

# STATS 225: Bayesian Analysis

## Gaussian Process Models

Babak Shahbaba

Department of Statistics, UCI

Winter, 2015

# Introduction

- A Gaussian process (GP) on the real line is a random real-valued function  $y(t)$ , which is completely determined by its mean function  $\mathbb{E}y(s)$  and covariance function (kernel)  $C_{st} = \text{Cov}(y(s), y(t))$ .
- A finite sample  $(y(t_1), \dots, y(t_n))$  has a multivariate Gaussian distribution with mean  $(\mathbb{E}y(t_1), \dots, \mathbb{E}y(t_n))$ , and covariance matrix  $(C_{t_i t_j})$ .
- Note that we are limited to kernels providing positive semi-definite covariance matrices.
- In this lecture, we discuss Gaussian process models for regression and classification, so the process is indexed by a set of predictors  $x$ .
- In this case, Gaussian process is used as a distribution over functions  $y(x)$ .
- We usually add an extra parameter to account for observation noise.

# Gaussian process models

- To introduce this concept, we start with a simple linear regression model.
- Recall that we presented a linear regression model as

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots \beta_p x_{ip} + \varepsilon_i$$

- Using normal priors (with mean zero, and in general, different variances) for  $\beta$ 's

$$\beta_j | \sigma_j \sim N(0, \sigma_j^2) \quad j = 0, \dots, p$$

# Gaussian process models

- In prior,  $\beta$  has a  $(p + 1)$  dimensional multivariate normal distribution

$$\beta | \Sigma_\beta \sim N(0, \Sigma_\beta)$$

- $\varepsilon$  also has an  $n$  dimensional multivariate normal distribution

$$\varepsilon | \Sigma_\varepsilon \sim N(0, \Sigma_\varepsilon)$$

- To obtain the distribution of  $y$  we multiply  $\beta$  by the matrix  $x$  and add  $\varepsilon$  to it.
- Based on the properties of multivariate normal distribution, the resulting distribution would still be multivariate normal  $N(0, C)$  where

$$C = x \Sigma_\beta x^T + \Sigma_\varepsilon$$

# Gaussian process models

- This gives us the prior distribution on the function  $y(x)$ .
- Since any finite subset of  $y(x)$  (e.g., for the  $n$  observed cases) would have a Gaussian distribution, the prior distribution on  $y(x)$  is a *Gaussian process*.
- As mentioned above, analogous to Gaussian distributions, Gaussian processes are defined by their mean (here, the mean is 0 in prior) and covariance function  $C$ .
- For the above linear model, the elements of  $C$  are

$$C_{ij} = \text{Cov}(y_i, y_j) = \sigma_0^2 + \sum_{u=1}^p x_{iu}x_{ju}\sigma_u^2 + \delta_{ij}\sigma_\epsilon^2$$

where  $\delta_{ij}$  is equal to 1 if  $i = j$ , and 0 otherwise.

# Gaussian process models

- Setting up the model this way, we are putting the prior directly on the relationship between  $x$  and  $y$  as opposed to on some parameters that represent this relationship (i.e., we cut out the middleman).
- This is specially useful if our objective is to predict future cases as opposed to making inference about the relationship between  $x$  and  $y$ .
- Note that the prior here is implicit and reflects our choice of the functional form.
- In the above example, we are assuming the relationship is linear. In general, we could use other covariance functions,  $C$ , to create nonlinear relationship.

# Gaussian process for nonlinear regression

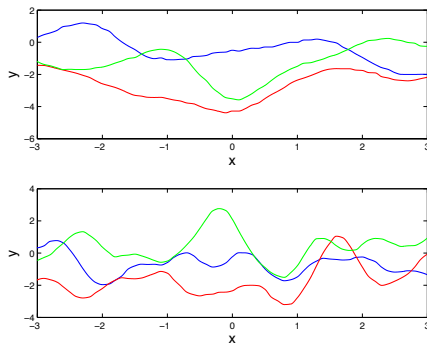
- For example, the following covariance function is very useful and includes a wide range of smooth nonlinear functions:

$$\text{Cov}(y_i, y_j) = \lambda^2 + \eta^2 \exp \left( - \sum_{u=1}^p \rho_u^2 (x_{iu} - x_{ju})^2 \right) + \delta_{ij} \sigma_\epsilon^2$$

- The constant part is used to make sure the model fit functions where the mean of  $y$  is not zero (the  $x$  matrix does not have a vector of 1's anymore). However, it is better to center  $y$  before analysis so we don't have to use a large constant.
- There is one  $\rho$  for each predictor.
- The noise parameter,  $\sigma_\epsilon^2$  (also called *jitter*), accounts for random variations and is essential to improve computation.

# The effect of parameters in the covariance function

- By using different  $\eta$ ,  $\rho$ 's,  $\lambda$  and  $\sigma_\varepsilon$ , we can generate a large variety of functions.



The top panel shows samples based on  $\eta = 1$ ,  $\rho = 1$ ,  $\lambda = 1$ , and  $\sigma_\varepsilon = 0.01$ . The bottom panel is based on the same priors except we set  $\rho = 2$ .



- As mentioned above, using a Gaussian process prior is especially useful if our goal is predicting future cases for which we only know the value of predictors,  $\tilde{x}$ .
- Assume that we have observed  $(x, y)$  for  $n$  cases, and we want to predict  $\tilde{y}$  for a new observation with predictor values  $\tilde{x}$ .
- Since the covariance function depends on  $x$ , we can find  $C_{n+1}$  for  $n$  the training cases and the new observation, i.e., for  $\begin{pmatrix} x \\ \tilde{x} \end{pmatrix}$ . To avoid confusion we denote the covariance matrix for just the training cases as  $C_n$ .
- We can write down  $C_{n+1}$  as follows:

$$C_{n+1} = \begin{pmatrix} C_n & K \\ K^T & v \end{pmatrix}$$

where  $K$  is the  $n \times 1$  covariance vector between  $\tilde{y}$  and the  $n$  observed  $y$ .  $v$  is the prior variance of  $\tilde{y}$  obtained based on the covariance function  $C$ .

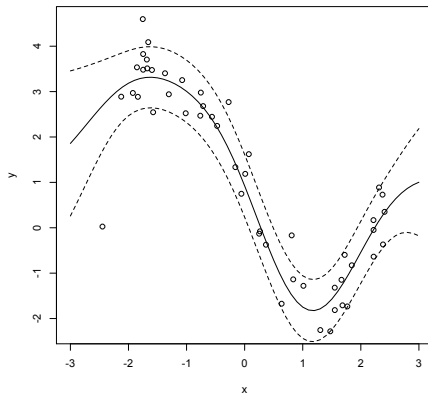
- Based the above setting, we can obtain the posterior predictive distribution for the new case.
- This distribution is also Gaussian with the following mean and variance:

$$\begin{aligned}E(\tilde{y}|y) &= K^T C_n^{-1} y \\ \text{Var}(\tilde{y}|y) &= v - K^T C_n^{-1} K\end{aligned}$$

- If we need a point estimate, we can use  $E(\tilde{y}|y)$ .

# Example

- The following example shows a Gaussian process model trained on 100 data points uniformly sampled from -2 to 2 .



# Example

- For the above model, we used the following covariance function:

$$\text{Cov}(y_i, y_j) = 2 + \exp(-0.5(x_i - x_j)^2) + \delta_{ij} \times 0.1$$

- The solid line is expected function based on a grid test points between -3 and 3.
- The dashed lines show the 95% interval for predictions.

# Hyperparameters

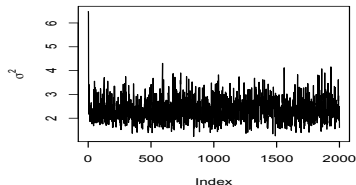
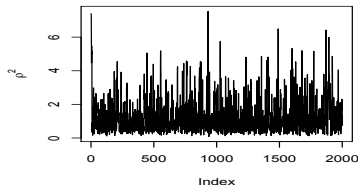
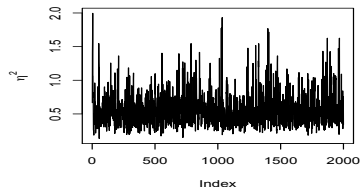
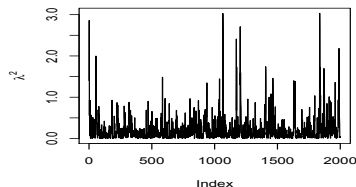
- In reality, we might not have enough information to fix the parameters of the covariance functions.
- In general, we would treat these parameters (e.g.,  $\eta$ ,  $\rho$ 's,  $\lambda$  and  $\sigma_\epsilon$ ) as hyperparameters.
- Therefore, we need to use MCMC simulations in order to obtain samples from the posterior distributions of these hyperparameters, and as usual, we integrate over these posterior distributions to obtain the posterior predictive probabilities.
- The log-likelihood function in this case is as follows:

$$\ell = -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log \det C - \frac{1}{2} x^T C^{-1} x$$

- Note that the computational cost of  $C^{-1}$  is in general  $\mathcal{O}(n^3)$ .

# Example

- For the following example,  $\eta^2$ ,  $\rho^2$ ,  $\lambda^2$ , and  $\sigma^2$  are hyperparameters with Gamma(1, 1) priors.



# Gaussian process models for classification

- For categorical outcome variables, we assume the Gaussian process prior over a continuous *latent function*,  $u(x)$ .
- We define the distribution of the response variable in terms of this latent function.
- For example, if the outcome variable  $y$  is binary, we can use the following logistic model:

$$P(y_i = 1|u(x_i)) = \frac{\exp(u(x_i))}{1 + \exp(u(x_i))}$$

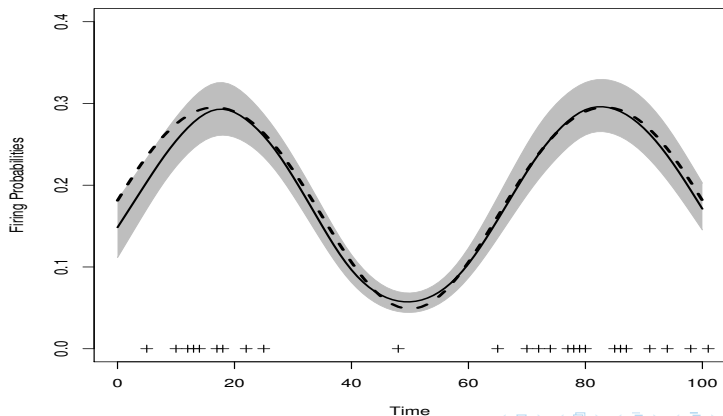
or alternatively,

$$P(y_i = 1|u(x_i)) = \frac{1}{1 + \exp(-u(x_i))}$$

- We can use a multinomial logit model for outcome variables with multiple categories.

# Gaussian process models for classification

- Here, we are using a GP model to estimate the underlying firing rates of a neuron (i.e.,  $y_t = 1$  when the neuron fires,  $y_t = 0$  otherwise).
- The dashed line shows the true firing probability and the plus signs show the firing time.





# Covariance functions (kernels)

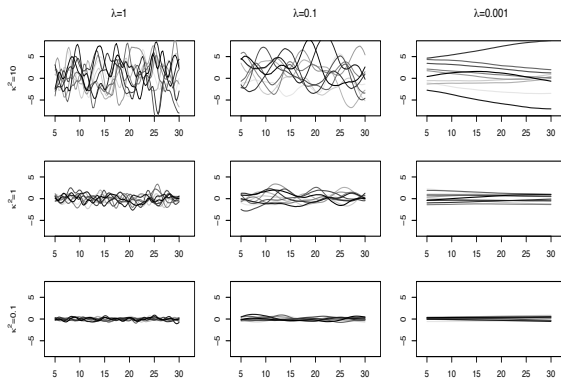
- Choosing an appropriate covariance function is an important step in setting Gaussian process models.
- Usually, from a class of valid kernels we choose a kernel that represents our beliefs regarding the underlying function,  $y(x)$ .
- Sometimes, we might choose a kernel that is computationally convenient computationally.
- In what follows, we discuss some of these alternative kernels.
- Note that we can create new kernels using their products and linear combinations.

# Squared exponential

- The covariance function we discussed earlier is called *squared exponential*, which has the following form:

$$C_{ij} = \kappa^2 \exp[-\lambda(x_i - x_j)^2]$$

- Here,  $\lambda$  controls the correlation length, while  $\kappa^2$  accounts for the height of oscillations in realizations of the GP.



# Matérn process

- The Matérn class of kernels is defined in terms of the smoothness parameter,  $\nu$ , length-scale,  $\rho$ , and variance  $\sigma^2$  as follows:

$$C_{ij} = \sigma^2 \frac{1}{\Gamma(\nu)2^{\nu-1}} \left( \sqrt{2\nu} \frac{|x_i - x_j|}{\rho} \right)^\nu K_\nu \left( \sqrt{2\nu} \frac{|x_i - x_j|}{\rho} \right)$$

- Here,  $\Gamma$  and  $K$  are the Gamma and modified Bessel functions respectively.
- For  $\nu = \frac{1}{2} + n$  and  $n = 0, 1, \dots$ , the corresponding GP is  $n$  times continuously differentiable.

# Ornstein-Uhlenbeck (OU) process

- Setting  $\nu = 1/2$ , we obtain a special case of GP called Ornstein-Uhlenbeck (OU) process,

$$C_{ij} = \sigma^2 \exp\left(-\frac{|x_i - x_j|}{\rho}\right)$$

- Compared to the squared exponential kernel, the resulting model is not very flexible.

# Brownian motion (Wiener process)

- Brownian motion (Wiener process) is still simpler (rougher) than the OU process

$$C_{ij} = \sigma^2 \min(x_i, x_j)$$

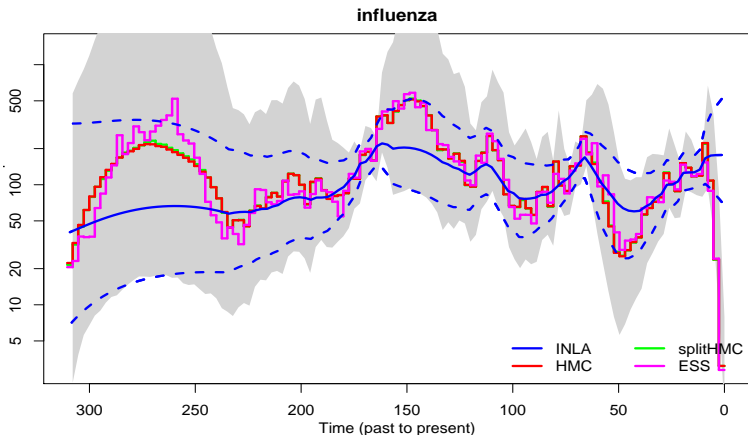
- In this case, the computational cost of inverting  $C$  is  $\mathcal{O}(n)$ .
- The OU process and the Wiener process are related through the following SDE:

$$dY_t = -\rho(Y_t - \mu)dt + dW_t,$$

where  $Y_t$  is an OU process and  $W_t$  is Brownian motion.

# Brownian motion (Wiener process)

- Here is an example of using Brownian motion to model population dynamics of influenza (Lan, et al., 2014) using different sampling algorithms, Hamiltonian Monte Carlo (HMC), SplitHMC, and Elliptical Slice Sampling (ESS), and comparing the results to Integrated Nested Laplace Approximation (INLA).



- To learn more about this topic, you could refer to
  - ▶ “Regression and classification using Gaussian process priors” (with discussion), by Neal, R. M. (1998).
  - ▶ “Gaussian Processes for Machine Learning,” by Rasmussen and Williams (2006)
  - ▶ “Hierarchical Modeling and Analysis for Spatial Data,” by Banerjee, Carlin, and Gelfand (2014).
  - ▶ “Dependent Matérn Processes for Multivariate Time Series, Vandenberg-Rodes, A. and Shahbaba, B. (2015).