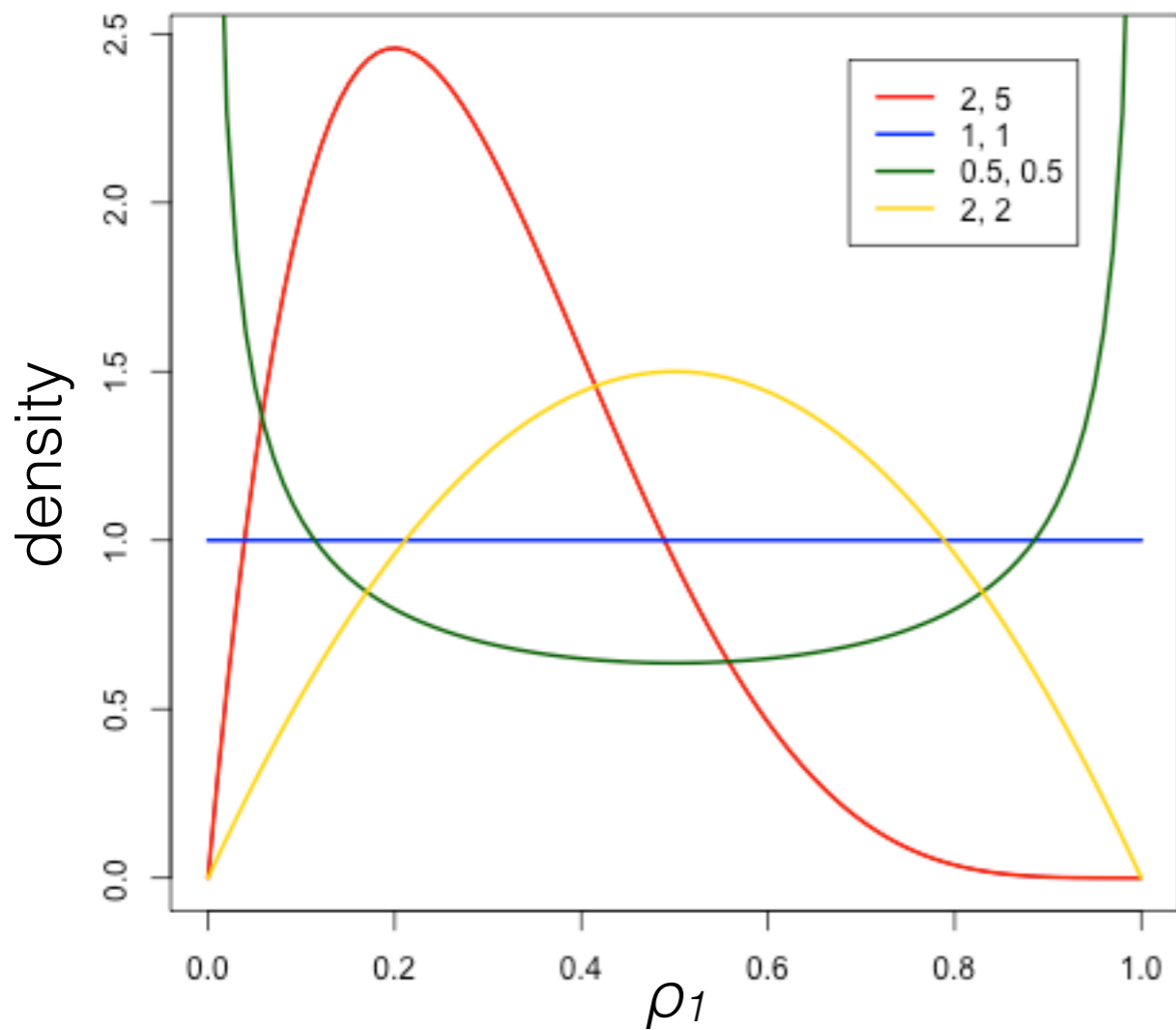


# Beta distribution review

$$\text{Beta}(\rho_1 | a_1, a_2) = \frac{\Gamma(a_1 + a_2)}{\Gamma(a_1)\Gamma(a_2)} \rho_1^{a_1-1} (1 - \rho_1)^{a_2-1}$$

$$\rho_1 \in (0, 1)$$

$$a_1, a_2 > 0$$

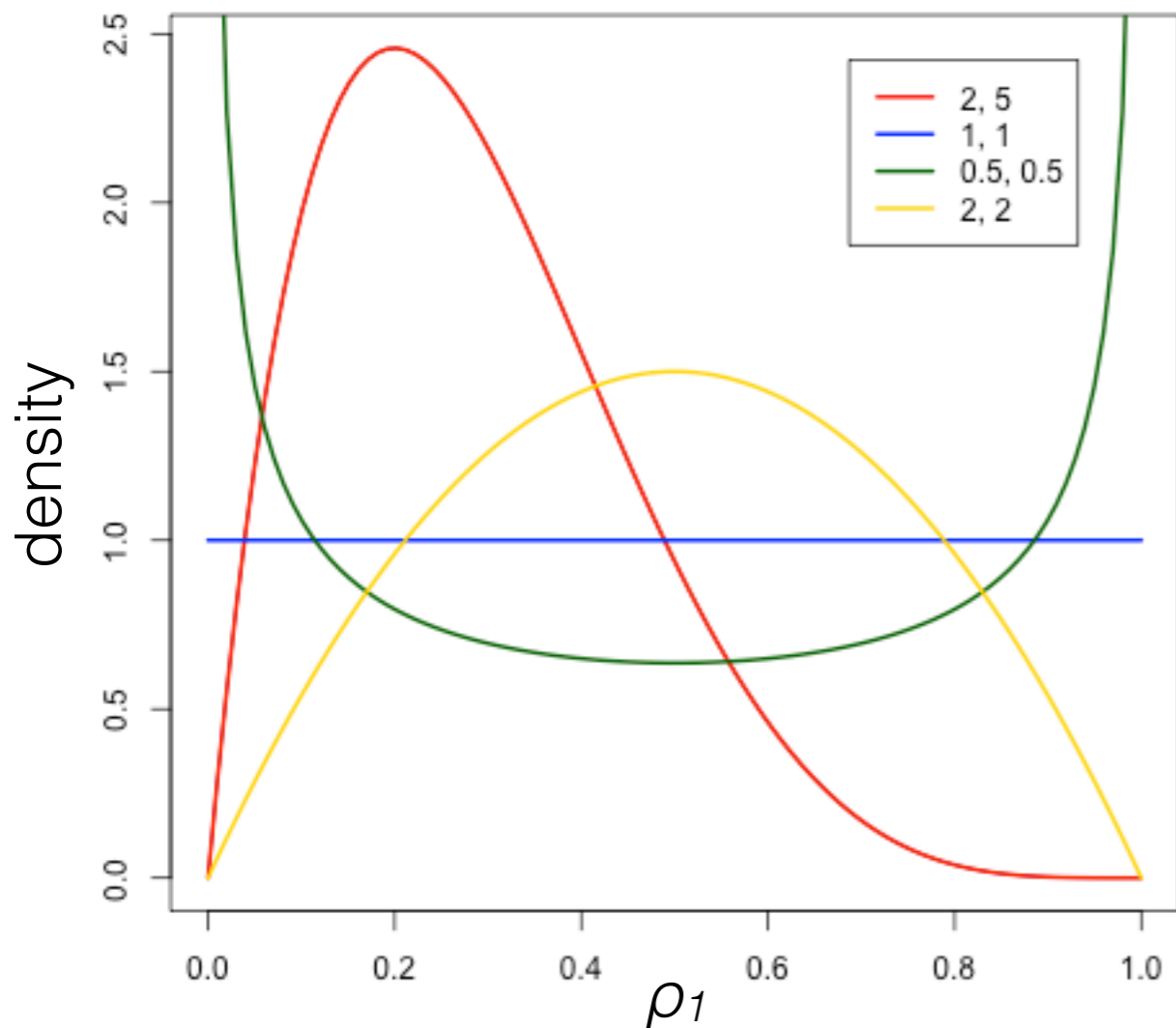


# Beta distribution review

$$\text{Beta}(\rho_1 | a_1, a_2) = \frac{\Gamma(a_1 + a_2)}{\Gamma(a_1)\Gamma(a_2)} \rho_1^{a_1-1} (1 - \rho_1)^{a_2-1}$$

$$\rho_1 \in (0, 1)$$

$$a_1, a_2 > 0$$



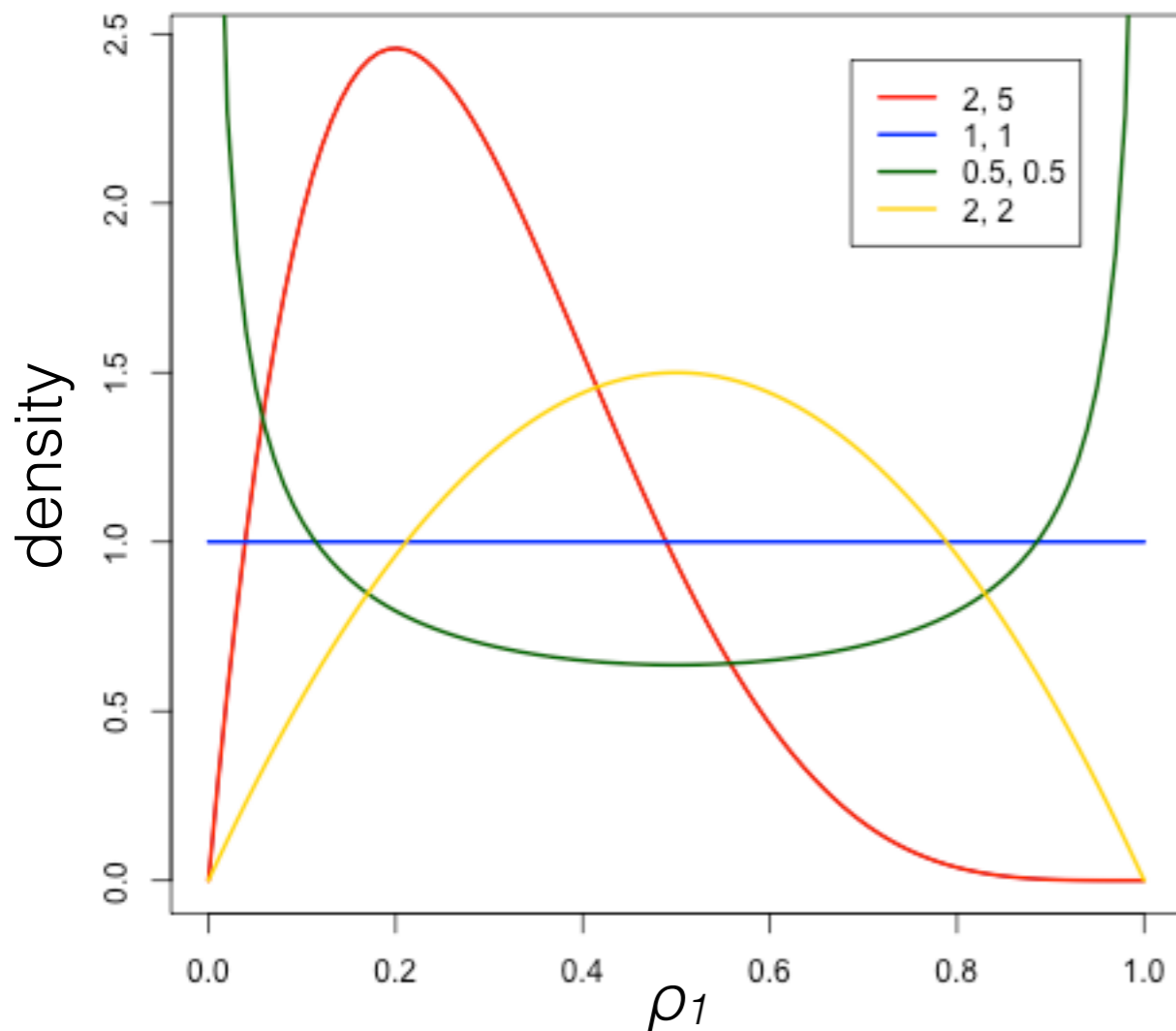
- What happens?

# Beta distribution review

$$\text{Beta}(\rho_1 | a_1, a_2) = \frac{\Gamma(a_1 + a_2)}{\Gamma(a_1)\Gamma(a_2)} \rho_1^{a_1-1} (1 - \rho_1)^{a_2-1}$$

$$\rho_1 \in (0, 1)$$

$$a_1, a_2 > 0$$



- What happens?

$$a = a_1 = a_2 \rightarrow 0$$

$$a = a_1 = a_2 \rightarrow \infty$$

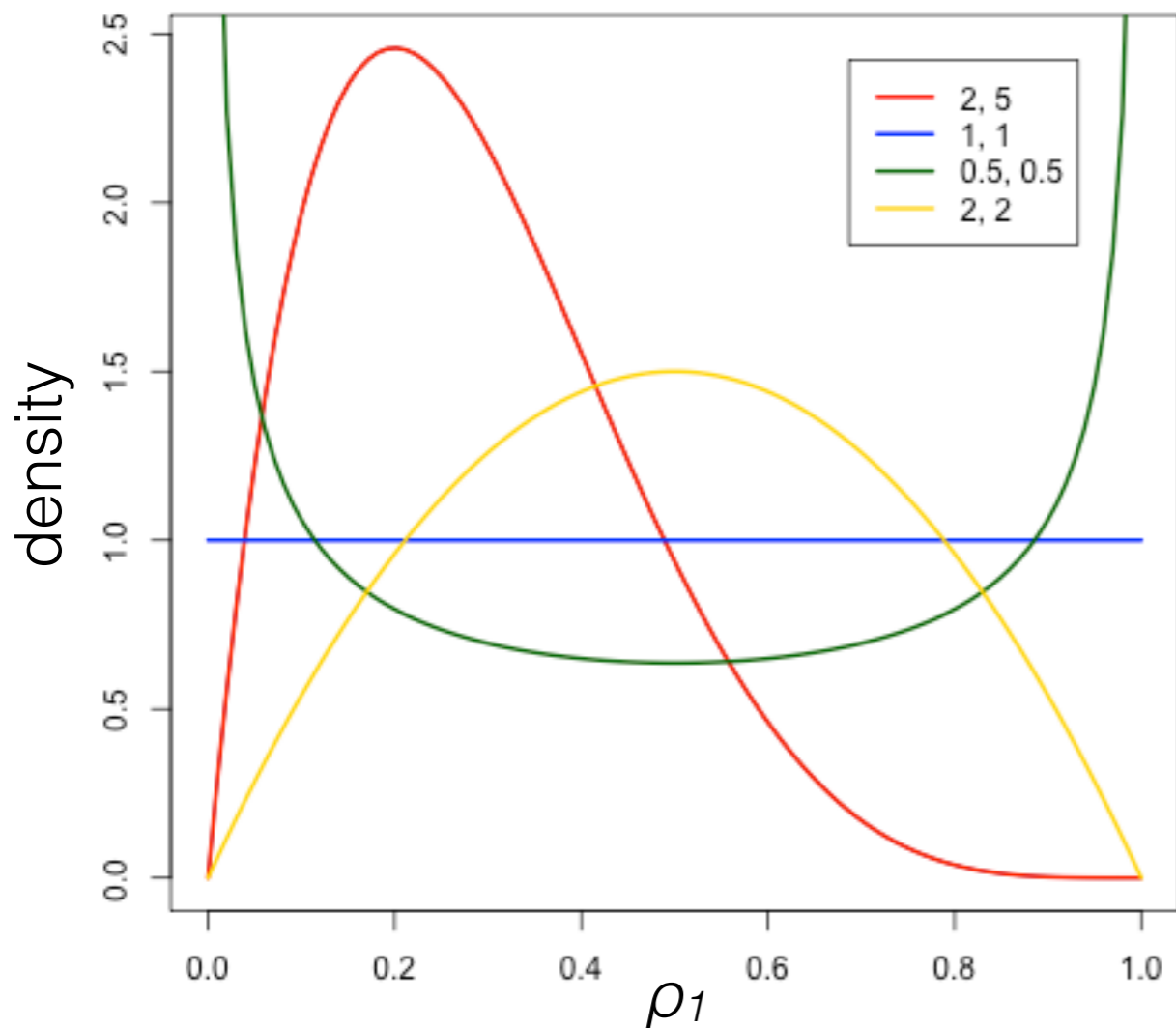
$$a_1 > a_2$$

# Beta distribution review

$$\text{Beta}(\rho_1 | a_1, a_2) = \frac{\Gamma(a_1 + a_2)}{\Gamma(a_1)\Gamma(a_2)} \rho_1^{a_1-1} (1 - \rho_1)^{a_2-1}$$

$$\rho_1 \in (0, 1)$$

$$a_1, a_2 > 0$$



- What happens?

$$a = a_1 = a_2 \rightarrow 0$$

$$a = a_1 = a_2 \rightarrow \infty$$

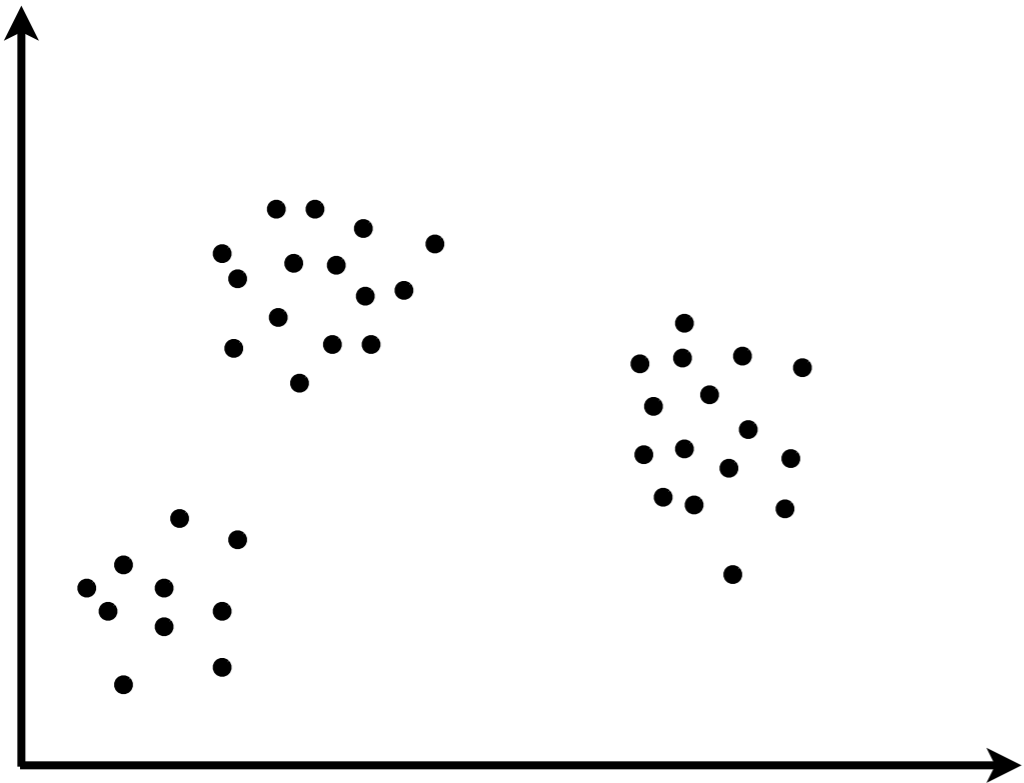
$$a_1 > a_2$$

[demo]

# Generative model

$$\mathbb{P}(\text{parameters}|\text{data}) \propto \mathbb{P}(\text{data}|\text{parameters})\mathbb{P}(\text{parameters})$$

- Finite Gaussian mixture model ( $K$  clusters)



# Generative model

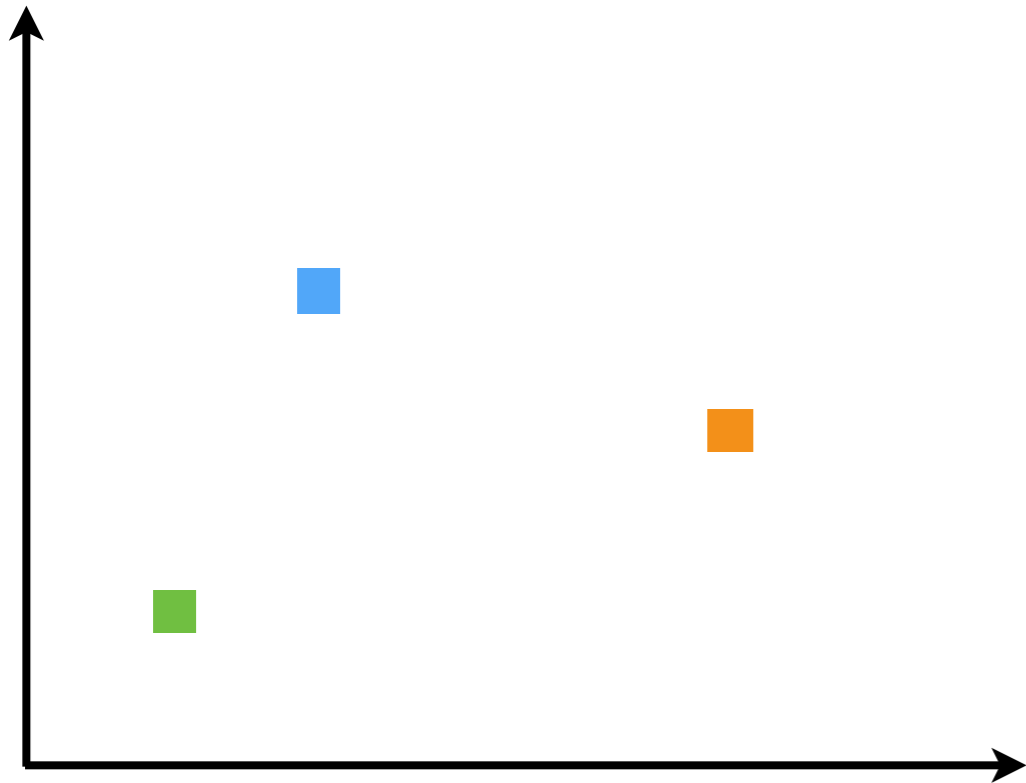
$$\mathbb{P}(\text{parameters}|\text{data}) \propto \mathbb{P}(\text{data}|\text{parameters})\mathbb{P}(\text{parameters})$$

- Finite Gaussian mixture model ( $K$  clusters)



# Generative model

$$\mathbb{P}(\text{parameters}|\text{data}) \propto \mathbb{P}(\text{data}|\text{parameters})\mathbb{P}(\text{parameters})$$



- Finite Gaussian mixture model ( $K$  clusters)

$$\mu_k \stackrel{iid}{\sim} \mathcal{N}(\mu_0, \Sigma_0)$$

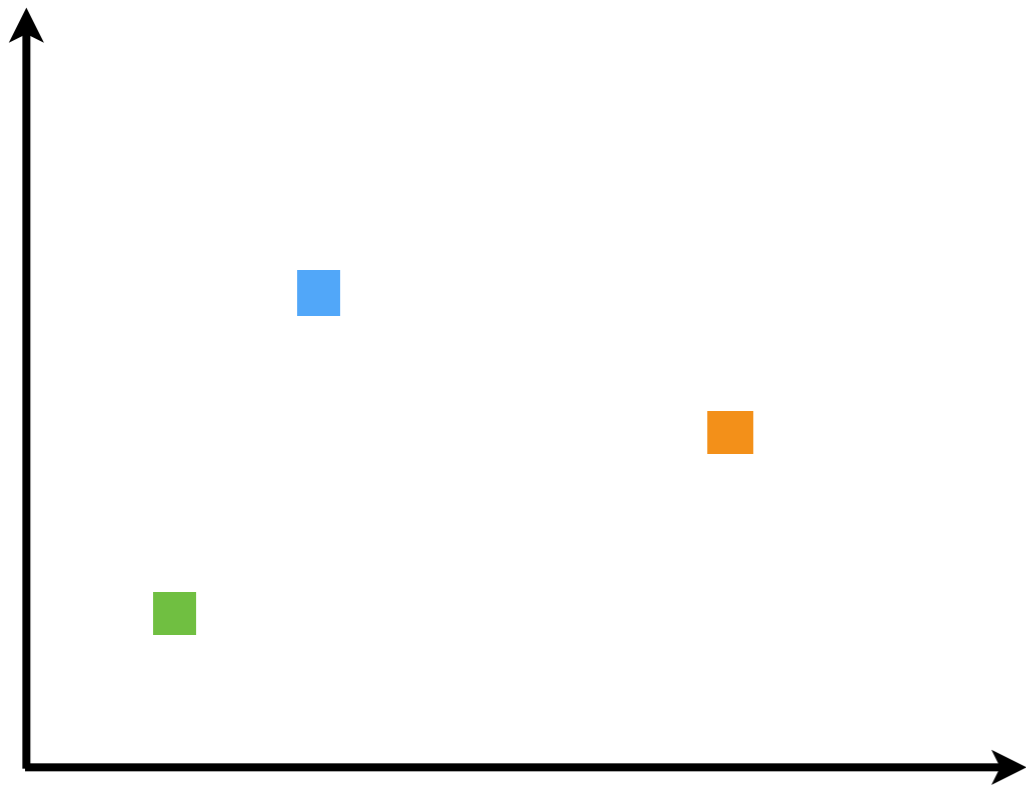
# Generative model

$$\mathbb{P}(\text{parameters}|\text{data}) \propto \mathbb{P}(\text{data}|\text{parameters})\mathbb{P}(\text{parameters})$$

- Finite Gaussian mixture model ( $K$  clusters)

$$\mu_k \stackrel{iid}{\sim} \mathcal{N}(\mu_0, \Sigma_0)$$

$$\rho_{1:K} \sim \text{Dirichlet}(a_{1:K})$$



$\rho_1$

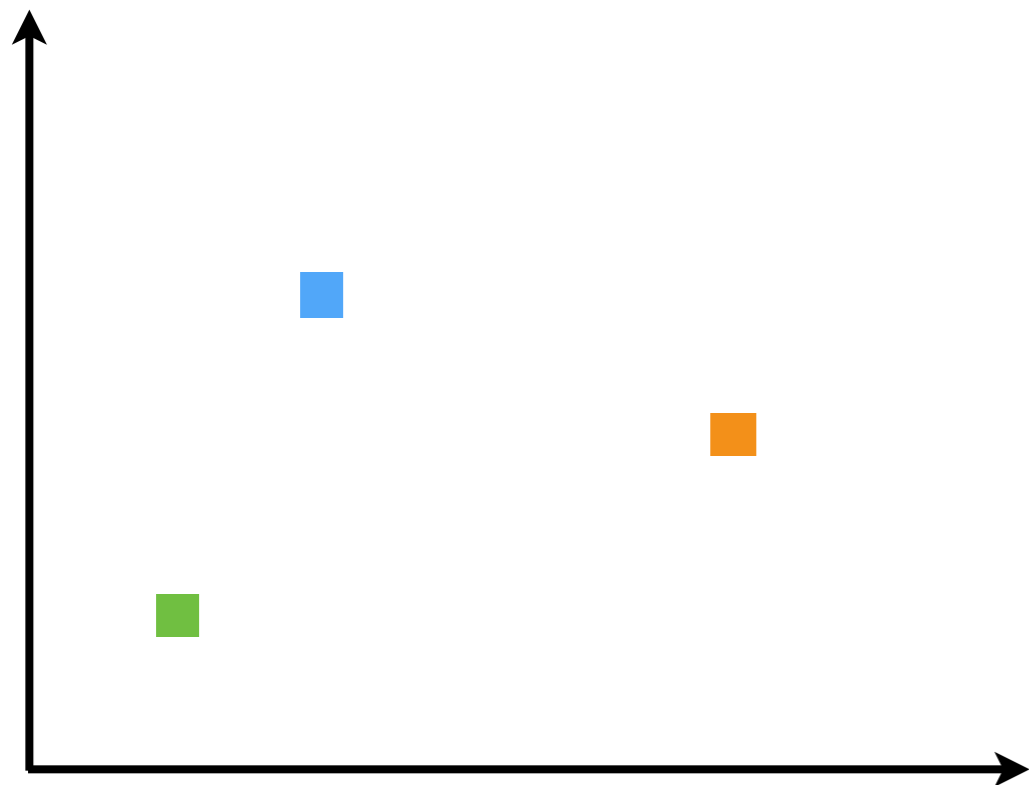
$\rho_2$

$\rho_3$



# Generative model

$$\mathbb{P}(\text{parameters}|\text{data}) \propto \mathbb{P}(\text{data}|\text{parameters})\mathbb{P}(\text{parameters})$$



- Finite Gaussian mixture model ( $K$  clusters)

$$\mu_k \stackrel{iid}{\sim} \mathcal{N}(\mu_0, \Sigma_0)$$

$$\rho_{1:K} \sim \text{Dirichlet}(a_{1:K})$$

$$z_n \stackrel{iid}{\sim} \text{Categorical}(\rho_{1:K})$$



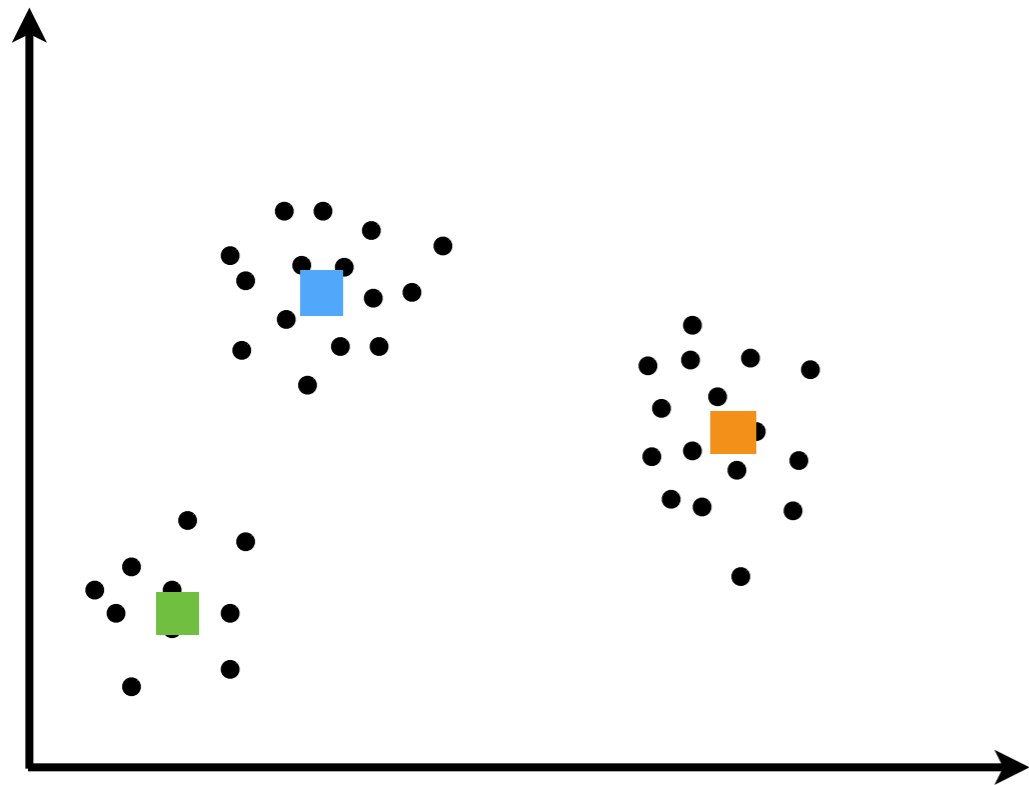
$\rho_1$

$\rho_2$

$\rho_3$

# Generative model

$$\mathbb{P}(\text{parameters}|\text{data}) \propto \mathbb{P}(\text{data}|\text{parameters})\mathbb{P}(\text{parameters})$$



- Finite Gaussian mixture model ( $K$  clusters)

$$\mu_k \stackrel{iid}{\sim} \mathcal{N}(\mu_0, \Sigma_0)$$

$$\rho_{1:K} \sim \text{Dirichlet}(a_{1:K})$$

$$z_n \stackrel{iid}{\sim} \text{Categorical}(\rho_{1:K})$$

$$x_n \stackrel{indep}{\sim} \mathcal{N}(\mu_{z_n}, \Sigma)$$



$\rho_1$

$\rho_2$

$\rho_3$

# Dirichlet distribution review

$$\text{Dirichlet}(\rho_{1:K} | a_{1:K}) = \frac{\Gamma(\sum_{k=1}^K a_k)}{\prod_{k=1}^K \Gamma(a_k)} \prod_{k=1}^K \rho_k^{a_k - 1} \quad a_k > 0$$

# Dirichlet distribution review

$$\text{Dirichlet}(\rho_{1:K} | a_{1:K}) = \frac{\Gamma(\sum_{k=1}^K a_k)}{\prod_{k=1}^K \Gamma(a_k)} \prod_{k=1}^K \rho_k^{a_k - 1}$$

$$a_k > 0$$

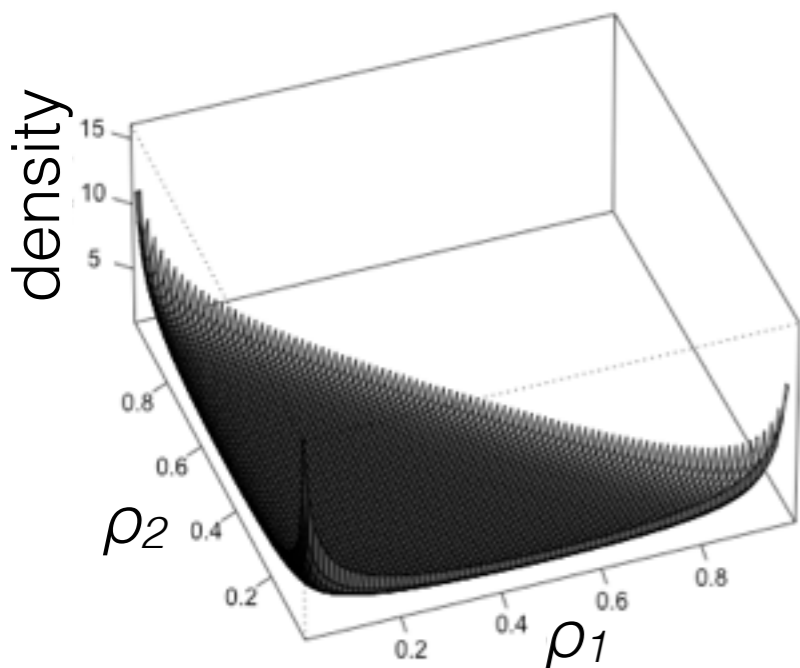
$$\rho_k \in (0, 1)$$

$$\sum_k \rho_k = 1$$

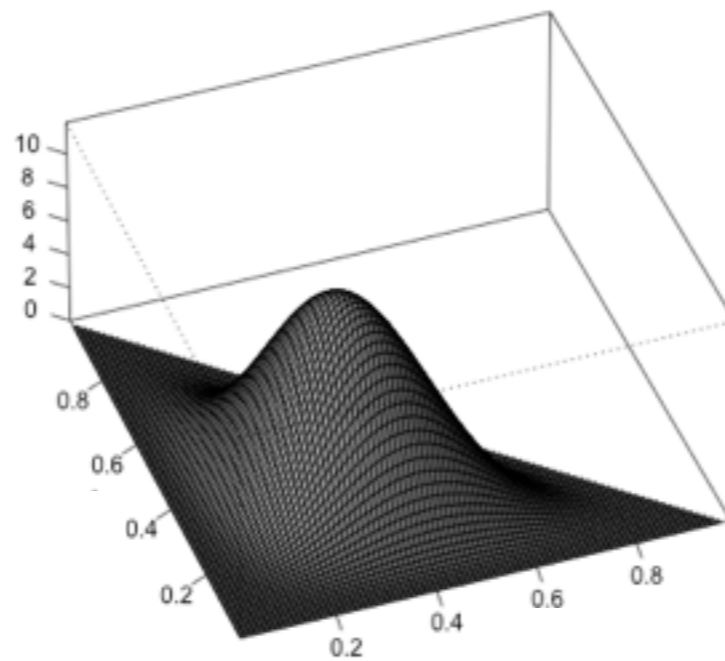
# Dirichlet distribution review

$$\text{Dirichlet}(\rho_{1:K} | a_{1:K}) = \frac{\Gamma(\sum_{k=1}^K a_k)}{\prod_{k=1}^K \Gamma(a_k)} \prod_{k=1}^K \rho_k^{a_k - 1} \quad \begin{array}{l} a_k > 0 \\ \rho_k \in (0, 1) \\ \sum_k \rho_k = 1 \end{array}$$

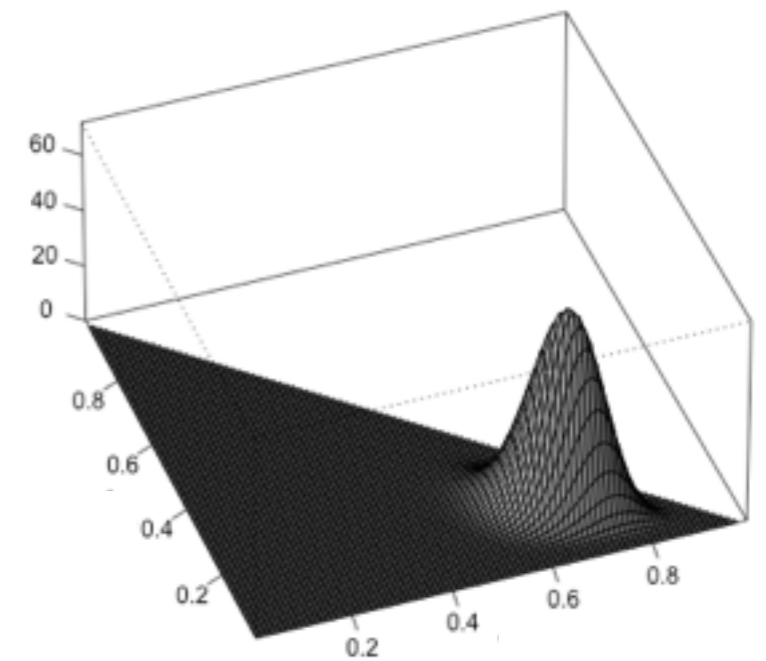
$a = (0.5, 0.5, 0.5)$



$a = (5, 5, 5)$



$a = (40, 10, 10)$

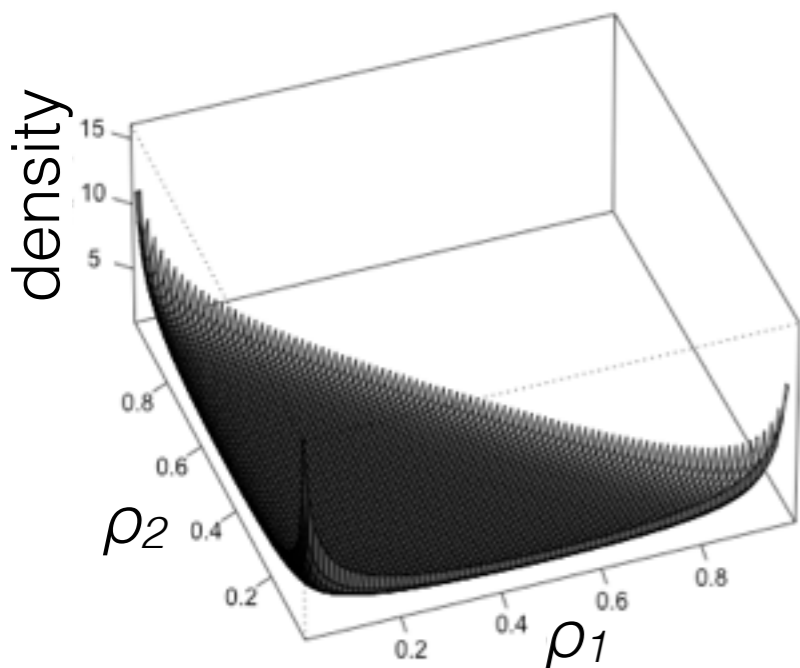


# Dirichlet distribution review

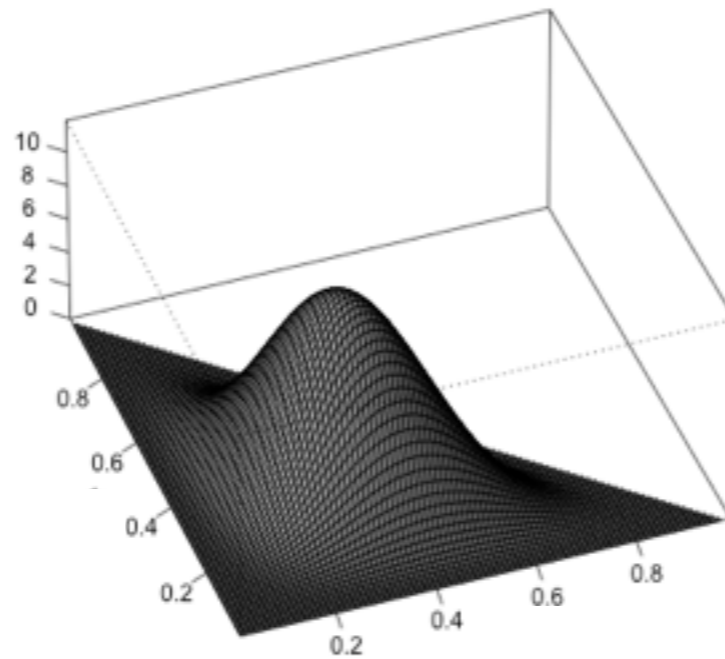
$$\text{Dirichlet}(\rho_{1:K} | a_{1:K}) = \frac{\Gamma(\sum_{k=1}^K a_k)}{\prod_{k=1}^K \Gamma(a_k)} \prod_{k=1}^K \rho_k^{a_k - 1}$$

$a_k > 0$   
 $\rho_k \in (0, 1)$   
 $\sum_k \rho_k = 1$

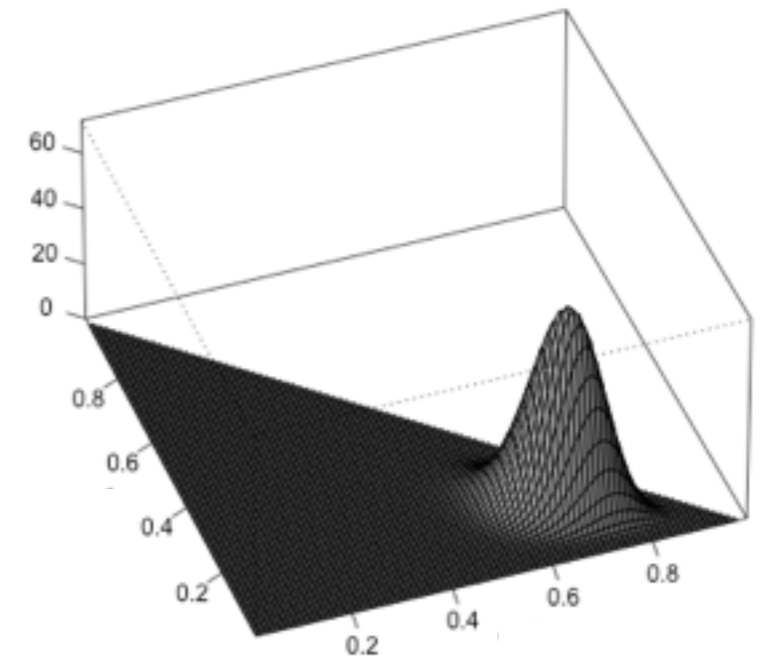
$a = (0.5, 0.5, 0.5)$



$a = (5, 5, 5)$



$a = (40, 10, 10)$



- What happens?

$$a = a_k = 1$$

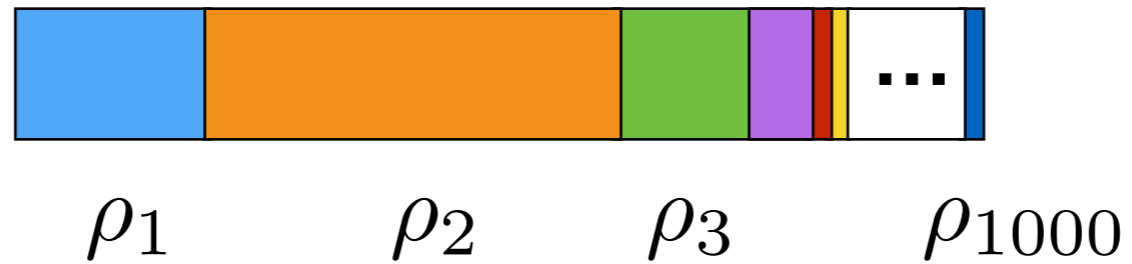
$$a = a_k \rightarrow 0$$

$$a = a_k \rightarrow \infty$$

[demo]

So far  $K \ll N$ . What if not?

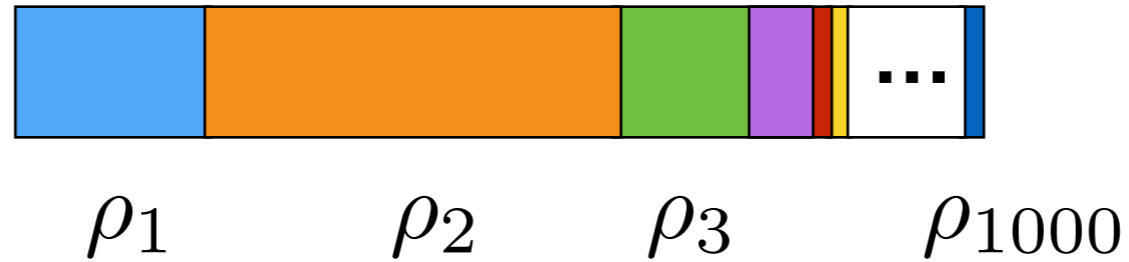
So far  $K \ll N$ . What if not?





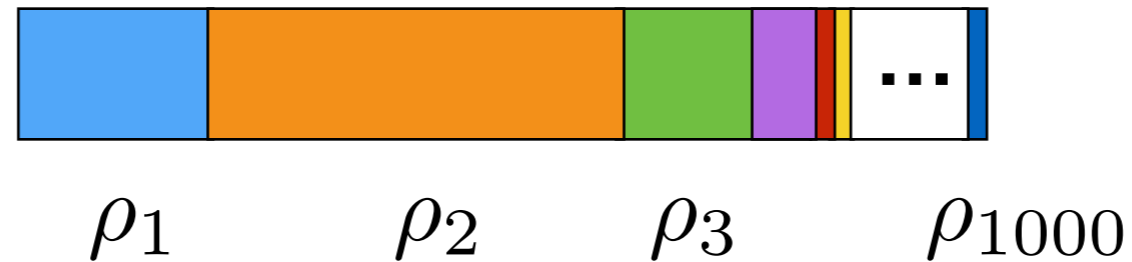
# So far $K \ll N$ . What if not?

- e.g. species sampling, topic modeling, groups on a social network, etc.



# So far $K \ll N$ . What if not?

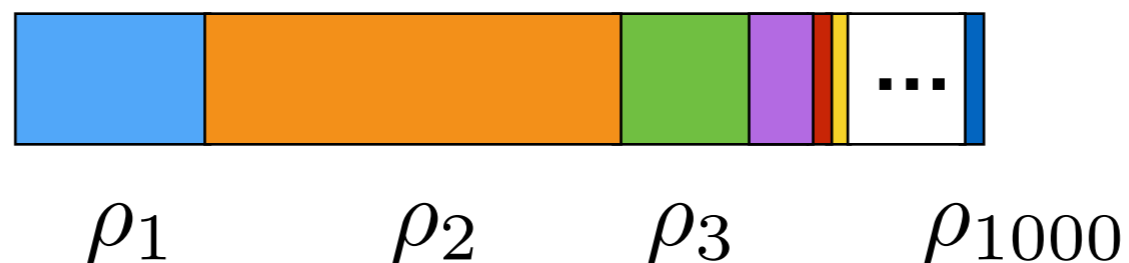
- e.g. species sampling, topic modeling, groups on a social network, etc.



- Components: number of latent groups

# So far $K \ll N$ . What if not?

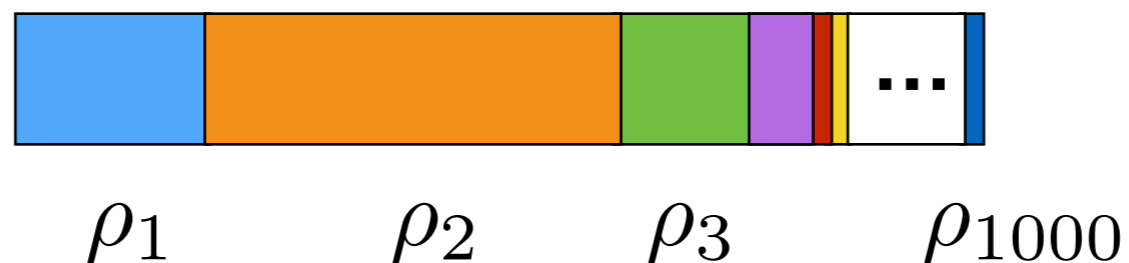
- e.g. species sampling, topic modeling, groups on a social network, etc.



- Components: number of latent groups
- Clusters: number of components represented in the data

# So far $K \ll N$ . What if not?

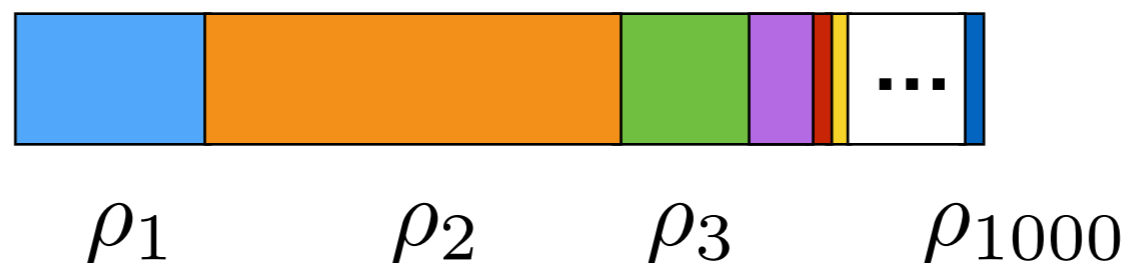
- e.g. species sampling, topic modeling, groups on a social network, etc.



- Components: number of latent groups
- Clusters: number of components represented in the data
- [demo 1, demo 2]

# So far $K \ll N$ . What if not?

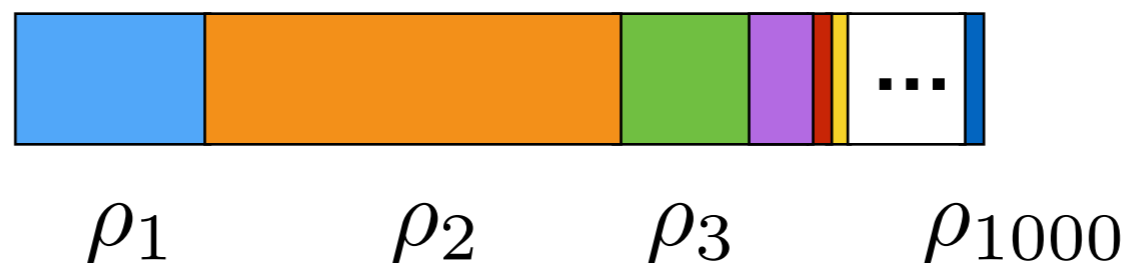
- e.g. species sampling, topic modeling, groups on a social network, etc.



- Components: number of latent groups
- Clusters: number of components represented in the data
- [demo 1, demo 2]
- Number of clusters is random

# So far $K \ll N$ . What if not?

- e.g. species sampling, topic modeling, groups on a social network, etc.



- Components: number of latent groups
- Clusters: number of components represented in the data
- [demo 1, demo 2]
- Number of clusters is random
- Number of clusters grows with  $N$

- Here, difficult to choose finite  $K$  in advance (contrast with small  $K$ ): don't know  $K$ , difficult to infer, streaming data

# Choosing $K = \infty$

- Here, difficult to choose finite  $K$  in advance (contrast with small  $K$ ): don't know  $K$ , difficult to infer, streaming data



# Choosing $K = \infty$

- Here, difficult to choose finite  $K$  in advance (contrast with small  $K$ ): don't know  $K$ , difficult to infer, streaming data
- How to generate  $K = \infty$  strictly positive frequencies that sum to one?

# Choosing $K = \infty$

- Here, difficult to choose finite  $K$  in advance (contrast with small  $K$ ): don't know  $K$ , difficult to infer, streaming data
- How to generate  $K = \infty$  strictly positive frequencies that sum to one?
  - Observation:  $\rho_{1:K} \sim \text{Dirichlet}(a_{1:K})$

# Choosing $K = \infty$

- Here, difficult to choose finite  $K$  in advance (contrast with small  $K$ ): don't know  $K$ , difficult to infer, streaming data
- How to generate  $K = \infty$  strictly positive frequencies that sum to one?
  - Observation:  $\rho_{1:K} \sim \text{Dirichlet}(a_{1:K})$



# Choosing $K = \infty$

- Here, difficult to choose finite  $K$  in advance (contrast with small  $K$ ): don't know  $K$ , difficult to infer, streaming data
- How to generate  $K = \infty$  strictly positive frequencies that sum to one?
  - Observation:  $\rho_{1:K} \sim \text{Dirichlet}(a_{1:K})$

$$\Leftrightarrow \rho_1 \stackrel{d}{=} \text{Beta}\left(a_1, \sum_{k=1}^K a_k - a_1\right)$$



# Choosing $K = \infty$

- Here, difficult to choose finite  $K$  in advance (contrast with small  $K$ ): don't know  $K$ , difficult to infer, streaming data
- How to generate  $K = \infty$  strictly positive frequencies that sum to one?
  - Observation:  $\rho_{1:K} \sim \text{Dirichlet}(a_{1:K})$

$$\Leftrightarrow \rho_1 \stackrel{d}{=} \text{Beta}\left(a_1, \sum_{k=1}^K a_k - a_1\right)$$



# Choosing $K = \infty$

- Here, difficult to choose finite  $K$  in advance (contrast with small  $K$ ): don't know  $K$ , difficult to infer, streaming data
- How to generate  $K = \infty$  strictly positive frequencies that sum to one?
  - Observation:  $\rho_{1:K} \sim \text{Dirichlet}(a_{1:K})$

$$\Leftrightarrow \rho_1 \stackrel{d}{=} \text{Beta}\left(a_1, \sum_{k=1}^K a_k - a_1\right) \perp \frac{(\rho_2, \dots, \rho_K)}{1 - \rho_1} \stackrel{d}{=} \text{Dirichlet}(a_2, \dots, a_K)$$



# Choosing $K = \infty$

- Here, difficult to choose finite  $K$  in advance (contrast with small  $K$ ): don't know  $K$ , difficult to infer, streaming data
- How to generate  $K = \infty$  strictly positive frequencies that sum to one?
  - Observation:  $\rho_{1:K} \sim \text{Dirichlet}(a_{1:K})$

$$\Leftrightarrow \rho_1 \stackrel{d}{=} \text{Beta}(a_1, \sum_{k=1}^K a_k - a_1) \perp\!\!\!\perp \frac{(\rho_2, \dots, \rho_K)}{1 - \rho_1} \stackrel{d}{=} \text{Dirichlet}(a_2, \dots, a_K)$$



# Choosing $K = \infty$

- Here, difficult to choose finite  $K$  in advance (contrast with small  $K$ ): don't know  $K$ , difficult to infer, streaming data
- How to generate  $K = \infty$  strictly positive frequencies that sum to one?
  - Observation:  $\rho_{1:K} \sim \text{Dirichlet}(a_{1:K})$

$$\Leftrightarrow \rho_1 \stackrel{d}{=} \text{Beta}\left(a_1, \sum_{k=1}^K a_k - a_1\right) \perp\!\!\!\perp \frac{(\rho_2, \dots, \rho_K)}{1 - \rho_1} \stackrel{d}{=} \text{Dirichlet}(a_2, \dots, a_K)$$





# Choosing $K = \infty$

- Here, difficult to choose finite  $K$  in advance (contrast with small  $K$ ): don't know  $K$ , difficult to infer, streaming data
- How to generate  $K = \infty$  strictly positive frequencies that sum to one?
  - Observation:  $\rho_{1:K} \sim \text{Dirichlet}(a_{1:K})$

$$\Leftrightarrow \rho_1 \stackrel{d}{=} \text{Beta}\left(a_1, \sum_{k=1}^K a_k - a_1\right) \perp\!\!\!\perp \frac{(\rho_2, \dots, \rho_K)}{1 - \rho_1} \stackrel{d}{=} \text{Dirichlet}(a_2, \dots, a_K)$$

- 
- “Stick breaking”

# Choosing $K = \infty$

- Here, difficult to choose finite  $K$  in advance (contrast with small  $K$ ): don't know  $K$ , difficult to infer, streaming data
- How to generate  $K = \infty$  strictly positive frequencies that sum to one?
  - Observation:  $\rho_{1:K} \sim \text{Dirichlet}(a_{1:K})$

$$\Leftrightarrow \rho_1 \stackrel{d}{=} \text{Beta}\left(a_1, \sum_{k=2}^K a_k\right) \perp\!\!\!\perp \frac{(\rho_2, \dots, \rho_K)}{1 - \rho_1} \stackrel{d}{=} \text{Dirichlet}(a_2, \dots, a_K)$$

- 
- “Stick breaking”

$$V_1 \sim \text{Beta}(a_1, a_2 + a_3 + a_4)$$

# Choosing $K = \infty$

- Here, difficult to choose finite  $K$  in advance (contrast with small  $K$ ): don't know  $K$ , difficult to infer, streaming data
- How to generate  $K = \infty$  strictly positive frequencies that sum to one?
  - Observation:  $\rho_{1:K} \sim \text{Dirichlet}(a_{1:K})$

$$\Leftrightarrow \rho_1 \stackrel{d}{=} \text{Beta}\left(a_1, \sum_{k=2}^K a_k\right) \perp\!\!\!\perp \frac{(\rho_2, \dots, \rho_K)}{1 - \rho_1} \stackrel{d}{=} \text{Dirichlet}(a_2, \dots, a_K)$$



- “Stick breaking”

$$V_1 \sim \text{Beta}(a_1, a_2 + a_3 + a_4)$$

$$\rho_1 = V_1$$

# Choosing $K = \infty$

- Here, difficult to choose finite  $K$  in advance (contrast with small  $K$ ): don't know  $K$ , difficult to infer, streaming data
- How to generate  $K = \infty$  strictly positive frequencies that sum to one?
  - Observation:  $\rho_{1:K} \sim \text{Dirichlet}(a_{1:K})$

$$\Leftrightarrow \rho_1 \stackrel{d}{=} \text{Beta}\left(a_1, \sum_{k=2}^K a_k\right) \perp\!\!\!\perp \frac{(\rho_2, \dots, \rho_K)}{1 - \rho_1} \stackrel{d}{=} \text{Dirichlet}(a_2, \dots, a_K)$$



- “Stick breaking”

$$V_1 \sim \text{Beta}(a_1, a_2 + a_3 + a_4)$$

$$\rho_1 = V_1$$

$$V_2 \sim \text{Beta}(a_2, a_3 + a_4)$$

# Choosing $K = \infty$

- Here, difficult to choose finite  $K$  in advance (contrast with small  $K$ ): don't know  $K$ , difficult to infer, streaming data
- How to generate  $K = \infty$  strictly positive frequencies that sum to one?
  - Observation:  $\rho_{1:K} \sim \text{Dirichlet}(a_{1:K})$

$$\Leftrightarrow \rho_1 \stackrel{d}{=} \text{Beta}\left(a_1, \sum_{k=2}^K a_k\right) \perp\!\!\!\perp \frac{(\rho_2, \dots, \rho_K)}{1 - \rho_1} \stackrel{d}{=} \text{Dirichlet}(a_2, \dots, a_K)$$



- “Stick breaking”

$$V_1 \sim \text{Beta}(a_1, a_2 + a_3 + a_4)$$

$$\rho_1 = V_1$$

$$V_2 \sim \text{Beta}(a_2, a_3 + a_4)$$

$$\rho_2 = (1 - V_1)V_2$$

# Choosing $K = \infty$

- Here, difficult to choose finite  $K$  in advance (contrast with small  $K$ ): don't know  $K$ , difficult to infer, streaming data
- How to generate  $K = \infty$  strictly positive frequencies that sum to one?
  - Observation:  $\rho_{1:K} \sim \text{Dirichlet}(a_{1:K})$

$$\Leftrightarrow \rho_1 \stackrel{d}{=} \text{Beta}\left(a_1, \sum_{k=2}^K a_k\right) \perp\!\!\!\perp \frac{(\rho_2, \dots, \rho_K)}{1 - \rho_1} \stackrel{d}{=} \text{Dirichlet}(a_2, \dots, a_K)$$



- “Stick breaking”

$$V_1 \sim \text{Beta}(a_1, a_2 + a_3 + a_4) \quad \rho_1 = V_1$$

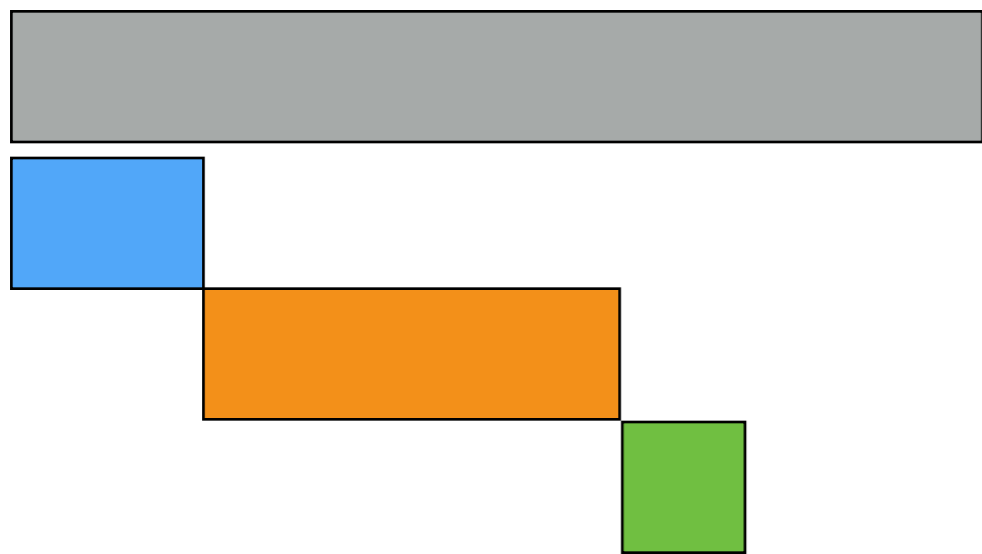
$$V_2 \sim \text{Beta}(a_2, a_3 + a_4) \quad \rho_2 = (1 - V_1)V_2$$

$$V_3 \sim \text{Beta}(a_3, a_4)$$

# Choosing $K = \infty$

- Here, difficult to choose finite  $K$  in advance (contrast with small  $K$ ): don't know  $K$ , difficult to infer, streaming data
- How to generate  $K = \infty$  strictly positive frequencies that sum to one?
  - Observation:  $\rho_{1:K} \sim \text{Dirichlet}(a_{1:K})$

$$\Leftrightarrow \rho_1 \stackrel{d}{=} \text{Beta}\left(a_1, \sum_{k=2}^K a_k\right) \perp \frac{(\rho_2, \dots, \rho_K)}{1 - \rho_1} \stackrel{d}{=} \text{Dirichlet}(a_2, \dots, a_K)$$



- “Stick breaking”

$$V_1 \sim \text{Beta}(a_1, a_2 + a_3 + a_4) \quad \rho_1 = V_1$$

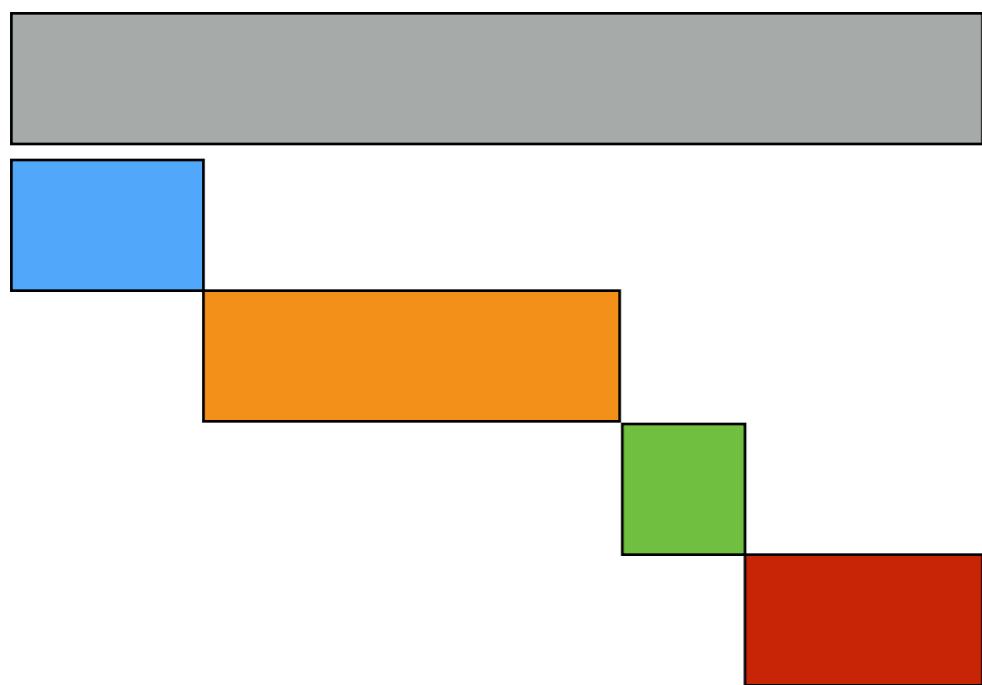
$$V_2 \sim \text{Beta}(a_2, a_3 + a_4) \quad \rho_2 = (1 - V_1)V_2$$

$$V_3 \sim \text{Beta}(a_3, a_4) \quad \rho_3 = (1 - V_1)(1 - V_2)V_3$$

# Choosing $K = \infty$

- Here, difficult to choose finite  $K$  in advance (contrast with small  $K$ ): don't know  $K$ , difficult to infer, streaming data
- How to generate  $K = \infty$  strictly positive frequencies that sum to one?
  - Observation:  $\rho_{1:K} \sim \text{Dirichlet}(a_{1:K})$

$$\Leftrightarrow \rho_1 \stackrel{d}{=} \text{Beta}\left(a_1, \sum_{k=2}^K a_k\right) \perp\!\!\!\perp \frac{(\rho_2, \dots, \rho_K)}{1 - \rho_1} \stackrel{d}{=} \text{Dirichlet}(a_2, \dots, a_K)$$



- “Stick breaking”

$$V_1 \sim \text{Beta}(a_1, a_2 + a_3 + a_4) \quad \rho_1 = V_1$$

$$V_2 \sim \text{Beta}(a_2, a_3 + a_4) \quad \rho_2 = (1 - V_1)V_2$$

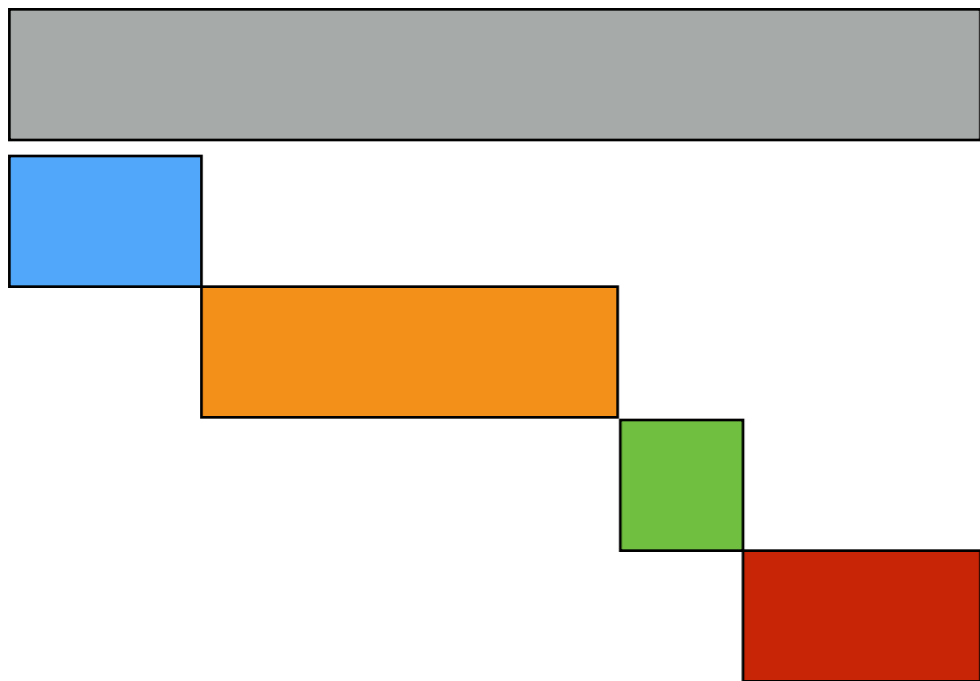
$$V_3 \sim \text{Beta}(a_3, a_4) \quad \rho_3 = (1 - V_1)(1 - V_2)V_3$$

$$\rho_4 = 1 - \sum_{k=1}^3 \rho_k$$



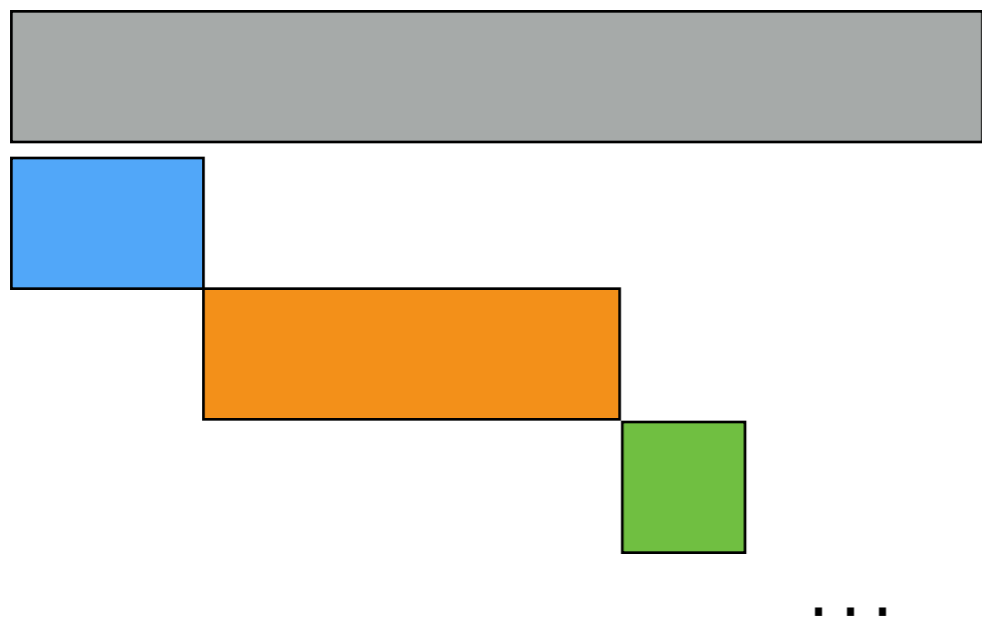
# Choosing $K = \infty$

- Here, difficult to choose finite  $K$  in advance (contrast with small  $K$ ): don't know  $K$ , difficult to infer, streaming data
- How to generate  $K = \infty$  strictly positive frequencies that sum to one?



# Choosing $K = \infty$

- Here, difficult to choose finite  $K$  in advance (contrast with small  $K$ ): don't know  $K$ , difficult to infer, streaming data
- How to generate  $K = \infty$  strictly positive frequencies that sum to one?



# Choosing $K = \infty$

- Here, difficult to choose finite  $K$  in advance (contrast with small  $K$ ): don't know  $K$ , difficult to infer, streaming data
- How to generate  $K = \infty$  strictly positive frequencies that sum to one?



# Choosing $K = \infty$

- Here, difficult to choose finite  $K$  in advance (contrast with small  $K$ ): don't know  $K$ , difficult to infer, streaming data
- How to generate  $K = \infty$  strictly positive frequencies that sum to one?



$$V_1 \sim \text{Beta}(a_1, b_1)$$

# Choosing $K = \infty$

- Here, difficult to choose finite  $K$  in advance (contrast with small  $K$ ): don't know  $K$ , difficult to infer, streaming data
- How to generate  $K = \infty$  strictly positive frequencies that sum to one?



$$V_1 \sim \text{Beta}(a_1, b_1) \quad \rho_1 = V_1$$

# Choosing $K = \infty$

- Here, difficult to choose finite  $K$  in advance (contrast with small  $K$ ): don't know  $K$ , difficult to infer, streaming data
- How to generate  $K = \infty$  strictly positive frequencies that sum to one?



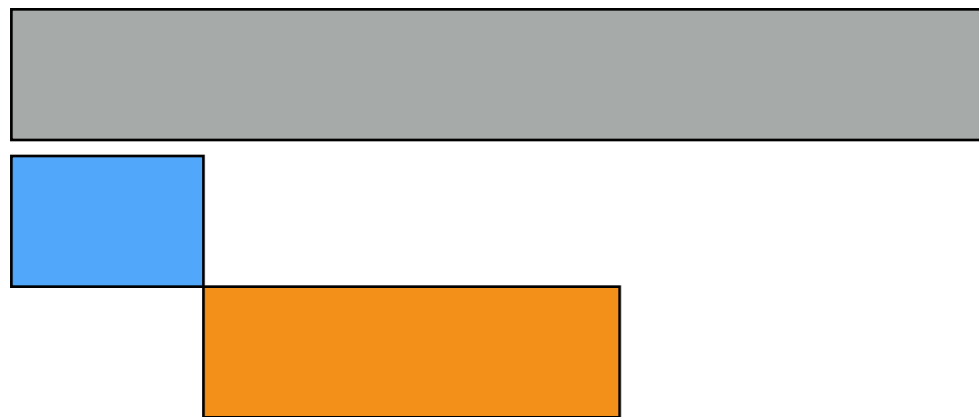
$$V_1 \sim \text{Beta}(a_1, b_1)$$

$$\rho_1 = V_1$$

$$V_2 \sim \text{Beta}(a_2, b_2)$$

# Choosing $K = \infty$

- Here, difficult to choose finite  $K$  in advance (contrast with small  $K$ ): don't know  $K$ , difficult to infer, streaming data
- How to generate  $K = \infty$  strictly positive frequencies that sum to one?



$$V_1 \sim \text{Beta}(a_1, b_1)$$

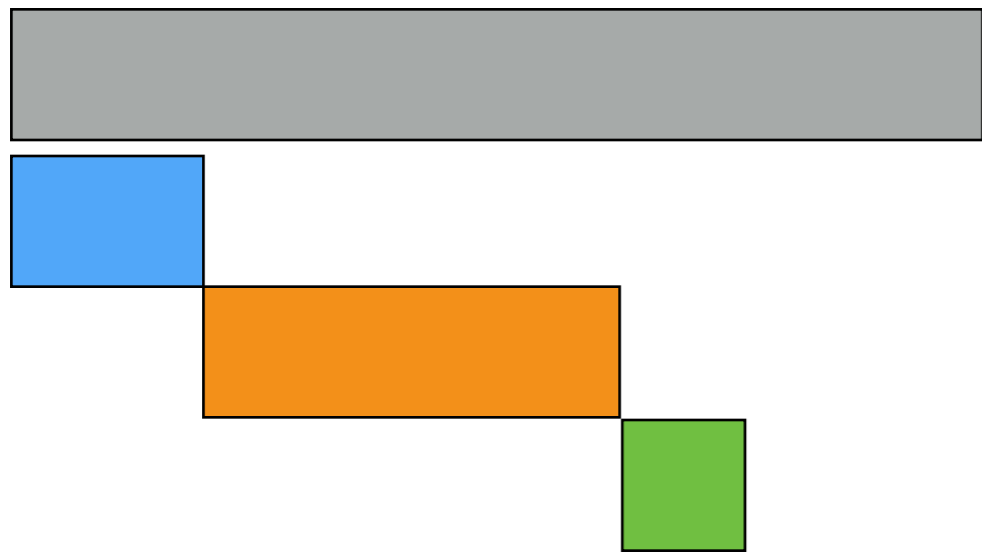
$$\rho_1 = V_1$$

$$V_2 \sim \text{Beta}(a_2, b_2)$$

$$\rho_2 = (1 - V_1)V_2$$

# Choosing $K = \infty$

- Here, difficult to choose finite  $K$  in advance (contrast with small  $K$ ): don't know  $K$ , difficult to infer, streaming data
- How to generate  $K = \infty$  strictly positive frequencies that sum to one?



$$V_1 \sim \text{Beta}(a_1, b_1)$$

$$\rho_1 = V_1$$

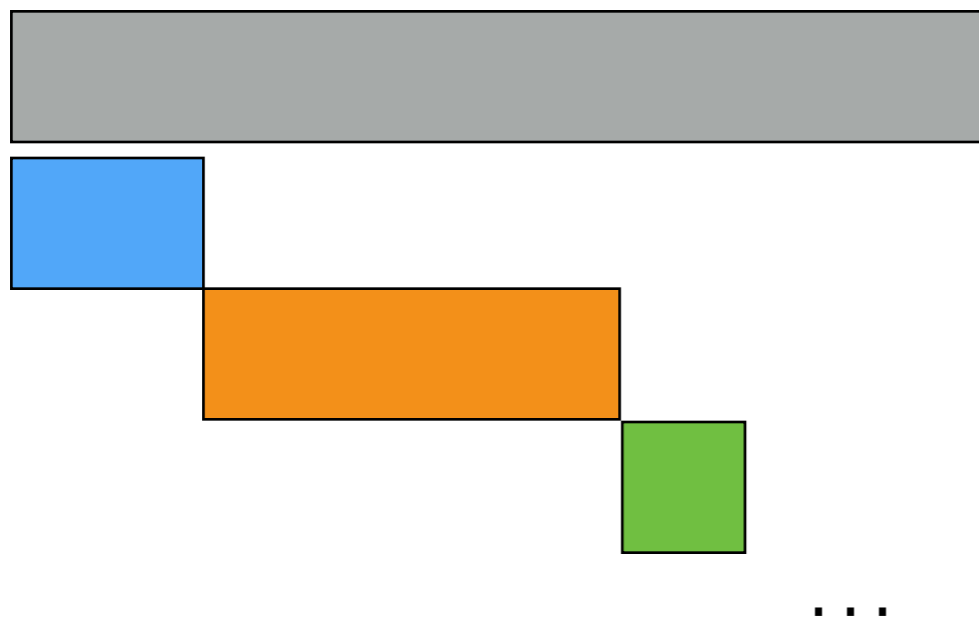
$$V_2 \sim \text{Beta}(a_2, b_2)$$

$$\rho_2 = (1 - V_1)V_2$$



# Choosing $K = \infty$

- Here, difficult to choose finite  $K$  in advance (contrast with small  $K$ ): don't know  $K$ , difficult to infer, streaming data
- How to generate  $K = \infty$  strictly positive frequencies that sum to one?



$$V_1 \sim \text{Beta}(a_1, b_1)$$

$$\rho_1 = V_1$$

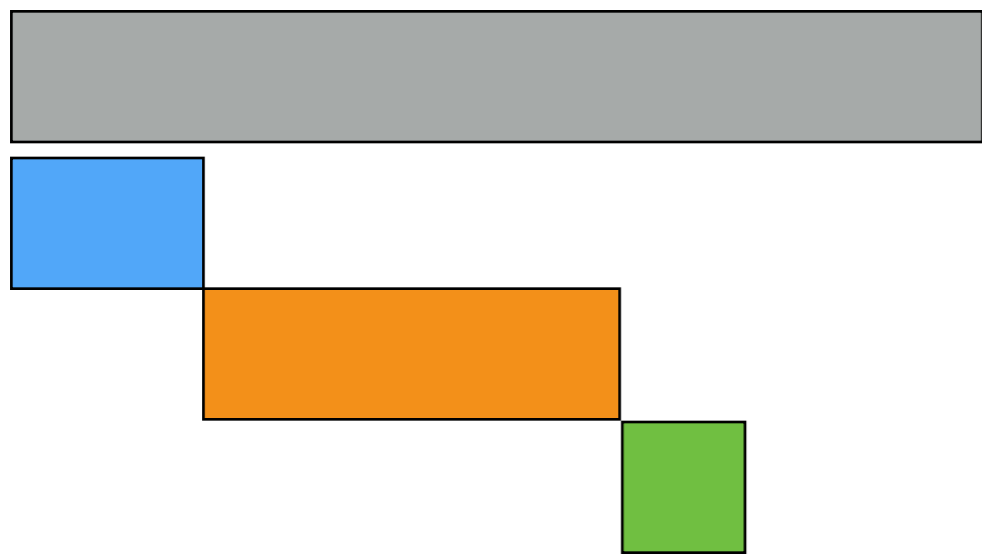
$$V_2 \sim \text{Beta}(a_2, b_2)$$

$$\rho_2 = (1 - V_1)V_2$$

...

# Choosing $K = \infty$

- Here, difficult to choose finite  $K$  in advance (contrast with small  $K$ ): don't know  $K$ , difficult to infer, streaming data
- How to generate  $K = \infty$  strictly positive frequencies that sum to one?



$$V_1 \sim \text{Beta}(a_1, b_1)$$

$$\rho_1 = V_1$$

$$V_2 \sim \text{Beta}(a_2, b_2)$$

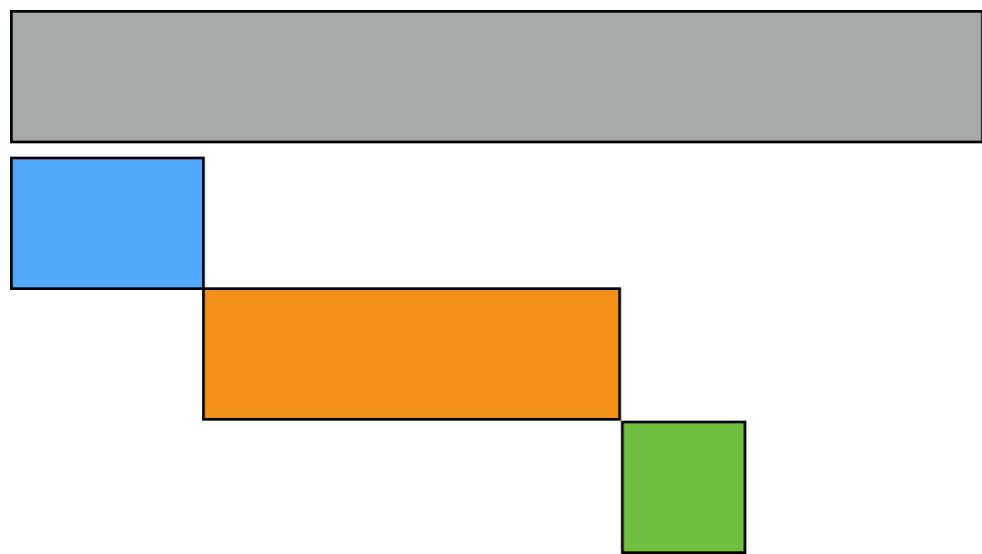
$$\rho_2 = (1 - V_1)V_2$$

...

$$V_k \sim \text{Beta}(a_k, b_k)$$

# Choosing $K = \infty$

- Here, difficult to choose finite  $K$  in advance (contrast with small  $K$ ): don't know  $K$ , difficult to infer, streaming data
- How to generate  $K = \infty$  strictly positive frequencies that sum to one?



$$V_1 \sim \text{Beta}(a_1, b_1)$$

$$\rho_1 = V_1$$

$$V_2 \sim \text{Beta}(a_2, b_2)$$

$$\rho_2 = (1 - V_1)V_2$$

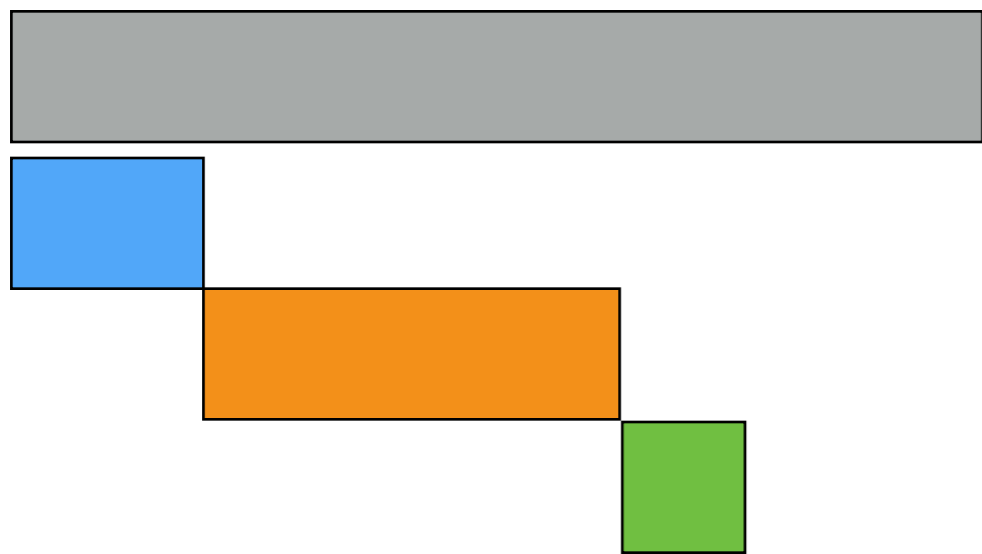
...

$$V_k \sim \text{Beta}(a_k, b_k)$$

$$\rho_k = \left[ \prod_{j=1}^{k-1} (1 - V_j) \right] V_k$$

# Choosing $K = \infty$

- Here, difficult to choose finite  $K$  in advance (contrast with small  $K$ ): don't know  $K$ , difficult to infer, streaming data
- How to generate  $K = \infty$  strictly positive frequencies that sum to one?



$$V_1 \sim \text{Beta}(a_1, b_1)$$

$$V_2 \sim \text{Beta}(a_2, b_2)$$

$$\dots \quad V_k \sim \text{Beta}(a_k, b_k)$$

$$\rho_1 = V_1$$

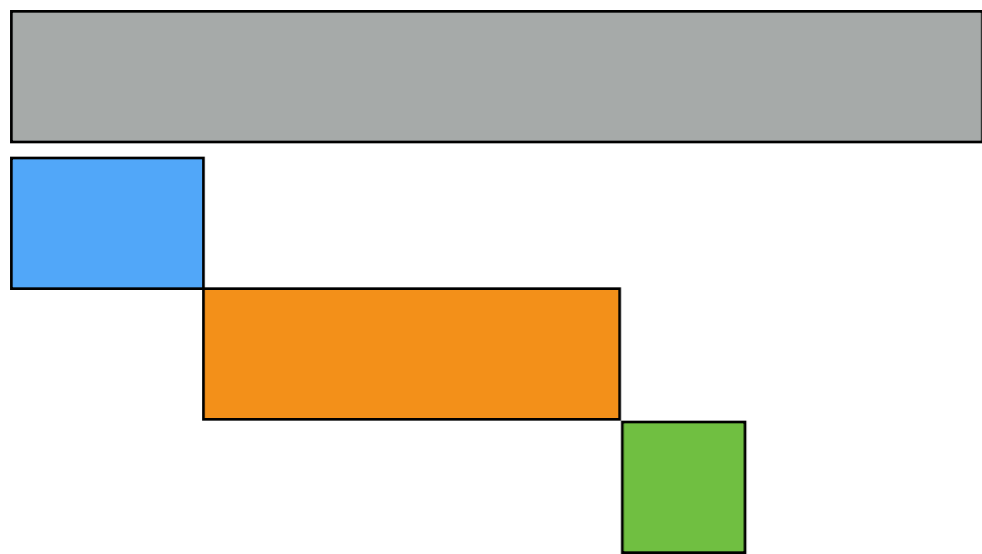
$$\rho_2 = (1 - V_1)V_2$$

$$\rho_k = \left[ \prod_{j=1}^{k-1} (1 - V_j) \right] V_k$$

[Ishwaran, James 2001]

# Choosing $K = \infty$

- Here, difficult to choose finite  $K$  in advance (contrast with small  $K$ ): don't know  $K$ , difficult to infer, streaming data
- How to generate  $K = \infty$  strictly positive frequencies that sum to one?
  - **Dirichlet process stick-breaking:**  $a_k = 1, b_k = \alpha > 0$



$$V_1 \sim \text{Beta}(a_1, b_1)$$

$$\rho_1 = V_1$$

$$V_2 \sim \text{Beta}(a_2, b_2)$$

$$\rho_2 = (1 - V_1)V_2$$

...

$$V_k \sim \text{Beta}(a_k, b_k)$$

$$\rho_k = \left[ \prod_{j=1}^{k-1} (1 - V_j) \right] V_k$$

[Ishwaran, James 2001]

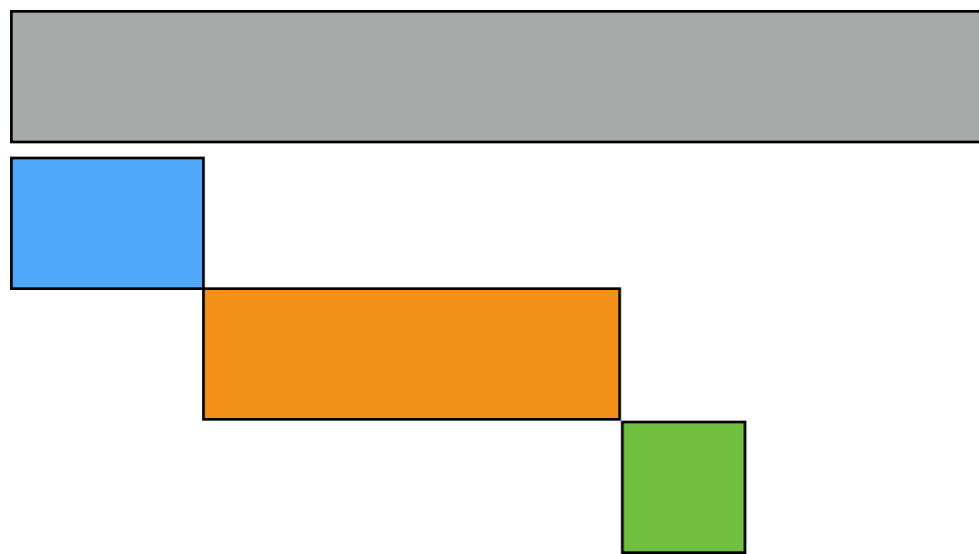
# Choosing $K = \infty$

- Here, difficult to choose finite  $K$  in advance (contrast with small  $K$ ): don't know  $K$ , difficult to infer, streaming data
- How to generate  $K = \infty$  strictly positive frequencies that sum to one?

- **Dirichlet process stick-breaking**:  $a_k = 1, b_k = \alpha > 0$

- Griffiths-Engen-McCloskey (**GEM**) distribution:

$$\rho = (\rho_1, \rho_2, \dots) \sim \text{GEM}(\alpha)$$



$$V_1 \sim \text{Beta}(a_1, b_1)$$

$$\rho_1 = V_1$$

$$V_2 \sim \text{Beta}(a_2, b_2)$$

$$\rho_2 = (1 - V_1)V_2$$

...

$$V_k \sim \text{Beta}(a_k, b_k)$$

$$\rho_k = \left[ \prod_{j=1}^{k-1} (1 - V_j) \right] V_k$$

# Choosing $K = \infty$

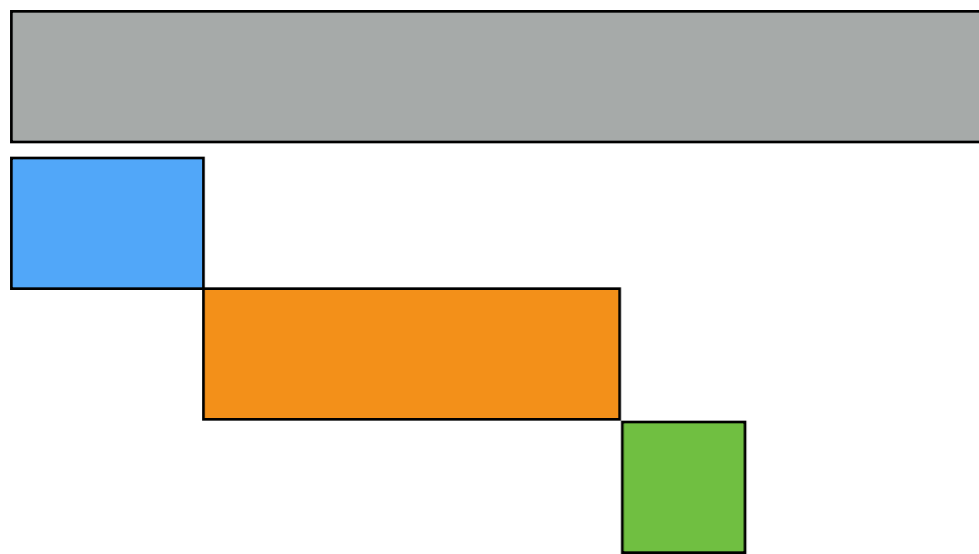
- Here, difficult to choose finite  $K$  in advance (contrast with small  $K$ ): don't know  $K$ , difficult to infer, streaming data
- How to generate  $K = \infty$  strictly positive frequencies that sum to one?

- **Dirichlet process stick-breaking**:  $a_k = 1, b_k = \alpha > 0$

- Griffiths-Engen-McCloskey (**GEM**) distribution:

$$\rho = (\rho_1, \rho_2, \dots) \sim \text{GEM}(\alpha)$$

[demo]



$$V_1 \sim \text{Beta}(a_1, b_1)$$

$$\rho_1 = V_1$$

$$V_2 \sim \text{Beta}(a_2, b_2)$$

$$\rho_2 = (1 - V_1)V_2$$

...

$$V_k \sim \text{Beta}(a_k, b_k)$$

$$\rho_k = \left[ \prod_{j=1}^{k-1} (1 - V_j) \right] V_k$$

# Dirichlet process mixture model



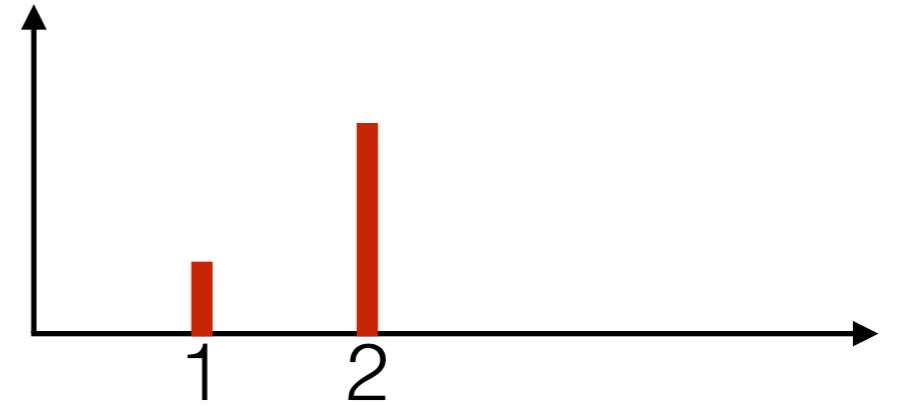
# Dirichlet process mixture model

- Gaussian mixture model

# Dirichlet process mixture model

- Gaussian mixture model

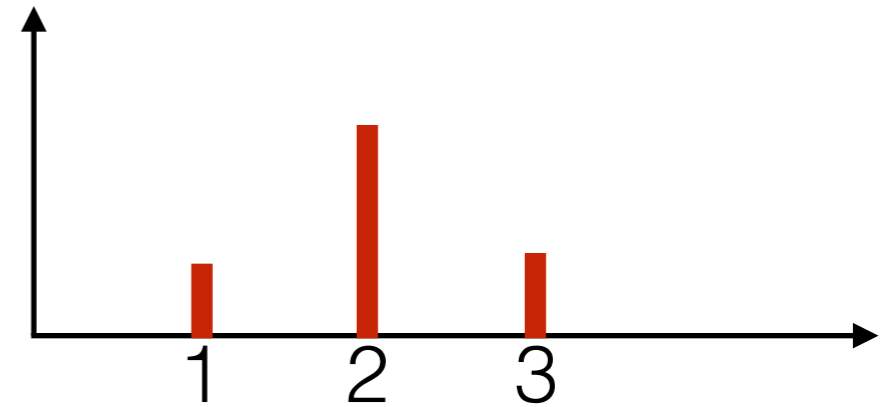
$$\rho = (\rho_1, \rho_2) \sim \text{Beta}(a_1, a_2)$$



# Dirichlet process mixture model

- Gaussian mixture model

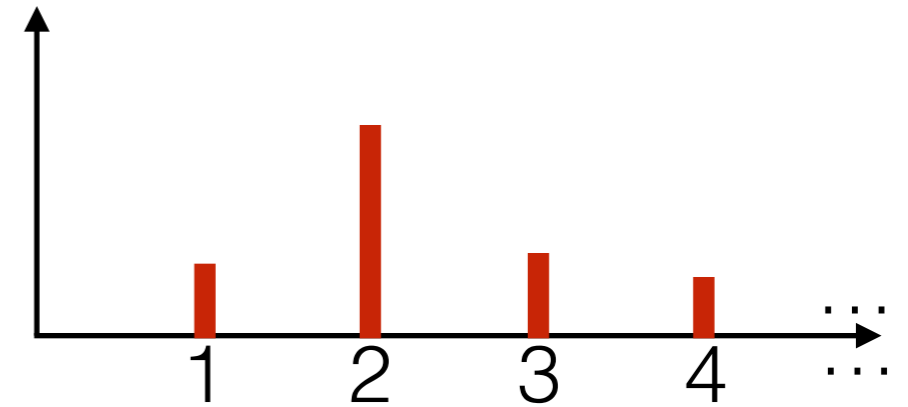
$$\rho = (\rho_1, \dots, \rho_K) \sim \text{Dir}(a_{1:K})$$



# Dirichlet process mixture model

- Gaussian mixture model

$$\rho = (\rho_1, \rho_2, \dots) \sim \text{GEM}(\alpha)$$

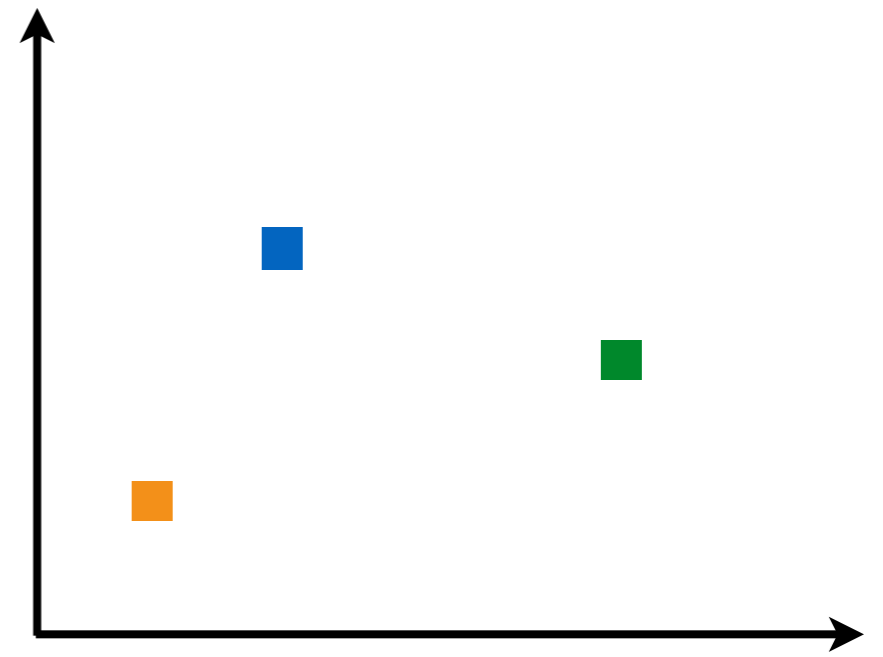
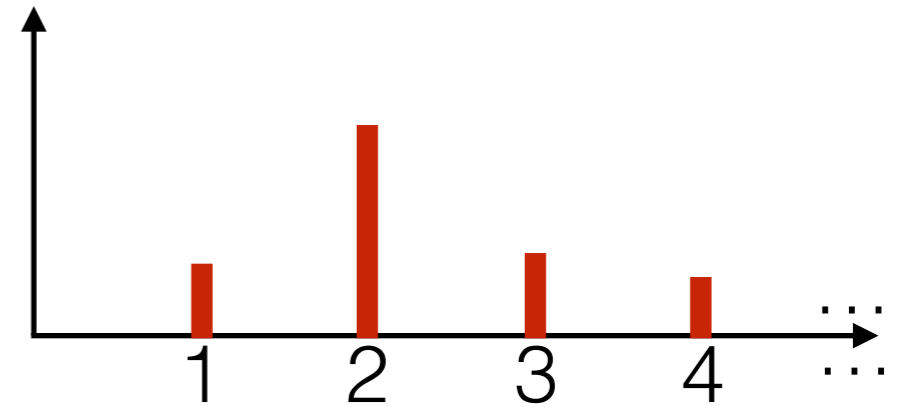


# Dirichlet process mixture model

- Gaussian mixture model

$$\rho = (\rho_1, \rho_2, \dots) \sim \text{GEM}(\alpha)$$

$$\mu_k \stackrel{iid}{\sim} \mathcal{N}(\mu_0, \Sigma_0), k = 1, 2, \dots$$

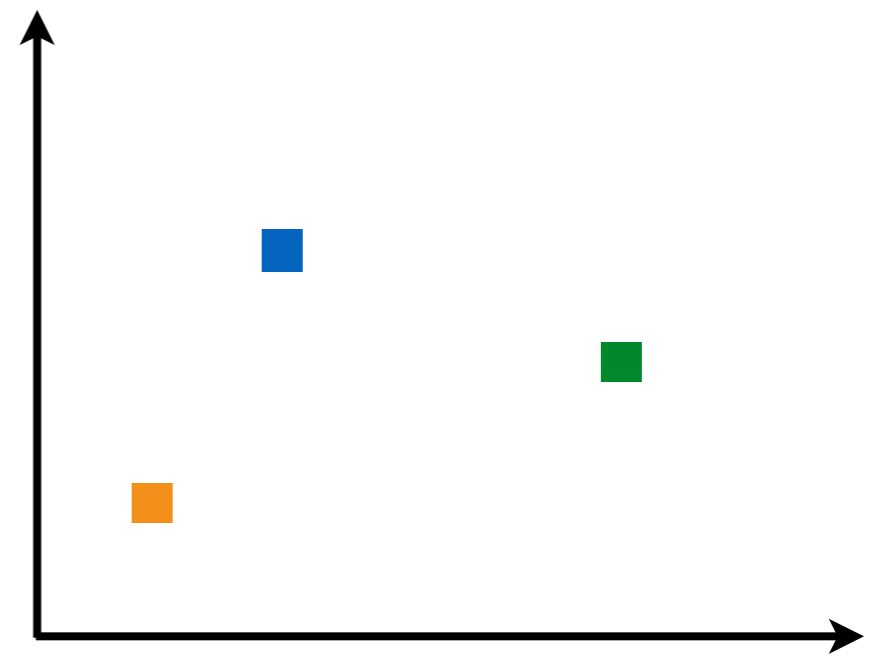
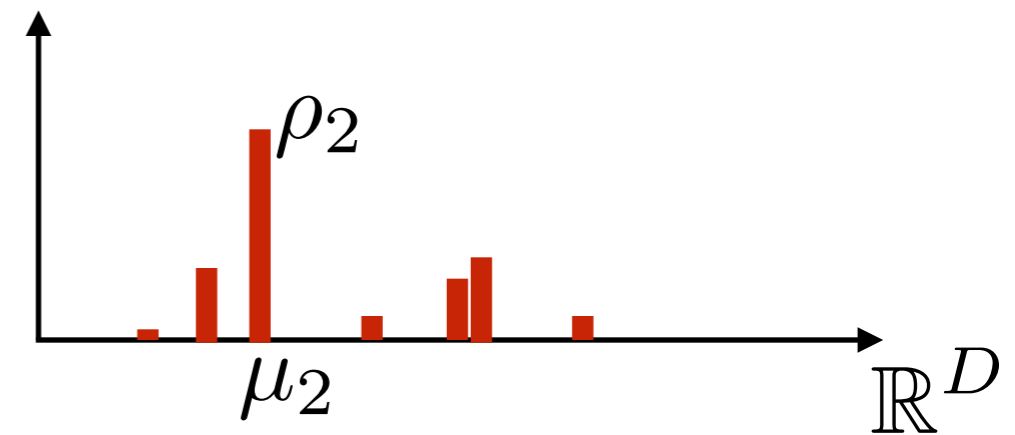
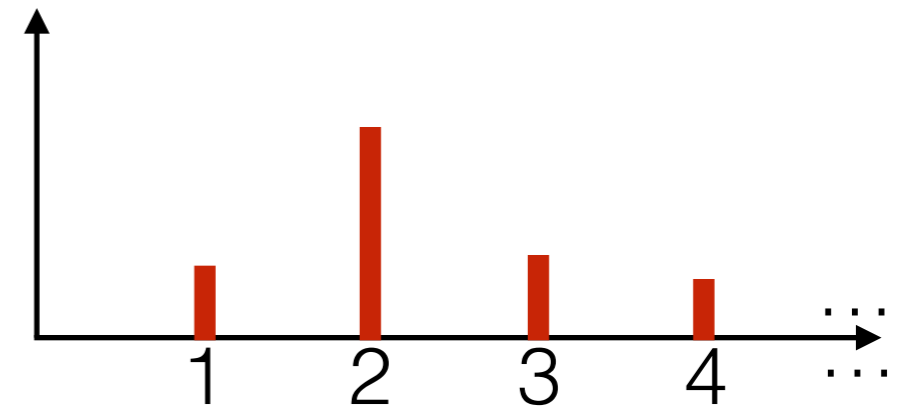


# Dirichlet process mixture model

- Gaussian mixture model

$$\rho = (\rho_1, \rho_2, \dots) \sim \text{GEM}(\alpha)$$

$$\mu_k \stackrel{iid}{\sim} \mathcal{N}(\mu_0, \Sigma_0), k = 1, 2, \dots$$



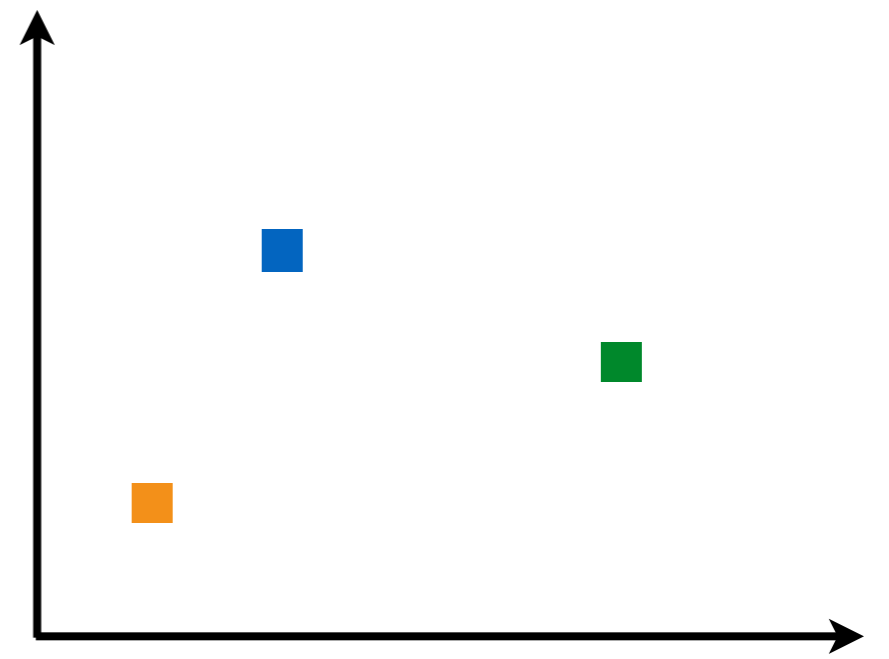
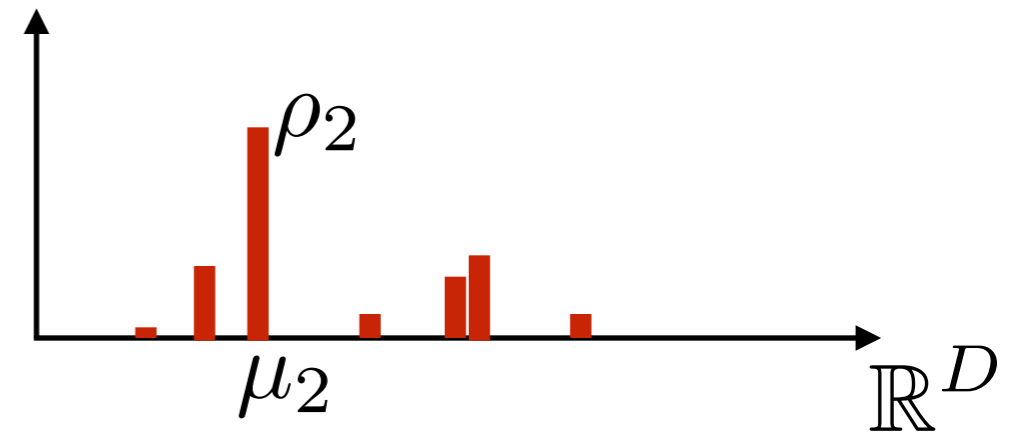
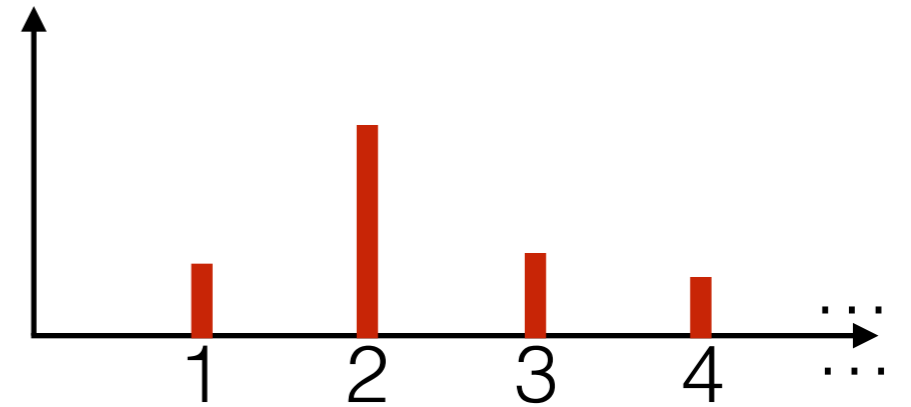
# Dirichlet process mixture model

- Gaussian mixture model

$$\rho = (\rho_1, \rho_2, \dots) \sim \text{GEM}(\alpha)$$

$$\mu_k \stackrel{iid}{\sim} \mathcal{N}(\mu_0, \Sigma_0), k = 1, 2, \dots$$

- i.e.  $G = \sum_{k=1}^{\infty} \rho_k \delta_{\mu_k}$



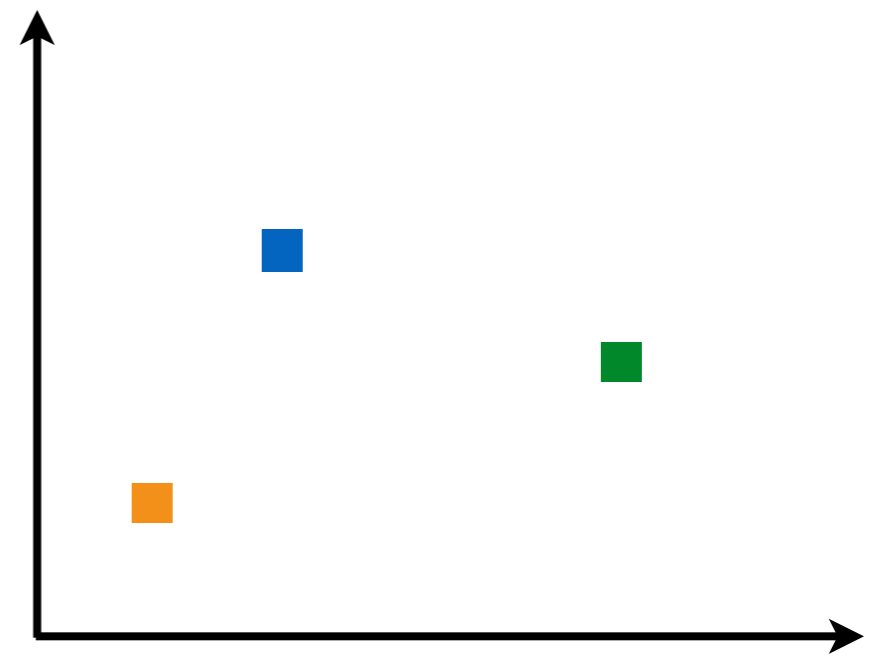
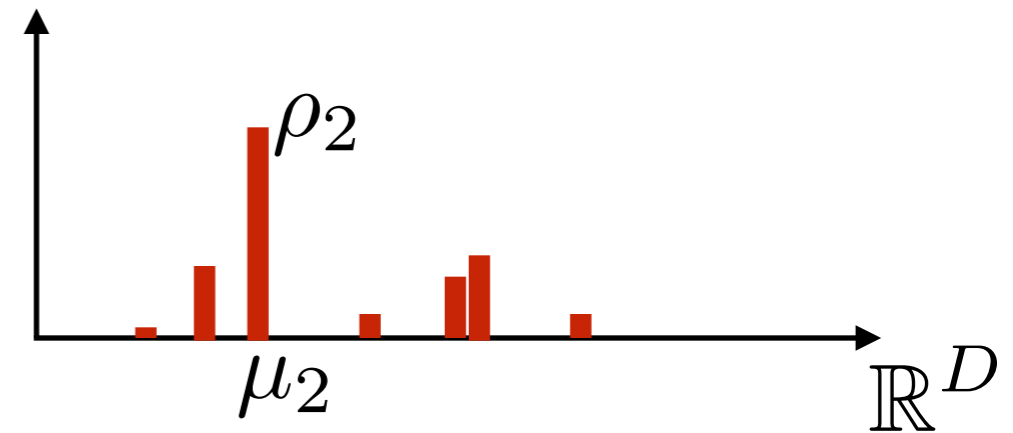
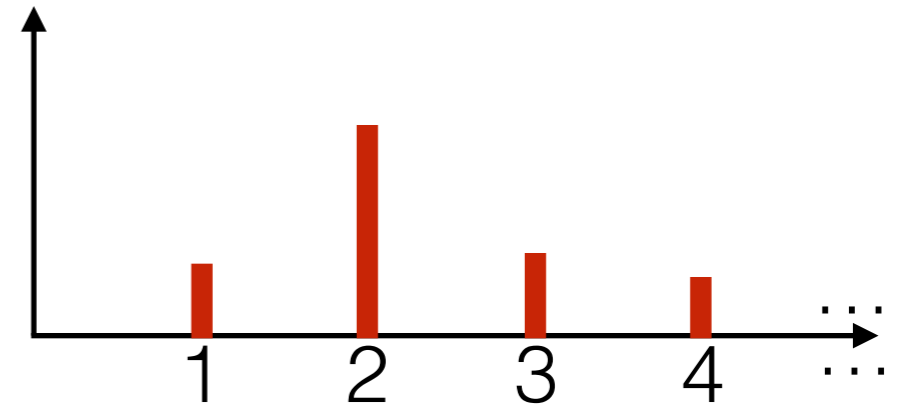
# Dirichlet process mixture model

- Gaussian mixture model

$$\rho = (\rho_1, \rho_2, \dots) \sim \text{GEM}(\alpha)$$

$$\mu_k \stackrel{iid}{\sim} \mathcal{N}(\mu_0, \Sigma_0), k = 1, 2, \dots$$

- i.e.  $G = \sum_{k=1}^{\infty} \rho_k \delta_{\mu_k} \stackrel{d}{=} \text{DP}(\alpha, \mathcal{N}(\mu_0, \Sigma_0))$





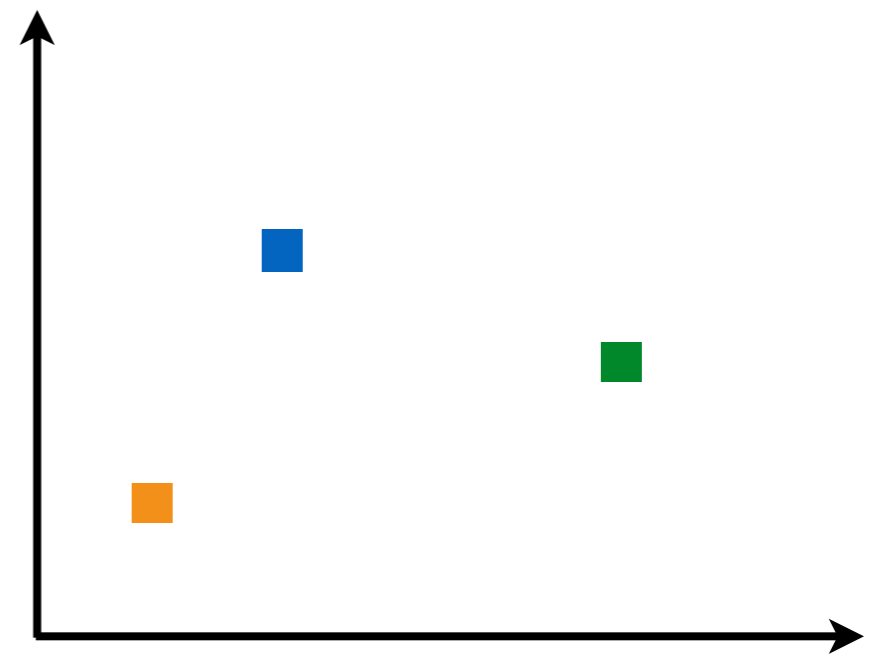
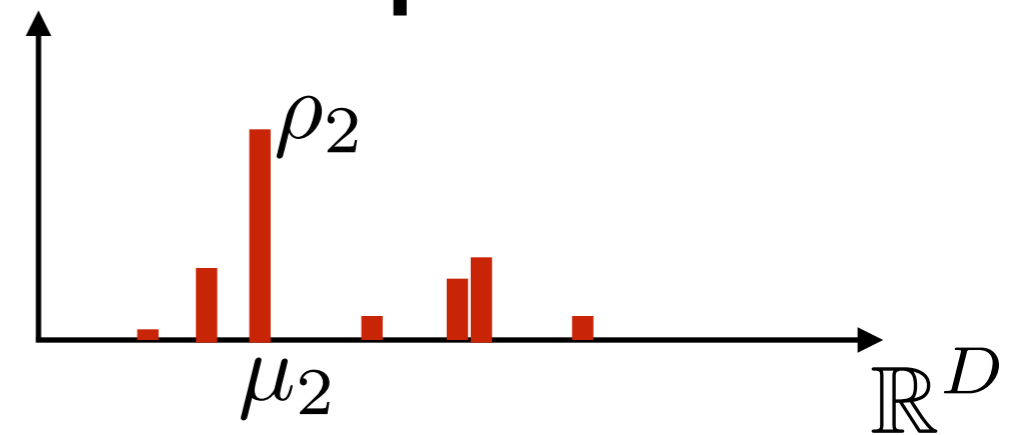
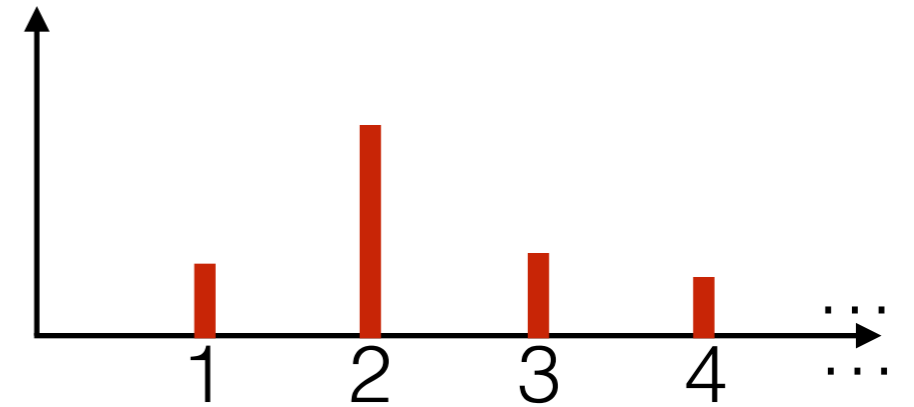
# Dirichlet process mixture model

- Gaussian mixture model

$$\rho = (\rho_1, \rho_2, \dots) \sim \text{GEM}(\alpha)$$

$$\mu_k \stackrel{iid}{\sim} \mathcal{N}(\mu_0, \Sigma_0), k = 1, 2, \dots$$

- i.e.  $G = \sum_{k=1}^{\infty} \rho_k \delta_{\mu_k} \stackrel{d}{=} \text{DP}(\alpha, \mathcal{N}(\mu_0, \Sigma_0))$  **Dirichlet process**



# Dirichlet process mixture model

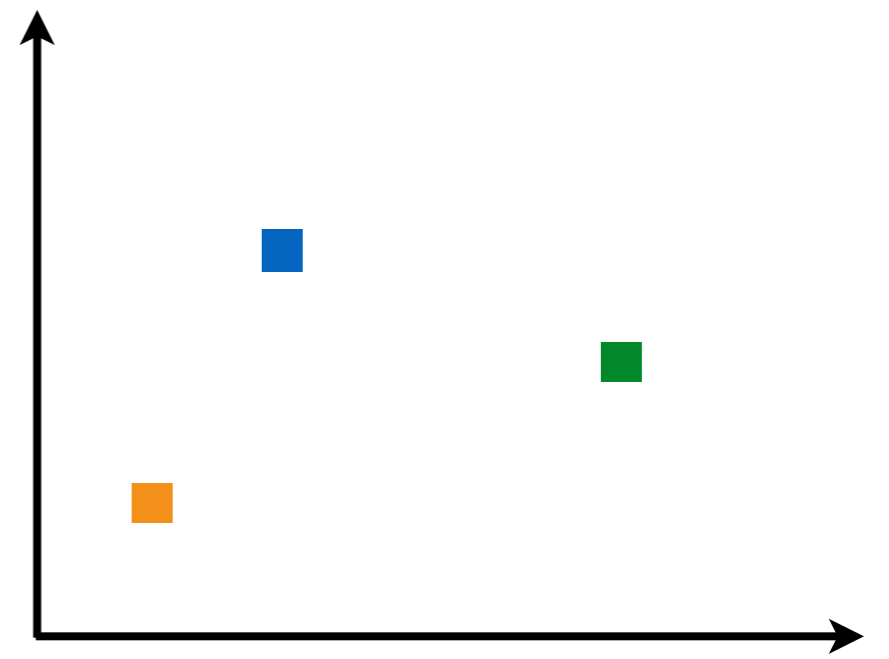
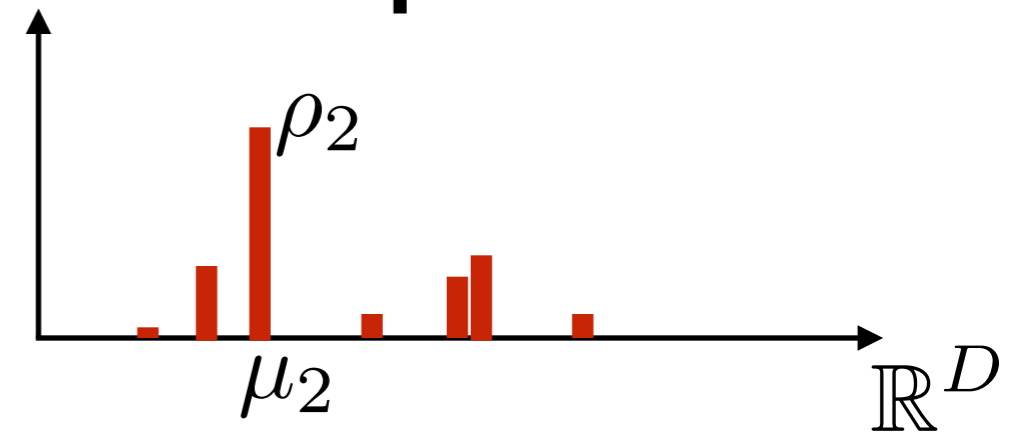
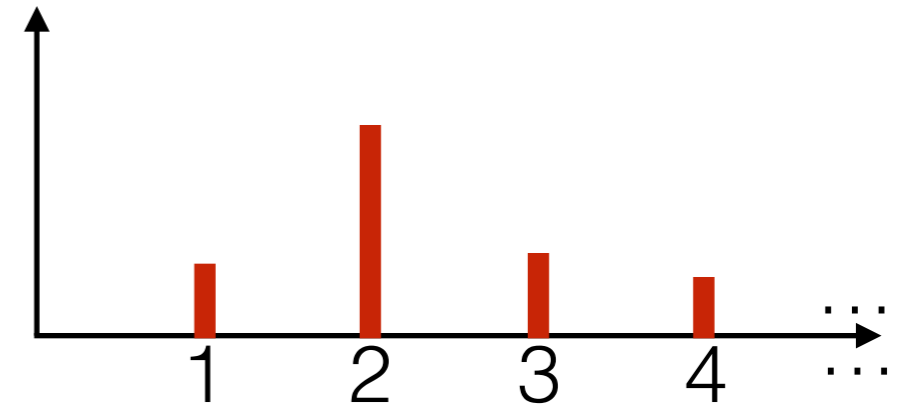
- Gaussian mixture model

$$\rho = (\rho_1, \rho_2, \dots) \sim \text{GEM}(\alpha)$$

$$\mu_k \stackrel{iid}{\sim} \mathcal{N}(\mu_0, \Sigma_0), k = 1, 2, \dots$$

- i.e.  $G = \sum_{k=1}^{\infty} \rho_k \delta_{\mu_k} \stackrel{d}{=} \text{DP}(\alpha, \mathcal{N}(\mu_0, \Sigma_0))$  **Dirichlet process**

$$z_n \stackrel{iid}{\sim} \text{Categorical}(\rho)$$



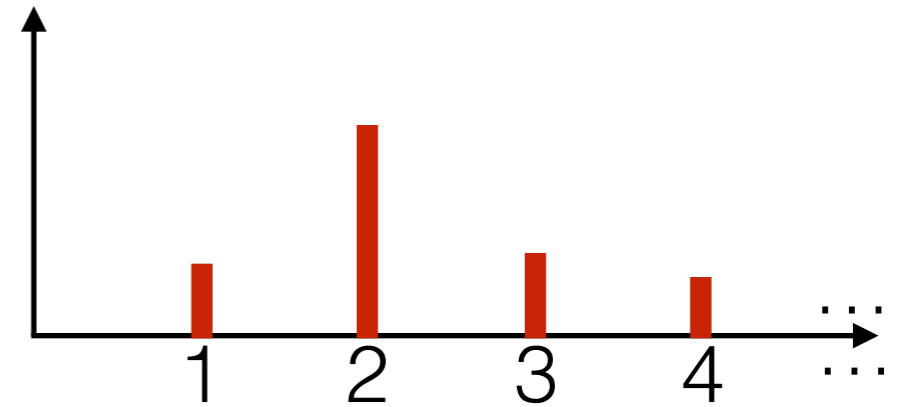
# Dirichlet process mixture model

- Gaussian mixture model

$$\rho = (\rho_1, \rho_2, \dots) \sim \text{GEM}(\alpha)$$

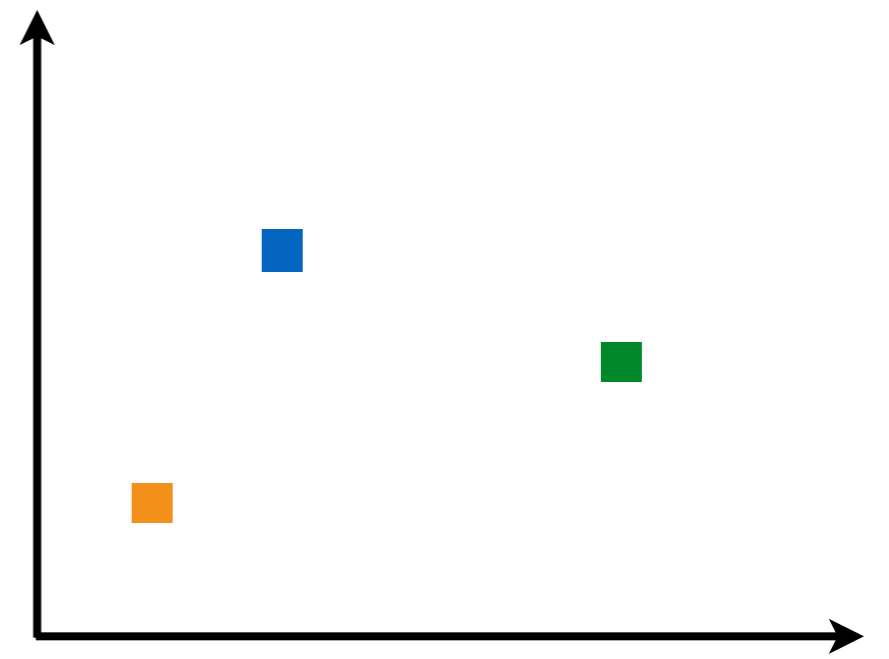
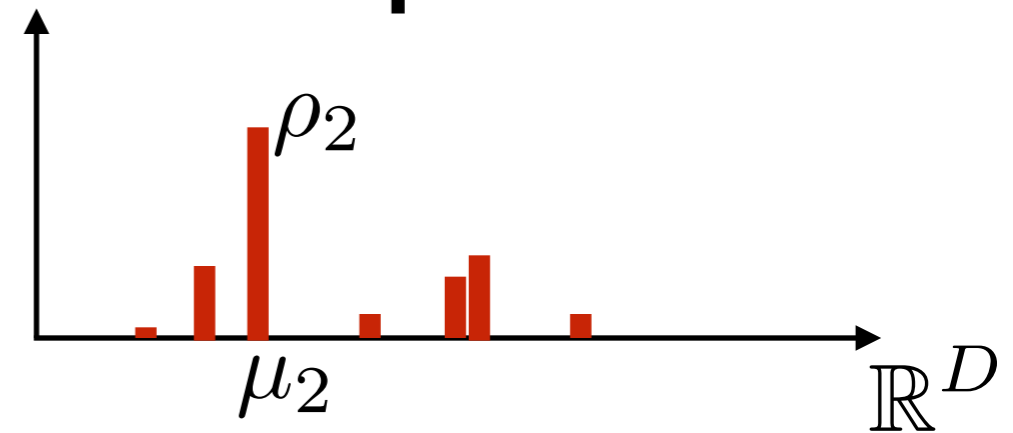
$$\mu_k \stackrel{iid}{\sim} \mathcal{N}(\mu_0, \Sigma_0), k = 1, 2, \dots$$

- i.e.  $G = \sum_{k=1}^{\infty} \rho_k \delta_{\mu_k} \stackrel{d}{=} \text{DP}(\alpha, \mathcal{N}(\mu_0, \Sigma_0))$  **Dirichlet process**



$$z_n \stackrel{iid}{\sim} \text{Categorical}(\rho)$$

$$\mu_n^* = \mu_{z_n}$$



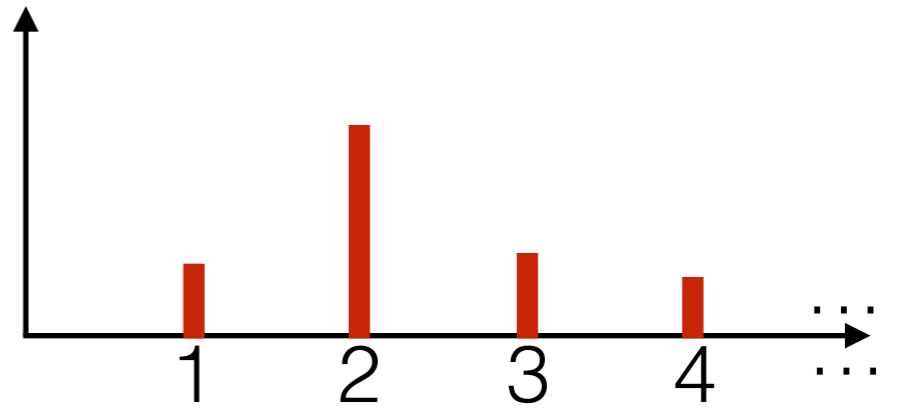
# Dirichlet process mixture model

- Gaussian mixture model

$$\rho = (\rho_1, \rho_2, \dots) \sim \text{GEM}(\alpha)$$

$$\mu_k \stackrel{iid}{\sim} \mathcal{N}(\mu_0, \Sigma_0), k = 1, 2, \dots$$

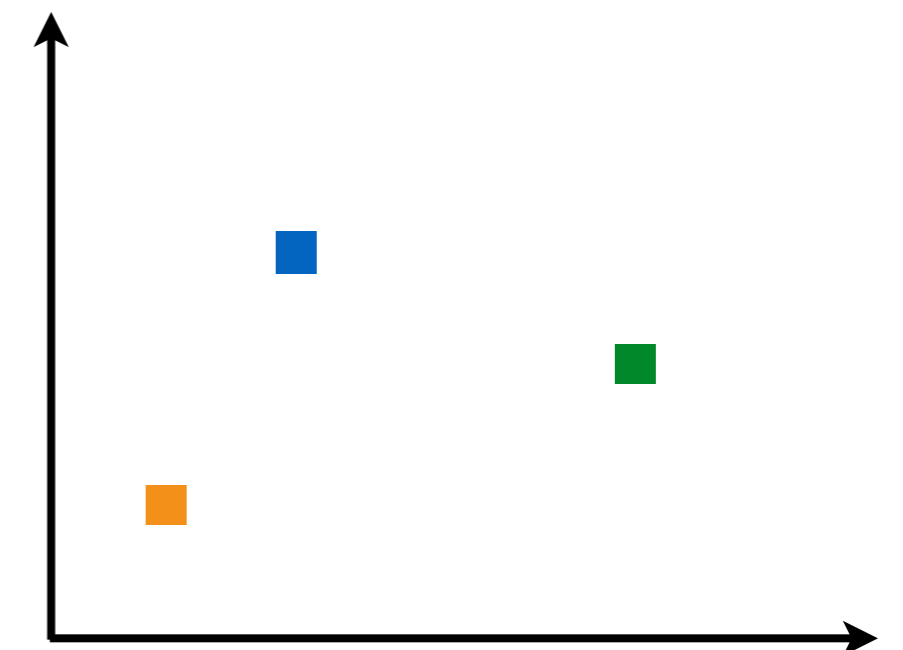
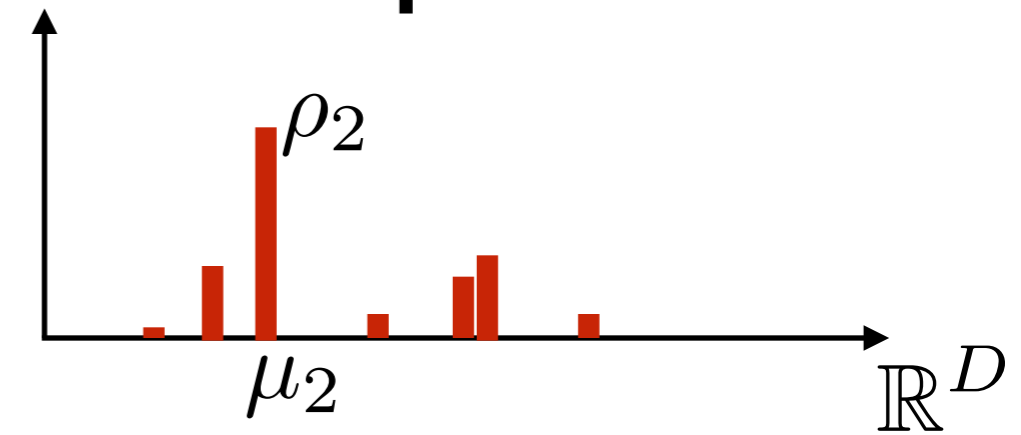
- i.e.  $G = \sum_{k=1}^{\infty} \rho_k \delta_{\mu_k} \stackrel{d}{=} \text{DP}(\alpha, \mathcal{N}(\mu_0, \Sigma_0))$  **Dirichlet process**



$$z_n \stackrel{iid}{\sim} \text{Categorical}(\rho)$$

$$\mu_n^* = \mu_{z_n}$$

- i.e.  $\mu_n^* \stackrel{iid}{\sim} G$



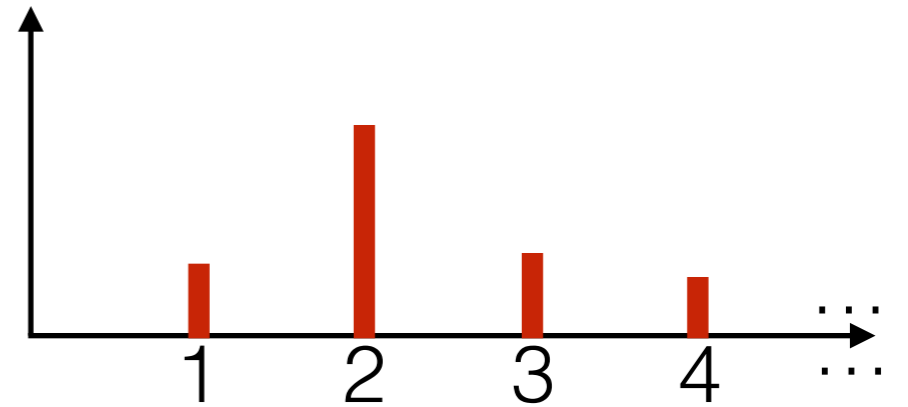
# Dirichlet process mixture model

- Gaussian mixture model

$$\rho = (\rho_1, \rho_2, \dots) \sim \text{GEM}(\alpha)$$

$$\mu_k \stackrel{iid}{\sim} \mathcal{N}(\mu_0, \Sigma_0), k = 1, 2, \dots$$

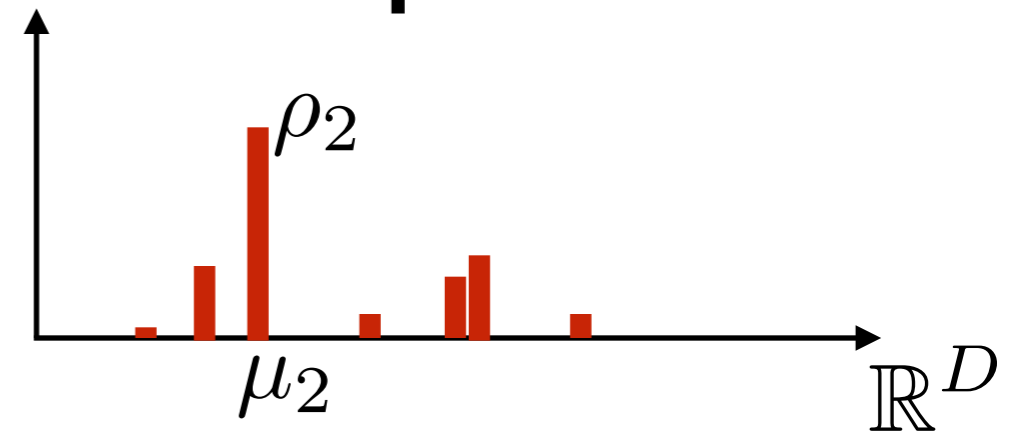
- i.e.  $G = \sum_{k=1}^{\infty} \rho_k \delta_{\mu_k} \stackrel{d}{=} \text{DP}(\alpha, \mathcal{N}(\mu_0, \Sigma_0))$  **Dirichlet process**



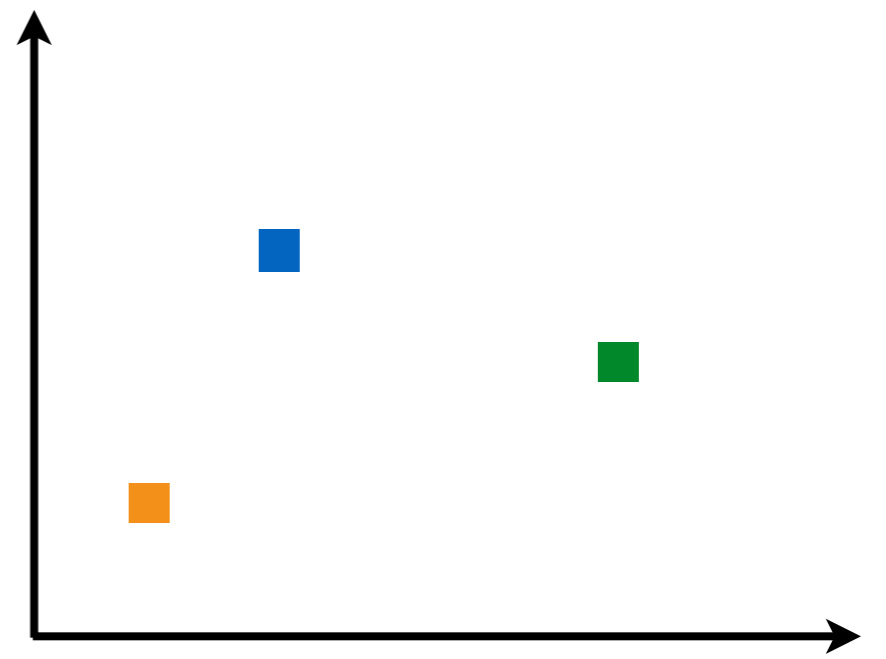
$$z_n \stackrel{iid}{\sim} \text{Categorical}(\rho)$$

$$\mu_n^* = \mu_{z_n}$$

- i.e.  $\mu_n^* \stackrel{iid}{\sim} G$



$$x_n \stackrel{indep}{\sim} \mathcal{N}(\mu_n^*, \Sigma)$$



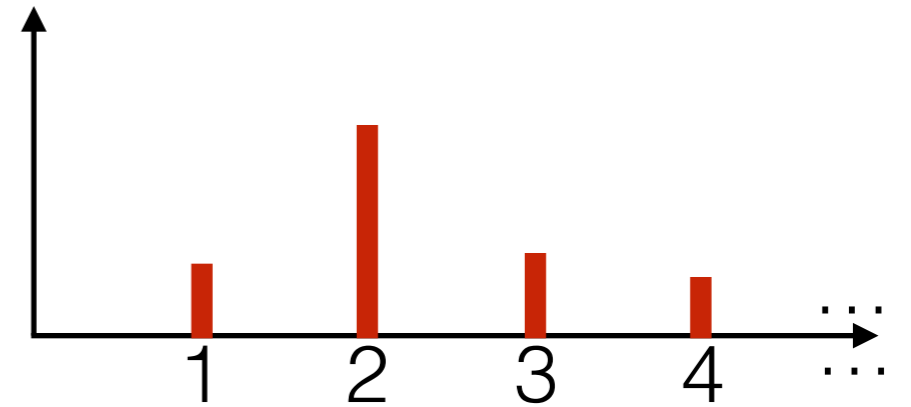
# Dirichlet process mixture model

- Gaussian mixture model

$$\rho = (\rho_1, \rho_2, \dots) \sim \text{GEM}(\alpha)$$

$$\mu_k \stackrel{iid}{\sim} \mathcal{N}(\mu_0, \Sigma_0), k = 1, 2, \dots$$

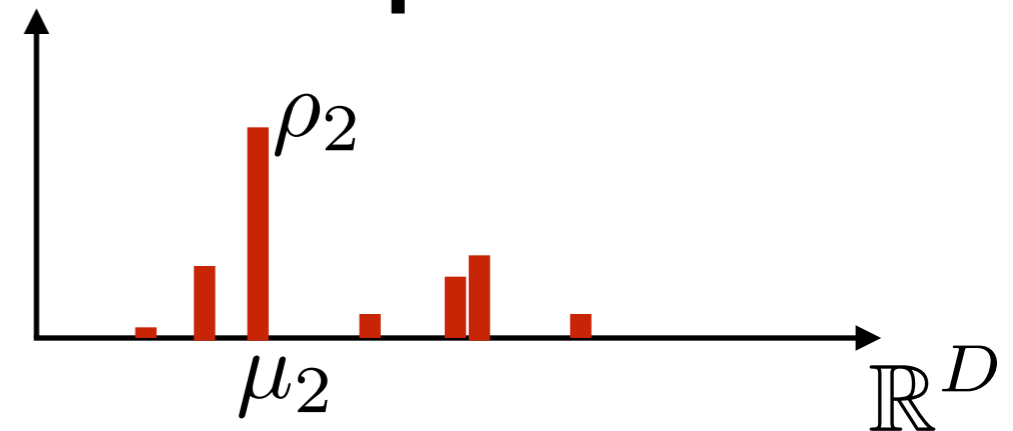
- i.e.  $G = \sum_{k=1}^{\infty} \rho_k \delta_{\mu_k} \stackrel{d}{=} \text{DP}(\alpha, \mathcal{N}(\mu_0, \Sigma_0))$  **Dirichlet process**



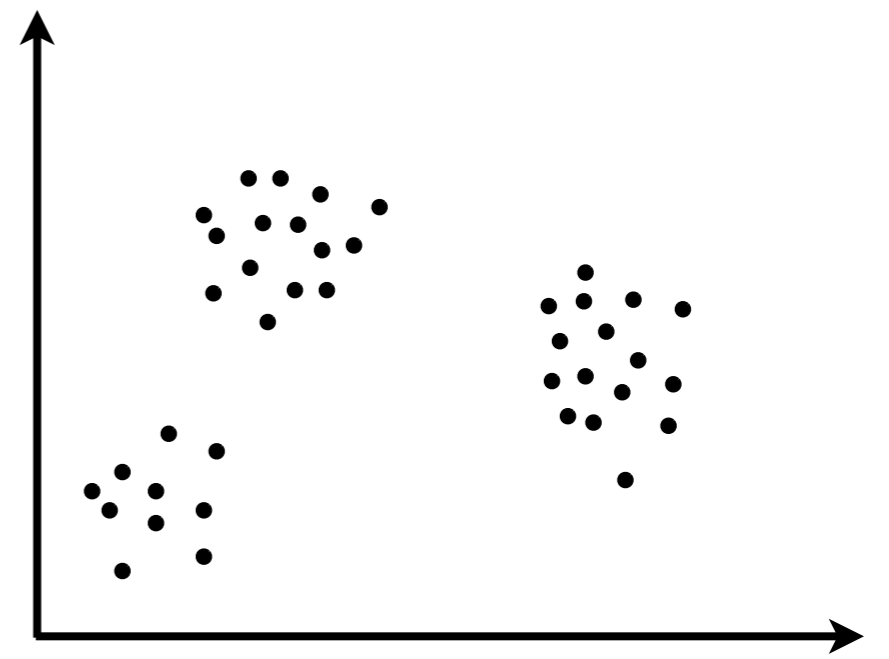
$$z_n \stackrel{iid}{\sim} \text{Categorical}(\rho)$$

$$\mu_n^* = \mu_{z_n}$$

- i.e.  $\mu_n^* \stackrel{iid}{\sim} G$



$$x_n \stackrel{indep}{\sim} \mathcal{N}(\mu_n^*, \Sigma)$$



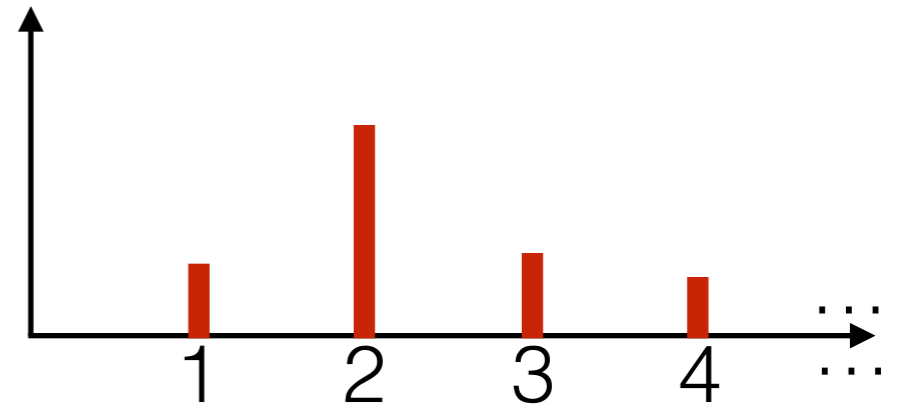
# Dirichlet process mixture model

- Gaussian mixture model

$$\rho = (\rho_1, \rho_2, \dots) \sim \text{GEM}(\alpha)$$

$$\mu_k \stackrel{iid}{\sim} \mathcal{N}(\mu_0, \Sigma_0), k = 1, 2, \dots$$

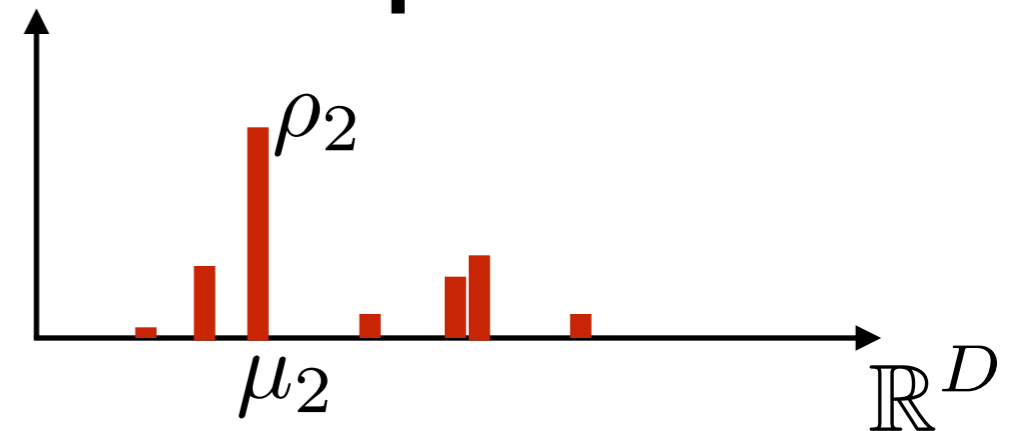
- i.e.  $G = \sum_{k=1}^{\infty} \rho_k \delta_{\mu_k} \stackrel{d}{=} \text{DP}(\alpha, \mathcal{N}(\mu_0, \Sigma_0))$  **Dirichlet process**



$$z_n \stackrel{iid}{\sim} \text{Categorical}(\rho)$$

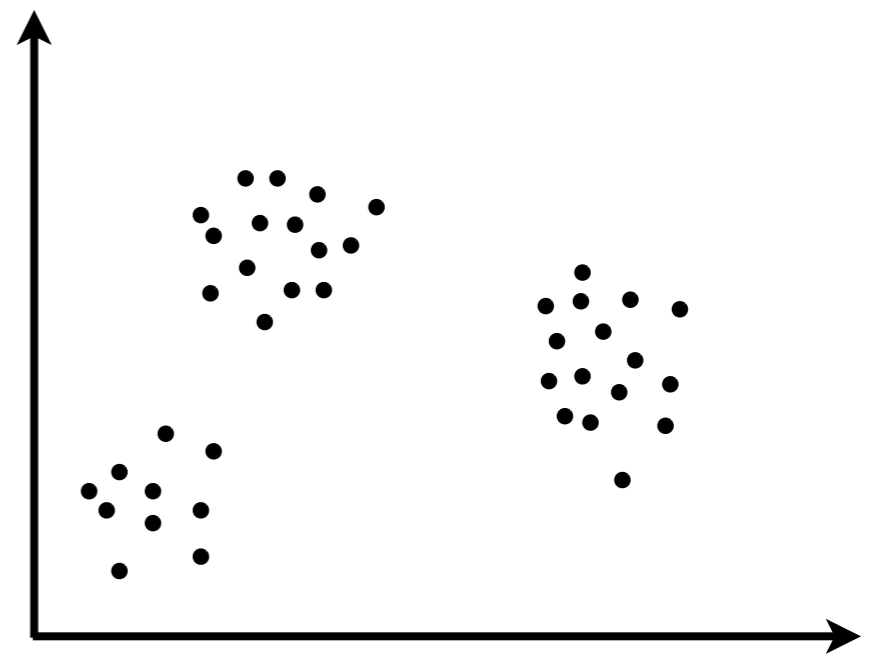
$$\mu_n^* = \mu_{z_n}$$

- i.e.  $\mu_n^* \stackrel{iid}{\sim} G$



$$x_n \stackrel{indep}{\sim} \mathcal{N}(\mu_n^*, \Sigma)$$

[demo]



# Roadmap

- Example problem: clustering
- Example NPBayes model: Dirichlet process
- Big questions
  - Why NPBayes?
  - What does a growing/infinite number of parameters really mean (in NPBayes)?
  - Why is NPBayes challenging but practical?



# Roadmap

- Example problem: clustering
- Example NPBayes model: Dirichlet process
- Big questions
  - Why NPBayes? **Learn more from more data**
  - What does a growing/infinite number of parameters really mean (in NPBayes)?
  - Why is NPBayes challenging but practical?

# Roadmap

- Example problem: clustering
- Example NPBayes model: Dirichlet process
- Big questions
  - Why NPBayes? **Learn more from more data**
  - What does a growing/infinite number of parameters really mean (in NPBayes)? **Components vs clusters; latent vs. realized**
  - Why is NPBayes challenging but practical?

# Roadmap

- Example problem: clustering
- Example NPBayes model: Dirichlet process
- Big questions
  - Why NPBayes? **Learn more from more data**
  - What does a growing/infinite number of parameters really mean (in NPBayes)? **Components vs clusters; latent vs. realized**
  - Why is NPBayes challenging but practical? **Infinite dimensional parameter, but finitely many realized**

# Roadmap

- Example problem: clustering
- Example NPBayes model: Dirichlet process
- Big questions
  - Why NPBayes? **Learn more from more data**
  - What does a growing/infinite number of parameters really mean (in NPBayes)? **Components vs clusters; latent vs. realized**
  - Why is NPBayes challenging but practical? **Infinite dimensional parameter, but finitely many realized**
- **Typical approaches:**

# Roadmap

- Example problem: clustering
- Example NPBayes model: Dirichlet process
- Big questions
  - Why NPBayes? **Learn more from more data**
  - What does a growing/infinite number of parameters really mean (in NPBayes)? **Components vs clusters; latent vs. realized**
  - Why is NPBayes challenging but practical? **Infinite dimensional parameter, but finitely many realized**
- **Typical approaches:**
  - **Integrate out the infinite parameter**

# Roadmap

- Example problem: clustering
- Example NPBayes model: Dirichlet process
- Big questions
  - Why NPBayes? **Learn more from more data**
  - What does a growing/infinite number of parameters really mean (in NPBayes)? **Components vs clusters; latent vs. realized**
  - Why is NPBayes challenging but practical? **Infinite dimensional parameter, but finitely many realized**
- **Typical approaches:**
  - **Integrate out the infinite parameter**
  - **Truncate the infinite parameter**

# References (page 1 of 4)

DJ Aldous. *Exchangeability and related topics*. Springer, 1983.

CE Antoniak. Mixtures of Dirichlet processes with applications to Bayesian nonparametric problems. *The Annals of Statistics*, 1974.

E Arjas and D Gasbarra. Nonparametric Bayesian inference from right censored survival data, using the Gibbs sampler. *Statistica Sinica*, 1994.

J Bertoin. *Random Fragmentation and Coagulation Processes*. Cambridge University Press, 2006.

D Blackwell and JB MacQueen. Ferguson distributions via Pólya urn schemes. *The Annals of Statistics*, 1973.

E Bowlby. NOAA/Olympic Coast NMS; NOAA/OAR/Office of Ocean Exploration - NOAA Photo Library. Retrieved from: [https://en.wikipedia.org/wiki/Opisthoteuthis\\_californiana#/media/File:Opisthoteuthis\\_californiana.jpg](https://en.wikipedia.org/wiki/Opisthoteuthis_californiana#/media/File:Opisthoteuthis_californiana.jpg)

T Broderick, MI Jordan, and J Pitman. Beta processes, stick-breaking, and power laws. *Bayesian Analysis*, 2012.

T Broderick, MI Jordan, and J Pitman. Cluster and feature modeling from combinatorial stochastic processes. *Statistical Science*, 2013.

T Broderick, J Pitman, and MI Jordan. Feature allocations, probability functions, and paintboxes. *Bayesian Analysis*, 2013.

T Broderick, AC Wilson, and MI Jordan. Posteriors, conjugacy, and exponential families for completely random measures. arXiv preprint arXiv:1410.6843, 2014.

T Campbell\*, JH Huggins\*, J How, T Broderick. Truncation random measures. arXiv preprint arXiv:1603.00861, 2016.

S Engen. A note on the geometric series as a species frequency model. *Biometrika*, 1975.

MD Escobar and M West. Bayesian density estimation and inference using mixtures. *Journal of the American Statistical Association*, 1995.

# References (page 2 of 4)

W Ewens. The sampling theory of selectively neutral alleles. *Theoretical Population Biology*, 1972.

W Ewens. Population genetics theory -- the past and the future. *Mathematical and Statistical Developments of Evolutionary Theory*, 1987.

TS Ferguson. A Bayesian analysis of some nonparametric problems. *The Annals of Statistics*, 1973.

TS Ferguson. Bayesian density estimation by mixtures of normal distributions. *Recent Advances in Statistics*, 1983.

EB Fox, personal website. Retrieved from: <http://www.stat.washington.edu/~ebfox/research.html> --- Associated paper: EB Fox, MC Hughes, EB Sudderth, and MI Jordan. *The Annals of Applied Statistics*, 2014.

S Ghosal, JK Ghosh, and RV Ramamoorthi. Posterior consistency of Dirichlet mixtures in density estimation. *The Annals of Statistics*, 1999.

S Goldwater, TL Griffiths, and M Johnson. Interpolating between types and tokens by estimating power-law generators. *NIPS*, 2005.

A Gnedin, B Hansen, and J Pitman. Notes on the occupancy problem with infinitely many boxes: general asymptotics and power laws. *Probability Surveys*, 2007.

TL Griffiths and Z Ghahramani. Infinite latent feature models and the Indian buffet process. *NIPS*, 2005.

DL Hartl and AG Clark. *Principles of Population Genetics, Fourth Edition*. 2003.

E Hewitt and LJ Savage. Symmetric measures on Cartesian products. *Transactions of the American Mathematical Society*, 1955.

NL Hjort. Nonparametric Bayes estimators based on beta processes in models for life history data. *The Annals of Statistics*, 1990.

DN Hoover. Relations on probability spaces and arrays of random variables, *Preprint, Institute for Advanced Study*, 1979.



# References (page 3 of 4)

- FM Hoppe. Pólya-like urns and the Ewens' sampling formula. *Journal of Mathematical Biology*, 1984.
- H Ishwaran and LF James. Gibbs sampling methods for stick-breaking priors. *Journal of the American Statistical Association*, 2001.
- L James. Poisson latent feature calculus for generalized Indian buffet processes. arXiv preprint arXiv:1411.2936, 2014.
- Y Kim. Nonparametric Bayesian estimators for counting processes. *The Annals of Statistics*, 1999.
- JFC Kingman. The representation of partition structures. *Journal of the London Mathematical Society*, 1978.
- JFC Kingman. On the genealogy of large populations. *Journal of Applied Probability*, 1982.
- JFC Kingman. *Poisson processes*, 1992.
- JR Lloyd, P Orbanz, Z Ghahramani, and DM Roy. Random function priors for exchangeable arrays with applications to graphs and relational data. *NIPS*, 2012.
- SN MacEachern and P Müller. Estimating mixtures of Dirichlet process models. *Journal of Computational and Graphical Statistics*, 1998.
- JW McCloskey. A model for the distribution of individuals by species in an environment. *Ph.D. thesis, Michigan State University*, 1965.
- K Miller, MI Jordan, and TL Griffiths. Nonparametric latent feature models for link prediction. *NIPS*, 2009.
- RM Neal. Density modeling and clustering using Dirichlet diffusion trees. *Bayesian Statistics*, 2003.
- P Orbanz. Construction of nonparametric Bayesian models from parametric Bayes equations. *NIPS*, 2009.
- P Orbanz. Conjugate Projective Limits. arXiv preprint arXiv:1012.0363, 2010.
- P Orbanz, DM Roy. Bayesian models of graphs, arrays and other exchangeable random structures. *IEEE TPAMI*, 2015.

# References (page 4 of 4)

GP Patil and C Taillie. Diversity as a concept and its implications for random communities. *Bulletin of the International Statistical Institute*, 1977.

J Pitman and M Yor. The two-parameter Poisson-Dirichlet distribution derived from a stable subordinator. *The Annals of Probability*, 1997.

A Rodríguez, DB Dunson & AE Gelfand. The nested Dirichlet process. *Journal of the American Statistical Association*, 2008.

S Saria, D Koller, and A Penn. Learning individual and population traits from clinical temporal data. *NIPS*, 2010.

J Sethuraman. A constructive definition of Dirichlet priors. *Statistica Sinica*, 1994.

EB Sudderth and MI Jordan. Shared segmentation of natural scenes using dependent Pitman-Yor processes. *NIPS*, 2009.

YW Teh. A hierarchical Bayesian language model based on Pitman-Yor processes. *ACL*, 2006.

YW Teh, C Blundell, and L Elliott. Modelling genetic variations using fragmentation-coagulation processes. *NIPS*, 2011.

YW Teh, MI Jordan, MJ Beal, and DM Blei. Hierarchical Dirichlet processes. *Journal of the American Statistical Association*, 2006.

R Thibaux and MI Jordan. Hierarchical beta processes and the Indian buffet process. *ICML*, 2007.

J Wakeley. *Coalescent Theory: An Introduction*, Chapter 3, 2008.

M West, P Müller, and MD Escobar. Hierarchical Priors and Mixture Models, With Application in Regression and Density Estimation. *Aspects of Uncertainty*, 1994.