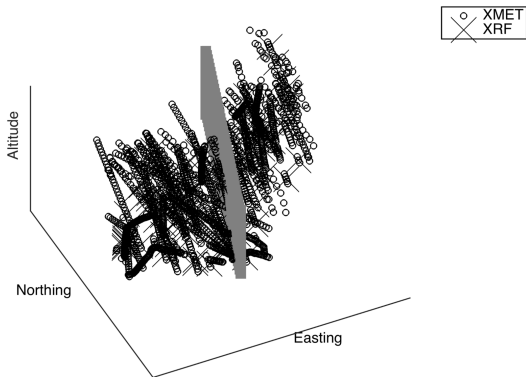


# Gaussian processes: Kriging and parameter estimation

Jo Eidsvik

Department of Mathematical Sciences, NTNU, Norway

# Grade prediction from boreholes



# Kriging interpolation

$$Y^*(\mathbf{s}_0) = \sum_{i=1}^n \alpha_i Y(\mathbf{s}_i) = \boldsymbol{\alpha}^t \mathbf{Y}$$

- ▶ Spatial interpolation from data  $Y(\mathbf{s}_1), \dots, Y(\mathbf{s}_n)$
- ▶ Best linear (spatial) predictor
- ▶ Kriging equals the optimal prediction for Gaussian model
- ▶ Unbiased and minimum prediction variance

## Versions

- ▶ Simple kriging  $E[Y(\mathbf{s})] = 0$ .
- ▶ Ordinary kriging  $E[Y(\mathbf{s})] = \mu$
- ▶ Universal kriging  $E[Y(\mathbf{s})] = \mathbf{h}(\mathbf{s})\boldsymbol{\beta}$
- ▶ Cokriging (Multivariate data  $Y_1(\mathbf{s}), \dots, Y_K(\mathbf{s})$ )

## Kriging derivation

$$\sigma_{s_0}^2 = \text{Var}[Y^*(s_0) - Y(s_0)] = E[Y^*(s_0) - Y(s_0)]^2 - \{E[Y^*(s_0)] - E[Y(s_0)]\}^2$$

'Mean Square Prediction Error' = 'Variance' + 'Bias squared'

$$\begin{aligned} \sigma_{s_0}^2 &= E[Y^*(s_0) - Y(s_0)]^2 = E\left[\sum_i \alpha_i Y(s_i) - Y(s_0)\right]^2 & (1) \\ &= E\left[\sum_i \sum_j \alpha_j \alpha_i Y(s_i) Y(s_j) - 2Y(s_0) \sum_i \alpha_i Y(s_i) + Y(s_0)^2\right] \\ &= \sum_i \sum_j \alpha_i \alpha_j C(i, j) - 2 \sum_i \alpha_i C(0, i) + C(0, 0) \\ &= \boldsymbol{\alpha}^t \mathbf{C} \boldsymbol{\alpha} - 2\boldsymbol{\alpha}^t \mathbf{C}_{0,\cdot} + C_0 \end{aligned}$$

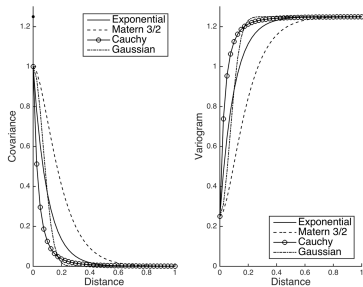
## Kriging derivation

Minimizing prediction error as a function of the weights  $\alpha_j$ .  
Optimal weights - derivative is 0 at the minimum.

$$\begin{aligned}\frac{d\sigma_{s_0}^2}{d\alpha} &= 2\mathbf{C}\alpha - 2\mathbf{C}_{0,\cdot} = 0 & (2) \\ \alpha &= \mathbf{C}^{-1}\mathbf{C}_{0,\cdot} \\ Y^*(s_0) &= \alpha^t \mathbf{Y} = \mathbf{C}_{0,\cdot}^t \mathbf{C}^{-1} \mathbf{Y}\end{aligned}$$

One can show the same with the variogram  
 $\gamma(\mathbf{h}) = \text{Var}(Y(\mathbf{s}) - Y(\mathbf{s} + \mathbf{h}))$ .

# Covariance functions and variograms



# Interpretation

$$Y^*(s_0) = \alpha^t \mathbf{Y} = \mathbf{C}_{0,\cdot}^t \mathbf{C}^{-1} \mathbf{Y}$$

- ▶ Unbiased linear predictor  $E[Y(s)] = 0$  for all  $s$ .
- ▶ Weights depend on  $\text{Cov}[Y(s_i), Y(s_0)]$ : Closer sites get larger weight
- ▶ Weights depend on  $\text{Cov}[Y(s_i), Y(s_j)]$ : Clustered sites get less weight



## Prediction variance

Plugging in optimal  $\alpha$  in  $\sigma_{s_0}^2$ .

$$\begin{aligned}
 \sigma_{s_0}^2 &= \alpha^t C \alpha - 2\alpha^t \mathbf{C}_{0,\cdot} + C_0 \\
 &= \mathbf{C}_{0,\cdot}^t C^{-1} C C^{-1} \mathbf{C}_{0,\cdot} - 2\mathbf{C}_{0,\cdot}^t C^{-1} \mathbf{C}_{0,\cdot} + C_0 \\
 &= C_0 - \mathbf{C}_{0,\cdot}^t C^{-1} \mathbf{C}_{0,\cdot}
 \end{aligned} \tag{3}$$

- ▶ Prediction variance is smaller than  $C_0$ .
- ▶ Decrease in prediction variance is larger close to data sites:  $\mathbf{C}_{0,\cdot}$  large.
- ▶ Prediction variance does not depend on data. It can be computed before the data acquisition.
- ▶ The spatial allocation of sites  $s_1, \dots, s_n$  is called 'spatial design'. This design impacts the prediction variance.

# Spatial regression model

Model:  $Y(\mathbf{s}) = \mathbf{h}(\mathbf{s})\boldsymbol{\beta} + w(\mathbf{s}) + \epsilon(\mathbf{s})$ .

1.  $Y(\mathbf{s})$  response variable at position  $\mathbf{s}$ .
2.  $\boldsymbol{\beta}$  regression effects.  $\mathbf{h}(\mathbf{s})$  covariates at  $\mathbf{s}$ .
3.  $w(\mathbf{s})$  zero-mean structured (spatially correlated) Gaussian process.
4.  $\epsilon(\mathbf{s})$  zero-mean unstructured (independent) Gaussian measurement noise.

## Gaussian model

Model:  $Y(\mathbf{s}) = \mathbf{h}(\mathbf{s})\boldsymbol{\beta} + w(\mathbf{s}) + \epsilon(\mathbf{s})$ .

Data at  $n$  locations:  $\mathbf{Y} = (Y(\mathbf{s}_1), \dots, Y(\mathbf{s}_n))'$ .

Likelihood:

$$l(\mathbf{Y}; \boldsymbol{\beta}, \boldsymbol{\theta}) = -\frac{1}{2} \log |\boldsymbol{\Sigma}| - \frac{1}{2} (\mathbf{Y} - \mathbf{H}\boldsymbol{\beta})' \boldsymbol{\Sigma}^{-1} (\mathbf{Y} - \mathbf{H}\boldsymbol{\beta})$$

$$\boldsymbol{\Sigma} = \text{Var}(\mathbf{w}) + \text{Var}(\boldsymbol{\epsilon}) = \mathbf{C} + \tau^2 \mathbf{I}$$

# Maximum likelihood

$$(\hat{\theta}, \hat{\beta}) = \operatorname{argmax}_{\theta, \beta} \{l(\mathbf{Y}; \beta, \theta)\}.$$

$$\hat{\theta}_{p+1} = \hat{\theta}_p - E \left( \frac{d^2 l(\mathbf{Y}; \hat{\beta}_p, \hat{\theta}_p)}{d\theta^2} \right)^{-1} \frac{dl(\mathbf{Y}; \hat{\beta}_p, \hat{\theta}_p)}{d\theta},$$

$$\hat{\beta}_p = \mathbf{A}^{-1} \mathbf{b}, \quad \mathbf{A} = \mathbf{A}(\hat{\theta}_p), \quad \mathbf{b} = \mathbf{b}(\hat{\theta}_p).$$

$$\mathbf{A} = \mathbf{H}' \boldsymbol{\Sigma}^{-1} \mathbf{H}$$

$$\mathbf{b} = \mathbf{H}' \boldsymbol{\Sigma}^{-1} \mathbf{Y}.$$

# Analytical derivatives

Formulas for matrix derivatives.

$$\begin{aligned}\mathbf{Q} &= \boldsymbol{\Sigma}^{-1} \\ \frac{d \log |\boldsymbol{\Sigma}|}{d\theta_r} &= \text{trace}\left(\mathbf{Q} \frac{d\boldsymbol{\Sigma}}{d\theta_r}\right) \\ \frac{d\mathbf{Z}'\mathbf{Q}\mathbf{Z}}{d\theta_r} &= -\mathbf{Z}'\mathbf{Q} \frac{d\boldsymbol{\Sigma}}{d\theta_r} \mathbf{Q}\mathbf{Z}.\end{aligned}$$

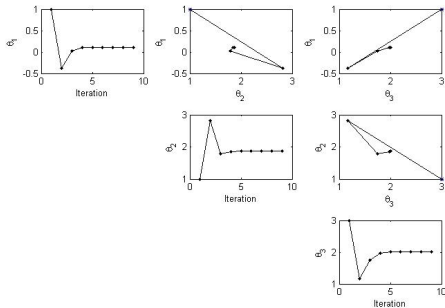
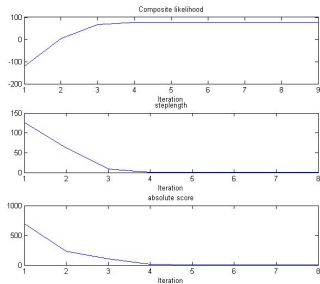
## Score and Hessian

$$\frac{dl}{d\theta_r} = -\frac{1}{2}\text{trace}\left(\mathbf{Q}\frac{d\boldsymbol{\Sigma}}{d\theta_r}\right) + \frac{1}{2}\mathbf{Z}'\mathbf{Q}\frac{d\boldsymbol{\Sigma}}{d\theta_r}\mathbf{Q}\mathbf{Z},$$

$$E\left(\frac{d^2l}{d\theta_r d\theta_s}\right) = -\frac{1}{2}\text{trace}\left(\mathbf{Q}\frac{d\boldsymbol{\Sigma}}{d\theta_s}\mathbf{Q}\frac{d\boldsymbol{\Sigma}}{d\theta_r}\right).$$

# Illustration maximization

Exponential covariance with nugget effect.  $\theta = (\theta_1, \theta_2, \theta_3)'$ : log **precision**, logistic **range**, log **nugget** precision.



# Properties

- ▶ maximum likelihood estimators are asymptotically unbiased.
- ▶ maximum likelihood estimators attain asymptotically minimum variance
- ▶ maximum likelihood estimators are asymptotically Gaussian distributed.



## Challenges

Model:  $Y(\mathbf{s}) = \mathbf{H}(\mathbf{s})\boldsymbol{\beta} + w(\mathbf{s}) + \epsilon(\mathbf{s})$ .

Data at  $n$  locations:  $\mathbf{Y} = (Y(\mathbf{s}_1), \dots, Y(\mathbf{s}_n))'$ .

Likelihood:

$$l(\mathbf{Y}; \boldsymbol{\beta}, \boldsymbol{\theta}) = -\frac{1}{2} \log |\boldsymbol{\Sigma}| - \frac{1}{2} (\mathbf{Y} - \mathbf{H}\boldsymbol{\beta})' \boldsymbol{\Sigma}^{-1} (\mathbf{Y} - \mathbf{H}\boldsymbol{\beta})$$

Challenges:

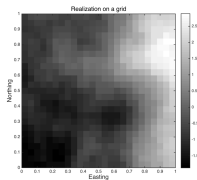
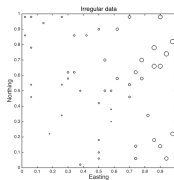
1. Build and store  $\boldsymbol{\Sigma}(\boldsymbol{\theta}) = \boldsymbol{\Sigma} = \mathbf{C} + \tau^2 \mathbf{I}_n$
2. Compute  $\log |\boldsymbol{\Sigma}|$
3. Compute  $\boldsymbol{\Sigma}^{-1}$  or  $(\mathbf{Y} - \mathbf{H}\boldsymbol{\beta})' \boldsymbol{\Sigma}^{-1} (\mathbf{Y} - \mathbf{H}\boldsymbol{\beta})$

## Possible solutions for large Gaussian models

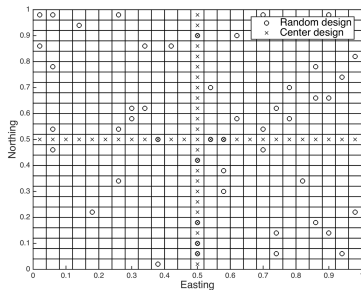
- ▶ Approximate likelihood (Fuentes 2007).
- ▶ Basis representation (Banerjee et al. 2008; Cressie and Johansson 2008).
- ▶ Markov representation (Lindgren et al. 2011).
- ▶ Tapered likelihood (Kaufman et al 2008).
- ▶ Composite likelihoods (Stein et al. 2004, Eidsvik et al 2014; Datta et al. 2016)
- ▶ Machine learning (Rasmussen and Williams 2006).
- ▶ Numerical linear algebra (Higham 2008, Aune et al., 2014).

# Example 1: Norwegian wood

$$Y(\mathbf{s}) = \mathbf{h}(\mathbf{s})\boldsymbol{\beta} + w(\mathbf{s}) + \epsilon(\mathbf{s})$$



# Data designs



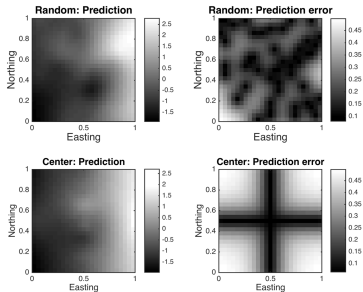
# Estimation: MLE

Table: Estimates(standard error).

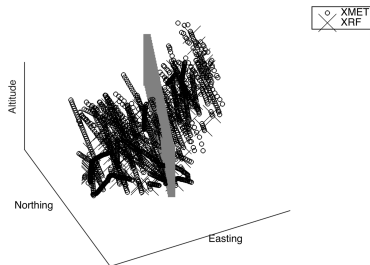
	$\beta_0$	$\beta_1$	$\beta_2$	$\sigma^2$	$\eta$	$\tau^2$
Center	-2.1 (0.6)	3.4 (0.7)	0.4 (0.7)	0.3 (0.14)	7.2 (2.0)	0.002 (0.001)
Random	-2.0 (0.5)	3.4 (0.6)	0.8 (0.5)	0.3 (0.12)	7.9 (2.0)	0.005 (0.007)
Truth	-2	3	1	0.25	9	0.0025

Matern covariance function.

# Predictions



## Example 2: Ore grade prediction in mining

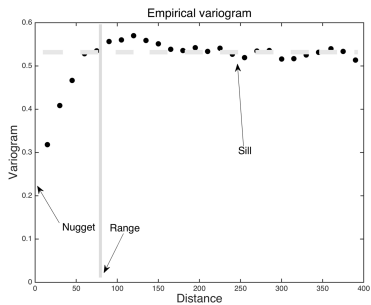


# Case

- ▶ Data at 1871 locations.
- ▶ Covariate is mineralization index (three possible classes)  
 $h(\mathbf{s}) = [1, \text{min.ind}(\mathbf{s})]$
- ▶ Spatial covariance is modeled by exponential covariance model.



# Variogram

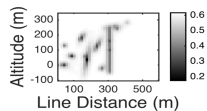
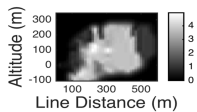


# Parameter estimation

Maximum likelihood (10 iterations of Fisher scoring.)

- ▶  $\beta_1 = 1.32$  (higher grades with mineralization index).
- ▶ Correlation range 50 m,  $\tau^2 = 0.45^2 = 0.2$ ,  $\sigma^2 = 0.62^2 = 0.38$ .

# Predictions

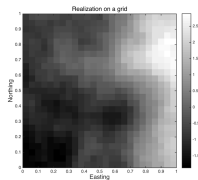
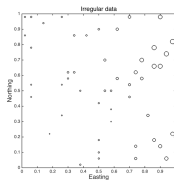


## Exercise: Norwegian wood

$$Y(\mathbf{s}) = \mathbf{h}(\mathbf{s})\boldsymbol{\beta} + w(\mathbf{s}) + \epsilon(\mathbf{s})$$

$$\mathbf{Y} \sim N(\mathbf{H}\boldsymbol{\beta}, \boldsymbol{\Sigma}(\boldsymbol{\theta}))$$

Different covariate.



## Spatio-temporal model

Model:  $Y(\mathbf{s}, t) = \mathbf{h}(\mathbf{s}, t)\boldsymbol{\beta} + w(\mathbf{s}, t) + \epsilon(\mathbf{s}, t)$ .

1.  $Y(\mathbf{s}, t)$  response variable at position  $\mathbf{s}$  at time  $t$ .
2.  $\boldsymbol{\beta}$  regression effects.  $\mathbf{h}(\mathbf{s}, t)$  covariates at  $\mathbf{s}$  at time  $t$ .
3.  $w(\mathbf{s}, t)$  zero-mean structured (spatio-temporally correlated) Gaussian process.
4.  $\epsilon(\mathbf{s}, t)$  zero-mean unstructured (independent) Gaussian measurement noise.

## Spatio-temporal statistics

Model:  $Y(\mathbf{s}, t) = \mathbf{h}(\mathbf{s}, t)\boldsymbol{\beta} + w(\mathbf{s}, t) + \epsilon(\mathbf{s}, t)$ .

Data at  $n_t$  locations at time  $t$ :  $\mathbf{Y}_t = (Y(\mathbf{s}_1, t), \dots, Y(\mathbf{s}_{n_t}, t))'$ ,

$t = t_1, t_2, \dots, t_n$ .

Goals could include:

- ▶ Estimate parameters: regression, noise structure in space and time, and noise of measurements.
- ▶ Characterize process in space and time: Smoothing, given all data. Filtering, given only current data. Prediction, given some data - look ahead in time. Interpolation (Kriging) in space.

## Common assumptions

Covariates  $h(\mathbf{s}, t)$  help include trends in space (say altitude, land-cover, etc.) or over time (hour, season, climate change, etc.), or coupling of space-time.

Covariance structure of  $w(\mathbf{s}, t)$  is

- ▶ Stationary in space and time:  $\text{Var}(w(\mathbf{s}, t)) = \text{Var}(w(\mathbf{s}', t'))$ ,  
 $\text{Corr}(w(\mathbf{s}, t), w(\mathbf{s}', t')) = \text{Corr}(w(\mathbf{s} + \mathbf{s}_0, t + t_0), w(\mathbf{s}' + \mathbf{s}_0, t' + t_0))$ .
- ▶ Separable in space and time:  
 $\text{Corr}(w(\mathbf{s}, t), w(\mathbf{s}', t')) =$   
 $\text{Corr}_s(w(\mathbf{s}, t), w(\mathbf{s}', t))\text{Corr}_t(w(\mathbf{s}, t), w(\mathbf{s}, t'))$ .

# Autoregressive spatial process

Markov in time and stationary:

$$w(\mathbf{s}, t) = \phi w(\mathbf{s}, t-1) + \delta(\mathbf{s}, t), \quad \text{Var}(w(\mathbf{s}, 0)) = \boldsymbol{\Sigma}_0, \quad \text{Var}(\delta(\mathbf{s}, t)) = (1 - \phi^2) \boldsymbol{\Sigma}_0$$

$\delta(\mathbf{s}, t)$  are independent in time.

This means that the one-step time correlation is  $\phi$  (and separable from space).



# Advection-diffusion equation

Structure process defined via a partial differential equation:

$$\frac{dw(\mathbf{s}, t)}{dt} = -\boldsymbol{\mu}^t \nabla w(\mathbf{s}, t) + \nabla \mathbf{D} \nabla w(\mathbf{s}, t) + \delta(\mathbf{s}, t),$$

$\boldsymbol{\mu}$  is advection (drift) term.  $\mathbf{D}$  is diffusion term,  $\delta$  is independent Gaussian noise.

Time-difference scheme gives Markovian process in time.

# Advection-diffusion equation

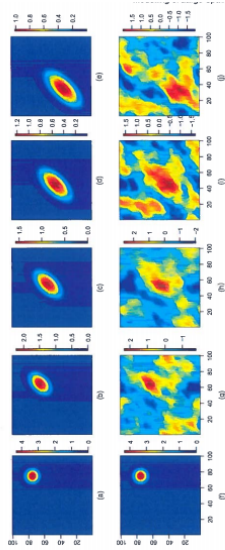


Fig. 1. Illustration of the SPDE (1) and the corresponding PDE (the drift vector points from north-west to south-west and the diffusive part exhibits anisotropy in the same direction; the same parameters are used for both the PDE and the SPDE, i.e.  $\zeta = -\log(0.99)$ ,  $\rho_1 = 0.06$ ,  $\gamma = 3$ ,  $\theta = 7/4$ ,  $\mu_x = -0.1$  and  $\mu_y = -0.1$ , and for the stochastic part of the SPDE, without stochastic term  $\beta(z, \mathbf{x})$ : (f)–(i) one sample from the distribution specified by the SPDE with a tiled initial condition; (a), (f)  $t = 1$ ; (b), (g)  $t = 2$ ; (c), (h)  $t = 3$ ; (d), (i)  $t = 4$ ; (e), (j)  $t = 5$ .

## Exercise Spatial case : Simulate and re-estimate

$$\mathbf{Y} = \mathbf{H}\boldsymbol{\beta} + \mathbf{LN}(0, \mathbf{I})$$

Cholesky matrix:

$$\mathbf{L}\mathbf{L}^t = \boldsymbol{\Sigma}(\boldsymbol{\theta})$$

- ▶ Fix parameters.
- ▶ Simulate a realization of  $\mathbf{Y}$  at data locations jointly with variables  $\mathbf{Y}_0$  at prediction locations.
- ▶ Fit parameters (maximum likelihood) from data  $\mathbf{Y}$ .
- ▶ Predict,  $\hat{\mathbf{Y}}_0$  (with uncertainty) given data and parameters.

## Exercise Spatial case : Joint Gaussian

$$\begin{pmatrix} \mathbf{Y}_0 \\ \mathbf{Y} \end{pmatrix} \sim N \left[ \begin{pmatrix} \mathbf{H}_0 \\ \mathbf{H} \end{pmatrix} \boldsymbol{\beta}, \begin{pmatrix} \boldsymbol{\Sigma}_0 & \boldsymbol{\Sigma}_{0,\cdot} \\ \boldsymbol{\Sigma}_{\cdot,0} & \boldsymbol{\Sigma} \end{pmatrix} \right]$$

$$[\mathbf{Y}_0 | \mathbf{Y}] \sim N(\mathbf{H}_0 \boldsymbol{\beta} + \boldsymbol{\Sigma}_{0,\cdot} \boldsymbol{\Sigma}^{-1} (\mathbf{Y} - \mathbf{H} \boldsymbol{\beta}), \boldsymbol{\Sigma}_0 - \boldsymbol{\Sigma}_{0,\cdot} \boldsymbol{\Sigma}^{-1} \boldsymbol{\Sigma}_{\cdot,0})$$