



Machine Learning and Nonparametric Bayes

Peter Orbanz
Assistant Professor
Statistics
Columbia

Tamara Broderick
ITT Career Development Assistant Professor
Electrical Engineering & Computer Science
MIT



Machine Learning and **Nonparametric Bayes**

Peter Orbanz

Assistant Professor
Statistics
Columbia

Tamara Broderick

ITT Career Development Assistant Professor
Electrical Engineering & Computer Science
MIT

Nonparametric Bayes

Nonparametric Bayes

- Bayesian statistics that is not parametric

Nonparametric Bayes

- Bayesian statistics that is not parametric (wait!)

Nonparametric Bayes

- Bayesian statistics that is not parametric
- Bayesian

Nonparametric Bayes

- Bayesian statistics that is not parametric
- Bayesian

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

Nonparametric Bayes

- Bayesian statistics that is not parametric
- 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 statistics that is not parametric
- 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 statistics that is not parametric
- 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 statistics that is not parametric
- 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

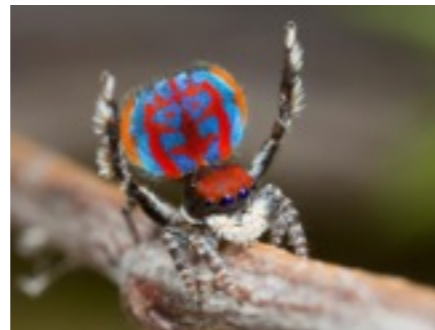
- Bayesian statistics that is not parametric
- Bayesian

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

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



[wikipedia.org]



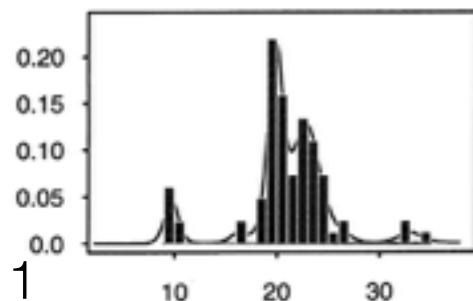
[Time Mag]

Nonparametric Bayes

- Bayesian statistics that is not parametric
- 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)



[Escobar,
West 1995;
Ghosal
et al 1999]

Nonparametric Bayes

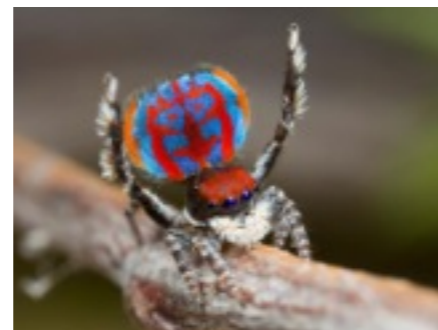
- Bayesian statistics that is not parametric
- Bayesian

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

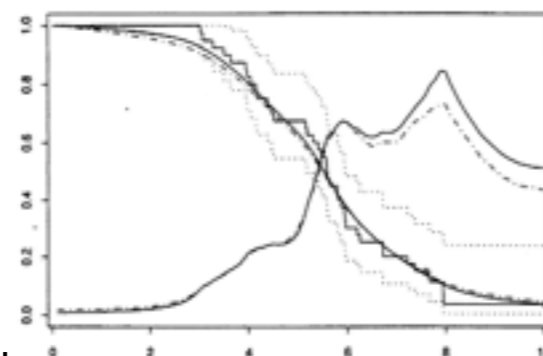
- Not parametric (i.e. not finite parameter, unbounded/growing/infinite number of parameters)



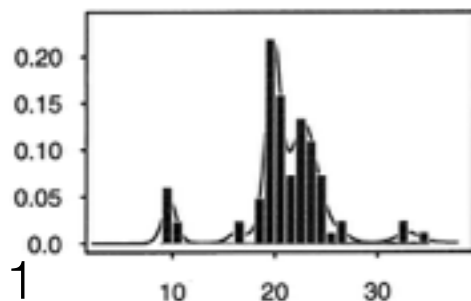
[wikipedia.org]



[Time Mag]



[Arjas, Gasbarra 1994]



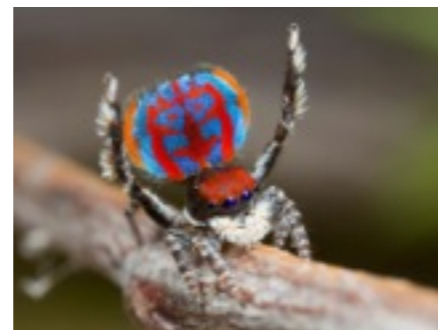
[Escobar, West 1995; Ghosal et al 1999]

Nonparametric Bayes

- Bayesian statistics that is not parametric
- Bayesian

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

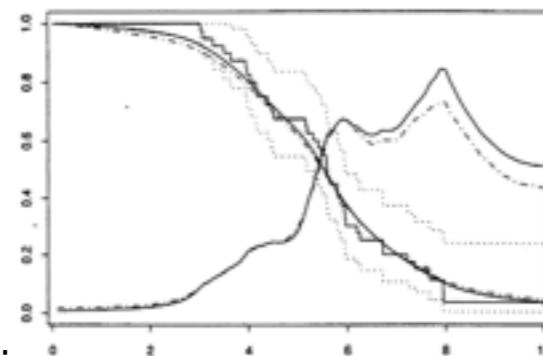
- Not parametric (i.e. not finite parameter, unbounded/growing/infinite number of parameters)



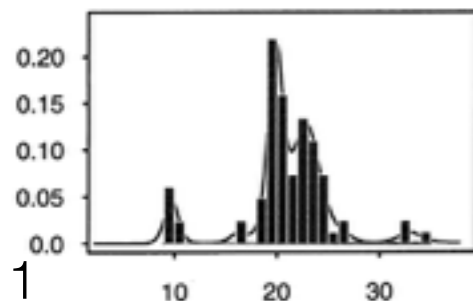
[Time Mag]



[Fox et al 2014]



[Arjas, Gasbarra 1994]



[Escobar, West 1995; Ghosal et al 1999]

Nonparametric Bayes

- Bayesian statistics that is not parametric
- Bayesian

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

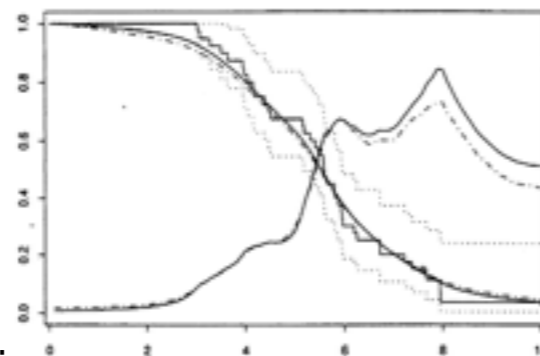
- Not parametric (i.e. not finite parameter, unbounded/growing/infinite number of parameters)



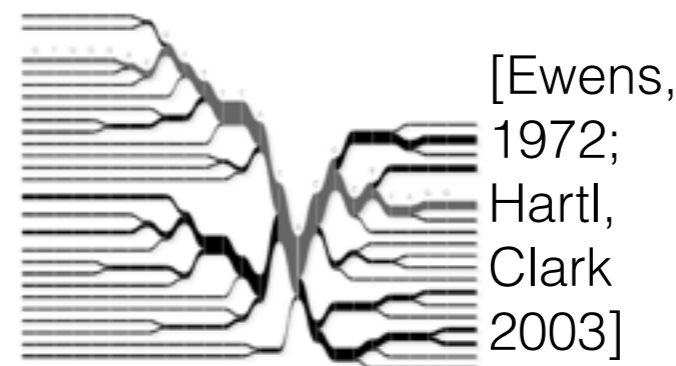
[Time Mag]



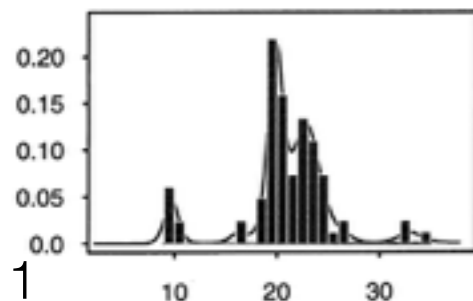
[Fox et al 2014]



[Arjas, Gasbarra 1994]



[Ewens, 1972; Hartl, Clark 2003]



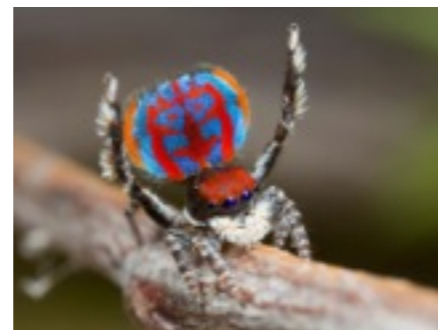
[Escobar, West 1995; Ghosal et al 1999]

Nonparametric Bayes

- Bayesian statistics that is not parametric
- Bayesian

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

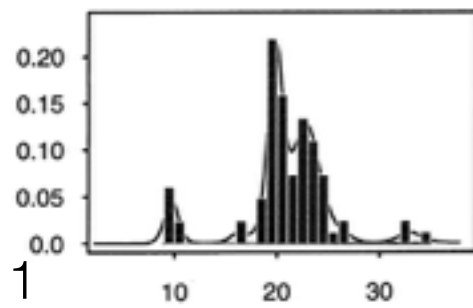
- Not parametric (i.e. not finite parameter, unbounded/growing/infinite number of parameters)



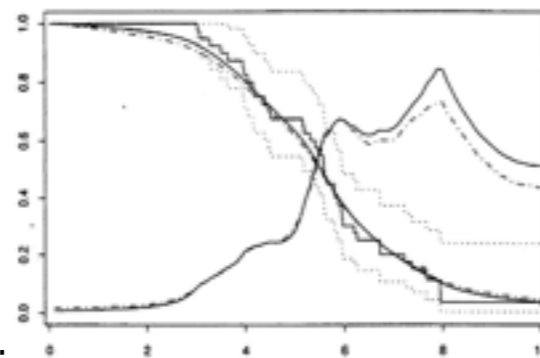
[Time Mag]



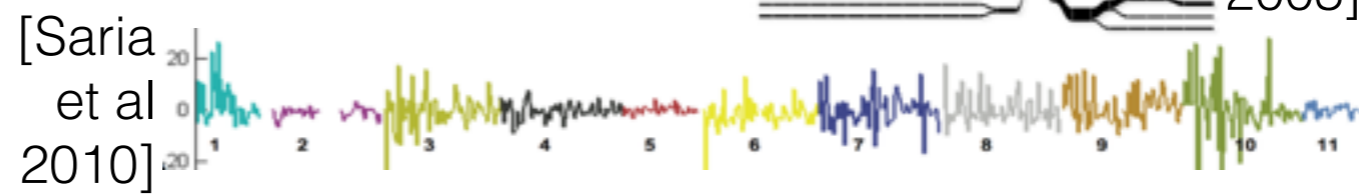
[Fox et al 2014]



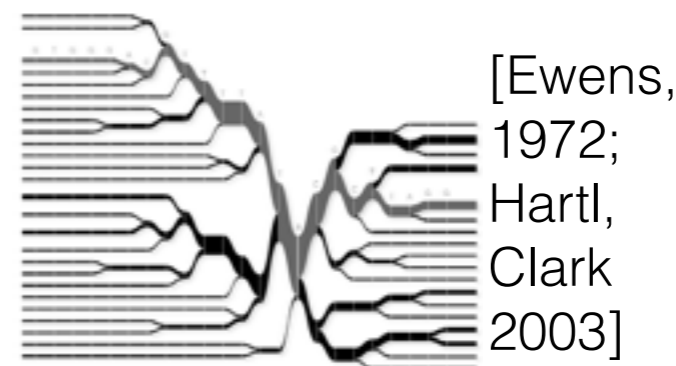
[Escobar, West 1995; Ghosal et al 1999]



[Arjas, Gasbarra 1994]



[Saria et al 2010]



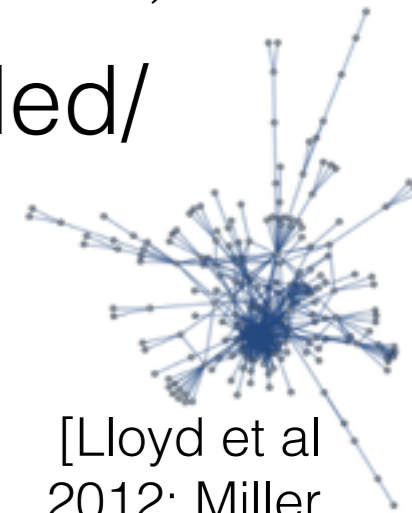
[Ewens, 1972; Hartl, Clark 2003]

Nonparametric Bayes

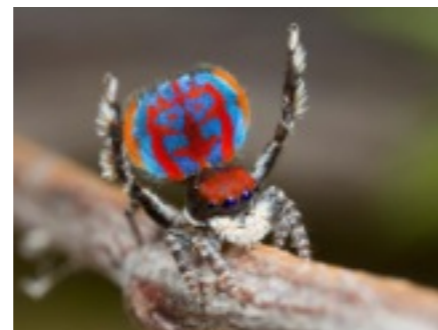
- Bayesian statistics that is not parametric
- Bayesian

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

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



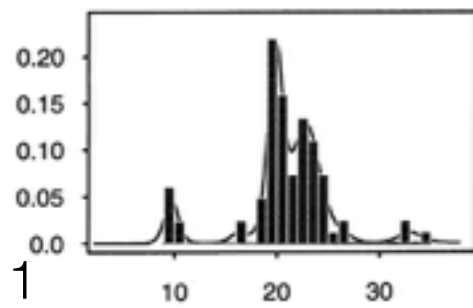
[Lloyd et al 2012; Miller et al 2010]



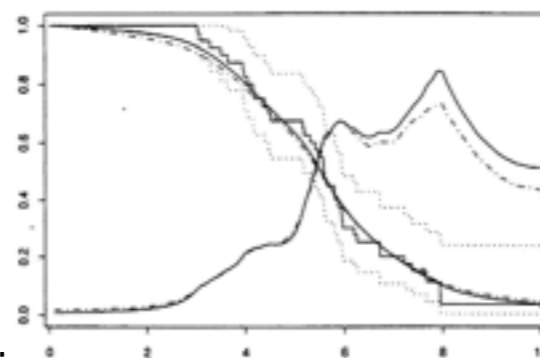
[Time Mag]



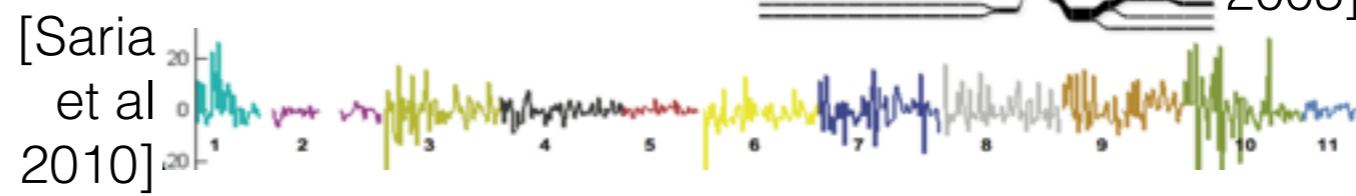
[Fox et al 2014]



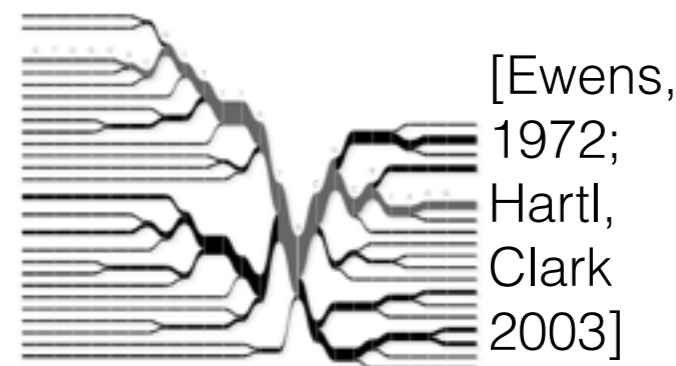
[Escobar, West 1995; Ghosal et al 1999]



[Arjas, Gasbarra 1994]



[Saria et al 2010]



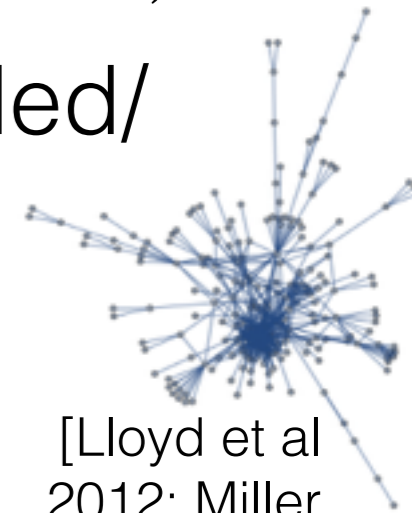
[Ewens, 1972; Hartl, Clark 2003]

Nonparametric Bayes

- Bayesian statistics that is not parametric
- 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)



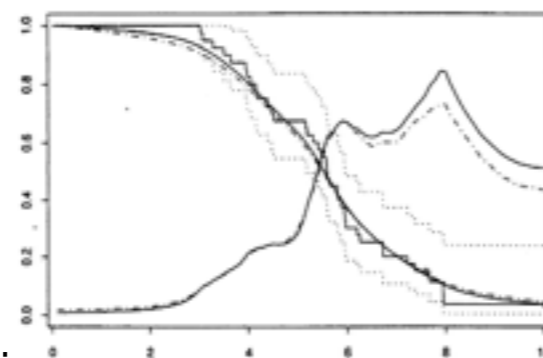
[Lloyd et al 2012; Miller et al 2010]



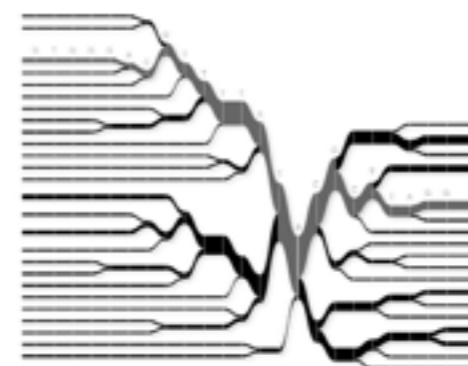
[Time Mag]



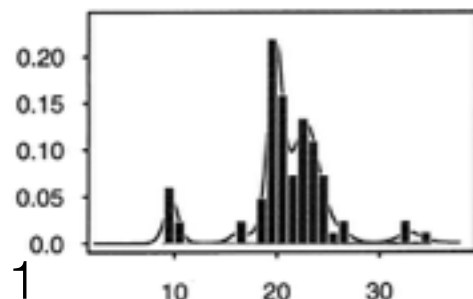
[Fox et al 2014]



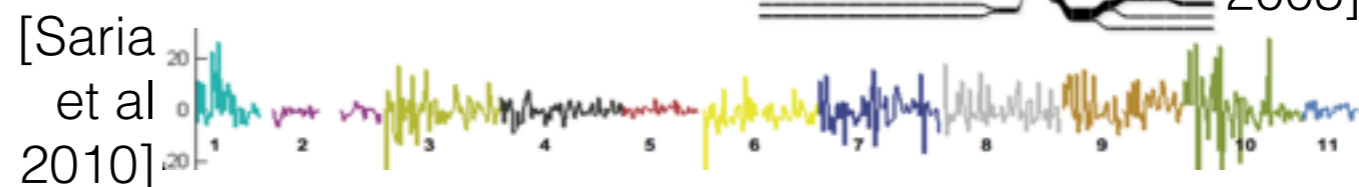
[Arjas, Gasbarra 1994]



[Ewens 1972; Hartl, Clark 2003]



[Escobar, West 1995; Ghosal et al 1999]



[Saria et al 2010]



[Sudderth, Jordan 2009]

Nonparametric Bayes

- A theoretical motivation: De Finetti's Theorem

Nonparametric Bayes

- A theoretical motivation: De Finetti's Theorem
- A data sequence is *infinitely exchangeable* if the distribution of any N data points doesn't change when permuted: $p(X_1, \dots, X_N) = p(X_{\sigma(1)}, \dots, X_{\sigma(N)})$

Nonparametric Bayes

- A theoretical motivation: De Finetti's Theorem
- A data sequence is *infinitely exchangeable* if the distribution of any N data points doesn't change when permuted: $p(X_1, \dots, X_N) = p(X_{\sigma(1)}, \dots, X_{\sigma(N)})$
- *De Finetti's Theorem* (roughly): A sequence X_1, X_2, \dots is infinitely exchangeable if and only if, for all N and some distribution P :

Nonparametric Bayes

- A theoretical motivation: De Finetti's Theorem
- A data sequence is *infinitely exchangeable* if the distribution of any N data points doesn't change when permuted: $p(X_1, \dots, X_N) = p(X_{\sigma(1)}, \dots, X_{\sigma(N)})$
- *De Finetti's Theorem* (roughly): A sequence X_1, X_2, \dots is infinitely exchangeable if and only if, for all N and some distribution P :

$$p(X_1, \dots, X_N) = \int_{\theta} \prod_{n=1}^N p(X_n | \theta) P(d\theta)$$

Nonparametric Bayes

- A theoretical motivation: De Finetti's Theorem
- A data sequence is *infinitely exchangeable* if the distribution of any N data points doesn't change when permuted: $p(X_1, \dots, X_N) = p(X_{\sigma(1)}, \dots, X_{\sigma(N)})$
- *De Finetti's Theorem* (roughly): A sequence X_1, X_2, \dots is infinitely exchangeable if and only if, for all N and some distribution P :

$$p(X_1, \dots, X_N) = \int_{\theta} \prod_{n=1}^N p(X_n | \theta) P(d\theta)$$

- Motivates:

Nonparametric Bayes

- A theoretical motivation: De Finetti's Theorem
- A data sequence is *infinitely exchangeable* if the distribution of any N data points doesn't change when permuted: $p(X_1, \dots, X_N) = p(X_{\sigma(1)}, \dots, X_{\sigma(N)})$
- *De Finetti's Theorem* (roughly): A sequence X_1, X_2, \dots is infinitely exchangeable if and only if, for all N and some distribution P :

$$p(X_1, \dots, X_N) = \int_{\theta} \prod_{n=1}^N p(X_n | \theta) P(d\theta)$$

- Motivates:
 - Parameters and likelihoods

Nonparametric Bayes

- A theoretical motivation: De Finetti's Theorem
- A data sequence is *infinitely exchangeable* if the distribution of any N data points doesn't change when permuted: $p(X_1, \dots, X_N) = p(X_{\sigma(1)}, \dots, X_{\sigma(N)})$
- *De Finetti's Theorem* (roughly): A sequence X_1, X_2, \dots is infinitely exchangeable if and only if, for all N and some distribution P :

$$p(X_1, \dots, X_N) = \int_{\theta} \prod_{n=1}^N p(X_n | \theta) P(d\theta)$$

- Motivates:
 - Parameters and likelihoods
 - Priors

Nonparametric Bayes

- A theoretical motivation: De Finetti's Theorem
- A data sequence is *infinitely exchangeable* if the distribution of any N data points doesn't change when permuted: $p(X_1, \dots, X_N) = p(X_{\sigma(1)}, \dots, X_{\sigma(N)})$
- *De Finetti's Theorem* (roughly): A sequence X_1, X_2, \dots is infinitely exchangeable if and only if, for all N and some distribution P :

$$p(X_1, \dots, X_N) = \int_{\theta} \prod_{n=1}^N p(X_n | \theta) P(d\theta)$$

- Motivates:
 - Parameters and likelihoods
 - Priors
 - “Nonparametric Bayesian” priors

Outline

Outline

- Dirichlet process

Outline

- Dirichlet process
 - Background for intuition

Outline

- Dirichlet process
 - Background for intuition
 - Generative model

Outline

- Dirichlet process
 - Background for intuition
 - Generative model
 - What does a growing/infinite number of parameters really mean (in Nonparametric Bayes)?

Outline

- Dirichlet process
 - Background for intuition
 - Generative model
 - What does a growing/infinite number of parameters really mean (in Nonparametric Bayes)?
- Chinese restaurant process

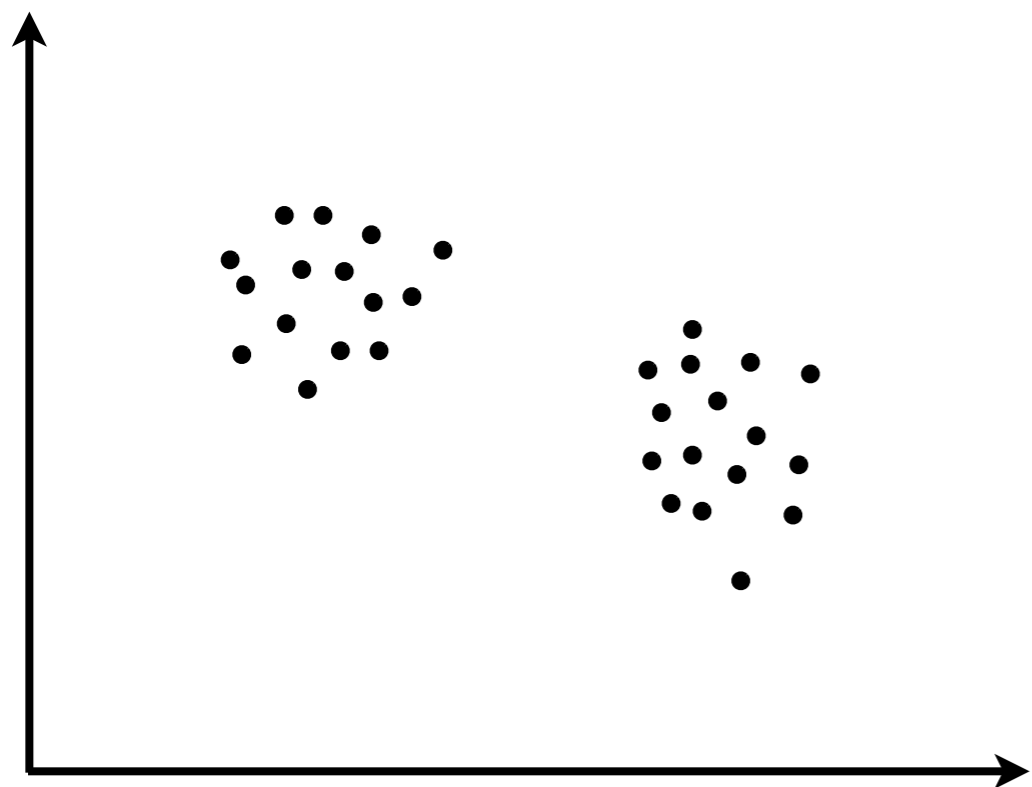
Outline

- Dirichlet process
 - Background for intuition
 - Generative model
 - What does a growing/infinite number of parameters really mean (in Nonparametric Bayes)?
- Chinese restaurant process
- Inference

Outline

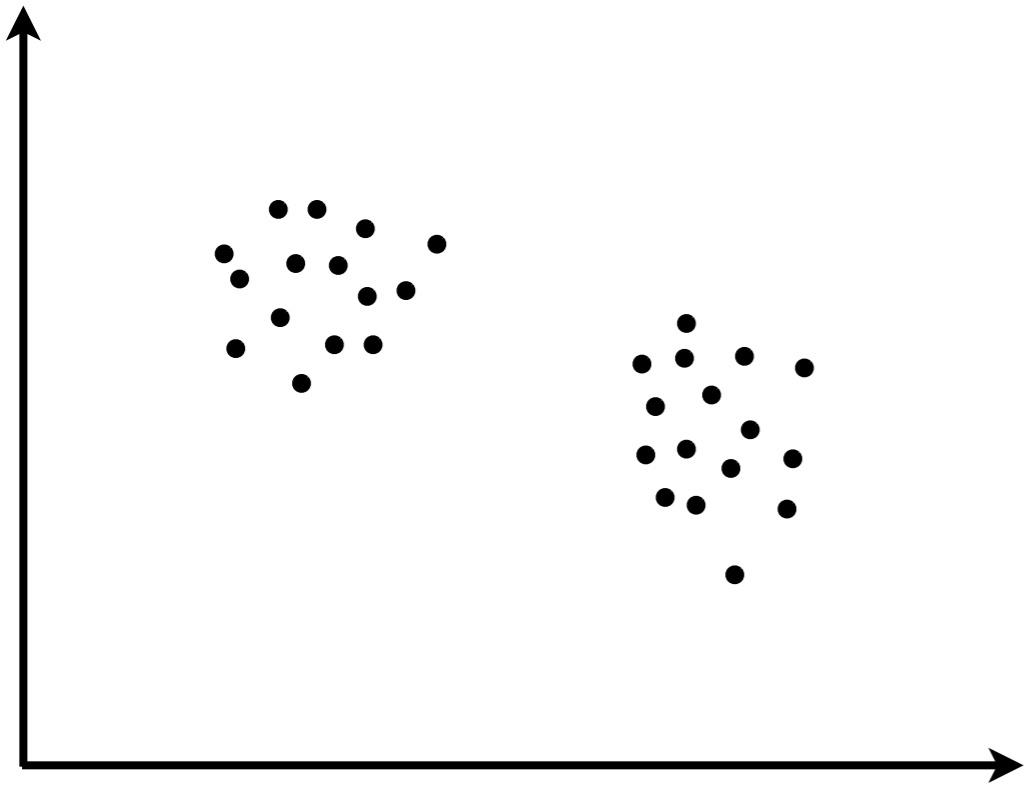
- Dirichlet process
 - Background for intuition
 - Generative model
 - What does a growing/infinite number of parameters really mean (in Nonparametric Bayes)?
- Chinese restaurant process
- Inference
- Venture further into the wild world of Nonparametric Bayesian statistics

Generative model



Generative model

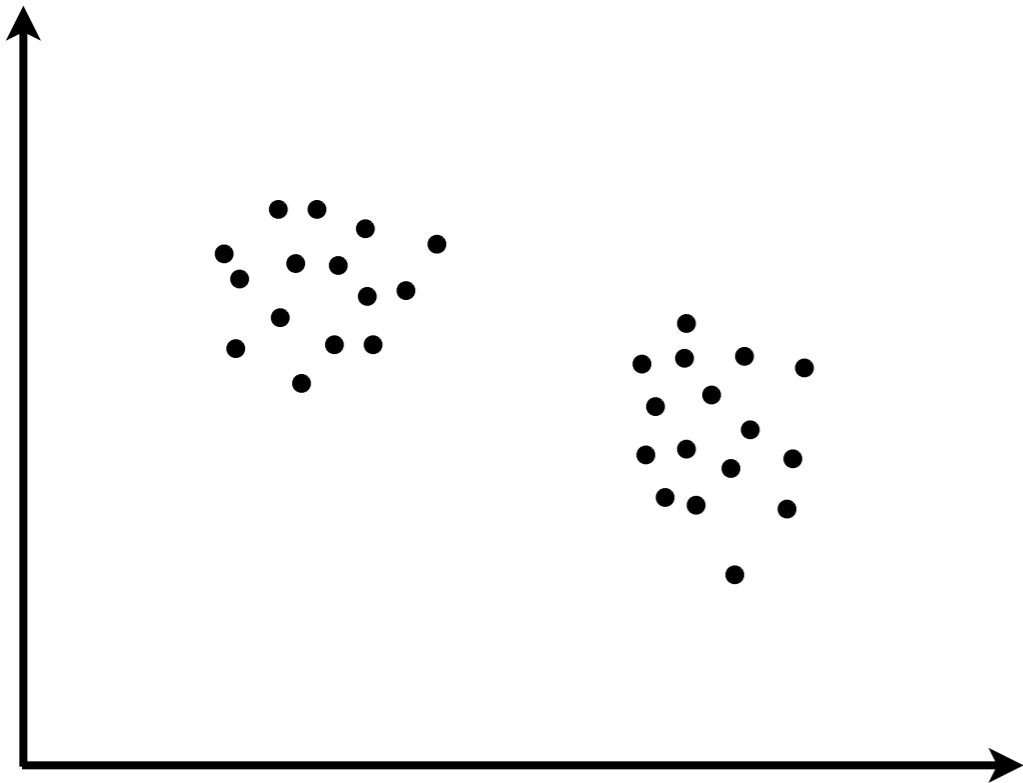
- Finite Gaussian mixture model ($K=2$ clusters)



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)

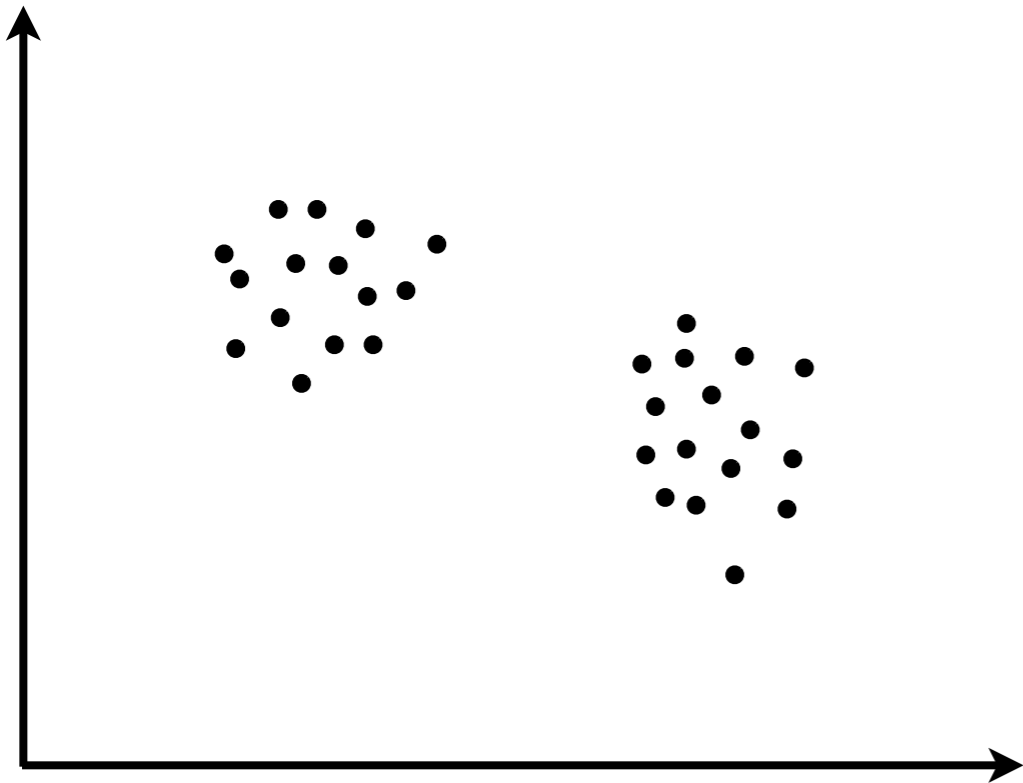


Generative model

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

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

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

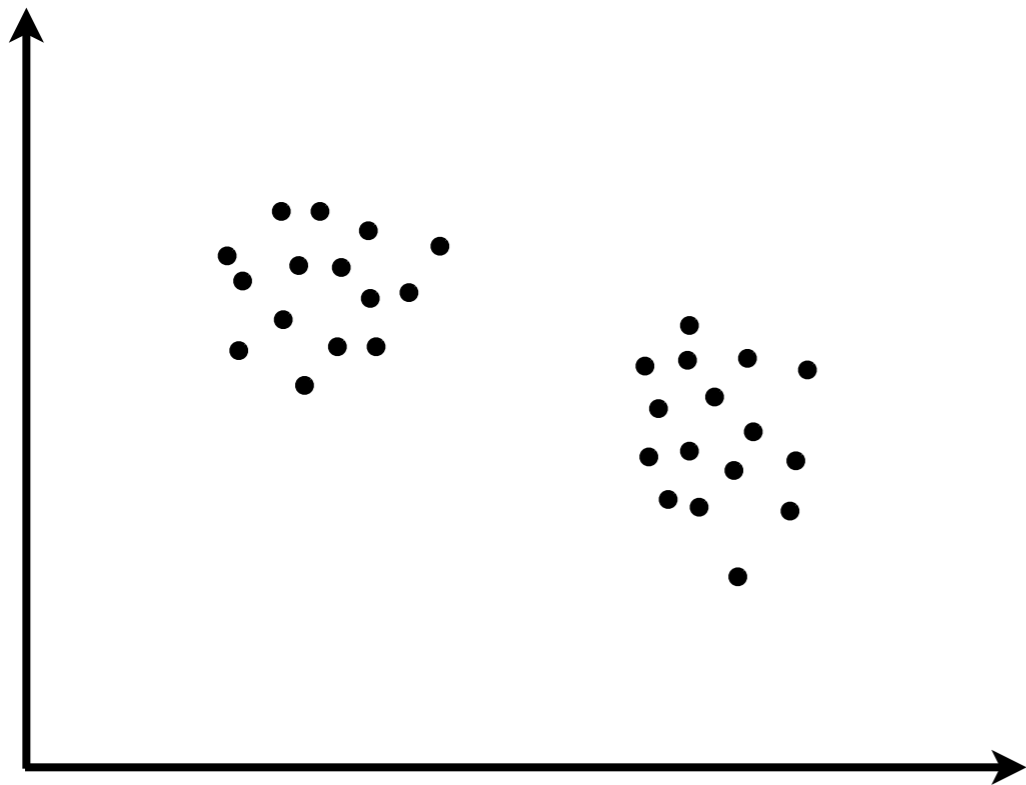


Generative model

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

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

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



ρ_1

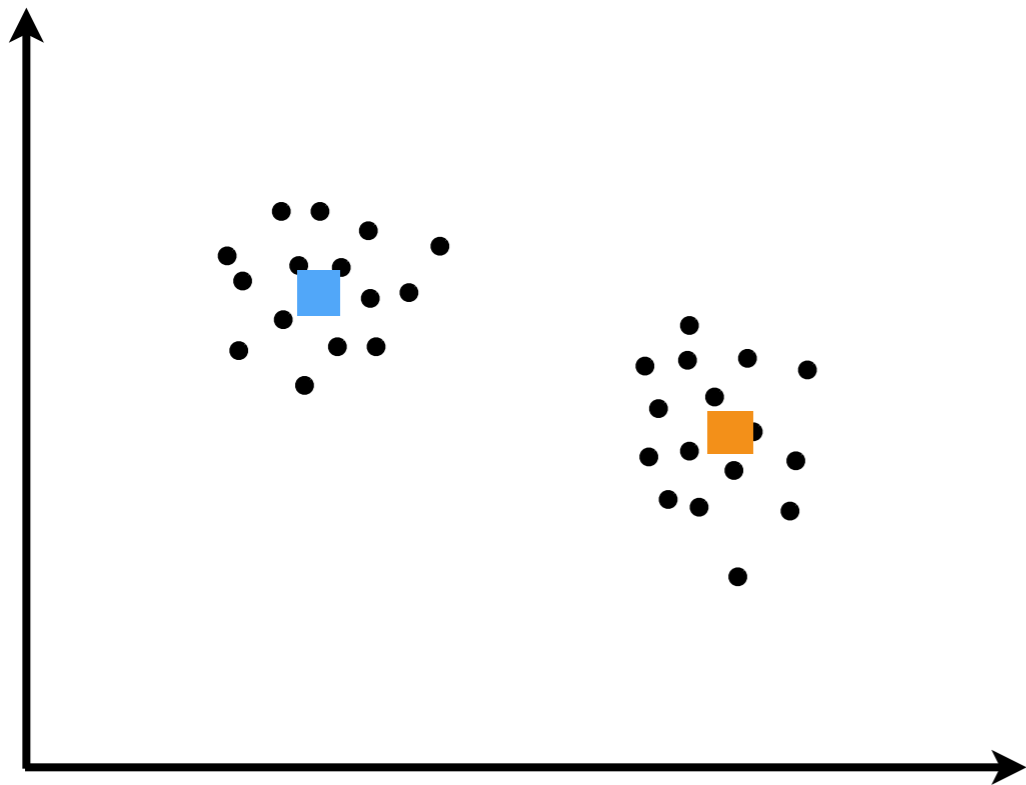
ρ_2

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)

$$z_n \stackrel{iid}{\sim} \text{Categorical}(\rho_1, \rho_2)$$
$$x_n \stackrel{indep}{\sim} \mathcal{N}(\mu_{z_n}, \Sigma)$$



ρ_1

ρ_2

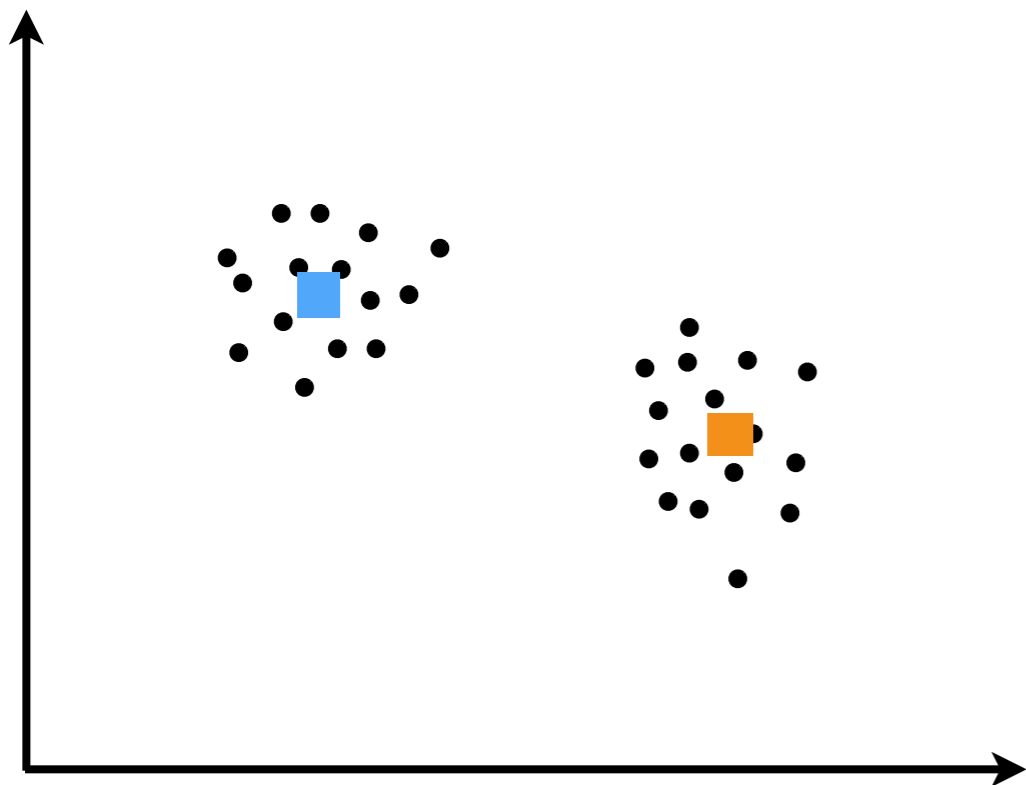
Generative model

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

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

$$z_n \stackrel{iid}{\sim} \text{Categorical}(\rho_1, \rho_2)$$
$$x_n \stackrel{indep}{\sim} \mathcal{N}(\mu_{z_n}, \Sigma)$$

- Don't know μ_1, μ_2

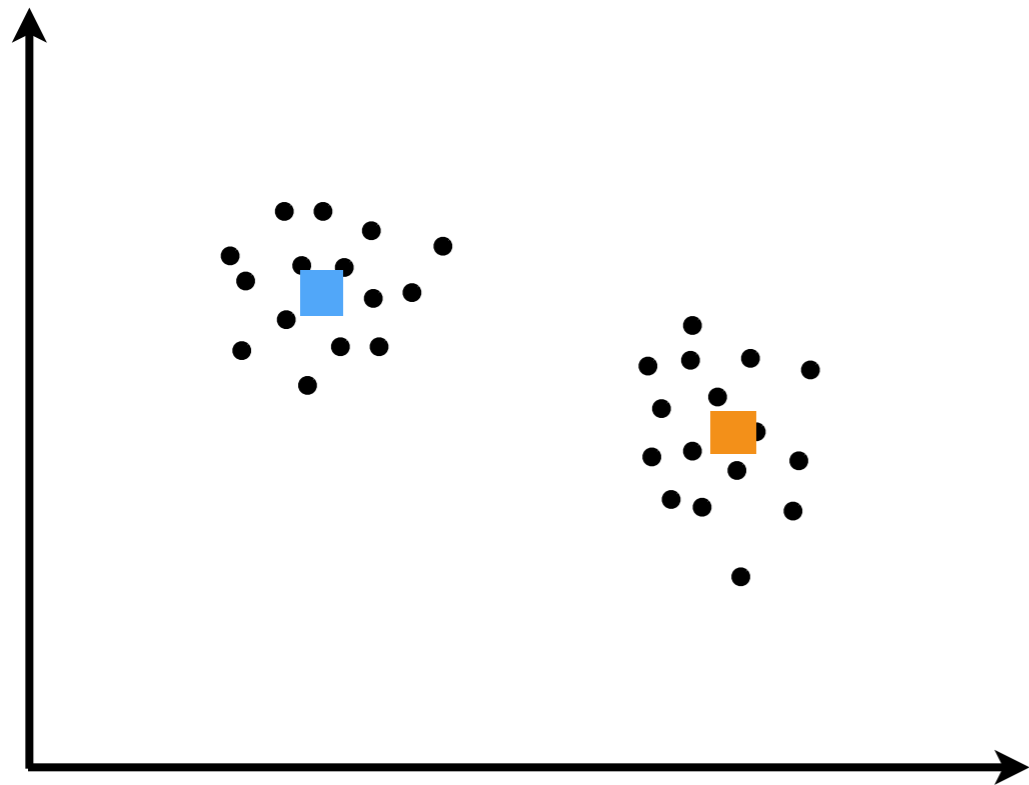


ρ_1

ρ_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)

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

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

- Don't know μ_1, μ_2

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

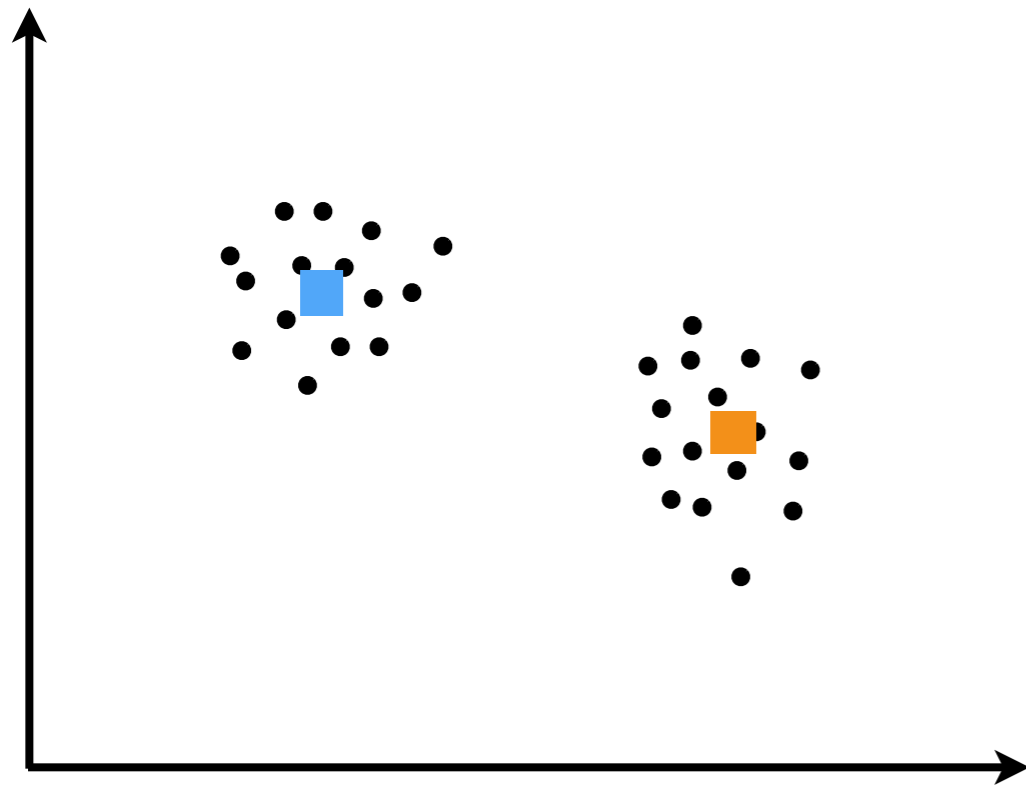


ρ_1

ρ_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)

$$z_n \stackrel{iid}{\sim} \text{Categorical}(\rho_1, \rho_2)$$
$$x_n \stackrel{indep}{\sim} \mathcal{N}(\mu_{z_n}, \Sigma)$$

- Don't know μ_1, μ_2

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

- Don't know ρ_1, ρ_2

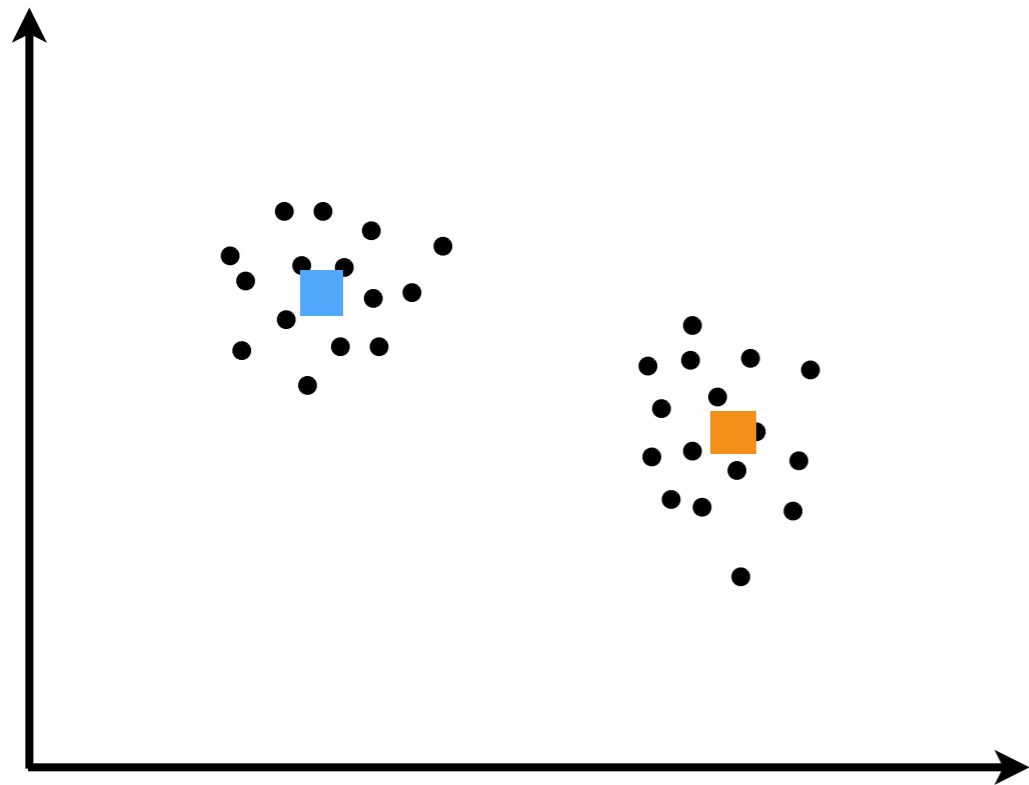


ρ_1

ρ_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)

$$z_n \stackrel{iid}{\sim} \text{Categorical}(\rho_1, \rho_2)$$
$$x_n \stackrel{indep}{\sim} \mathcal{N}(\mu_{z_n}, \Sigma)$$

- Don't know μ_1, μ_2

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

- Don't know ρ_1, ρ_2

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

$$\rho_2 = 1 - \rho_1$$

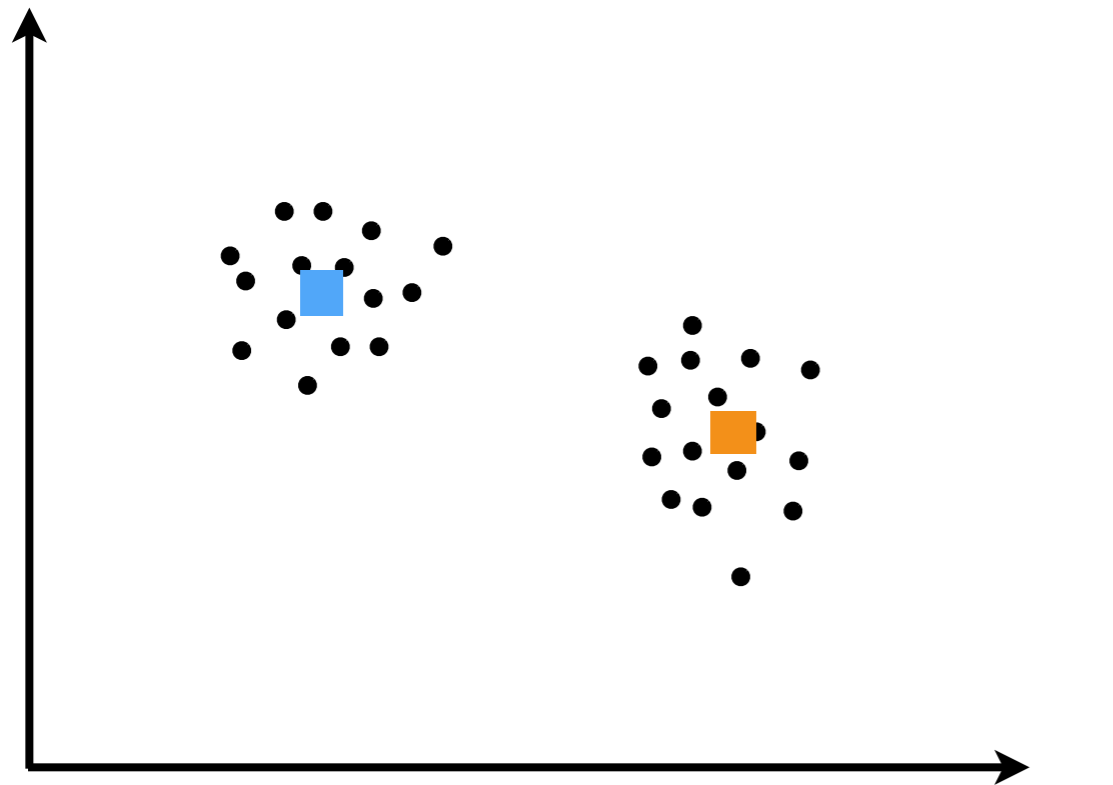


ρ_1

ρ_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)

$$z_n \stackrel{iid}{\sim} \text{Categorical}(\rho_1, \rho_2)$$
$$x_n \stackrel{indep}{\sim} \mathcal{N}(\mu_{z_n}, \Sigma)$$

- Don't know μ_1, μ_2

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

- Don't know ρ_1, ρ_2

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

$$\rho_2 = 1 - \rho_1$$

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



ρ_1

ρ_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)

$$z_n \stackrel{iid}{\sim} \text{Categorical}(\rho_1, \rho_2)$$
$$x_n \stackrel{indep}{\sim} \mathcal{N}(\mu_{z_n}, \Sigma)$$

- Don't know μ_1, μ_2

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

- Don't know ρ_1, ρ_2

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

$$\rho_2 = 1 - \rho_1$$

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



Generative model

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



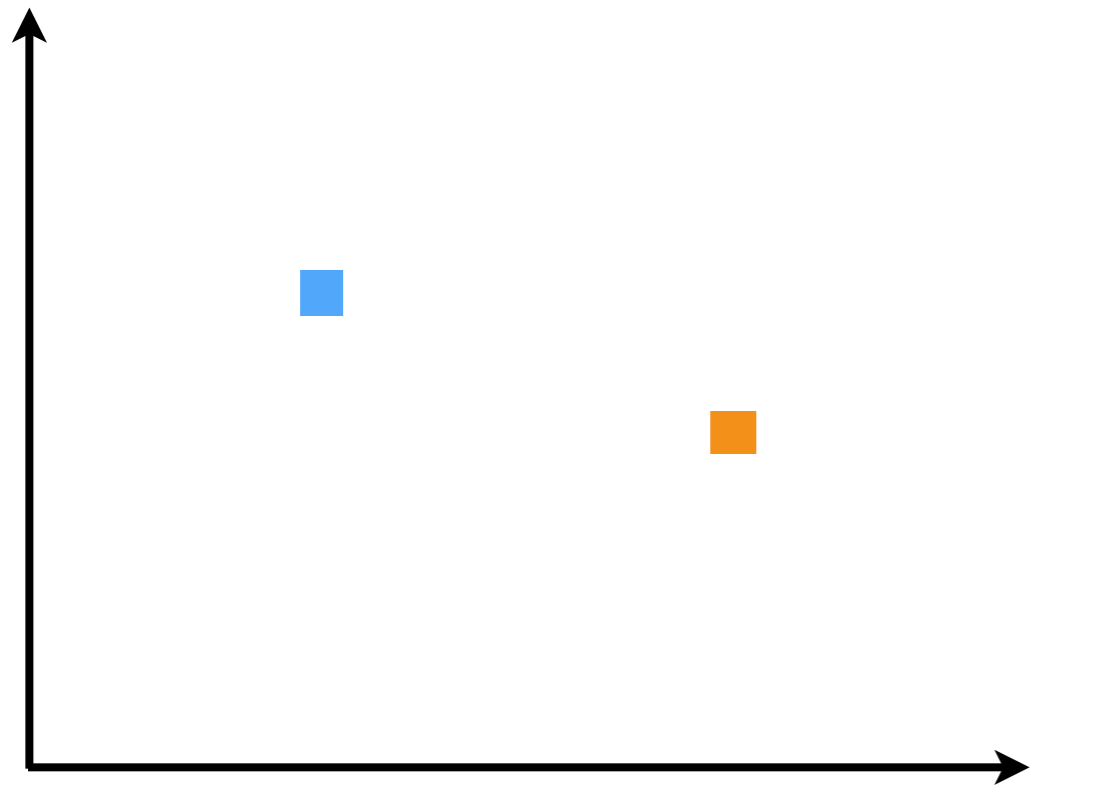
ρ_1

ρ_2

- Finite Gaussian mixture model ($K=2$ clusters)
$$z_n \stackrel{iid}{\sim} \text{Categorical}(\rho_1, \rho_2)$$
$$x_n \stackrel{indep}{\sim} \mathcal{N}(\mu_{z_n}, \Sigma)$$
- Don't know μ_1, μ_2
$$\mu_k \stackrel{iid}{\sim} \mathcal{N}(\mu_0, \Sigma_0)$$
- Don't know ρ_1, ρ_2
$$\rho_1 \sim \text{Beta}(a_1, a_2)$$
$$\rho_2 = 1 - \rho_1$$
- Inference goal: assignments of data points to clusters, cluster parameters

Generative model

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



ρ_1

ρ_2

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

$$z_n \stackrel{iid}{\sim} \text{Categorical}(\rho_1, \rho_2)$$
$$x_n \stackrel{indep}{\sim} \mathcal{N}(\mu_{z_n}, \Sigma)$$

- Don't know μ_1, μ_2

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

- Don't know ρ_1, ρ_2

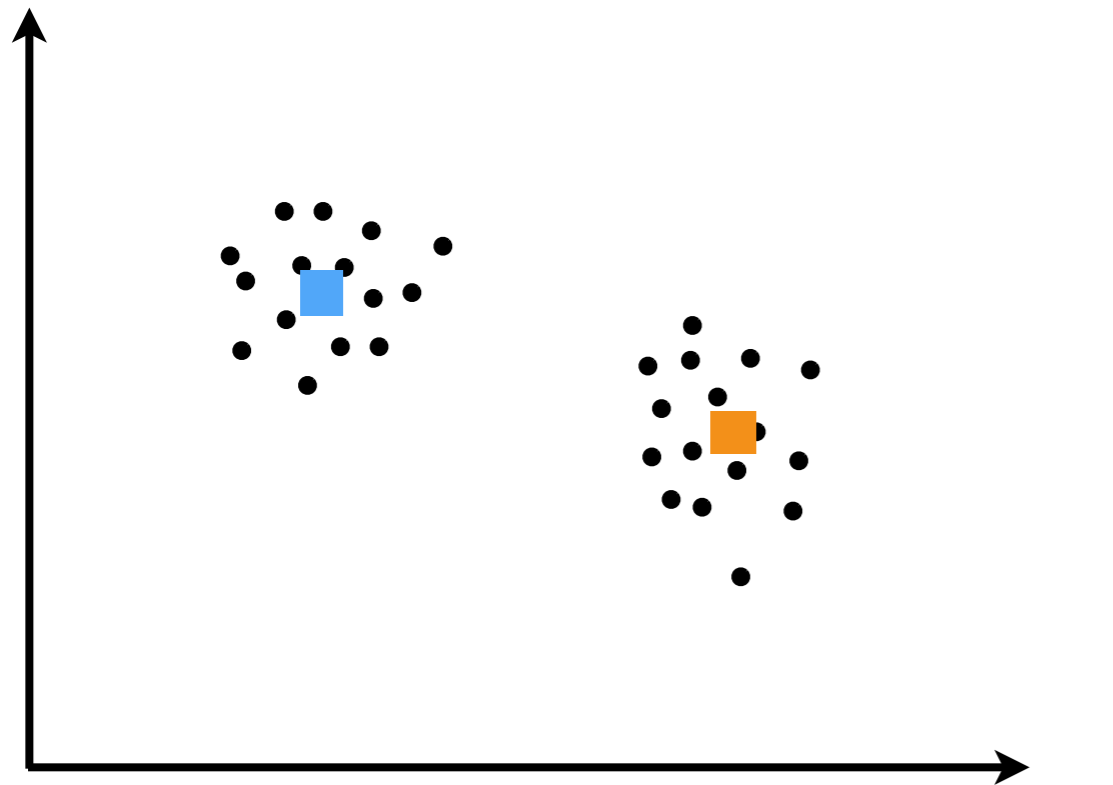
$$\rho_1 \sim \text{Beta}(a_1, a_2)$$

$$\rho_2 = 1 - \rho_1$$

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

Generative model

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



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

$$z_n \stackrel{iid}{\sim} \text{Categorical}(\rho_1, \rho_2)$$
$$x_n \stackrel{indep}{\sim} \mathcal{N}(\mu_{z_n}, \Sigma)$$

- Don't know μ_1, μ_2

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

- Don't know ρ_1, ρ_2

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

$$\rho_2 = 1 - \rho_1$$

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



ρ_1

ρ_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

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

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

$$a_1, a_2 > 0$$

- Gamma function Γ

Beta distribution review

$$\text{Beta}(\rho_1 | a_1, a_2) = \frac{\Gamma(a_1 + a_2)}{\Gamma(a_1)\Gamma(a_2)} \rho_1^{a_1-1} (1 - \rho_1)^{a_2-1} \quad \begin{array}{l} \rho_1 \in (0, 1) \\ a_1, a_2 > 0 \end{array}$$

- Gamma function Γ
 - integer m : $\Gamma(m + 1) = m!$

Beta distribution review

$$\text{Beta}(\rho_1 | a_1, a_2) = \frac{\Gamma(a_1 + a_2)}{\Gamma(a_1)\Gamma(a_2)} \rho_1^{a_1-1} (1 - \rho_1)^{a_2-1} \quad \begin{array}{l} \rho_1 \in (0, 1) \\ a_1, a_2 > 0 \end{array}$$

- Gamma function Γ
 - integer m : $\Gamma(m + 1) = m!$
 - for $x > 0$: $\Gamma(x + 1) = x\Gamma(x)$

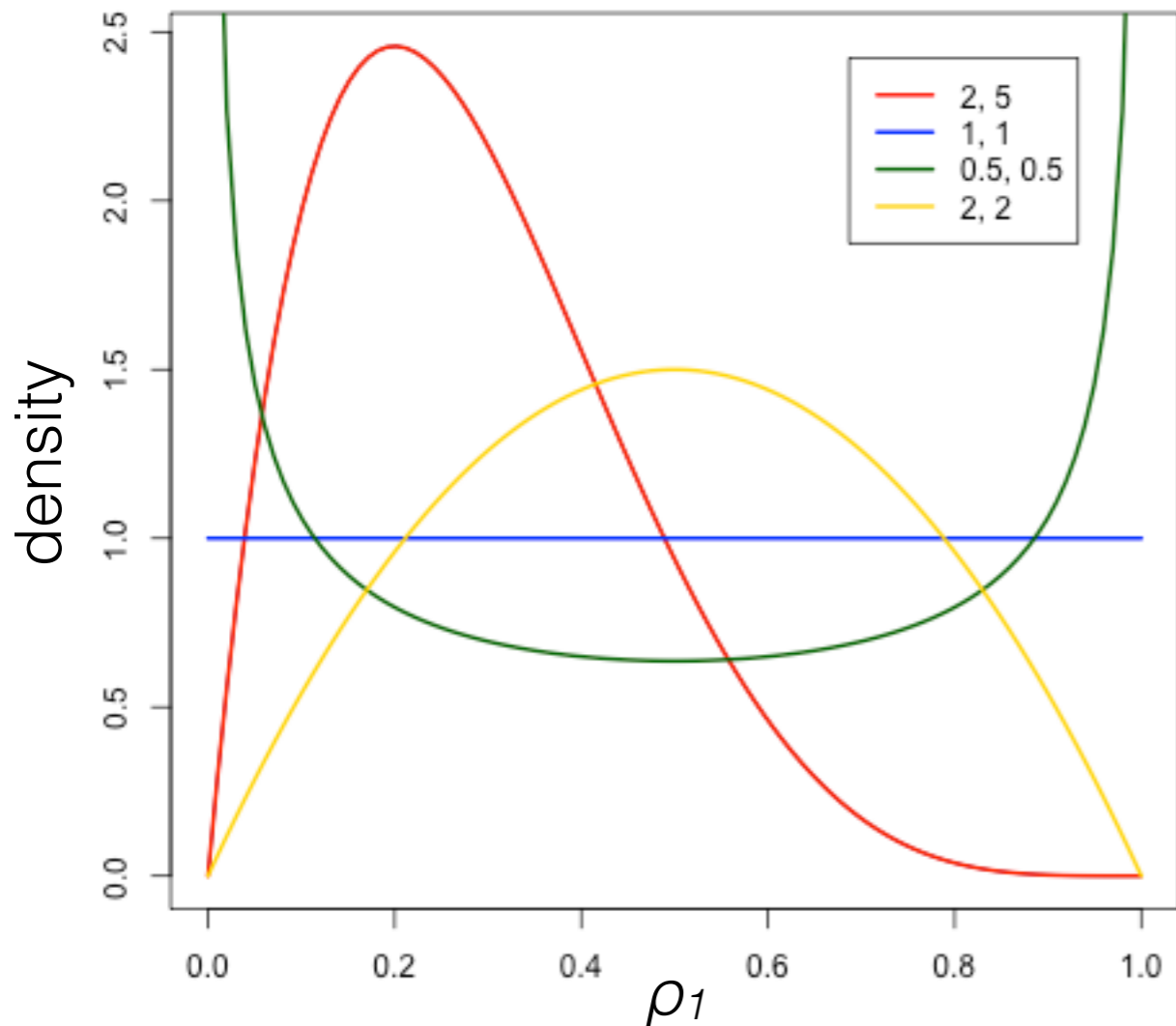
Beta distribution review

$$\text{Beta}(\rho_1 | a_1, a_2) = \frac{\Gamma(a_1 + a_2)}{\Gamma(a_1)\Gamma(a_2)} \rho_1^{a_1-1} (1 - \rho_1)^{a_2-1} \quad \begin{array}{l} \rho_1 \in (0, 1) \\ a_1, a_2 > 0 \end{array}$$

- Gamma function Γ
 - integer m : $\Gamma(m + 1) = m!$
 - for $x > 0$: $\Gamma(x + 1) = x\Gamma(x)$
- What happens?

Beta distribution review

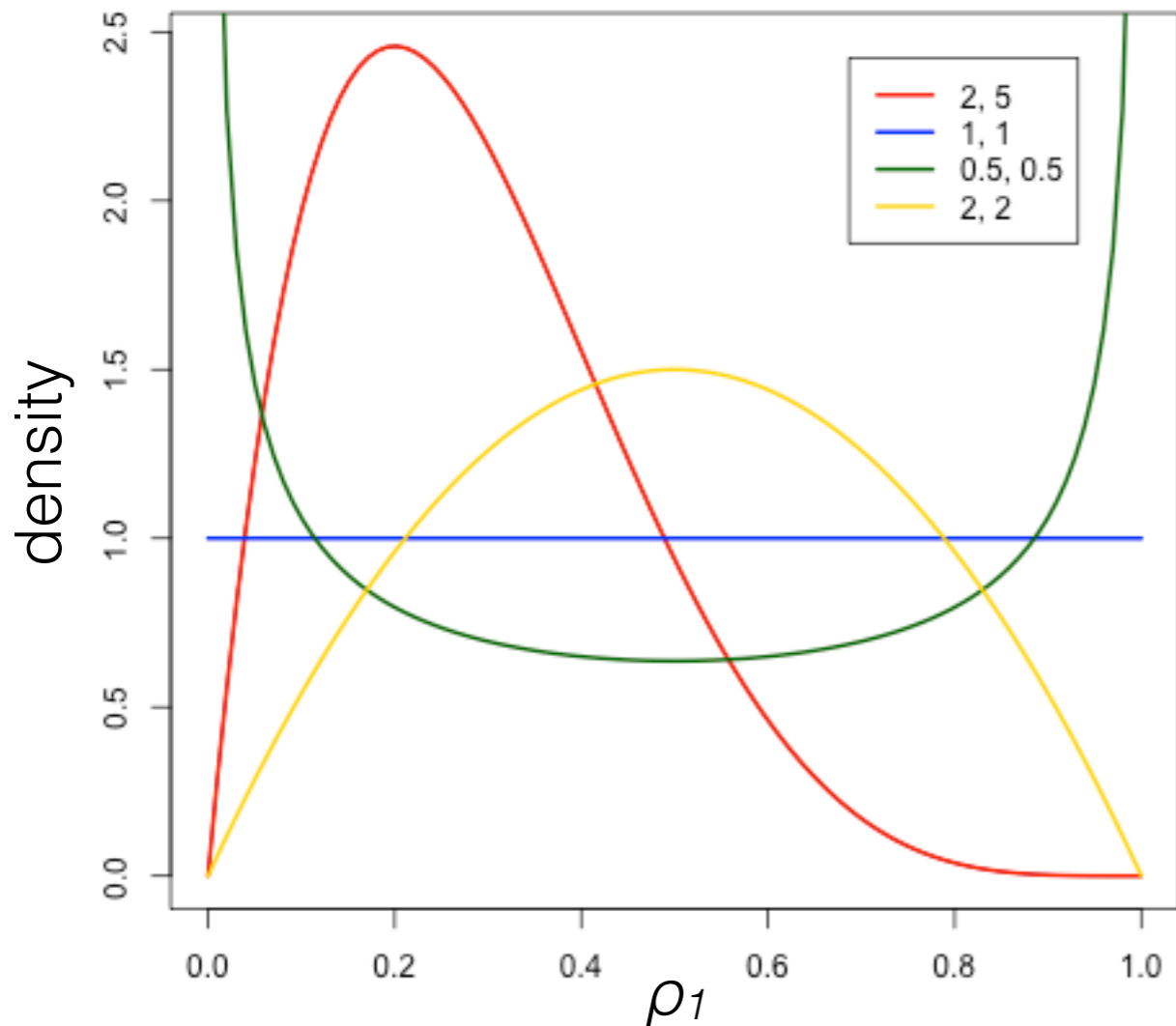
$$\text{Beta}(\rho_1 | a_1, a_2) = \frac{\Gamma(a_1 + a_2)}{\Gamma(a_1)\Gamma(a_2)} \rho_1^{a_1-1} (1 - \rho_1)^{a_2-1} \quad \begin{array}{l} \rho_1 \in (0, 1) \\ a_1, a_2 > 0 \end{array}$$



- Gamma function Γ
 - integer m : $\Gamma(m + 1) = m!$
 - for $x > 0$: $\Gamma(x + 1) = x\Gamma(x)$
- What happens?

Beta distribution review

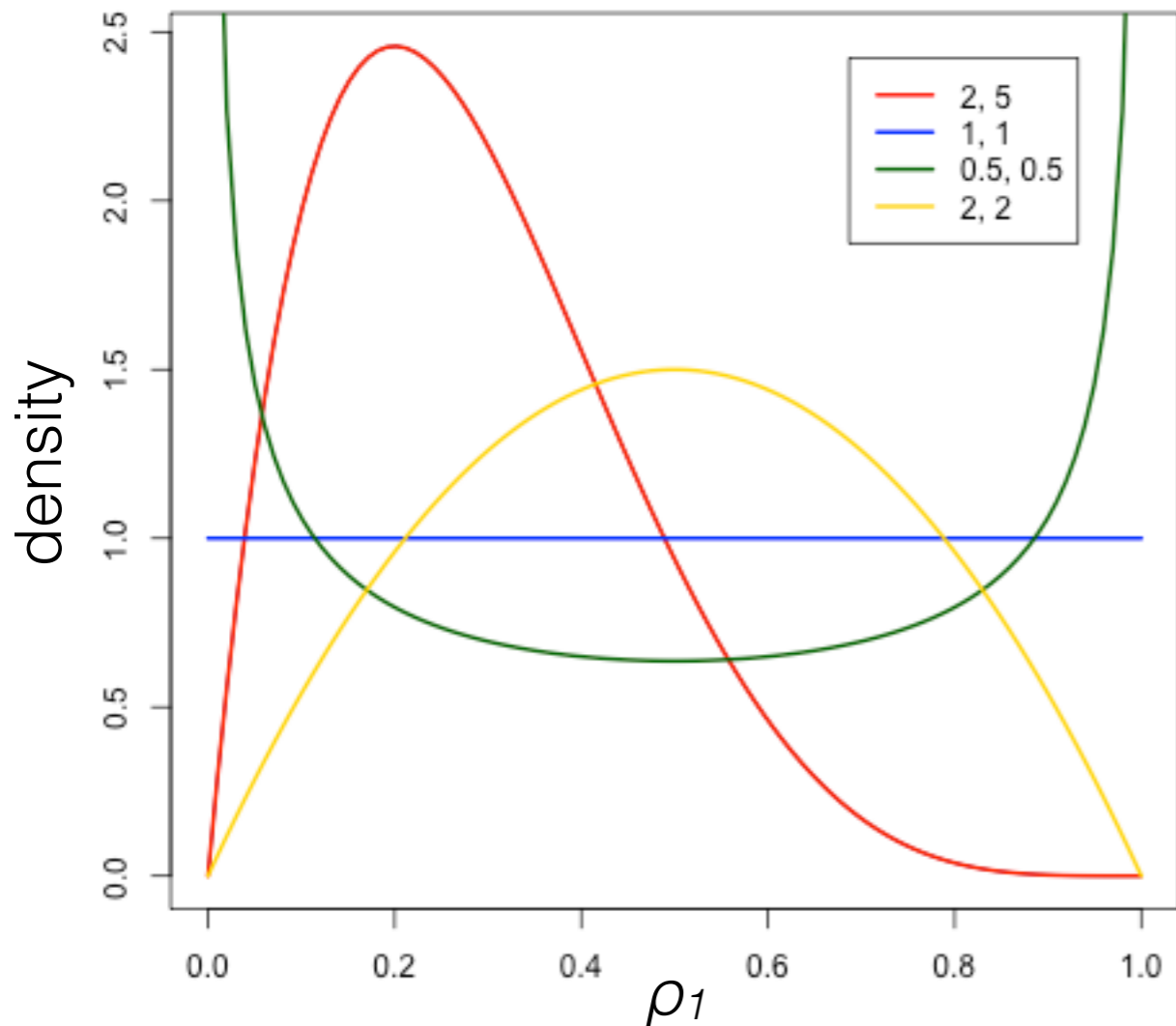
$$\text{Beta}(\rho_1 | a_1, a_2) = \frac{\Gamma(a_1 + a_2)}{\Gamma(a_1)\Gamma(a_2)} \rho_1^{a_1-1} (1 - \rho_1)^{a_2-1} \quad \begin{array}{l} \rho_1 \in (0, 1) \\ a_1, a_2 > 0 \end{array}$$



- Gamma function Γ
 - integer m : $\Gamma(m + 1) = m!$
 - for $x > 0$: $\Gamma(x + 1) = x\Gamma(x)$
- What happens?
 - $a = a_1 = a_2 \rightarrow 0$

Beta distribution review

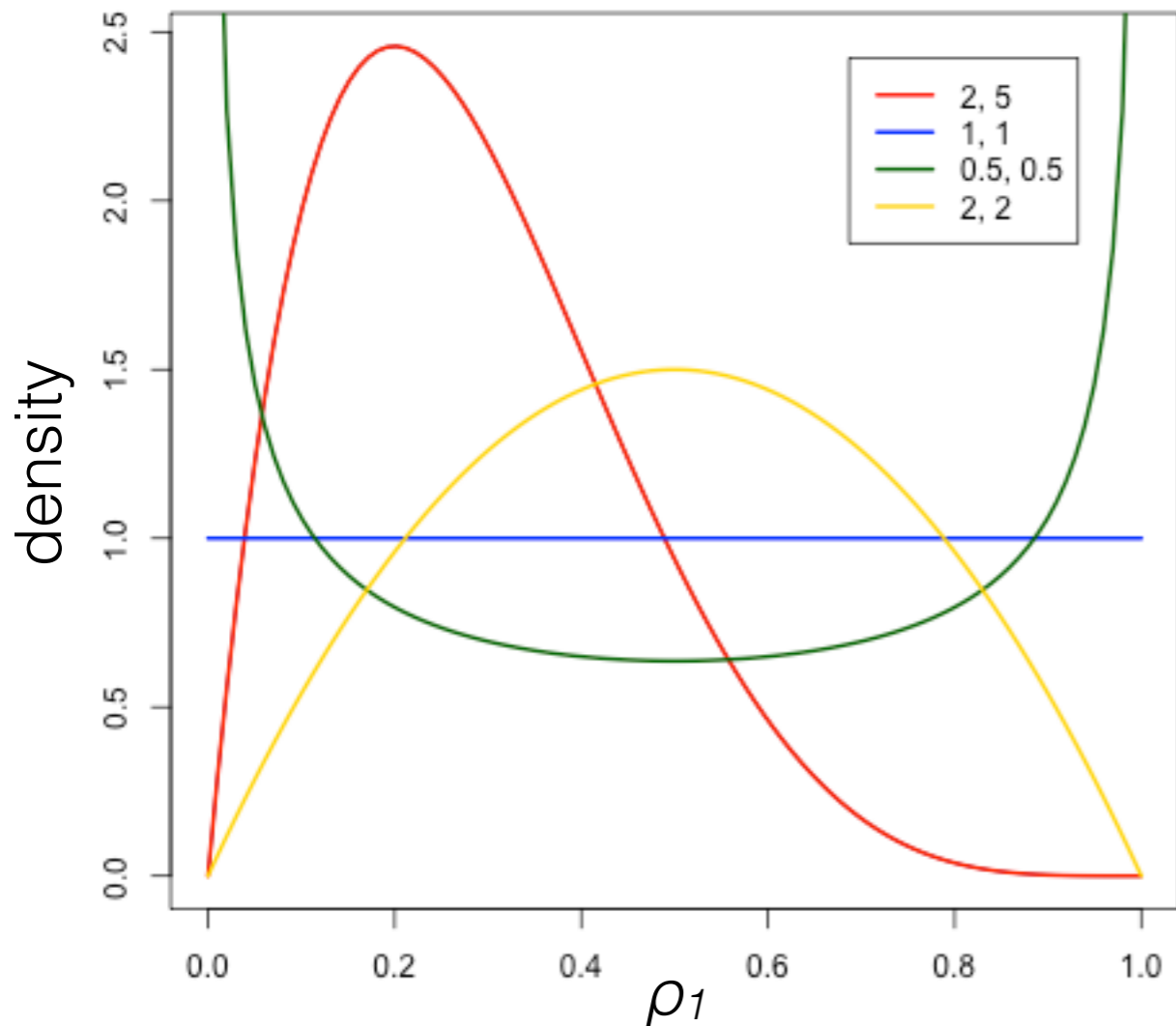
$$\text{Beta}(\rho_1 | a_1, a_2) = \frac{\Gamma(a_1 + a_2)}{\Gamma(a_1)\Gamma(a_2)} \rho_1^{a_1-1} (1 - \rho_1)^{a_2-1} \quad \begin{array}{l} \rho_1 \in (0, 1) \\ a_1, a_2 > 0 \end{array}$$



- Gamma function Γ
 - integer m : $\Gamma(m + 1) = m!$
 - for $x > 0$: $\Gamma(x + 1) = x\Gamma(x)$
- What happens?
 - $a = a_1 = a_2 \rightarrow 0$
 - $a = a_1 = a_2 \rightarrow \infty$

Beta distribution review

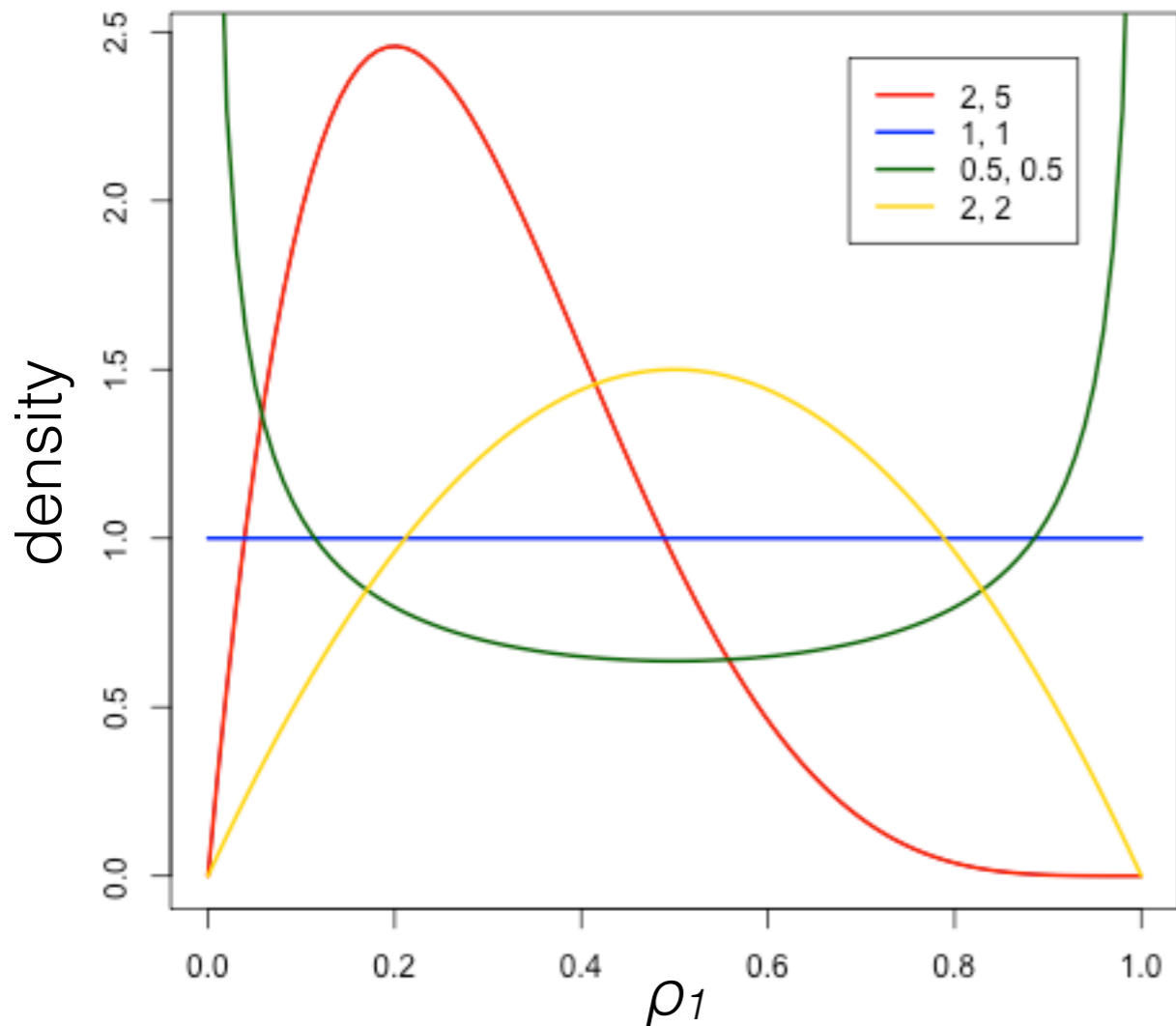
$$\text{Beta}(\rho_1 | a_1, a_2) = \frac{\Gamma(a_1 + a_2)}{\Gamma(a_1)\Gamma(a_2)} \rho_1^{a_1-1} (1 - \rho_1)^{a_2-1} \quad \begin{array}{l} \rho_1 \in (0, 1) \\ a_1, a_2 > 0 \end{array}$$



- Gamma function Γ
 - integer m : $\Gamma(m + 1) = m!$
 - for $x > 0$: $\Gamma(x + 1) = x\Gamma(x)$
- What happens?
 - $a = a_1 = a_2 \rightarrow 0$
 - $a = a_1 = a_2 \rightarrow \infty$
 - $a_1 > a_2$

Beta distribution review

$$\text{Beta}(\rho_1 | a_1, a_2) = \frac{\Gamma(a_1 + a_2)}{\Gamma(a_1)\Gamma(a_2)} \rho_1^{a_1-1} (1 - \rho_1)^{a_2-1} \quad \begin{array}{l} \rho_1 \in (0, 1) \\ a_1, a_2 > 0 \end{array}$$

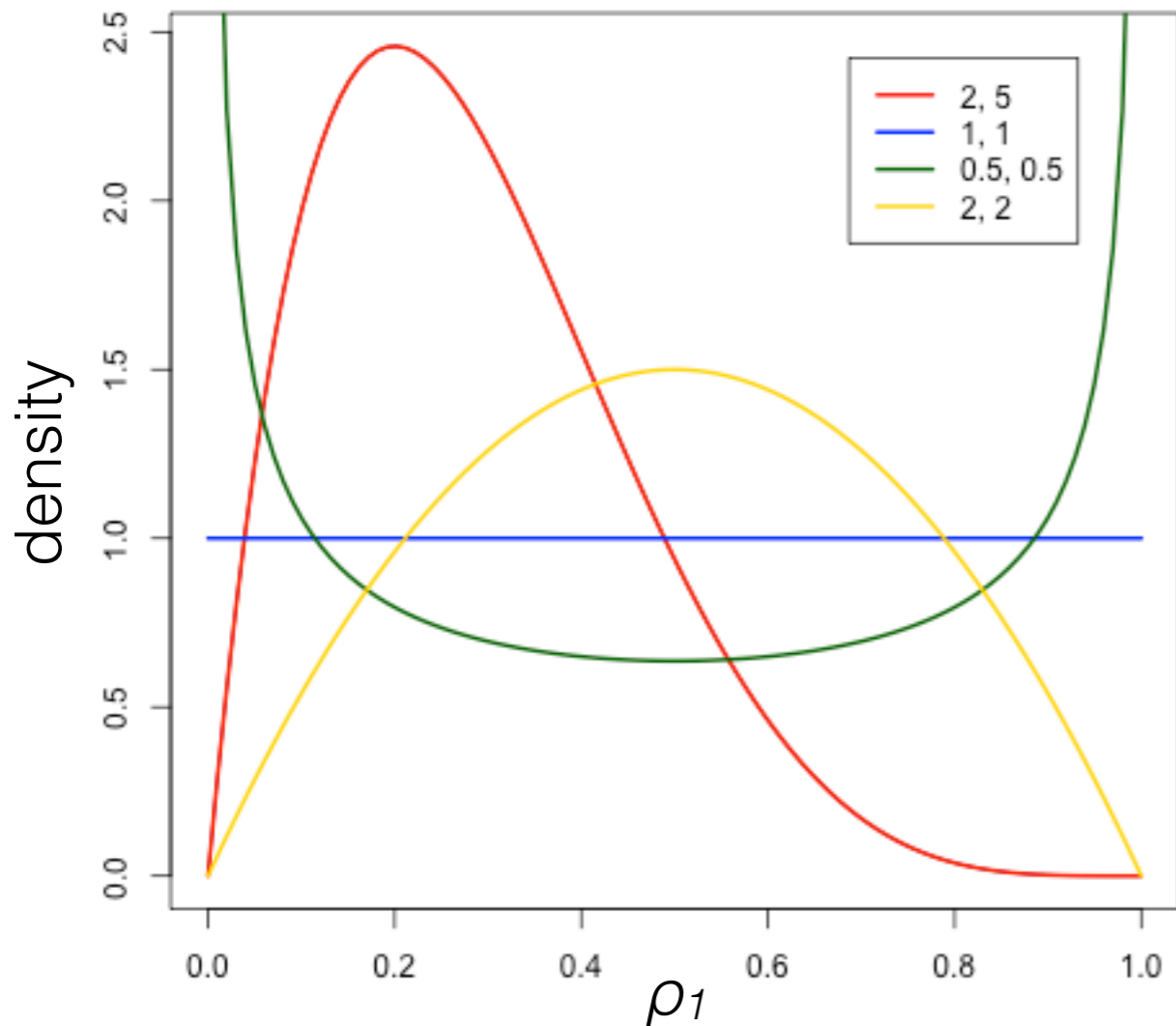


- Gamma function Γ
 - integer m : $\Gamma(m + 1) = m!$
 - for $x > 0$: $\Gamma(x + 1) = x\Gamma(x)$
- What happens?
 - $a = a_1 = a_2 \rightarrow 0$
 - $a = a_1 = a_2 \rightarrow \infty$
 - $a_1 > a_2$

[demo]

Beta distribution review

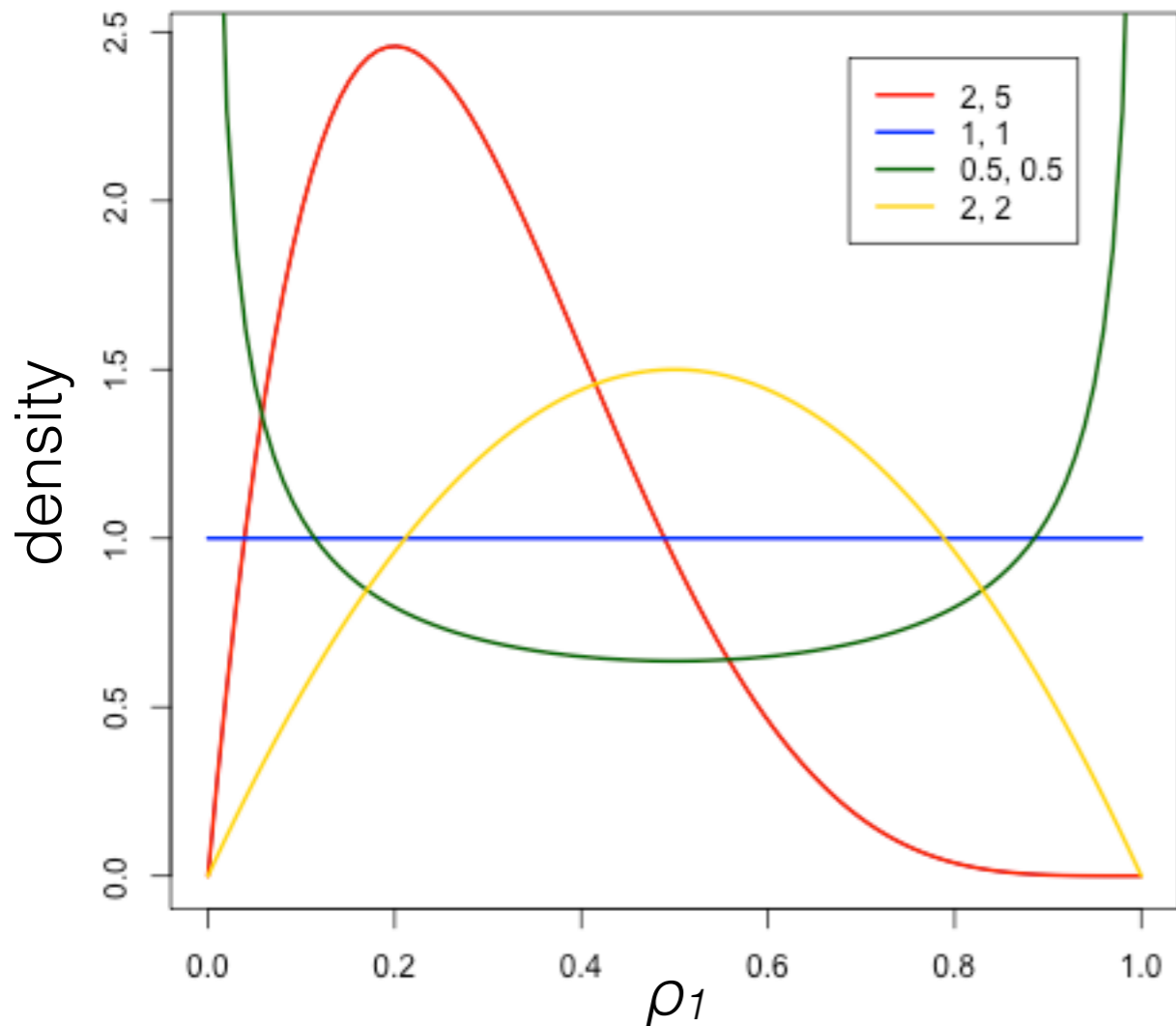
$$\text{Beta}(\rho_1 | a_1, a_2) = \frac{\Gamma(a_1 + a_2)}{\Gamma(a_1)\Gamma(a_2)} \rho_1^{a_1-1} (1 - \rho_1)^{a_2-1} \quad \begin{array}{l} \rho_1 \in (0, 1) \\ a_1, a_2 > 0 \end{array}$$



- Gamma function Γ
 - integer m : $\Gamma(m + 1) = m!$
 - for $x > 0$: $\Gamma(x + 1) = x\Gamma(x)$
- What happens?
 - $a = a_1 = a_2 \rightarrow 0$
 - $a = a_1 = a_2 \rightarrow \infty$
 - $a_1 > a_2$ [demo]
- Beta is conjugate to Cat

Beta distribution review

$$\text{Beta}(\rho_1 | a_1, a_2) = \frac{\Gamma(a_1 + a_2)}{\Gamma(a_1)\Gamma(a_2)} \rho_1^{a_1-1} (1 - \rho_1)^{a_2-1} \quad \begin{array}{l} \rho_1 \in (0, 1) \\ a_1, a_2 > 0 \end{array}$$

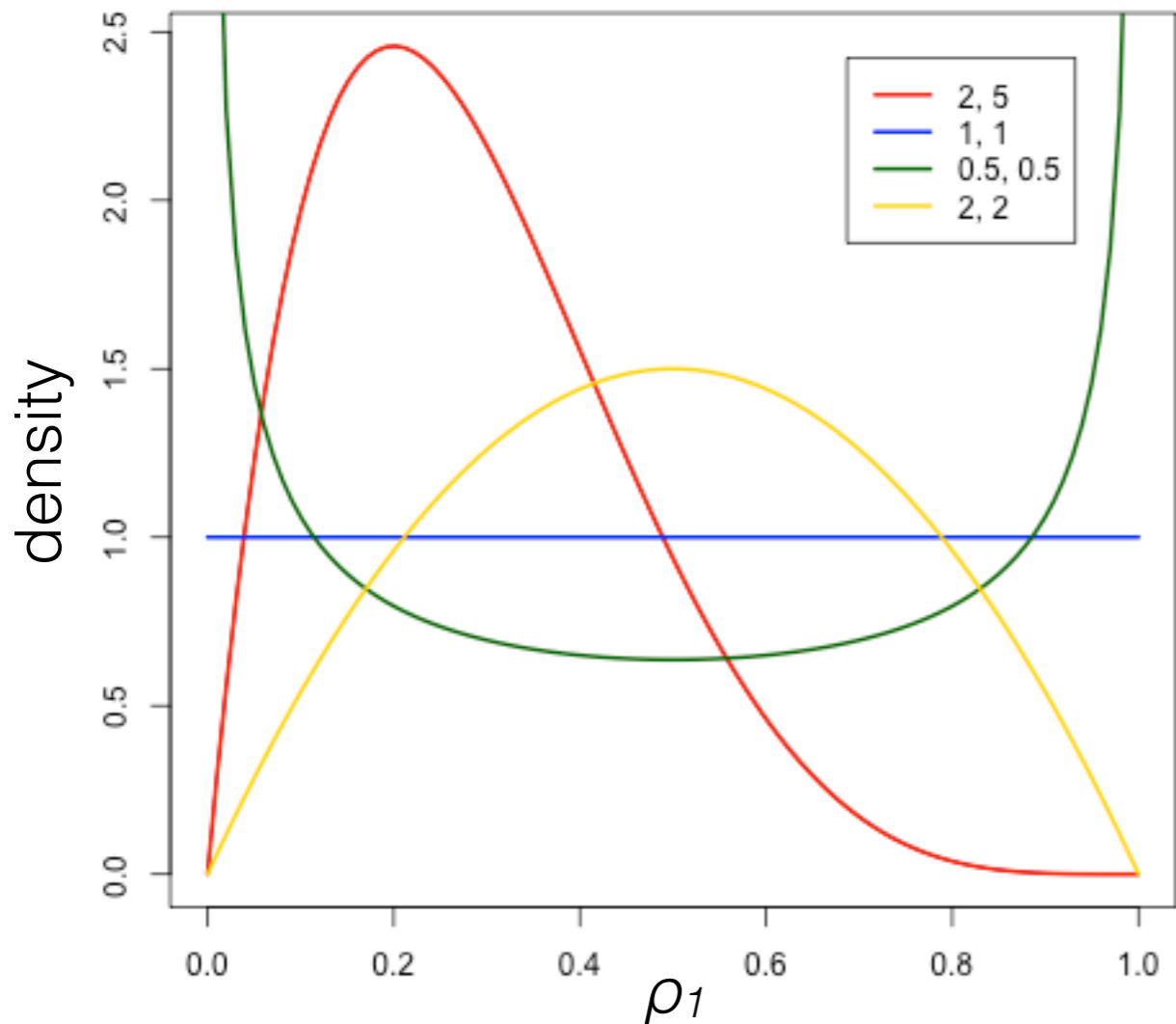


- Gamma function Γ
 - integer m : $\Gamma(m + 1) = m!$
 - for $x > 0$: $\Gamma(x + 1) = x\Gamma(x)$
- What happens?
 - $a = a_1 = a_2 \rightarrow 0$
 - $a = a_1 = a_2 \rightarrow \infty$
 - $a_1 > a_2$ [demo]
- Beta is conjugate to Cat

$\rho_1 \sim \text{Beta}(a_1, a_2), z \sim \text{Cat}(\rho_1, \rho_2)$

Beta distribution review

$$\text{Beta}(\rho_1 | a_1, a_2) = \frac{\Gamma(a_1 + a_2)}{\Gamma(a_1)\Gamma(a_2)} \rho_1^{a_1-1} (1 - \rho_1)^{a_2-1} \quad \begin{array}{l} \rho_1 \in (0, 1) \\ a_1, a_2 > 0 \end{array}$$



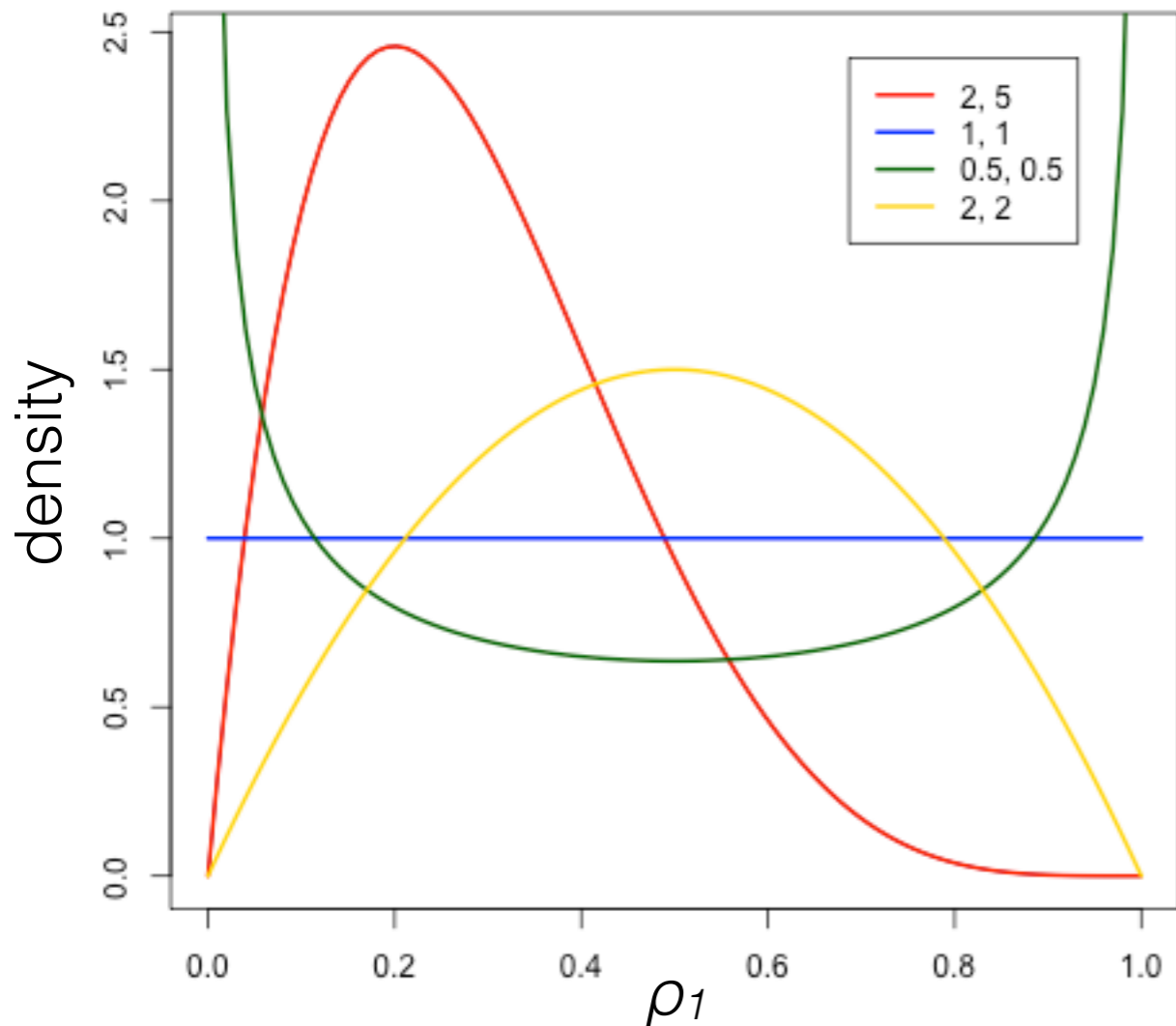
- Gamma function Γ
 - integer m : $\Gamma(m + 1) = m!$
 - for $x > 0$: $\Gamma(x + 1) = x\Gamma(x)$
- What happens?
 - $a = a_1 = a_2 \rightarrow 0$
 - $a = a_1 = a_2 \rightarrow \infty$
 - $a_1 > a_2$ [demo]
- Beta is conjugate to Cat

$\rho_1 \sim \text{Beta}(a_1, a_2), z \sim \text{Cat}(\rho_1, \rho_2)$

$$p(\rho_1, z) \propto$$

Beta distribution review

$$\text{Beta}(\rho_1 | a_1, a_2) = \frac{\Gamma(a_1 + a_2)}{\Gamma(a_1)\Gamma(a_2)} \rho_1^{a_1-1} (1 - \rho_1)^{a_2-1} \quad \begin{array}{l} \rho_1 \in (0, 1) \\ a_1, a_2 > 0 \end{array}$$

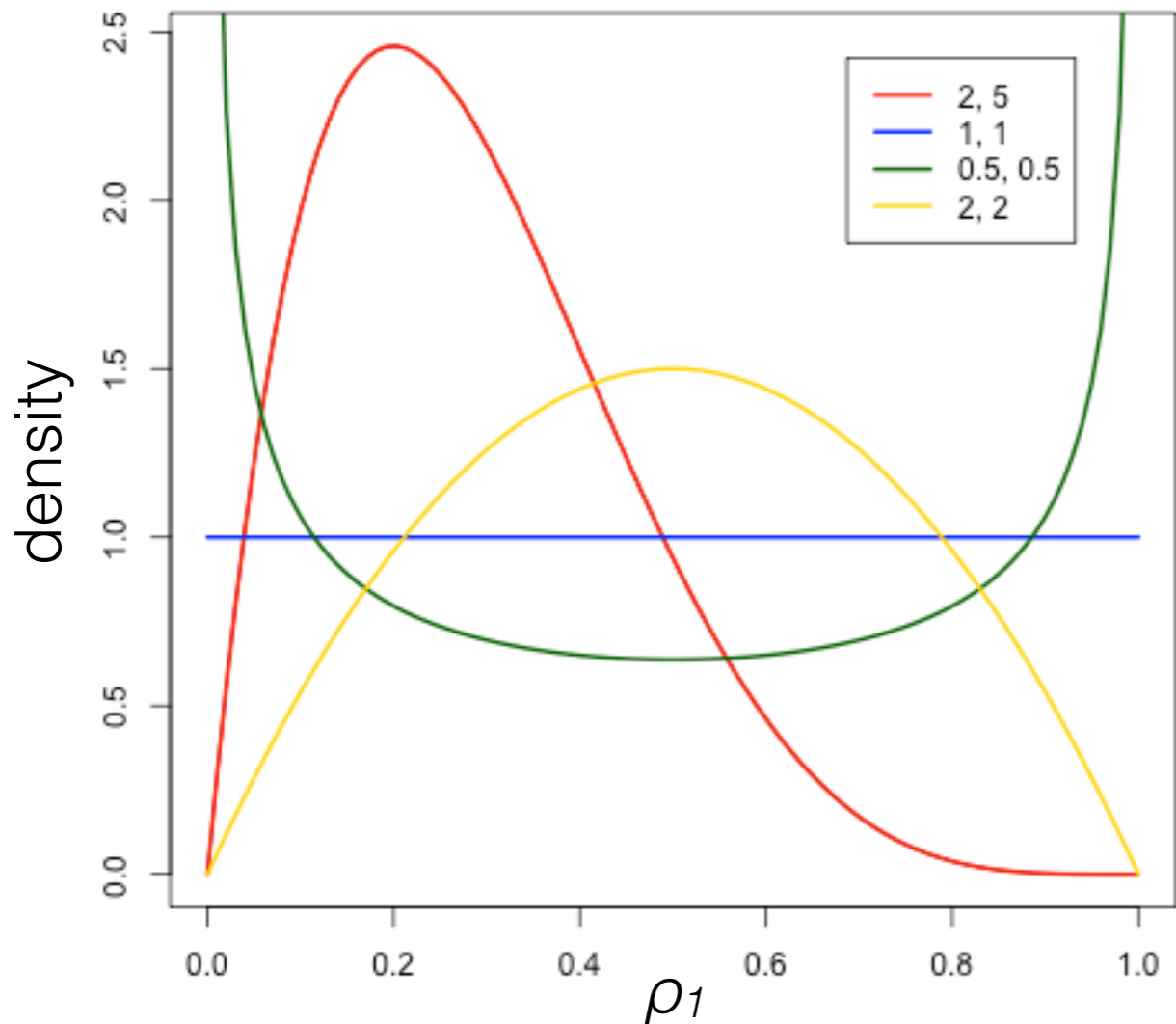


- Gamma function Γ
 - integer m : $\Gamma(m + 1) = m!$
 - for $x > 0$: $\Gamma(x + 1) = x\Gamma(x)$
 - What happens?
 - $a = a_1 = a_2 \rightarrow 0$
 - $a = a_1 = a_2 \rightarrow \infty$
 - $a_1 > a_2$ [demo]
 - Beta is conjugate to Cat
- $\rho_1 \sim \text{Beta}(a_1, a_2), z \sim \text{Cat}(\rho_1, \rho_2)$

$$p(\rho_1, z) \propto \rho_1^{\mathbf{1}\{z=1\}} (1 - \rho_1)^{\mathbf{1}\{z=2\}}$$

Beta distribution review

$$\text{Beta}(\rho_1 | a_1, a_2) = \frac{\Gamma(a_1 + a_2)}{\Gamma(a_1)\Gamma(a_2)} \rho_1^{a_1-1} (1 - \rho_1)^{a_2-1} \quad \begin{array}{l} \rho_1 \in (0, 1) \\ a_1, a_2 > 0 \end{array}$$



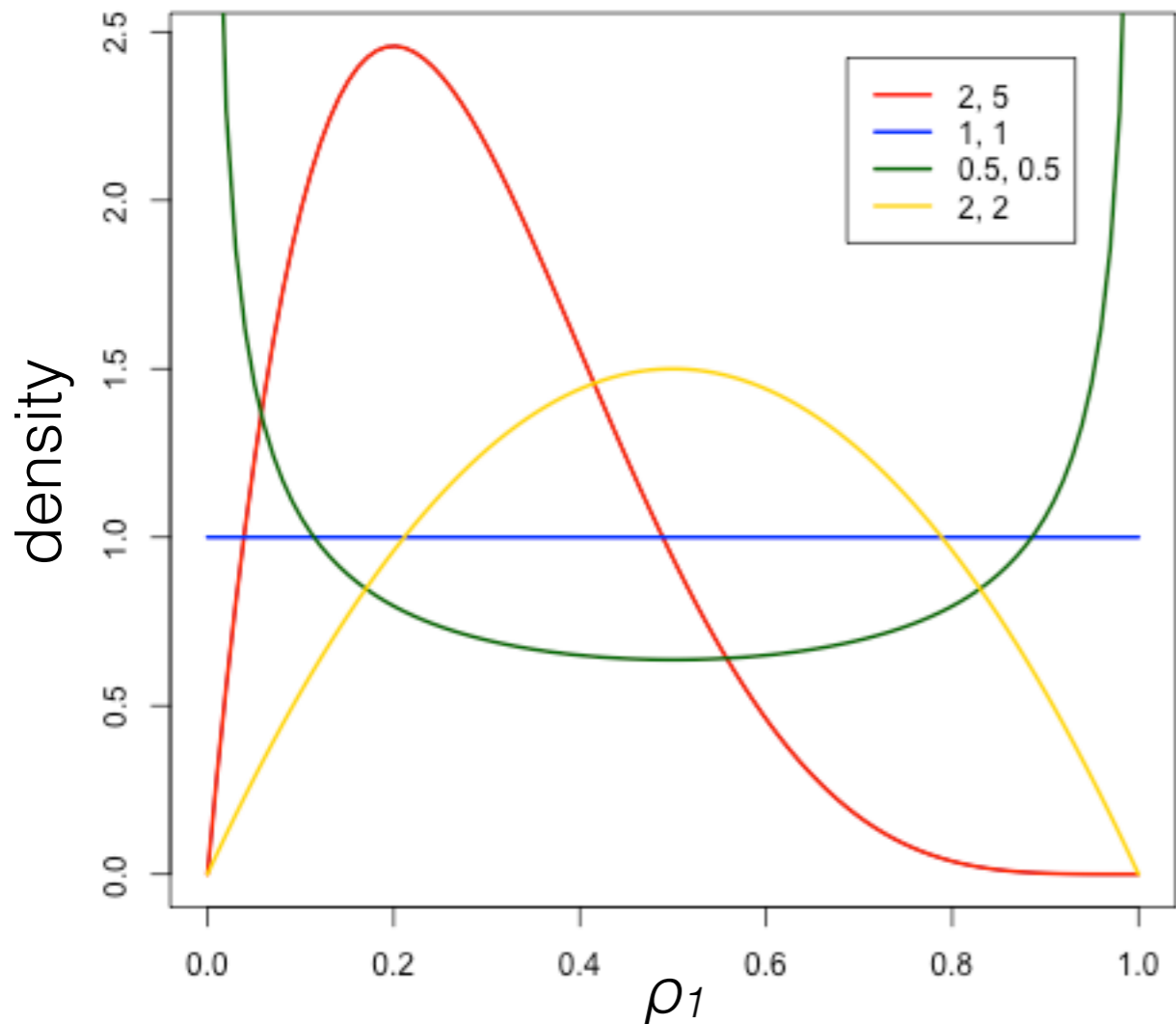
- Gamma function Γ
 - integer m : $\Gamma(m + 1) = m!$
 - for $x > 0$: $\Gamma(x + 1) = x\Gamma(x)$
- What happens?
 - $a = a_1 = a_2 \rightarrow 0$
 - $a = a_1 = a_2 \rightarrow \infty$
 - $a_1 > a_2$ [demo]
- Beta is conjugate to Cat

$$\rho_1 \sim \text{Beta}(a_1, a_2), z \sim \text{Cat}(\rho_1, \rho_2)$$

$$p(\rho_1, z) \propto \rho_1^{\mathbf{1}\{z=1\}} (1 - \rho_1)^{\mathbf{1}\{z=2\}} \rho_1^{a_1-1} (1 - \rho_1)^{a_2-1}$$

Beta distribution review

$$\text{Beta}(\rho_1 | a_1, a_2) = \frac{\Gamma(a_1 + a_2)}{\Gamma(a_1)\Gamma(a_2)} \rho_1^{a_1-1} (1 - \rho_1)^{a_2-1} \quad \begin{array}{l} \rho_1 \in (0, 1) \\ a_1, a_2 > 0 \end{array}$$



- Gamma function Γ
 - integer m : $\Gamma(m + 1) = m!$
 - for $x > 0$: $\Gamma(x + 1) = x\Gamma(x)$
- What happens?
 - $a = a_1 = a_2 \rightarrow 0$
 - $a = a_1 = a_2 \rightarrow \infty$
 - $a_1 > a_2$ [demo]
- Beta is conjugate to Cat

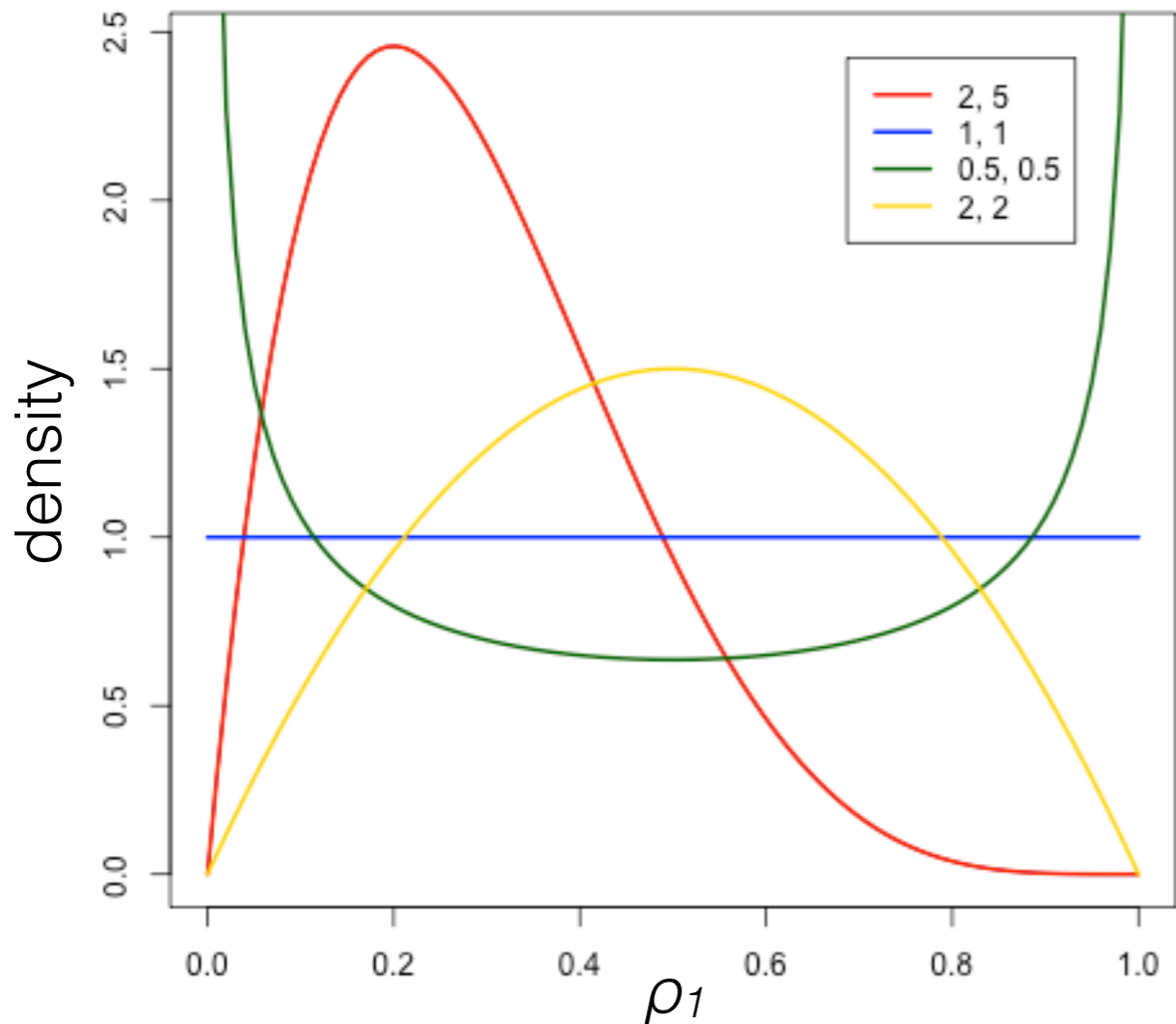
$$\rho_1 \sim \text{Beta}(a_1, a_2), z \sim \text{Cat}(\rho_1, \rho_2)$$

$$p(\rho_1, z) \propto \rho_1^{\mathbf{1}\{z=1\}} (1 - \rho_1)^{\mathbf{1}\{z=2\}} \rho_1^{a_1-1} (1 - \rho_1)^{a_2-1}$$

$$p(\rho_1 | z) \propto$$

Beta distribution review

$$\text{Beta}(\rho_1 | a_1, a_2) = \frac{\Gamma(a_1 + a_2)}{\Gamma(a_1)\Gamma(a_2)} \rho_1^{a_1-1} (1 - \rho_1)^{a_2-1} \quad \begin{array}{l} \rho_1 \in (0, 1) \\ a_1, a_2 > 0 \end{array}$$



- Gamma function Γ
 - integer m : $\Gamma(m + 1) = m!$
 - for $x > 0$: $\Gamma(x + 1) = x\Gamma(x)$
- What happens?
 - $a = a_1 = a_2 \rightarrow 0$
 - $a = a_1 = a_2 \rightarrow \infty$
 - $a_1 > a_2$ [demo]
- Beta is conjugate to Cat

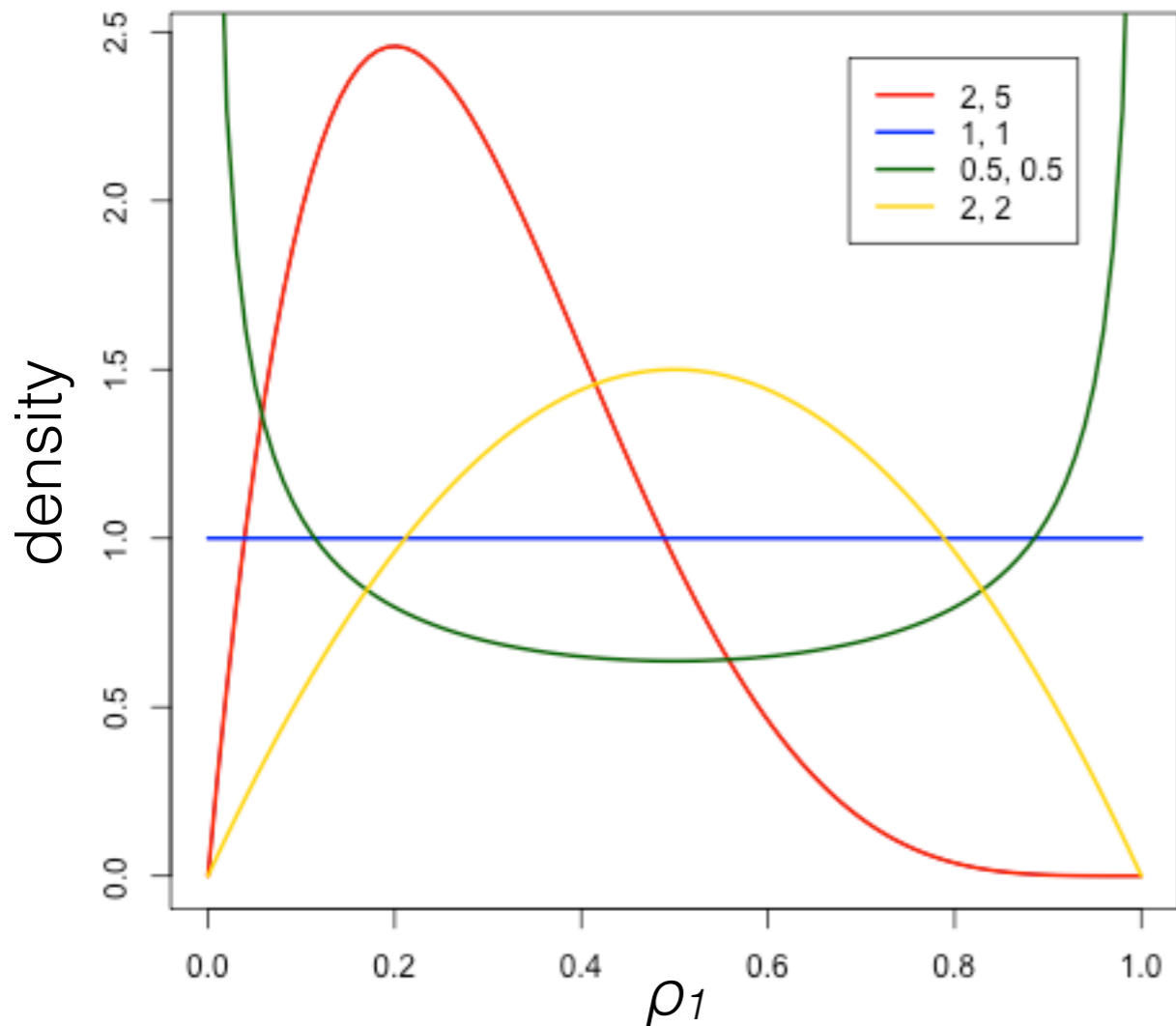
$$\rho_1 \sim \text{Beta}(a_1, a_2), z \sim \text{Cat}(\rho_1, \rho_2)$$

$$p(\rho_1, z) \propto \rho_1^{\mathbf{1}\{z=1\}} (1 - \rho_1)^{\mathbf{1}\{z=2\}} \rho_1^{a_1-1} (1 - \rho_1)^{a_2-1}$$

$$p(\rho_1 | z) \propto \rho_1^{a_1 + \mathbf{1}\{z=1\} - 1} (1 - \rho_1)^{a_2 + \mathbf{1}\{z=2\} - 1}$$

Beta distribution review

$$\text{Beta}(\rho_1 | a_1, a_2) = \frac{\Gamma(a_1 + a_2)}{\Gamma(a_1)\Gamma(a_2)} \rho_1^{a_1-1} (1 - \rho_1)^{a_2-1} \quad \begin{array}{l} \rho_1 \in (0, 1) \\ a_1, a_2 > 0 \end{array}$$



- Gamma function Γ
 - integer m : $\Gamma(m + 1) = m!$
 - for $x > 0$: $\Gamma(x + 1) = x\Gamma(x)$
- What happens?
 - $a = a_1 = a_2 \rightarrow 0$
 - $a = a_1 = a_2 \rightarrow \infty$
 - $a_1 > a_2$ [demo]
- Beta is conjugate to Cat

$$\rho_1 \sim \text{Beta}(a_1, a_2), z \sim \text{Cat}(\rho_1, \rho_2)$$

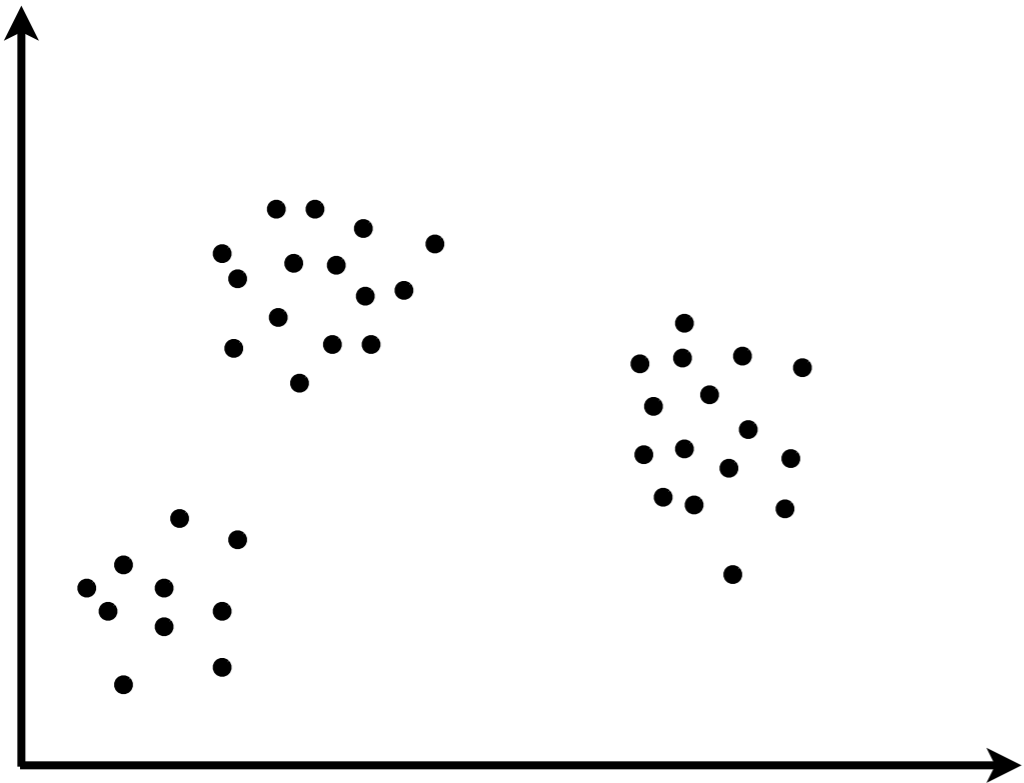
$$p(\rho_1, z) \propto \rho_1^{\mathbf{1}\{z=1\}} (1 - \rho_1)^{\mathbf{1}\{z=2\}} \rho_1^{a_1-1} (1 - \rho_1)^{a_2-1}$$

$$p(\rho_1 | z) \propto \rho_1^{a_1 + \mathbf{1}\{z=1\} - 1} (1 - \rho_1)^{a_2 + \mathbf{1}\{z=2\} - 1} \propto \text{Beta}(\rho_1 | a_1 + \mathbf{1}\{z = 1\}, a_2 + \mathbf{1}\{z = 2\})$$

Generative model

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

- Finite Gaussian mixture model (K 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)

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



ρ_1

ρ_2

ρ_3

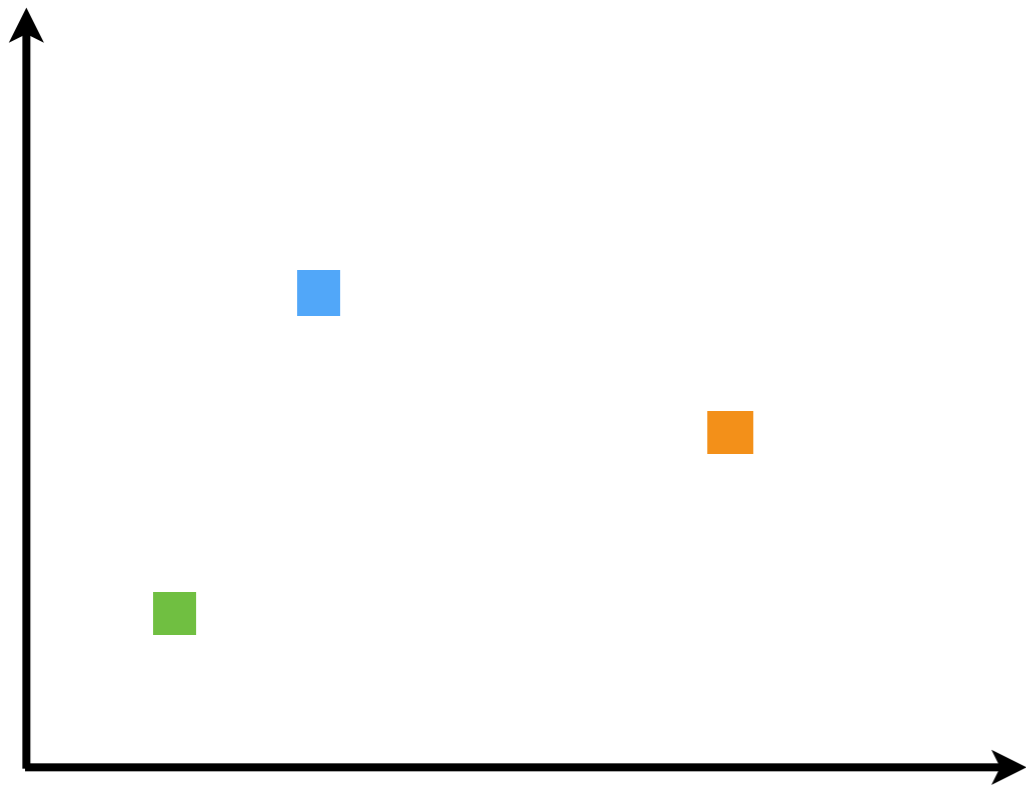
Generative model

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

- Finite Gaussian mixture model (K clusters)

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

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



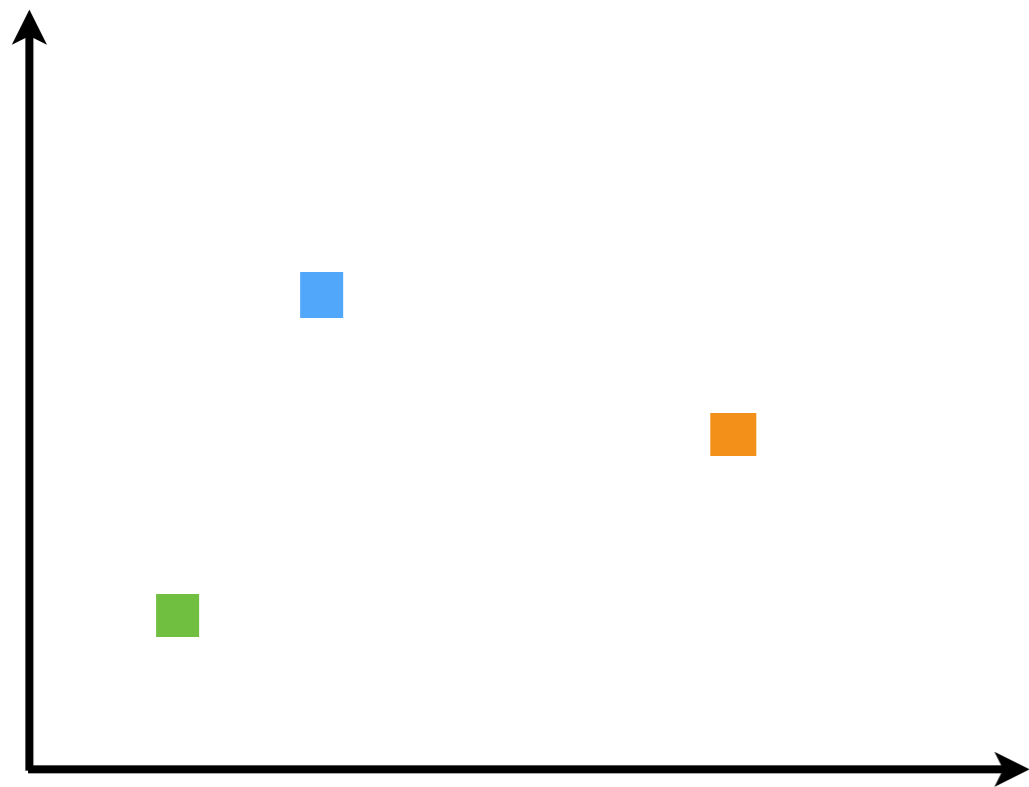
ρ_1

ρ_2

ρ_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)

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

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

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



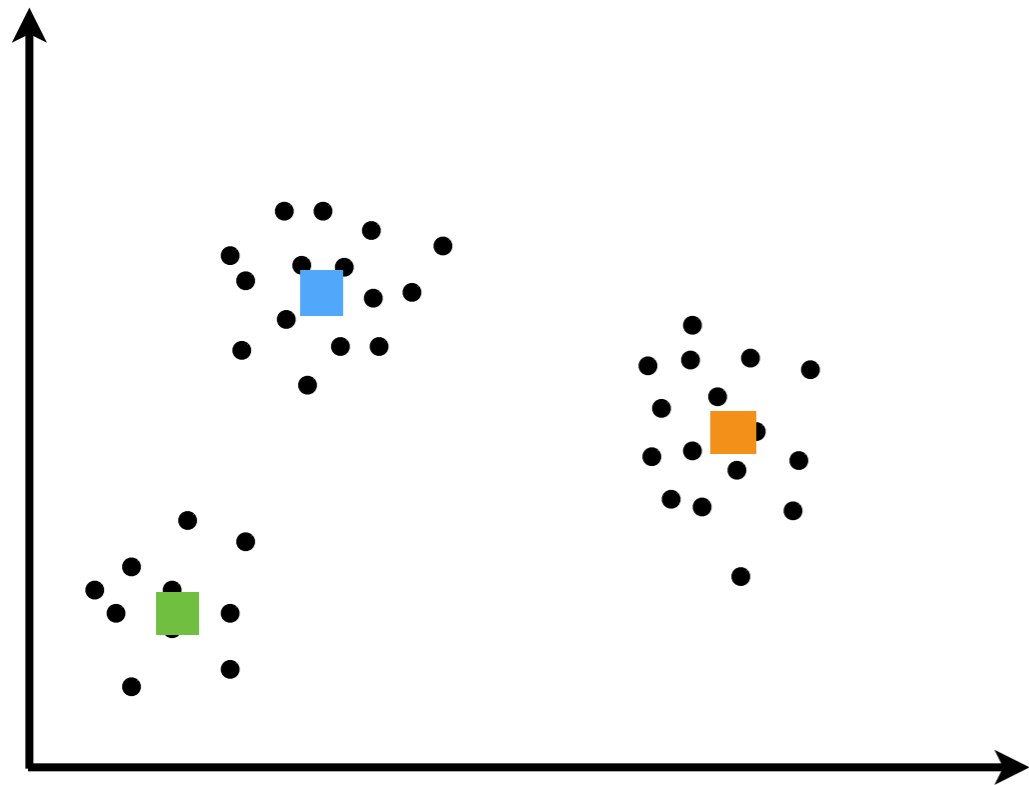
ρ_1

ρ_2

ρ_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)

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

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

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

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



ρ_1

ρ_2

ρ_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

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

Dirichlet distribution review

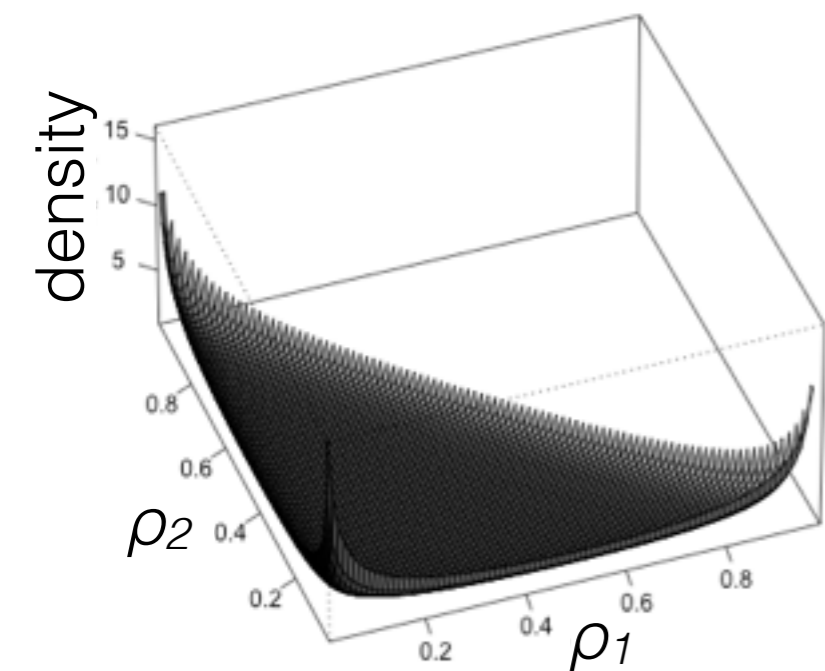
$$\text{Dirichlet}(\rho_{1:K} | a_{1:K}) = \frac{\Gamma(\sum_{k=1}^K a_k)}{\prod_{k=1}^K \Gamma(a_k)} \prod_{k=1}^K \rho_k^{a_k - 1} \quad a_k > 0$$

- What happens?

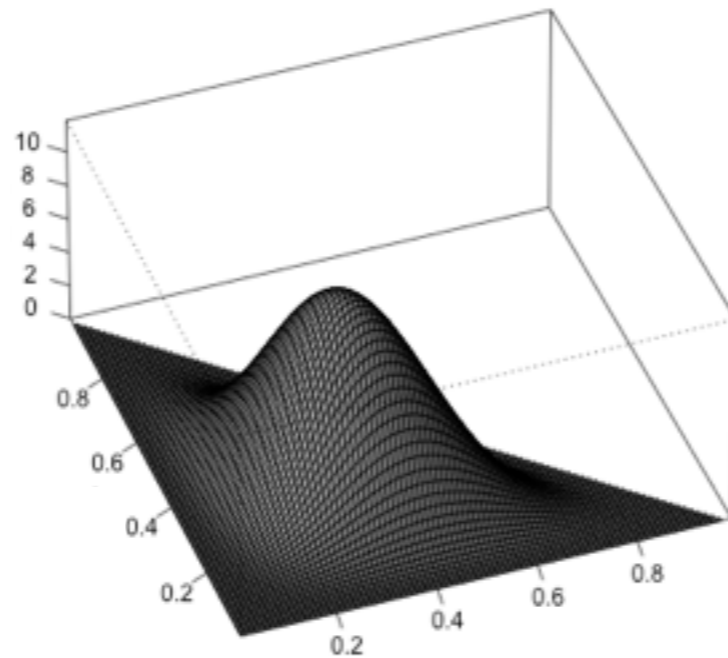
Dirichlet distribution review

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

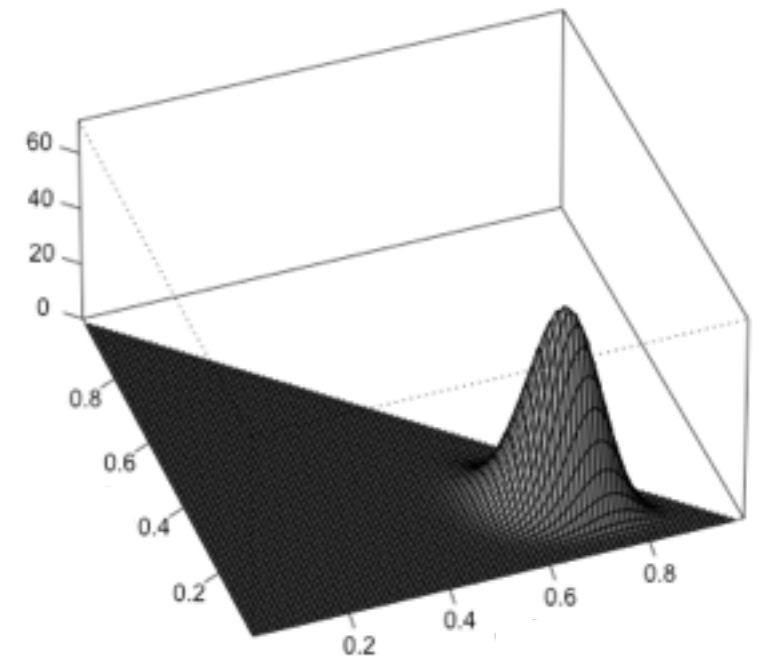
$a = (0.5, 0.5, 0.5)$



$a = (5, 5, 5)$



$a = (40, 10, 10)$

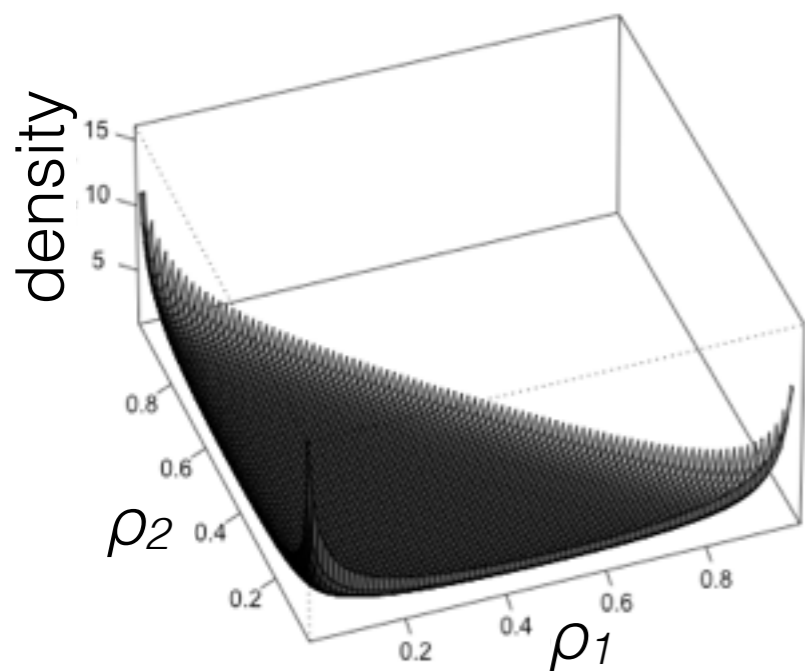


- What happens?

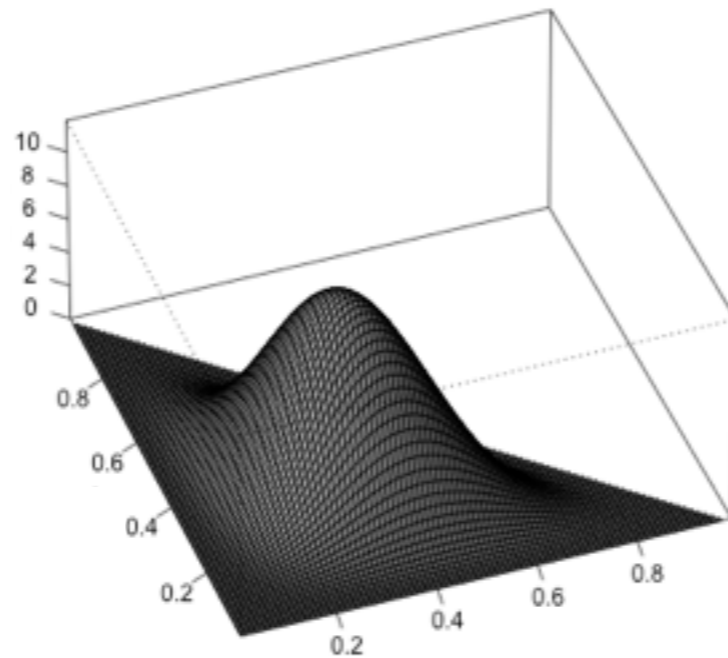
Dirichlet distribution review

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

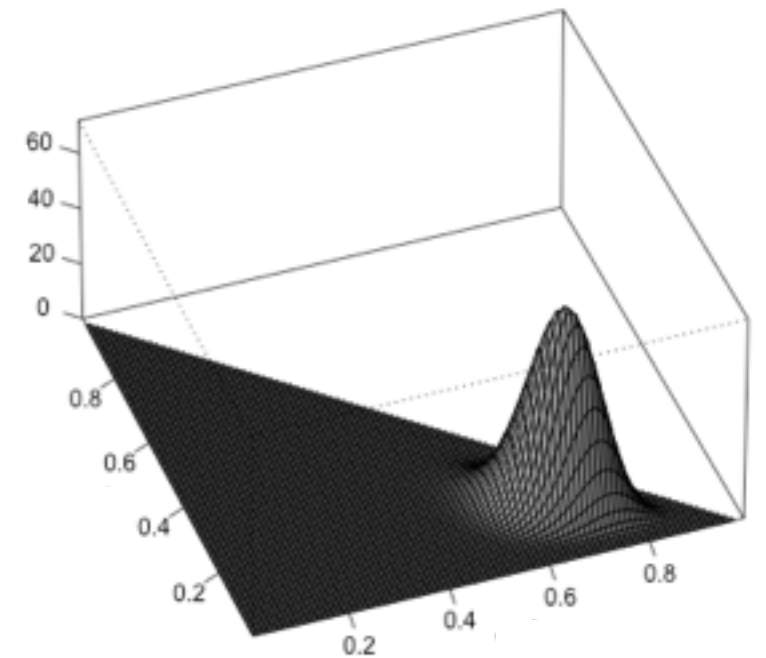
$a = (0.5, 0.5, 0.5)$



$a = (5, 5, 5)$



$a = (40, 10, 10)$

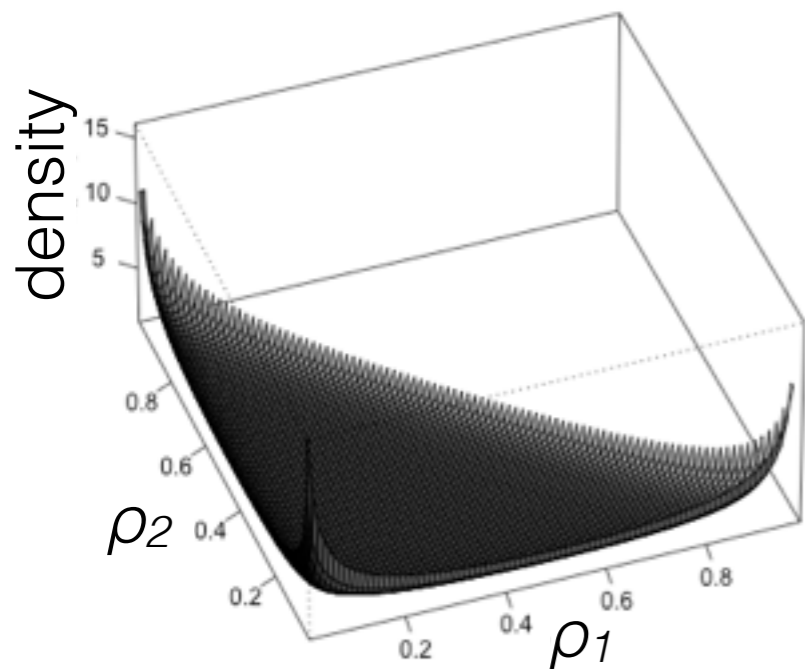


- What happens? $a = a_k = 1$

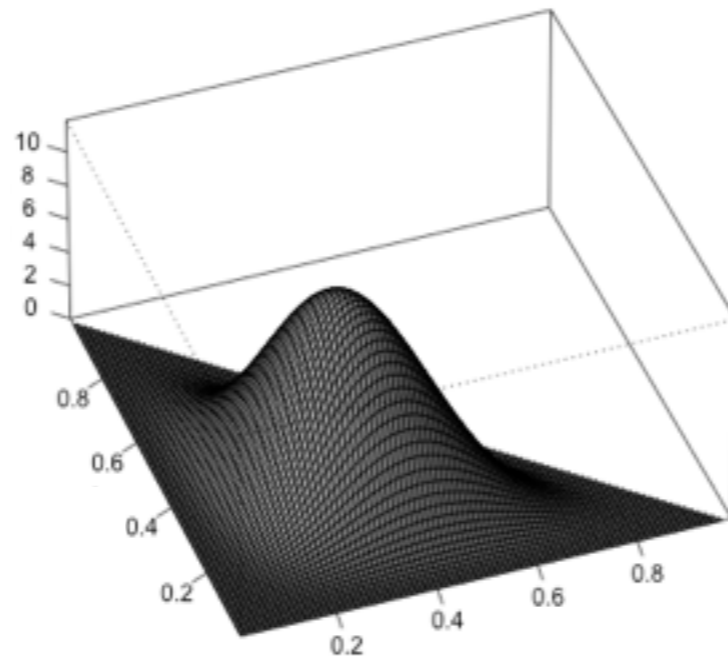
Dirichlet distribution review

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

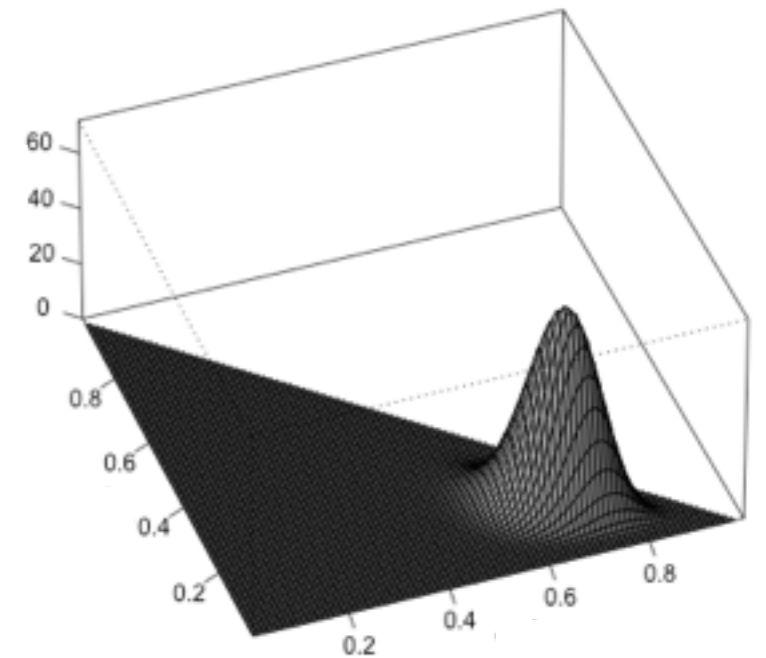
$a = (0.5, 0.5, 0.5)$



$a = (5, 5, 5)$



$a = (40, 10, 10)$

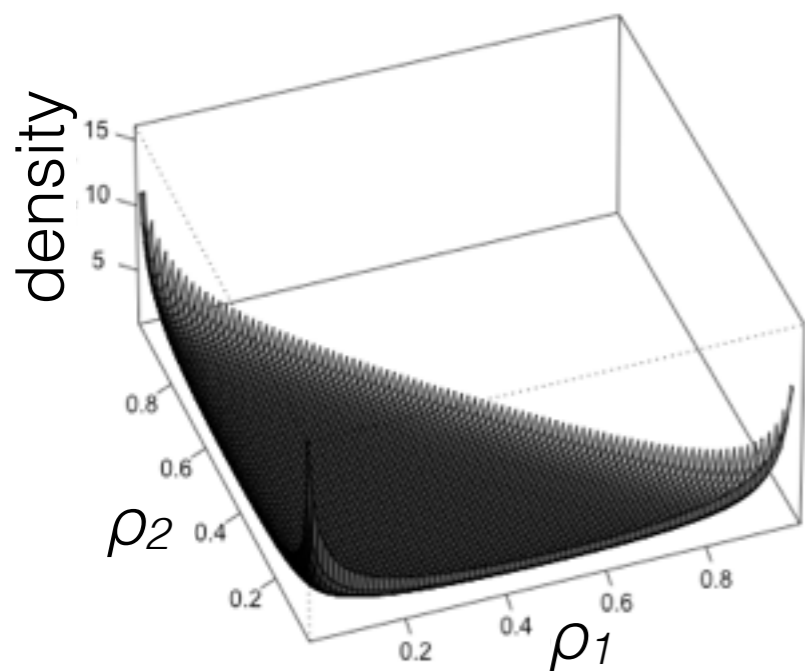


- What happens? $a = a_k = 1$ $a = a_k \rightarrow 0$

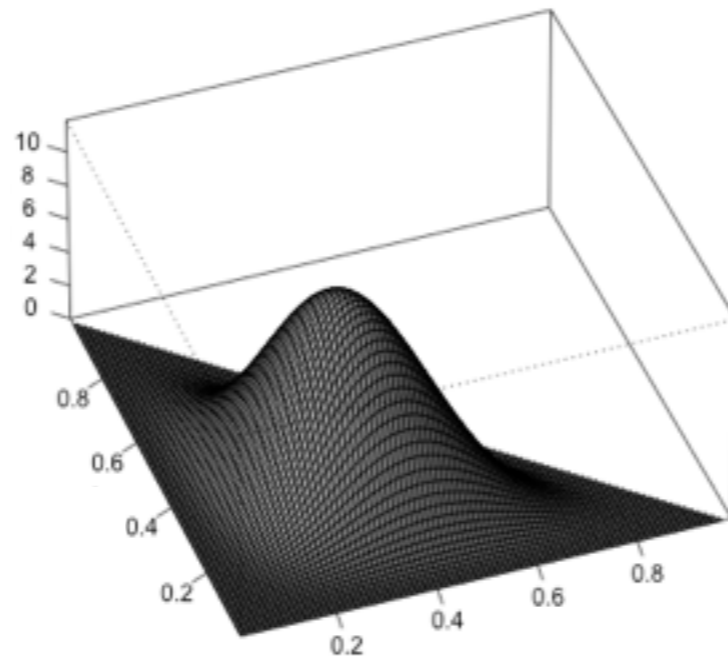
Dirichlet distribution review

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

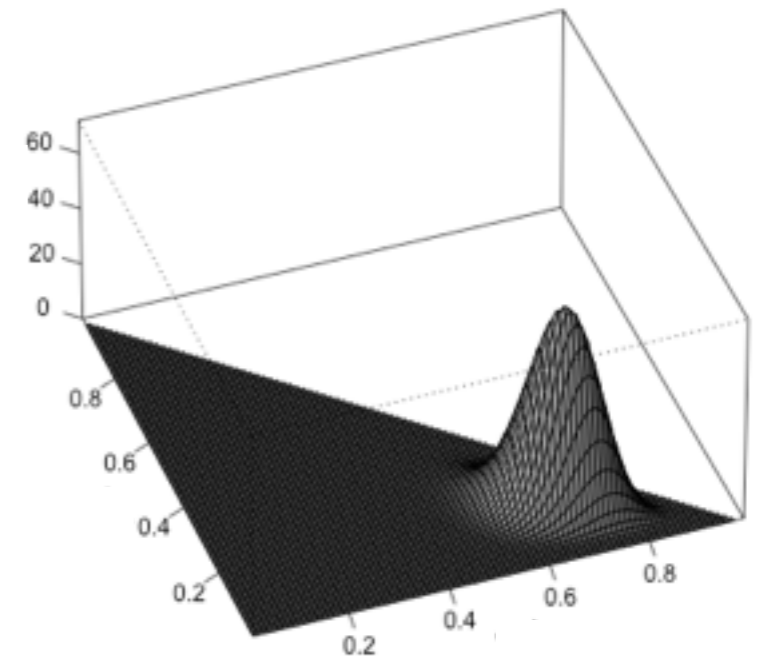
$a = (0.5, 0.5, 0.5)$



$a = (5, 5, 5)$



$a = (40, 10, 10)$

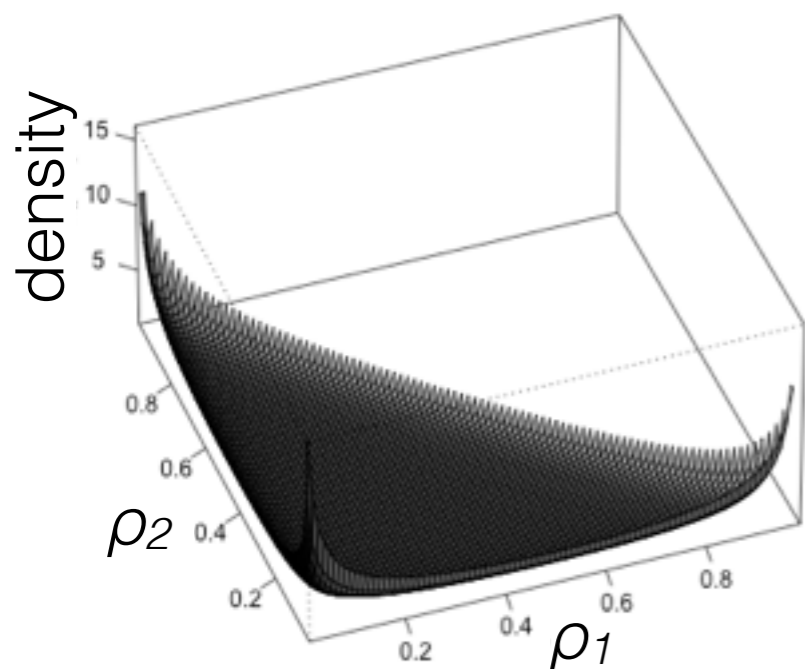


- What happens? $a = a_k = 1$ $a = a_k \rightarrow 0$ $a = a_k \rightarrow \infty$

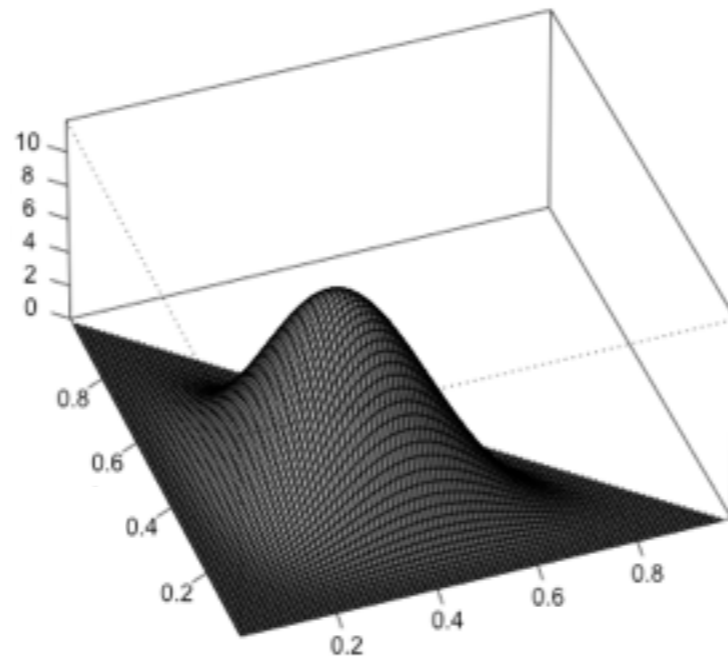
Dirichlet distribution review

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

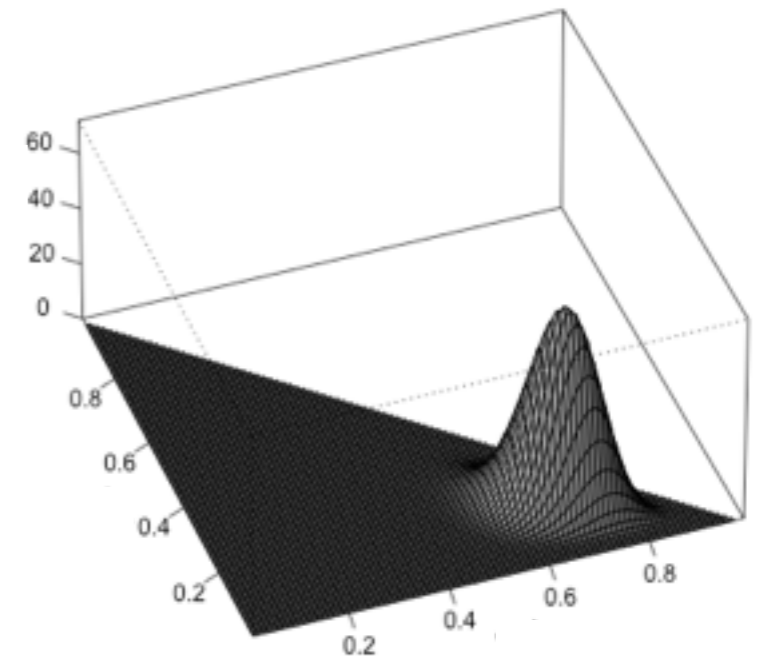
$a = (0.5, 0.5, 0.5)$



$a = (5, 5, 5)$



$a = (40, 10, 10)$



- What happens?

$$a = a_k = 1$$

$$a = a_k \rightarrow 0$$

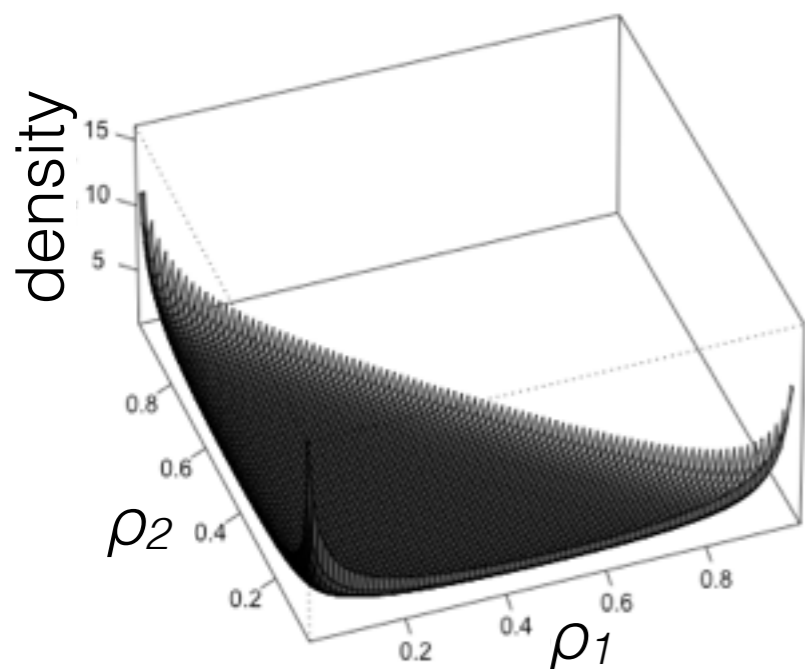
$$a = a_k \rightarrow \infty$$

[demo]

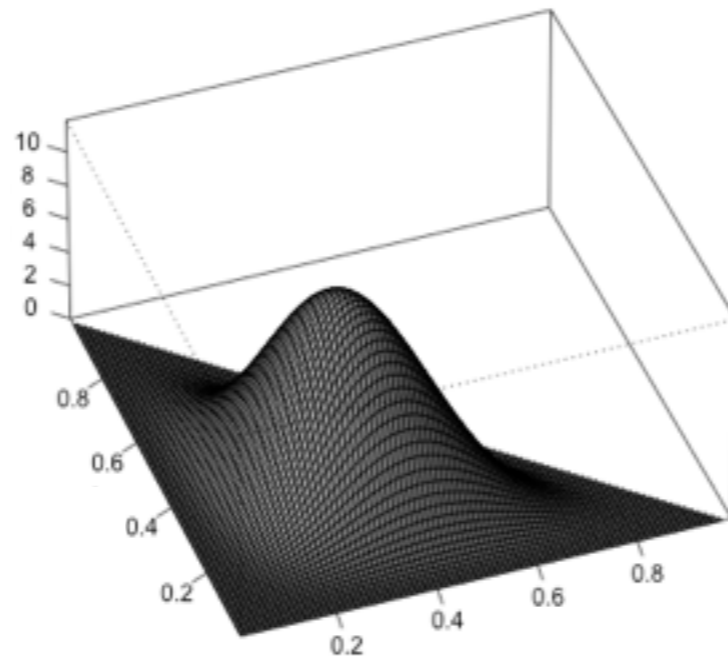
Dirichlet distribution review

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

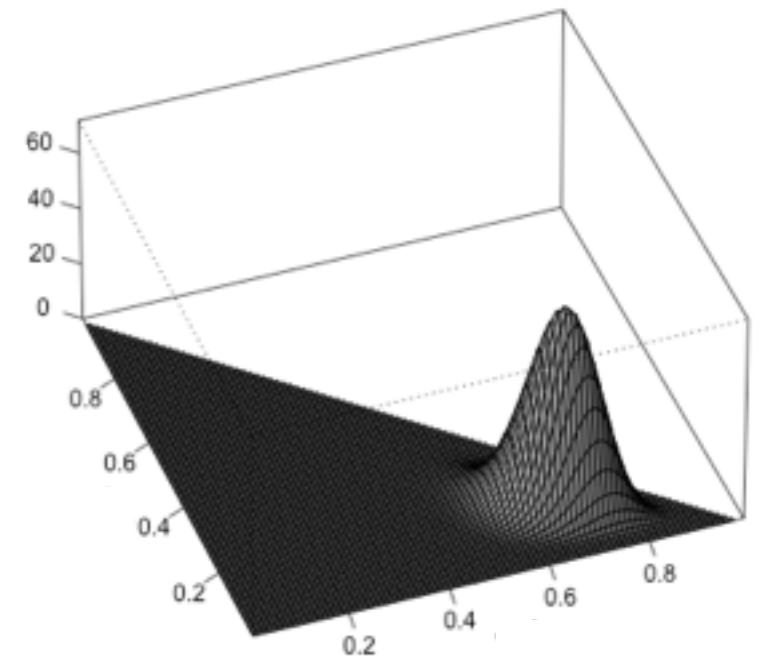
$a = (0.5, 0.5, 0.5)$



$a = (5, 5, 5)$



$a = (40, 10, 10)$

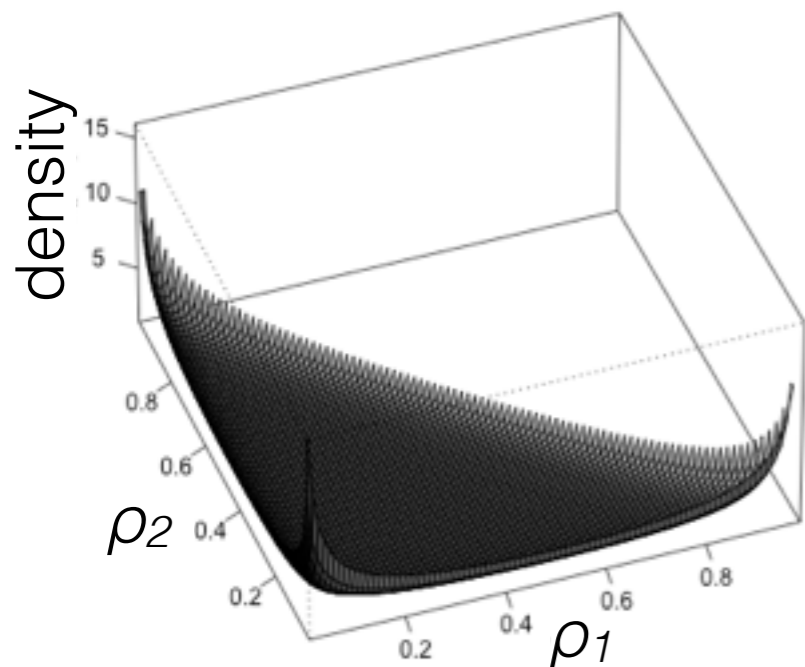


- What happens? $a = a_k = 1$ $a = a_k \rightarrow 0$ $a = a_k \rightarrow \infty$
- Dirichlet is conjugate to Categorical [demo]

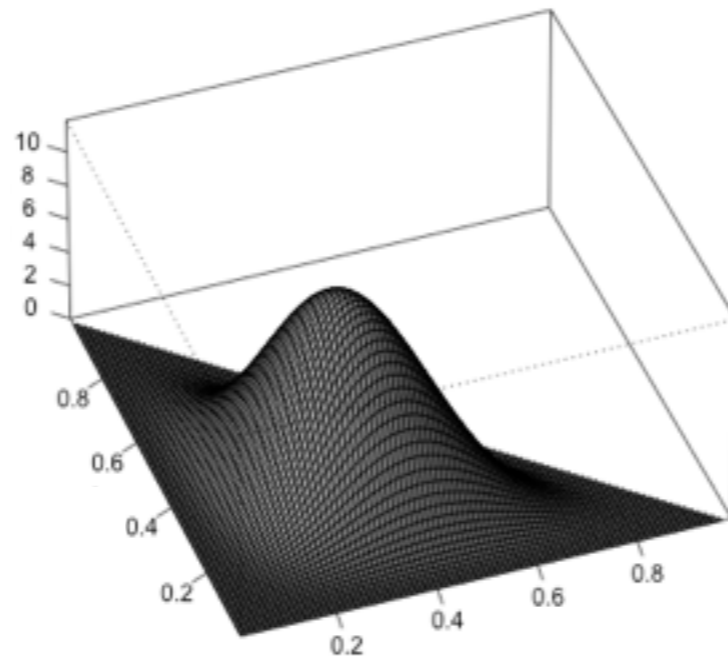
Dirichlet distribution review

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

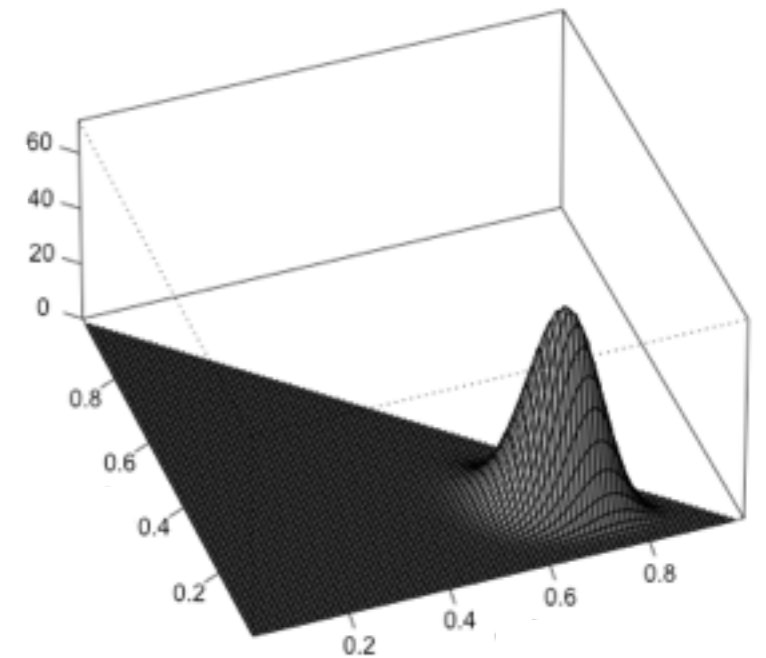
$a = (0.5, 0.5, 0.5)$



$a = (5, 5, 5)$



$a = (40, 10, 10)$

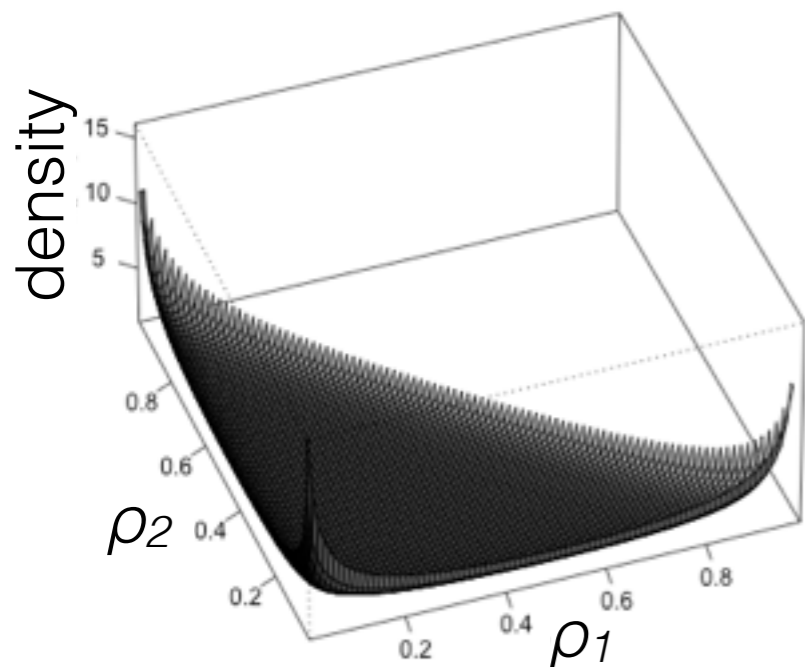


- What happens? $a = a_k = 1$ $a = a_k \rightarrow 0$ $a = a_k \rightarrow \infty$
- Dirichlet is conjugate to Categorical
 $\rho_{1:K} \sim \text{Dirichlet}(a_{1:K}), z \sim \text{Cat}(\rho_{1:K})$ [demo]

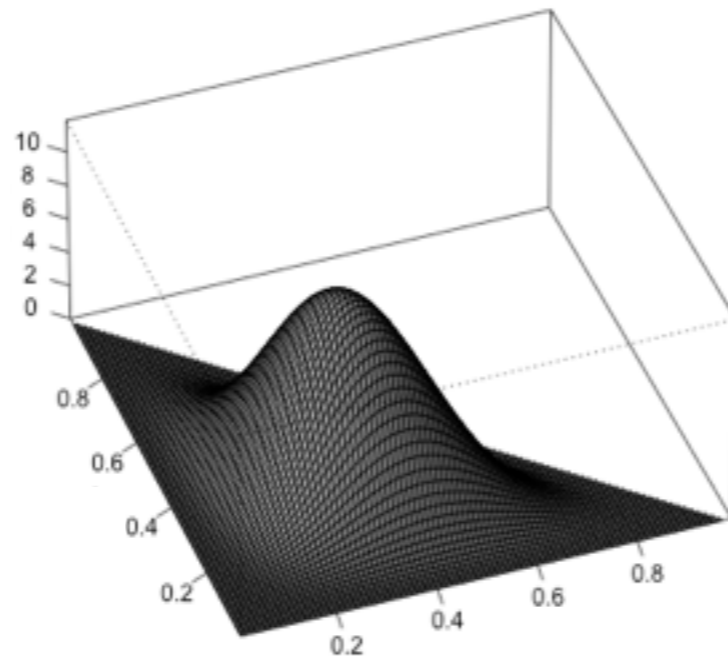
Dirichlet distribution review

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

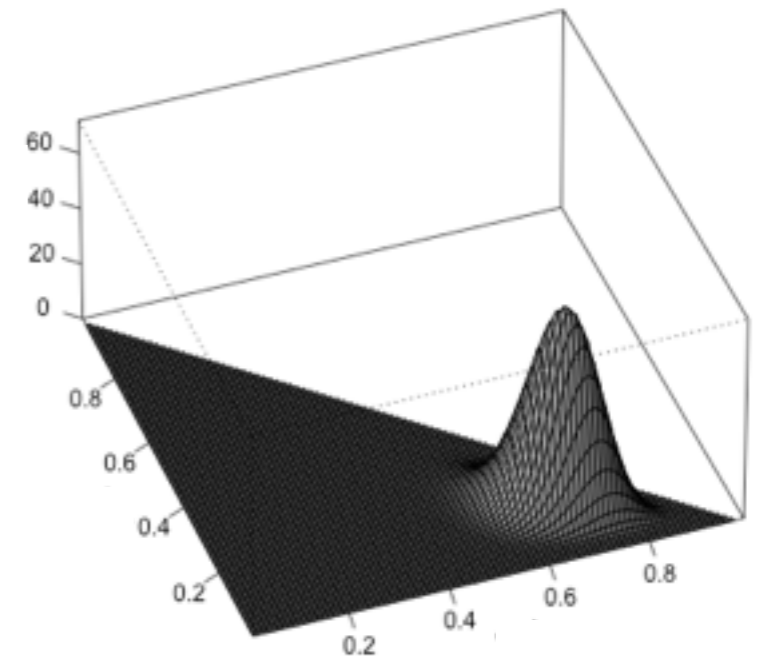
$$a = (0.5, 0.5, 0.5)$$



$$a = (5, 5, 5)$$



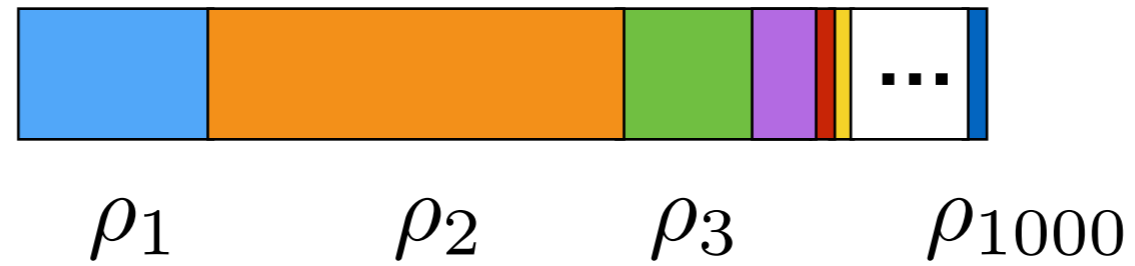
$$a = (40, 10, 10)$$



- What happens? $a = a_k = 1$ $a = a_k \rightarrow 0$ $a = a_k \rightarrow \infty$
 - Dirichlet is conjugate to Categorical [demo]
- $$\rho_{1:K} \sim \text{Dirichlet}(a_{1:K}), z \sim \text{Cat}(\rho_{1:K})$$
- $$\rho_{1:K} | z \stackrel{d}{=} \text{Dirichlet}(a'_{1:K}), a'_k = a_k + \mathbf{1}\{z = k\}$$

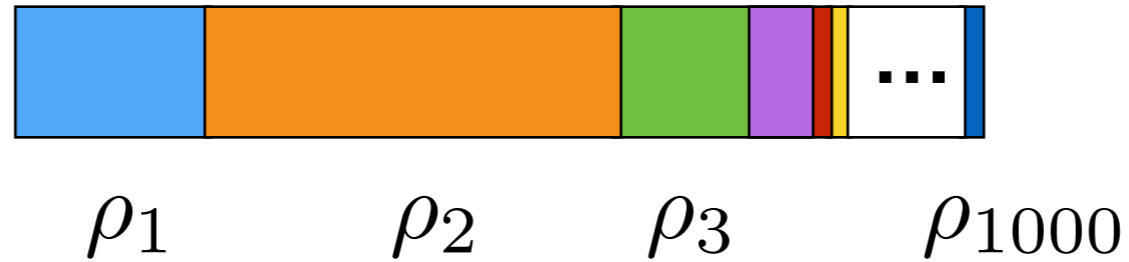
What if $K > N$?

What if $K > N$?



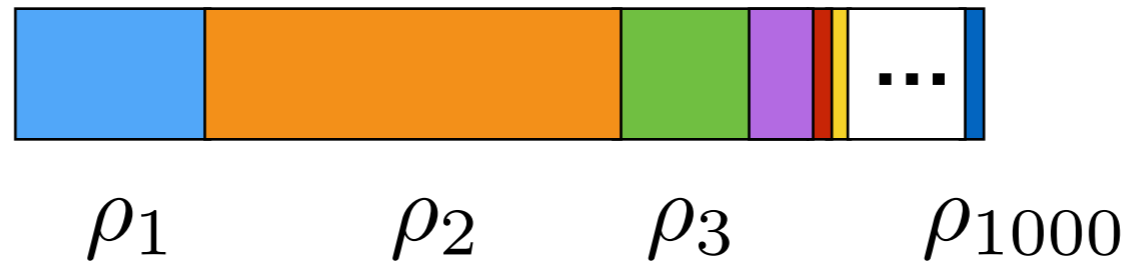
What if $K > N$?

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



What if $K > N$?

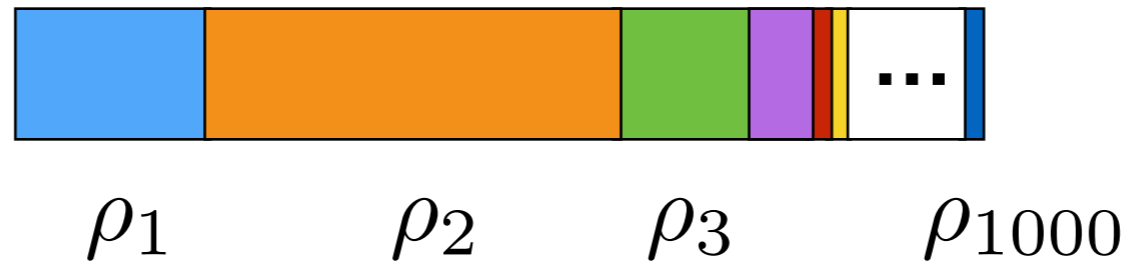
- e.g. species sampling, topic modeling, groups on a social network, etc.



- Components: number of latent groups

What if $K > N$?

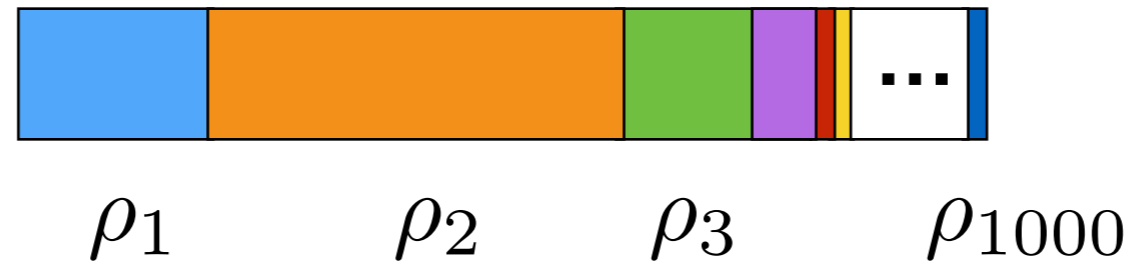
- e.g. species sampling, topic modeling, groups on a social network, etc.



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

What if $K > N$?

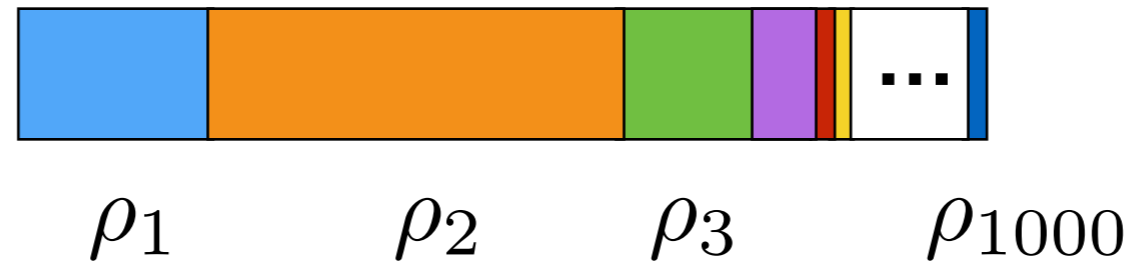
- e.g. species sampling, topic modeling, groups on a social network, etc.



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

What if $K > N$?

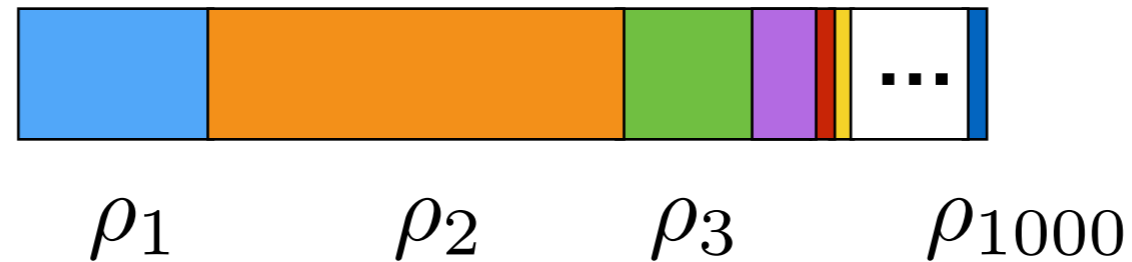
- e.g. species sampling, topic modeling, groups on a social network, etc.



- Components: number of latent groups
- Clusters: number of components represented in the data
- [demo 1, demo 2]
- Number of clusters for N data points is $< K$ and random

What if $K > N$?

- 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 for N data points is $< K$ and 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: $\rho_{1:K} \sim \text{Dirichlet}(a_{1:K})$

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



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

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



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

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



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

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



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

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



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

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

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

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

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

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



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

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



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

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



- “Stick breaking”

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

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

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

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



- “Stick breaking”

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

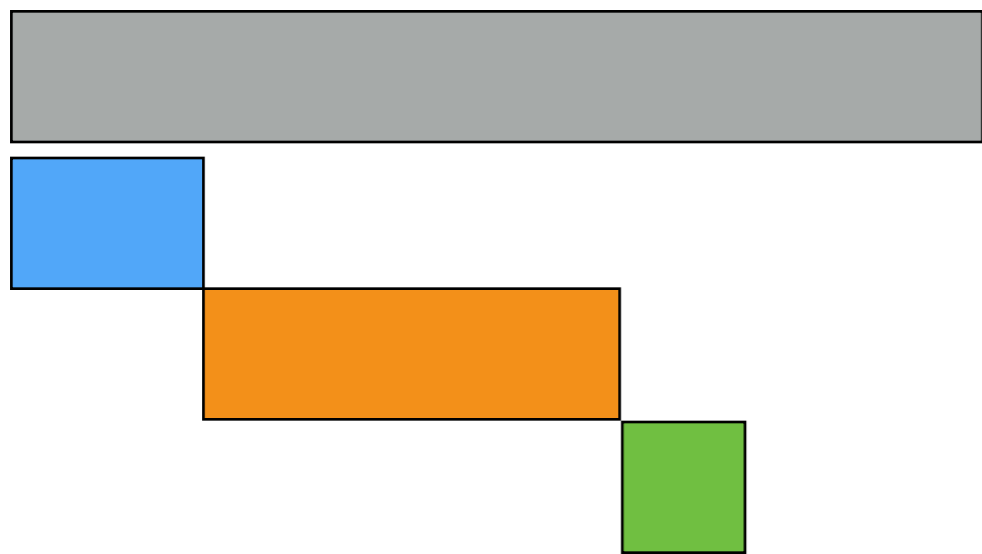
$$V_2 \sim \text{Beta}(a_2, a_3 + a_4) \quad \rho_2 = (1 - V_1)V_2$$

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

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

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



- “Stick breaking”

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

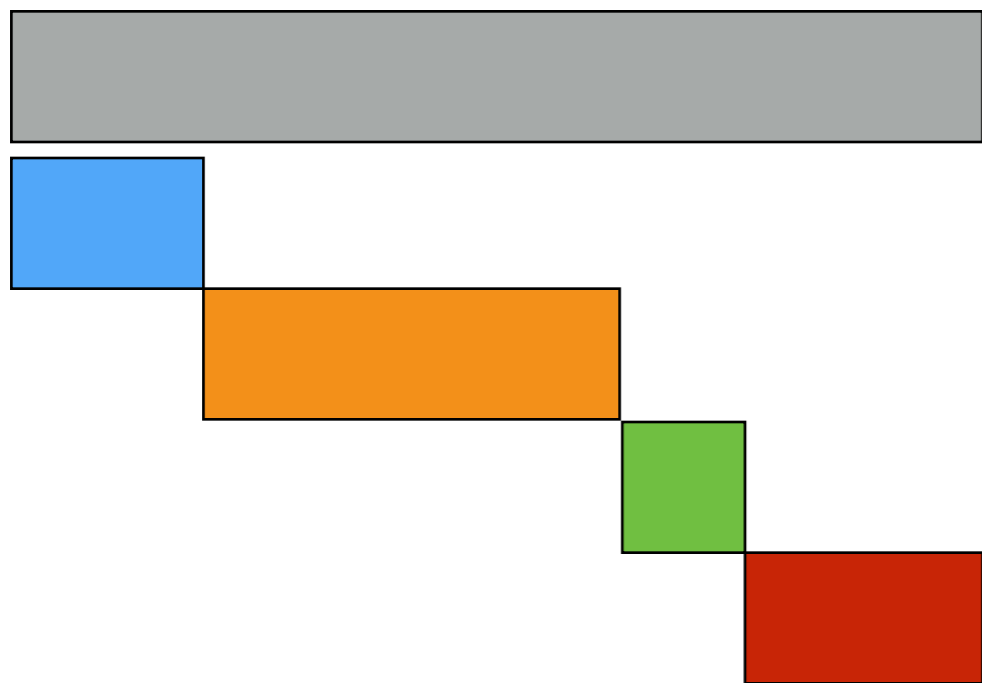
$$V_2 \sim \text{Beta}(a_2, a_3 + a_4) \quad \rho_2 = (1 - V_1)V_2$$

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

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

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



- “Stick breaking”

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

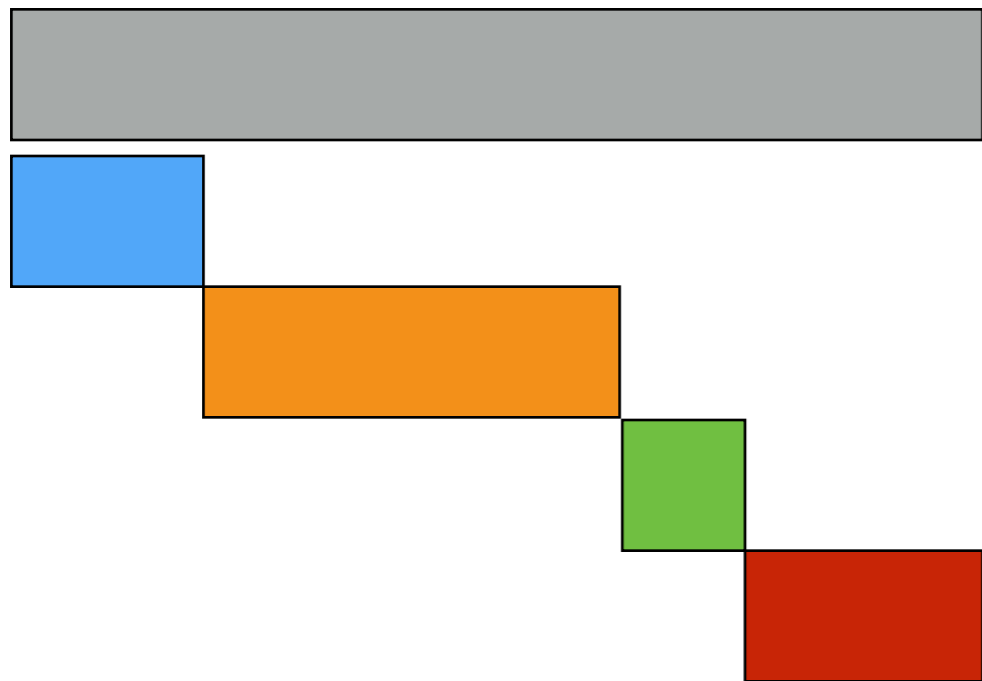
$$V_2 \sim \text{Beta}(a_2, a_3 + a_4) \quad \rho_2 = (1 - V_1)V_2$$

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

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

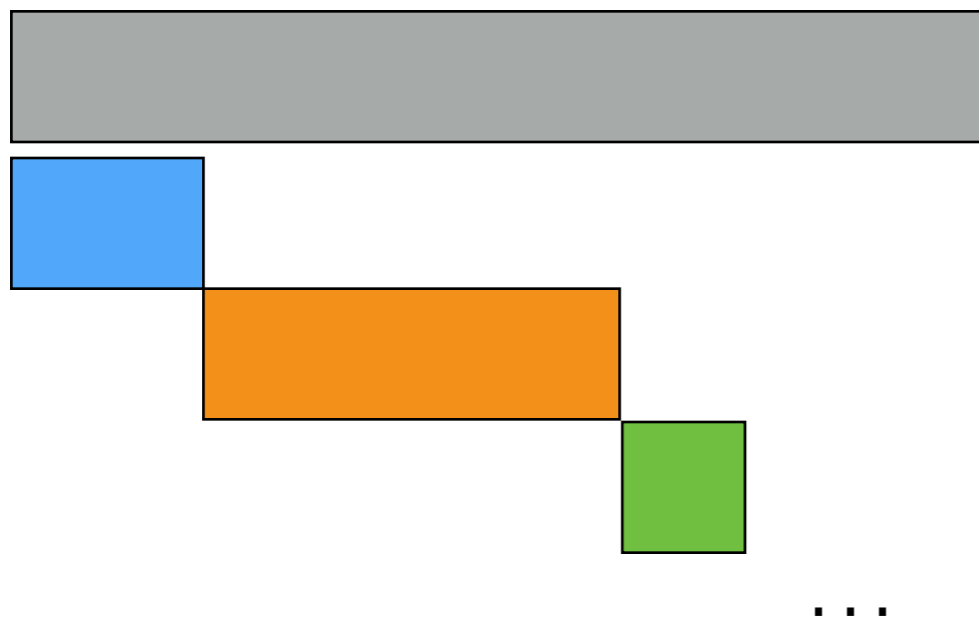
Choosing $K = \infty$

- 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?



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

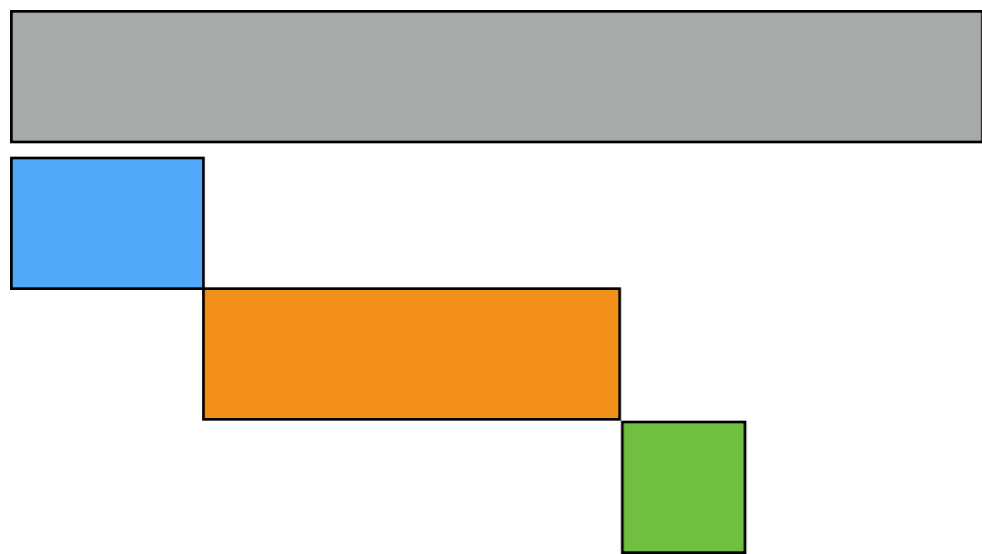
$$\rho_1 = V_1$$

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

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

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?



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

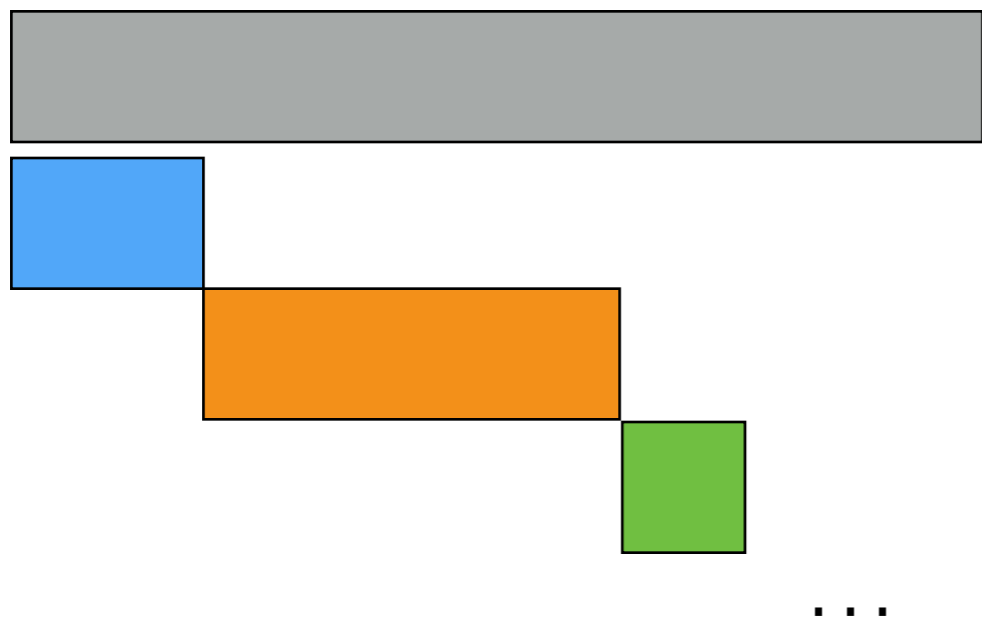
$$\rho_1 = V_1$$

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

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

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?



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

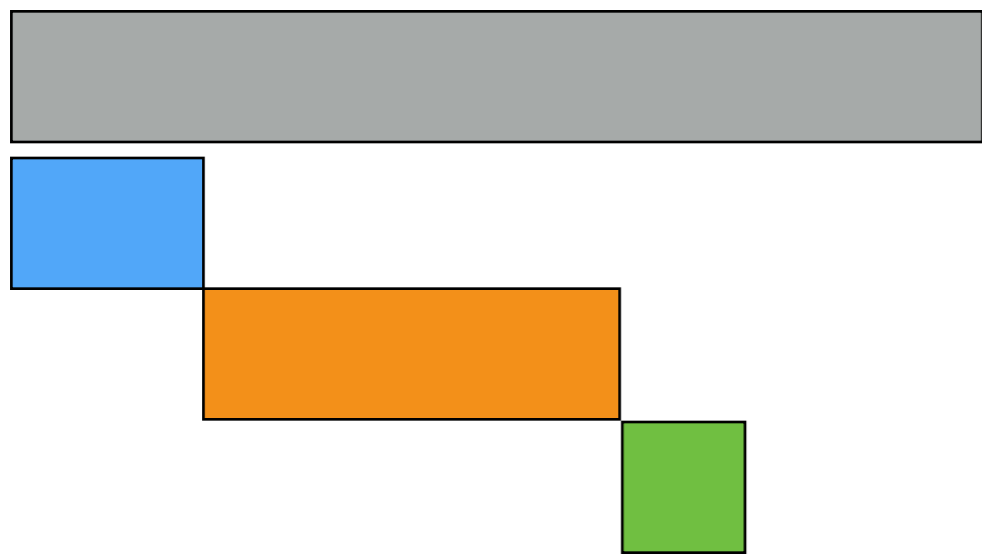
$$\rho_1 = V_1$$

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

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

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?



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

$$\rho_1 = V_1$$

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

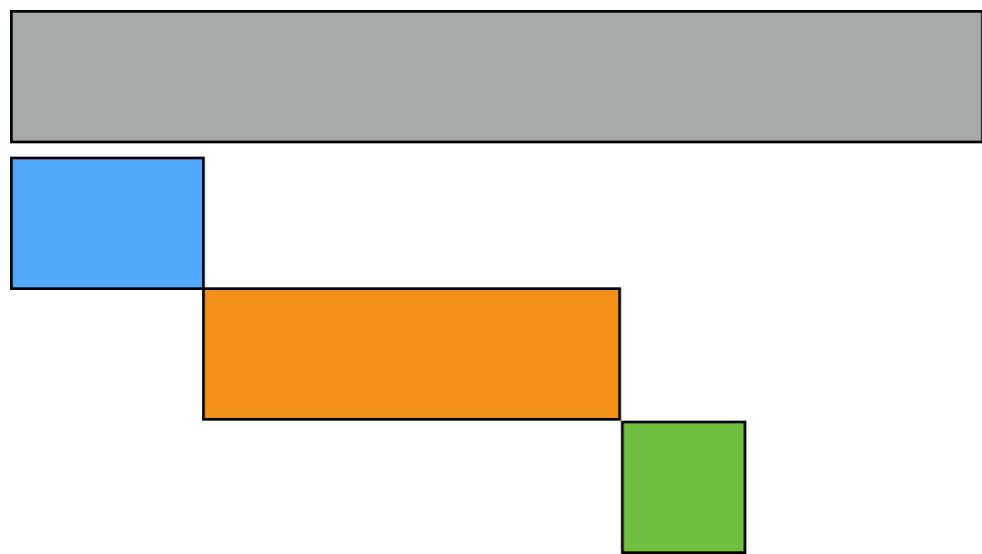
$$\rho_2 = (1 - V_1)V_2$$

...

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

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?



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

$$\rho_1 = V_1$$

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

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

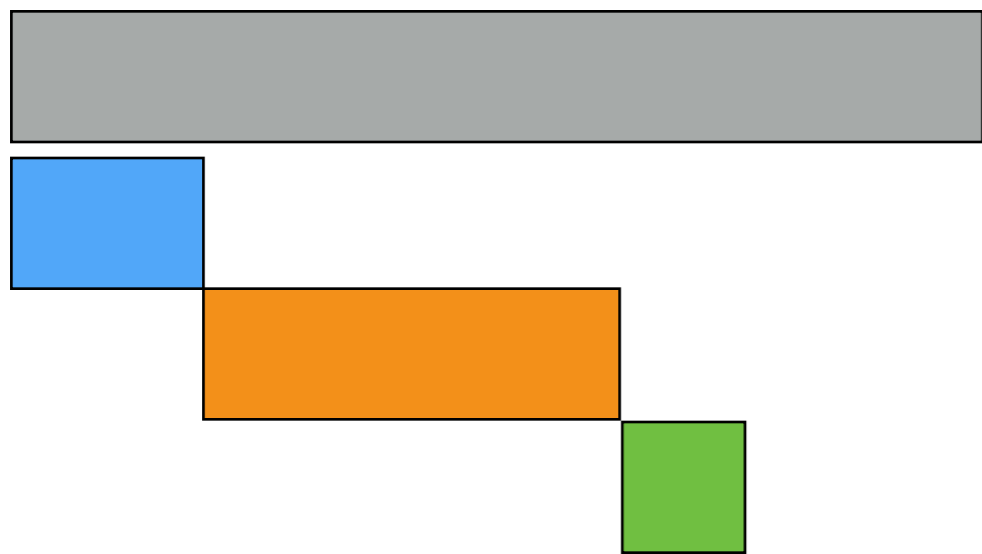
...

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

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

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?



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

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

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

$$\rho_1 = V_1$$

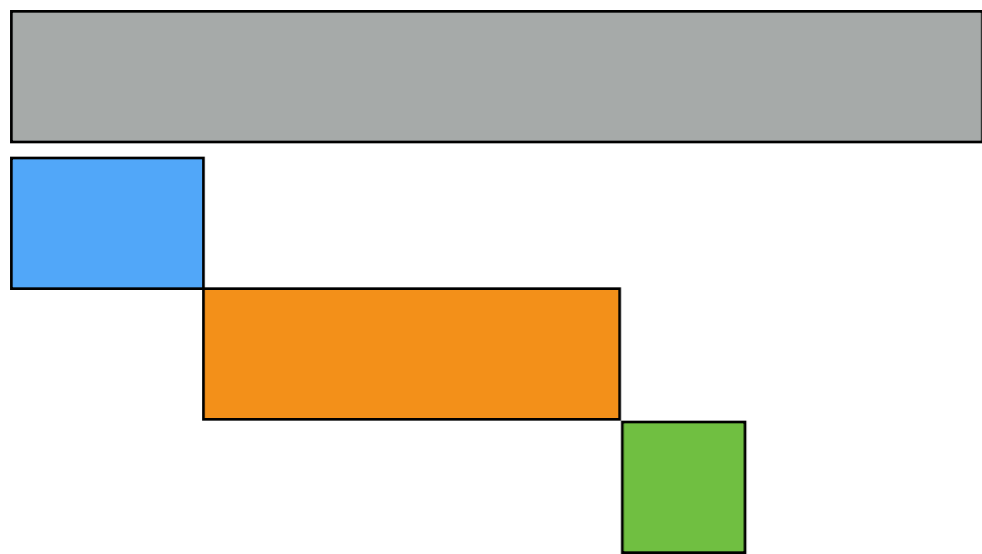
$$\rho_2 = (1 - V_1)V_2$$

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

[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]

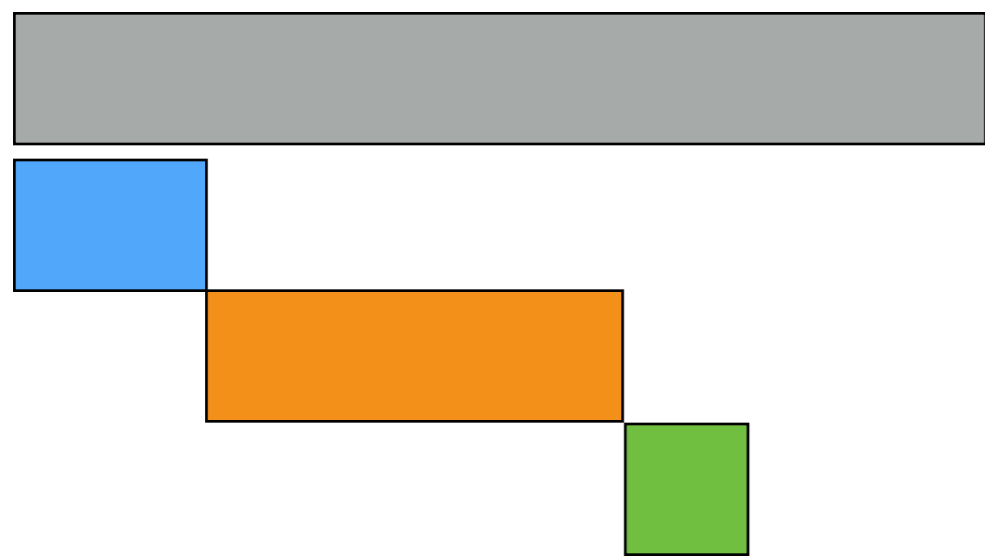
Choosing $K = \infty$

- 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:

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



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

$$\rho_1 = V_1$$

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

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

...

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

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

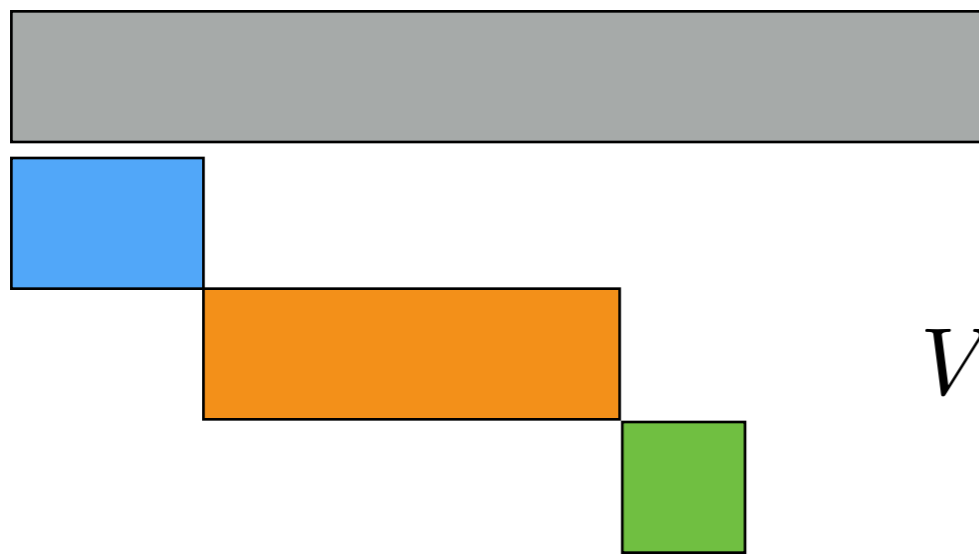
Choosing $K = \infty$

- 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:

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



$$V_k \stackrel{iid}{\sim} \text{Beta}(1, \alpha)$$

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

...

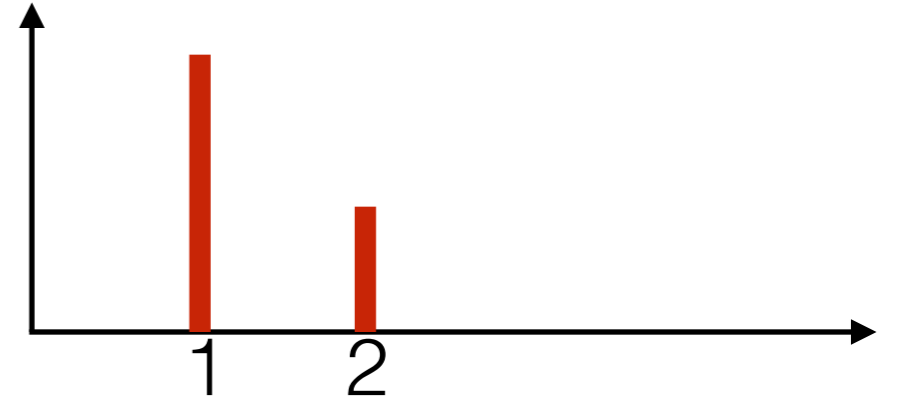
Distributions

Distributions

- Beta \rightarrow random distribution over 1, 2

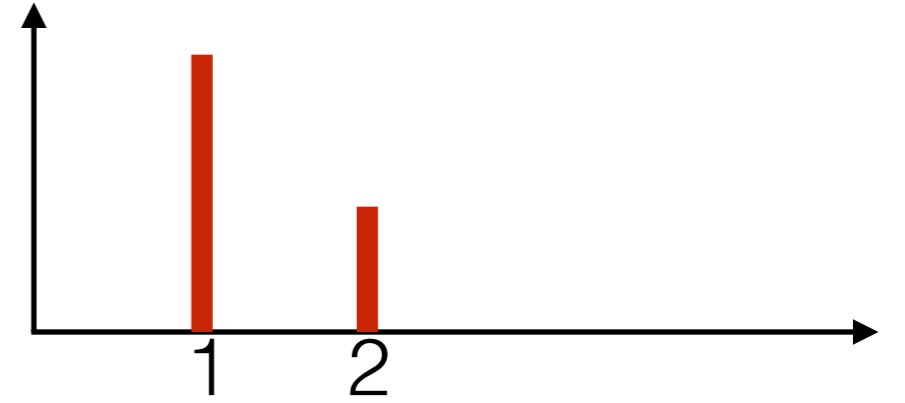
Distributions

- Beta \rightarrow random distribution over 1, 2



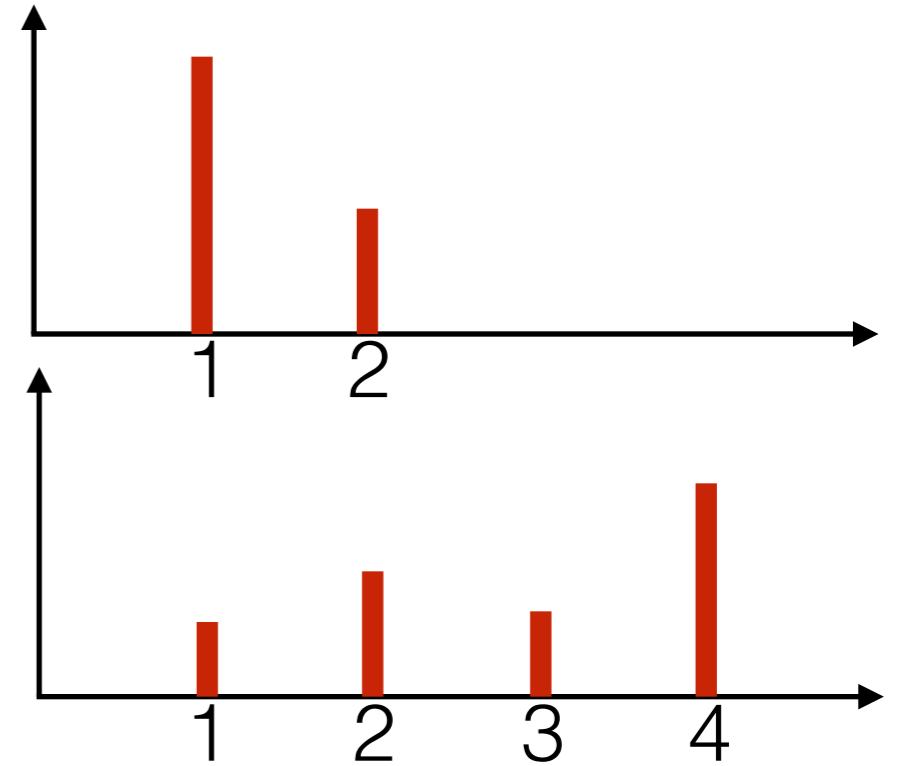
Distributions

- Beta \rightarrow random distribution over 1, 2
- Dirichlet \rightarrow random distribution over 1, 2, \dots , K



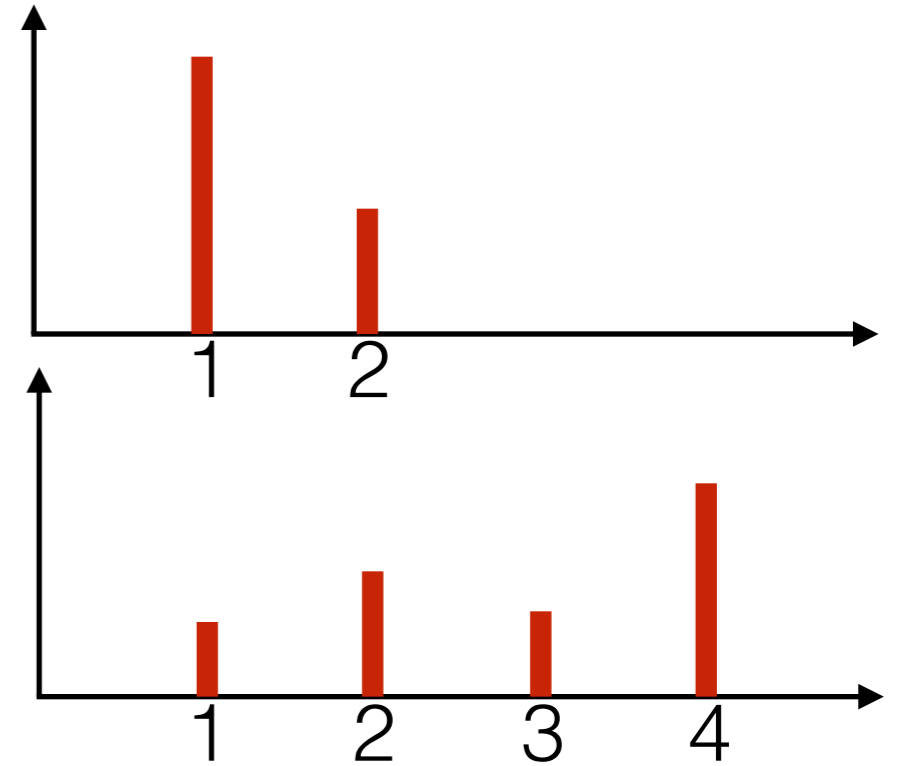
Distributions

- Beta \rightarrow random distribution over 1, 2
- Dirichlet \rightarrow random distribution over 1, 2, ..., K



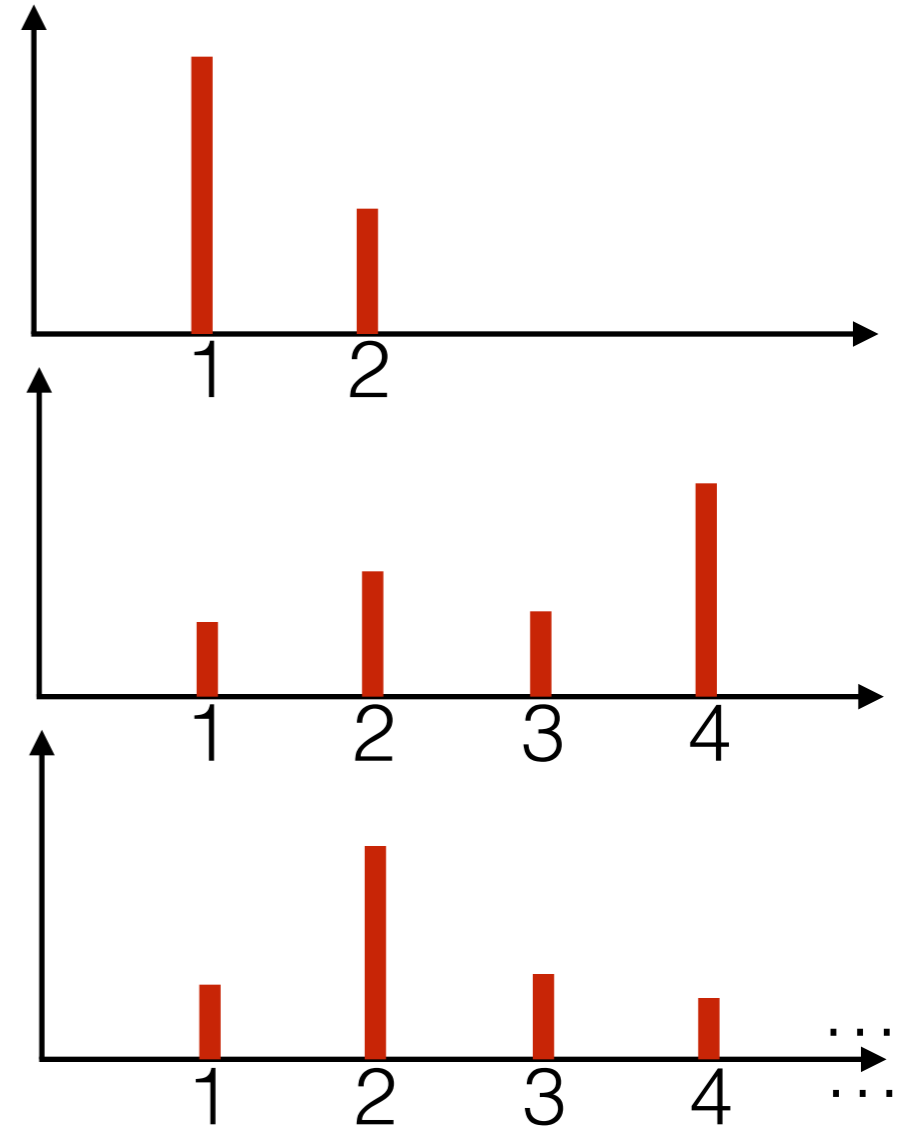
Distributions

- Beta \rightarrow random distribution over 1, 2
- Dirichlet \rightarrow random distribution over 1, 2, ..., K
- GEM / Dirichlet process stick-breaking \rightarrow random distribution over 1, 2, ...



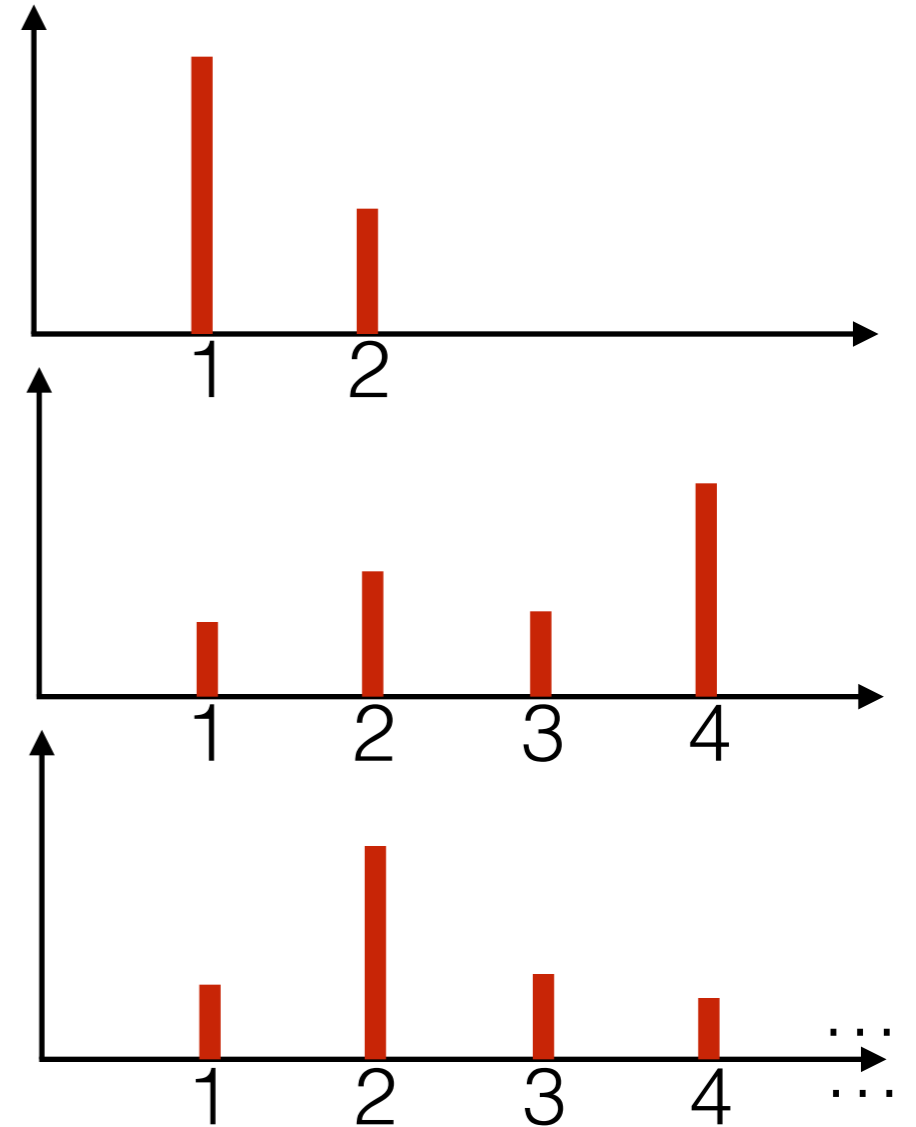
Distributions

- Beta \rightarrow random distribution over 1, 2
- Dirichlet \rightarrow random distribution over 1, 2, ..., K
- GEM / Dirichlet process stick-breaking \rightarrow random distribution over 1, 2, ...



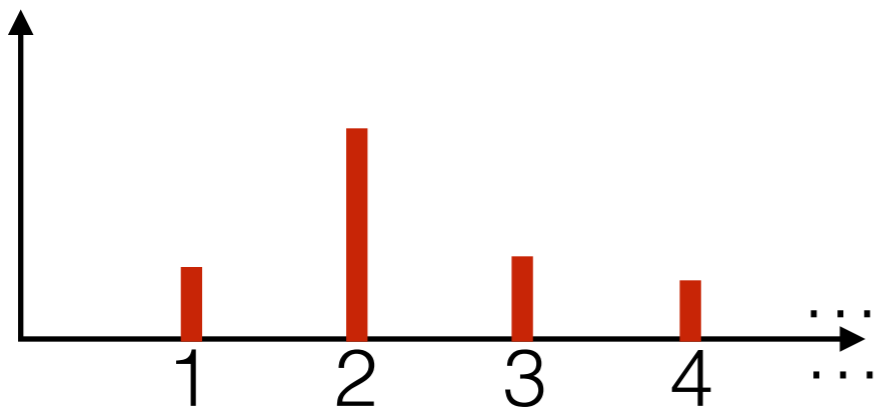
Distributions

- Beta \rightarrow random distribution over 1, 2
- Dirichlet \rightarrow random distribution over 1, 2, ..., K
- GEM / Dirichlet process stick-breaking \rightarrow random distribution over 1, 2, ...



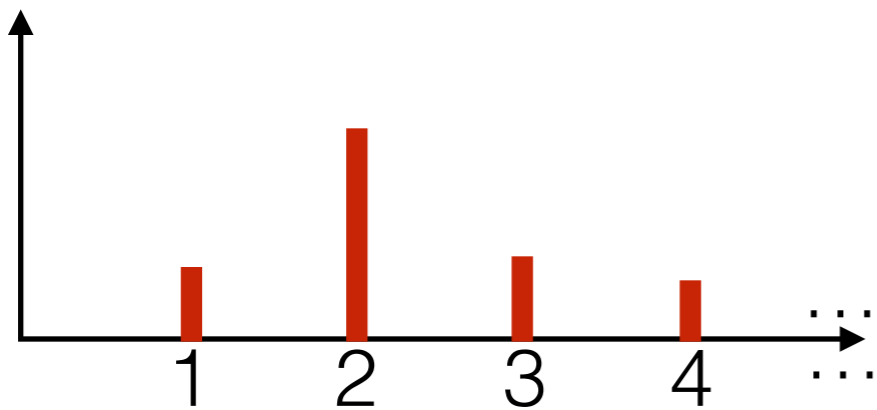
- Infinity of parameters: components
- Growing number of parameters: clusters

Exercises



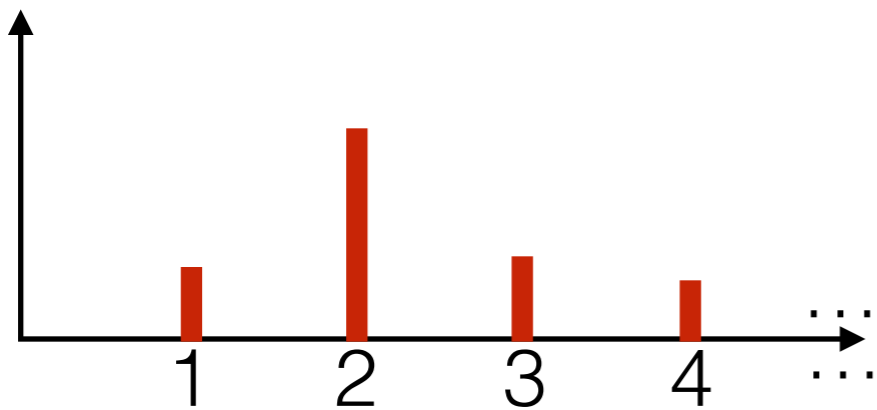
Exercises

- Prove the beta (Dirichlet) is conjugate to the categorical



Exercises

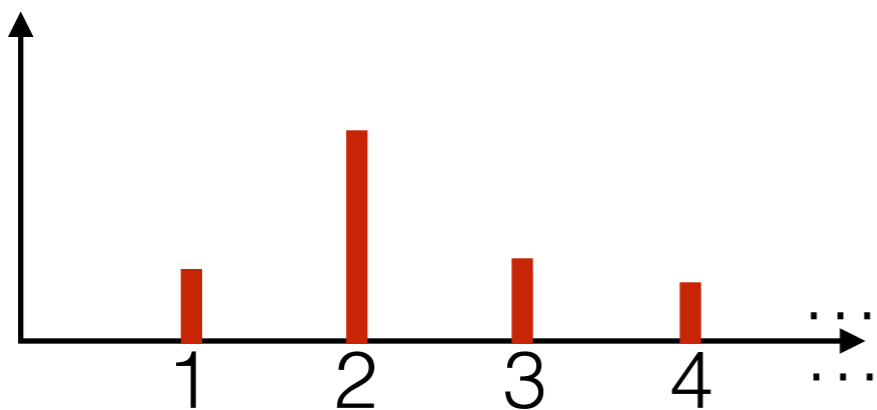
- Prove the beta (Dirichlet) is conjugate to the categorical
 - What is the posterior after N data points?



Exercises

- Prove the beta (Dirichlet) is conjugate to the categorical
 - What is the posterior after N data points?
- Suppose $\rho_{1:K} \sim \text{Dirichlet}(a_{1:K})$; prove that

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

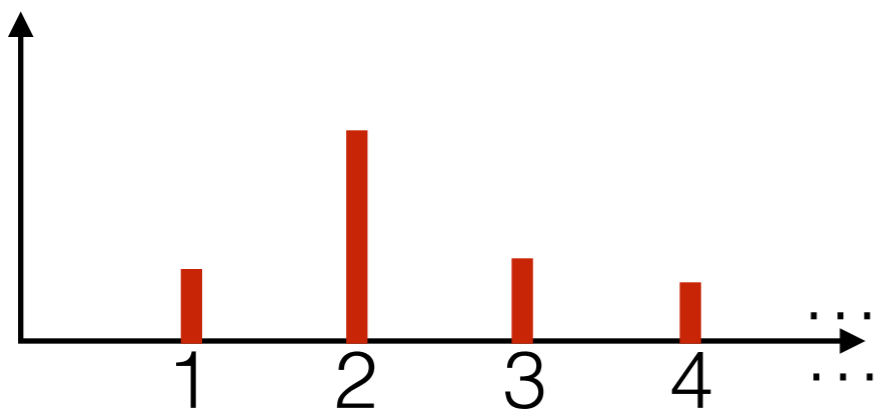


Exercises

- Prove the beta (Dirichlet) is conjugate to the categorical
 - What is the posterior after N data points?
- Suppose $\rho_{1:K} \sim \text{Dirichlet}(a_{1:K})$; prove that

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

- Code your own GEM simulator for ρ ; why is this hard?

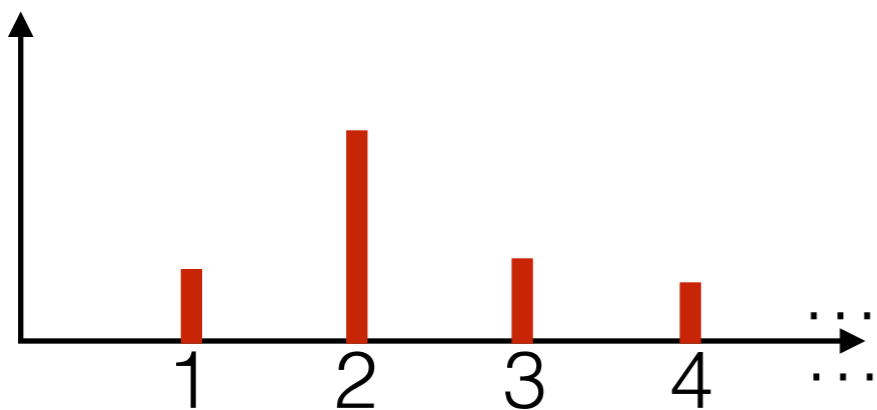


Exercises

- Prove the beta (Dirichlet) is conjugate to the categorical
 - What is the posterior after N data points?
- Suppose $\rho_{1:K} \sim \text{Dirichlet}(a_{1:K})$; prove that

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

- Code your own GEM simulator for ρ ; why is this hard?
- Simulate drawing cluster indicators (z) from your ρ

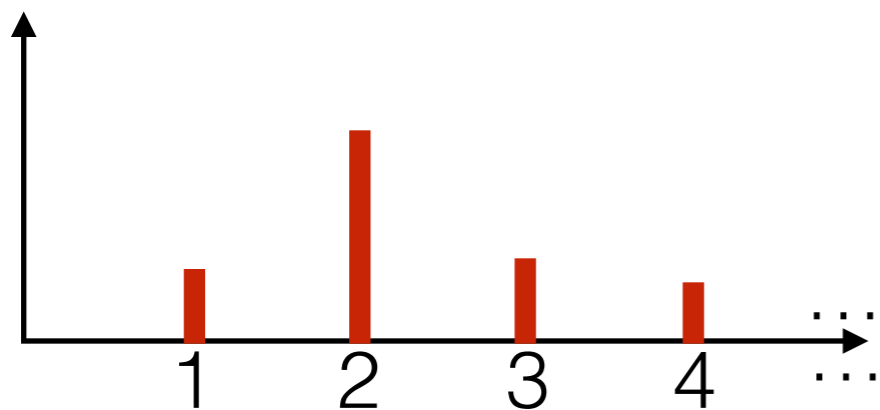


Exercises

- Prove the beta (Dirichlet) is conjugate to the categorical
 - What is the posterior after N data points?
- Suppose $\rho_{1:K} \sim \text{Dirichlet}(a_{1:K})$; prove that

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

- Code your own GEM simulator for ρ ; why is this hard?
- Simulate drawing cluster indicators (z) from your ρ



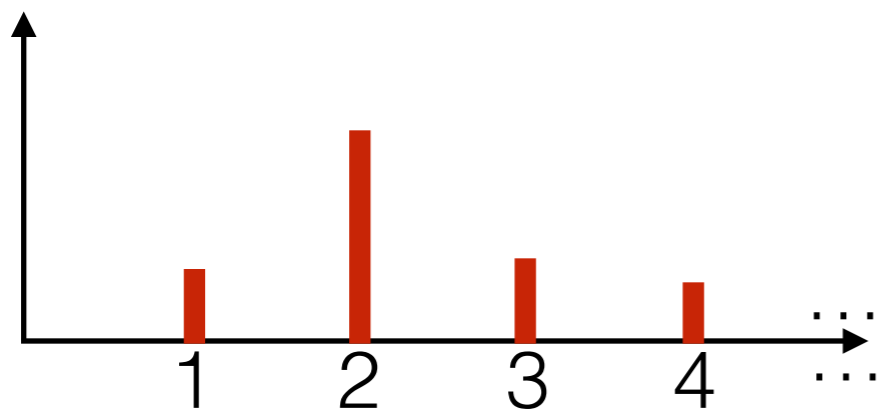
- Compare the number of clusters as N changes in the GEM case with the growth in the $K=1000$ case

Exercises

- Prove the beta (Dirichlet) is conjugate to the categorical
 - What is the posterior after N data points?
- Suppose $\rho_{1:K} \sim \text{Dirichlet}(a_{1:K})$; prove that

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

- Code your own GEM simulator for ρ ; why is this hard?
- Simulate drawing cluster indicators (z) from your ρ



- Compare the number of clusters as N changes in the GEM case with the growth in the $K=1000$ case
- How does the growth in N change when you change α ?

References

A full reference list is provided at the end of the “Part II” slides.