Stochastic Processes, Kernel Regression, Infinite Mixture Models



Gabriel Huang

(TA for Simon Lacoste-Julien)

IFT 6269 : Probabilistic Graphical Models - Fall 2018

Stochastic Process

Random Function

Today

- Motivate Gaussian and Dirichlet distribution in Bayesian Framework.
- Kolmogorov's extension theorem.
- Define Gaussian Process and Dirichlet Process from finite-dimensional marginals.
- Gaussian Process:
 - Motivating Applications: Kriging, Hyperparameter optimization.
 - Properties: Conditioning/Posterior distribution.
 - Demo.
- Dirichlet Process:
 - Motivating Application: Clustering with unknown number of clusters.
 - Construction: stick-breaking, Polya urn, Chinese Restaurant Process.
 - De Finetti theorem.
 - How to use.
 - Demo.

Disclaimer

I will be skipping the more theoretical building blocks of stochastic processes (e.g. measure theory) in order to be able to cover more material.

Recall some distributions

Gaussian Distribution

Samples X in \mathbb{R}^d .

Dirichlet Distribution

Samples π in Simplex Δ_{d-1} verifies $\pi_1 + \dots + \pi_d = 1$





Why Gaussian and Dirichlet? They are often used as priors

Bayesians like to use those distributions as priors over model parameters p(x)

Why?

Because they are very convenient to represent/update.

Conjugate Priors

 $p(\theta|x) \propto p(x|\theta)p(\theta)$ Likelihood Posterior Prior model

Conjugate Prior means: Posterior in same family as prior

$$p(\theta|x) \propto p(x|\theta)p(\theta)$$

$$Posterior \qquad Prior \\ \theta|x \sim Gaussian(\mu', \Sigma') \qquad \theta|x \sim Gaussian(\mu, \Sigma)$$

$$Likelihood \\ x|\theta \sim Gaussian(\theta, \Sigma_L)$$
Gaussian is conjugate prior for

Gaussian likelihood model.



Dirichlet is conjugate prior for **Categorical/Multinoulli** likelihood model.

So taking the posterior is simply a matter of updating the parameters of the prior.

Back to Gaussian and Dirichlet

Gaussian Distribution

Samples X in \mathbb{R}^d .

Dirichlet Distribution

Samples π in Simplex Δ_{d-1} verifies $\pi_1 + \dots + \pi_d = 1$





Gaussian and Dirichlet are indexed with a <u>finite set of integers</u> {1, ..., d}

They are random vectors.

$$(X_1, X_2, \dots, X_d)$$

 $(\pi_1, \pi_2, \dots, \pi_d)$

Can we index random variables with <u>infinite</u> sets as well?

In other words, define **random functions**.

Defining stochastic processes from their marginals.

Suppose we want to define a random function (stochastic process)

$X: t \in T \to \mathbf{R},$

where T is an infinite set of indices.

Imagine a joint distribution over all the (X_t) .

Kolmogorov Extension Theorem

informal statement

Assume that for any $k \ge 1$, and every finite subset of indices $(t_1, t_2, ..., t_k)$, we can define a marginal probability (finite-dimensional distribution) $p_{t_1,t_2,...,t_k}(X_{t_1}, X_{t_2}, ..., X_{t_k})$ Then, if all marginal probabilities agree, there **exists** a

unique stochastic process $X: t \in T \rightarrow \mathbf{R}$ which satisfies the given marginals.

So Kolmogorov's extension theorem gives us a way to <u>implicitly define stochastic</u>

processes.

(However it does not tell us how to <u>construct</u> them.)

Defining Gaussian Process from finite-dimensional marginals.

Characterizing Gaussian Process

Samples $X \sim GP(\mu, \Sigma)$ of a Gaussian Process are random functions $X: T \rightarrow \mathbf{R}$ defined on the domain T (such as time $T = \mathbf{R}$, or vectors $T = \mathbf{R}^d$).

We can also see them as an infinite collection $(X_t)_{t \in T}$ indexed by T.

Parameters are the **Mean** function $\mu(t)$ and **Covariance** function $\Sigma(t, t')$.

For any $t_1, t_2, ..., t_k \in T$ we define the following finite-dimensional distributions $p(X_{t_1}, X_{t_2}, ..., X_{t_k})$.

$$X_{t_1}, X_{t_2}, \dots, X_{t_k} \sim \mathcal{N}((\boldsymbol{\mu}(\boldsymbol{t_i})_i, (\boldsymbol{\Sigma}(\boldsymbol{t_i}, \boldsymbol{t_j}))_{i,j}))$$

Since they are consistent with each other, Kolmogorov's extension theorem states that they define a unique stochastic process, we will call **Gaussian Process**:

$$X \sim GP(\mu, \Sigma)$$

Characterizing Gaussian Process

Some properties are immediate consequences of definition:

• $\mathbf{E}[X_t] = \mu(t)$

•
$$\mathbf{Cov}(X_t, X_{t'}) = E[(X_t - \mu(t))(X_{t'} - \mu(t'))^T] = \Sigma(t, t')$$

Any linear combination of distinct dimensions is still a Gaussian:

$$\sum_{i=1}^{k} \alpha_i * X_{t_i} \sim \mathcal{N}(\cdot, \cdot)$$

Characterizing Gaussian Process

Some properties are immediate consequences of definition:

- <u>Stationarity</u>: $\Sigma(t, t') = \Sigma(t t')$ does not depend on the positions
- <u>Continuity:</u> $\lim_{t' \to t} \Sigma(t, t') = \Sigma(t, t)$
- Any linear combination is still a Gaussian:

$$\sum_{i=1}^{k} \alpha_i * X_{t_i} \sim \mathcal{N}(\cdot, \cdot)$$

Example Samples



Posteriors of Gaussian Process. How to use them for regression?

Interactive Demo

need a volunteer

http://chifeng.scripts.mit.edu/stuff/gp-demo/

Gaussian processes are very useful for doing regression on an unknown function f: y = f(x).

Say we don't know anything about that function, except the fact that it is smooth.

Before observing any data, we represent our belief on the unknown function f with the following prior:



WARNING: Change of notation! x is now the index and f(x) is the random function Controls smoothness (bandwidth/length-scale)

Now, assume we observe a training set

$$X_n = (x_1, x_2, \dots, x_n), y_n = (y_1, y_2, \dots, y_n)$$

and we want to predict the value $y^* = f(x^*)$ associated with a new test point x^* .

One way to do that is to compute the **posterior** $f | X_n, y_n$ after observing the evidence (training set).

Bayes' Rule $p(f|X_n, y_n) \propto p(y_n|f, X_n) p(f)$

• Gaussian Process Prior:

$$p(f) = GP(\mu(x), \Sigma(x, x'))$$

• Gaussian Likelihood:

$$p(y_n|f, X_n) = \mathcal{N}(f(X), \epsilon^2 I_n)$$

• -> Gaussian Process Posterior:

 $p(f|X_n, y_n) = GP(\mu'(x), \Sigma'(x, x'))$ for some $\mu'(x), \Sigma'(x, x')$.

<u>Remember</u>: Gaussian Process is conjugate prior for Gaussian likelihood model.

Bayes' Rule $p(f|X_n, y_n) \propto p(y_n|f, X_n) p(f)$

• Gaussian Process Prior:

$$p(f) = GP(\mu(x), \Sigma(x, x'))$$

• Dirac Likelihood: $(\epsilon \rightarrow 0)$

$$p(y_n|f,X_n) = \delta(y_n - f(X_n))$$

that is, y_n is now deterministic after observing f, X_n .

$$y_n = f(X_n)$$

• -> Gaussian Process Posterior:

$$p(f|X_n, y_n) = GP(\mu'(x), \Sigma'(x, x'))$$
for some $\mu'(x), \Sigma'(x, x')$.

The problem is there is no easy way to represent the parameters of the posterior $\mu(x), \Sigma(x, x')$ efficiently.

Instead of computing the full posterior f, we will just evaluate the posterior at one point $y^* = f(x^*)$.

We want:
$$p(y^*|X_n, y_n, x^*)$$

We want: $p(y^*|X_n, y_n, x^*)$

The finite-dimensional marginals of the Gaussian process give that:



Theorem: For a Gaussian vector with distribution

$$\frac{X_1}{X_2} \sim \mathcal{N}\left(\frac{\mu_1}{\mu_2}, \frac{\Sigma_{1,1}}{\Sigma_{2,1}}, \frac{\Sigma_{1,2}}{\Sigma_{2,2}}\right)$$

the conditional distribution $X_2 | X_1$ is given by

$$X_2 \sim \mathcal{N}(\mu_2 + \Sigma_{2,1}\Sigma_{1,1}^{-1}(x_1 - \mu_1)),$$

This Theorem will be useful for the Kalman filter, later on ...

$$\Sigma_{2,2} - \Sigma_{2,1} \Sigma_{1,1}^{-1} \Sigma_{1,2}$$
)

[Schur's complement] ³⁵

Applying the previous theorem gives us the posterior y^* .

$$\mu' = \mu(x^*) + \Sigma(x^*, X_n) \Sigma(X_n, X_n)^{-1}(X_n - \mu(X_n))$$


Active Learning with Gaussian Process.

Active Learning

Active Learning is iterative process:

- Generate a question x^* .
- Query the world with the question (by acting, can be costly)
- Obtain an answer $y^* = f(x^*)$.
- Improve model by learning from the answer.
- Repeat.

Active Learning

Gaussian process is good for cases where it is expensive to evaluate $y^* = f(x^*)$.

- Kriging. y* is the amount of natural resource, x* is new 2D/3D location to dig. Every evaluation is mining and can cost millions.
- Hyperparameter optimization (Bayesian optimization).
 y* is the validation loss, x* is set of hyperparameters to test. Every evaluation is running an experiment and can take hours.

Back to the demo

(Talk about utility function)

http://chifeng.scripts.mit.edu/stuff/gp-demo/

Formal equivalence with Kernelized Linear Regression. [blackboard if time]

<u>Rasmussen & Williams (2006)</u> <u>http://www.gaussianprocess.org/gpml/chapters/RW2.pdf</u>

Dirichlet Processes. Stick Breaking Construction

$$\pi = (\pi_1, \pi_2, ...) \sim GEM(\alpha)$$
scalar weights
sum up to 1
$$G = \sum_{k=1}^{+\infty} \pi_k \delta_{\theta_k}$$

$$G = \sum_{k=1}^{+\infty} \pi_k \delta_{\theta_k}$$
Diracs concentrate
probability mass π_k
at θ_k

G is a **random probability measure**:

- **random**: both π and θ are random
- probability measure: it is a convex combination of Diracs, which are probability measures

Dirichlet Process

Consider Gaussian G₀



Two independent samples G from $DP(\alpha, G_0)$



Each sample G is a probability distribution (e.g. over parameters) and can be written as a mixture of diracs.

 $+\infty$

k=1

 $\pi_k o$

Measuring is counting $G(A) = \sum_{k=1}^{+\infty} \pi_k * 1\{\theta_k \in A\}$



For a fixed subset A, notice how G(A) is random. In fact even the π_k change value for each sample.

To generate a **finite** sequence of (mixture) weights $\pi = \pi_1, \pi_2, ..., \pi_k$ that sum up to 1, we can use the **Dirichlet** distribution $\pi \sim Dirichlet(\alpha_1, ..., \alpha_k)$

How to generate an **infinite** sequence of (mixture) weights $\pi = \pi_1, \pi_2, \dots$ which sum up to 1? we can use **stick-breaking** $\pi \sim GEM(1, \alpha)$

Beta Distribution



 $\alpha, \beta \rightarrow +\infty$ gives peaked distribution around $\alpha/(\alpha + \beta)$



Defining Dirichlet Process from finite-dimensional marginals.

Dirichlet Process

Samples $G \sim DP(\alpha, G_0)$ of a Dirichlet Process are themselves probability measures (i.e. distributions) over a measurable space (Ω, \mathcal{F}) .

 $G: \mathcal{F} \to \mathbf{R}_+$

which associate a probability to every measurable subset $A \in \mathcal{F}$.

Note: \mathcal{F} is the set of all measurable subsets $A \subseteq \Omega$.

Parameters are the **base probability distribution** G_0 (over Ω) and the parameter $\alpha > 0$.



Kolmogorov Consistency Construction

For any $k \ge 0$, consider any partition $A_1, A_2, ..., A_k$ of the space Ω . We define the following finite-dimensional distributions

$$G(A_1), \dots, G(A_k) \sim Dirichlet(\alpha * G_0(A_1), \dots, \alpha * G_0(A_k))$$

Since they can be proved* to be consistent with each other, Kolmogorov's extension theorem states that they define a unique stochastic process, we will call **Dirichlet Process**:

$$G \sim DP(\alpha, G_0)$$

Here A_1, A_2, A_3 is a partition of the parameter space Ω . Assume $\alpha = 10, G_0 = \mathcal{N}(0, I_2)$. π_1 Draw two distributions $G_1, G_2 \sim_{iid} DP(\alpha, G_0)$. A_3 First sample π_5 $G_1(A_1) = \pi_5 + \pi_6 + \pi_8$ π_6 $G_1(A_2) = \pi_2 + \pi_4$ $G_1(A_3) = \pi_1$ Second sample $G_2(A_1) = \pi_3 + \pi_4 + \pi_5$ $G_2(A_2) = \pi_2$ π_2 $G_2(A_3) = \pi_1$ *Probability masses for base distribution (deterministic)* $G_0(A_1) = \mathcal{N}(0, I_2)(A_1) = 0.8$ A_3 A_2 $G_0(A_2) = \mathcal{N}(0, I_2)(A_2) = 0.2$ π_5 $G_0(A_3) = \mathcal{N}(0, I_2)(A_3) = 0.2$ Then we have that $G(A_1), G(A_2), G(A_3) \sim \text{Dirichlet}(8,2,2)$ 53

All constructions match.

It can be shown that Stick-Breaking and Kolmogorov consistency definitions match. Defining Dirichlet Process from Chinese Restaurant Process /BlackWell-McQueen Urn.

Infinity of Tables

 $\bullet \bullet \bullet$



$G_0 = Uniform(\{Fish, Pork, Tofu\})$





Infinity of Tables

- Customer 1 arrives.
- Takes any free table.
- Sample a dish $\theta_1 \sim G_0$ -> Tofu
- state={{1}}, n=1 customers





Infinity of Tables

- Customer 2 arrives.
- P(new table) $\propto \alpha$
- $P(table \{1\}) \propto |\{1\}| = 1$
- Decides to sit at {1}
- Share dish: $\theta_2 = \theta_1 = Tofu$
- {{1,2}}, n=2 customers





Infinity of Tables



- Customer 3 arrives.
- P(new table) $\propto \alpha$
- $P(table \{1,2\}) \propto |\{1,2\}| = 2$
- Decides to sit at new table
- Sample a dish $\theta_3 \sim G_0$ -> Pork
- {{1,2},{3}}, n=3 customers





Infinity of Tables

• • •



- Customer 4 arrives.
- P(new table) $\propto \alpha$
- $P(table \{1,2\}) \propto |\{1,2\}| = 2$
- $P(table \{3\}) \propto |\{3\}| = 1$
- Share dish, $\theta_4 = \theta_1 =$ Tofu
- {{1,2,4},{3}}, n=4 customers





Infinity of Tables

 $\bullet \bullet \bullet$

- Customer 5 arrives.
- P(new table) $\propto \alpha$

Fish

- P(table $\{1,2,4\} \propto |\{1,2,4\}| = 3$
- $P(table \{3\}) \propto |\{3\}| = 1$
- Pick new table
- Sample new dish $\theta_5 = Fish$
- {{1,2,4},{3},{5}}, n=5 customers



Pork

2



Infinity of Tables

- Customer 6 arrives.
- P(new table) $\propto \alpha$
- P(table $\{1,2,4\}) \propto |\{1,2,4\}| = 3$
- $P(table \{3\}) \propto |\{3\}| = 1$
- $P(table \{5\}) \propto |\{5\}| = 1$
- Pick table {1,2,4}
- Share dish $\theta_6 = \theta_1 = Tofu$
- {{1,2,4,6},{3},{5}}, n=6 customers

We can look at the sequence of dishes

$$\begin{array}{l} \theta_1 = Tofu \\ \theta_2 = Tofu \\ \theta_3 = Pork \\ \theta_4 = Tofu \\ \theta_5 = Fish \\ \theta_6 = Tofu \end{array}$$

It can be shown that the distribution of $(\theta_t)_t$ is **exchangeable**. That is: $p(\theta_1 = u_1, \theta_2 = u_2, ...) = p(\theta_1 = u_{\sigma(1)}, \theta_2 = u_{\sigma(2)}, ...)$

The **order** in which the customers arrive is actually **not important**.

De Finetti's Theorem

- Suppose that we agree that if our data are reordered, it doesn't matter
 - this is generally **not** an assertion of "independent and identically distributed"; rather, it is an assertion of "exchangeability"
- *Exchangeability*: the joint probability distribution underlying the data is invariant to permutation

Theorem (De Finetti, 1935). If $(x_1, x_2, ...)$ are infinitely exchangeable, then the joint probability $p(x_1, x_2, ..., x_N)$ has a representation as a mixture:

$$p(x_1, x_2, \dots, x_N) = \int \left(\prod_{i=1}^N p(x_i \mid \theta)\right) dP(\theta)$$

M. I. Jordan NIPS 2017 Tutorial http://faculty.dbmi.pitt. edu/day/Bioinf2132advanced-Bayes-and-R/Bioinf2132documents-2017/2017-11-30/nipstutorial05.pdf

for some random variable θ .

• I.e., if you assert exchangeability, it is reasonable to act as if there is an underlying parameter, there is a prior on that parameter, and the data are conditionally IID given that parameter

De Finetti's Theorem

Applied to the CRP, it means there **exist** a **unique**^{*} random variable G, such that all θ become independent conditionally to G.

We can show that $G \sim DP(\alpha, G_0)!$

Here: $G_0 = \frac{1}{3}\delta_{fish} + \frac{1}{3}\delta_{pork} + \frac{1}{3}\delta_{tofu}$, α is the same (\propto new table) Let $\pi = (\pi_1, \pi_2, ...) \sim \text{GEM}(\alpha)$ stick-breaking. Sample $\theta = (\theta_{k=1}, \theta_{k=2}, ...) \sim_{iid} G_0$ Now we can form our random measure $\mathbf{G} = \sum_{k=1}^{+\infty} \pi_k * \delta_{\theta_k}$. And we sample $\theta_{i=1}, \theta_{i=2}, ... \sim_{iid} \mathbf{G}$

*Unique in distribution.

 θ_i is the parameter for data point i (customer i) θ_k is the parameter for component k (table k)

Blackwell-McQueen Urn Polya Urn

Same process, different story.

Each dish is a set of unique ball colors.

Each customer is a successive draw.

Using Dirichlet Process for infinite mixture models.

Infinity of **Components**

 $\bullet \bullet \bullet$

 $G_0 = \mathcal{N}(\mu_0, \Sigma_0)$





Infinity of Tables

- Sample parameter for data point 1.
- Takes any free table.
- Sample a parameter $\theta_1 \sim G_0$
- state={{1}}, n=1 customers



Infinity of Tables

- Sample parameter for data point 2.
- P(new table) $\propto \alpha$
- $P(table \{1\}) \propto |\{1\}| = 1$
- Decides to sit at {1}
- Share dish: $\theta_2 = \theta_1 = Tofu$
- {{1,2}}, n=2 customers





Infinity of Tables



- Sample parameter for data point 3.
- P(new table) $\propto \alpha$
- $P(table \{1,2\}) \propto |\{1,2\}| = 2$
- Decides to sit at new table
- Sample a dish $\theta_3 \sim G_0$ -> Pork
- {{1,2},{3}}, n=3 customers
Chinese Restaurant Process (CRP)



 $\boldsymbol{\theta}_{k=1} = (\mathbf{0}, \mathbf{1})$



Infinity of Tables

 $\bullet \bullet \bullet$

- Customer 4 arrives.
- P(new table) $\propto \alpha$
- $P(table \{1,2\}) \propto |\{1,2\}| = 2$
- $P(table \{3\}) \propto |\{3\}| = 1$
- Share dish, $\theta_4 = \theta_1 =$ Tofu
- {{1,2,4},{3}}, n=4 customers

Chinese Restaurant Process (CRP)



 $\boldsymbol{\theta}_{k=1} = (0, 1)$



Infinity of Tables

 \bullet \bullet \bullet

- Customer 5 arrives.
- P(new table) $\propto \alpha$
- P(table $\{1,2,4\}) \propto |\{1,2,4\}| = 3$
- $P(table \{3\}) \propto |\{3\}| = 1$
- Pick new table
- Sample new dish $\theta_5 = Fish$
- {{1,2,4},{3},{5}}, n=5 customers

Chinese Restaurant Process (CRP)



 $\boldsymbol{\theta}_{k=2} = (\mathbf{0}, \mathbf{1})$



Infinity of Tables

 $\bullet \bullet \bullet$

- Customer 6 arrives.
- P(new table) $\propto \alpha$
- P(table $\{1,2,4\} \propto |\{1,2,4\}| = 3$
- $P(table \{3\}) \propto |\{3\}| = 1$
- $P(table \{5\}) \propto |\{5\}| = 1$
- Pick table {1,2,4}
- Share dish $\theta_6 = \theta_1 = Tofu$
- {{1,2,4,6},{3},{5}}, n=6 customers



We can describe a generative process of data points.

(but first let's recall the generative process for GMM)

