

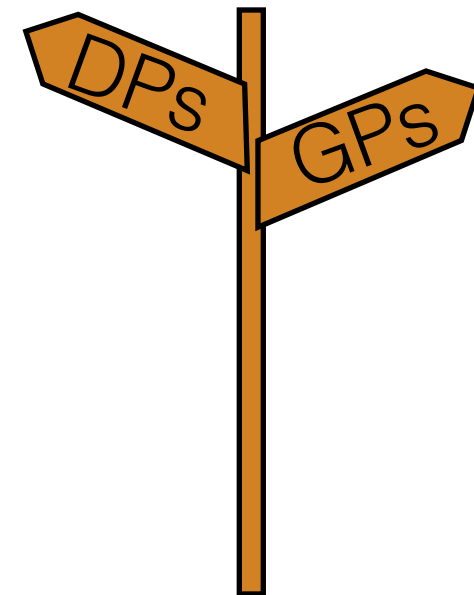
# Nonparametric Bayesian Methods: Models, Algorithms, and Applications (Day 5)

Tamara Broderick

ITT Career Development Assistant Professor  
Electrical Engineering & Computer Science  
MIT

# Roadmap

- Bayes Foundations
- Unsupervised Learning
  - Example problem: clustering
  - Example BNP model: Dirichlet process (DP)
  - Chinese restaurant process
- Supervised Learning
  - Example problem: regression
  - Example BNP model: Gaussian process (GP)
- Venture further into the wild world of Nonparametric Bayes
- Big questions
  - Why BNP?
  - What does an infinite/growing number of parameters really mean (in BNP)?
  - Why is BNP challenging but practical?

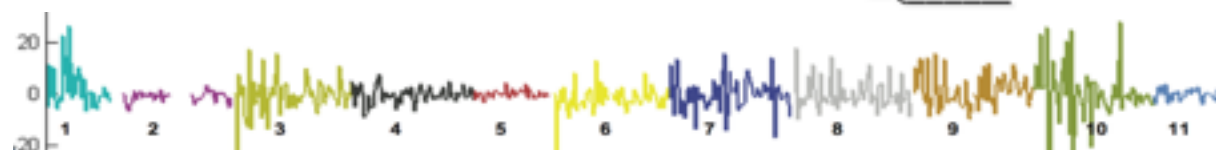


# Applications

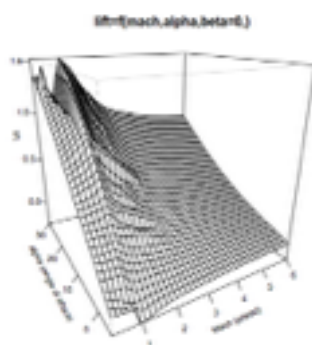


[wikipedia.org]

[Saria  
et al  
2010]



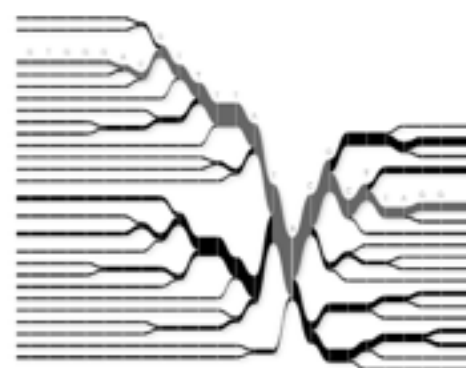
[US CDC PHIL;  
Futoma, Hariharan,  
Heller 2017]



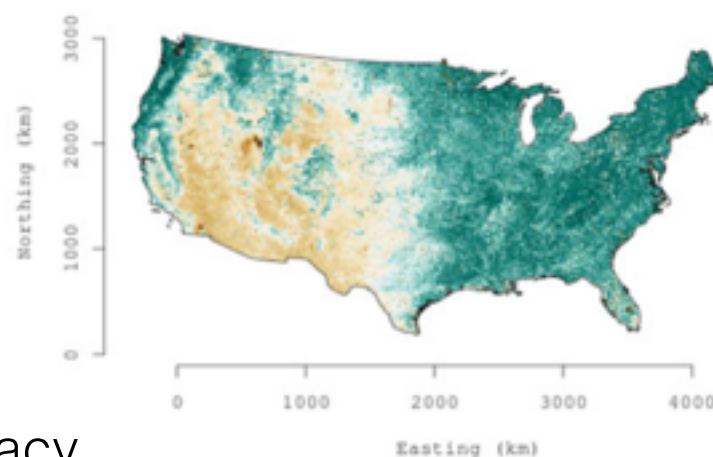
[Gramacy,  
Lee 2009]



[Ed Bowlby, NOAA]



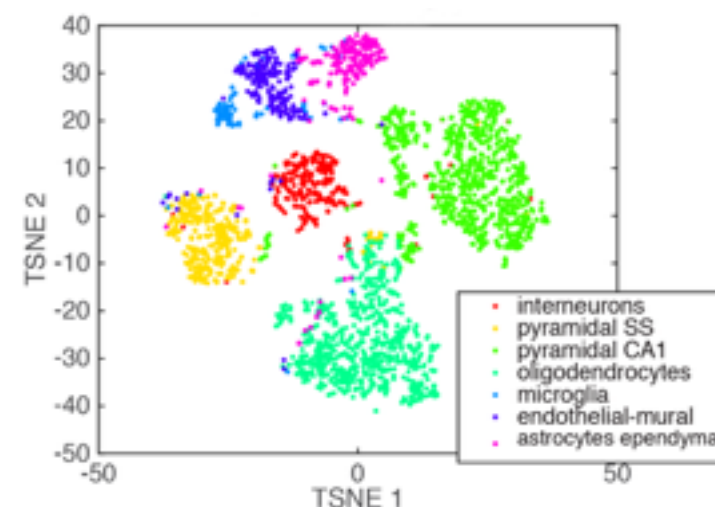
[Ewens  
1972;  
Hartl,  
Clark  
2003]



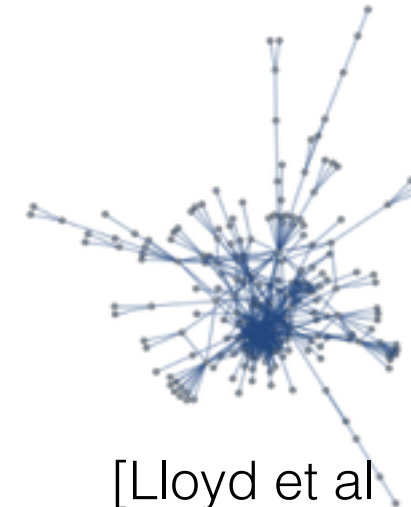
[Datta,  
Banerjee,  
Finley,  
Gelfand  
2016]



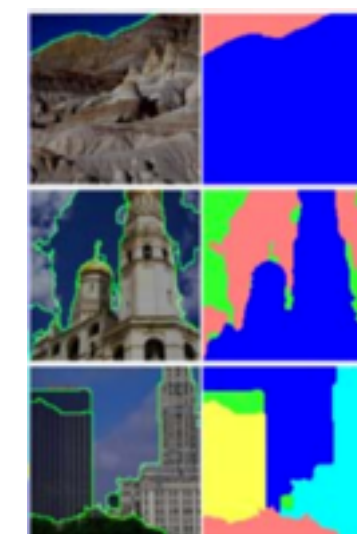
[Fox et al 2014]



[Prabhakaran, Azizi, Carr,  
Pe'er 2016]

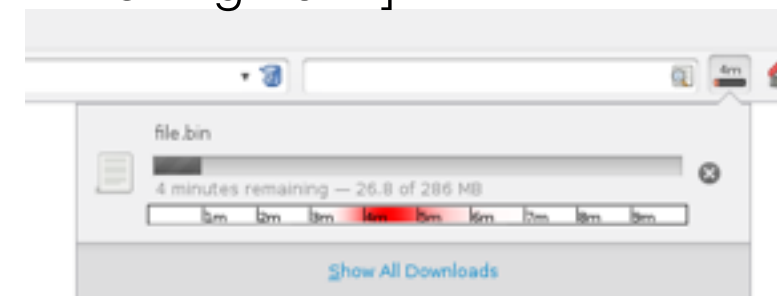


[Lloyd et al  
2012; Miller  
et al 2010]



[Sudderth,  
Jordan 2009]

[Kiefel,  
Schuler,  
Hennig 2014]



[Deisenroth, Fox, Rasmussen 2015]



[Chati,  
Balakrishnan  
2017]



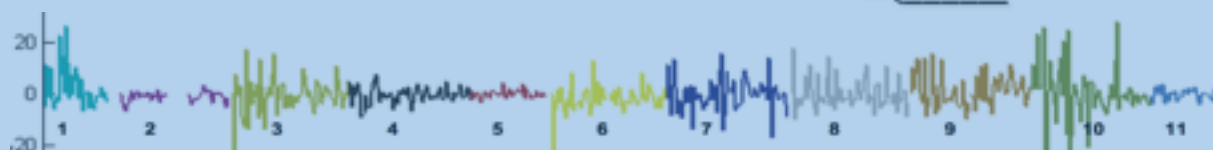


# Applications

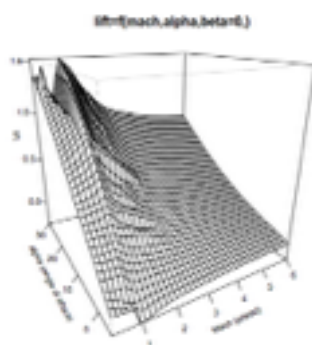


[wikipedia.org]

[Saria et al 2010]



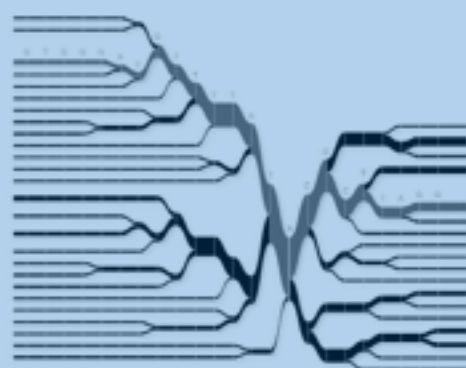
[US CDC PHIL; Futoma, Hariharan, Heller 2017]



[Gramacy, Lee 2009]



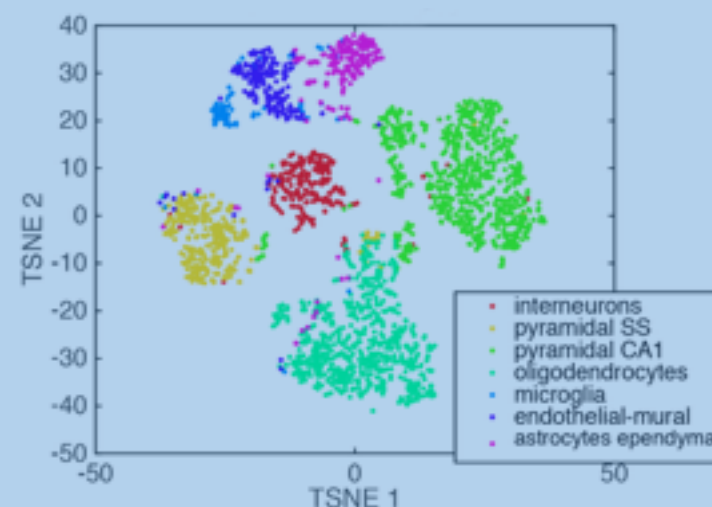
[Ed Bowlby, NOAA]



[Ewens 1972; Hartl, Clark 2003]



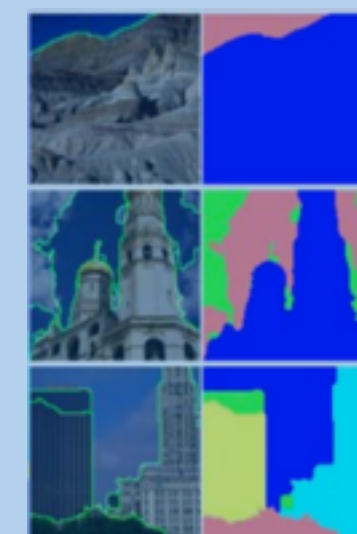
[Fox et al 2014]



[Prabhakaran, Azizi, Carr, Pe'er 2016]

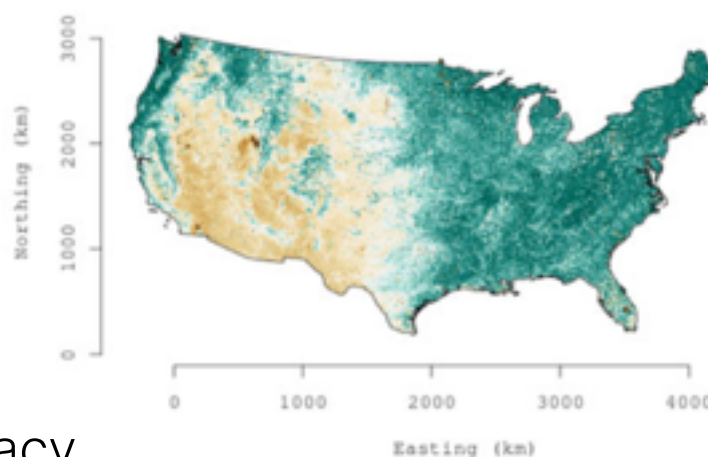


[Lloyd et al 2012; Miller et al 2010]

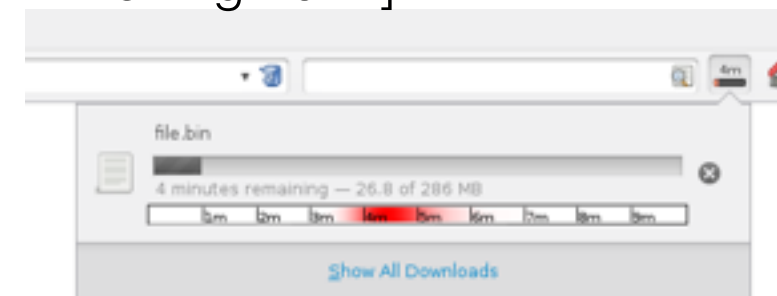


[Sudderth, Jordan 2009]

[Datta, Banerjee, Finley, Gelfand 2016]



[Kiefel, Schuler, Hennig 2014]



[Deisenroth, Fox, Rasmussen 2015]



[Chati, Balakrishnan 2017]



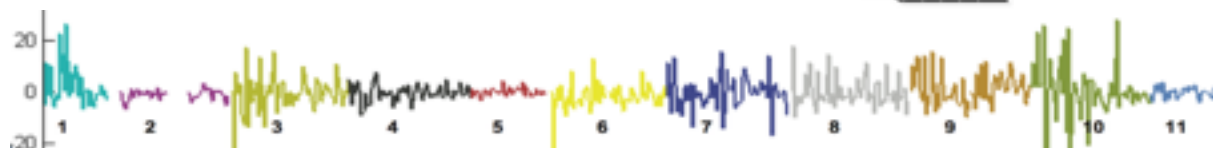


# Applications

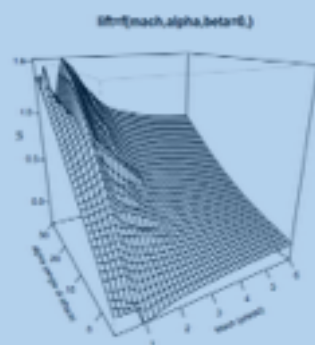


[wikipedia.org]

[Saria  
et al  
2010]



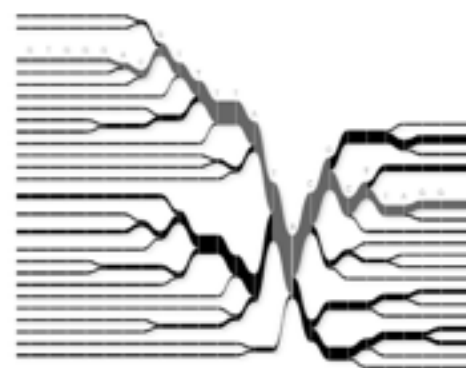
[US CDC PHIL;  
Futoma, Hariharan,  
Heller 2017]



[Gramacy,  
Lee 2009]



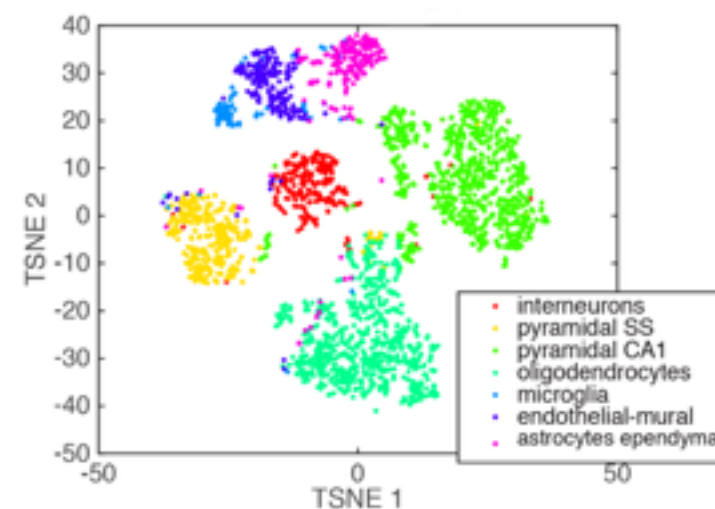
[Ed Bowlby, NOAA]



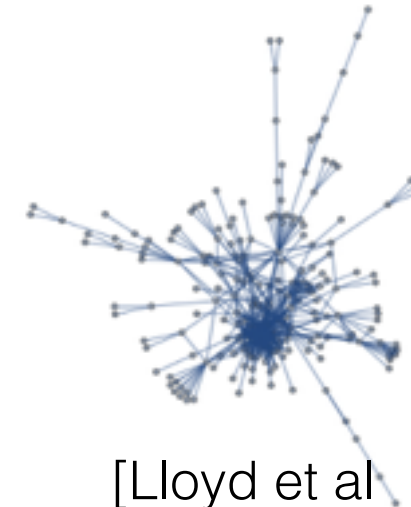
[Ewens  
1972;  
Hartl,  
Clark  
2003]



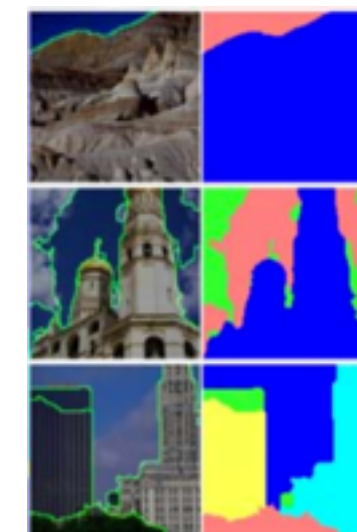
[Fox et al 2014]



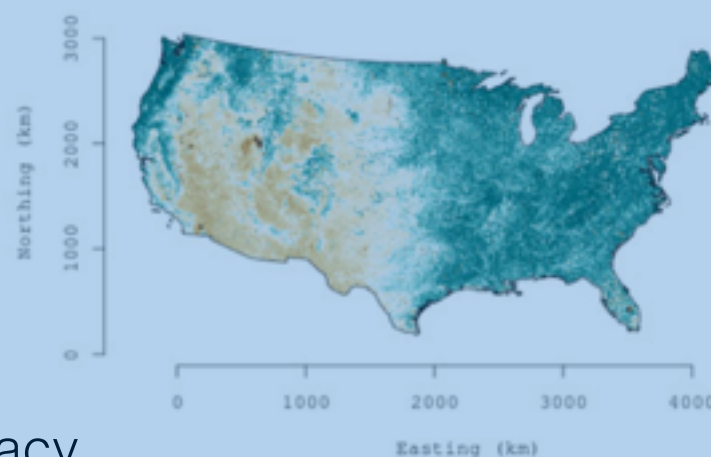
[Prabhakaran, Azizi, Carr,  
Pe'er 2016]



[Lloyd et al  
2012; Miller  
et al 2010]

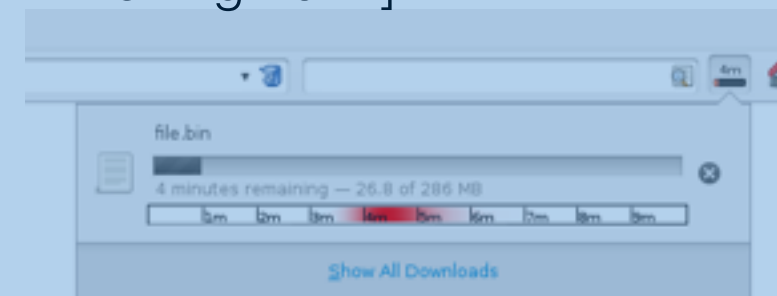


[Sudderth,  
Jordan 2009]



[Datta,  
Banerjee,  
Finley,  
Gelfand  
2016]

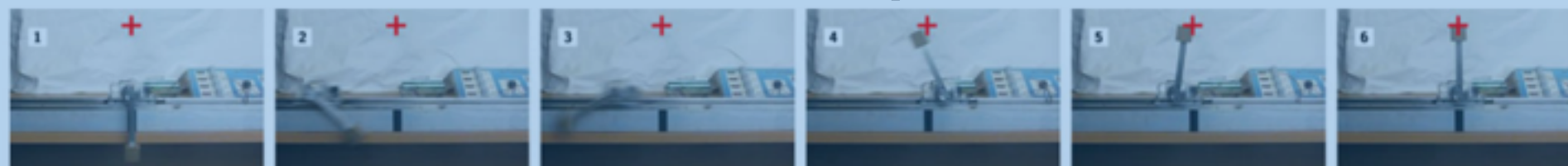
[Kiefel,  
Schuler,  
Hennig 2014]



[Deisenroth, Fox, Rasmussen 2015]



[Chati,  
Balakrishnan  
2017]



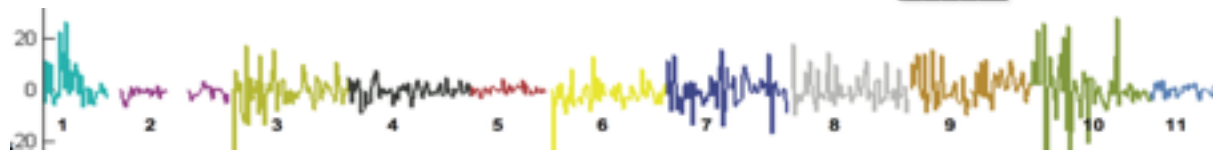


# Applications

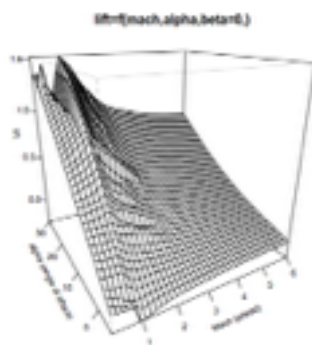


[wikipedia.org]

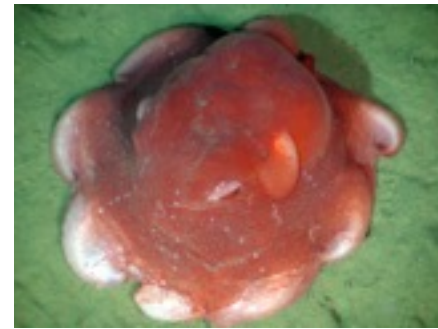
[Saria et al 2010]



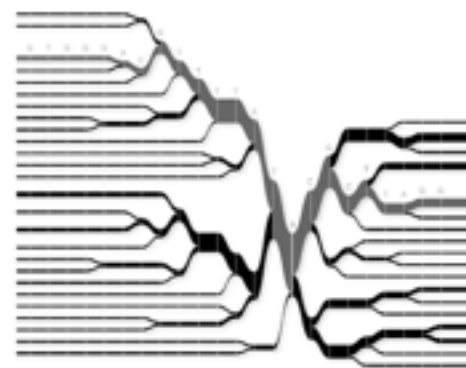
[US CDC PHIL; Futoma, Hariharan, Heller 2017]



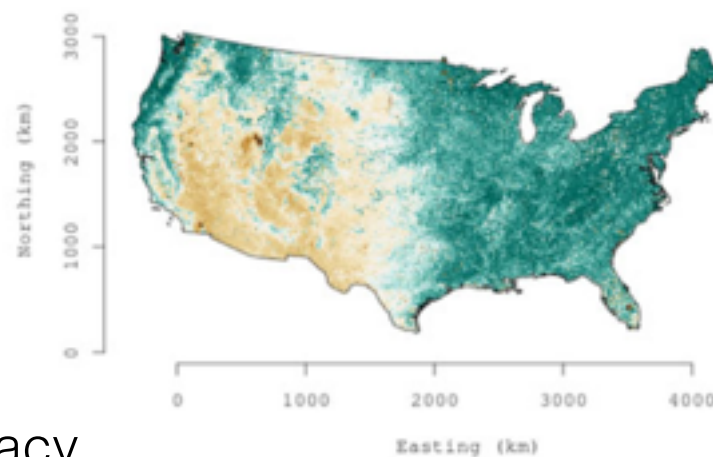
[Gramacy, Lee 2009]



[Ed Bowlby, NOAA]



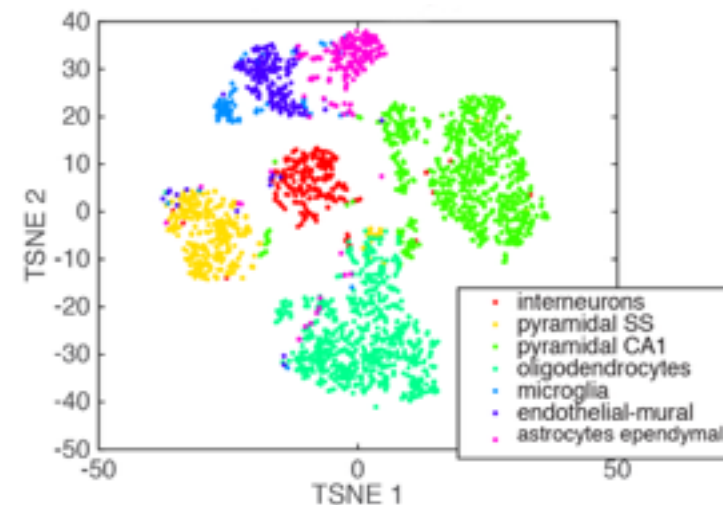
[Ewens 1972; Hartl, Clark 2003]



[Datta, Banerjee, Finley, Gelfand 2016]

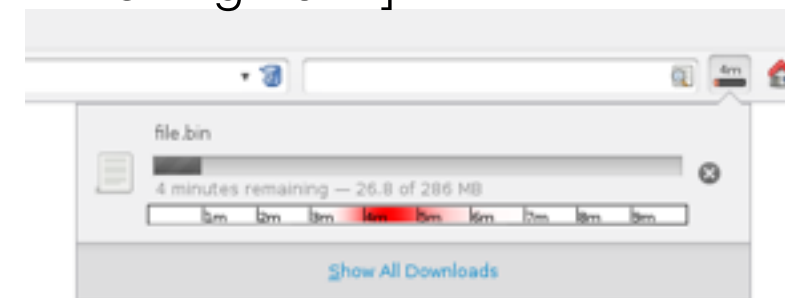


[Fox et al 2014]

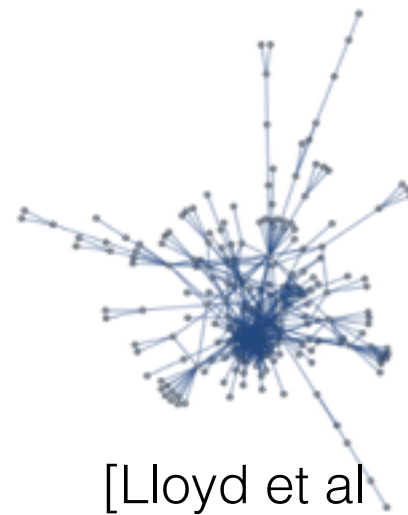


[Prabhakaran, Azizi, Carr, Pe'er 2016]

[Kiefel, Schuler, Hennig 2014]



[Deisenroth, Fox, Rasmussen 2015]



[Lloyd et al 2012; Miller et al 2010]



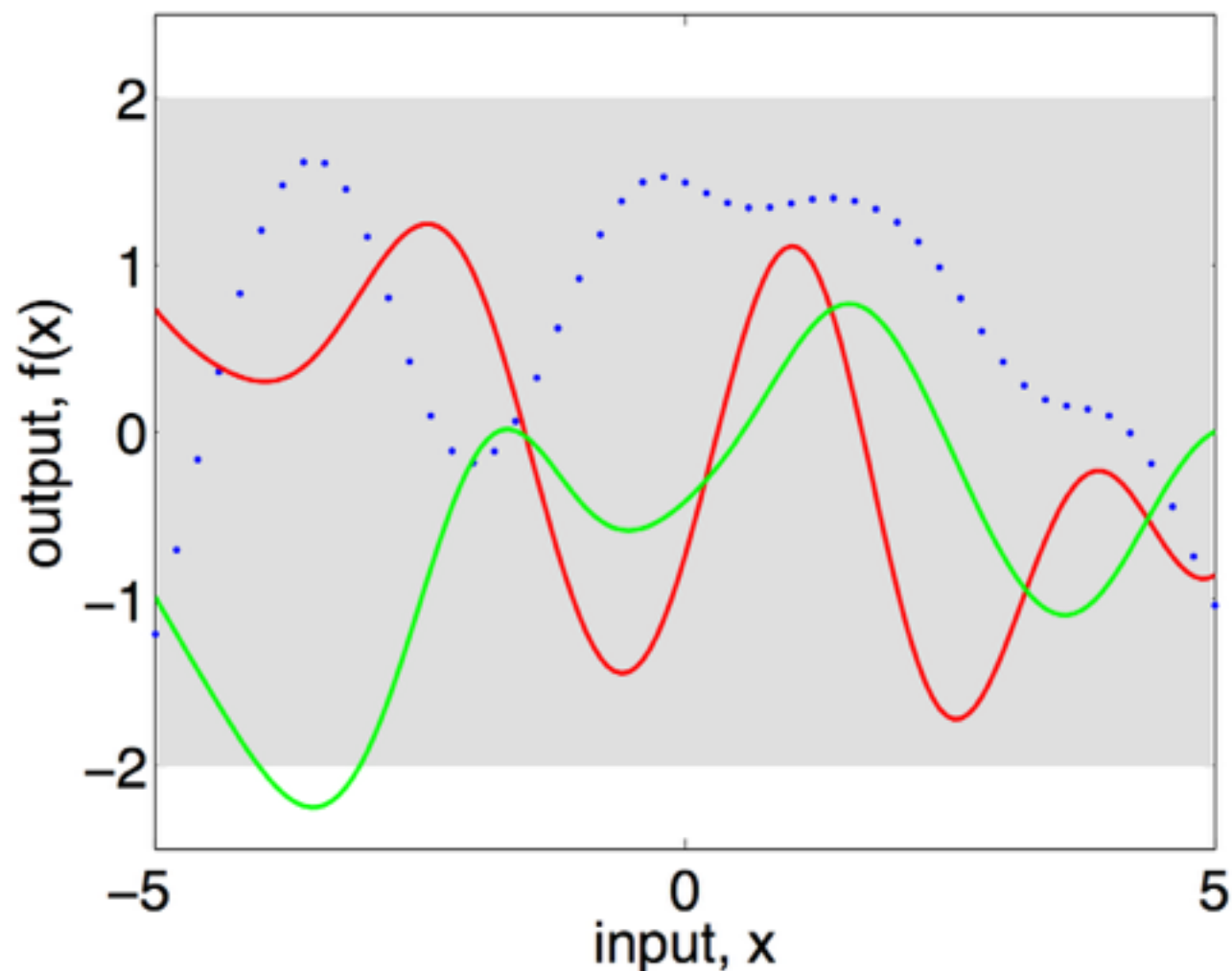
[Sudderth, Jordan 2009]



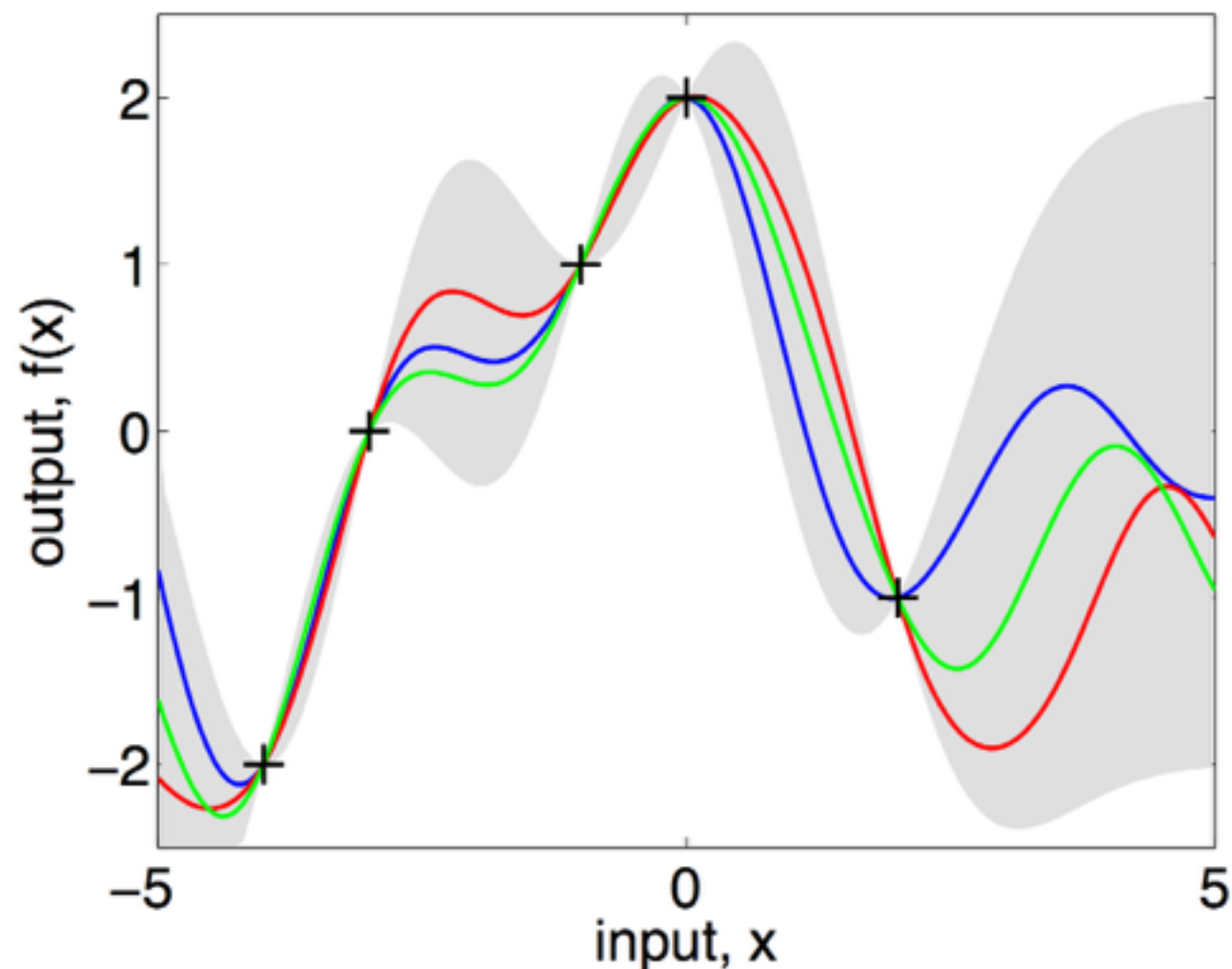
[Chati, Balakrishnan 2017]



# Regression



(a), prior



(b), posterior





Fast inference

Power laws

Hierarchies

Feature allocations

Coalescents/  
Diffusions/Trees

Networks/graphs

Poisson  
processes

de Finetti

*Here be Dragons*



# More Markov Chain Monte Carlo



# More Markov Chain Monte Carlo

- Slice sampling



# More Markov Chain Monte Carlo

- Slice sampling

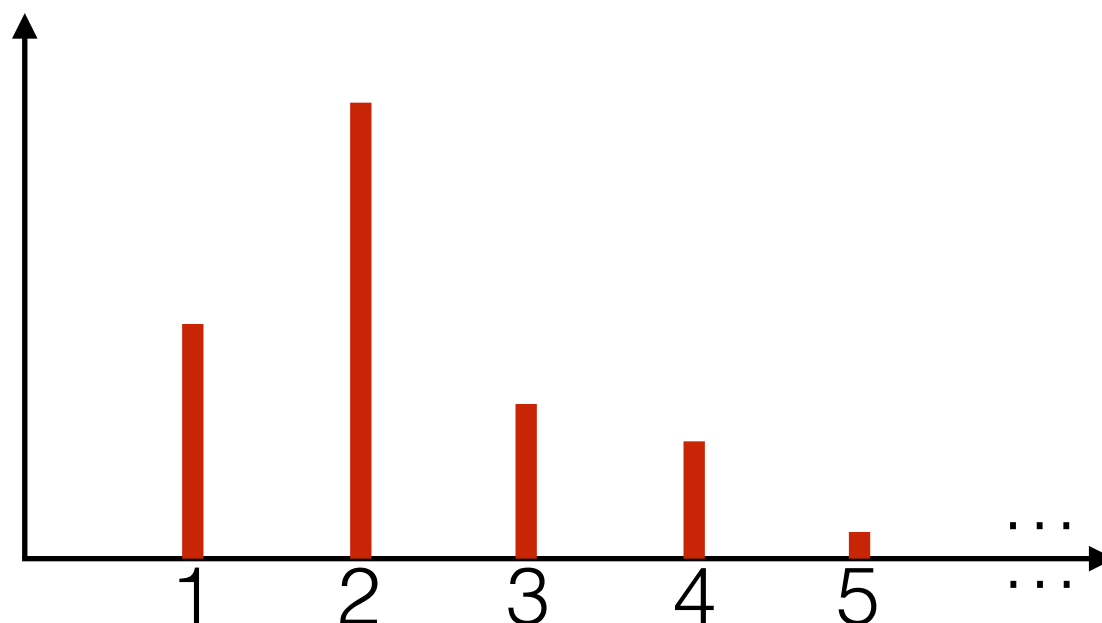


# More Markov Chain Monte Carlo

- Slice sampling
  - auxiliary variable  $\rightarrow$  finite conditionals

# More Markov Chain Monte Carlo

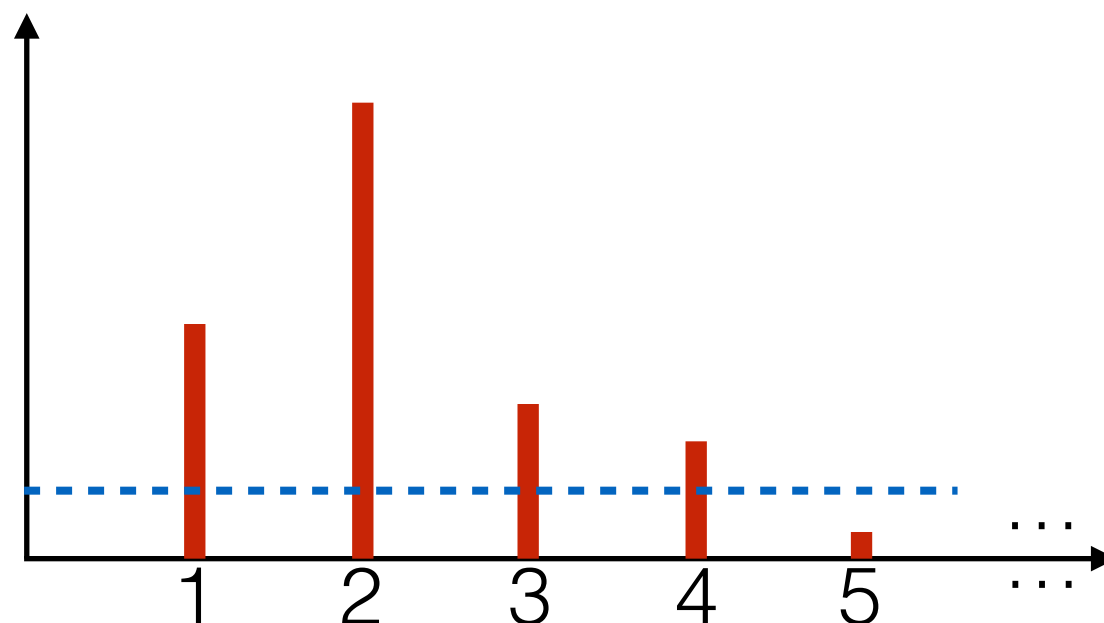
- Slice sampling
  - auxiliary variable  $\rightarrow$  finite conditionals





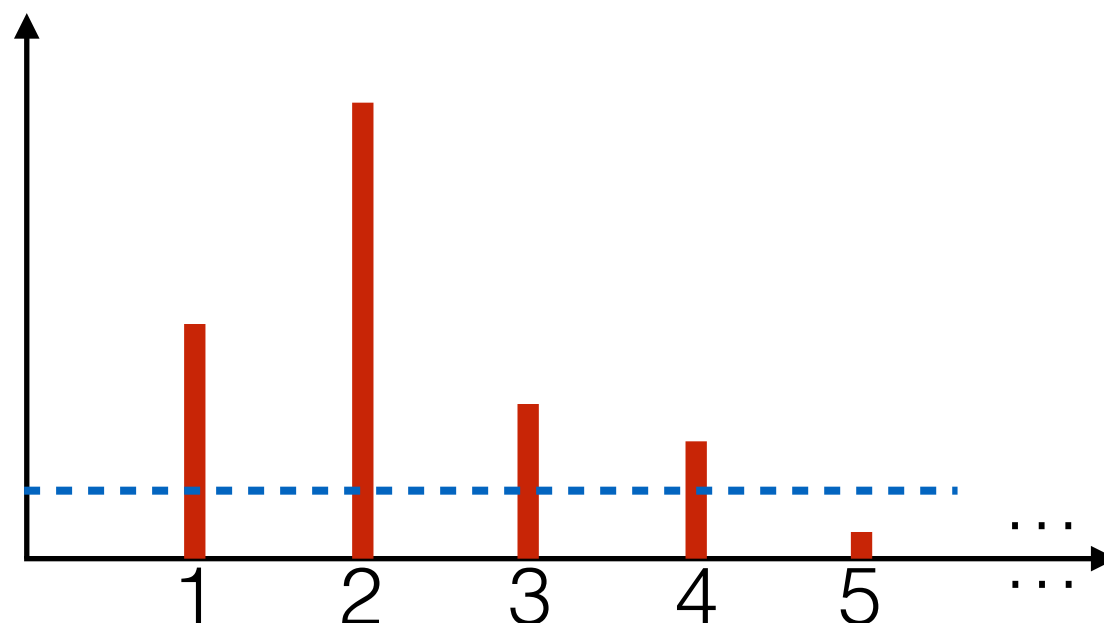
# More Markov Chain Monte Carlo

- Slice sampling
  - auxiliary variable  $\rightarrow$  finite conditionals



# More Markov Chain Monte Carlo

- Slice sampling
  - auxiliary variable  $\rightarrow$  finite conditionals

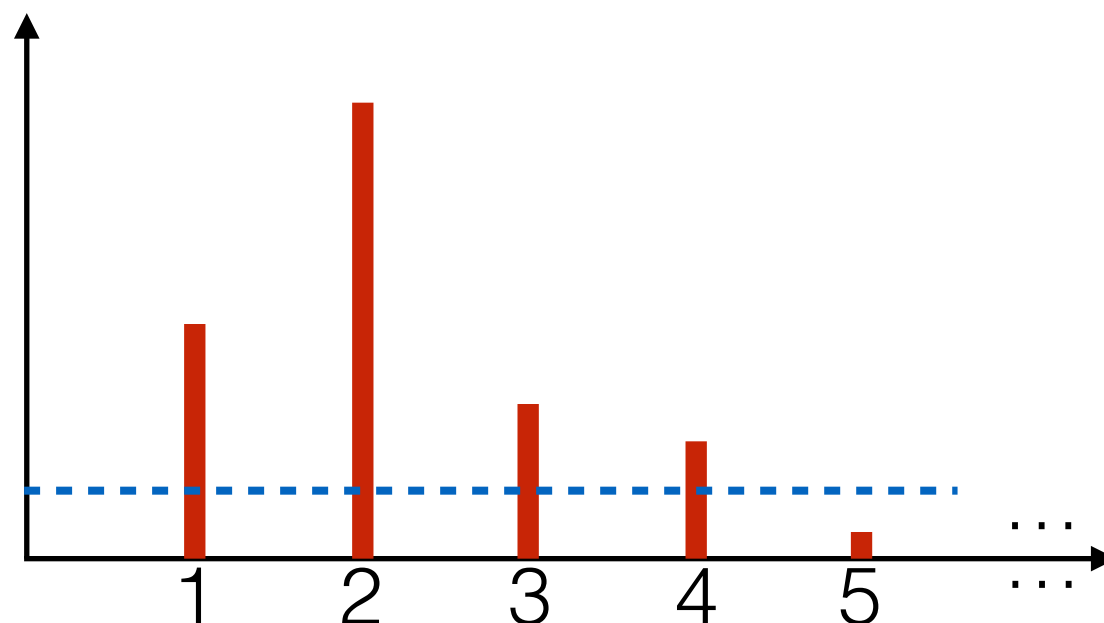


- Approximate with truncated distribution



# More Markov Chain Monte Carlo

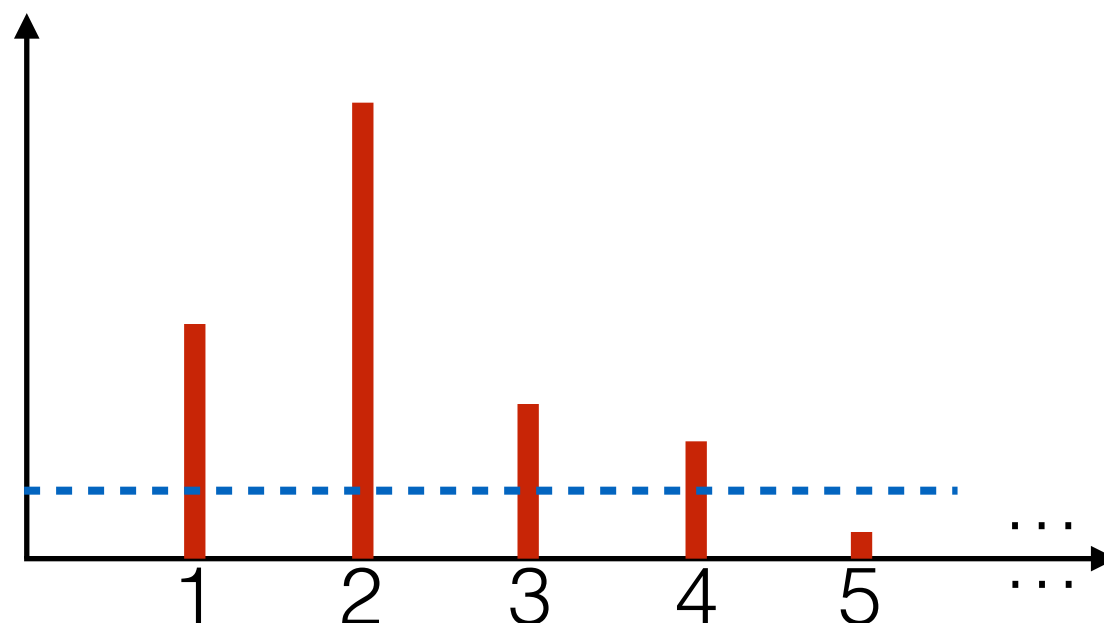
- Slice sampling
  - auxiliary variable  $\rightarrow$  finite conditionals



- Approximate with truncated distribution

# More Markov Chain Monte Carlo

- Slice sampling
  - auxiliary variable  $\rightarrow$  finite conditionals



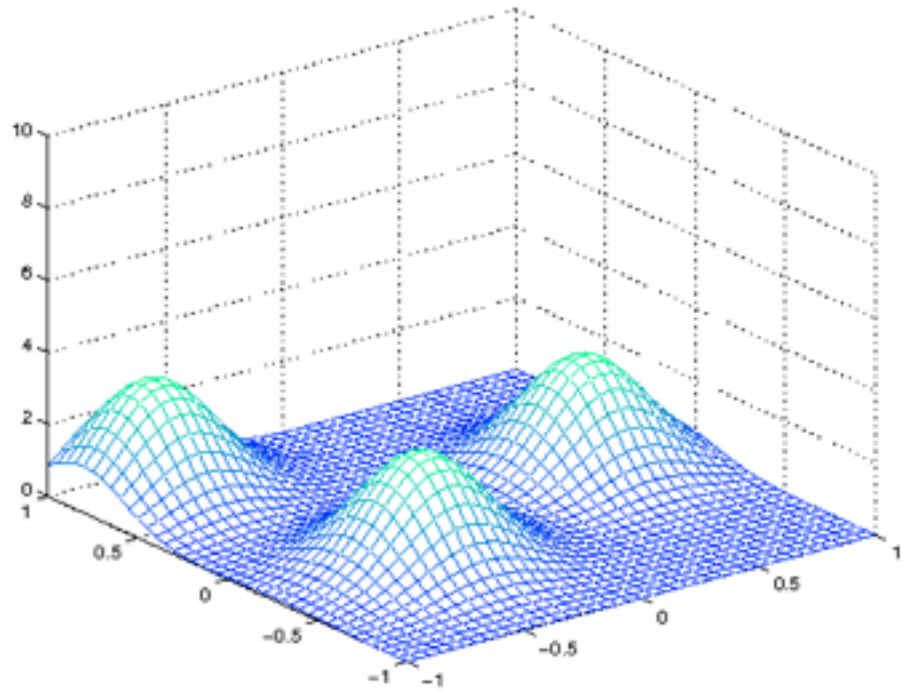
- Approximate with truncated distribution
  - E.g., Hamiltonian Monte Carlo



# Variational Bayes

# Variational Bayes

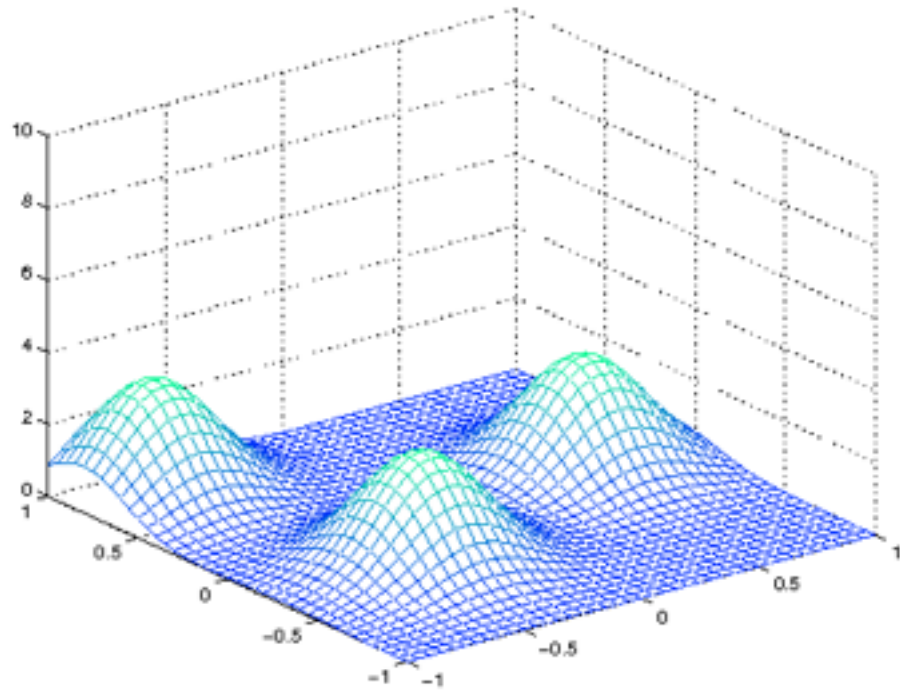
- Variational Bayes (VB)





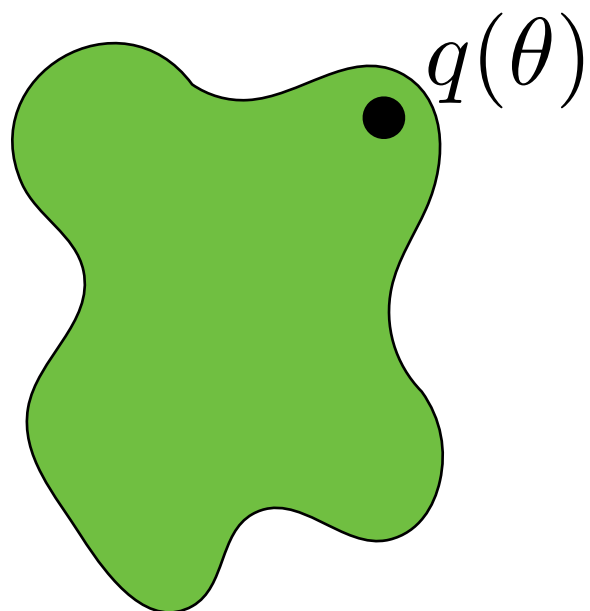
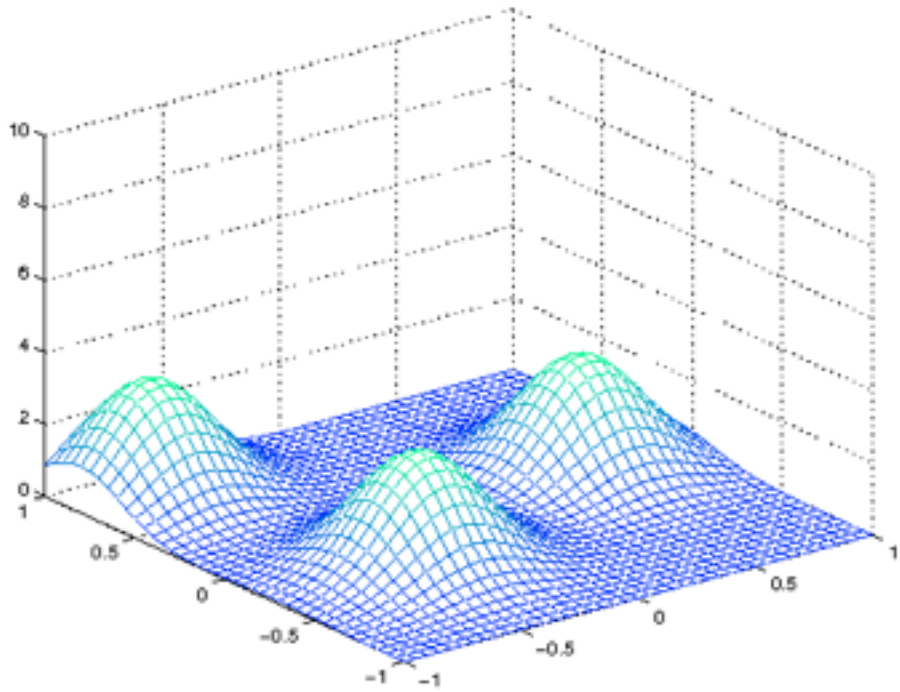
# Variational Bayes

- Variational Bayes (VB)
  - Approximation  $q^*(\theta)$  for posterior  $p(\theta|x)$



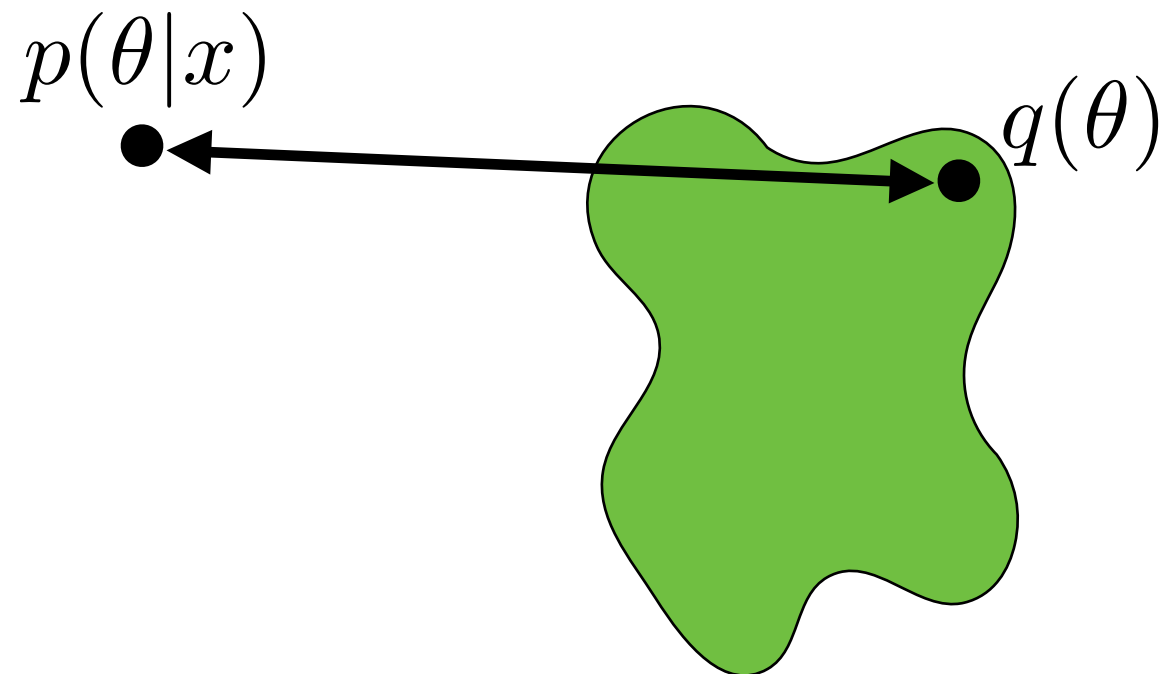
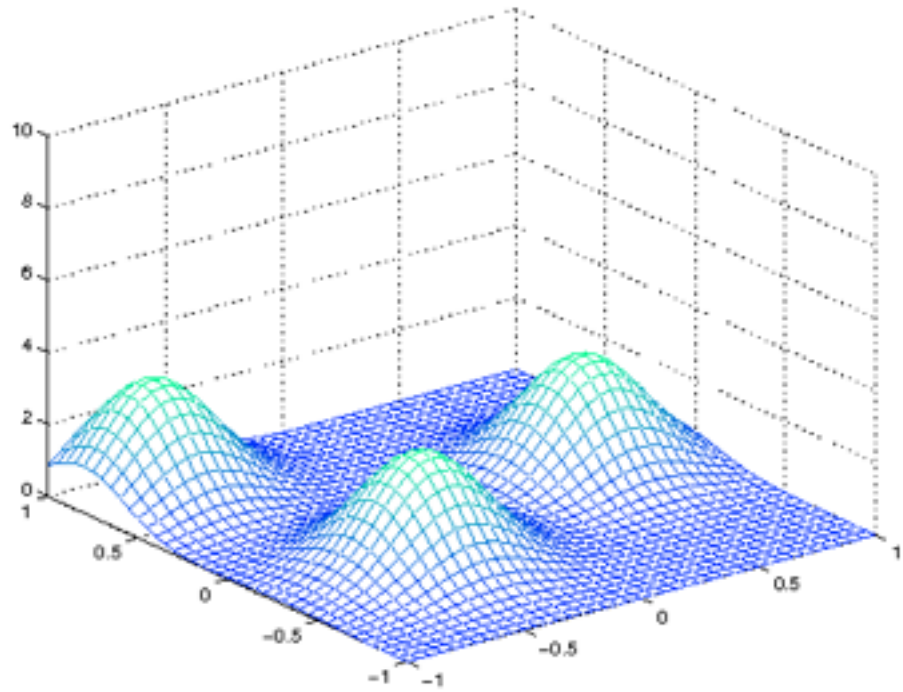
# Variational Bayes

- Variational Bayes (VB)
  - Approximation  $q^*(\theta)$  for posterior  $p(\theta|x)$



# Variational Bayes

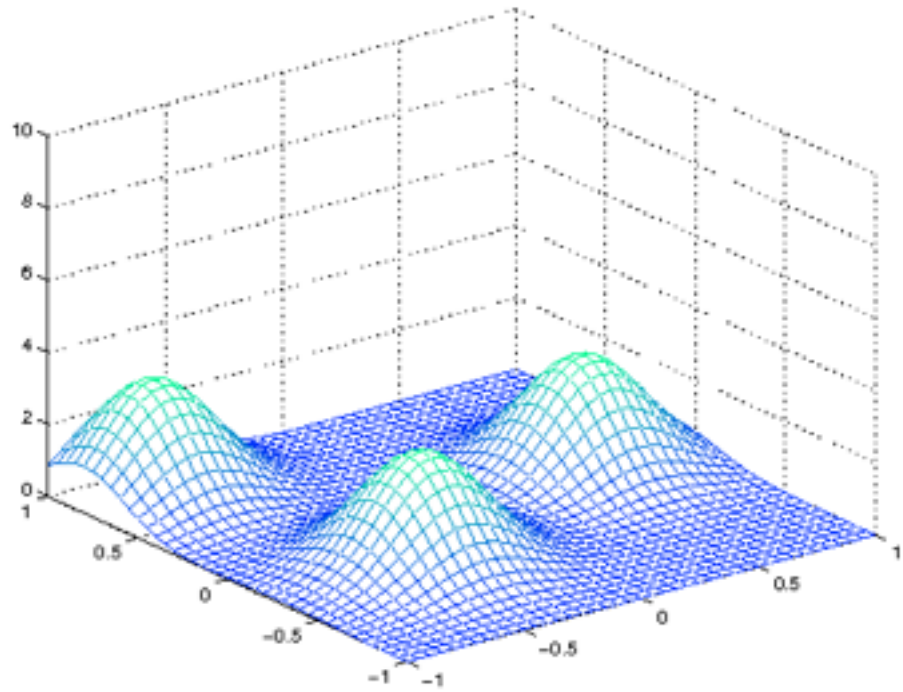
- Variational Bayes (VB)
  - Approximation  $q^*(\theta)$  for posterior  $p(\theta|x)$



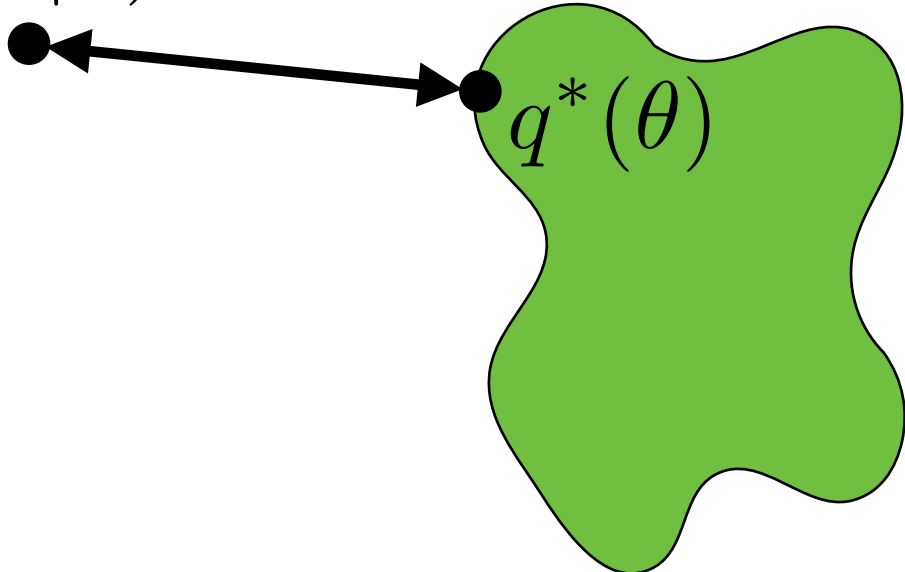


# Variational Bayes

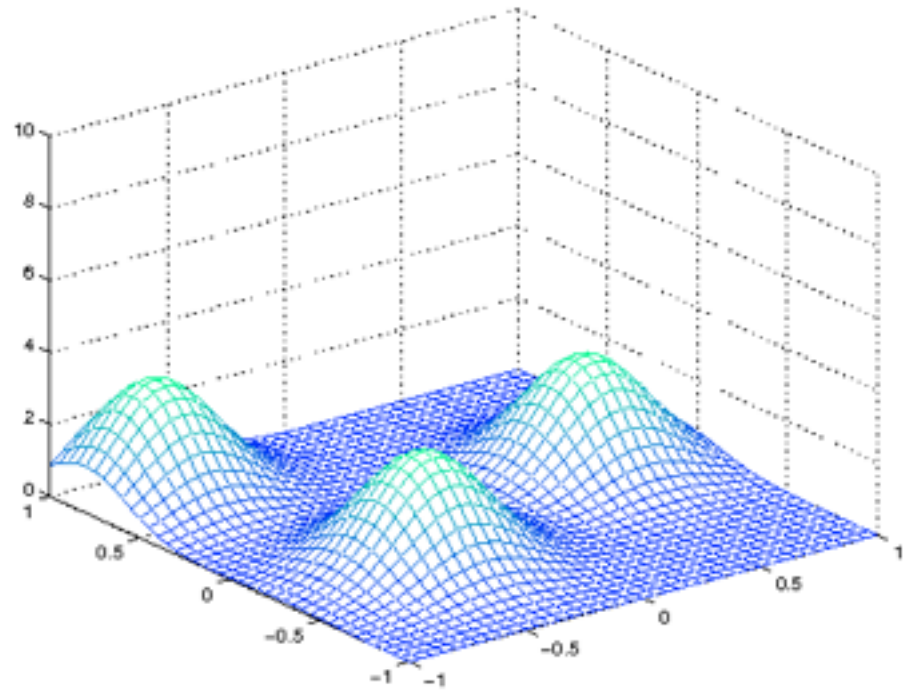
- Variational Bayes (VB)
  - Approximation  $q^*(\theta)$  for posterior  $p(\theta|x)$



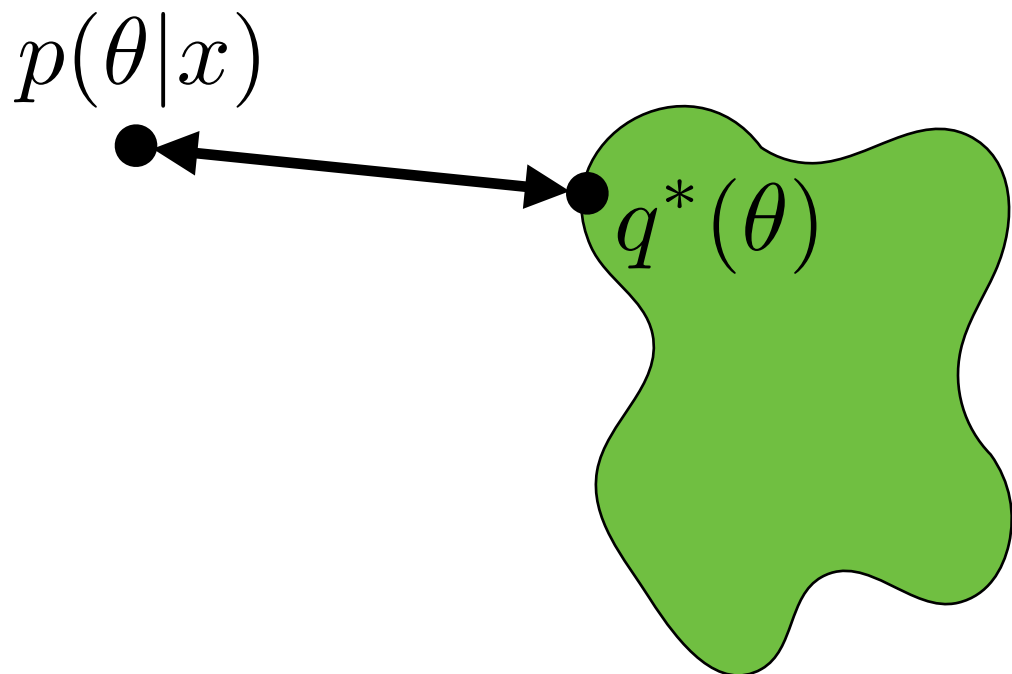
$p(\theta|x)$



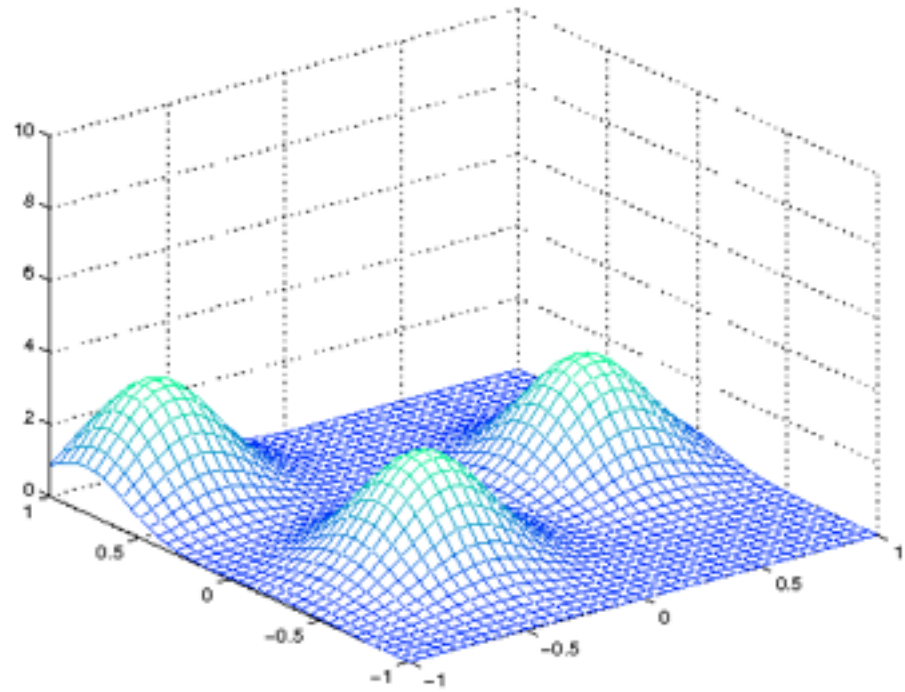
# Variational Bayes



- Variational Bayes (VB)
  - Approximation  $q^*(\theta)$  for posterior  $p(\theta|x)$
  - “Close”: Minimize Kullback-Liebler (KL) divergence:
$$KL(q||p(\cdot|x))$$

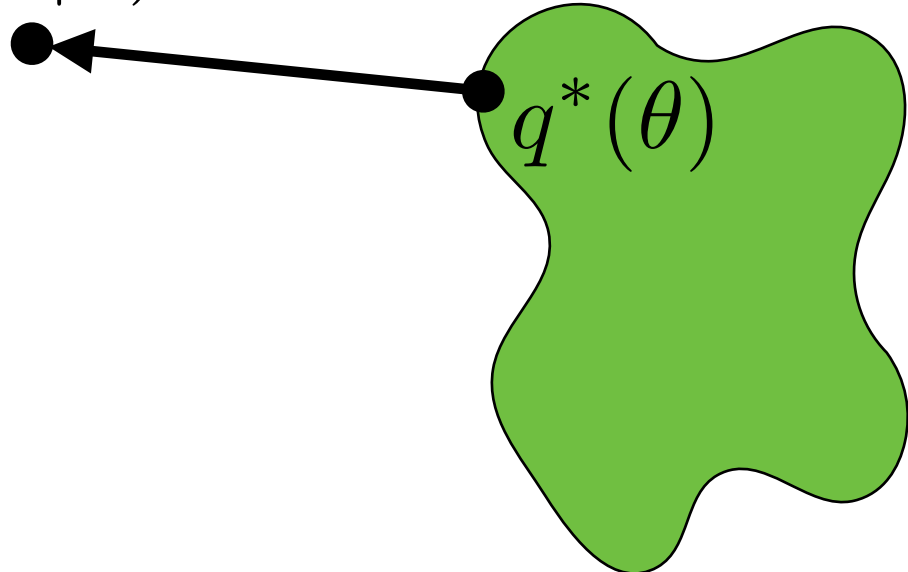


# Variational Bayes



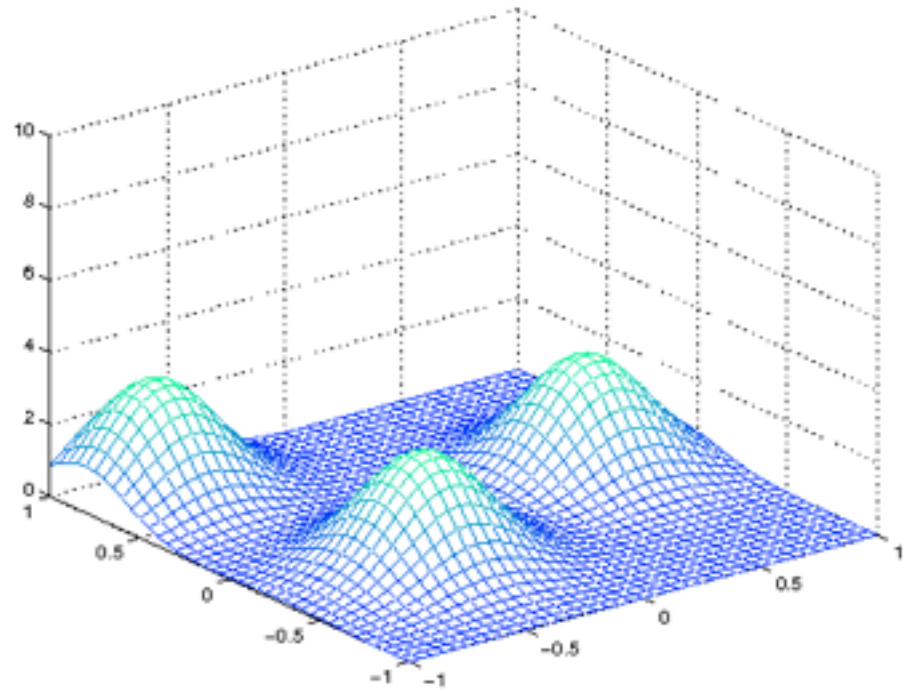
- Variational Bayes (VB)
  - Approximation  $q^*(\theta)$  for posterior  $p(\theta|x)$
  - “Close”: Minimize Kullback-Liebler (KL) divergence:
$$KL(q||p(\cdot|x))$$

$p(\theta|x)$



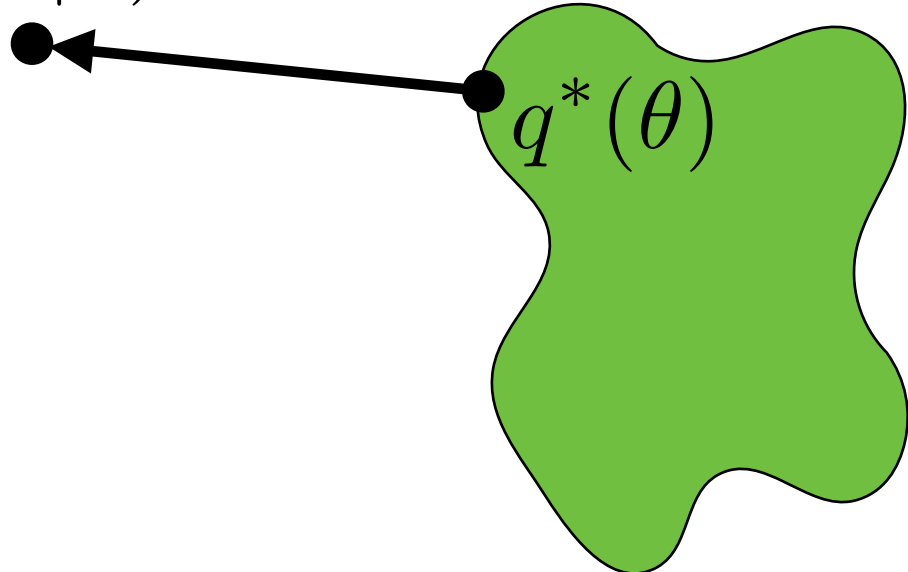


# Variational Bayes

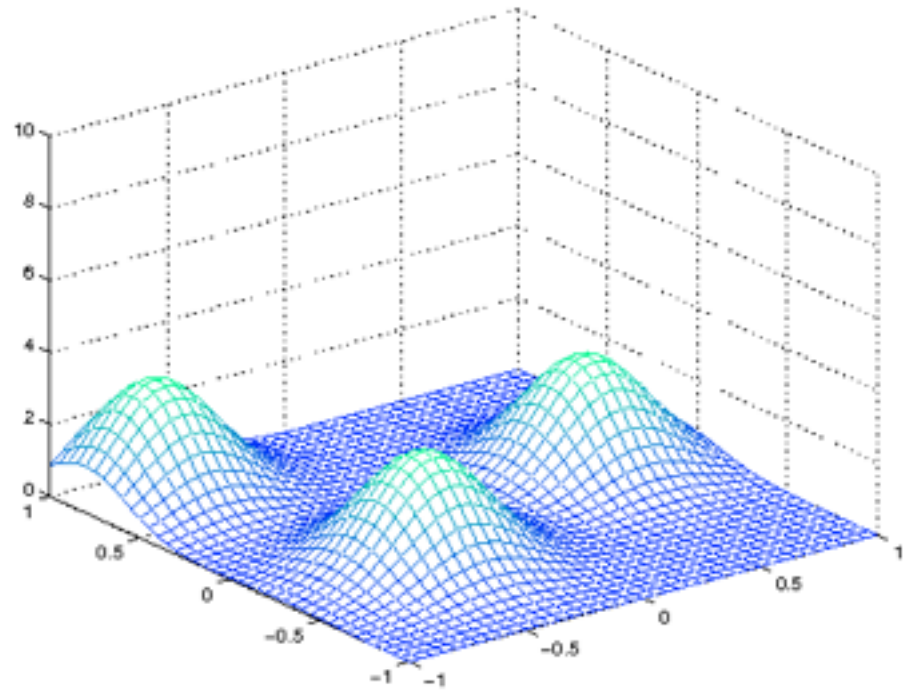


- Variational Bayes (VB)
  - Approximation  $q^*(\theta)$  for posterior  $p(\theta|x)$
  - “Close”: Minimize Kullback-Liebler (KL) divergence:
$$KL(q||p(\cdot|x))$$
  - “Nice”: factorizes, exponential family, truncation

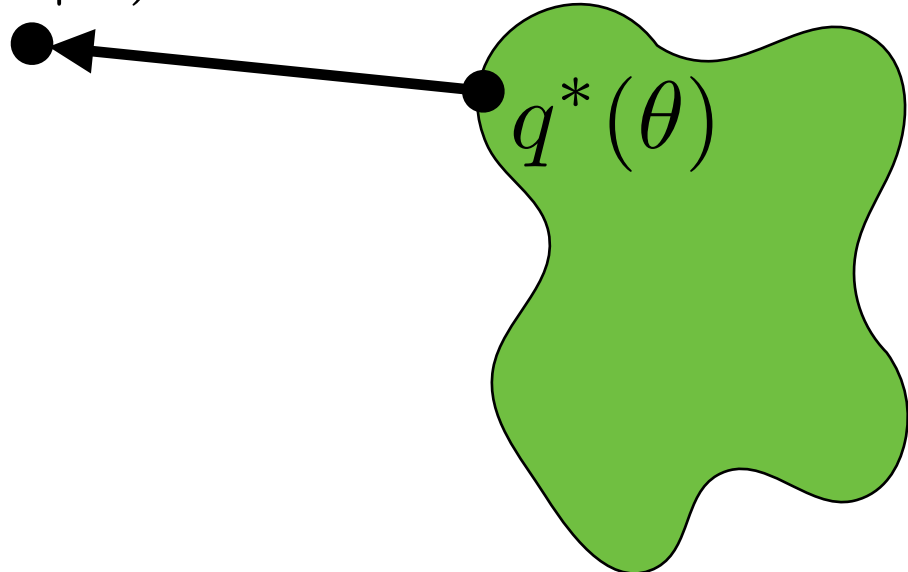
$p(\theta|x)$



# Variational Bayes

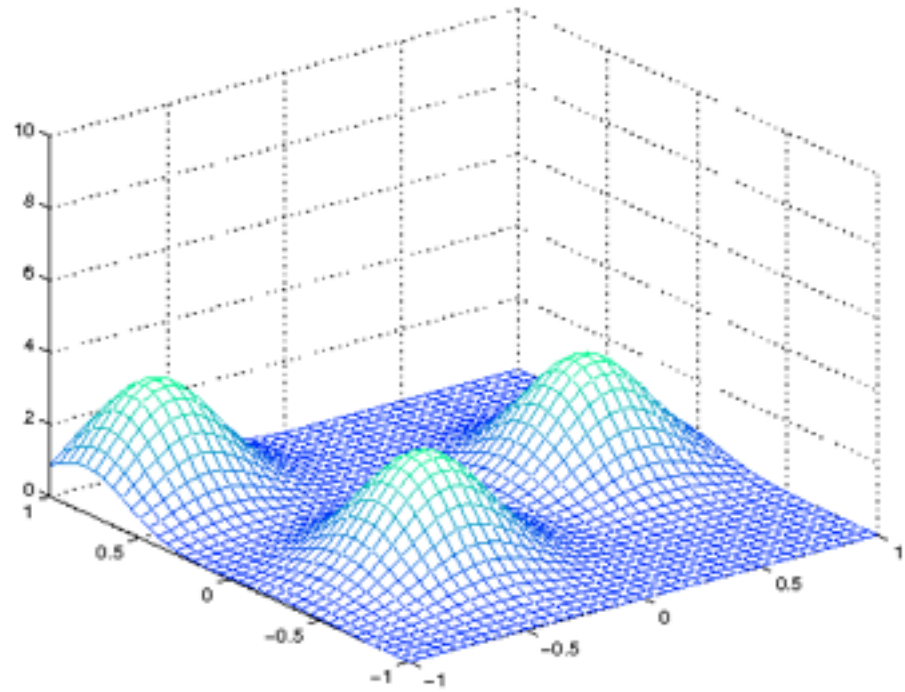


$p(\theta|x)$

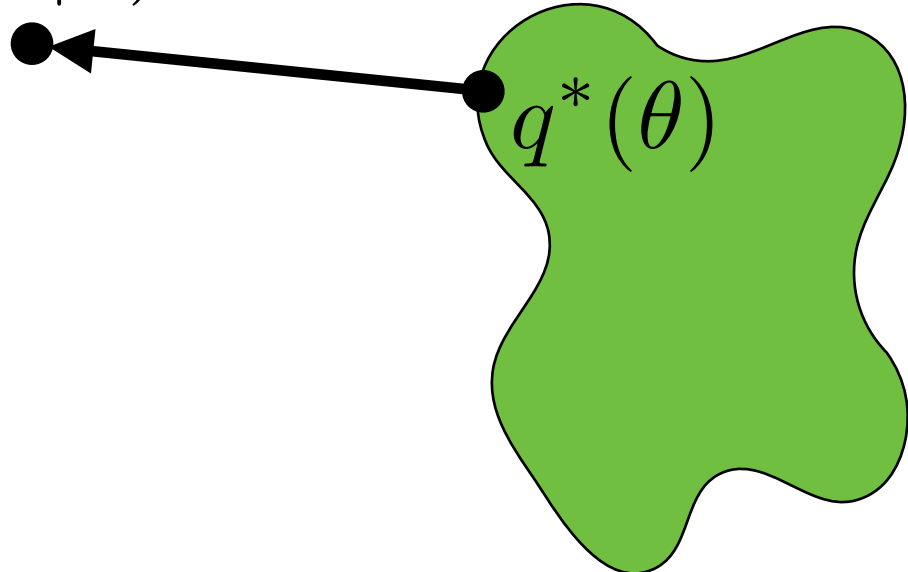


- Variational Bayes (VB)
  - Approximation  $q^*(\theta)$  for posterior  $p(\theta|x)$
  - “Close”: Minimize Kullback-Liebler (KL) divergence:
$$KL(q||p(\cdot|x))$$
  - “Nice”: factorizes, exponential family, truncation
- VB practical success

# Variational Bayes



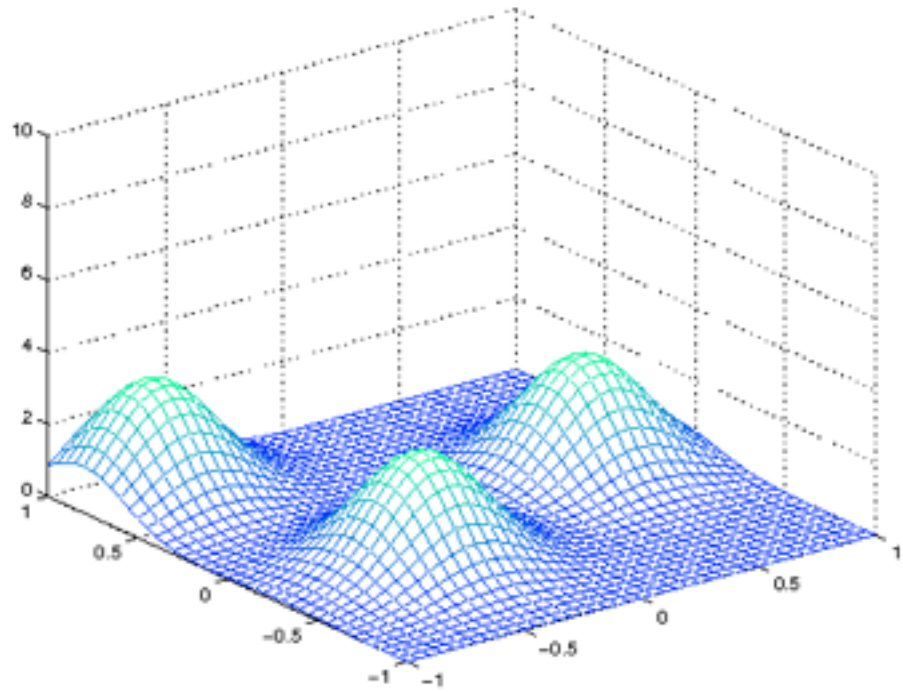
$p(\theta|x)$



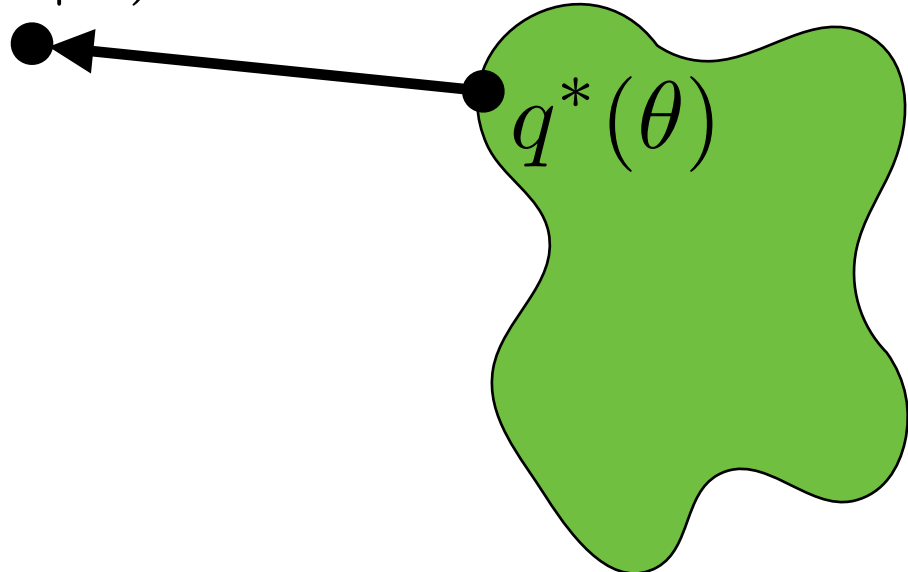
- Variational Bayes (VB)
  - Approximation  $q^*(\theta)$  for posterior  $p(\theta|x)$
  - “Close”: Minimize Kullback-Liebler (KL) divergence:
$$KL(q||p(\cdot|x))$$
  - “Nice”: factorizes, exponential family, truncation
- VB practical success
  - point estimates and prediction



# Variational Bayes

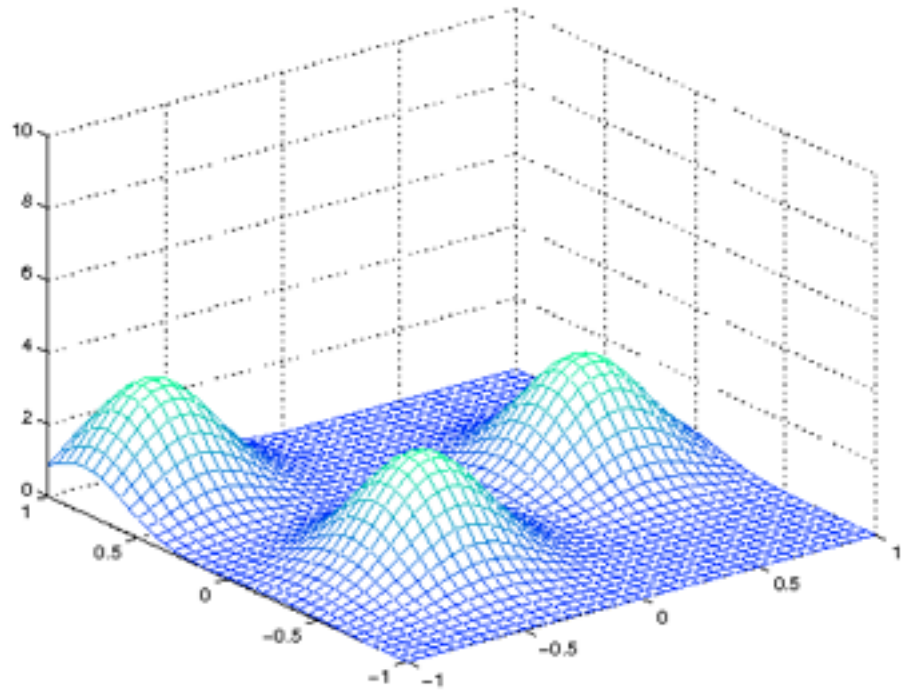


$p(\theta|x)$

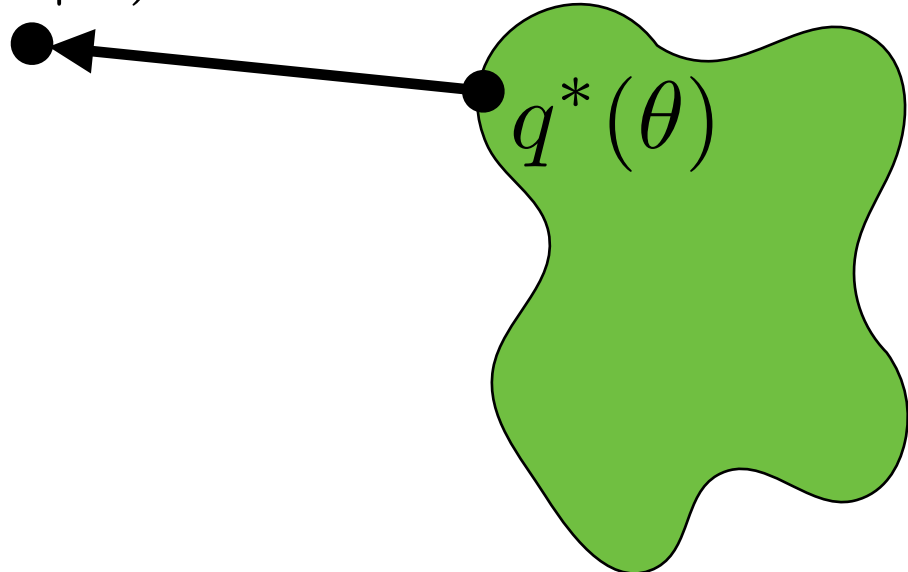


- Variational Bayes (VB)
  - Approximation  $q^*(\theta)$  for posterior  $p(\theta|x)$
  - “Close”: Minimize Kullback-Liebler (KL) divergence:
$$KL(q||p(\cdot|x))$$
  - “Nice”: factorizes, exponential family, truncation
- VB practical success
  - point estimates and prediction
  - fast, streaming, distributed

# Variational Bayes



$p(\theta|x)$



- Variational Bayes (VB)
  - Approximation  $q^*(\theta)$  for posterior  $p(\theta|x)$
  - “Close”: Minimize Kullback-Liebler (KL) divergence:
$$KL(q||p(\cdot|x))$$
  - “Nice”: factorizes, exponential family, truncation
- VB practical success
  - point estimates and prediction
  - fast, streaming, distributed
  - can underestimate uncertainties

[Broderick, Boyd, Wibisono, Wilson, Jordan 2013;

Giordano, Broderick, Jordan 2015; Huggins, Campbell, Broderick 2016]

# Clustering

	Arts	Econ	Sports	Health	Technology
Document 1					
Document 2					
Document 3					
Document 4					
Document 5					
Document 6					
Document 7					

# Feature allocation

	Arts	Econ	Sports	Health	Technology
Document 1					
Document 2					
Document 3					
Document 4					
Document 5					
Document 6					
Document 7					



# Feature allocation

	Arts	Econ	Sports	Health	Technology
Document 1					
Document 2					
Document 3					
Document 4					
Document 5					
Document 6					
Document 7					

- Indian buffet process

# Feature allocation

	Arts	Econ	Sports	Health	Technology
Document 1					
Document 2					
Document 3					
Document 4					
Document 5					
Document 6					
Document 7					

- Indian buffet process

# Feature allocation

	Arts	Econ	Sports	Health	Technology
Document 1					
Document 2					
Document 3					
Document 4					
Document 5					
Document 6					
Document 7					

- Indian buffet process
- Beta process

# Feature allocation

	Arts	Econ	Sports	Health	Technology
Document 1					
Document 2					
Document 3					
Document 4					
Document 5					
Document 6					
Document 7					

- Indian buffet process
- Beta process

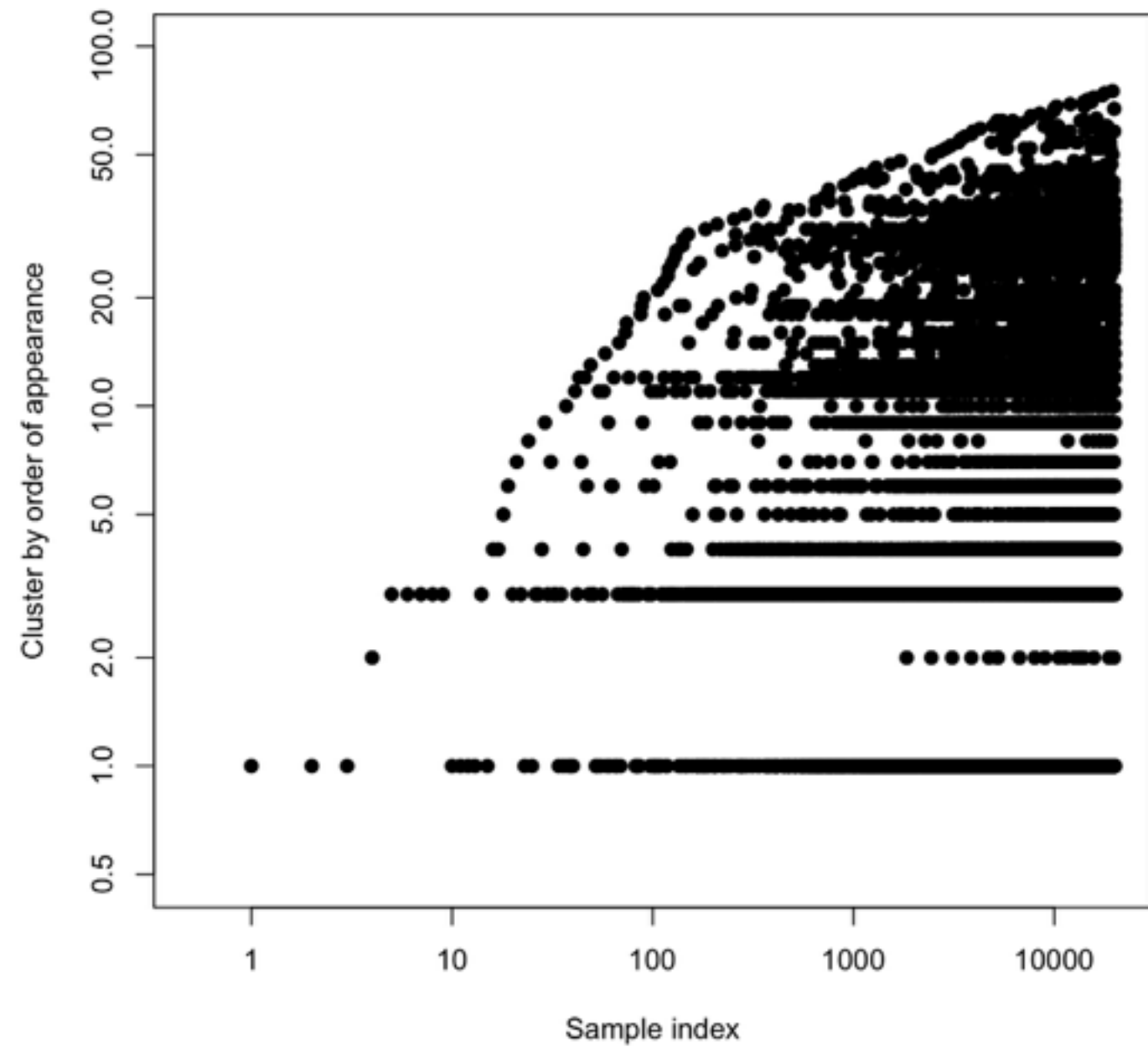


# Feature allocation

	Arts	Econ	Sports	Health	Technology
Document 1					
Document 2					
Document 3					
Document 4					
Document 5					
Document 6					
Document 7					

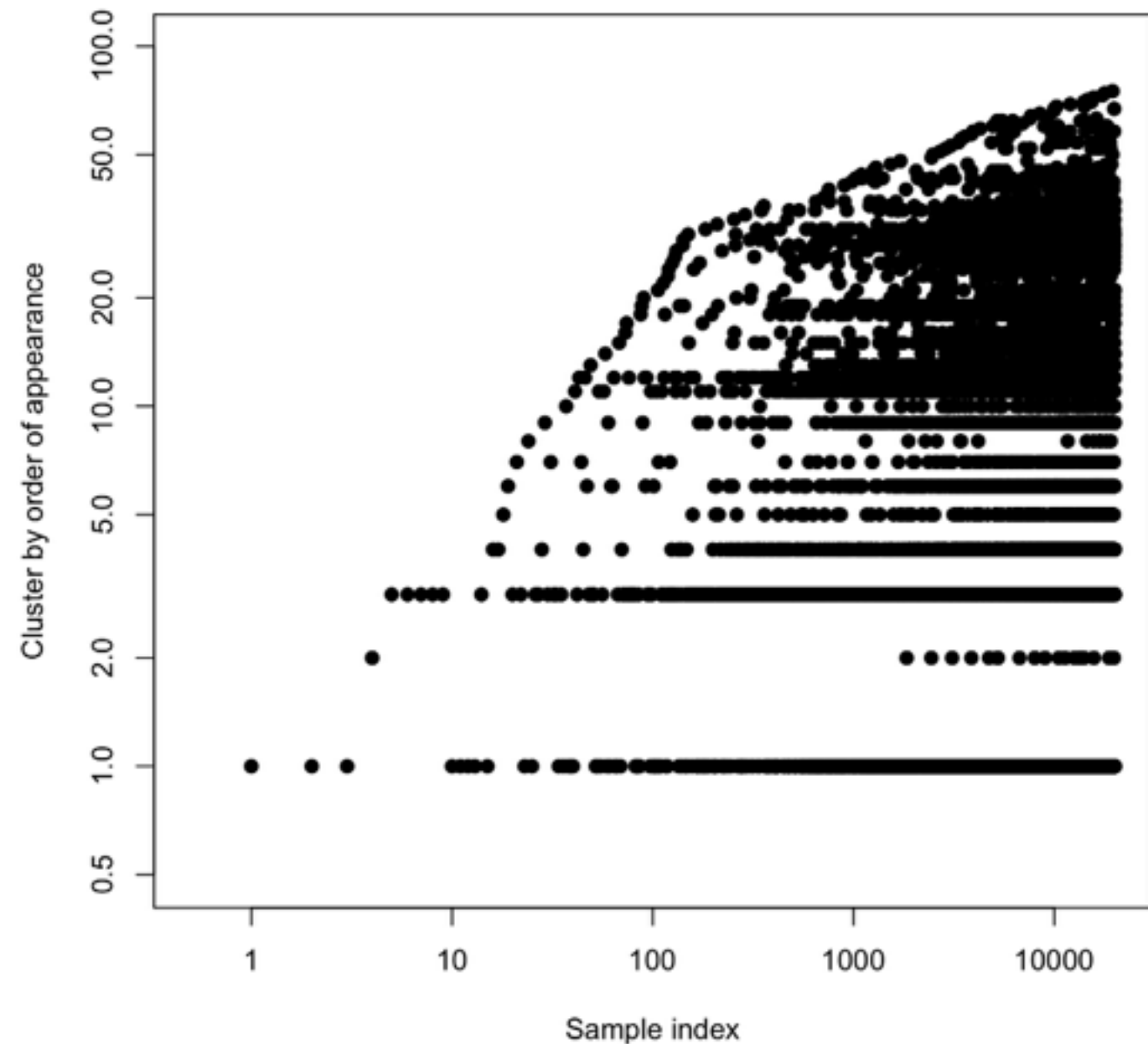
- Indian buffet process
- Beta process

# Power laws



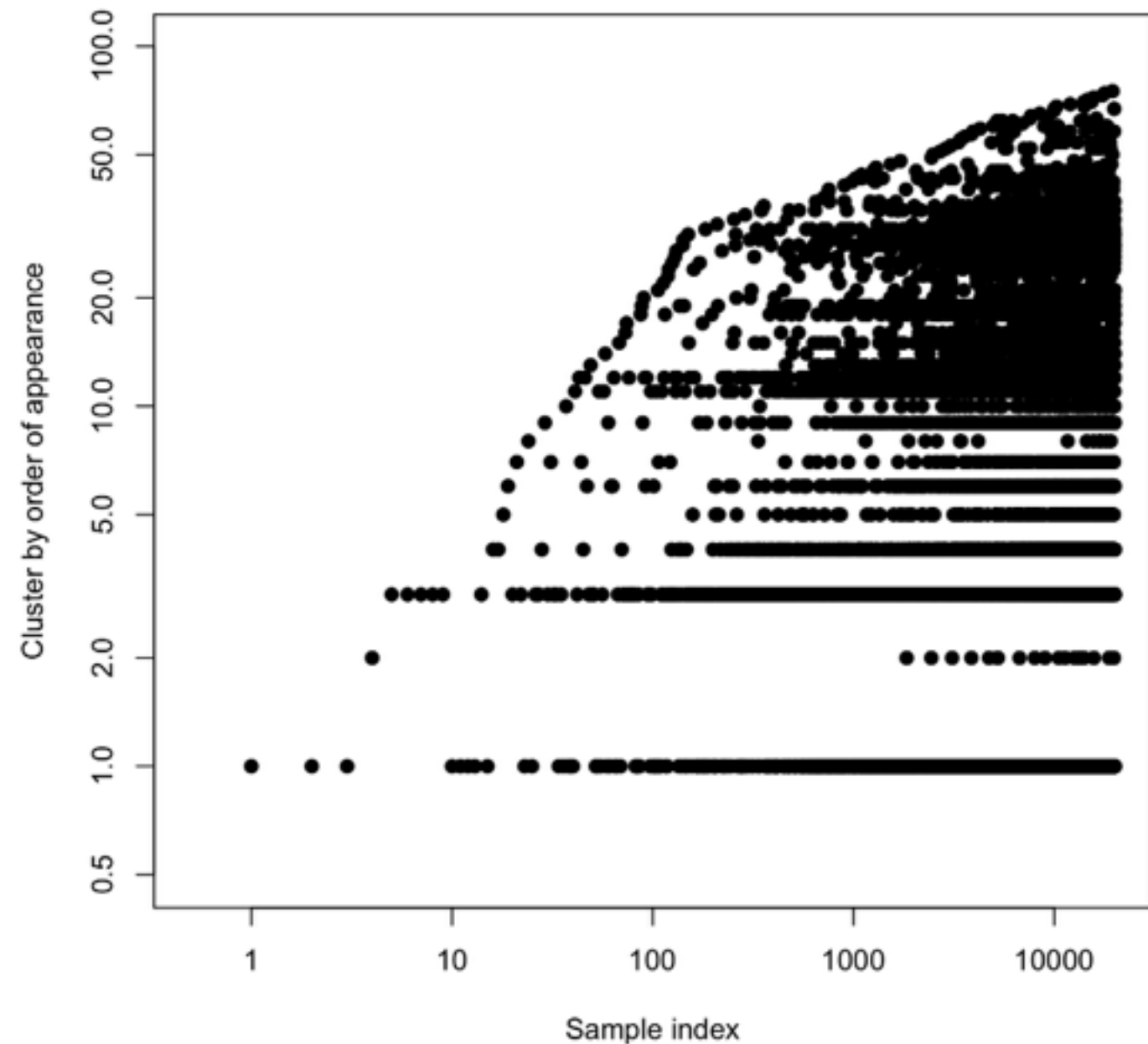
# Power laws

- $K_N := \#$  clusters occupied by  $N$  data points



# Power laws

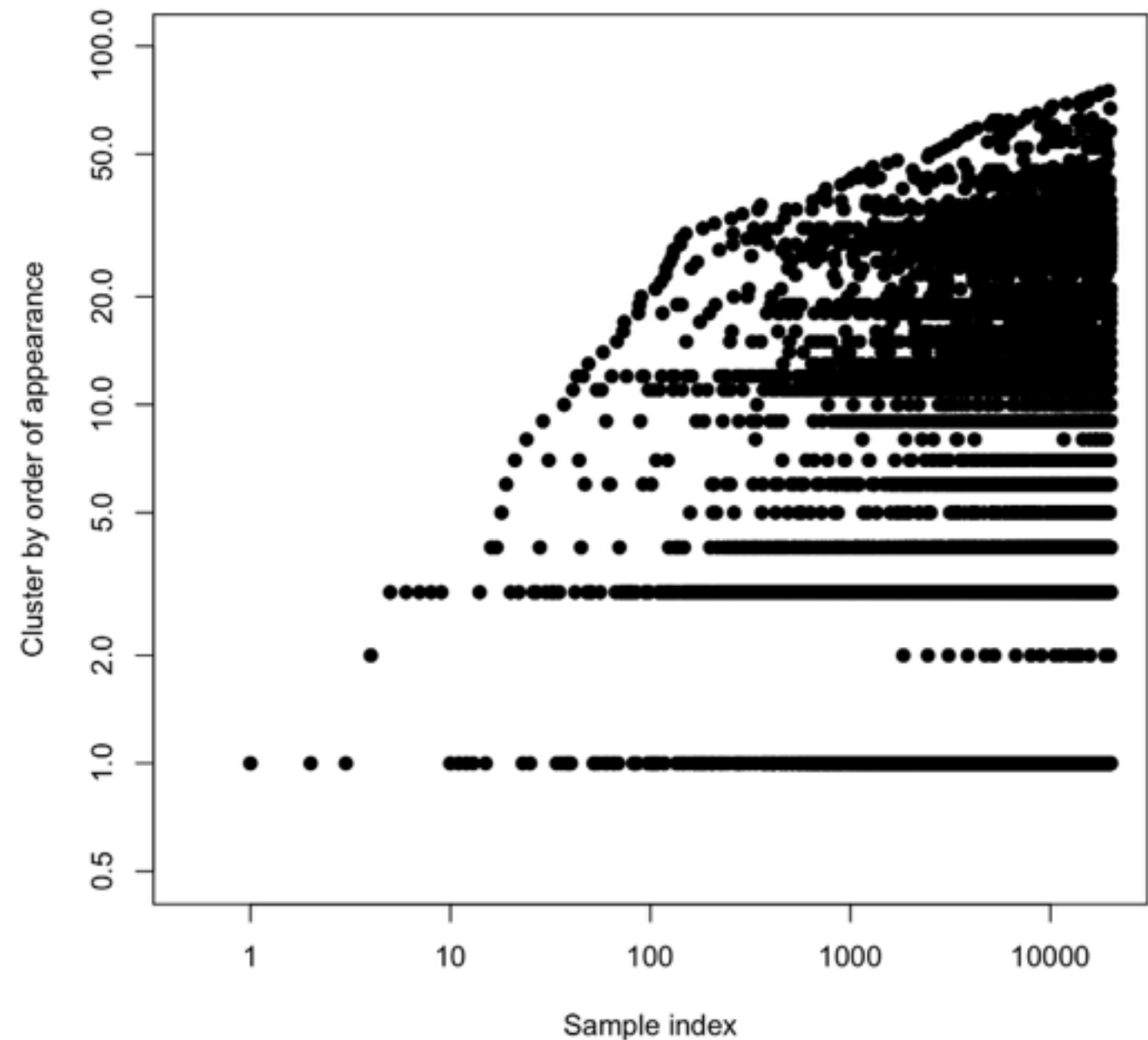
- $K_N := \#$  clusters occupied by  $N$  data points
- CRP:  $K_N \sim \alpha \log N$  w.p. 1





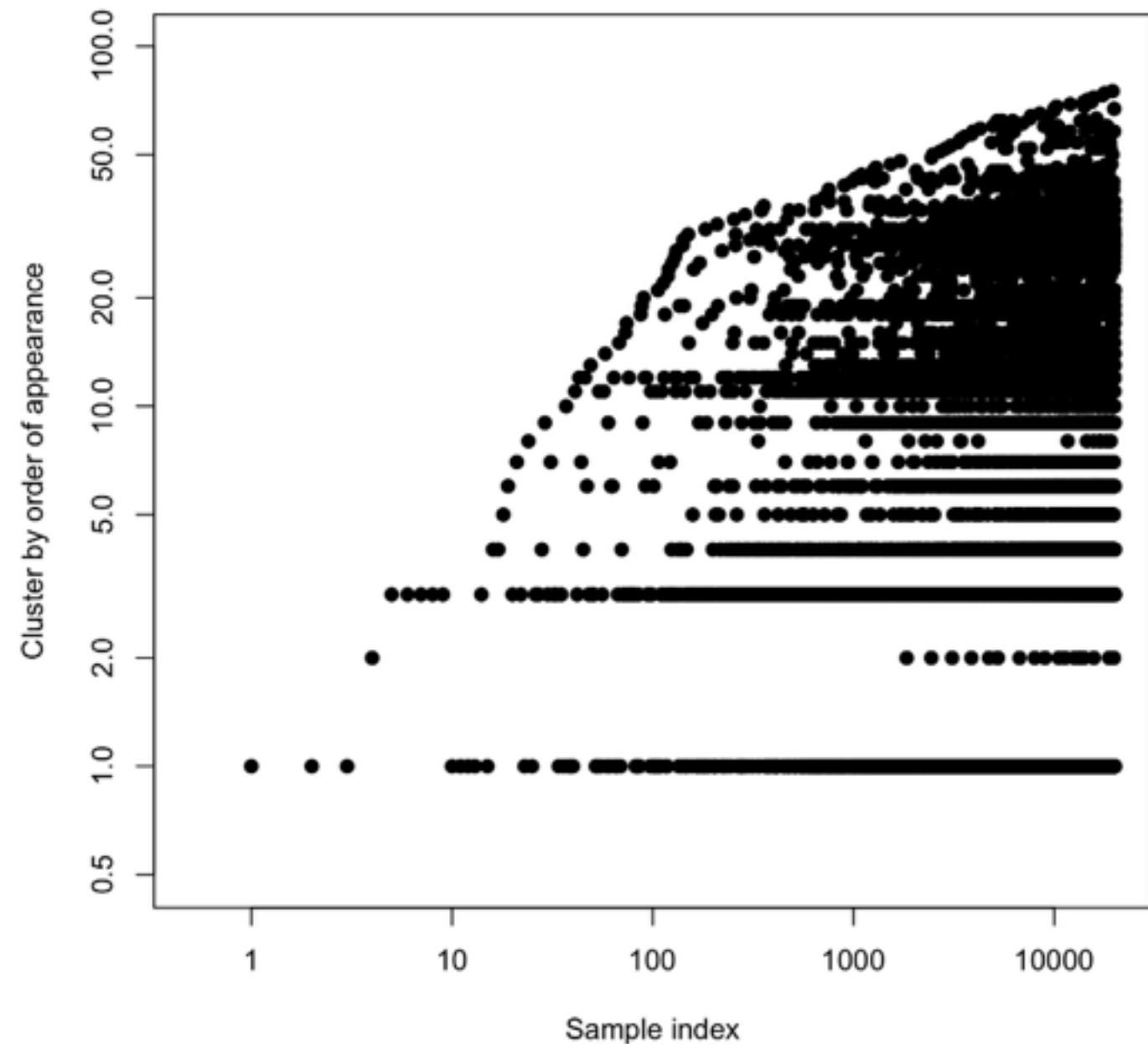
# Power laws

- $K_N := \#$  clusters occupied by  $N$  data points
- CRP:  $K_N \sim \alpha \log N$  w.p. 1
  - vs. Heaps' law, Herdan's law, etc



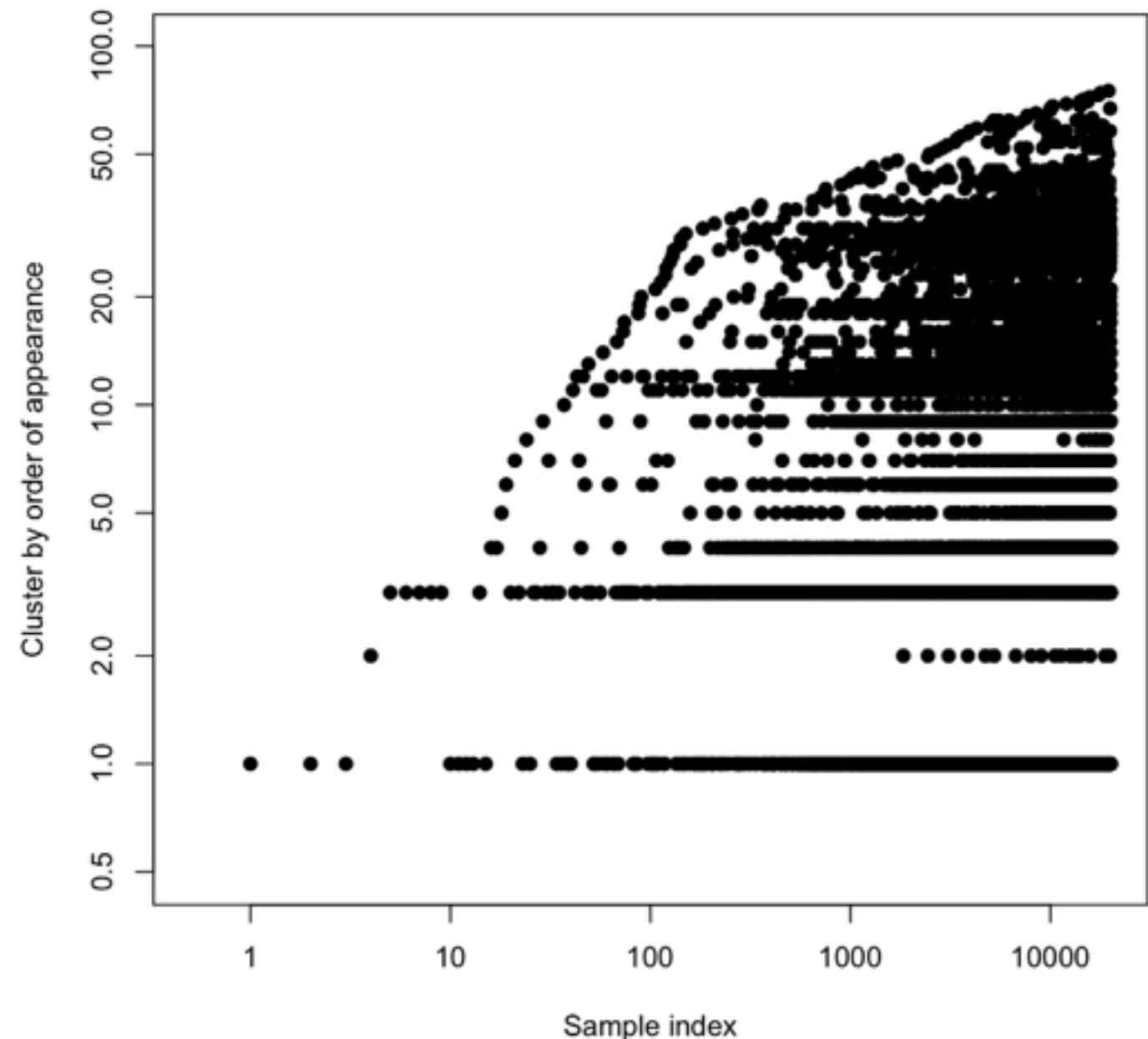
# Power laws

- $K_N := \#$  clusters occupied by  $N$  data points
- CRP:  $K_N \sim \alpha \log N$  w.p. 1
  - vs. Heaps' law, Herdan's law, etc



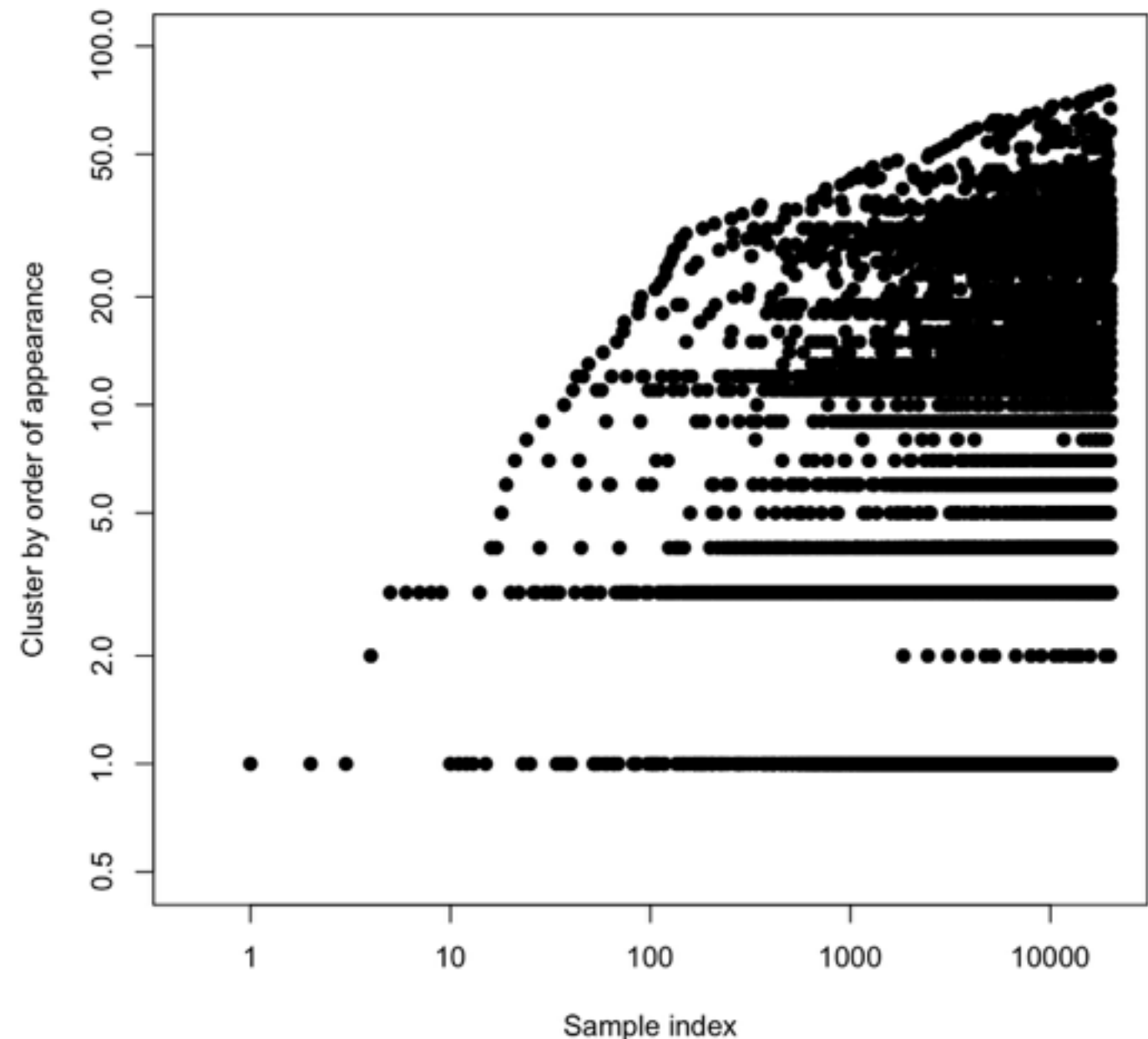
# Power laws

- $K_N := \#$  clusters occupied by  $N$  data points
- CRP:  $K_N \sim \alpha \log N$  w.p. 1
  - vs. Heaps' law, Herdan's law, etc
- Pitman-Yor process:



# Power laws

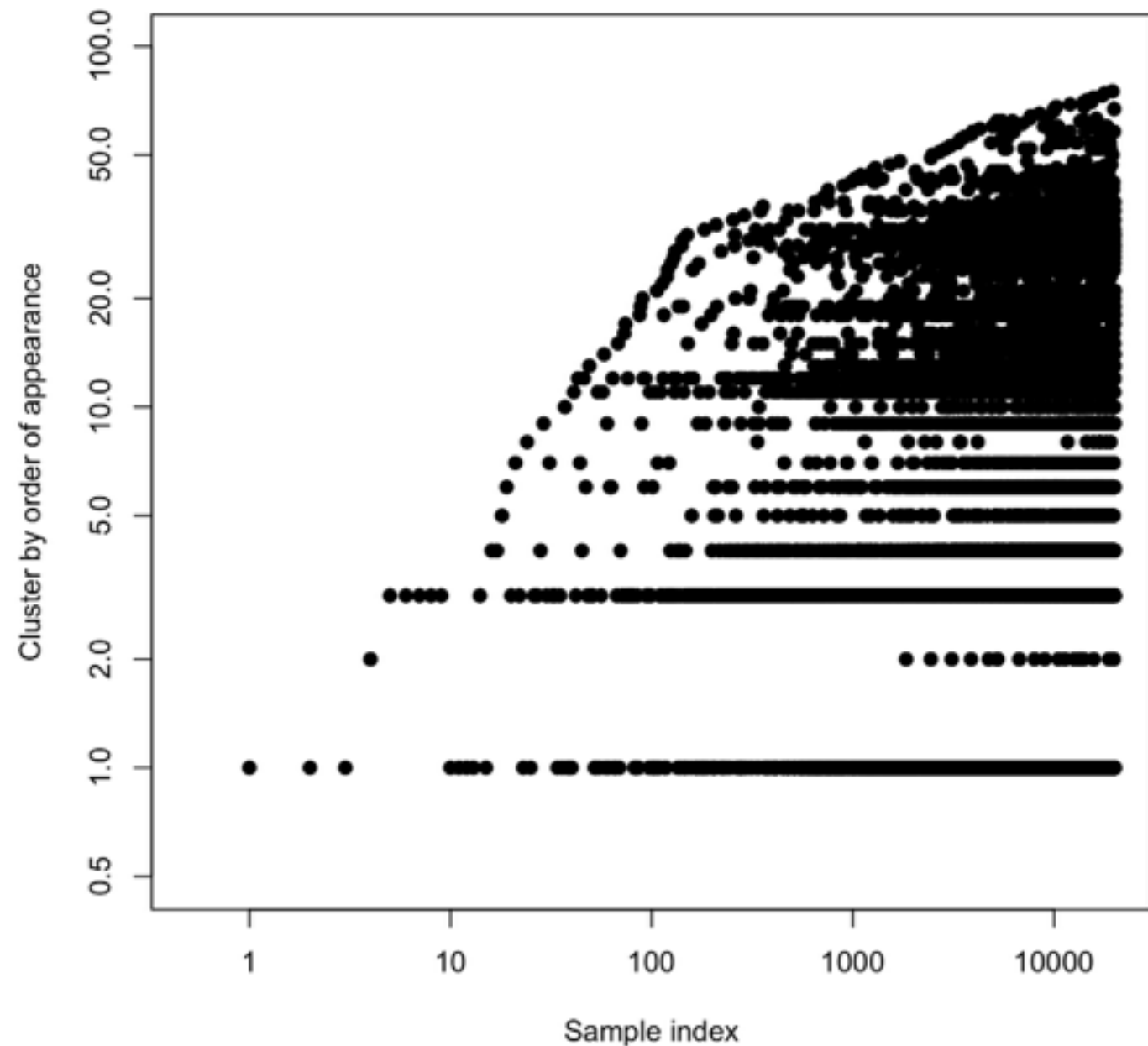
- $K_N := \#$  clusters occupied by  $N$  data points
- CRP:  $K_N \sim \alpha \log N$  w.p. 1
  - vs. Heaps' law, Herdan's law, etc
- Pitman-Yor process:





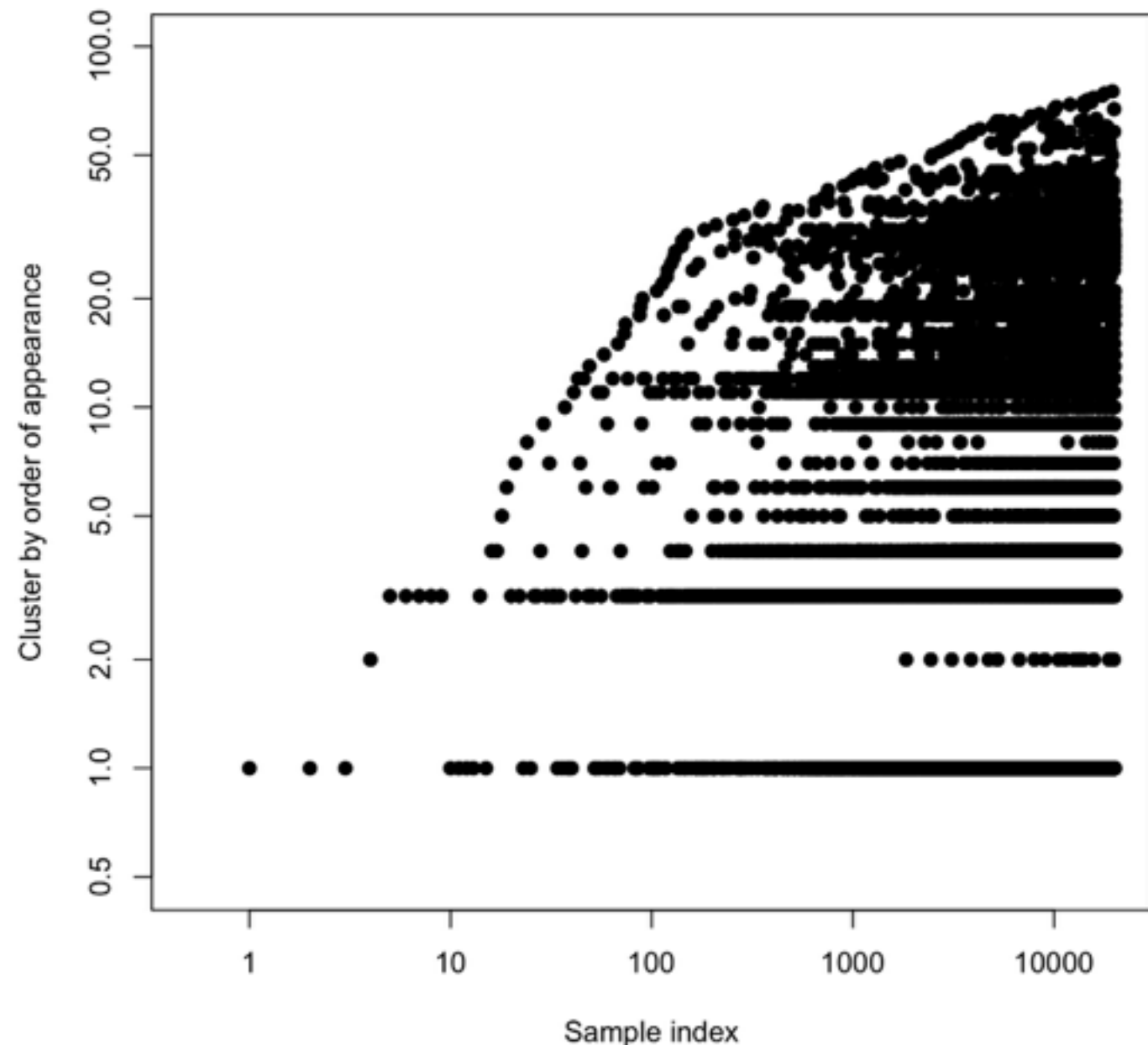
# Power laws

- $K_N := \#$  clusters occupied by  $N$  data points
- CRP:  $K_N \sim \alpha \log N$  w.p. 1
  - vs. Heaps' law, Herdan's law, etc
- Pitman-Yor process:  
 $K_N \sim S_\alpha N^\sigma$  w.p. 1



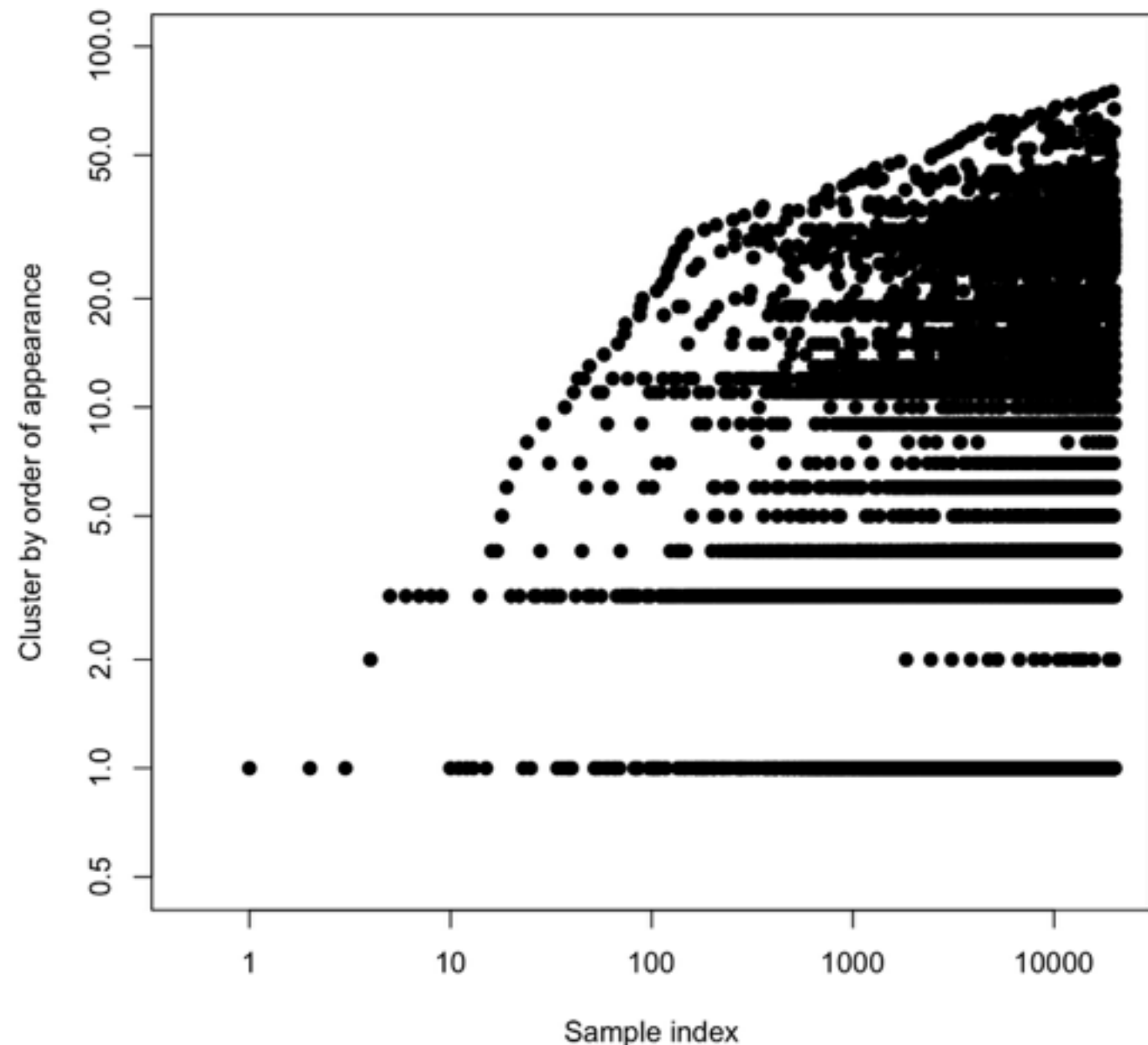
# Power laws

- $K_N := \#$  clusters occupied by  $N$  data points
- CRP:  $K_N \sim \alpha \log N$  w.p. 1
  - vs. Heaps' law, Herdan's law, etc
- Pitman-Yor process:  
 $K_N \sim S_\alpha N^\sigma$  w.p. 1



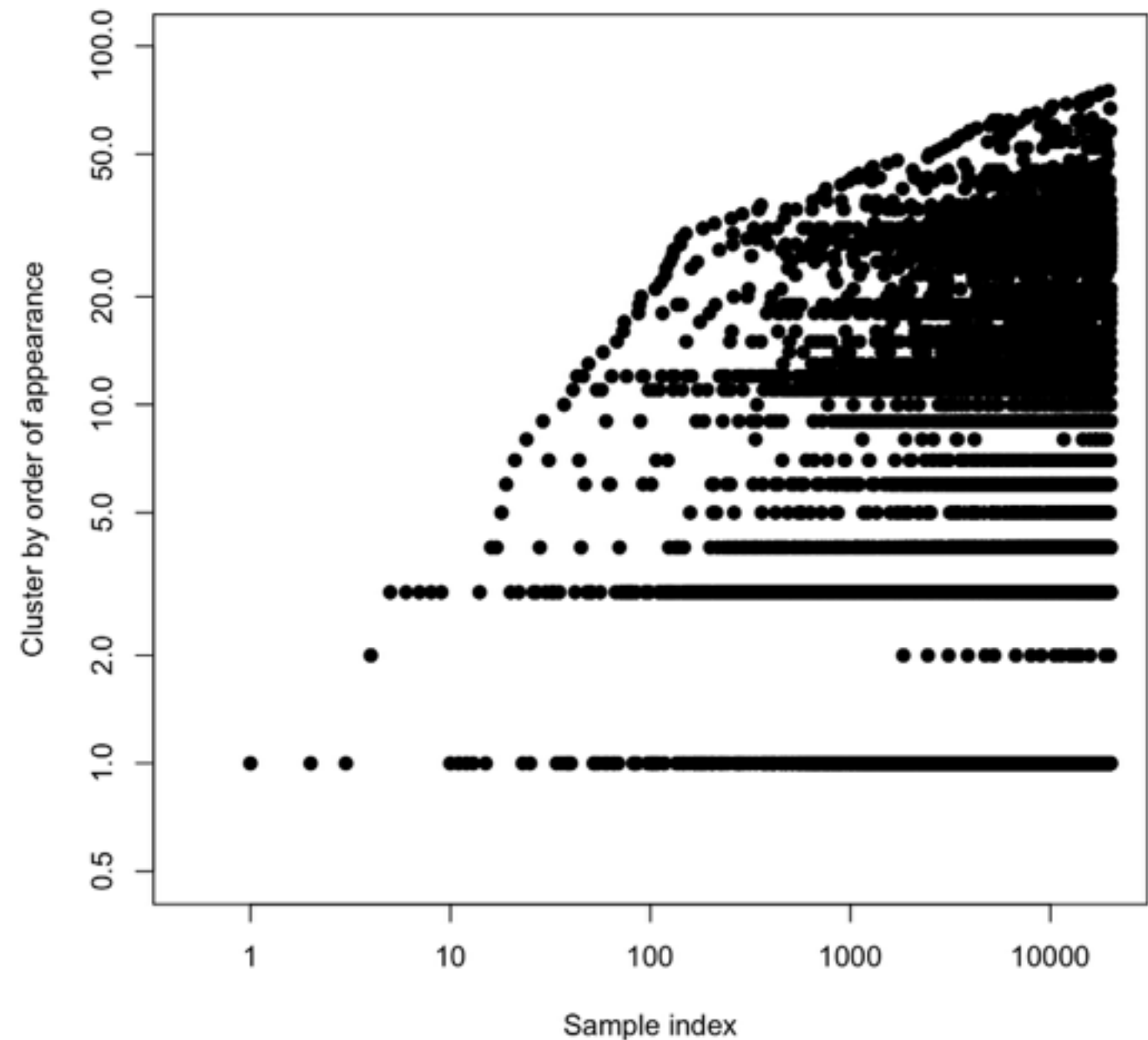
# Power laws

- $K_N := \#$  clusters occupied by  $N$  data points
- CRP:  $K_N \sim \alpha \log N$  w.p. 1
  - vs. Heaps' law, Herdan's law, etc
- Pitman-Yor process:
$$K_N \sim S_\alpha N^\sigma \text{ w.p. } 1$$
  - related to Zipf's law (ranked frequencies)



# Power laws

- $K_N := \#$  clusters occupied by  $N$  data points
- CRP:  $K_N \sim \alpha \log N$  w.p. 1
  - vs. Heaps' law, Herdan's law, etc
- Pitman-Yor process:
$$K_N \sim S_\alpha N^\sigma \text{ w.p. } 1$$
  - related to Zipf's law (ranked frequencies)
- Not just clusters



# Hierarchies



# Hierarchies

- Hierarchical  
Dirichlet process

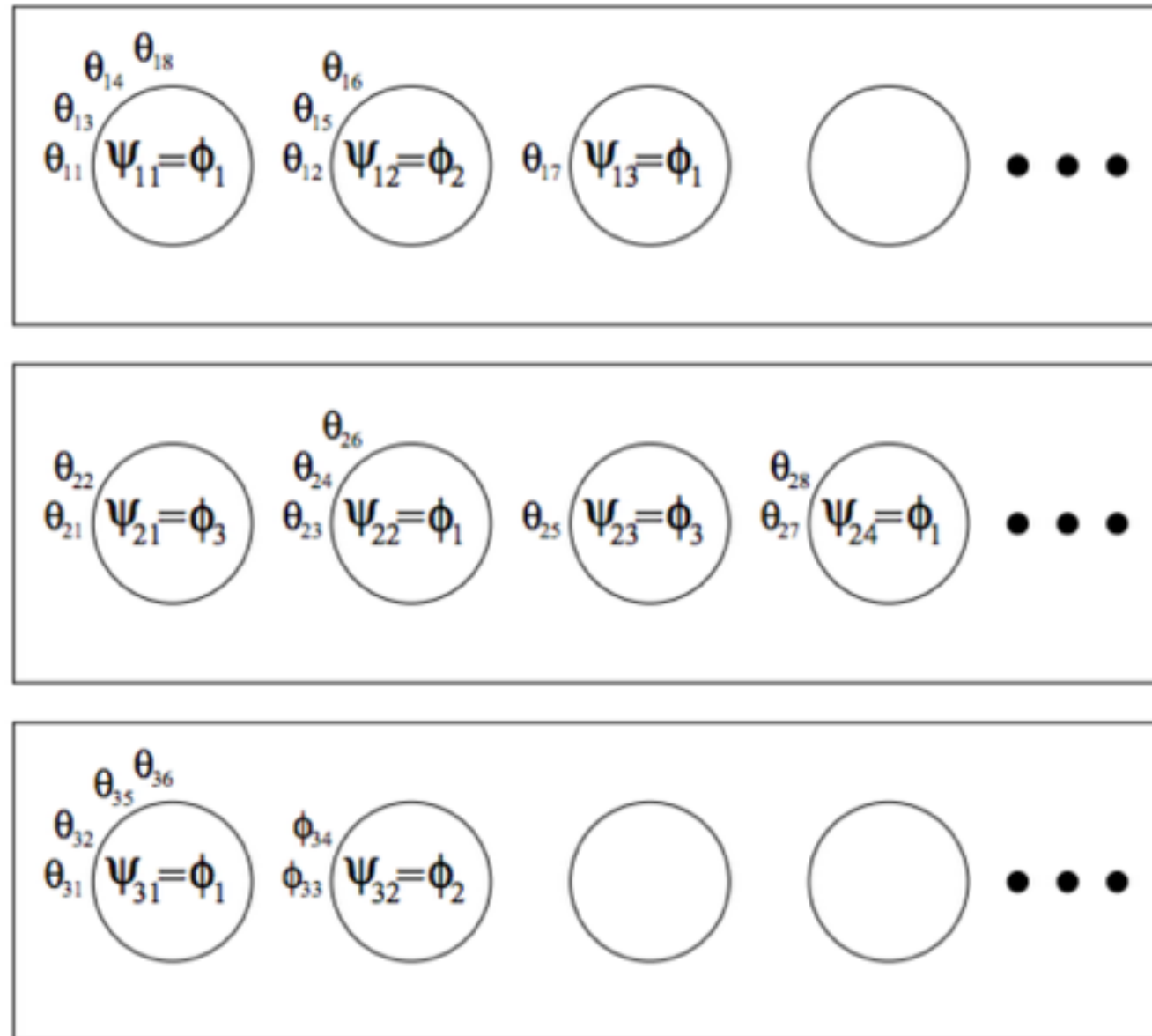
# Hierarchies

- Hierarchical Dirichlet process

# Hierarchies

- Hierarchical Dirichlet process
- Chinese restaurant franchise

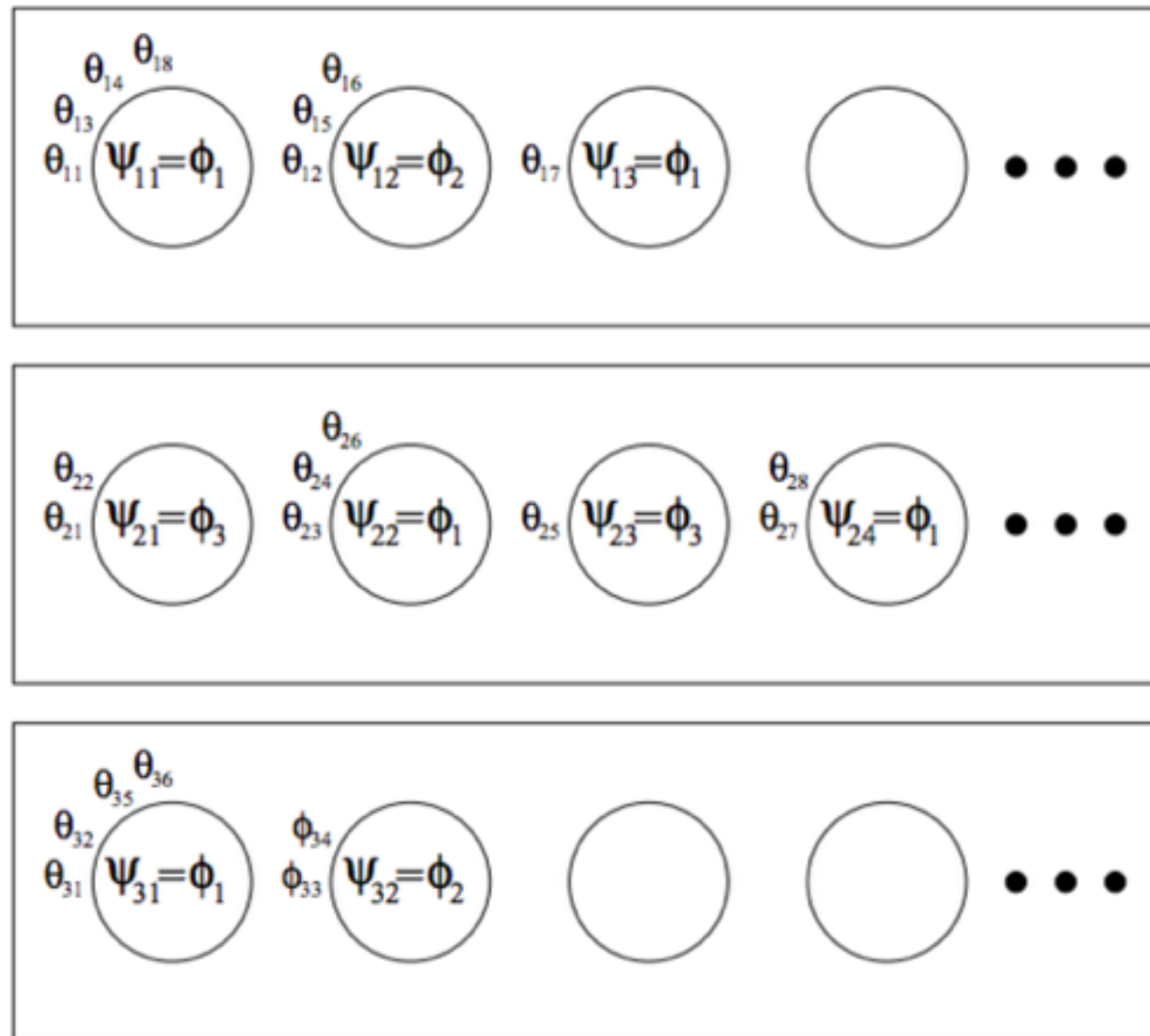
# Hierarchies



- Hierarchical Dirichlet process
- Chinese restaurant franchise

[Teh et al 2006]

# Hierarchies

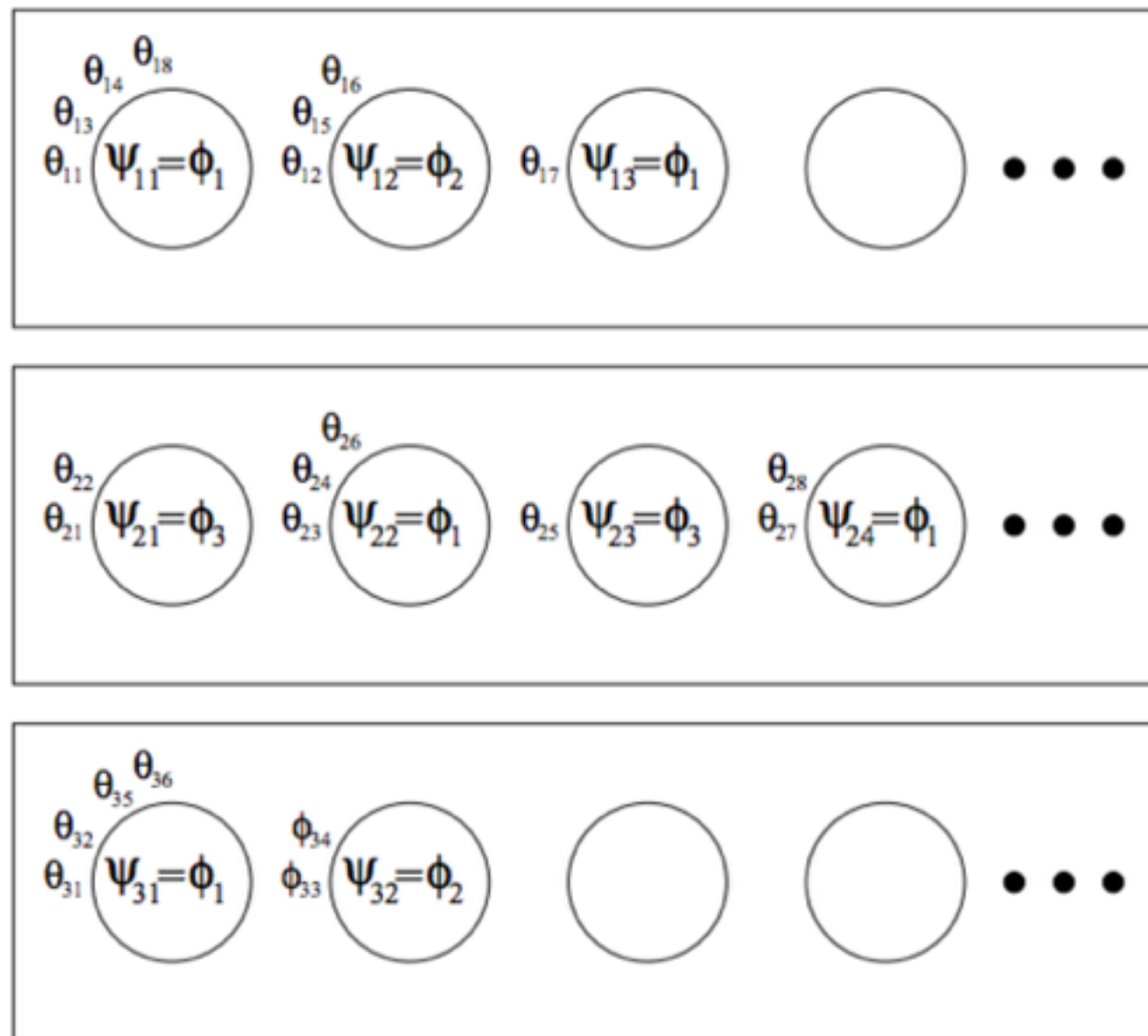


- Hierarchical Dirichlet process
- Chinese restaurant franchise
- Hierarchical beta process

[Teh et al 2006]



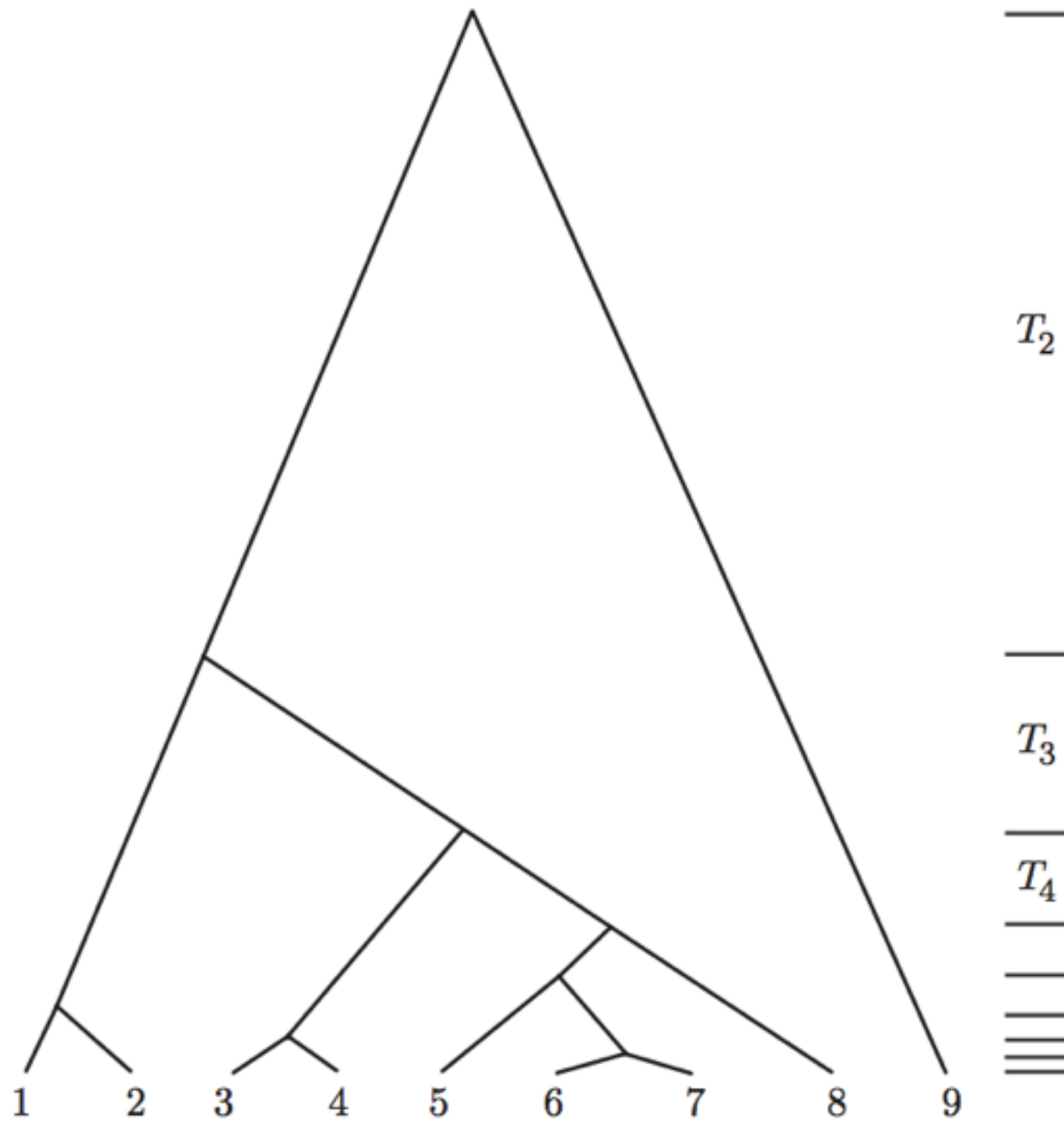
# Hierarchies



- Hierarchical Dirichlet process
- Chinese restaurant franchise
- Hierarchical beta process

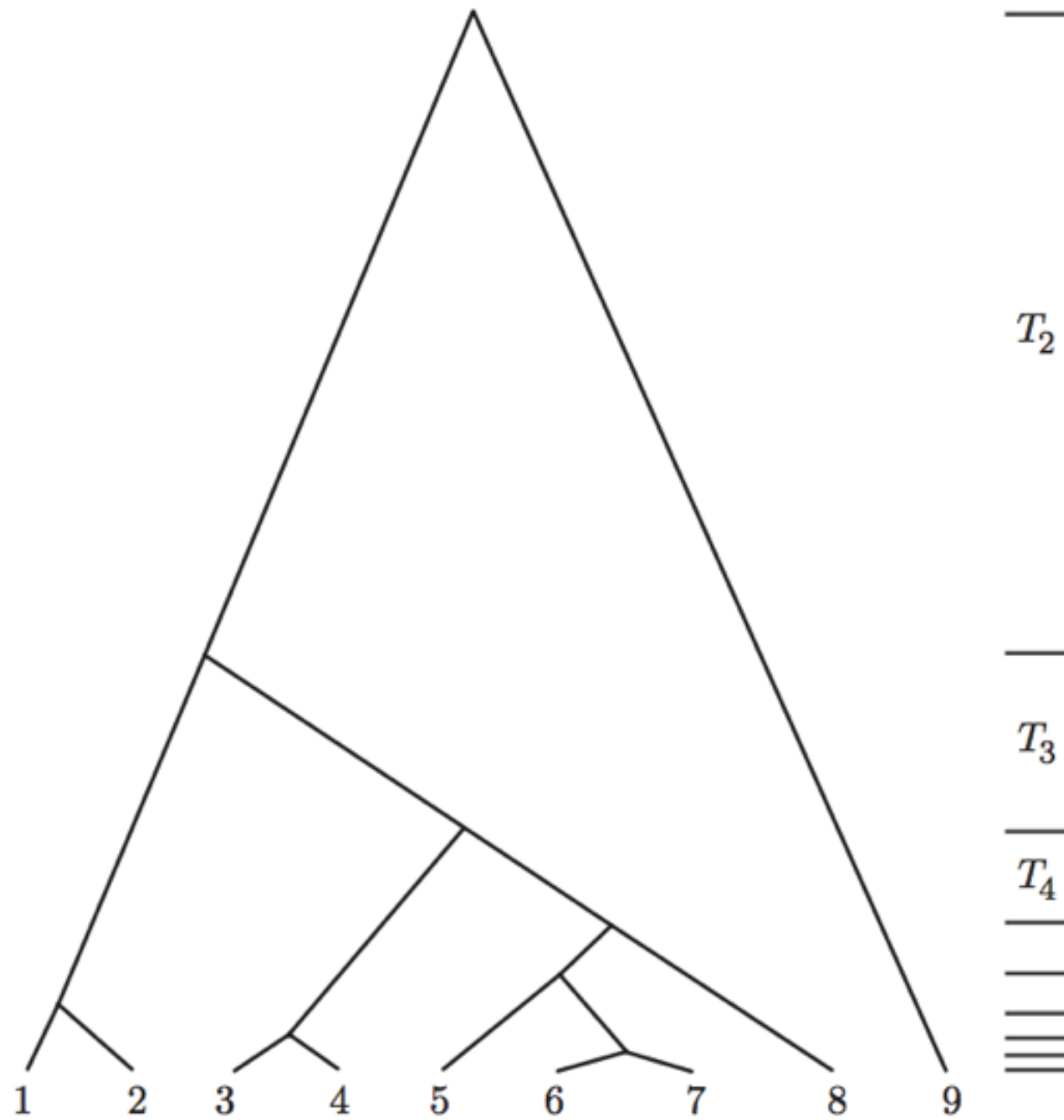
[Teh et al 2006]

# Genealogy, trees, beyond trees



[Wakeley 2008]

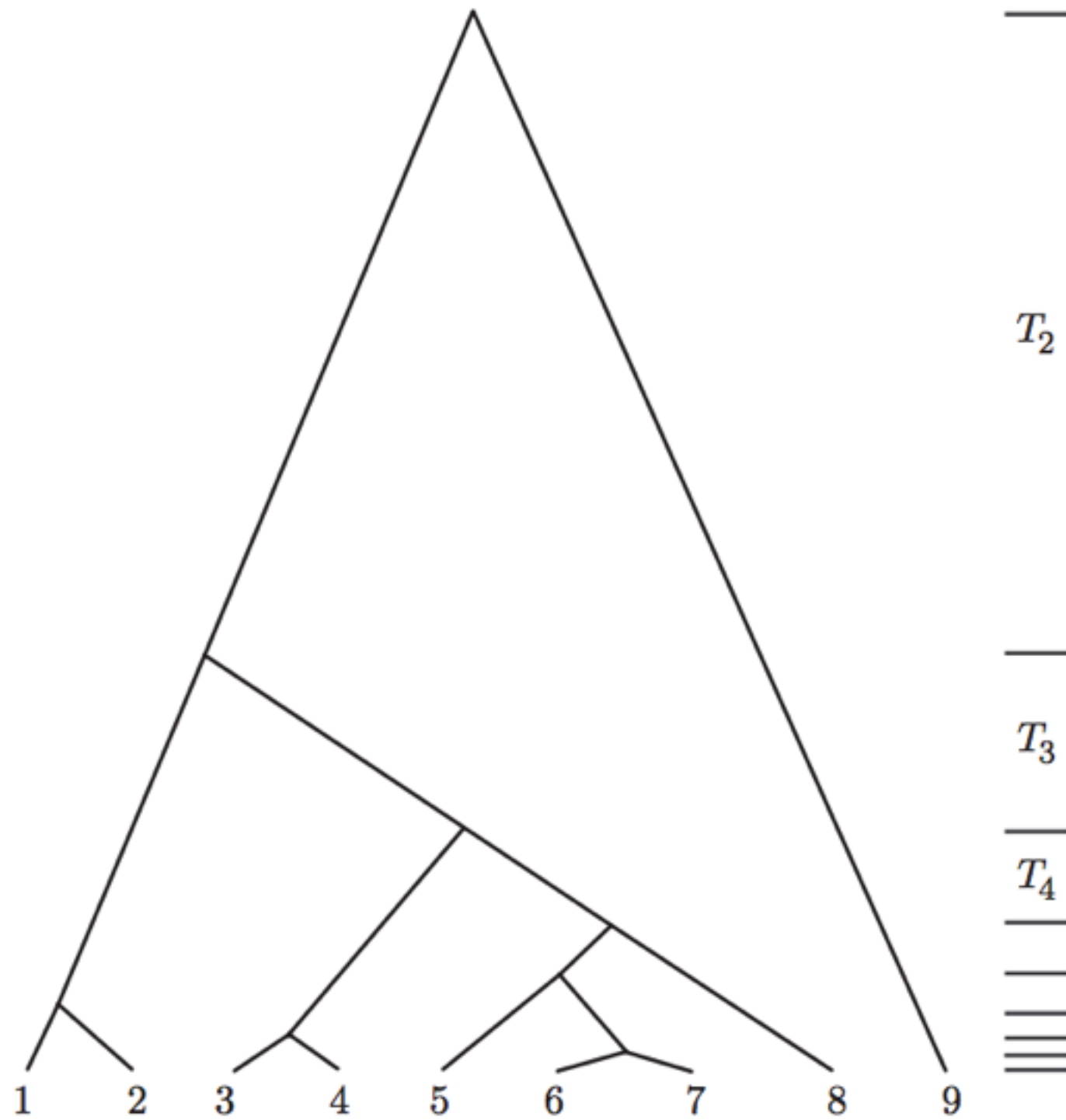
# Genealogy, trees, beyond trees



- Kingman coalescent

[Wakeley 2008]

# Genealogy, trees, beyond trees

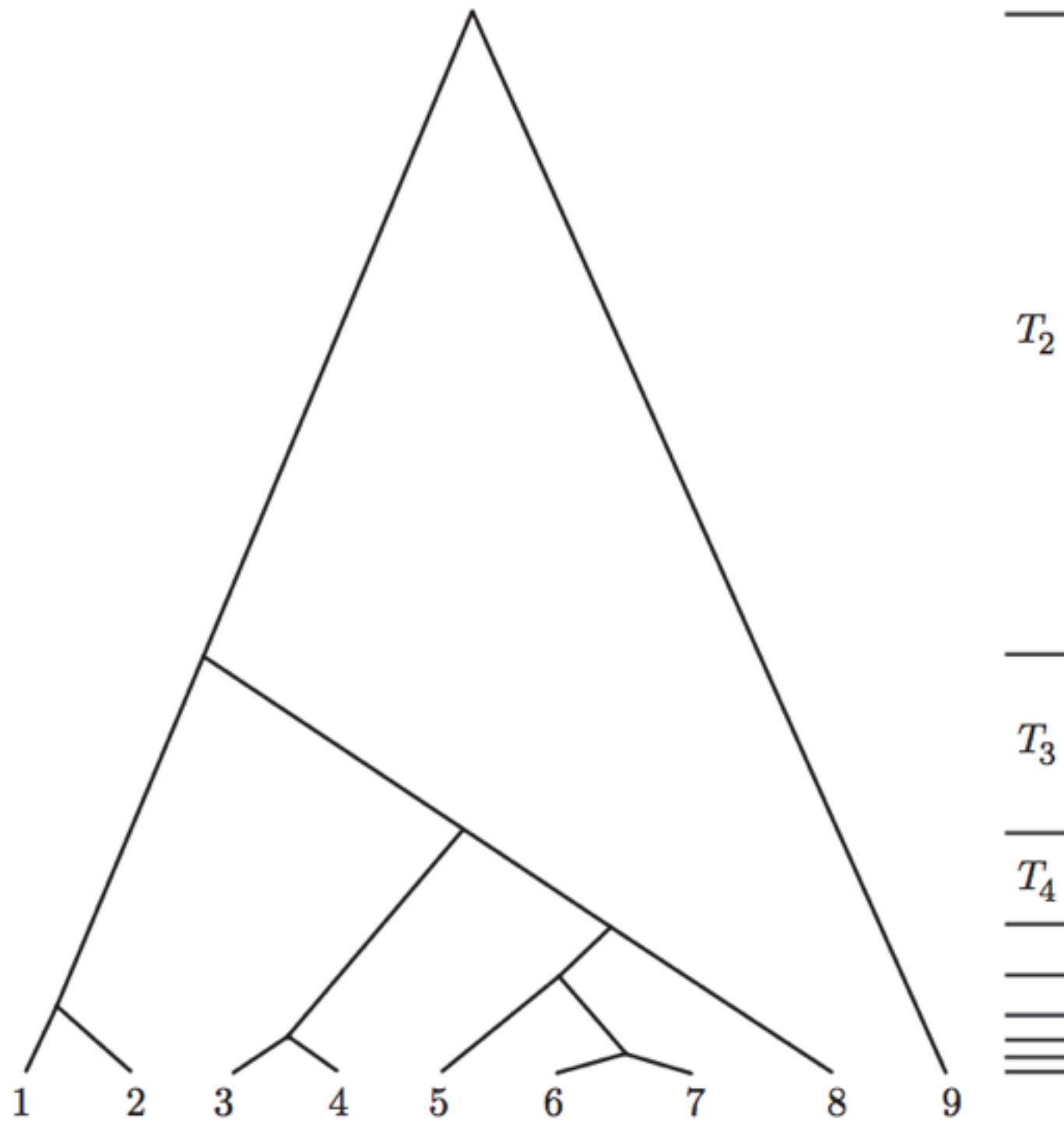


- Kingman coalescent

[Wakeley 2008]

[Kingman 1982]

# Genealogy, trees, beyond trees

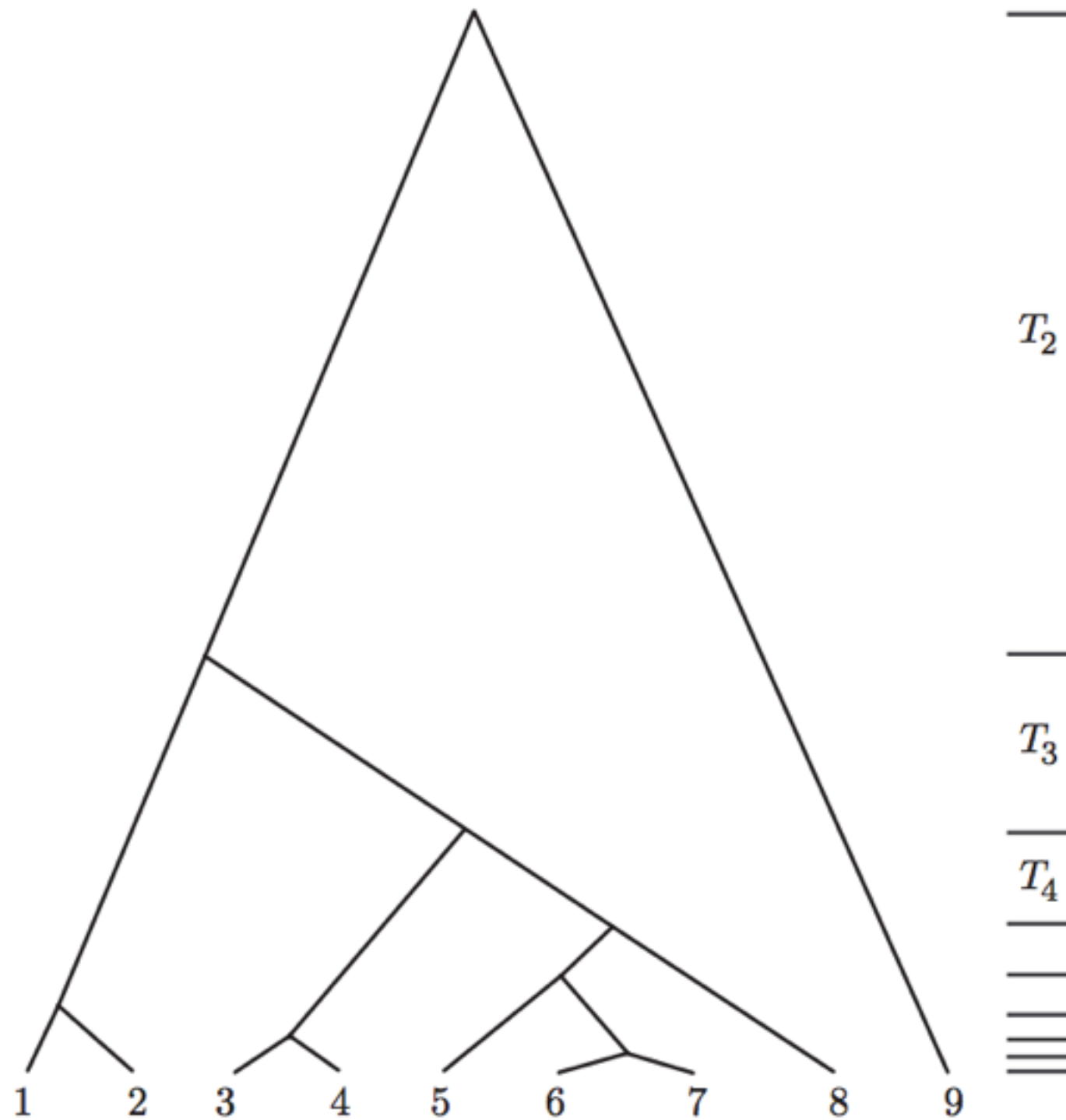


[Wakeley 2008]

[Kingman 1982]

- Kingman coalescent
- Fragmentation
- Coagulation

# Genealogy, trees, beyond trees



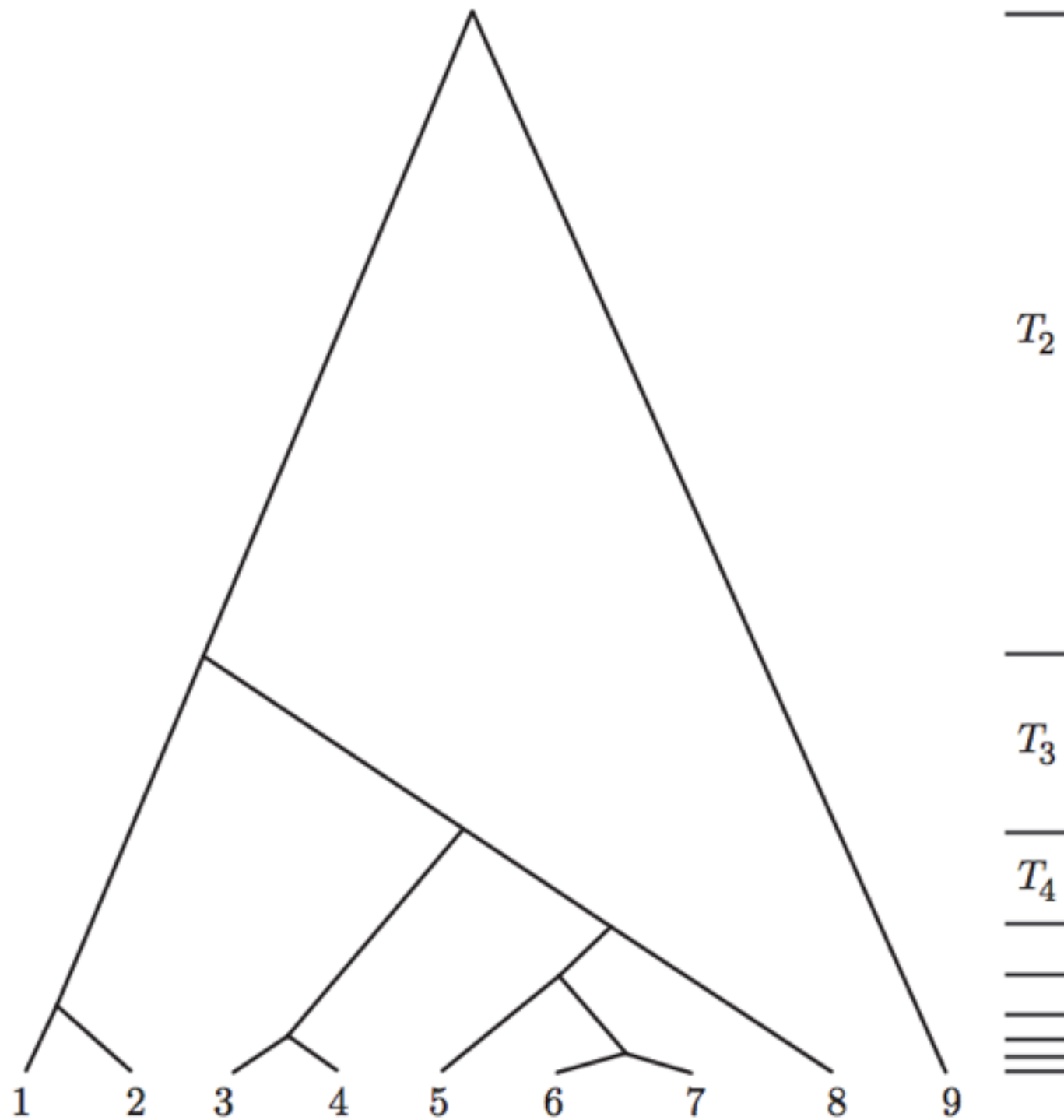
- Kingman coalescent
- Fragmentation
- Coagulation

[Wakeley 2008]

[Kingman 1982, Bertoin 2006, Teh et al 2011]



# Genealogy, trees, beyond trees

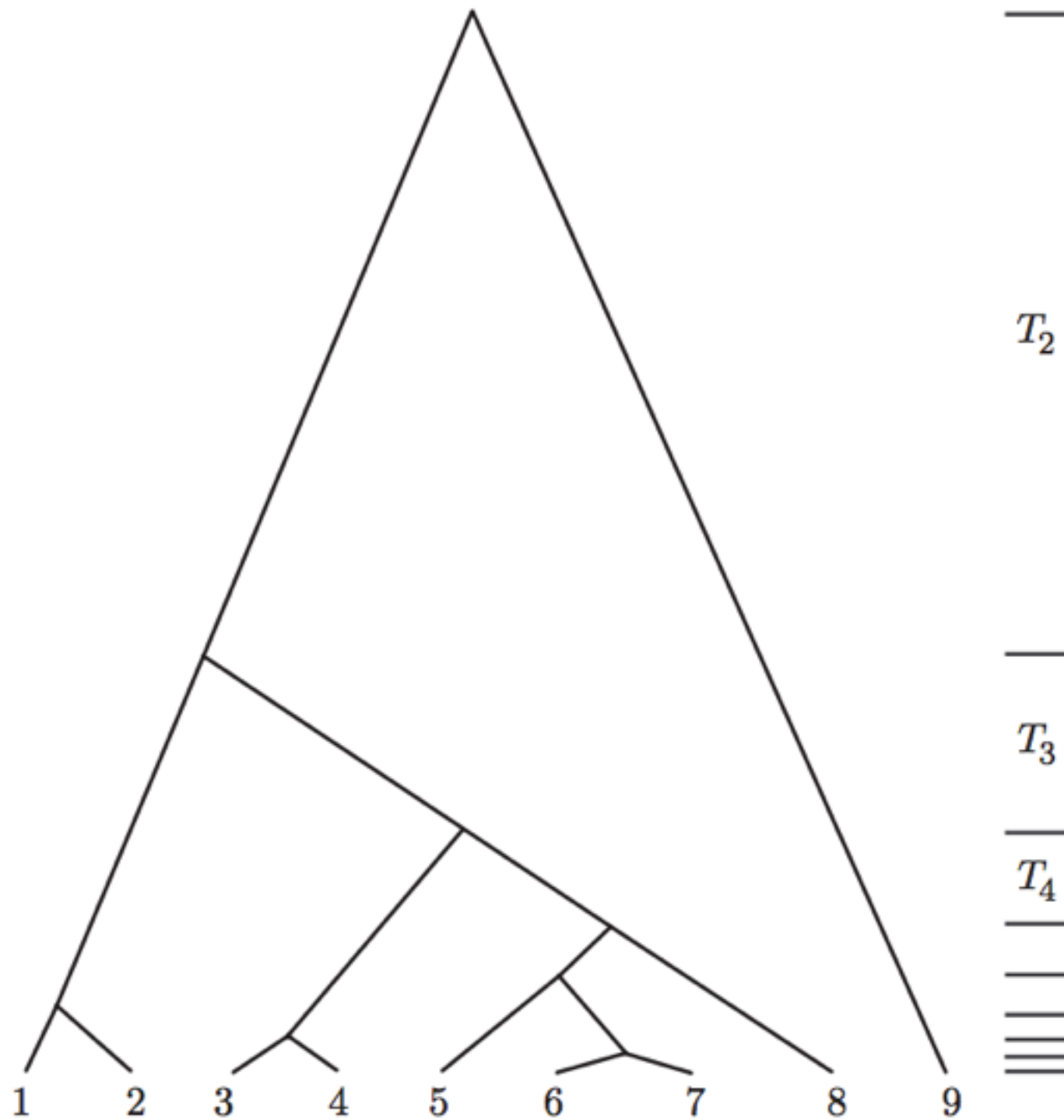


[Wakeley 2008]

[Kingman 1982, Bertoin 2006, Teh et al 2011]

- Kingman coalescent
- Fragmentation
- Coagulation
- Dirichlet diffusion tree

# Genealogy, trees, beyond trees



[Wakeley 2008]

[Kingman 1982, Bertoin 2006, Teh et al 2011, Neal 2003]

- Kingman coalescent
- Fragmentation
- Coagulation
- Dirichlet diffusion tree

# Conjugacy & Poisson point processes

# Conjugacy & Poisson point processes

- Beta process, Bernoulli process (Indian buffet)

# Conjugacy & Poisson point processes

- Beta process, Bernoulli process (Indian buffet)
- Gamma process, Poisson likelihood process (DP, CRP)

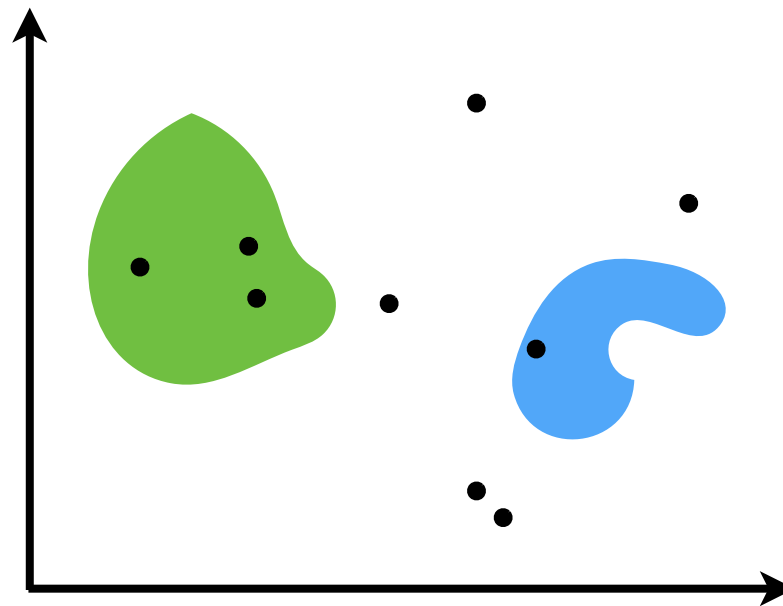
# Conjugacy & Poisson point processes

- Beta process, Bernoulli process (Indian buffet)
- Gamma process, Poisson likelihood process (DP, CRP)
- Beta process, negative binomial process



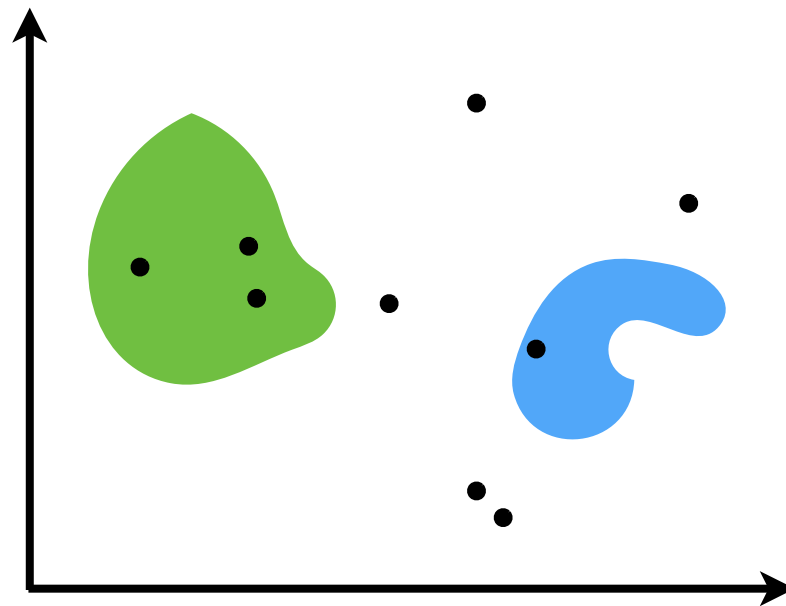
# Conjugacy & Poisson point processes

- Beta process, Bernoulli process (Indian buffet)
- Gamma process, Poisson likelihood process (DP, CRP)
- Beta process, negative binomial process



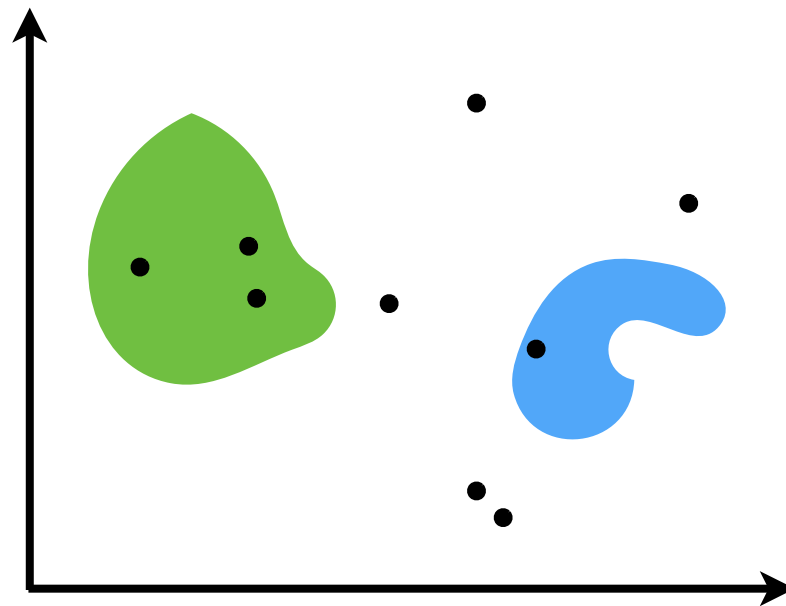
# Conjugacy & Poisson point processes

- Beta process, Bernoulli process (Indian buffet)
- Gamma process, Poisson likelihood process (DP, CRP)
- Beta process, negative binomial process



# Conjugacy & Poisson point processes

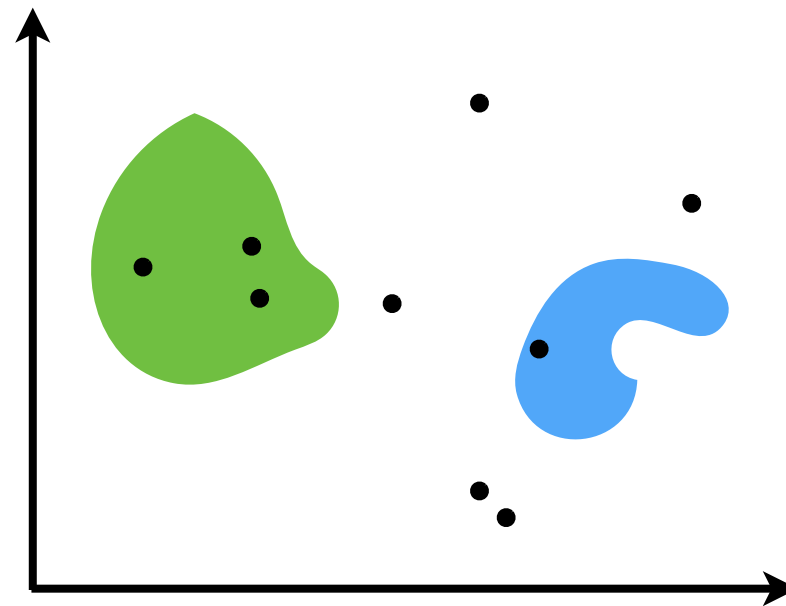
- Beta process, Bernoulli process (Indian buffet)
- Gamma process, Poisson likelihood process (DP, CRP)
- Beta process, negative binomial process



- Posteriors, conjugacy, and exponential families for completely random measures

# Conjugacy & Poisson point processes

- Beta process, Bernoulli process (Indian buffet)
- Gamma process, Poisson likelihood process (DP, CRP)
- Beta process, negative binomial process



- Posteriors, conjugacy, and exponential families for completely random measures

# De Finetti mixing measures

# De Finetti mixing measures

- Clustering: Kingman paintbox



# De Finetti mixing measures

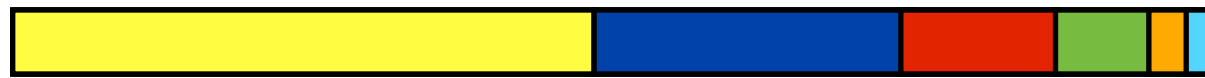
- Clustering: Kingman paintbox



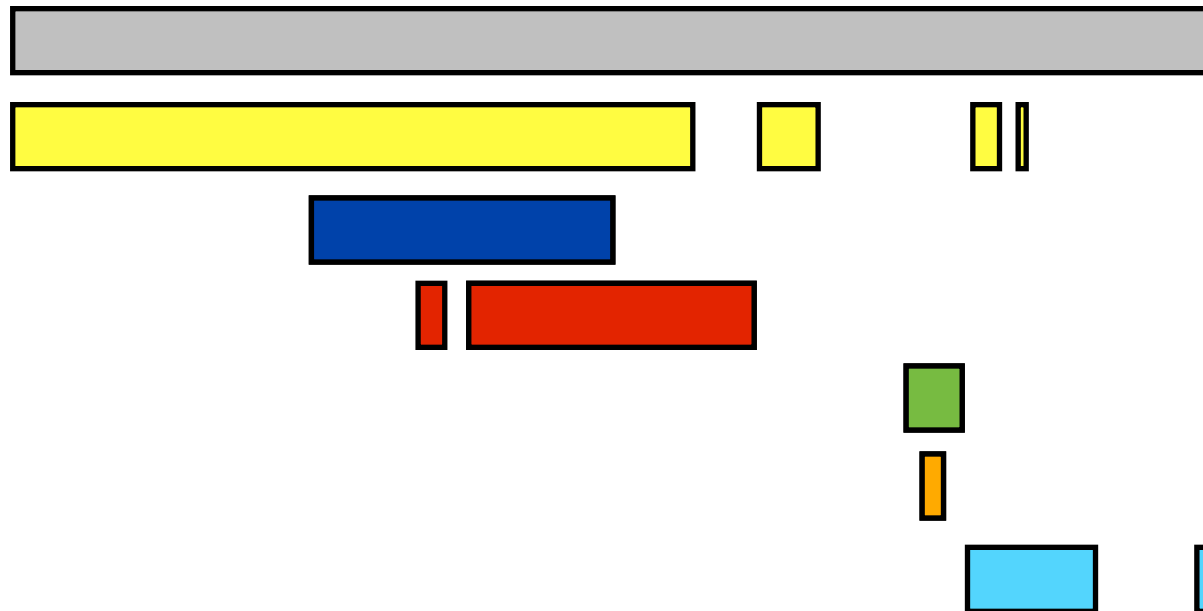


# De Finetti mixing measures

- Clustering: Kingman paintbox



- Feature allocation: Feature paintbox

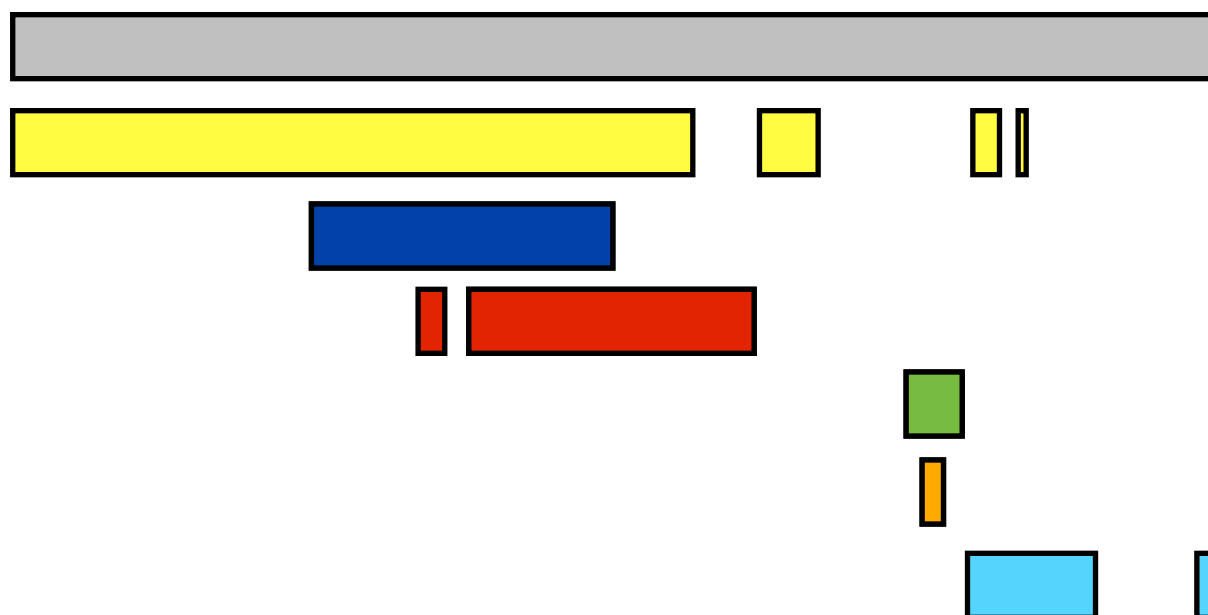


# De Finetti mixing measures

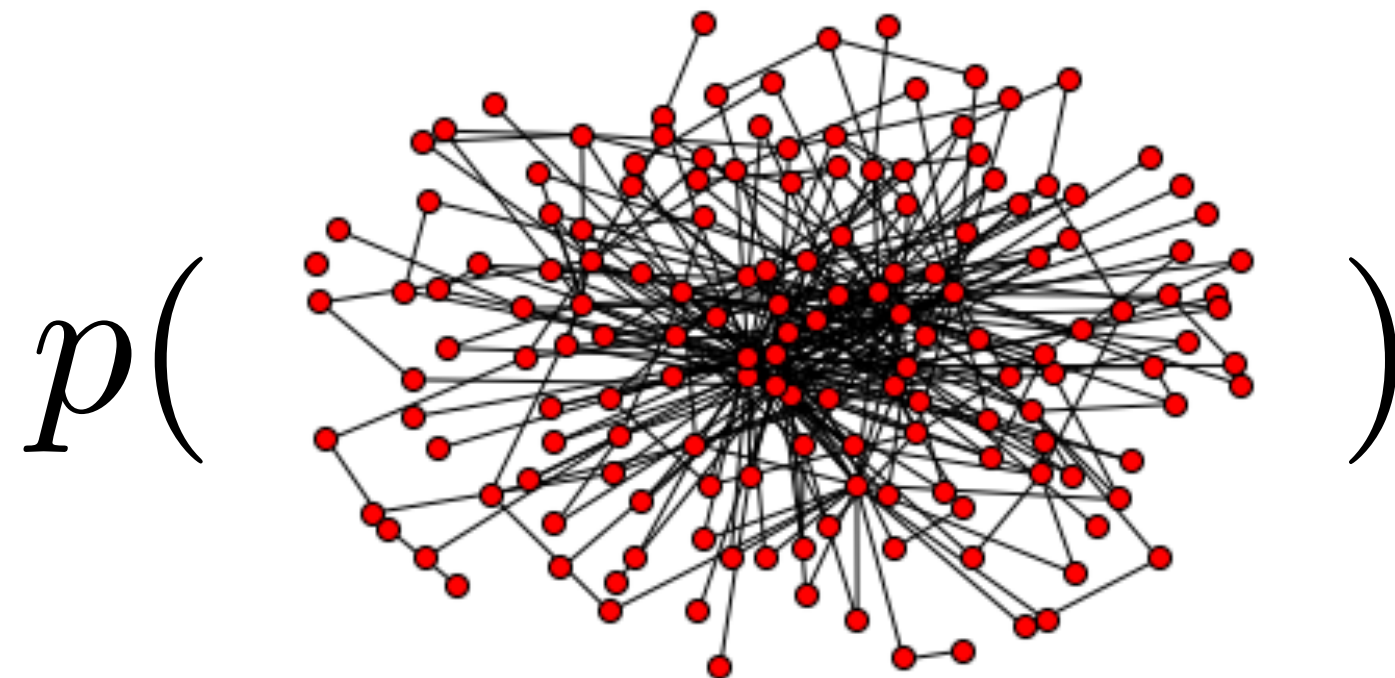
- Clustering: Kingman paintbox



- Feature allocation: Feature paintbox



# Probabilistic models for graphs



E.g. online social networks,  
biological networks,  
communication networks,  
transportation networks

- Rich relationships, coherent uncertainties, prior info
- Stochastic block model, mixed membership stochastic block model, infinite relational model, and many more
- Assume: Adding more data doesn't change distribution of earlier data (*projectivity*)
- **Problem:** model misspecification, dense graphs

# Edge exchangeability



**Thm. A wide range of edge-exchangeable graph sequences are sparse**

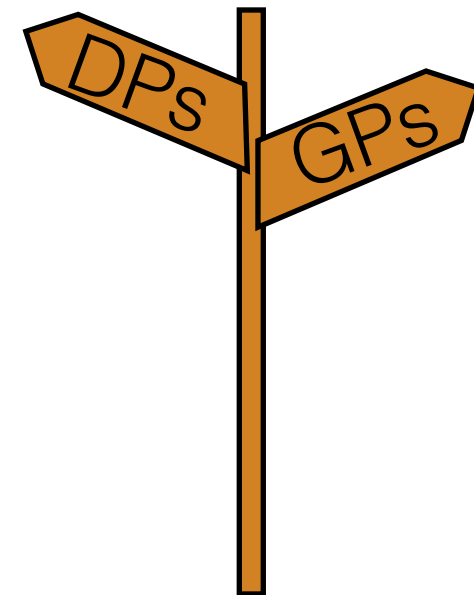
$G_1$        $G_2$        $G_3$        $G_4$

**Thm. A paintbox-style characterization for edge-exchangeable graph sequences**

$$p\left( \begin{array}{c} \text{1} \quad \text{2} \\ \text{3} \end{array} \right) = p\left( \begin{array}{c} \text{2} \quad \text{4} \\ \text{1} \end{array} \right)$$

# Roadmap

- Bayes Foundations
- Unsupervised Learning
  - Example problem: clustering
  - Example BNP model: Dirichlet process (DP)
  - Chinese restaurant process
- Supervised Learning
  - Example problem: regression
  - Example BNP model: Gaussian process (GP)
- Venture further into the wild world of Nonparametric Bayes
- Big questions
  - Why BNP?
  - What does an infinite/growing number of parameters really mean (in BNP)?
  - Why is BNP challenging but practical?



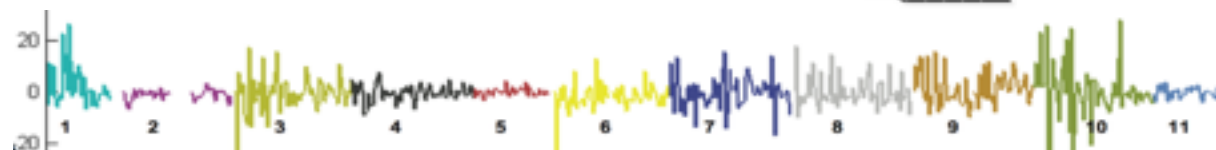


# Applications

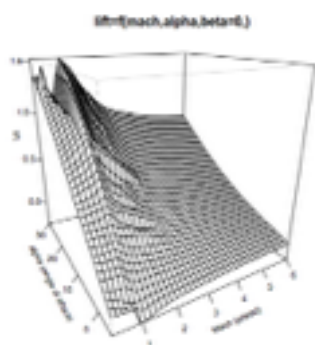


[wikipedia.org]

[Saria  
et al  
2010]



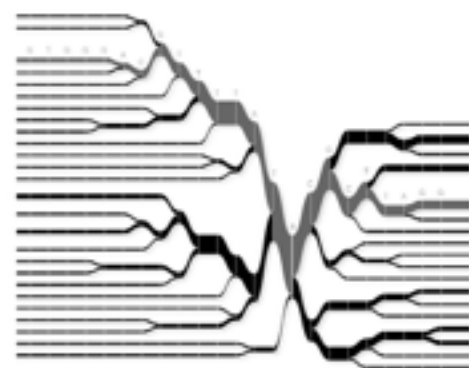
[US CDC PHIL;  
Futoma, Hariharan,  
Heller 2017]



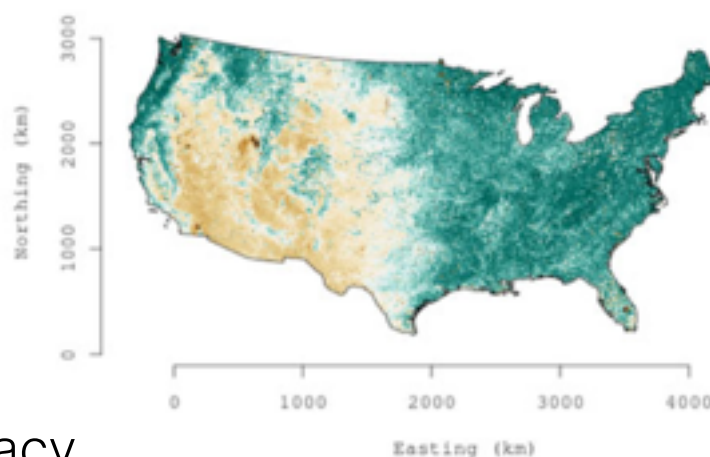
[Gramacy,  
Lee 2009]



[Ed Bowlby, NOAA]



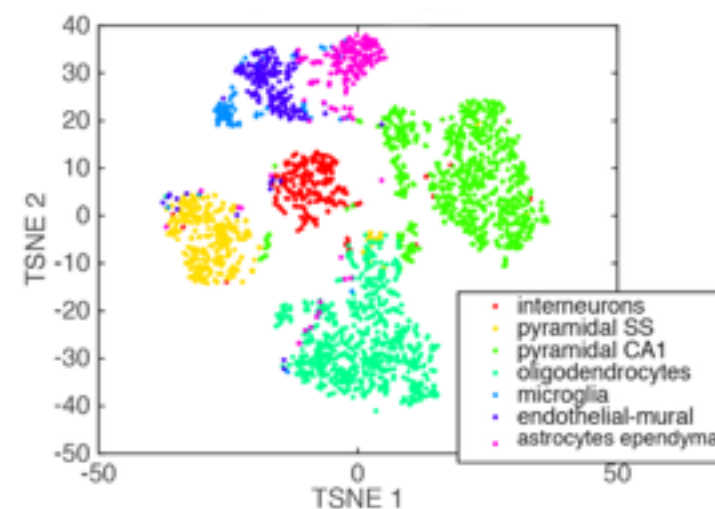
[Ewens  
1972;  
Hartl,  
Clark  
2003]



[Datta,  
Banerjee,  
Finley,  
Gelfand  
2016]

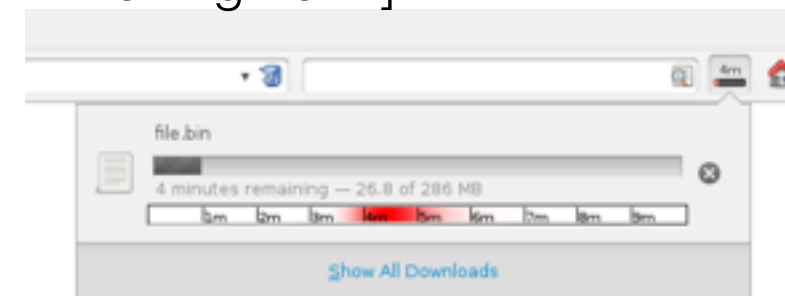


[Fox et al 2014]

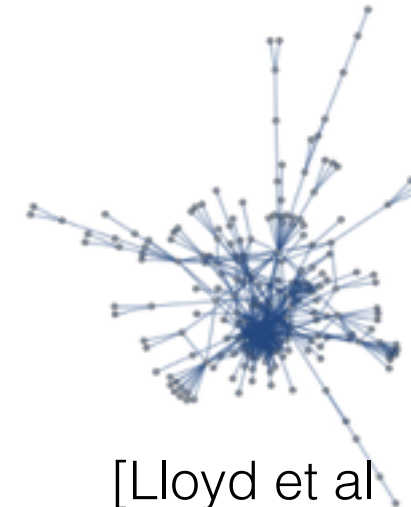


[Prabhakaran, Azizi, Carr,  
Pe'er 2016]

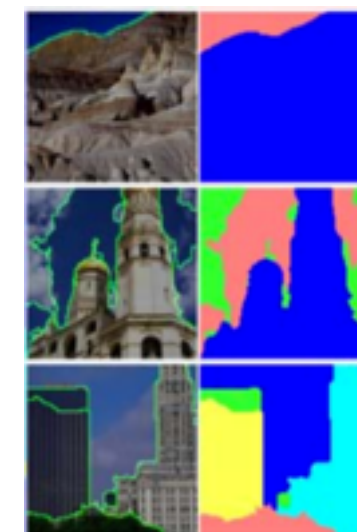
[Kiefel,  
Schuler,  
Hennig 2014]



[Deisenroth, Fox, Rasmussen 2015]



[Lloyd et al  
2012; Miller  
et al 2010]



[Sudderth,  
Jordan 2009]



[Chati,  
Balakrishnan  
2017]

