

Lesson 8:

Generalised Least Squares

Assumptions of linear models

$$y_{ijk} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + e_{ijk}$$

Linear models make many assumptions, including:

1. The model makes biological sense/ physical sense
2. Additivity (terms are added together)
3. Linearity
4. Independence of errors
5. Homoscedasticity – equal variance of errors
6. Normality of errors.

Two situations where GLS is useful


1. Data heteroscedasticity (unclear mean-variance relationships)
2. Autocorrelation between data points (non-independence)

Linear model with 'good' errors (IID)

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \text{ where } \boldsymbol{\varepsilon} \sim N(\mathbf{0}, \mathbf{V})$$

Where $\mathbf{V} = \text{Cov}[\boldsymbol{\varepsilon} | \mathbf{X}]$, the covariance matrix of errors. Ideally:

1. ε_i have even range across predicted y (homoscedasticity)
2. off-diagonal elements of \mathbf{V} are zero (no autocorrelation)

$$\mathbf{Y} = \mathbf{X} * \boldsymbol{\beta} + \boldsymbol{\varepsilon} \text{ where } \boldsymbol{\varepsilon} \sim N(0, \sigma^2 \mathbf{I})$$


The following matrix equation illustrates the linear model with its components:

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_{n-1} \\ y_n \end{bmatrix} = \begin{bmatrix} \text{Int} & \text{X1} & \text{X2} \\ 1 & x_{1,1} & x_{2,1} \\ 1 & x_{1,2} & x_{2,2} \\ 1 & x_{1,3} & x_{2,3} \\ \vdots & \vdots & \vdots \\ 1 & x_{1,n-1} & x_{2,n-1} \\ 1 & x_{1,n} & x_{2,n} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \varepsilon_3 \\ \vdots \\ \varepsilon_{n-1} \\ \varepsilon_n \end{bmatrix}$$

where

$$\begin{bmatrix} \varepsilon_{1,1} & 0 & 0 & \dots & 0 & 0 \\ 0 & \varepsilon_{2,2} & 0 & \dots & 0 & 0 \\ 0 & 0 & \varepsilon_{3,3} & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ 0 & 0 & 0 & \dots & \varepsilon_{n-1,n-1} & 0 \\ 0 & 0 & 0 & \dots & 0 & \varepsilon_{n,n} \end{bmatrix}$$

Linear model with heteroscedasticity and/ or error covariance

In reality $\underline{\varepsilon} \sim N(0, \mathbf{V})$, where :

1. ε_i are heteroscedastic
2. Off-diagonal elements are non-zero (covariance)

$$\mathbf{Y} = \mathbf{X} * \boldsymbol{\beta} + \boldsymbol{\varepsilon} \quad \text{where} \quad \underline{\varepsilon} \sim N(0, \mathbf{V})$$

X

Int X1 X2

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_{n-1} \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_{1,1} & x_{2,1} \\ 1 & x_{1,2} & x_{2,2} \\ 1 & x_{1,3} & x_{2,3} \\ \vdots & \vdots & \vdots \\ 1 & x_{1,n-1} & x_{2,n-1} \\ 1 & x_{1,n} & x_{2,n} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \varepsilon_3 \\ \vdots \\ \varepsilon_{n-1} \\ \varepsilon_n \end{bmatrix}$$

where

$$\begin{bmatrix} w_1 \varepsilon_{1,1} & \varepsilon_{1,2} & \varepsilon_{1,3} & \dots & \varepsilon_{1,n-1} & \varepsilon_{1,n} \\ \varepsilon_{2,1} & w_2 \varepsilon_{2,2} & \varepsilon_{2,3} & \dots & \varepsilon_{2,n-1} & \varepsilon_{2,n} \\ \varepsilon_{3,1} & \varepsilon_{3,2} & w_3 \varepsilon_{3,3} & \dots & \varepsilon_{3,n-1} & \varepsilon_{3,n} \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ \varepsilon_{n-1,1} & \varepsilon_{n-1,2} & \varepsilon_{n-1,3} & \dots & \varepsilon_{n-1,n-1} & \varepsilon_{n-1,n} \\ \varepsilon_{n,1} & \varepsilon_{n,2} & \varepsilon_{n,3} & \dots & \varepsilon_{n,n-1} & w_n \varepsilon_{n,n} \end{bmatrix}$$

Generalised least squares models

- The problem: $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ where $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \mathbf{V})$
- OLS solution: $\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} \quad \because \mathbf{V} = \mathbf{I}\sigma^2$
- The GLS solution: $\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{V}^{-1}\mathbf{y}$
- You might recall this from the mixed models lecture.

- We will deal with heteroscedasticity and error covariance separately in the remaining problems

To be clear:

- Heteroscedasticity is about the diagonal elements (homoscedastic residuals should be ~equal in size)
- Error covariance is about the off-diagonal elements (should be close to zero i.e. uncorrelated)

1. Heteroscedascity

$$\begin{aligned}
 & \mathbf{Y} = \mathbf{X} * \boldsymbol{\beta} + \boldsymbol{\varepsilon} \quad \text{where} \quad \boldsymbol{\varepsilon} \sim N(0, w\sigma^2\mathbf{I}) \quad \text{X} \\
 & \begin{array}{c} \text{Int} \quad \text{X1} \quad \text{X2} \\ \left[\begin{array}{c} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_{n-1} \\ 1 \\ y_n \end{array} \right] = \left[\begin{array}{ccc} 1 & x_{1,1} & x_{2,1} \\ 1 & x_{1,2} & x_{2,2} \\ 1 & x_{1,3} & x_{2,3} \\ \vdots & \vdots & \vdots \\ 1 & x_{1,n-1} & x_{2,n-1} \\ 1 & x_{1,n} & x_{2,n} \end{array} \right] \left[\begin{array}{c} \beta_0 \\ \beta_1 \\ \beta_2 \end{array} \right] + \left[\begin{array}{c} \varepsilon_1 \\ \varepsilon_2 \\ \varepsilon_3 \\ \vdots \\ \varepsilon_{n-1} \\ \varepsilon_n \end{array} \right] \end{array} \quad \text{where}
 \end{aligned}$$

$$\left[\begin{array}{cccccc} w_1\varepsilon_{1,1} & 0 & 0 & \dots & 0 & 0 \\ 0 & w_2\varepsilon_{2,2} & 0 & \dots & 0 & 0 \\ 0 & 0 & w_3\varepsilon_{3,3} & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & w_{n-1}\varepsilon_{n-1,n-1} & 0 \\ 0 & 0 & 0 & \dots & 0 & w_n\varepsilon_{n,n} \end{array} \right]$$

1. Dealing with heteroscedascity

- Usually we try to define relationships between means and variances, e.g. Poisson, Binomial or Gamma distributions.
- Sometimes we might need more flexibility if we are uncertain about the mean-variance relationship
- CAVEAT: I find it quite difficult to decide when to apply these methods UNLESS I already know that the mean-variance relationship is systematically skewed and cannot be described by a known distribution

Generalised Least Squares

- GLS extends the linear model to address complex error relationships

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \text{ where } \boldsymbol{\varepsilon} \sim \mathbf{N}(\mathbf{0}, \mathbf{I}w_i\sigma^2)$$

- We aim to apply inverse weights, $1/w_i$, that are the inverse of the relationship $\mathbf{I}w_i\sigma^2$
- Can model these w_i using the weights argument in the `gls()` function

The gls() function

- `mod <- gls(y~x, data=dat,
weights=varFixed(form=~1/w2))`

Variance function. Options
include:

- 1) `varFixed`: weights are linear
- 2) `varPower`: weights are power-law
- 3) `varIdent`: weights differ by factor level

Express the weights
in terms of the
covariate

The lme () function

- The lme() function allows one to include random effects into the mix as well
- `mod <- lme(y~x, random=~1|group, data=dat, weights=varFixed(form=~1/w2))`

Code 8.1

Heteroscedastic models

Exercise 8.1

- Return to the plantdamage data.
- Fit a random effects model with lme to test whether damage, light and their interaction affects growth.
- Use growth in untransformed form.
- First, run the model assuming no covariance structure (mod 1) and plot the residuals.
- Inspect the lme model for evidence of heteroscedasticity. Try this command
Ps: `plot(mod, resid(.)~fitted(.)|light, abline=0, layout=c(2, 1))`
- Then try varIdent and varExp models (mod 1b, mod 1c) and compare them against your first model using AIC and plot their residuals.
- Second, assume growth is exponentially distributed. Transform it and run the model assuming no co-variance structure (mod 2), and plot the residuals.
- Then try varIdent and varExp models with the transformed growth and compare them using AIC. (mod 2b, mod 2c)
- Finally compare the summaries from mod 1c and mod 2.

2. Autocorrelation

$$\begin{aligned}
 & \mathbf{Y} = \mathbf{X} * \boldsymbol{\beta} + \boldsymbol{\varepsilon} \quad \text{where} \quad \boldsymbol{\varepsilon} \sim N(0, \mathbf{V})
 \end{aligned}$$

X

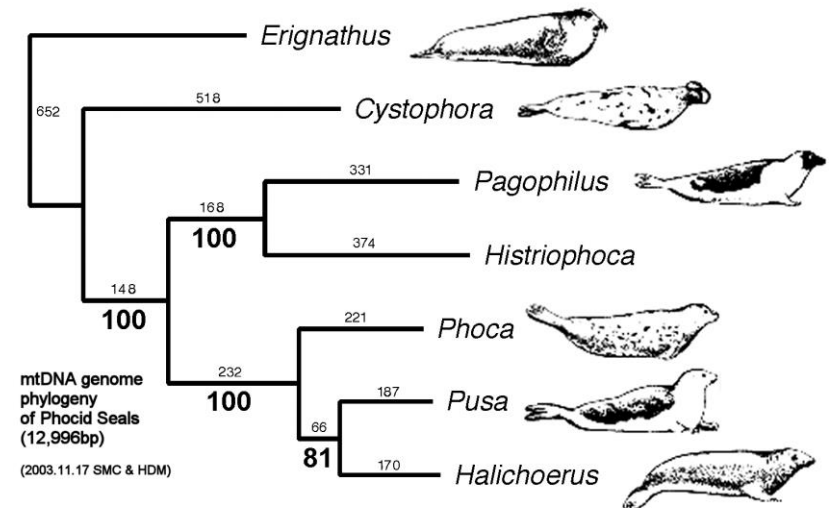
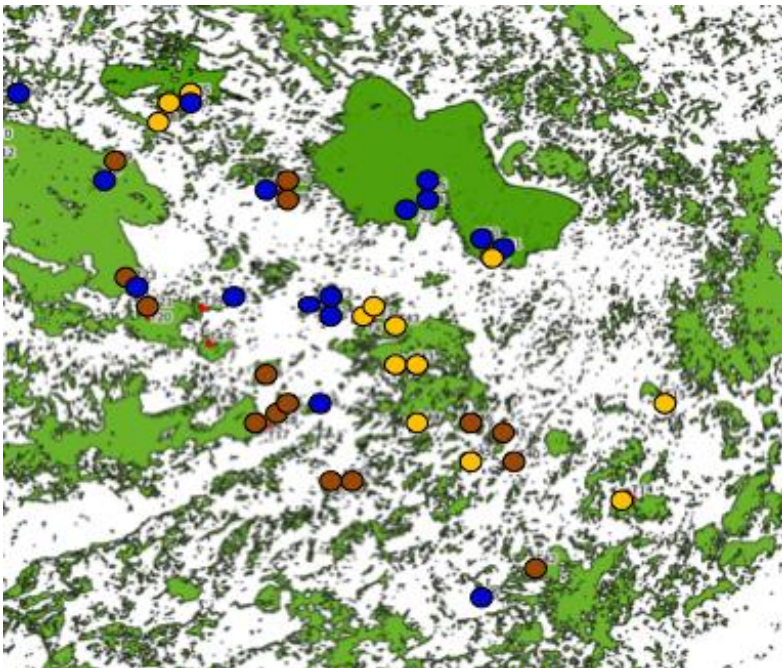
	Int	X1	X2
y_1	1	$x_{1,1}$	$x_{2,1}$
y_2	1	$x_{1,2}$	$x_{2,2}$
y_3	1	$x_{1,3}$	$x_{2,3}$
\vdots		\vdots	\vdots
y_{n-1}	1	$x_{1,n-1}$	$x_{2,n-1}$
y_n	1	$x_{1,n}$	$x_{2,n}$

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_{n-1} \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_{1,1} & x_{2,1} \\ 1 & x_{1,2} & x_{2,2} \\ 1 & x_{1,3} & x_{2,3} \\ \vdots & \vdots & \vdots \\ 1 & x_{1,n-1} & x_{2,n-1} \\ 1 & x_{1,n} & x_{2,n} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \varepsilon_3 \\ \vdots \\ \varepsilon_{n-1} \\ \varepsilon_n \end{bmatrix} \quad \text{where}$$

$\varepsilon_{1,1}$	$\varepsilon_{1,2}$	$\varepsilon_{1,3}$...	$\varepsilon_{1,n-1}$	$\varepsilon_{1,n}$
$\varepsilon_{2,1}$	$\varepsilon_{2,2}$	$\varepsilon_{2,3}$...	$\varepsilon_{2,n-1}$	$\varepsilon_{2,n}$
$\varepsilon_{3,1}$	$\varepsilon_{3,2}$	$\varepsilon_{3,3}$...	$\varepsilon_{3,n-1}$	$\varepsilon_{3,n}$
\vdots	\vdots	\vdots		\vdots	\vdots
\vdots	\vdots	\vdots		\vdots	\vdots
$\varepsilon_{n-1,1}$	$\varepsilon_{n-1,2}$	$\varepsilon_{n-1,3}$...	$\varepsilon_{n-1,n-1}$	$\varepsilon_{n-1,n}$
$\varepsilon_{n,1}$	$\varepsilon_{n,2}$	$\varepsilon_{n,3}$...	$\varepsilon_{n,n-1}$	$\varepsilon_{n,n}$

2. Dealing with auto-correlation

- Data points that are close together in space, time or phylogeny are not independent.



2. Dealing with auto-correlation

- Data points that are close together in space, time or phylogeny are not independent.
- Difficult to turn these relationships into discrete blocks (as we did for LMMs).
- Possible to model the auto-correlation directly with `gls()` and `lme()` using the correlation argument.

Covariance structures in gls

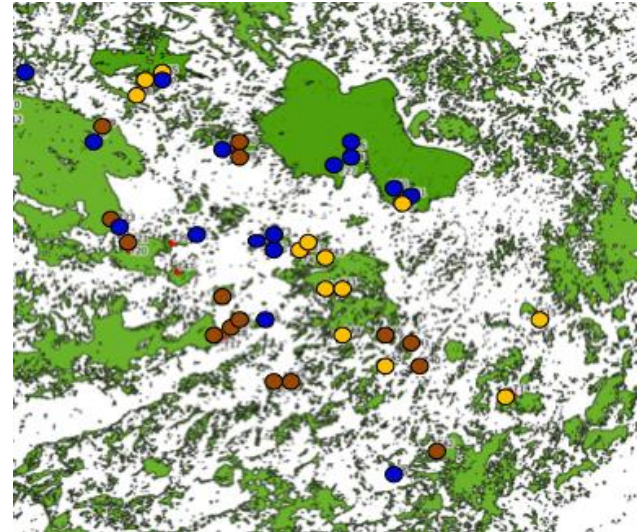
- All linear models have the form
- $Y = X\beta + \varepsilon$ such that $\varepsilon \sim N(0, V)$

Where Σ is a $n \times n$ matrix.

- Each value, from row i and column j is the covariance of the errors of replicates i and j
- OLS models assume that $V \sim I\sigma^2$ (in other words: off-diagonal covariances are close to zero)
- In GLS we essentially describe the relationships in V using an assumed covariance structure

Covariance between data

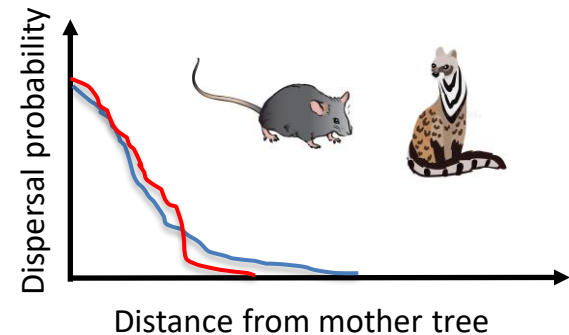
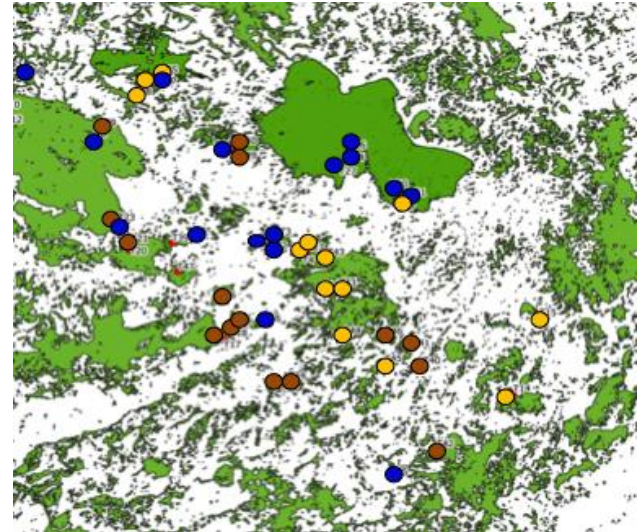
- Consider the correlation between plant communities in space.
- In the figure, it is clear that closer communities are more similar.
- This could be because of geology
- But it could also be because plants have a higher chance of establishing offspring closer to a mother tree
- This relatedness effect falls off as two points move further apart..
- i.e. the correlation between points decreases with distance



Covariance between data

Take homes:

1. Spatial autocorrelation is due to some process that we may find difficult to diagnose properly
2. The autocorrelation process often has a spatial limit
e.g. max distance of seed dispersal
3. The autocorrelation-dist relationship has a shape!



Describing covariance between data points: Semivariograms

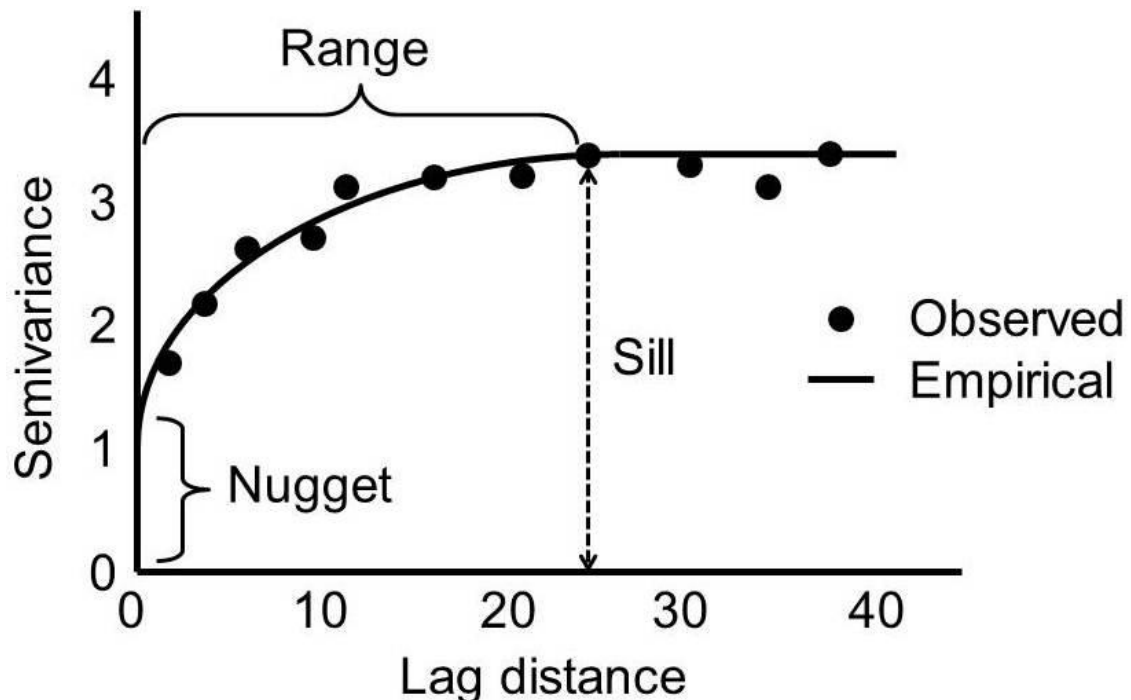
Semi-variograms describe differences (variance) between data points in terms of some distance measure (space, time, DNA).

$$\hat{\gamma}(h \pm \delta) := \frac{1}{2|N(h \pm \delta)|} \sum_{(i,j) \in N(h \pm \delta)} |z_i - z_j|^2$$

Range: distance over which points are correlated in space/ time

Sill: average distance between uncorrelated data points = variance

Nugget: some variance at small scale (stochasticity?/ measurement error?)



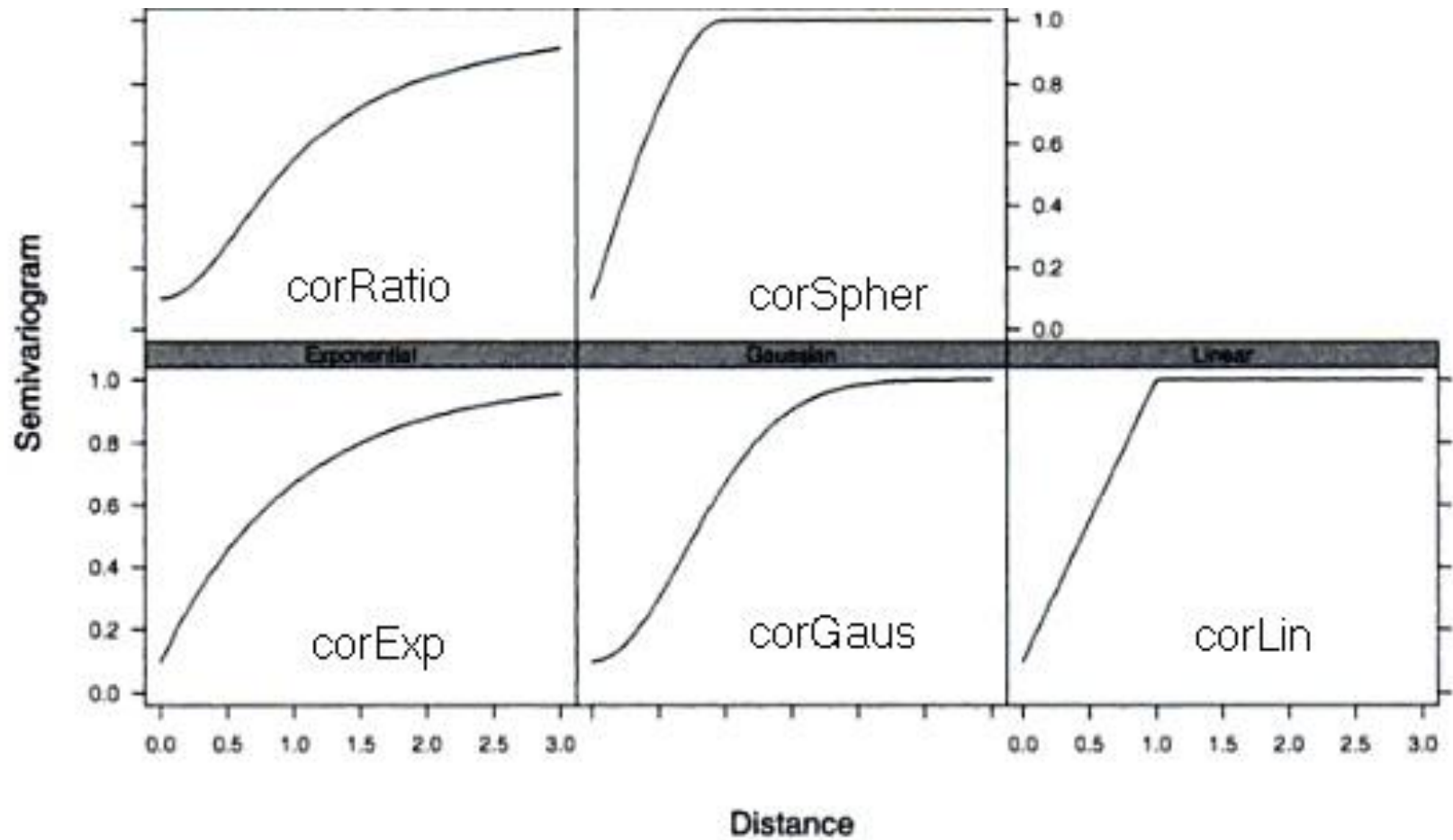
Fitting a correlation structure

- Distance correlation relationships have pattern;

```
mod <- gls(y~x, data=dat, correlation=corExp(form=~x+y))
```

- Dependency decreases as an exponential function of distance.
- Other options:
 - corSpher
 - corGaus
 - corLin
 - corRatio

Correlation structures for shapes

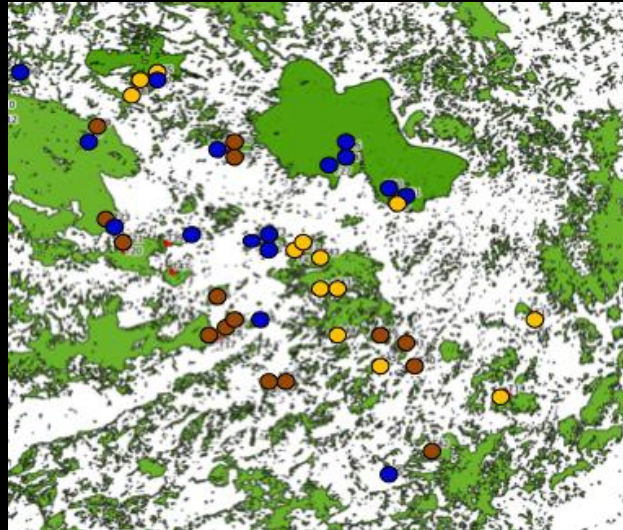


Approach

- 1) fit a basic model without a correlation structure,
- 2) examine a variogram of the residuals (i.e., their correlation with respect to distance),
- 3) choose an appropriate corStruct class based on the variogram shape (maybe test a few..), and
- 4) test whether addition of the correlation term significantly improves model fit via AICs or LRTs (**anova** function).
- Use **Variogram** function in **nlme** to diagnose the type:
`plot(Variogram(gls1,form=~lon+lat))`

Code 8.2

Spatial correlation structures



Exercise 8.2

- The data in the file trees.csv come from an analysis of tree growth in an Austrian forest.
- The response of Relative growth rate (RGR) to light was compared between two species (beech and spruce)
- Determine if there is an interaction between species and light in their effects on RGR
- Are the residuals independent?
- Modify the models to deal with any spatial dependence