

Advanced Bayesian Computation Weeks 5

Rajarshi Guhaniyogi
Winter 2018

February 9, 2018

Choice of Covariance Functions

- There are certain realistic assumptions often employed on the covariance function $C_{\theta}(t_i, t_j)$.
- Stationary: $C_{\theta}(t_i, t_j) = C_{\theta}(t_i - t_j)$; Isotropic: $C_{\theta}(t_i, t_j) = C_{\theta}(\|t_i - t_j\|)$.
- The covariance function of a stationary process can be represented as the Fourier transform of a positive finite measure.
- This is the famous Bochner's theorem.
- Let $h = t_i - t_j$, A real-valued function $C_{\theta}(h)$ on \mathbb{R}^D is the covariance function of a stationary real valued random process on \mathbb{R}^D if and only if it can be represented as

$$C_{\theta}(h) = \int \cos(2\pi h \cdot s) dH(s),$$

where $H(s)$ is a positive finite measure.

Choice of Covariance Functions

- If $H(s)$ has a density $S(s)$, then $S(s)$ is called the spectral density.
- Note that C_θ and $S(s)$ are Fourier-dual of each other, i.e.

$$C_\theta(h) = \int \cos(2\pi h \cdot s) dS(s) ds, \quad S(s) = \int \cos(-2\pi h \cdot s) C_\theta(h) dh.$$

- Most of the practical applications we take the covariance kernel as an isotropic kernel.
- **Squared Exponential Covariance:** $C_\theta(r) = \exp\left(-\frac{r^2}{\phi}\right)$.
- The spectral density is given by $S(s) = (\sqrt{2\pi\phi})^D \exp(-2\pi^2\phi s^2)$.
- The sample path is infinitely differentiable.
- This is the most popularly used covariance kernel in the machine learning literature.

- **Matern Covariance:**

$$C_{\theta}(r) = \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu}r}{\phi} \right) K_{\nu} \left(\frac{\sqrt{2\nu}r}{\phi} \right),$$

K_{ν} is the modified Bessel function.

- The spectral density is
- Here ν is called the smoothness parameter which determines the smoothness of the sample path.
- As ν is increased, the sample paths are more smooth.
- As $\nu \rightarrow \infty$, Matern covariance kernel converges to the squared exponential covariance kernel.
- If k is the greatest integer less than ν , then the Gaussian process is k times mean square differentiable.

- **Exponential Covariance:**

$$C_{\theta}(r) = \exp\left(-\frac{r}{2\phi}\right).$$

- This is a special case of the Matern covariance kernel with $\nu = 1/2$.
- The sample path is only continuous, not even differentiable once.
- In the one dimensional case this is the covariance function of the Ornstein-Uhlenbeck (OU) process.
- The OU process [Uhlenbeck process and Ornstein, 1930] was introduced as a mathematical model of the velocity of a particle undergoing Brownian motion.
- Gaussian process with the exponential kernel is not even once mean square differentiable.

- **Rational Quadratic Covariance:**

$$C_{\theta}(r) = \left(1 + \frac{r^2}{2\alpha\phi}\right)^{-\alpha}$$

- This is a scale mixture of squared exponential kernel.
- Sometimes used for a greater flexibility over squared exponential.
- The Gaussian process is infinitely mean squared differentiable for every α .
- As $\alpha \rightarrow \infty$, the rational quadratic function takes more and more the shape of a squared exponential covariance function.

Constructing More Complicated Covariance Functions

- Rather than using an isotropic function, one may want to use stationary covariance function.
- Define the distance metric $r(t_i, t_j)^2 = (t_i - t_j)'M(t_i - t_j)$, M is a positive definite matrix.
- Now replace any of the already defined covariance kernels by $C_{\theta}(t_i - t_j) = C_{\theta}(r(t_i, t_j))$.
- Some non-stationary covariance functions are used in some applications.
- For example, dot product covariance function

$$C_{\theta}(t_i, t_j) = \sigma^2 + t_i \cdot t_j$$

- Neural network covariance function is used sometimes.

$$C_{\theta}(t_i, t_j) = \frac{2}{\pi} \sin^{-1} \left(\frac{2\tilde{t}_i' \Sigma \tilde{t}_j}{\sqrt{(1 + 2\tilde{t}_i' \Sigma \tilde{t}_i)(1 + 2\tilde{t}_j' \Sigma \tilde{t}_j)}} \right),$$

where $\tilde{t}_i = (1, t_i)'$.

Model and Gaussian process prior

- The non-linear relationship between y and \mathbf{x} is given by

$$y = f(\mathbf{x}) + \epsilon, \quad \epsilon \sim N(0, \tau^2)$$

- Need prior distributions on $\{f(\cdot), \tau^2\}$.
- $f(\cdot)$ is assigned a Gaussian process prior distribution $GP(\mu, C_{\theta}(\cdot, \cdot))$
- Lets demonstrate everything with the exponential covariance function. Let $\theta = (\sigma^2, \phi)$,

$$C_{\sigma^2, \phi}(t_i, t_j) = \sigma^2 \exp(-\|t_i - t_j\|/\phi).$$

- Note that this is going to create sample paths which are only continuous and not even once differentiable while the function $f(\cdot)$ may be more smooth.

- One may use Matern covariance kernel with a prior on the smoothness parameter ν .
- It was not possible to learn ν .
- Thus it is a common practice to assign the prior distribution with fixing ν .
- For exponential kernel we are fixing $\nu = 2$.
- Prior on σ^2 and τ^2 are assigned inverse-gamma(a,b) prior.
- A Normal prior on μ .
- We need to be careful in assigning prior on ϕ .

- Suppose we have the data $(y_1, \mathbf{x}_1), \dots, (y_n, \mathbf{x}_n)$.
- Thus $y_i = f(\mathbf{x}_i) + \epsilon_i$, $\epsilon_i \sim N(0, \tau^2)$.
- Suppose $\mathbf{y} = (y_1, \dots, y_n)'$, $\boldsymbol{\epsilon} = (\epsilon_1, \dots, \epsilon_n)'$ and $\mathbf{f} = (f(\mathbf{x}_1), \dots, f(\mathbf{x}_n))'$.
- Thus $\mathbf{y} = \mathbf{f} + \boldsymbol{\epsilon}$, $\boldsymbol{\epsilon} \sim N(\mathbf{0}, \tau^2 \mathbf{I})$.
- From the Gaussian process specification, a priori $\mathbf{f} \sim N(\mu \mathbf{1}_n, \mathbf{C}_{\sigma^2, \phi})$.
- The posterior distribution

$$p(\phi, \sigma^2, \mu, \tau^2 | \mathbf{y}) \propto N(\mathbf{y} | \mu \mathbf{1}_n, \mathbf{C}_{\sigma^2, \phi} + \tau^2 \mathbf{I}) \times N(\mu | \mu_\mu, \sigma_\mu^2) \\ \times IG(\tau^2 | a, b) \times IG(\sigma^2 | a, b) \times p(\phi).$$

Model fitting proceeds through MCMC steps. We run Gibbs within Metropolis.

- $\mu|-$ follows a normal distribution.
- σ^2, τ^2, ϕ are updated using Metropolis steps.
- In spatial statistics, one uses the classical geostatistical Gaussian process model

$$y(\mathbf{s}) = \mathbf{x}(\mathbf{s})'\boldsymbol{\beta} + f(\mathbf{s}) + \epsilon(\mathbf{s}), \quad \epsilon(\mathbf{s}) \sim N(0, \tau^2).$$

- $\mathbf{y} = (y(\mathbf{s}_1), \dots, y(\mathbf{s}_n))'$, $\mathbf{X} = [\mathbf{x}(\mathbf{s}_1) : \dots : \mathbf{x}(\mathbf{s}_n)]'$,
 $\mathbf{f} = (f(\mathbf{s}_1), \dots, f(\mathbf{s}_n))'$.
- Let $\mathbf{C}_{\sigma^2, \phi} = ((C_{\sigma^2, \phi}(\mathbf{s}_i, \mathbf{s}_j)))_{i,j=1}^n$

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{f} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim N(\mathbf{0}, \tau^2 \mathbf{I}).$$

Posterior distribution of the function

- Note that

$$p(\mathbf{f}|\mathbf{y}) \propto N(\mathbf{y}|\mathbf{f}, \tau^2 \mathbf{I}) \times N(\mathbf{f}|\mu \mathbf{1}_n, \mathbf{C}_{\sigma^2, \phi})$$

- Thus $\mathbf{f}|\mathbf{y}$, – is a multivariate normal distribution.
- For $\{\mu^{(l)}, \sigma^{2(l)}, \tau^{2(l)}, \phi^{(l)}\}_{l=1}^L$ L post burn-in MCMC samples, we draw $\mathbf{f}^{(1)}, \dots, \mathbf{f}^{(L)}$ L MCMC samples of the posterior realization of the function at n data points.
- For inference at an arbitrary point \mathbf{x} , we calculate the distribution of $f(\mathbf{x})|\mathbf{y}$.

$$p(f(\mathbf{x})|\mathbf{y}) = \int [p(f(\mathbf{x})|f(\mathbf{x}_1), \dots, f(\mathbf{x}_n), \sigma^2, \tau^2, \phi, \mu) \\ p(f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)|\sigma^2, \tau^2, \phi, \mu, \mathbf{y})p(\sigma^2, \tau^2, \phi, \mu|\mathbf{y})].$$

- We already know how to draw samples from $p(f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)|\sigma^2, \tau^2, \phi, \mu, \mathbf{y})$ and $p(\sigma^2, \tau^2, \phi, \mu|\mathbf{y})$.

- $f(\mathbf{x})|f(\mathbf{x}_1), \dots, f(\mathbf{x}_n), \sigma^2, \tau^2, \phi, \mu \sim N(\mu_f, \sigma_f^2)$.

$$\sigma_f^2 = C_{\sigma^2, \phi}(\mathbf{x}, \mathbf{x}) - \mathbf{c}_{\sigma^2, \phi}(\mathbf{x})' \mathbf{C}_{\sigma^2, \phi}^{-1} \mathbf{c}_{\sigma^2, \phi}(\mathbf{x}).$$

$$\mu_f = \mu + \mathbf{c}_{\sigma^2, \phi}(\mathbf{x})' \mathbf{C}_{\sigma^2, \phi}^{-1} (\mathbf{f} - \mu \mathbf{1}_n)$$

$$\mathbf{c}_{\sigma^2, \phi}(\mathbf{x}) = (C_{\sigma^2, \phi}(\mathbf{x}, \mathbf{x}_1), \dots, C_{\sigma^2, \phi}(\mathbf{x}, \mathbf{x}_n))'$$

- For spatial process models, finding posterior distribution of the function is also similar.

- Suppose the prediction of response is required at \mathbf{x} .
- Note that, $y \sim N(f(\mathbf{x}), \tau^2)$.
- We have already seen how to draw post burn-in samples $f(\mathbf{x})^{(1)}, \dots, f(\mathbf{x})^{(L)}$ from $f(\mathbf{x})|y_1, \dots, y_n$.
- Posterior predictive samples $y^{(1)}, \dots, y^{(L)}$ are drawn from $y^{(l)} \sim N(f(\mathbf{x})^{(l)}, \tau^{2(l)})$.
- In sample prediction can be similarly performed.

Multivariate Gaussian Process model

- Multivariate Gaussian process models are most often used in the geostatistical analysis.
- Let $\mathbf{y}(\mathbf{s}) = (y_1(\mathbf{s}), \dots, y_m(\mathbf{s}))'$, $\mathbf{w}(\mathbf{s}) = (w_1(\mathbf{s}), \dots, w_m(\mathbf{s}))'$, $\boldsymbol{\epsilon}(\mathbf{s}) = (\epsilon_1(\mathbf{s}), \dots, \epsilon_m(\mathbf{s}))'$.
- The multivariate model is given by

$$\mathbf{y}(\mathbf{s}) = \mathbf{B}\mathbf{x}(\mathbf{s}) + \mathbf{w}(\mathbf{s}) + \boldsymbol{\epsilon}(\mathbf{s}),$$

$$\mathbf{B} = ((\beta_{ij}))_{i,j=1}^{m,p}, \quad \boldsymbol{\epsilon}(\mathbf{s}) \sim N(\mathbf{0}, \boldsymbol{\Psi}).$$

- When data observed at $\mathbf{s}_1, \dots, \mathbf{s}_n$,

$$\mathbf{y} = \text{vec}(\mathbf{B}\mathbf{X}) + \mathbf{w} + \boldsymbol{\epsilon},$$

$$\mathbf{w} = (\mathbf{w}(\mathbf{s}_1)', \dots, \mathbf{w}(\mathbf{s}_n)')', \quad \boldsymbol{\epsilon} = (\boldsymbol{\epsilon}(\mathbf{s}_1)', \dots, \boldsymbol{\epsilon}(\mathbf{s}_n)')', \\ \mathbf{X} = [\mathbf{x}(\mathbf{s}_1) : \dots : \mathbf{x}(\mathbf{s}_n)].$$

- $\boldsymbol{\epsilon} \sim N(\mathbf{0}, \mathbf{I} \otimes \boldsymbol{\Psi})$.

Modeling Multivariate Gaussian Process

- **Linear Model Coregionalization:** $\mathbf{w}(\mathbf{s}) = \mathbf{A}\mathbf{v}(\mathbf{s})$, where $\mathbf{v}(\mathbf{s}) = (v_1(\mathbf{s}), \dots, v_m(\mathbf{s}))'$, \mathbf{A} is an $m \times m$ matrix.
- Is \mathbf{A} identifiable?
- Popular specification is \mathbf{A} is a lower triangular matrix with diagonal entries all positive.
- $v_1(\mathbf{s}), \dots, v_m(\mathbf{s})$ are assigned independent Gaussian process.
- **Multivariate Matern Kernel:** Specification of Matern kernel for the multivariate Gaussian process so that marginally each component follows a Gaussian process with a univariate Matern kernel.