

Kriging

In my view Kriging can be summarized as random field prediction in the presence of functional covariates with unknown coefficients.

"classic" Model where Kriging is applied.

Data:

d_1, \dots, d_n at spatial locations $x_1, \dots, x_n \in \mathbb{R}^d$.

Model:

$$(K1) \quad d_i = Y(x_i) + \varepsilon_i$$

$$(K2) \quad Y(x) = \sum_{p=1}^m \beta_p f_p(x) + Z(x)$$

(K3) $f_p: \mathbb{R}^d \rightarrow \mathbb{R}$ are known "functional covariates".

(K4) β_p are unknown coefficients

$$(K5) \quad (\varepsilon_1, \dots, \varepsilon_n) \sim N(0, \sigma^2 I).$$

$$(K6) \quad Z \sim GRF_{\mathbb{R}^d}(0, K)$$

$\hookrightarrow K: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$.

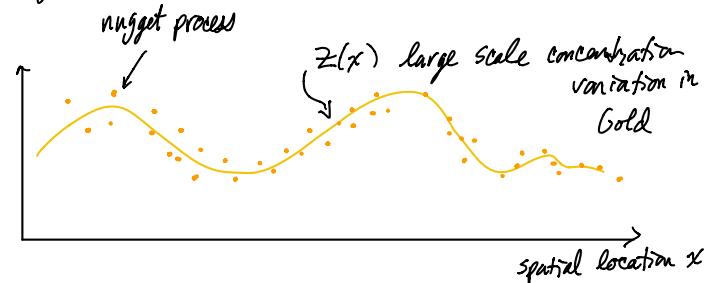
Goal:

Predict $Y(x_0)$ at some $x_0 \in \mathbb{R}^d$

This is a natural model in the sense that the linear model can explain large scale (deterministic) variation, ε_i 's account for spatially uncorrelated error and $Z(x)$ accounts for "unknown spatial covariates" which induce spatially correlated errors using the linear model $\sum_p \beta_p f_p(x)$.

Note: The additive errors $\varepsilon_1, \dots, \varepsilon_n$ correspond to what is called a nugget effect. It can be used to model instrumental noise or a microscale process that has correlation length scale much smaller than $Z(x)$.

The name comes from mining applications e.g.



There are a few equivalent characterizations of Kriging, each giving a different perspective but yield equivalent predictions.

Kriging with Generalized least squares
Regression

I think this is the most intuitive and direct but makes it difficult to see why one can use generalized cov funs (later).

Note that

$$\begin{bmatrix} d_1 \\ \vdots \\ d_n \\ Y(x_0) \end{bmatrix} = \begin{bmatrix} \sum_p \beta_p f_p(x_1) \\ \vdots \\ \sum_p \beta_p f_p(x_n) \\ \sum_p \beta_p f_p(x_0) \end{bmatrix} + \begin{bmatrix} Z(x_1) \\ \vdots \\ Z(x_n) \\ Z(x_0) \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \\ 0 \end{bmatrix}$$

$$\text{So letting } f(x_0) = \begin{bmatrix} f_1(x_0) \\ \vdots \\ f_m(x_0) \end{bmatrix} \text{ & } f(x_1, \dots, x_n) = \begin{bmatrix} 1 & \dots & 1 \\ f(x_1) & \dots & f(x_n) \\ \vdots & \ddots & \vdots \end{bmatrix}$$

we have

$$\begin{bmatrix} d_1 \\ \vdots \\ d_n \\ Y(x_0) \end{bmatrix} \sim N \left(\begin{bmatrix} f(x_1, \dots, x_n)^T \beta \\ \vdots \\ f(x_n)^T \beta \\ Y(x_0) \end{bmatrix}, \begin{bmatrix} K(x_i, x_j)_{i,j=1}^n + \sigma^2 I & K(x_j, x_0) \\ \vdots & \vdots \\ K(x_0, x_i) & K(x_0, x_0) \end{bmatrix} \right)$$

" " "

$$\begin{bmatrix} X\beta \\ \vdots \\ f(x_n)^T \beta \\ Y(x_0) \end{bmatrix} \quad \begin{bmatrix} \Sigma_{11} & \Sigma_{10} \\ \vdots & \vdots \\ \Sigma_{01} & \Sigma_{00} \end{bmatrix}$$

so $d = X\beta + \varepsilon$ is simply a GLS regression problem:

$$\hat{\beta} = (X^T \Sigma_{11}^{-1} X)^{-1} X^T \Sigma_{11}^{-1} d \quad (1a)$$

To predict $Y(x_0)$ we simply use the multivariate model for $(d^T, Y(x_0))^T$ as above and write

$$E(Y(x_0)|d) = \sum_{01} \sum_{11}^{-1} (d - X\beta) + f(x_0)^T \beta$$

which requires known β so the Kriging prediction simply plugs in $\hat{\beta}$:

$$\hat{Y}(x_0) := \sum_{01} \sum_{11}^{-1} (d - X\hat{\beta}) + f(x_0)^T \hat{\beta} \quad (1b)$$

At this point it's easy to remember and compute but it's not clear the impact of plugging in $\hat{\beta}$ for β in $E(Y(x_0)|d)$.

The next characterization shows it is optimal among "linear unbiased predictions".

Kriging as a best linear unbiased predictor

The idea predict $\hat{Y}(x_0)$ via $d^T \lambda$

$$\text{where } \lambda = \arg \min_{\lambda \in \mathcal{U}B} E(Y(x_0) - d^T \lambda)^2$$

where the "unbias" constraint $\lambda \in \mathcal{U}B$ is:

$$E(\hat{Y}(x_0)) = f(x_0)^T \beta$$

Note the LHS \uparrow equals

$$\begin{aligned} E(d^T \lambda) &= (f(x_1, \dots, x_n)^T \cdot \beta)^T \cdot \lambda \\ &= \beta^T \cdot f(x_1, \dots, x_n) \cdot \lambda \end{aligned}$$

so the unbiased constraint becomes

$$E(\hat{Y}(x_0)) = f(x_0)^T \beta$$

\uparrow

$$\forall \beta \in \mathbb{R}^m, \beta^T \cdot f(x_1, \dots, x_n) \cdot \lambda = \beta^T \cdot f(x_0) \cdot \lambda$$

\updownarrow

$$\boxed{f(x_1, \dots, x_n)} \cdot \boxed{\lambda} = \boxed{f(x_0)}$$

$\underbrace{}$ length $n \dots$
i.e. the number of
obs points

$m \times 1$

\updownarrow

$$\lambda \in \mathcal{U}B$$

Remark: Notice that if we define

$$\lambda^{\text{ext}} = (\lambda_1, \lambda_2, \dots, \lambda_n, -1)^T$$

$$\text{then } f(x_1, \dots, x_n, x_0) \cdot \lambda^{\text{ext}} = 0$$

To summarize we predict $\hat{Y}(x_0) = d^T \cdot \lambda$
where λ satisfies

$$f(x_1, \dots, x_n) \cdot \lambda = f(x_0)$$

$$\text{i.e. } f(x_1, \dots, x_n, x_0) \cdot \lambda^{\text{ext}} = 0$$

and minimizes

$$\begin{aligned} & E(Y(x_0) - d^T \cdot \lambda)^2 \\ &= E\left(\begin{bmatrix} d \\ Y(x_0) \end{bmatrix}^T \cdot \lambda^{\text{ext}}\right)^2 \\ &= (\lambda^{\text{ext}})^T \begin{bmatrix} [K(x_i, x_j)]_{i,j=1}^n + \sigma^2 I & K(x_j, x_0) \\ \vdots & \vdots \\ \dots & K(x_0, x_j) \dots & K(x_0, x_0) \end{bmatrix} \lambda^{\text{ext}} \\ &\quad \text{This term is } (f(x_1, \dots, x_n, x_0) \cdot \lambda^{\text{ext}})^2 \\ &\quad \text{but is zero due to the unbiased constraint} \\ &= (\lambda^{\text{ext}})^T \begin{bmatrix} \Sigma_{11} & \Sigma_{10} \\ \Sigma_{01} & \Sigma_{00} \end{bmatrix} \lambda^{\text{ext}} \end{aligned}$$

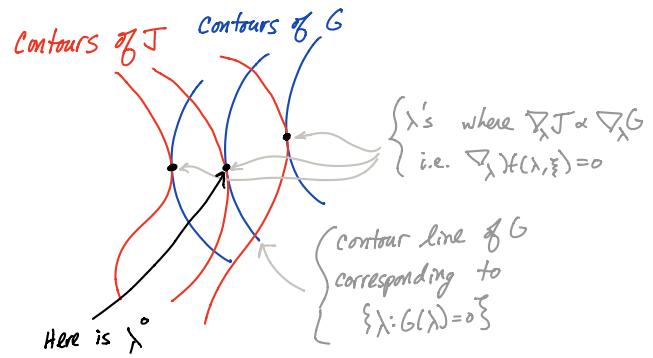
Recall the technique of using Lagrange multipliers for constrained minimization:

$$\lambda^* \in \arg \min_{\{\lambda: G(\lambda)=0\}} J(\lambda) \quad \begin{array}{l} \text{where} \\ \dim(G(\lambda)) = m \\ \dim(\lambda) = n > m \end{array}$$

$$\Updownarrow \quad \lambda^* \in \left\{ \lambda : \exists \xi \in \mathbb{R}^m \text{ s.t. } \begin{array}{l} \nabla J(\lambda, \xi) = 0 \\ \text{and } G(\lambda) = 0 \end{array} \right\}$$

$$\text{with } \mathcal{H}(\lambda, \xi) = J(\lambda) + \xi^T \cdot G(\lambda)$$

The picture looks like this



In our case we have

$$\begin{aligned} \mathcal{H}(\lambda, \xi) &= J(\lambda) + \xi^T \cdot G(\lambda) \\ \frac{1}{2} \begin{bmatrix} \lambda \\ -1 \end{bmatrix}^T \begin{bmatrix} \Sigma_{11} & \Sigma_{10} \\ \Sigma_{01} & \Sigma_{00} \end{bmatrix} \begin{bmatrix} \lambda \\ -1 \end{bmatrix} & \quad \xi^T \cdot (f(x_1, \dots, x_n) \lambda - f(x_0)) \end{aligned}$$

Also

$$\frac{\partial J(\lambda)}{\partial \lambda} = \begin{bmatrix} \lambda \\ -1 \end{bmatrix}^T \begin{bmatrix} \Sigma_{11} & \Sigma_{10} \\ \Sigma_{01} & \Sigma_{00} \end{bmatrix} \begin{bmatrix} I_{n \times n} \\ 0_{1 \times n} \end{bmatrix}$$

so that

$$\begin{aligned} \nabla_\lambda J(\lambda) &= \begin{bmatrix} I_{n \times n} & 0_{n \times 1} \\ \Sigma_{11} & \Sigma_{10} \end{bmatrix} \begin{bmatrix} \lambda \\ -1 \end{bmatrix} \\ &= \Sigma_{11} \lambda - \Sigma_{10} \end{aligned}$$

and therefore

$$\nabla_\lambda \mathcal{H}(\lambda, \xi) = 0 \text{ and}$$

$$G(\lambda) = 0$$

\Updownarrow

$$\Sigma_{11} \lambda - \Sigma_{10} + f(x_1, \dots, x_n)^T \cdot \xi = 0 \text{ and}$$

$$f(x_1, \dots, x_n) \lambda - f(x_0) = 0$$

\Updownarrow

$$\begin{bmatrix} \Sigma_{11} & f(x_1, \dots, x_n)^T \\ f(x_1, \dots, x_n) & 0_{m \times m} \end{bmatrix} \begin{bmatrix} \lambda \\ \xi \end{bmatrix} = \begin{bmatrix} \Sigma_{10} \\ f(x_0) \end{bmatrix}$$

Therefore the Kriging predictor via (unbias) constrained minimizing MSE is:

$$\hat{Y}(x_0) = d^T \lambda \quad (2a)$$

where λ satisfies

$$\begin{bmatrix} \Sigma_{11} & f(x_1, \dots, x_n)^T \\ f(x_1, \dots, x_n) & 0_{m \times m} \end{bmatrix} \begin{bmatrix} \lambda \\ \xi \end{bmatrix} = \begin{bmatrix} \Sigma_{10} \\ f(x_0) \end{bmatrix} \quad (2b)$$

with $\Sigma_{ii} = [K(x_i, x_j)]_{i,j=1}^n + \sigma^2 I$ and $\Sigma_{10} = \begin{bmatrix} \vdots \\ K(x_0, x_0) \\ \vdots \end{bmatrix}$

Remark: Notice that in this characterization we automatically get an estimate of the variance of prediction error by

$$\begin{aligned} E(Y(x_0) - \hat{Y}(x_0))^2 \\ = \begin{bmatrix} \lambda \\ -1 \end{bmatrix}^T \begin{bmatrix} \Sigma_{11} & \Sigma_{10} \\ \Sigma_{01} & \Sigma_{00} \end{bmatrix} \begin{bmatrix} \lambda \\ -1 \end{bmatrix} \end{aligned} \quad (2c)$$

where $(\lambda^T, \xi^T)^T$ solves the above Lagrange multiplier problem.

Kriging as a Spline interpolator

A typical spline interpolator has the form:

$$\hat{Y}(x_0) = \sum_{k=1}^n c_k K(x_0, x_k) + \sum_{p=1}^m b_p f_p(x_0)$$

↑ ↑
reproducing kernel knots

where $c = (c_1, \dots, c_n)^T$ & $b = (b_1, \dots, b_m)^T$ are given by minimizing "bending energy" and "data fidelity" as follows

$$(c^*, b^*) \in \underset{(c, b)}{\operatorname{argmin}} \left\{ \mathcal{D}(c, b) + \mathcal{E}(c) \right\}$$

where

$\mathcal{D}(c, b) \equiv$ data fidelity

$$= \frac{1}{2\sigma^2} \| d - [\hat{Y}(x_1), \dots, \hat{Y}(x_n)]^T \|^2$$

$$= \frac{1}{2\sigma^2} \| d - (Kc + F^T b) \|^2$$

↑ ↑
[$K(x_i, x_j)$]_{i,j=1}ⁿ $F := f(x_1, \dots, x_n)$

and

$$\begin{aligned} \mathcal{E}(c) &\equiv \text{bending energy} \\ &= \frac{1}{2} c^T K c \end{aligned}$$

with $\sum_{p=1}^m b_p f_p$
unpenalized

$\frac{1}{2} \langle Kc | K^{-1} | Kc \rangle$

Now the minimizer is characterized by

$$\nabla_c \mathcal{D}(c, b) + \nabla_c E(c) = 0$$

$$\nabla_b \mathcal{D}(c, b) = 0$$

\Updownarrow

$$- (d - Kc - F^T b) + \sigma^2 c = 0$$

after a left multiply with $\sigma^2 K^{-1}$

$$- F(d - Kc - F^T b) = 0$$

\Updownarrow this is $\sigma^2 c$ from the line above

$$(K + \sigma^2 I)c + F^T b = d$$

$$Fc = 0$$

In summary the Kriging predictor (defined via a spline approach) is defined as:

$$\hat{Y}(x_0) = \sum_{k=1}^n c_k K(x_0, x_k) + \sum_{p=1}^m b_p f_p(x_0) \quad (3a)$$

where

$$\begin{bmatrix} \Sigma_{11} & f(x_1, \dots, x_n)^T \\ f(x_1, \dots, x_n) & 0_{m \times m} \end{bmatrix} \begin{bmatrix} c \\ b \end{bmatrix} = \begin{bmatrix} d \\ 0 \end{bmatrix} \quad (3b)$$

$$\text{with } \Sigma_{11} = \left[K(x_i, x_j) \right]_{i,j=1}^n + \sigma^2 I.$$

All three Kriging characterizations are the same

$\hat{Y}(x_0)$ from (2a)

$$= d^T \gamma$$

$$= \begin{bmatrix} \gamma \\ \xi \end{bmatrix}^T \begin{bmatrix} d \\ 0 \end{bmatrix}$$

$$(2b) = \begin{bmatrix} \Sigma_{10} \\ f(x_0) \end{bmatrix}^T \cdot \begin{bmatrix} \Sigma_{11} & f(x_1, \dots, x_n)^T \\ f(x_1, \dots, x_n) & 0_{m \times m} \end{bmatrix}^{-1} \begin{bmatrix} d \\ 0 \end{bmatrix}$$

$$(3b) = \begin{bmatrix} \Sigma_{10} \\ f(x_0) \end{bmatrix}^T \begin{bmatrix} c \\ d \end{bmatrix}$$

$= \hat{Y}(x_0)$ from (3a)

$\Rightarrow \hat{Y}(x_0)$ from (1a)

verified by noticing that
 $b = \hat{\beta}$ and $c = \Sigma^{-1}(d - f(x_1, \dots, x_n)^T \hat{\beta})$
satisfy (3b)

Now

- $\hat{Y}(x_0)$ from (1a) is probably easiest to remember and intuitive
- $\hat{Y}(x_0)$ from (2a) shows this predictor is optimal (in m.s. sense) among all unbiased predictors and gives prediction uncertainty.
- $\hat{Y}(x_0)$ from (3a) is used to establish all three are equal and gives a "bayesian" interpretation (viewing $E(c)$ as a "prior").

Replacing the nugget variance with a Noise-to-Signal parameter

By examining the Kriging prediction $\hat{Y}(x_0)$ given in the three different forms above, one discovers a number of simplifications that reduce the number of parameters needed. Here is one example.

Start by writing the Kriging model with an overall variance parameter σ_z^2 for Z and the nugget variance σ_ϵ^2 for ϵ_k

$$(*) \quad \begin{cases} d_k = \sum_{p=1}^m b_p f_p(x_p) + Z(x_p) + \epsilon_k \\ Z \sim GRF_p(0, \sigma_z^2 K) \\ \epsilon_1, \dots, \epsilon_n \stackrel{iid}{\sim} N(0, \sigma_\epsilon^2) \end{cases} \quad \text{so } K(x, x) = 1$$

by simply factoring out σ_z^2 one notices $\hat{Y}(x_0)$ in (1a) from model (*) is the same prediction as $\hat{Y}(x_0)$ in (1a) from model (**):

$$(**) \quad \begin{cases} d_k = \sum_{p=1}^m b_p f_p(x_p) + Z(x_p) + \epsilon_k \\ Z \sim GRF_p(0, K) \\ \epsilon_1, \dots, \epsilon_n \stackrel{iid}{\sim} N(0, \frac{\sigma_\epsilon^2}{\sigma_z^2}) \end{cases} \quad \begin{matrix} K(x, x) = 1 \\ \text{still} \end{matrix}$$

\curvearrowleft ratio of nugget var to Z var.

So once the correlation function K and the covariates f_p are defined the behavior of $\hat{Y}(x_0)$ is completely determined by a single parameter: $\sigma_\epsilon^2 / \sigma_z^2 = \frac{N}{S}$

Warning: To generate error bars for $\hat{Y}(x_0)$ with (2a) and (2b) you will need to know both σ_ϵ^2 and σ_z^2 separately.

Kriging with the variogram

Recall the characterization of $\hat{Y}(x_0)$ given in (2a) and (2b) defined $\hat{Y}(x_0) = d^T \lambda$ where λ is defined by

$$\begin{aligned} & \underset{\lambda \in \mathbb{R}^n}{\operatorname{argmin}} E(Y(x_0) - d^T \lambda)^2 \\ &= \underset{\lambda \in \mathbb{R}^n}{\operatorname{argmin}} \begin{bmatrix} \lambda \\ -1 \end{bmatrix}^T \left[\begin{array}{c|c} \Sigma_{11} & \Sigma_{10} \\ \hline \Sigma_{01} & \Sigma_{00} \end{array} \right] \begin{bmatrix} \lambda \\ -1 \end{bmatrix} \\ &= \underset{\lambda \in \mathbb{R}^n}{\operatorname{argmin}} \left\{ \begin{bmatrix} \lambda \\ -1 \end{bmatrix}^T \left(K(z_i, z_j) \right)_{i,j=1}^{n+1} \begin{bmatrix} \lambda \\ -1 \end{bmatrix} + \sigma_e^2 \|\lambda\|_2^2 \right\} \end{aligned}$$

where $(z_1, \dots, z_{n+1}) = (x_1, \dots, x_n, x_0)$

Also recall that

$$\begin{aligned} \lambda \in \mathbb{R}^n &\Leftrightarrow f(x_1, \dots, x_n, x_0) \cdot \begin{bmatrix} \lambda \\ -1 \end{bmatrix} = 0 \\ &\Leftrightarrow f(z_1, \dots, z_{n+1}) \begin{bmatrix} \lambda \\ -1 \end{bmatrix} = 0 \end{aligned}$$

Therefore we get the same λ , i.e. same $\hat{Y}(x_0)$ and predicted sd, if we replace $K(x, y)$ with $\tilde{K}(x, y)$ so long as the following holds:

$$\begin{bmatrix} \lambda \\ -1 \end{bmatrix}^T \left(K(z_i, z_j) \right)_{i,j=1}^{n+1} \begin{bmatrix} \lambda \\ -1 \end{bmatrix} = \begin{bmatrix} \lambda \\ -1 \end{bmatrix}^T \left(\tilde{K}(z_i, z_j) \right)_{i,j=1}^{n+1} \begin{bmatrix} \lambda \\ -1 \end{bmatrix}$$

Whenever $f(z_1, \dots, z_{n+1}) \begin{bmatrix} \lambda \\ -1 \end{bmatrix} = 0$

Or more generally if $\forall n, \forall y_1, \dots, y_n \in \mathbb{R}^d$, and $\forall c_1, \dots, c_n \in \mathbb{R}$ we have

$$C^T \left(K(y_i, y_j) \right)_{i,j=1}^n c = C^T \left(\tilde{K}(y_i, y_j) \right)_{i,j=1}^n c$$

$$\text{whenever } f(y_1, \dots, y_n)^T c = 0$$

Example

Suppose $m=1$ and the only covariate we have is constant, i.e.

$$f_i(x) \equiv 1.$$

Then

$$f(y_1, \dots, y_n)^T c = 0 \Leftrightarrow c_1 + \dots + c_n = 0$$

Now consider cov fun $K(x, y)$ and

$$\begin{aligned} & \text{the variogram of } K(x, y) \\ \gamma(x, y) &= \overbrace{K(x, x) + K(y, y) - 2K(xy)} \\ K(x, y) &= \frac{1}{2} (K(xx) + K(yy) - \gamma(x, y)) \\ \tilde{K}(x, y) &= \frac{1}{2} (\gamma(x, a_0) + \gamma(y, a_0) - \gamma(x, y)) \end{aligned}$$

If $c_1 + \dots + c_n = 0$ then

$$\begin{aligned} & \sum_{k,p=1}^n c_k c_p K(y_k, y_p) \quad \text{Both } y_k \\ &= \frac{1}{2} \sum_{k,p=1}^n c_k c_p K(y_k, y_k) \quad \text{Both } y_p \\ &+ \frac{1}{2} \sum_{k,p=1}^n c_k c_p K(y_p, y_p) \\ &- \frac{1}{2} \sum_{k,p=1}^n c_k c_p \gamma(y_k, y_p) \end{aligned}$$

But since

$$\begin{aligned} \sum_{k,p=1}^n c_k c_p K(y_k, y_p) \\ = \frac{1}{2} \left(\sum_{p=1}^n c_p \right) \left(\sum_{k=1}^n c_k K(y_k, y_p) \right) \\ = 0 \end{aligned}$$

and similarly $\frac{1}{2} \sum_{k,p=1}^n c_k c_p K(y_p, y_p) = 0$

we have

$$\begin{aligned} \sum_{k,p=1}^n c_k c_p K(y_k, y_p) &= \sum_{k,p=1}^n c_k c_p \frac{(-1)}{2} \gamma(y_k, y_p) \\ &= \sum_{k,p=1}^n c_k c_p \tilde{K}(y_k, y_p) \end{aligned}$$

↑
same reasoning

Whenever $c_1 + \dots + c_n = 0$.

i.e. We just need to use the $\frac{(-1)}{2}$ variogram in place of $K(x, y)$ for Kriging and s.d. prediction, when one of the covariates $f_p(x) = 1$.

Notice also that if one has scale and variance parameters (ρ, σ^2) ... and use the variogram $\gamma(x, y) = \|x - y\|^{2\rho}$ (with $\rho = 1, p = 1$) then

$$\sigma^2 \gamma\left(\frac{x}{\rho}, \frac{y}{\rho}\right) = \frac{\sigma^2}{\rho^{2\rho}} \gamma(x, y)$$

↑

The separate ρ, σ merge into a single parameter

$$\sigma^2 / \rho^{2\rho}$$

Kriging with a generalized covariance function

It is pretty hard to extend the technique for variograms with $f_i(x) = 1$ to general covariates $f_p(x)$... with one exception:

$$\left\{ \begin{array}{l} \text{polynomials} \\ \text{of order } \leq q \end{array} \right\} \subseteq \text{span} \{f_1(x), \dots, f_m(x)\}. \quad (i)$$

For convenience let the monomials of order $\leq q$ be denoted

$$\underbrace{m_0(x), \dots, m_{M_q}(x)}_{=1}^{M_q}$$

So we seek functions $\tilde{K}(x, y)$ s.t.

$$f(y_1, \dots, y_n) \cdot c = 0 \text{ and } (i)$$

$$\downarrow \quad \sum_{k,p=1}^n c_k c_p K(y_k, y_p) = \sum_{k,p=1}^n c_k c_p \tilde{K}(y_k, y_p) \quad (ii)$$

holds for any $y_1, \dots, y_n \in \mathbb{R}^d$ and $c_1, \dots, c_n \in \mathbb{R}$.

A $\tilde{K}(x, y)$ with this property (for some positive definite K) is called a generalized covariance function of order q .

The official (but equivalent) definition is usually given as follows:

Definition: $\tilde{Y}(x, y): \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is a generalized covariance function on \mathbb{R}^d of order q

$$\left. \begin{array}{l} x_1, \dots, x_n \in \mathbb{R}^d, c_1, \dots, c_n \in \mathbb{R} \\ \text{and } \sum_{i=1}^n c_i m(x_i) = 0 \\ \text{for all monomials } m(\cdot) \\ \text{on } \mathbb{R}^d \text{ of order } \leq q \end{array} \right\} \Rightarrow \sum_{i,j=1}^n c_i c_j \tilde{Y}(x_i, x_j) \geq 0$$

Remark: Smaller $g \Rightarrow$ more stringent, in above definition. If \mathcal{G} satisfies above for g , then also for any $g' > g$.

To find a $\tilde{K} \neq K$ that satisfies (ii) recall how it worked in the previous section where $f_i(x) \equiv 1$. In this case we defined

$$\tilde{Z}(x) := Z(x) - Z(a_0)$$

which removed an overall constant and discovered $\tilde{K}(x,y) = \text{cov}(\tilde{Z}(x), \tilde{Z}(y))$ can be used instead of $K(x,y)$ to yield equivalent kriging predictions.

In the present case this suggests to define \tilde{Z} by removing all monomials of order $\leq g$ from Z so that at some fixed points

$$a_0, a_1, \dots, a_{M_g} \in \mathbb{R}^d$$

we have $\tilde{Z}(a_i) = 0$.

To this end define

$$\tilde{Z}(x) := Z(x) - \underbrace{\varphi^i(x)Z(a_i)}_{\substack{\text{Short for} \\ \sum_{i=0}^{M_g} \varphi^i(x)Z(a_i)}} \quad \left(\begin{array}{c} \\ \end{array} \right)$$

where $\varphi^i(x)Z(a_i)$ is given by

$$\left[\begin{array}{c} 1 \\ m_1(x) \\ \vdots \\ m_{M_g}(x) \end{array} \right]^T \left[\begin{array}{cc|c} 1 & m_1(a_0) & m_{M_g}(a_0) \\ | & | & | \\ 1 & m_1(a_n) & m_{M_g}(a_n) \end{array} \right]^{-1} \left[\begin{array}{c} Z(a_0) \\ | \\ Z(a_{M_g}) \end{array} \right]$$

$$= \left[\varphi^0(x), \varphi^1(x), \dots, \varphi^{M_g}(x) \right]$$

This is what we want since

$$\varphi^i(a_j)Z(a_i) = Z(a_j), \forall j$$

$$\therefore \tilde{Z}(a_j) = 0 \quad \forall j$$

Note: there is an implicit assumption here that the a_i 's are "spread out" enough so the above matrix inverse is defined.

Claim: $\tilde{K}(x,y) = \text{cov}(\tilde{Z}(x), \tilde{Z}(y))$ defined as above satisfies (ii).

Proof:

$$\begin{aligned} \tilde{K}(x,y) &= K(x,y) - \varphi^i(x)K(y, a_i) - \varphi^i(y)K(x, a_i) \\ &\quad + \varphi^i(x)\varphi^j(y)K(a_i, a_j) \end{aligned}$$

which can be written in operator notation

$$\tilde{K}(x,y) = \iint ds dt [s_{x-s} + \varphi^i(x)s_{a_i-s}] [t_{y-t} + \varphi^j(y)t_{a_j-t}] K(s,t) \quad \left(\begin{array}{c} \\ \end{array} \right)$$

where $\int ds s_{x-s} f(s) = f(x)$

Now suppose (i) and $f(y_1, \dots, y_n) \cdot c = 0$.

$$\therefore \sum_{i=1}^n m_j(y_i) c_i = 0 \quad \text{for any } m_j(x) \in \{m_0(x), \dots, m_{M_g}(x)\}$$

$$\therefore \sum_{i=1}^n \varphi^j(y_i) c_i = 0 \quad \text{for any } j \in \{0, 1, \dots, M_g\}$$

Therefore

$$\begin{aligned} \sum_{k,p=1}^n c_k c_p \tilde{K}(y_k, y_p) &= \sum_{k,p=1}^n c_k c_p K(y_k, y_p) \\ &\quad - \sum_{k,p=1}^n c_k c_p \varphi^i(y_k) K(y_p, a_i) \\ &\quad - \sum_{k,p=1}^n c_k c_p \varphi^i(y_p) K(y_k, a_i) \\ &\quad + \sum_{k,p=1}^n c_k c_p \varphi^i(y_k) \varphi^j(y_p) K(a_i, a_j) \end{aligned}$$

where

$$\sum_{k,p=1}^n c_k c_p \varphi^i(y_k) \varphi^j(y_p) K(a_i, a_j) \\ = K(a_i, a_j) \left(\underbrace{\sum_{k=1}^n c_k \varphi^i(y_k)}_{=0} \right) \left(\underbrace{\sum_{p=1}^n c_p \varphi^j(y_p)}_{=0} \right) = 0$$

$$\sum_{k,p=1}^n c_k c_p \varphi^i(y_k) K(y_p, a_i) \\ = \sum_{p=1}^n c_p K(y_p, a_i) \left(\underbrace{\sum_{k=1}^n c_k \varphi^i(y_k)}_{=0} \right) = 0$$

and similarly $\sum_{k,p=1}^n c_k c_p \varphi^i(y_p) K(y_k, a_i) = 0$.

$$\therefore \sum_{k,p=1}^n c_k c_p \tilde{K}(y_k, y_p) = \sum_{k,p=1}^n c_k c_p K(y_k, y_p)$$

as was to be shown.

□□□

Remark:

What is interesting about $\sum_{k,p=1}^n c_k c_p \tilde{K}(y_k, y_p)$ is that it can be simplified further by dropping any term in $\tilde{K}(x, y)$ that has the form $\|x - y\|^{2j}$ where $j \leq q$

the max monomial order.

This follows since

$$\|x - y\|^{2j} = \dots + m(x)m'(y) + \dots$$

where m and m' are monomials, one of which has has order $\leq j$.

$$\therefore \sum_{k,p=1}^n c_k c_p \|y_k - y_p\|^{2j} \\ = \dots + \left(\underbrace{\sum_{k=1}^n c_k m(y_k)}_{\text{either } 0 \text{ or } \text{is zero}} \right) \left(\underbrace{\sum_{p=1}^n c_p m'(y_p)}_{\text{is zero}} \right) + \dots \\ = 0 \quad (\text{iii})$$

This suggests that there exists other bivariate functions $\tilde{K}(x, y)$ which satisfy (ii) but which are not positive definite (except on increments that annihilate monomials of order $\leq q$)

A collection of generalized auto-cov fans

Suppose $K(t): \mathbb{R} \rightarrow \mathbb{R}$ is a symmetric function which has $2q$ finite derivatives at $t=0$.
Let

$$\mathcal{D}_g K(t) := K(t) - \sum_{k=0}^q \frac{K^{(2k)}(0)}{(2k)!} t^{2k}$$

Claim: If $K(t)$ is an auto-covariance function on \mathbb{R}^d then $\mathcal{D}_g K(t)$ is a generalized auto-covariance function on \mathbb{R}^d of order $m \geq q$

Proof: Let $x_1, \dots, x_n \in \mathbb{R}^d$ and $c_1, \dots, c_n \in \mathbb{R}$ such that

$$\sum_{i=1}^n c_i m(x_i) = 0$$

for all $m(\cdot) \in \{\text{monomials on } \mathbb{R}^d \text{ of order } \leq q\}$.

By (iii) we have that

$$\sum_{k,p=1}^n c_k c_p K(\|x_k - x_p\|) = \sum_{k,p=1}^n c_k c_p \mathcal{D}_g K(\|x_k - x_p\|)$$

≥ 0 by positive definiteness

≥ 0

□□□

Example:

$$\begin{aligned} -\frac{1}{2} \underbrace{\gamma(\|x-y\|)}_{\text{variogram}} &= K(\|x-y\|) - K(0) \\ &= D_g K(\|x-y\|) \text{ for } g=0 \end{aligned}$$

$\therefore -\frac{1}{2} \gamma(\|x-y\|)$ is a generalized auto-covariance of order 0.

Example:

Consider the Matérn auto-cov M_ν .

By properties of M_ν

$$|M_\nu^{(2g)}(0)| < \infty, \forall \text{ integer } g \leq \nu$$

$\Rightarrow D_{L\nu} M_\nu(\|x-y\|)$ is a generalized covariance function of order g for any $g \geq L\nu$

Also recall from Lecture 1 that

$$M_\nu(t) = \sum_{k=0}^{L\nu} b_k t^{2k} + c_0 \mathcal{Y}_\nu(t) + \mathcal{O}(t^{2L\nu+2})$$

$\underbrace{\phantom{b_k t^{2k}}}_{\text{order } 2L\nu}$ $\underbrace{\phantom{c_0 \mathcal{Y}_\nu(t)}}_{\text{order } 2L\nu+2}$

as $t \rightarrow \infty$ where $c_0 > 0$ and

$$\mathcal{Y}_\nu(t) := (-1)^{L\nu+1} \times \begin{cases} t^{2\nu} \log t & \text{if } \nu = 1, 2, \dots \\ t^{2\nu} & \text{otherwise} \end{cases}$$

These \mathcal{Y}_ν are also generalized auto-covariance functions. To see why let $x_1, \dots, x_n \in \mathbb{R}^d$ and $c_1, \dots, c_n \in \mathbb{R}$ s.t.

$$\sum_{i=1}^n c_i m(x_i) = 0$$

for all $m(\cdot) \in \{\text{monomials in } \mathbb{R}^d \text{ of order } \leq L\nu\}$.

Now

$$\begin{aligned} 0 &\leq \lim_{\varepsilon \downarrow 0} \sum_{i,j=1}^n c_i c_j \frac{M_\nu(\varepsilon \|x_i - x_j\|)}{\varepsilon^{2\nu}} \\ &= \lim_{\varepsilon \downarrow 0} \sum_{i,j=1}^n c_i c_j D_{L\nu,ij} \frac{M_\nu(\varepsilon \|x_i - x_j\|)}{\varepsilon^{2\nu}} \\ &\quad \text{since } \sum_{i,j=1}^n c_i c_j \|x_i - x_j\|^{2L\nu} = 0 \\ &\quad \text{for any } k \leq L\nu \end{aligned}$$

$$\begin{aligned} &= \lim_{\varepsilon \downarrow 0} \sum_{i,j=1}^n c_i c_j c_0 \frac{\mathcal{Y}_\nu(\varepsilon \|x_i - x_j\|)}{\varepsilon^{2\nu}} \\ &\quad + \lim_{\varepsilon \downarrow 0} \mathcal{O}\left(\frac{\varepsilon^{2L\nu+2}}{\varepsilon^{2\nu}}\right) \end{aligned}$$

Notice that $\lim_{\varepsilon \downarrow 0} \mathcal{O}\left(\frac{\varepsilon^{2L\nu+2}}{\varepsilon^{2\nu}}\right) = 0$.

Also if $\nu \notin \{1, 2, \dots\}$ then

$$\begin{aligned} \sum_{i,j=1}^n c_i c_j \frac{\mathcal{Y}_\nu(\varepsilon \|x_i - x_j\|)}{\varepsilon^{2\nu}} &= \sum_{i,j=1}^n c_i c_j (-1)^{L\nu+1} \|x_i - x_j\|^{2\nu} \\ &= \sum_{i,j=1}^n \mathcal{Y}_\nu(\|x_i - x_j\|) \end{aligned}$$

and if $\nu \in \{1, 2, \dots\}$ then

$$\begin{aligned} \sum_{i,j=1}^n c_i c_j \frac{\mathcal{Y}_\nu(\varepsilon \|x_i - x_j\|)}{\varepsilon^{2\nu}} &= \sum_{i,j=1}^n c_i c_j (-1)^{L\nu+1} \|x_i - x_j\|^{2\nu} \log \|x_i - x_j\| \\ &\quad + \sum_{i,j=1}^n c_i c_j (-1)^{L\nu+1} \|x_i - x_j\|^{2\nu} \log \varepsilon \end{aligned}$$

This is zero since $2\nu = 2L\nu$ and x_i, c_i annihilate polys of order $L\nu$.

$$= \sum_{i,j=1}^n c_i c_j \mathcal{Y}_\nu(\|x_i - x_j\|)$$

Putting this all together we get

$$0 \leq \lim_{\varepsilon \downarrow 0} \sum_{i,j=1} c_i c_j \frac{M_\nu(\varepsilon \|x_i - x_j\|)}{\varepsilon^{2\nu}}$$

$$= \sum_{i,j=1}^n c_i c_j c_\nu \underbrace{M_\nu(\|x_i - x_j\|)}_{\substack{\uparrow \\ \text{positive}}}$$

This establishes the following

Claim: For any $\nu > 0$ & $d \in \{1, 2, 3, \dots\}$

$$M_\nu(t) := (-1)^{\lfloor \nu \rfloor + 1} \times \begin{cases} t^{\lfloor \nu \rfloor} \log t & \text{if } \nu = 1, 2, \dots \\ t^{\lfloor \nu \rfloor} & \text{otherwise} \end{cases}$$

are generalized auto-covariance functions on \mathbb{R}^d of order $g \geq \lfloor \nu \rfloor$.

Pre-Kriging REML parameter estimation

As is usually the case, the covariance function $K(x, y)$ used for Kriging has unknown parameters that need to be estimated from the data.

Write θ for a generic parameter vector and $K_\theta(x, y)$ for the dependence of $K(x, y)$ on θ .

Recall the regression characterization of the data d_1, \dots, d_n

$$d = \begin{bmatrix} d_1 \\ \vdots \\ d_n \end{bmatrix} \sim N\left(f(x_1, \dots, x_n)^T \beta, \Sigma_\theta + \sigma_e^2 I\right)$$

↑
the matrix with
 $(i,j)^{\text{th}}$ entry $K_\theta(x_i, x_j)$

If we construct a matrix M such that

$$M f(x_1, \dots, x_n)^T = \text{Zero matrix}$$

then

$$Md \sim N\left(\underbrace{M f(x_1, \dots, x_n)^T \beta}_{=0}, M(\Sigma_\theta + \sigma_e^2 I)M^T\right)$$

so the likelihood of Md only depends on θ and not β . The REML (i.e. restricted max likelihood) estimate of θ is given by

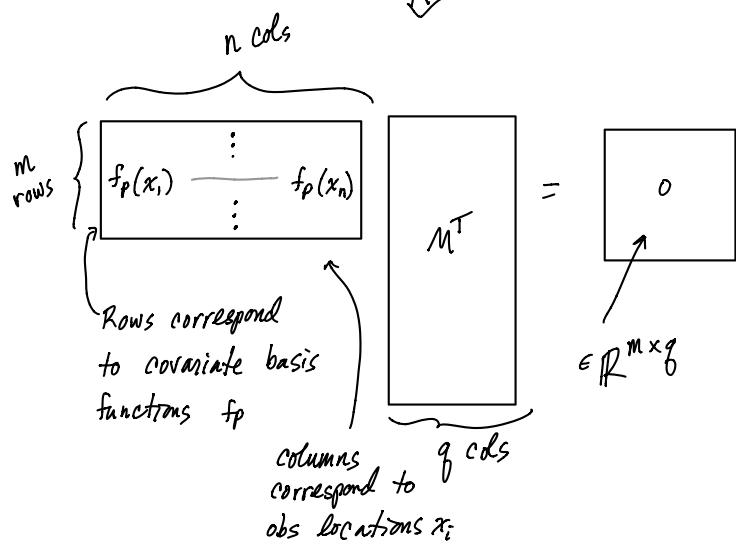
$$\hat{\theta}_{\text{REML}} = \arg \max_{\theta} \log P(Md | \theta).$$

Notice Md will necessarily be a lossy compression. Let's see how many degrees of freedom we drop:

$$M f(x_1, \dots, x_n)^T = \text{zero matrix}$$



$$f(x_1, \dots, x_n) M^T = \text{zero matrix}$$



So the best we can hope for is an M where the columns of M^T spans the null space of $f(x_1, \dots, x_n)$ which has dimension at least $n-m$, i.e.

$$M f(x_1, \dots, x_n)^T = \text{zero matrix}$$



$$\dim(Md) \geq n-m$$

Remark: A theme of the previous sections was that the final form of $\widehat{Y}(x_0)$ or $\text{var}(\widehat{Y}(x_0))$ doesn't require specifying everything about $K(x,y)$... just the parts that are not annihilated by linear combinations that satisfy

$$f(x_1, \dots, x_n) \cdot C = 0$$

or

$$f(x_1, \dots, x_n, x_0) \cdot \begin{bmatrix} x \\ -1 \end{bmatrix} = 0$$

Notice that REML estimation is no different since we are using the likelihood of

$$Md \sim N\left(0, M \underbrace{\left(\Sigma_0 + \sigma^2 \mathbb{I}\right)}_{\Sigma} M^T\right)$$

where $f(x_1, \dots, x_n) M^T = 0$

\Downarrow
M $\Sigma_0 M^T$ annihilates parts of K_0

REML + Kriging with \hat{Y}_j and parameter compression $\left(\frac{\sigma^2}{\rho^{2v}}, \sigma_\epsilon^2\right)$

Now we have a complete package:

- (cov) parameter estimation
- interpolation, $\hat{Y}(x_0)$
- standard errors, $sd(\hat{Y}(x_0))$

which is especially simplified using

- Generalized covariance Modeling

$$\text{with } K_G(x, y) = \sigma^2 C_v M_v \left(\frac{\|x-y\|}{\rho} \right)$$

$$= \frac{\sigma^2}{\rho^{2v}} C_v (2v)^v M_v(\|x-y\|)$$

when v is an integer the increments annihilate the lower order term.

So the generalized covariance parameter vector is given by $\theta = (\theta_1, \theta_2)$

$$\theta_1 = v$$

$$\theta_2 = \frac{\sigma^2}{\rho^{2v}}$$

- $\{f_1(x), \dots, f_m(x)\}$ contain all monomials (in x) of order $\leq [v]$

In this case there are then 3 parameters:

$$\theta_1 = v, \quad \theta_2 = \frac{\sigma^2}{\rho^{2v}}, \quad \sigma_\epsilon^2 = \text{nugget variance.}$$

Warning about interpreting $\hat{\beta}$ when using generalized covariance functions with Kriging

Recall that Kriging is a method for estimating

$$Y(x) = \sum_{p=1}^m f_p(x) \beta_p + Z(x)$$

where $Z \sim \text{GRF}(\sigma, K_\theta)$.

The fact that our goal is the sum of the covariate term plus $Z(x)$, in essence, allows us to conflate the fixed effects:

$$\sum_{p=1}^m f_p(x) \beta_p$$

with the random effects

$$Z(x)$$

This means that $\hat{\beta}$ computed from K_θ will be different than $\hat{\beta}$ computed from a generalized cov G_θ , even though the predictions of \hat{Y} yield the same result.

Explicit form of the Principle irregular term for $M_v(t)$

Recall that

$$M_v(t) = \sum_{k=0}^{\lfloor v \rfloor} b_k t^{2k} + c_v Y_v(t) + O(t^{\lfloor v \rfloor + 2})$$

$\underbrace{\phantom{b_k t^{2k}}}_{\text{order } 2k}$ $\underbrace{Y_v(t)}_{\text{order } \lfloor v \rfloor + 2}$

as $t \rightarrow 0$ where $c_v > 0$ and

$$Y_v(t) := (-1)^{\lfloor v \rfloor + 1} \times \begin{cases} t^{2v} \log t & \text{if } v = 1, 2, \dots \\ t^{2v} & \text{otherwise.} \end{cases}$$

It would be nice to know the constant term c_v so one can compare $M_v(t)$ to its principle irregular term $c_v Y_v(t)$.

The derivation is based on 9.6.2, 9.6.10 and 9.6.11 of Abramowitz and Stegun

case 1: $v \notin \mathbb{Z}$

$$t^v K_v(t)$$

$$= \frac{1}{2} \frac{\pi t^v}{\sin(v\pi)} (I_{-v}(z) - I_v(z)) \quad (9.6.2)$$

$$= \frac{1}{2} \frac{\pi t^v}{\sin(v\pi)} \left(\left(\frac{t}{2}\right)^{-1} \sum_{k=0}^{\infty} \frac{\left(\frac{t}{2}\right)^{2k}}{k! \Gamma(-v+k+1)} - \left(\frac{t}{2}\right)^v \sum_{k=0}^{\infty} \frac{\left(\frac{t}{2}\right)^{2k}}{k! \Gamma(v+k+1)} \right) \quad (9.6.10)$$

$$= \frac{1}{2} \frac{\pi}{\sin(v\pi)} \left(2^v \sum_{k=0}^{\infty} \frac{\left(\frac{t}{2}\right)^{2k}}{k! \Gamma(-v+k+1)} \right)$$

$$+ \frac{1}{2} \frac{\pi}{\sin(v\pi)} \left(-\frac{t^{2v}}{2^v \Gamma(v+1)} \right) \quad \text{Principle irregular term (PIT)}$$

$$\frac{1}{2} \frac{\pi}{\sin(v\pi)} \left(-\frac{1}{2^v} \sum_{k=1}^{\infty} \frac{\left(\frac{t}{2}\right)^{2k-2v}}{k! \Gamma(v+k+1)} \right)$$

\therefore

The PIT of $M_v(t)$

$$= \text{the PIT of } \frac{2^{1-v}}{\Gamma(v)} (\sqrt{v} t)^v K_v(\sqrt{v} t)$$

$$= \left(\frac{2^{1-v}}{\Gamma(v)} \right) \left(\frac{1}{2} \frac{\pi}{\sin(v\pi)} \right) \left(-\frac{(\sqrt{v} t)^{2v}}{2^v \Gamma(v+1)} \right)$$

$$= \frac{\pi}{\sin(v\pi)} \left(\frac{-(v/2)^v}{\Gamma(v) \Gamma(v+1)} \right) t^{2v}$$

Case 2: $v \in \mathbb{Z}$

$$t^v K_v(t) = (\text{low order} + \text{higher order}) \\ + t^v (-1)^{v+1} \log\left(\frac{t}{z}\right) \left(\frac{t}{z}\right)^v \frac{1}{r(v+1)}$$

by (9.6.11).

∴

The PIT of $M_v(t)$

$$= \text{the PIT of } \frac{2^{1-v}}{\Gamma(v)} (\sqrt{v}t)^v K_v(\sqrt{v}t)$$

$$= 2(-1)^v \left(\frac{-(v)_v}{\Gamma(v)\Gamma(v+1)} \right) \log(t) t^{2v}$$

In Summary

$$c_v Y_v(t) = \begin{cases} 2(-1)^v b_v \log(t) t^{2v}, & v \in \mathbb{Z} \\ \frac{\pi}{\sin(v\pi)} b_v t^{2v}, & v \notin \mathbb{Z} \end{cases}$$

where

$$b_v = \frac{-(v)_v}{\Gamma(v)\Gamma(v+1)}$$

