

Machine Learning

Lecture 11 – Bayesian nonparametrics



Fredrik Lindsten¹, Thomas Schön²

¹Department of Engineering
University of Cambridge.
[www-sigproc.eng.cam.ac.uk/
Main/FSML2](http://www-sigproc.eng.cam.ac.uk/Main/FSML2)

²Department of Information Technology
Uppsala University.
user.it.uu.se/~thosc112

Contents – lecture 11

2(21)

1. Summary of lecture 10
2. Parametric vs. nonparametric models
3. Mixture models – the standard example
4. The Dirichlet process
 - Stick-breaking
 - The Blackwell-MacQueen and Chinese restaurant processes
 - Dirichlet process mixture models
5. Beyond DP mixture models

Summary of lecture 10 (I/II)

3(21)

The **idea** underlying Monte Carlo is to generate samples $\{z^i\}_{i=1}^M$ according to some proposal distribution $q(z)$ and possibly compute a weight for each sample, resulting in an **empirical estimate**

$$\hat{\pi}(z) = \sum_{i=1}^M w^i \delta_{z^i}(z)$$

of the target distribution $\pi(z)$. This allows for approximations of general integrals according to

$$\mathbb{E}[g(z)] = \int g(z) \pi(z) dz \approx \sum_{i=1}^M w^i g(z^i)$$

Summary of lecture 10 (II/II)

4(21)

Two “basic Monte Carlo samplers” were introduced; rejection sampling and importance sampling.

A **Markov chain Monte Carlo (MCMC)** method allows us to generate samples from an arbitrary target distribution $\pi(z)$ by simulating a Markov chain whose stationary distribution is $\pi(z)$.

A **Markov chain** $\{z^m\}_{m \geq 1}$ is a stochastic process specified by

1. Initial distribution: $z^1 \sim \mu_1(z^1)$
2. Transition kernel: $z^{m+1} \mid z^m \sim K(z^{m+1} \mid z^m)$

Two **constructive** ways of building Markov chains with a particular user-defined stationary distribution were introduced:

1. Metropolis Hastings (MH) sampler
2. Gibbs sampler

Parametric model,

$$Y \sim p(Y | \theta)$$

for some finite dimensional parameter θ .

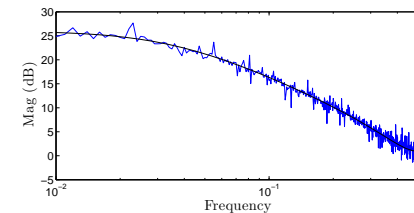
1. Complexity/flexibility of model \approx dimension of θ .
2. Can lead to over- or underfitting when there is a mismatch between the model complexity and the amount of available data!
3. Order selection is often a hard problem.

Nonparametric model – flexible model in which the complexity increases with the amount of data.

1. Attempts to avoid order selection.
2. The number of “parameters” increases with the number of data points.

Ex. (Empirical transfer function estimate, ETFE)

$$\hat{G}_N(e^{i\omega}) = \frac{Y_N(\omega)}{U_N(\omega)}$$



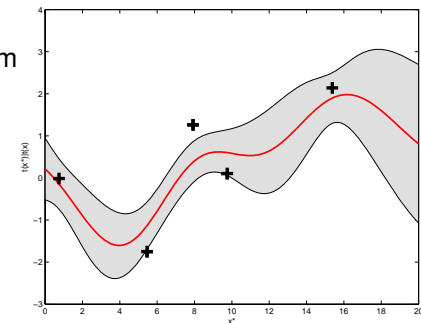
- Bayesian parametric model = latent random variables (parameters).
- Bayesian nonparametric model = **latent stochastic process**.

Recall Gaussian processes from lecture 5,

$$f(\cdot) \sim \text{GP}(m(x), k(x, x')),$$

$$y_n = f(x_n) + e_n,$$

for $n = 1, \dots, N$.



Many possibilities, depending on what we want to capture with the model,

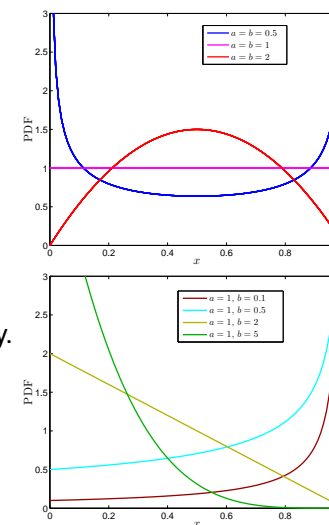
1. Gaussian process
2. Dirichlet process, Chinese restaurant process
3. Pitman-Yor process
4. Beta process, Indian buffet process
5. ...

Beta distribution:

$$x \sim \text{Be}(a, b)$$

Parameters: $a, b > 0$.

- Support: $0 \leq x \leq 1$.
- Often used as prior for a probability.
- Conjugate prior for Bernoulli, binomial and geometric distr.

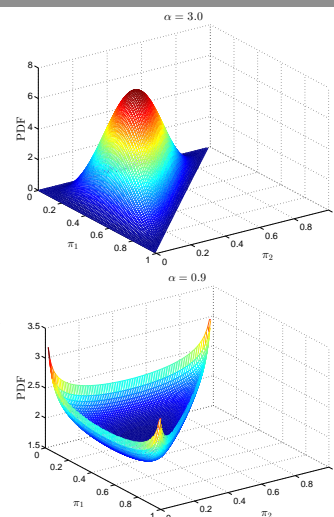


Dirichlet distribution:

$$\pi \sim \text{Dir}(\alpha_1, \dots, \alpha_K)$$

Parameters: $\alpha_k > 0$.

- Support: $0 \leq \pi_k \leq 1$ and $\sum_{k=1}^K \pi_k = 1$.
- A draw $\pi = (\pi_1, \dots, \pi_M)$ can be interpreted as a discrete probability distribution.
- The Dirichlet distribution is a “distribution over distributions”!
- Conjugate prior for discrete and multinomial distributions.



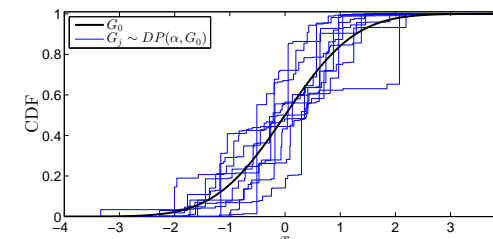
Dirichlet process:

$$G \sim \text{DP}(\alpha, G_0),$$

with base distribution G_0 and concentration parameter α .

A draw from the DP is a discrete probability distribution!

- $\mathbb{E}[G] = G_0$
- $\mathbb{V}[G] \propto (1 + \alpha)^{-1}$



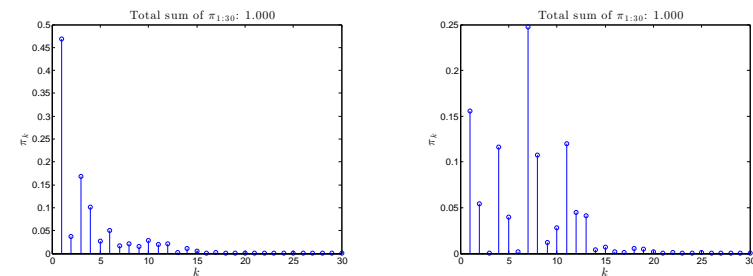
Can we make this construction nonparametric? Define an **infinite mixture model** by letting $K \rightarrow \infty$, i.e. we get $G = \sum_{k=1}^{\infty} \pi_k \delta_{\phi_k}$, where,

$$\begin{aligned} \phi_k &\stackrel{\text{i.i.d.}}{\sim} G_0, \\ \pi &\sim \text{Dir}(\alpha/K, \dots, \alpha/K), \quad K \rightarrow \infty. \end{aligned}$$

1. Will π have a proper distribution as $K \rightarrow \infty$?
2. Will $\sum_{k=1}^{\infty} \pi_k = 1$?
3. Will the model have clustering properties?

A better way – use a constructive definition.

Illustration of the stick-breaking construction of the Dirichlet process



Explicit representations (like stick-breaking) are of great practical importance for the derivation of algorithms.
(MATLAB code available on web site)

Assume that $G \sim \text{DP}(\alpha, G_0)$ and $\theta_1 \sim G$.

What can be said about the *posterior* distribution “ $G \mid \theta_1$ ”?

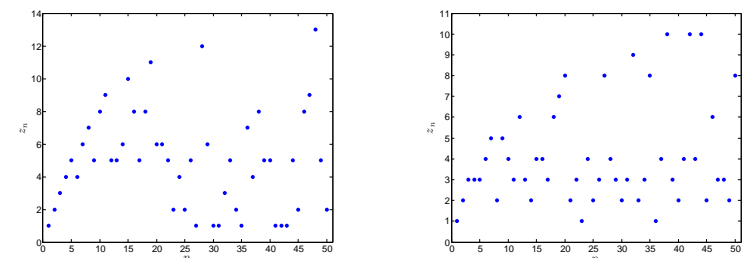
Discrete-Dirichlet-conjugacy carries over to DP!

$$G \mid \theta_1 \sim \text{DP} \left(\alpha + 1, \frac{\alpha G_0 + \delta_{\theta_1}}{\alpha + 1} \right).$$

Iterating the posterior update we get,

$$G \mid \theta_1, \dots, \theta_n \sim \text{DP} \left(\alpha + n, \frac{\alpha G_0 + \sum_{i=1}^n \delta_{\theta_i}}{\alpha + n} \right).$$

Under the assumption of a DP mixture, the CRP directly describes the prior over how data is clustered.



This implicit representation is practically useful as they can lead to simple and efficient inference algorithms.

(MATLAB code available on web site)

Note that

$$p(z_{1:n}) = \prod_{i=1}^n p(z_i | z_{1:i-1}).$$

We can thus write the DP mixture model in terms of the CRP,

$$z_{n+1} | z_1, \dots, z_n = \begin{cases} k & \text{w.p. } \frac{m_k}{\alpha + n}, \\ K+1 & \text{w.p. } \frac{\alpha}{\alpha + n}, \end{cases}$$

$$\phi_k \stackrel{\text{i.i.d.}}{\sim} G_0, \quad k = 1, 2, \dots$$

$$y_n | \{z_n = k\}, \phi_k \sim p(y_n | \phi_k).$$

Similar to a finite mixture model with latent variables, but now the latent variables are given by the CRP.

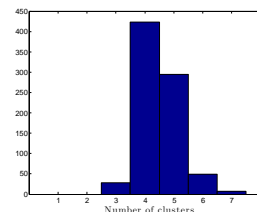
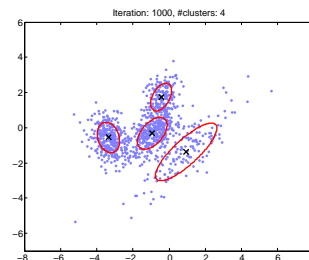
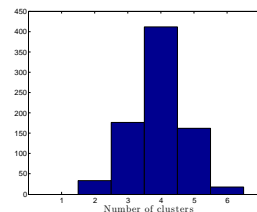
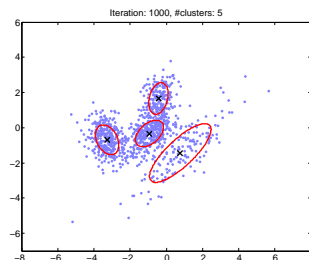
Inference for DP mixture models

- Different representations (stick-breaking, Blackwell-MacQueen, CRP) give rise to different algorithms.
- Different inference tools – MCMC (Gibbs/split-merge), VB, Particle filters, Maximization-Expectation, ...

ex) Gibbs sampler using CRP (Neal, 2000).

Given $\{y_n\}_{n=1}^N$, iterate:

- For $n = 1, \dots, N$ draw: $z_n | z_{-n}, y_n, \{\phi_k\}_{k=1}^K$;
- For $k \in \{z_1, \dots, z_n\}$ draw: $\phi_k | \{\text{all } y_n \text{ s.t. } z_n = k\}$.



(MATLAB code available on web site)

S. J. Gershman and D. M. Blei. **A tutorial on Bayesian nonparametric models.** *Journal of Mathematical Psychology* (56):1-12, 2012.

E. B. Sudderth. **Graphical Models for Visual Object Recognition and Tracking.** PhD thesis, MIT, 2006.

R. M. Neal, **Markov chain sampling methods for Dirichlet process mixture models,** *Journal of computational and graphical statistics*, vol. 9, no. 2, pp. 249-265, 2000.

- Nonparametric models allow the complexity to increase with the amount of data
- Bayesian nonparametrics = latent stochastic processes
- Dirichlet process,
 - A draw from the Dirichlet process is a (random) discrete probability distribution
 - Dirichlet process mixture model for clustering
 - Hierarchical DPs can be used to construct an “infinite” HMM
- Indian buffet process for feature models