

Non-parametric Bayesian Models

Jie Tang

Department of Computer Science and Technology

Tsinghua University

Contents

- 1 Introduction
- 2 Gaussian Process
 - Multivariate Gaussian Distribution
 - Definition
 - GP Regression
- 3 Dirichlet Process
 - Definition
 - Stick Breaking Construction
 - Blackwell-MacQueen Urn Scheme
 - Chinese Restaurant Process
 - DP Mixture

Contents

- 1 Introduction
- 2 Gaussian Process
 - Multivariate Gaussian Distribution
 - Definition
 - GP Regression
- 3 Dirichlet Process
 - Definition
 - Stick Breaking Construction
 - Blackwell-MacQueen Urn Scheme
 - Chinese Restaurant Process
 - DP Mixture

Introduction

We'll talk about *nonparametric Bayesian methods* in following slides

- Regression, classification – Gaussian Process

- Clustering – Dirichlet Process

What is *nonparametric* ?

What is *Bayesian* ?

What is Nonparametric Model?

Nonparametric doesn't mean there are no parameters.

Parametric model assumes a finite set of parameters.

- Finite model parameters captures everything about the data
- Predictions for future data are made merely through parameters, independent from observed data.

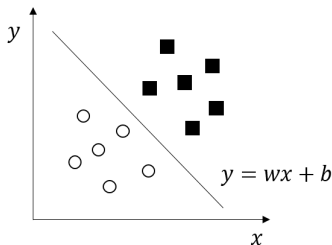
- Model complexity are bounded, even with infinite data

Nonparametric model assumes the model can't be specified by a finite set of parameters. But often a infinite-dimensional parameter space will work.

- Model complexity can grow with data

Example: Classification

Consider a linear model

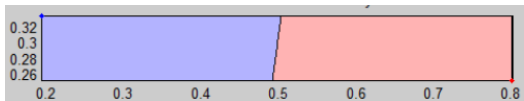


No matter how many samples are trained, number of parameters are fixed (w, b)

Example: Classification

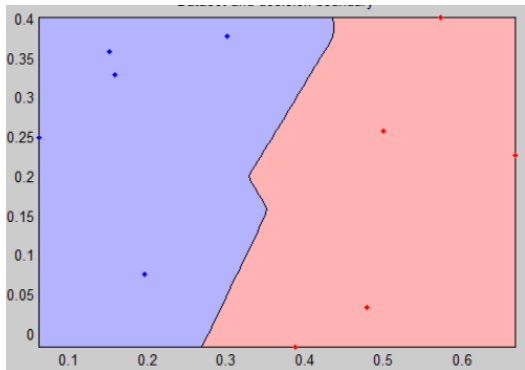
How about k-NN?

2 samples



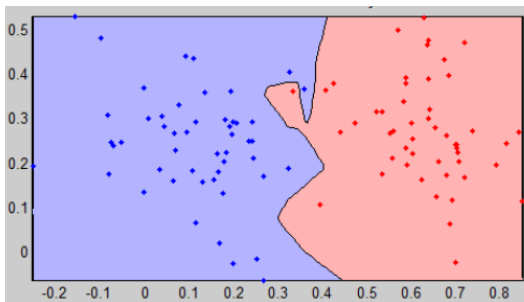
Example: Classification

10 samples



Example: Classification

100 samples



Parameters are unbounded as dataset grow larger

Nonparametric!

What about SVM?

Linear SVM

$$y = \text{sign}(w^T x)$$

What about SVM?

Linear SVM

$$y = \text{sign}(w^T x)$$

Kernel SVM

$$y = \text{sign}\left(\sum_i \alpha_i y_i \kappa(x_i, x)\right)$$

Parametric vs. Nonparametric

Parametric Models:

- Inflexible. Simple and easy to interpret.

- Performs well and fast, if model captures the data correctly

- Otherwise will give inconsistent result

Nonparametric Models:

- Flexible.

- Complex and hard to interpret, thus giving inaccurate estimation

Question: which of the following algorithm belongs to parametric models?

- (A) k-NN
- (B) SVM
- (C) Logistic Regression
- (D) Decision Tree

What is Bayesian?

Consider a regression problem

Given training data D , we aim to find a model f . For each sample x , return a prediction y

How to give an optimal y ?

What is Bayesian?

Consider a regression problem

Given training data D , we aim to find a model f . For each sample x , return a prediction y

How to give an optimal y ?

- MLE

- MAP

- Bayesian

Maximum Likelihood Estimation

In training, we maximize a likelihood

$$f^* = \arg \max_f L(f|D) = \arg \max_f p(D|f)$$

In prediction, just feed f with x

$$p(y|x, D) = p(y|f^*, x)$$

Maximum A Posteriori

To avoid overfitting, we place a prior on the model

$$f^* = \arg \max_f \frac{p(D|f)p(f)}{p(D)}$$

Prediction process remains the same

$$p(y|x, D) = p(y|f^*, x)$$

Bayesian Model

In training, we just estimate the posterior distribution of f (same as MAP)

$$p(f|D) = \frac{p(D|f)p(f)}{p(D)}$$

In prediction, we no longer want a point estimation. Instead, we consider all possible models and take the expectation

$$p(y|x, D) = \int_f p(y|f, x)p(f|D) df$$

Why Bayesian?

Infinite Exchangeability:

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

De Finetti's Theorem: If (x_1, x_2, \dots) are infinite exchangeable, then there exists a random variable f , $\forall n$,

$$p(x_1, \dots, x_n) = \int_f \prod_{i=1}^n p(x_i|f) p(f) df$$

Why Bayesian?

Infinite Exchangeability:

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

De Finetti's Theorem: If (x_1, x_2, \dots) are infinite exchangeable, then there exists a random variable f , $\forall n$,

$$p(x_1, \dots, x_n) = \int_f \prod_{i=1}^n p(x_i|f) p(f) df$$

How to define the prior $p(f)$?

Parametric Approach

In previous courses, we focused on parametric representation of f , e.g.

$$f(x) = w^T x$$

Instead of defining **prior** on function itself, we define prior on the parameters

$$p(w) = \mathcal{N}(0, \sigma^2 I)$$

Now, we define the prior on the function f directly

Contents

- 1 Introduction
- 2 Gaussian Process
 - Multivariate Gaussian Distribution
 - Definition
 - GP Regression
- 3 Dirichlet Process
 - Definition
 - Stick Breaking Construction
 - Blackwell-MacQueen Urn Scheme
 - Chinese Restaurant Process
 - DP Mixture

Warmup: Multivariate Gaussian

Before introducing Gaussian Process, let's recall multivariate Gaussian distribution first.

We say a n -dimensional random variable \mathbf{x} follows multivariate Gaussian

$$\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

if

$$p(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^n |\boldsymbol{\Sigma}|}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right\}$$

where

$$\mathbb{E}[\mathbf{x}] = \boldsymbol{\mu}$$

$$\mathbb{E}[(x_i - \mu_i)(x_j - \mu_j)] = \Sigma_{ij}$$

Marginal Distribution

Suppose

$$\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix} \right)$$

Integrating \mathbf{x}_2 out,

$$\mathbf{x}_1 \sim \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11})$$

Just drop irrelevant components

Conditional Distribution

Suppose

$$\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix} \right)$$

Then

$$p(\mathbf{x}_1 | \mathbf{x}_2) = \mathcal{N}(\mathbf{x}_1 | \boldsymbol{\mu}_{1|2}, \boldsymbol{\Sigma}_{1|2})$$

$$\boldsymbol{\mu}_{1|2} = \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} (\mathbf{x}_2 - \boldsymbol{\mu}_2)$$

$$\boldsymbol{\Sigma}_{1|2} = \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{21}$$

We will use it later

How to get that?

Some Details

First, we diagonalize the variance matrix to ease the inverse matrix calculation

$$\begin{bmatrix} \Sigma_{11}/\Sigma_{22} & O \\ O & \Sigma_{22} \end{bmatrix} = \begin{bmatrix} I_1 & -\Sigma_{12}\Sigma_{22}^{-1} \\ O & I_2 \end{bmatrix} \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \begin{bmatrix} I_1 & O \\ -\Sigma_{22}^{-1}\Sigma_{21} & I_2 \end{bmatrix}$$

Define $\Sigma/\Sigma_{22} = \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}$ here. Then

$$\begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}^{-1} = \begin{bmatrix} I_1 & O \\ -\Sigma_{22}^{-1}\Sigma_{21} & I_2 \end{bmatrix} \begin{bmatrix} (\Sigma/\Sigma_{22})^{-1} & O \\ O & \Sigma_{22}^{-1} \end{bmatrix} \begin{bmatrix} I_1 & -\Sigma_{12}\Sigma_{22}^{-1} \\ O & I_2 \end{bmatrix}$$

More Details

Check the exponent of $p(\mathbf{x}_1, \mathbf{x}_2)$:

$$\begin{aligned} & (x_1 - \mu_1, x_2 - \mu_2) \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}^{-1} \begin{pmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{pmatrix} \\ &= \begin{pmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{pmatrix}^T \begin{bmatrix} I_1 & O \\ -\Sigma_{22}^{-1} \Sigma_{21} & I_2 \end{bmatrix} \begin{bmatrix} (\Sigma/\Sigma_{22})^{-1} & O \\ O & \Sigma_{22}^{-1} \end{bmatrix} \\ & \quad \begin{bmatrix} I_1 - \Sigma_{12} \Sigma_{22}^{-1} \\ O & I_2 \end{bmatrix} \begin{pmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{pmatrix} \\ &= (x_1 - \mu_1 - \Sigma_{12} \Sigma_{22}^{-1} (x_2 - \mu_2))^T (\Sigma/\Sigma_{22})^{-1} \\ & \quad (x_1 - \mu_1 - \Sigma_{12} \Sigma_{22}^{-1} (x_2 - \mu_2)) + C \end{aligned}$$

Note that $p(\mathbf{x}_1|\mathbf{x}_2) = \frac{p(\mathbf{x}_1, \mathbf{x}_2)}{p(\mathbf{x}_2)}$, and C is just the exponent of $p(\mathbf{x}_2)$.

From the equation we can get the mean and variance of the conditional distribution.

Contents

- 1 Introduction
- 2 Gaussian Process
 - Multivariate Gaussian Distribution
 - Definition
 - GP Regression
- 3 Dirichlet Process
 - Definition
 - Stick Breaking Construction
 - Blackwell-MacQueen Urn Scheme
 - Chinese Restaurant Process
 - DP Mixture

Gaussian Process

Gaussian Process is a generalization of multivariate Gaussian to infinite variables – *stochastic process*

$\forall \mathbf{x}_1, \dots, \mathbf{x}_n, (f(\mathbf{x}_1), \dots, f(\mathbf{x}_n))$ is jointly Gaussian

A Gaussian Process is fully specified by mean function $m(\mathbf{x})$ and covariance function $\kappa(\mathbf{x}, \mathbf{x}')$

$$m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})]$$

$$\kappa(\mathbf{x}, \mathbf{x}') = \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))]$$

Then,

$$(f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_n))^T \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{K})$$

where

$$\mu_i = m(\mathbf{x}_i), K_{ij} = \kappa(\mathbf{x}_i, \mathbf{x}_j)$$

We require κ to be a positive definite kernel.

Contents

- 1 Introduction
- 2 Gaussian Process
 - Multivariate Gaussian Distribution
 - Definition
 - GP Regression
- 3 Dirichlet Process
 - Definition
 - Stick Breaking Construction
 - Blackwell-MacQueen Urn Scheme
 - Chinese Restaurant Process
 - DP Mixture

GP Regression with noise-free observations

For simplicity, let's start with predictions using noise free observations

Training set $D = \{(x_i, f_i)\}$, where $f_i = f(x_i)$

Noise free! – No uncertainty.

Test set \mathbf{X}_* . We want to predict outputs f_*

From definition of GP,

$$\begin{pmatrix} \mathbf{f} \\ \mathbf{f}_* \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} \boldsymbol{\mu} \\ \boldsymbol{\mu}_* \end{pmatrix}, \begin{bmatrix} \mathbf{K} & \mathbf{K}_* \\ \mathbf{K}_*^T & \mathbf{K}_{**} \end{bmatrix} \right)$$

GP Regression with noise-free observations

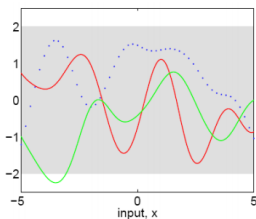
Using results of marginal multivariate Gaussian,

$$\begin{aligned}p(\mathbf{f}_*|\mathbf{X}_*, \mathbf{X}, \mathbf{f}) &= \mathcal{N}(\boldsymbol{\mu}_*, \boldsymbol{\Sigma}_*) \\ \boldsymbol{\mu}_* &= \boldsymbol{\mu}(\mathbf{X}_*) + \mathbf{K}_*^T \mathbf{K}^{-1}(\mathbf{f} - \boldsymbol{\mu}(\mathbf{X})) \\ \boldsymbol{\Sigma}_* &= \mathbf{K}_{**} - \mathbf{K}_*^T \mathbf{K}^{-1} \mathbf{K}_*\end{aligned}$$

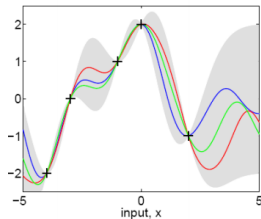
We usually set $\boldsymbol{\mu}(\mathbf{x}) = 0$, because GP can model the mean arbitrarily well.

Visualization

Let's see some samples from GP prior and posterior.



(a) prior

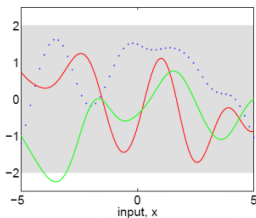


(b) posterior

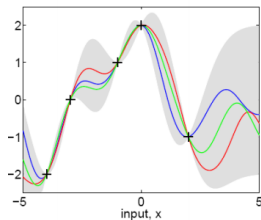
Training data points are marked with '+'.
No uncertainty at training data. Why?

Visualization

Let's see some samples from GP prior and posterior.



(c) prior



(d) posterior

Training data points are marked with '+'.
No uncertainty at training data. Why?

$\Sigma_* = \mathbf{K}_{**} - \mathbf{K}_*^T \mathbf{K}^{-1} \mathbf{K}_* = 0$ when \mathbf{K}_* and \mathbf{K}_{**} are sub-matrices of \mathbf{K} – Noise-free

GP Regression with noisy observations

When observations are noisy, the true f is hidden.

Instead, we observe f plus a Gaussian noise

$$y = f(\mathbf{x}) + \epsilon$$

where

$$\epsilon \sim \mathcal{N}(0, \sigma_y^2)$$

Now

$$\text{cov}[y_p, y_q] = \kappa(\mathbf{x}_p, \mathbf{x}_q) + \sigma_y^2 \delta_{pq}$$

in other words

$$\text{cov}[\mathbf{y}|\mathbf{X}] = \mathbf{K} + \sigma_y^2 \mathbf{I} \triangleq \mathbf{K}_y$$

GP Regression with noisy observations

Take $m(\mathbf{x}) = 0$ here

Similar to noise-free result

$$\begin{pmatrix} \mathbf{y} \\ \mathbf{f}_* \end{pmatrix} \sim \mathcal{N} \left(0, \begin{bmatrix} \mathbf{K}_y & \mathbf{K}_* \\ \mathbf{K}_*^T & \mathbf{K}_{**} \end{bmatrix} \right)$$

$$p(\mathbf{f}_* | \mathbf{X}_*, \mathbf{X}, \mathbf{y}) = \mathcal{N}(\boldsymbol{\mu}_*, \boldsymbol{\Sigma}_*)$$

$$\boldsymbol{\mu}_* = \mathbf{K}_*^T \mathbf{K}_y^{-1} \mathbf{y}$$

$$\boldsymbol{\Sigma}_* = \mathbf{K}_{**} - \mathbf{K}_*^T \mathbf{K}_y^{-1} \mathbf{K}_*$$

Deeper Insight

Feed the result with a single sample \mathbf{x}_*

$$p(f_*|\mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \mathcal{N}(\mathbf{k}_*^T \mathbf{K}_y^{-1} \mathbf{y}, k_{**} - \mathbf{k}_*^T \mathbf{K}_y^{-1} \mathbf{k}_*)$$
$$\bar{f}_* = \sum_i \alpha_i \kappa(\mathbf{x}_i, \mathbf{x}_*) \quad \alpha = \mathbf{K}_y^{-1} \mathbf{y}$$

What can we find?

Deeper Insight

Feed the result with a single sample \mathbf{x}_*

$$p(f_*|\mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \mathcal{N}(\mathbf{k}_*^T \mathbf{K}_y^{-1} \mathbf{y}, k_{**} - \mathbf{k}_*^T \mathbf{K}_y^{-1} \mathbf{k}_*)$$
$$\bar{f}_* = \sum_i \alpha_i \kappa(\mathbf{x}_i, \mathbf{x}_*) \quad \alpha = \mathbf{K}_y^{-1} \mathbf{y}$$

What can we find?

Although we define GP on infinite dimension, it's suffices to consider only $n + 1$ dimension.

In terms of prediction mean, GP regression can be considered as linear predictor when features are projected to high dimension space using feature map $\phi(\cdot)$ implied by κ . (Recall $\kappa(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$)

Effect of Kernel Parameters

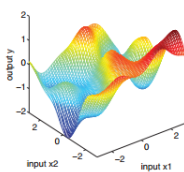
Kernel parameters are crucial in GP

RBF kernel: $\kappa(\mathbf{x}_i, \mathbf{x}_j) = \sigma_f^2 \exp\{(\mathbf{x}_i - \mathbf{x}_j)^T \mathbf{M}(\mathbf{x}_i - \mathbf{x}_j)\} + \sigma_y^2 \delta_{ij}$

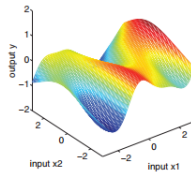
$$\mathbf{M} = \mathbf{I}$$

$$\mathbf{M} = \text{diag}(1, 3)^{-2}$$

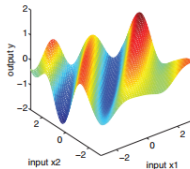
$$\mathbf{M} = (1, -1; -1, 1) + \text{diag}(6, 6)^{-2}$$



(a)



(b)



(c)

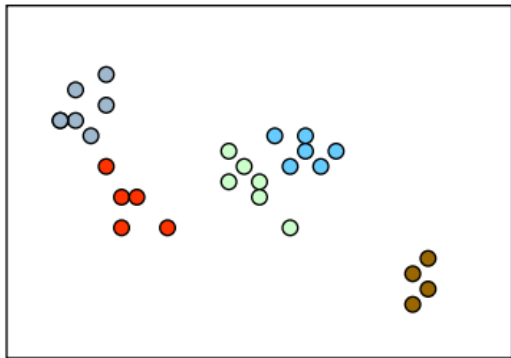
Contents

- 1 Introduction
- 2 Gaussian Process
 - Multivariate Gaussian Distribution
 - Definition
 - GP Regression
- 3 Dirichlet Process
 - Definition
 - Stick Breaking Construction
 - Blackwell-MacQueen Urn Scheme
 - Chinese Restaurant Process
 - DP Mixture

Example: Clustering

We're give a data set which are generated by a mixture of Gaussian.

How to model the data and cluster them?



Gaussian Mixture Clustering

Suppose the dataset is generated by
 K Gaussian components

Generating process:

$$\phi_k = (\mu_k, \Sigma_k) \sim \mathcal{NIW}(\nu)$$

$$\pi \sim \text{Dirichlet}(\alpha/K, \dots, \alpha/K)$$

$$z_i \sim \text{Multinomial}(\pi)$$

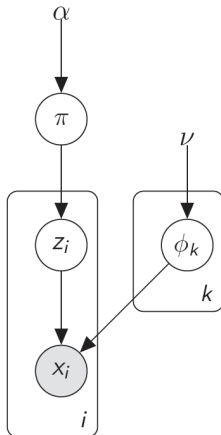
$$x_i \sim \mathcal{N}(\mu_{z_i}, \Sigma_{z_i})$$

where

$\mathcal{NIW}(\nu)$: conjugate prior of Gaussian

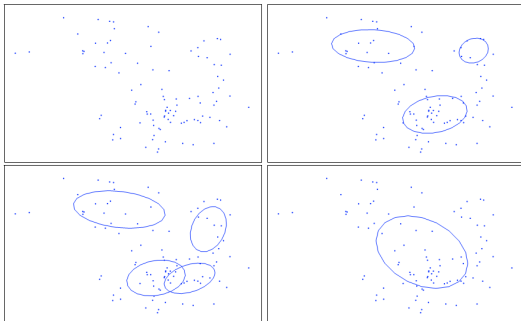
$\text{Dirichlet}(\dots)$: conjugate prior of multinomial

z_i : mixture indicator of x_i



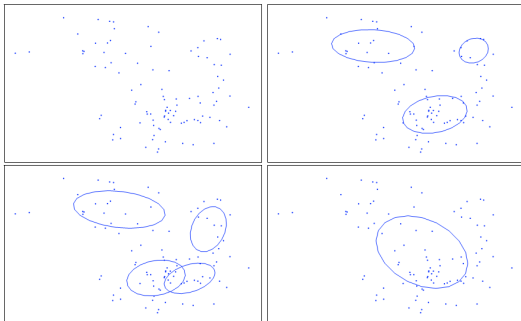
How to choose K?

How many clusters?



How to choose K?

How many clusters?



What if let $K \rightarrow \infty$?

Infinity Mixture Models

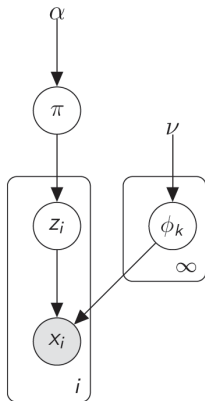
Imagine $K \gg 0$,

In Bayesian inference, we integrate ϕ_k and π out. Number of latent variables would not grow with K – **No overfitting**

At most n **active components** are associated with data

This is an *infinite mixture model*

Issue: can we take $K \rightarrow \infty$? What's the limiting model?



An Alternative View of Mixture Model

Mixture parameter for each data point can be sampled directly

z_i is absorbed into a *discrete probability measure* G

δ_{ϕ_k} is an *atom* positioned at ϕ_k

θ_i is the **sampled Gaussian parameter** to generate x_i

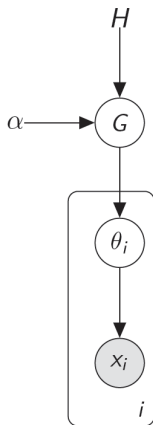
$$\phi_k \sim H$$

$$\pi_k \sim \text{Dirichlet}(\alpha/K, \dots, \alpha/K)$$

$$G = \sum_{k=1}^K \pi_k \delta_{\phi_k}$$

$$\theta_i \sim G$$

$$x_i \sim p(x_i | \theta_i)$$

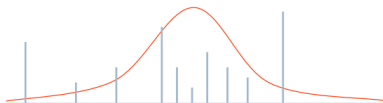


An Alternative View of Mixture Model

H: continuous



G: discrete, with mass at ϕ_k



Recall

$$G = \sum_{k=1}^K \pi_k \delta_{\phi_k}$$

Take $K \rightarrow \infty$,

$$G = \sum_{k=1}^{\infty} \pi_k \delta_{\phi_k}$$

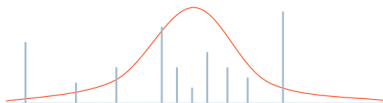
To execute Bayesian inference, we need to **find a prior** for G

An Alternative View of Mixture Model

H: continuous



G: discrete, with mass at ϕ_k



Recall

$$G = \sum_{k=1}^K \pi_k \delta_{\phi_k}$$

Take $K \rightarrow \infty$,

$$G = \sum_{k=1}^{\infty} \pi_k \delta_{\phi_k}$$

To execute Bayesian inference, we need to find a prior for G

Dirichlet Process

Warmup: Dirichlet Distribution

Dirichlet distribution is a distribution over the K -dim probability simplex

$$\Omega_K = \{(\pi_1, \dots, \pi_K) : \pi_k \geq 0, \sum_k \pi_k = 1\}$$

We say (π_1, \dots, π_K) is Dirichlet distributed with parameters $(\alpha_1, \dots, \alpha_K)$, if

$$p(\pi_1, \dots, \pi_K) = \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_k \pi_k^{\alpha_k - 1}$$

Warmup: Properties of Dirichlet Distribution

Aggregation Property:

If

$$(\pi_1, \dots, \pi_i, \pi_{i+1}, \dots, \pi_K) \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_i, \alpha_{i+1}, \dots, \alpha_K)$$

then

$$(\pi_1, \dots, \pi_i + \pi_{i+1}, \dots, \pi_K) \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_i + \alpha_{i+1}, \dots, \alpha_K)$$

Inverse of Aggregation Property:

If

$$(\pi_1, \dots, \pi_K) \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_K)$$

$$t \sim \text{Beta}(\alpha_i \beta, \alpha_i (1 - \beta))$$

$$(aka(t, 1 - t) \sim \text{Dirichlet}(\beta \alpha_i, (1 - \beta) \alpha_i))$$

then

$$\begin{aligned} &(\pi_1, \dots, t\pi_i, (1 - t)\pi_{i+1}, \dots, \pi_K) \\ &\sim \text{Dirichlet}(\alpha_1, \dots, \beta \alpha_i, (1 - \beta) \alpha_{i+1}, \dots, \alpha_K) \end{aligned}$$

Warmup: Dirichlet-Multinomial conjugacy

Dirichlet-Multinomial conjugacy:

If

$$\begin{aligned}(\pi_1, \dots, \pi_K) &\sim \text{Dirichlet}(\alpha_1, \dots, \alpha_K) \\ z &\sim \text{Multinomial}(\pi_1, \dots, \pi_K)\end{aligned}$$

then

$$(\pi_1, \dots, \pi_K) | z \sim \text{Dirichlet}(\alpha_1 + \delta_z(1), \dots, \alpha_K + \delta_z(K))$$

Contents

- 1 Introduction
- 2 Gaussian Process
 - Multivariate Gaussian Distribution
 - Definition
 - GP Regression
- 3 Dirichlet Process
 - Definition
 - Stick Breaking Construction
 - Blackwell-MacQueen Urn Scheme
 - Chinese Restaurant Process
 - DP Mixture

Dirichlet Process

Starting from π , slice the space finer and finer, what do we get?

$$(\pi) \sim \textit{Dirichlet}(\alpha)$$

$$(\pi_1, \pi_2) \sim \textit{Dirichlet}(\alpha_1, \alpha_2)$$

$$(\pi_{11}, \pi_{12}, \pi_{21}, \pi_{22}) \sim \textit{Dirichlet}(\alpha_{11}, \alpha_{12}, \alpha_{21}, \alpha_{22}) \dots$$

Because the dividing schema is arbitrary, finally, any finite partition of π should follow Dirichlet distribution.

Dirichlet Process

Starting from π , slice the space finer and finer, what do we get?

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

$$(\pi_1, \pi_2) \sim \text{Dirichlet}(\alpha_1, \alpha_2)$$

$$(\pi_{11}, \pi_{12}, \pi_{21}, \pi_{22}) \sim \text{Dirichlet}(\alpha_{11}, \alpha_{12}, \alpha_{21}, \alpha_{22}) \dots$$

Because the dividing schema is arbitrary, finally, any finite partition of π should follow Dirichlet distribution.

Dirichlet Process is a generalization of Dirichlet distribution to infinite many components – Recall Gaussian Process and multivariate Gaussian distribution.

Dirichlet Process

Let G be a probability measure over Ω . For any measurable subset $A \subset \Omega$, $G(A) = p(x \in A), x \in \Omega$

Dirichlet process (DP) is a distribution over G .

We write

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

if for any finite partition (A_1, \dots, A_n) of Ω ,

$$(G(A_1), \dots, G(A_n)) \sim \text{Dirichlet}(\alpha H(A_1), \dots, \alpha H(A_n))$$

H is the **base distribution**

$$E(G(A)) = H(A)$$

α is the **concentration parameter**

$$\text{Var}(G(A)) = \frac{H(A)(1 - H(A))}{\alpha + 1}$$

Dirichlet Process

Our definition is all about properties of DP. How to construct a concrete sample from DP?

Why is DP useful for clustering?

We'll see later

Posterior Dirichlet Process

We'll need the posterior Dirichlet Process later

Suppose $G \sim DP(\alpha, H)$, $\theta \sim G$

What is $p(\theta)$ and $G|\theta$?

For any $\theta \in A \subset \Omega$,

$$p(A) = \int_G G(A)p(G) \, dG = H(A)$$

So $\theta \sim H$

Posterior Dirichlet Process

Using Dirichlet-multinomial conjugacy, for any partition (A_1, \dots, A_n) ,

$$\begin{aligned} & (G(A_1), \dots, G(A_n)) | \theta \\ & \sim \text{Dirichlet}(\alpha H(A_1) + \delta_\theta(A_1), \dots, \alpha H(A_n) + \delta_\theta(A_n)) \end{aligned}$$

where $\delta_\theta(A) = 1$ iff $\theta \in A$

So, the posterior is also a DP

$$G | \theta \sim DP \left(\alpha + 1, \frac{\alpha H + \delta_\theta}{\alpha + 1} \right)$$

Generalization for n observations:

$$G | \theta_1, \dots, \theta_n \sim DP \left(\alpha + n, \frac{\alpha H + \sum_{i=1}^n \delta_{\theta_i}}{\alpha + n} \right)$$

We will construct a sample from DP based on the posterior

Contents

- 1 Introduction
- 2 Gaussian Process
 - Multivariate Gaussian Distribution
 - Definition
 - GP Regression
- 3 Dirichlet Process
 - Definition
 - Stick Breaking Construction
 - Blackwell-MacQueen Urn Scheme
 - Chinese Restaurant Process
 - DP Mixture

Stick Breaking

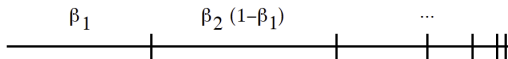
Define an infinite sequence of Beta random variables:

$$\beta_k \sim \text{Beta}(1, \alpha_0) \quad k = 1, 2, \dots$$

And then define an infinite sequence of mixing proportions as:

$$\begin{aligned} \pi_1 &= \beta_1 \\ \pi_k &= \beta_k \prod_{l=1}^{k-1} (1 - \beta_l) \quad k = 2, 3, \dots \end{aligned}$$

This can be viewed as breaking off portions of a stick:



Stick Breaking Construction with Dirichlet Process

Sample θ_1 from H

Consider partition $(\theta_1, \Omega - \{\theta_1\})$

$$(G(\theta_1), G(\Omega - \{\theta_1\})) | \theta_1 \sim \text{Dir}(\alpha H(\theta_1) + 1, \alpha H(\Omega - \{\theta_1\}))$$

Noting that $G(\Omega - \{\theta_1\}) = 1 - G(\theta_1)$, $H(\theta) = 0$, $H(\Omega) = 1$, we get

$$G(\theta_1) \sim \text{Beta}(1, \alpha)$$

Let $\beta_1 = G(\theta_1)$

$$G = \beta_1 \delta_{\theta_1} + (1 - \beta_1) G_1$$

where G_1 is a renormalized probability measure on $\Omega - \{\theta_1\}$

What is G_1 ?

Stick Breaking Construction with Dirichlet Process

Consider arbitrary partition $(\theta_1, A_1, \dots, A_n)$, we have

$$\begin{aligned} & (G(\theta_1), G(A_1), \dots, G(A_n)) \\ &= (\beta_1, (1 - \beta_1)G_1(A_1), \dots, (1 - \beta_1)G_1(A_n)) \\ &\sim (1, H(A_1), \dots, H(A_n)) \end{aligned}$$

Using conditional distribution of Dirichlet distribution,

$$\frac{1}{1 - \beta_1} ((1 - \beta_1)G_1(A_1), \dots, (1 - \beta_1)G_1(A_n)) \sim (H(A_1), \dots, H(A_n))$$

So, we construct a new DP:

$$G_1 \sim DP(\alpha, H)$$

We can continue the breaking process likewise

Stick Breaking Construction with Dirichlet Process

Continue the process

$$G_2 = \beta_1 \delta_{\theta_1} + (1 - \beta_1) G_1$$

$$G_3 = \beta_1 \delta_{\theta_1} + (1 - \beta_1)(\beta_2 + (1 - \beta_2) G_2)$$

\vdots

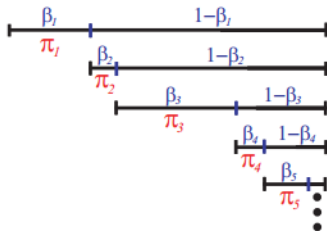
$$G_k = \sum_{k=1}^{\infty} \pi_k \delta_{\theta_k}$$

where

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

$$\pi_k = \beta_k \prod_{i=1}^{k-1} (1 - \beta_i)$$

$$\theta_k \sim H$$



Stick Breaking Construction

The above process is called the *stick breaking construction*

What can we find about DP from this construction?

A sample G from DP is **discrete** with probability 1

If we keep sampling from G , we'll get more and more **repeated** values

That's reasonable for clustering (why?)

Contents

- 1 Introduction
- 2 Gaussian Process
 - Multivariate Gaussian Distribution
 - Definition
 - GP Regression
- 3 Dirichlet Process
 - Definition
 - Stick Breaking Construction
 - Blackwell-MacQueen Urn Scheme
 - Chinese Restaurant Process
 - DP Mixture

Blackwell-MacQueen Urn Scheme

But for clustering, we're actually interested in sampling mixture parameters $\theta_1, \theta_2, \dots$ from Ω , integrating G out.

Blackwell-MacQueen Urn Scheme

But for clustering, we're actually interested in sampling mixture parameters $\theta_1, \theta_2, \dots$ from Ω , integrating G out.

It seems easy. If $\theta_1, \theta_2, \dots, \theta_n$ are independent given G , then we can simply sample each θ from H independently. Right?

Blackwell-MacQueen Urn Scheme

But for clustering, we're actually interested in sampling mixture parameters $\theta_1, \theta_2, \dots$ from Ω , integrating G out.

It seems easy. If $\theta_1, \theta_2, \dots, \theta_n$ are independent given G , then we can simply sample each θ from H independently. Right?

Unfortunately, the answer is No! θ s are no longer independent if we integrate G out

$$p(\theta_1, \theta_2, \dots, \theta_n) = \int_G \prod_{i=1}^n p(\theta_i | G) p(G) dG$$

Blackwell-MacQueen Urn Scheme

Consider two samples θ_1, θ_2 .

$$\begin{aligned} p(\theta_2|\theta_1) &= \int_G p(\theta_2, G|\theta_1) dG \\ &= \int_G p(\theta_2|G)p(G|\theta_1) dG \end{aligned}$$

We have already known the posterior distribution of G given θ_1 ,

$$G|\theta_1 \sim \text{DP}(\alpha + 1, \frac{\alpha H + \delta_{\theta_1}}{\alpha + 1})$$

Then, given θ_1 , the second sample θ_2 can be sampled from the posterior base distribution.

The sampling process becomes quite simple

Blackwell-MacQueen Urn Scheme

First sample:

$$\theta_1 \sim H, \quad G|\theta_1 \sim \text{DP}(\alpha + 1, \frac{\alpha H + \delta_{\theta_1}}{\alpha + 1})$$

Second sample:

$$\theta_2|\theta_1 \sim \frac{\alpha H + \delta_{\theta_1}}{\alpha + 1}, \quad G|\theta_1, \theta_2 \sim \text{DP}(\alpha + 2, \frac{\alpha H + \delta_{\theta_1} + \delta_{\theta_2}}{\alpha + 2})$$

n^{th} sample:

$$\theta_n|\theta_1, \dots, \theta_{n-1} \sim \frac{\alpha H + \sum_{k=1}^{n-1} \delta_{\theta_k}}{\alpha + n - 1}, \quad G|\theta_1, \dots, \theta_n \sim \text{DP}(\alpha + n, \frac{\alpha H + \sum_{k=1}^n \delta_{\theta_k}}{\alpha + n})$$

Blackwell-MacQueen Urn Scheme

One issue, how to sample from $\frac{\alpha H + \sum_{k=1}^{n-1} \delta_{\theta_k}}{\alpha + n - 1}$?

With probability $\propto \alpha$, sample from H

With probability $\propto n - 1$, randomly pick one from $\theta_1, \dots, \theta_{n-1}$

The above sample process is called *Blackwell-MacQueen Urn Scheme*

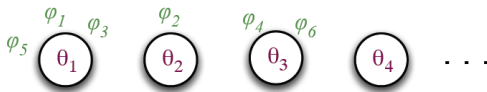
Samples $\theta_1, \dots, \theta_n$ induces a partition over $1, \dots, n$, if we cluster repeated values together

Working with float value is annoying. We can postpone sampling from H – Chinese Restaurant Process

Contents

- 1 Introduction
- 2 Gaussian Process
 - Multivariate Gaussian Distribution
 - Definition
 - GP Regression
- 3 Dirichlet Process
 - Definition
 - Stick Breaking Construction
 - Blackwell-MacQueen Urn Scheme
 - Chinese Restaurant Process
 - DP Mixture

Chinese Restaurant Process



Generating process of **CRP**:

customer 1 sits at the table 1.

the n^{th} customer

With **probability** $\frac{n_k}{\alpha + n - 1}$, sits at existing table k , where n_k is number of customers at table k .

With **probability** $\frac{\alpha}{\alpha + n - 1}$, sits at a new table $K + 1$, where K is number of existing tables.

CRP teases apart clustering property of DP from base distribution H

To get back to Blackwell-MacQueen Urn Scheme, sample θ_k for each table k

CRP and DP

We can prove that the order of customer coming in doesn't affect partition probability

Infinite Exchangeability:

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

De Finetti's Theorem: If (x_1, x_2, \dots) are infinite exchangeable, then there exists a random variable θ , $\forall n$,

$$p(x_1, \dots, x_n) = \int_{\theta} \prod_{i=1}^n p(x_i | \theta) p(\theta) d\theta$$

Here, DP is the distribution of θ

CRP and Clustering

We can view *data points* as customers, *clusters* as tables.

CRP defines a prior distribution on partition (cluster) of data.

Number of clusters can be induced

Combined with a parameterized probability distribution for each cluster, we can get a posterior

That's all we need in cluster settings

Contents

- 1 Introduction
- 2 Gaussian Process
 - Multivariate Gaussian Distribution
 - Definition
 - GP Regression
- 3 Dirichlet Process
 - Definition
 - Stick Breaking Construction
 - Blackwell-MacQueen Urn Scheme
 - Chinese Restaurant Process
 - DP Mixture

Gibbs Sampling DP Mixtures

Go back to [original Gaussian Mixture clustering](#) model. Let's place $DP(\alpha, H)$ as prior of G – DP mixture

Introduce z_i as cluster indicator for each x_i

We use [collapsed Gibbs sampling](#) for inference

$$p(z_i = k | \mathbf{z}_{-i}, \mathbf{x}, \boldsymbol{\nu}, \alpha) \propto p(z_i = k | \mathbf{z}_{-i}, \alpha) p(x_i | \mathbf{x}_{-i}, z_i = k, \mathbf{z}_{-i}, \boldsymbol{\nu})$$

where

$$p(z_i = k | \mathbf{z}_{-i}, \alpha) = \frac{n_k}{\alpha + n - 1}$$

for existing component k among \mathbf{z}_{-i} , and

$$p(z_i = K + 1 | \mathbf{z}_{-i}, \alpha) = \frac{\alpha}{\alpha + n - 1}$$

for new component $K + 1$, $\forall j \neq i, z_j \neq K + 1$

Gibbs Sampling for DP Mixtures

What is $p(x_i | \mathbf{x}_{-i}, z_i = k, \mathbf{z}_{-i}, \nu)$?

For new component $K + 1$, x_i is conditional independent from other data points

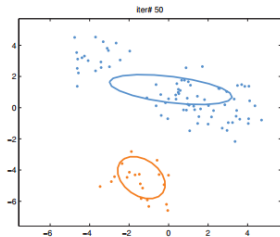
$$p(x_i | \mathbf{x}_{-i}, z_i = K + 1, \mathbf{z}_{-i}, \nu) = \int_{\phi} p(x_i | \phi) p(\phi | \nu) d\phi$$

For existing component k , x_i is conditional independent from data points in different cluster

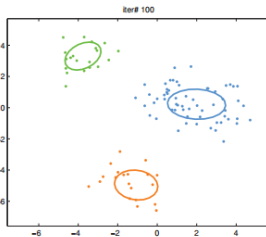
$$p(x_i | \mathbf{x}_{-i}, z_i = k, \mathbf{z}_{-i}, \nu) = \frac{p(x_i, \mathbf{x}_{-i}^{(k)} | \nu)}{p(\mathbf{x}_{-i}^{(k)} | \nu)}$$

It's model specific.

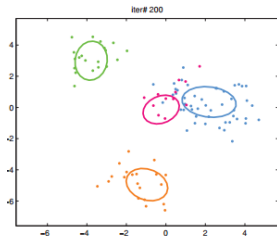
Visualization



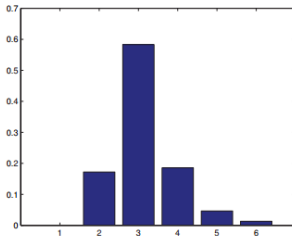
(a)



(b)



(c)



(d)

Other Methods

Hierarchical Dirichlet Processes

Infinite Hidden Markov Models

Polya Trees

Dirichlet (Diffusion) Trees

Dirichlet Forest

India Buffet Processes

(Zoubin Ghahramani, Non-parametric Bayesian Methods. 2005)

Outline

- 1 Introduction
- 2 Gaussian Process
 - Multivariate Gaussian Distribution
 - Definition
 - GP Regression
- 3 Dirichlet Process
 - Definition
 - Stick Breaking Construction
 - Blackwell-MacQueen Urn Scheme
 - Chinese Restaurant Process
 - DP Mixture

Thanks.

HP: <http://keg.cs.tsinghua.edu.cn/jietang/>

Email: jietang@tsinghua.edu.cn