

# Statistical Machine Learning

## Lecture 01 Non-Parametric Bayesian Introduction

Spring 2021  
Sharif University of Technology

# Bayesian Nonparametrics

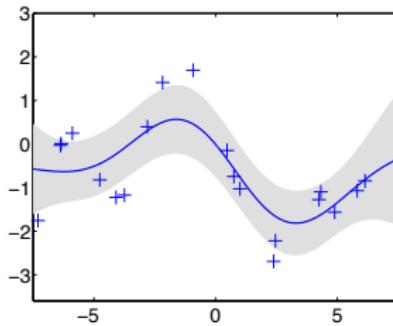
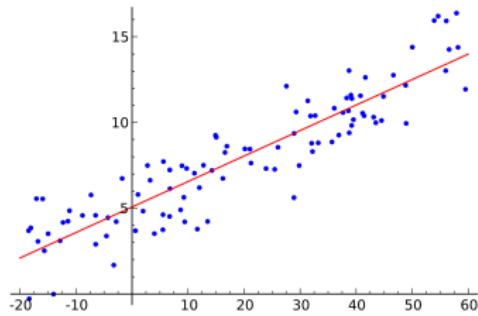
## Part I

Peter Orbanz

# PARAMETERS AND PATTERNS

## Parameters

$$P(X|\theta) = \text{Probability}[\text{data}|\text{pattern}]$$



## Inference idea

$$\text{data} = \text{underlying pattern} + \text{independent noise}$$

# TERMINOLOGY

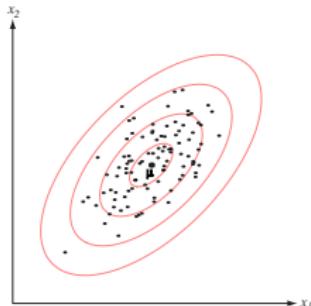
## Parametric model

- ▶ Number of parameters fixed (or constantly bounded) w.r.t. sample size

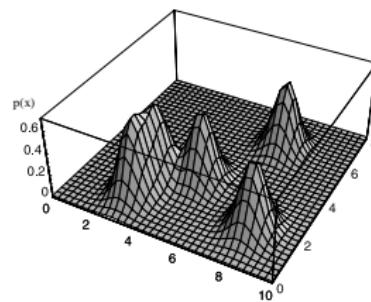
## Nonparametric model

- ▶ Number of parameters grows with sample size
- ▶  $\infty$ -dimensional parameter space

## Example: Density estimation



Parametric



Nonparametric

# NONPARAMETRIC BAYESIAN MODEL

## Definition

A nonparametric Bayesian model is a Bayesian model on an  $\infty$ -dimensional parameter space.

## Interpretation

Parameter space  $\mathcal{T}$  = set of possible patterns, for example:

Problem	$\mathcal{T}$
Density estimation	Probability distributions
Regression	Smooth functions
Clustering	Partitions

Solution to Bayesian problem = posterior distribution on patterns

# EXCHANGEABILITY

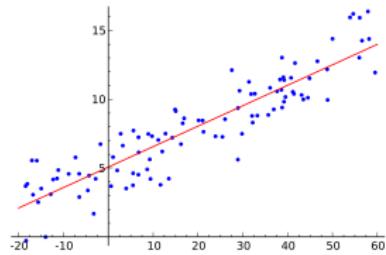
Can we justify our assumptions?

Recall:

$$\text{data} = \text{pattern} + \text{noise}$$

In Bayes' theorem:

$$Q(d\theta|x_1, \dots, x_n) = \frac{\prod_{j=1}^n p(x_j|\theta)}{p(x_1, \dots, x_n)} Q(d\theta)$$



## Definition

$X_1, X_2, \dots$  are *exchangeable* if  $P(X_1, X_2, \dots)$  is invariant under any permutation  $\sigma$ :

$$P(X_1 = x_1, X_2 = x_2, \dots) = P(X_1 = x_{\sigma(1)}, X_2 = x_{\sigma(2)}, \dots)$$

In words:

Order of observations does not matter.

# EXCHANGEABILITY AND CONDITIONAL INDEPENDENCE

## De Finetti's Theorem

$$P(X_1 = x_1, X_2 = x_2, \dots) = \int_{\mathbf{M}(\mathcal{X})} \left( \prod_{j=1}^{\infty} \theta(X_j = x_j) \right) Q(d\theta)$$

$\Updownarrow$

$X_1, X_2, \dots$  exchangeable

where:

- ▶  $\mathbf{M}(\mathcal{X})$  is the set of probability measures on  $\mathcal{X}$
- ▶  $\theta$  are values of a random probability measure  $\Theta$  with distribution  $Q$

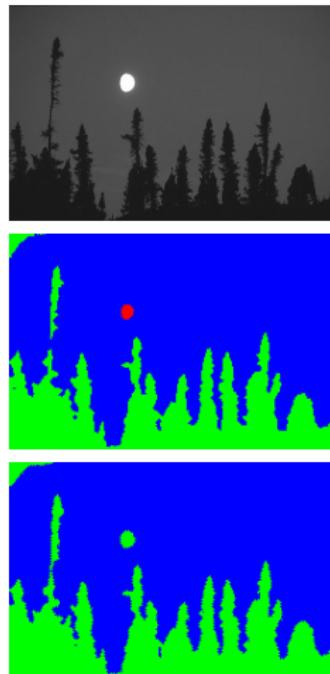
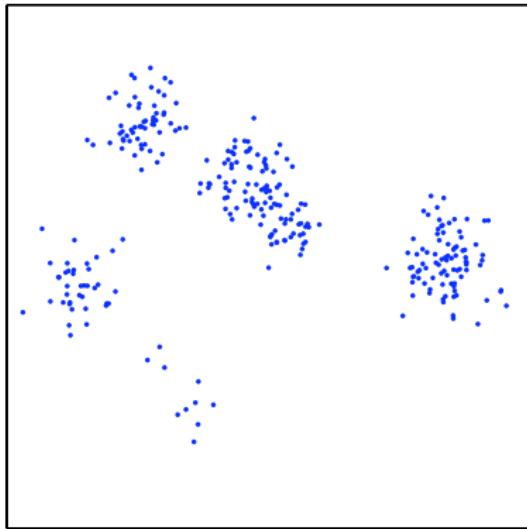
## Implications

- ▶ Exchangeable data decomposes into pattern and noise
- ▶ More general than i.i.d.-assumption
- ▶ Caution:  $\theta$  is in general an  $\infty$ -dimensional quantity

# CLUSTERING

# CLUSTERING

- ▶ Observations  $X_1, X_2, \dots$
- ▶ Each observation belongs to exactly one cluster
- ▶ Unknown pattern = partition of  $\{1, \dots, n\}$  or  $\mathbb{N}$



# MIXTURE MODELS

## Mixture models

$$p(x|m) = \int_{\Omega_\theta} p(x|\theta)m(d\theta)$$

$m$  is called the *mixing measure*

## Two-stage sampling

Sample  $X \sim p(\cdot|m)$  as:

1.  $\Theta \sim m$
2.  $X \sim p(\cdot|\theta)$

## Finite mixture model

$$p(x|\boldsymbol{\theta}, \mathbf{c}) = \int_{\Omega_\theta} p(x|\theta)m(d\theta) \quad \text{with} \quad m(\cdot) = \sum_{k=1}^K c_k \delta_{\theta_k}(\cdot)$$

# BAYESIAN MM

## Random mixing measure

$$M(\cdot) = \sum_{k=1}^K C_k \delta_{\Theta_k}(\cdot)$$

## Conjugate priors

A Bayesian model is *conjugate* if the posterior is an element of the same class of distributions as the prior ("closure under sampling").

$p(x \theta)$	conjugate prior
$\frac{1}{Z(\theta)} h(x) \exp(\langle S(x), \theta \rangle)$	$\frac{1}{K(\lambda, y)} \exp(\langle \theta, y \rangle - \lambda \log Z(\theta))$
Gaussian	Gaussian/inverse Wishart
multinomial	Dirichlet
...	...

## Choice of priors in BMM

- ▶ Choose conjugate prior for each parameter
- ▶ In particular: Dirichlet prior on  $(C_1, \dots, C_k)$

# DIRICHLET PROCESS MIXTURES

## Dirichlet process

A Dirichlet process is a distribution on random probability measures of the form

$$M(\cdot) = \sum_{k=1}^{\infty} C_k \delta_{\Theta_k}(\cdot) \quad \text{where} \quad \sum_{k=1}^{\infty} C_k = 1$$

## Constructive definition of DP $(\alpha, G_0)$

$$\Theta_k \sim_{\text{iid}} G_0$$

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

Compute  $C_k$  as

$$C_k := V_k \prod_{i=1}^{k-1} (1 - V_i)$$

"Stick-breaking construction"

# POSTERIOR DISTRIBUTION

## DP Posterior

$$\theta_{n+1} | \theta_1, \dots, \theta_n \sim \frac{1}{n + \alpha} \sum_{j=1}^n \delta_{\theta_j}(\theta_{n+1}) + \frac{\alpha}{n + \alpha} G_0(\theta_{n+1})$$

## Mixture Posterior

$$p(x_{n+1} | x_1, \dots, x_n) = \sum_{k=1}^{K_n} \frac{n_k}{n + \alpha} p(x_{n+1} | \theta_k^*) + \frac{\alpha}{n + \alpha} \int p(x_{n+1} | \theta) G_0(\theta) d\theta$$

## Conjugacy

- ▶ The posterior of DP  $(\alpha, G_0)$  is DP  $\left( \alpha + n, \frac{1}{n + \alpha} (\sum_k n_k \delta_{\theta_k^*} + \alpha G_0) \right)$
- ▶ Hence: The Dirichlet process is conjugate.

# INFERENCE

## Latent variables

$$p(x_{n+1}|x_1, \dots, x_n) = \sum_{k=1}^{K_n} \frac{n_k}{n + \alpha} p(x_{n+1}|\theta_k^*) + \frac{\alpha}{n + \alpha} \int p(x_{n+1}|\theta) G_0(\theta) d\theta$$

We do not actually observe the  $\Theta_j$  (they are latent). We observe  $X_j$ .

## Assignment probabilities

$$\begin{pmatrix} q_{10} & q_{11} & \dots & q_{1K_n} \\ \vdots & \vdots & & \vdots \\ q_{n0} & q_{n1} & \dots & q_{nK_n} \end{pmatrix}$$

Where:

- ▶  $q_{jk} \propto n_k p(x_j|\theta_k^*)$
- ▶  $q_{j0} \propto \alpha \int p(x_j|\theta) G_0(\theta) d\theta$

## Gibbs Sampling

Uses an assignment variable  $\phi_j$  for each observation  $X_j$ .

- ▶ Assignment step: Sample  $\phi_j \sim \text{Multinomial}(q_{j0}, \dots, q_{jK_n})$
- ▶ Parameter sampling:  $\theta_k^* \sim G_0(\theta_k^*) \prod_{x_j \in \text{Cluster}_k} p(x_j|\theta_k^*)$

# NUMBER OF CLUSTERS

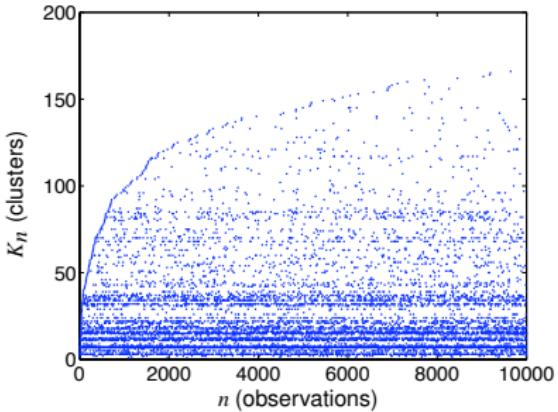
## Dirichlet process

$K_n = \#$  clusters in sample of size  $n$

$$\mathbb{E}[K_n] = O(\log(n))$$

## Modeling assumption

- ▶ Parametric clustering:  $K_\infty$  is *finite* (possibly unknown, but fixed).
- ▶ Nonparametric clustering:  $K_\infty$  is *infinite*



## Rephrasing the question

- ▶ Estimate of  $K_n$  is controlled by distribution of the cluster sizes  $C_k$  in  $\sum_k C_k \delta_{\Theta_k}$ .
- ▶ Ask instead: What should we assume about the distribution of  $C_k$ ?

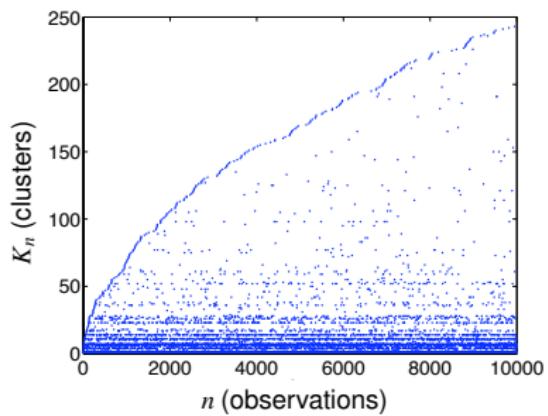
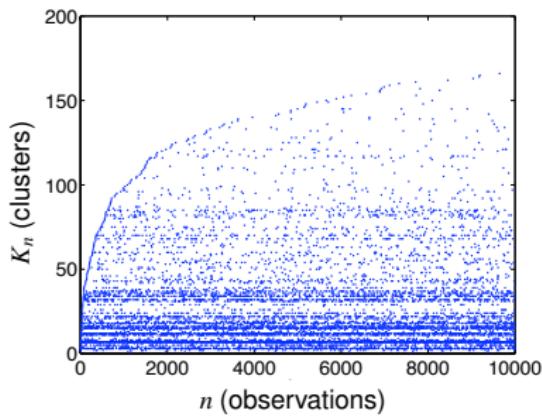
# GENERALIZING THE DP

## Pitman-Yor process

$$p(x_{n+1}|x_1, \dots, x_n) = \sum_{k=1}^{K_n} \frac{n_k - d}{n + \alpha} p(x_{n+1}|\theta_k^*) + \frac{\alpha + K_n \cdot d}{n + \alpha} \int p(x_{n+1}|\theta) G_0(\theta) d\theta$$

Discount parameter  $d \in [0, 1]$ .

## Cluster sizes



# POWER LAWS

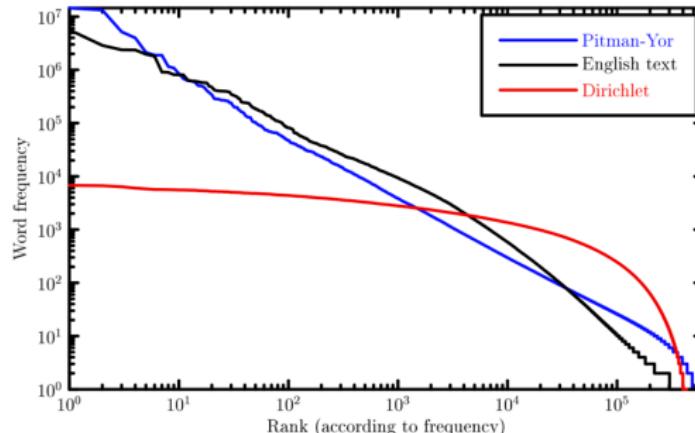
The distribution of cluster sizes is called a *power law* if

$$C_j \sim \gamma(\beta) \cdot j^{-\beta} \quad \text{for some } \beta \in [0, 1] .$$

## Examples of power laws

- ▶ Word frequencies
- ▶ Popularity (number of friends) in social networks

## Pitman-Yor language model



# RANDOM PARTITIONS

## Discrete measures and partitions

Sampling from a discrete measure determines a *partition* of  $\mathbb{N}$  into blocks  $b_k$ :

$$\Theta_n \sim_{\text{iid}} \sum_{k=1}^{\infty} c_k \delta_{\theta_k^*} \quad \text{and set} \quad n \in b_k \Leftrightarrow \Theta_n = \theta_k^*$$

As  $n \rightarrow \infty$ , the block proportions converge:  $\frac{|b_k|}{n} \rightarrow c_k$

## Induced random partition

The distribution of a random discrete measure  $M = \sum_{k=1}^{\infty} C_k \delta_{\Theta_k}$  induces the distribution of a *random partition*  $\Pi = (B_1, B_2, \dots)$ .

## Exchangeable random partitions

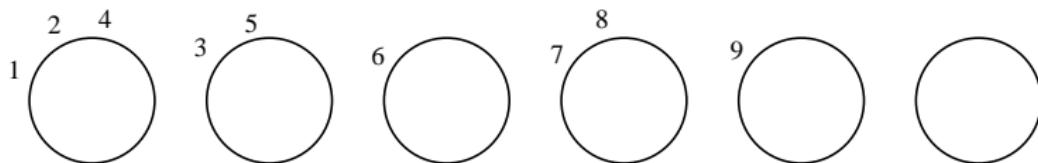
- $\Pi$  is called *exchangeable* if its distribution depends only on the sizes of its blocks.
- All exchangeable random partitions, and only those, can be represented by a random discrete distribution as above (Kingman's theorem).

# CHINESE RESTAURANT PROCESS

## Chinese Restaurant Process

The distribution of the random partition induced by the Dirichlet process is called the *Chinese Restaurant Process*.

"Customers and tables" analogy



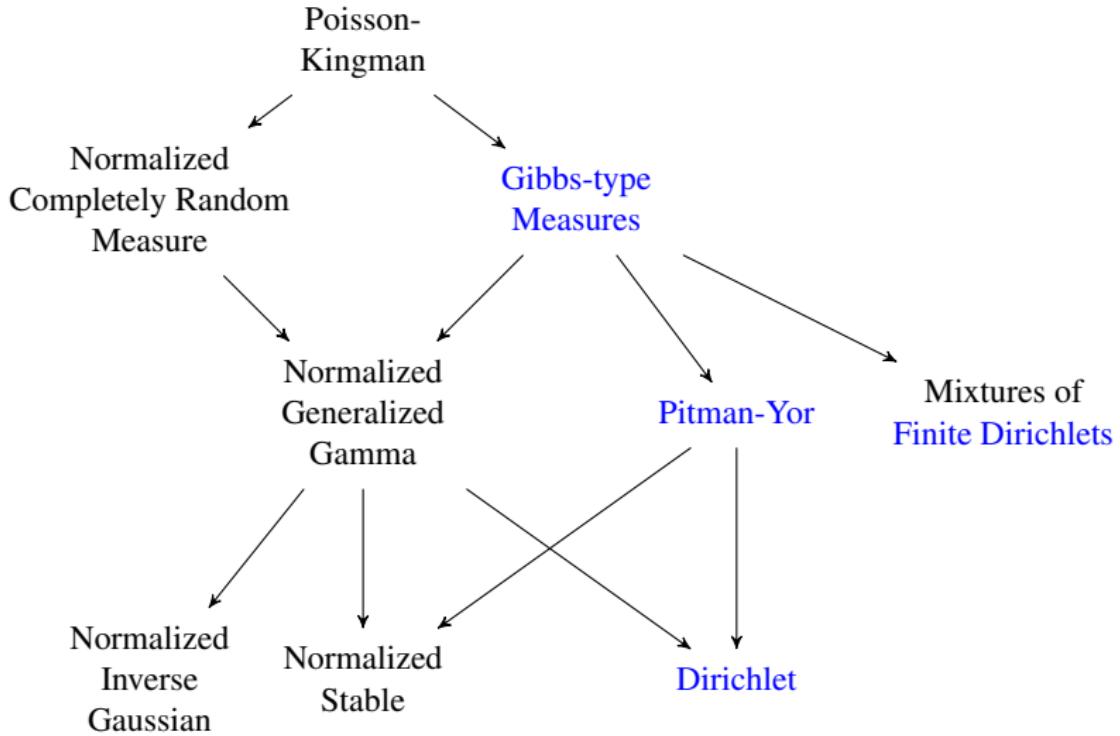
Customers = observations (indices in  $\mathbb{N}$ )

Tables = clusters (blocks)

## Historical remark

- ▶ Originally introduced by Dubins & Pitman as a distribution on infinite permutations
- ▶ A permutation of  $n$  items defines a partition of  $\{1, \dots, n\}$  (regard cycles of permutation as blocks of partition)
- ▶ The induced distribution on partitions is the CRP we use in clustering

# FAMILIES OF EXCHANGEABLE RANDOM PARTITIONS



# RANDOM DISCRETE MEASURES

## Classification (due to Prünster)

class	probability of new cluster	prior class
I	$\mathbb{P}\{\Theta_{n+1} \in \text{new cluster}   \Theta^{(n)}\} = f(n)$	Dirichlet processes
II	$\mathbb{P}\{\Theta_{n+1} \in \text{new cluster}   \Theta^{(n)}\} = f(n, K_n)$	Gibbs-type measures
III	$\mathbb{P}\{\Theta_{n+1} \in \text{new cluster}   \Theta^{(n)}\} = f(n, K_n, \mathbf{n})$	

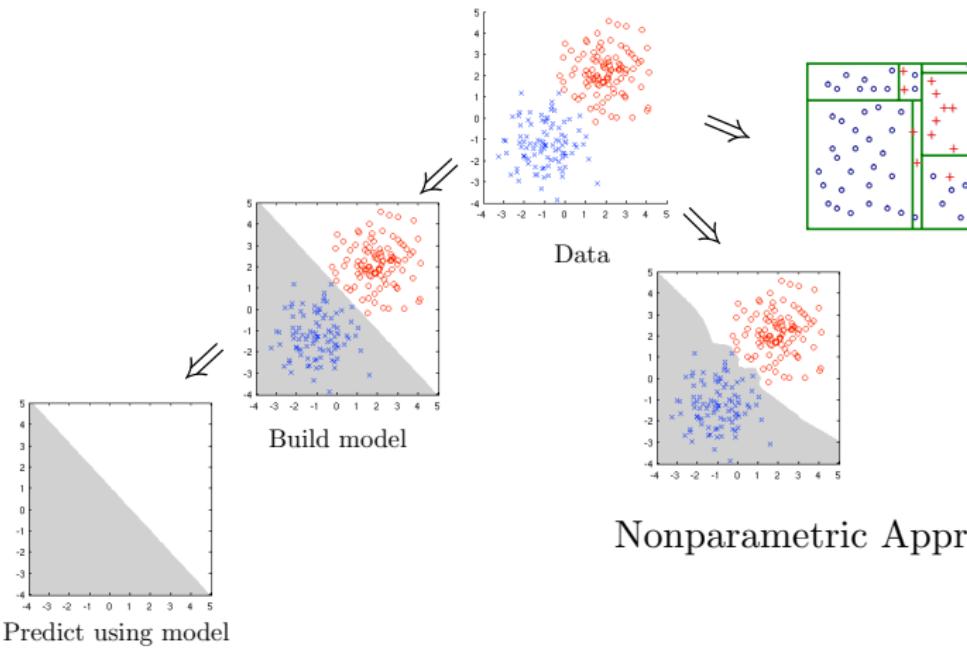
## General partition priors

- ▶ Gibbs-type measures are completely classified [GP06b]
- ▶ Properties of some cases well-studied, e.g.:
  - ▶ Dirichlet process
  - ▶ Pitman-Yor process
  - ▶ Normalized inverse Gaussian process [LMP05b]
- ▶ In the future: We will have a range of models which express different prior assumptions on the distribution of cluster sizes.

## Nonparametric

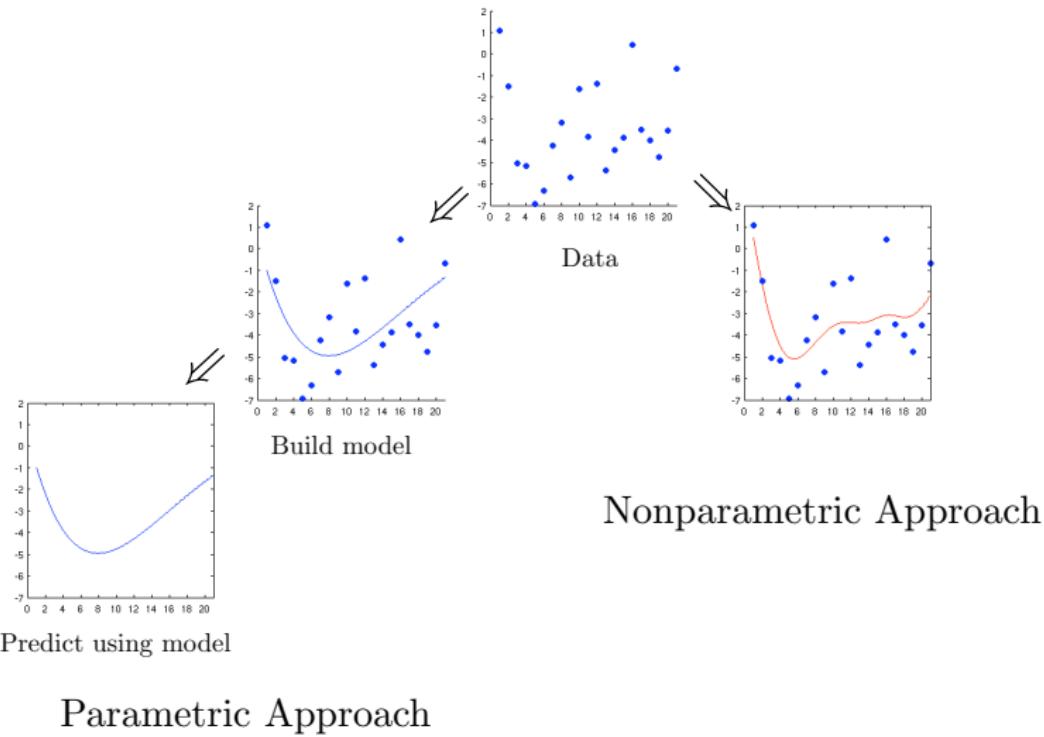
**Nonparametric**: Does NOT mean there are no parameters.

# Example: Classification

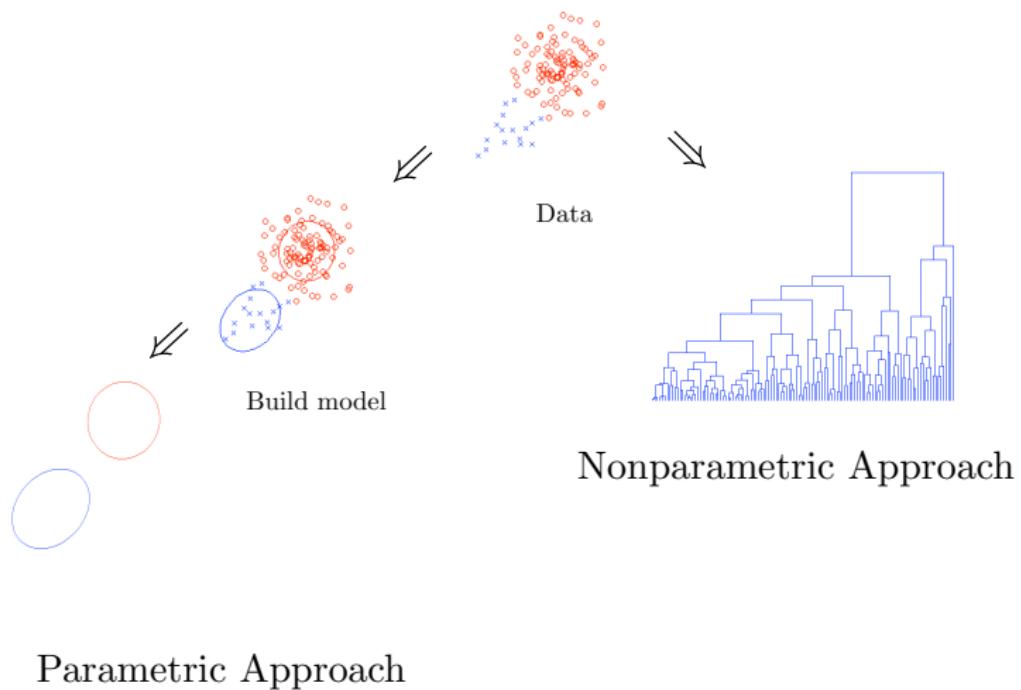


Parametric Approach

# Example: Regression



# Example: Clustering



# Why Be Bayesian?

You can take a course on this question.

# Why Be Bayesian?

You can take a course on this question. One answer:

**Infinite Exchangeability:**  $\forall n \ p(x_1, \dots, x_n) = p(x_{\sigma(1)}, \dots, x_{\sigma(n)})$

# Why Be Bayesian?

You can take a course on this question. One answer:

**Infinite Exchangeability:**  $\forall n \ p(x_1, \dots, x_n) = p(x_{\sigma(1)}, \dots, x_{\sigma(n)})$

**De Finetti's Theorem (1955):** If  $(x_1, x_2, \dots)$  are *infinitely exchangeable*, then  $\forall n$

$$p(x_1, \dots, x_n) = \int \left( \prod_{i=1}^n p(x_i | \theta) \right) dP(\theta)$$

for some random variable  $\theta$ .

## Simple Example

Task: Toss a (potentially biased) coin  $N$  times. Compute  $\theta$ , the probability of heads.

Suppose we observe: {T, H, H, T}. What do we think  $\theta$  is?

## Simple Example

Task: Toss a (potentially biased) coin  $N$  times. Compute  $\theta$ , the probability of heads.

Suppose we observe: {T, H, H, T}. What do we think  $\theta$  is? The maximum likelihood estimate is  $\theta = 1/2$ . Seems reasonable.

## Simple Example

Task: Toss a (potentially biased) coin  $N$  times. Compute  $\theta$ , the probability of heads.

Suppose we observe: {T, H, H, T}. What do we think  $\theta$  is? The maximum likelihood estimate is  $\theta = 1/2$ . Seems reasonable.

Now suppose we observe: {H, H, H, H}. What do we think  $\theta$  is?

## Simple Example

Task: Toss a (potentially biased) coin  $N$  times. Compute  $\theta$ , the probability of heads.

Suppose we observe: {T, H, H, T}. What do we think  $\theta$  is? The maximum likelihood estimate is  $\theta = 1/2$ . Seems reasonable.

Now suppose we observe: {H, H, H, H}. What do we think  $\theta$  is? The maximum likelihood estimate is  $\theta = 1$ . Seem reasonable?

## Simple Example

Task: Toss a (potentially biased) coin  $N$  times. Compute  $\theta$ , the probability of heads.

Suppose we observe: {T, H, H, T}. What do we think  $\theta$  is? The maximum likelihood estimate is  $\theta = 1/2$ . Seems reasonable.

Now suppose we observe: {H, H, H, H}. What do we think  $\theta$  is? The maximum likelihood estimate is  $\theta = 1$ . Seem reasonable?

Not really. Why?

## Simple Example

When we observe  $\{H, H, H, H\}$ , why does  $\theta = 1$  seem unreasonable?

## Simple Example

When we observe {H, H, H, H}, why does  $\theta = 1$  seem unreasonable?

Prior knowledge! We believe coins generally have  $\theta \approx 1/2$ . How to encode this? By using a Beta *prior on  $\theta$* .

## Bayesian Approach to Estimating $\theta$

Place a  $\text{Beta}(a, b)$  prior on  $\theta$ . This prior has the form

$$p(\theta) \propto \theta^{a-1} (1-\theta)^{b-1}.$$

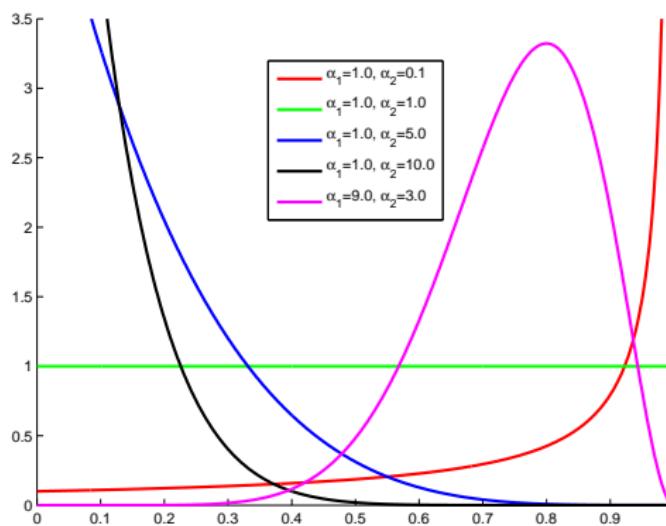
What does this distribution look like?

# Bayesian Approach to Estimating $\theta$

Place a Beta( $a, b$ ) prior on  $\theta$ . This prior has the form

$$p(\theta) \propto \theta^{a-1} (1-\theta)^{b-1}.$$

What does this distribution look like?



## Bayesian Approach to Estimating $\theta$

After observing  $X$ , a sequence with  $n$  heads and  $m$  tails, the posterior on  $\theta$  is:

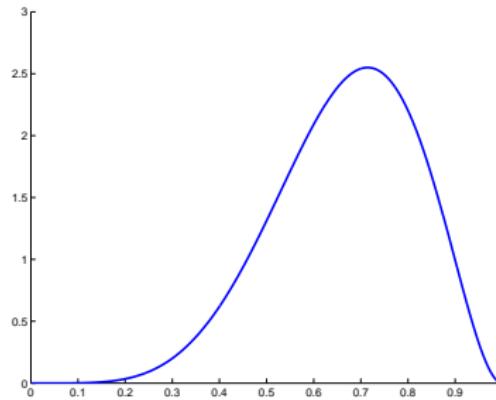
$$\begin{aligned} p(\theta|X) &\propto p(X|\theta)p(\theta) \\ &\propto \theta^{a+n-1}(1-\theta)^{b+m-1} \\ &\sim \text{Beta}(a+n, b+m). \end{aligned}$$

## Bayesian Approach to Estimating $\theta$

After observing  $X$ , a sequence with  $n$  heads and  $m$  tails, the posterior on  $\theta$  is:

$$\begin{aligned} p(\theta|X) &\propto p(X|\theta)p(\theta) \\ &\propto \theta^{a+n-1}(1-\theta)^{b+m-1} \\ &\sim \text{Beta}(a+n, b+m). \end{aligned}$$

If  $a = b = 1$  and we observe 5 heads and 2 tails, Beta(6, 3) looks like



# Nonparametric Bayesian Methods

Now we know what **nonparametric** and **Bayesian** mean. What should we expect from **nonparametric Bayesian** methods?

- Complexity of our model should be allowed to grow as we get more data.
- Place a prior on an unbounded number of parameters.