

- I will start with a simple model, using species diversity data
- Strong spatial dependence, $\hat{\rho} = 0.79$
- what is the mean diversity? How precise is our estimate?
- Sampling discussion:
 - The 64 squares are a systematic sample: Only have 1 sample
 - If treat as a simple random sample,
 - usual se probably wrong
- So shift to model-based inference
 - If assume each square independent, $\hat{\mu} = 6.09$, $se \hat{\mu} = 0.36$
 - But we know that model is wrong
 - If we assume spatially correlated, what changes?

Spatial inference

- If recognize the spatial correlation between nearby squares, $\hat{\mu} = 5.69$, $se \hat{\mu} = 1.69$
 - Difference in estimates is not too big (7% change)
 - Difference in se's is huge (470% change)
- This data set has 64 obs
- I find it very helpful to compare sample sizes for equal se's, rather than directly compare se's.
 - So, under the correlated data model, how many obs. needed to get se = 0.36?
 - A: 1410! Huge number! Spatial correlation matters!
 - (BTW, also for geostatistical models)
- Computed as $64 * (1.69/0.36)^2$
 - Correlated model: Variance from 64 obs is $1.69^2 = 2.856$
 - Variance from 128 obs will be $2.856/2$.
 - Variance from N obs will be $2.856/(N/64)$. Want = 0.36^2 . Solve for N

Accounting for spatial correlation

- How did I get estimates and se's "recognizing the spatial correlation"?
- Three common approaches
 - 1) Simultaneous Autoregressive (SAR) Model
 - 2) Conditional Autoregressive (CAR) Model
 - 3) geostatistical model
- Geostatistical model can be used for either point or areal data
- We will consider each in turn.
- First need to talk about:
 - Describing correlated data (multivariate normal distributions)
 - Regression using matrices
 - Regression with correlated data

Multivariate normal distributions

- Usual 401 setup: $Y_i \stackrel{iid}{\sim} N(\mu_i, \sigma^2)$
 - Mean, μ_i , for each observation may be constant, dependent on treatment, or unique to each observation (e.g., $\beta_0 + \beta_1 X_i$)
 - Variance same for each obs. (no subscript on σ^2)
 - Independent observations
- Now, need to describe correlations among pairs of observations
- $Y_i \sim N(\mu_i, \sigma^2)$ is not sufficient
- Towards that goal: collect all the observations in a vector \mathbf{Y}

$$\mathbf{Y} = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}, \mathbf{Y}' = [Y_1 \ Y_2 \ \cdots \ Y_n]$$

- Collect the means into a vector: $\mu' = [\mu_1 \mu_2 \cdots \mu_n]$
- Variance now becomes the variance-covariance matrix
- Consider 3 random variables, X, Y and Z
- VC matrix is a 3 x 3 matrix

$$\Sigma = \begin{bmatrix} \sigma_X^2 & \sigma_{XY} & \sigma_{XZ} \\ \sigma_{XY} & \sigma_Y^2 & \sigma_{YZ} \\ \sigma_{XZ} & \sigma_{YZ} & \sigma_Z^2 \end{bmatrix}$$

- When n values of Y , VC matrix is $n \times n$

Variances and covariances

- Diagonal values are variances:

$$\sigma_X^2 = \text{Var } X = E (X - \mu_X)^2$$

- Off diagonal values are covariances:

$$\sigma_{XY} = \text{Cov } X, Y = E (X - \mu_X)(Y - \mu_Y)$$

- $\text{Cov } X, Y = \text{Cov } Y, X$, so VC matrix is symmetric
- variance X is covariance of X with itself:
 $E(X - \mu_X)^2 = E(X - \mu_X)(X - \mu_X)$
- correlation between X and Y is

$$\text{Cor } X, Y = \frac{\text{Cov } X, Y}{\sqrt{\text{Var } X} \sqrt{\text{Var } Y}}$$

- $\text{Cov} = 0$ means $\text{Cor} = 0$
- In general, independent means $\text{Cov} = 0$.
- When observations have Multivariate normal distribution, $\text{Cov} = 0$ means independent

Multivariate normal distribution

- So can write VC matrix for "401" observations (independent, constant variance)

$$\Sigma = \begin{bmatrix} \sigma^2 & 0 & 0 & \cdots & 0 \\ 0 & \sigma^2 & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & \sigma^2 \end{bmatrix} = \sigma^2 \begin{bmatrix} 1 & 0 & 0 & \cdots & 0 \\ 0 & 1 & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & 1 \end{bmatrix} = \sigma^2 I$$

- I is the identity matrix.
 - matrix equivalent of the scalar 1

Least squares regression, again

- Model: $Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$
- Can write down (but won't) equations to estimate β_0 and β_1
- Do not generalize easily to more parameters, e.g.,
 $Z_i = \beta_0 + \beta_1 X_i + \beta_2 Y_i + \varepsilon_i$
- Can use matrices to simplify everything

$$\mathbf{X} = \begin{bmatrix} 1 & X_1 \\ 1 & X_2 \\ 1 & X_3 \\ \vdots & \vdots \\ 1 & X_n \end{bmatrix}, \beta = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix}, \mathbf{X}\beta = \begin{bmatrix} 1\beta_0 + X_1\beta_1 \\ 1\beta_0 + X_2\beta_1 \\ 1\beta_0 + X_3\beta_1 \\ \vdots \\ 1\beta_0 + X_n\beta_1 \end{bmatrix}$$

- Model: $\mathbf{Y} = \mathbf{X}\beta + \varepsilon$
 - Can have many columns in \mathbf{X} , or just 1

Least squares regression, again

- For any regression model with independent errors:

$$\hat{\beta}_{ols} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y}$$

- $()^{-1}$ is the matrix inverse, equivalent of scalar reciprocal
- $\mathbf{X}\mathbf{X}^{-1} = \mathbf{I}$ and $\mathbf{X}^{-1}\mathbf{X} = \mathbf{I}$
- labeled OLS for ordinary least squares
 - will see another type of LS soon (for correlated obs)
- and $\text{Var } \hat{\beta}_{ols} = \sigma^2 (\mathbf{X}'\mathbf{X})^{-1}$

Least squares regression, again

- Example: constant mean, $Y_i = \mu + \varepsilon_i$

$$\mathbf{X} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}, \beta = [\mu]$$

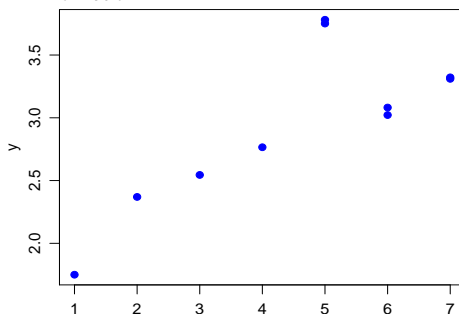
$$\mathbf{X}'\mathbf{X} = \mathbf{1}'\mathbf{1} = 1 + 1 + \dots = n, (\mathbf{X}'\mathbf{X})^{-1} = 1/n$$

$$\mathbf{X}'\mathbf{Y} = \mathbf{1}'\mathbf{Y} = Y_1 + Y_2 + \dots + Y_n = \Sigma Y$$

- so $\hat{\beta}_{ols} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y} = \Sigma Y / n$ and $\text{Var } \hat{\mu} = \sigma^2 / n$

Generalized Least Squares

- When errors are independent, each obs. is an additional piece of information
- Two positively correlated observations are less than 2 pieces of information



Generalized Least Squares

- Model:

$$\mathbf{Y} = \mathbf{X}\beta + \varepsilon, \text{Var } \varepsilon = \Sigma$$

- Estimates:

$$\hat{\beta}_{gls} = (\mathbf{X}'\Sigma^{-1}\mathbf{X})^{-1} \mathbf{X}'\Sigma^{-1}\mathbf{Y}$$

- What happens when $\Sigma = \sigma^2 \mathbf{I}$?

$$\begin{aligned} \hat{\beta}_{gls} &= (\mathbf{X}'\Sigma^{-1}\mathbf{X})^{-1} \mathbf{X}'\Sigma^{-1}\mathbf{Y} \\ &= (\mathbf{X}'(\sigma^2 \mathbf{I})^{-1}\mathbf{X})^{-1} \mathbf{X}'(\sigma^2 \mathbf{I})^{-1}\mathbf{Y} \\ &= (\mathbf{X}'(1/\sigma^2)\mathbf{X})^{-1} \mathbf{X}'(1/\sigma^2)\mathbf{Y} \\ &= \sigma^2(1/\sigma^2) (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y} = \hat{\beta}_{ols} \end{aligned}$$

- the GLS estimator simplifies to the usual OLS estimator
 - when observations are independent with constant variance

Generalized Least Squares

- Another way to think about GLS
- For most Σ , can find a "square-root" matrix, so that $\Sigma = C' C$
- Apply this idea to Σ^{-1} to get a square root matrix for the inverse, we will call this B , so $\Sigma^{-1} = B' B$
- Multiply all terms in the model by B

$$BY = BX\beta + B\epsilon$$

- Look what happens to the errors:

$$\text{Var } B\epsilon = B\Sigma B' = BC'(CB') = II = I$$

- Pre multiplying by square root of the inverse covariance matrix removes the correlation - transformed errors are now independent!

Generalized Least Squares

- So, if know the correlation matrix, can compute new $Y^* = BY$ and new $X^* = BX$
- and use OLS on X^* and Y^* : $Y^* = X^*\beta + \epsilon^*$
 - The regression coefficients from OLS on transformed variables are:
 - $\beta = [X^{*'} X^*]^{-1} (X^{*'} Y^*) = [(X' B')(BX)]^{-1} [(X' B')(BY)]$
 - $= (X \Sigma^{-1} X)^{-1} (X' \Sigma^{-1} Y)$
 - which are the GLS estimates for the original model

Generalized Least Squares

- Take home: If you know Σ , can compute the "best" estimates for any regression model
- using GLS
- A generalization (the Aitken model) says that all you need is the correlation part of Σ
 - Write Σ as $\sigma^2 C$ where C is (now) a correlation matrix
 - only need to know C , can estimate σ^2 .
- Practical problem is that the correlation part almost always has to be estimated
- Σ (or C) depends on unknown parameters.

SAR models

- General regression / ANOVA model
- $Y = X\beta + \epsilon$, if $X = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}$, then $Y = \mu + \epsilon$
- model spatial correlation by allowing ϵ to depend on error values in neighboring regions

$$\epsilon(s_i) = \sum_{j=1}^N b_{ij} \epsilon_j + \nu(s_i)$$

- b_{ij} are the elements of the spatial dependence matrix expressing dependence among regions
- $\nu(s_i)$ is an independent random disturbance for each region.
Usually assume $\nu(s_i) \stackrel{iid}{\sim} N(0, \sigma_\nu^2)$

- What this model “means”. Some examples:
- In all:
 - Focus on center observation (location s_5)
 - assume row standardized rook's neighbors, so $\{b_{ij}\}$ is

$$\begin{matrix} & 0 & 0.25 & 0 \\ 0.25 & \mathbf{0} & 0.25 \\ 0 & 0.25 & 0 \end{matrix}$$
- $\varepsilon(s_5) = 10$ Is this value large or small?
 - when b_{ij} is not zero, depends on neighbors
 - so look at $\nu(s_5)$ for $\varepsilon(s_5) = 10$
 - and different values for neighbors ε

- $$\begin{matrix} & 7 & \mathbf{8} & 10 \\ \mathbf{12} & 10 & \mathbf{8} \\ 11 & \mathbf{10} & 9 \end{matrix}$$
 - Only the bolded neighbors matter, $\sum b_{ij}\varepsilon(s_j) = 9.5$, $\nu(s_i) = 0.5$
 - large error but similar to neighbors, so independent contribution, $\nu(s_i)$ is small

- $$\begin{matrix} & 7 & \mathbf{16} & 10 \\ \mathbf{14} & 10 & \mathbf{12} \\ 11 & \mathbf{16} & 9 \end{matrix}$$
 - $\sum b_{ij}\varepsilon(s_j) = 14.5$, $\nu(s_i) = -4.5$
 - large error, but neighbors are larger, so independent contribution, $\nu(s_i)$ is negative

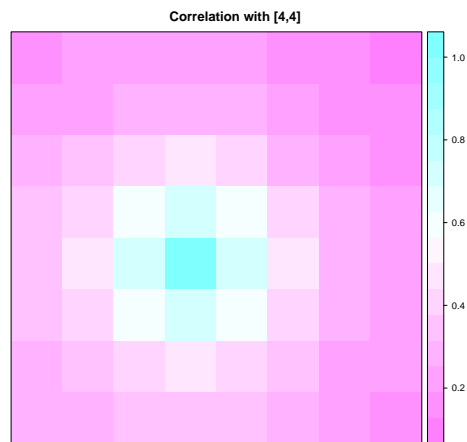
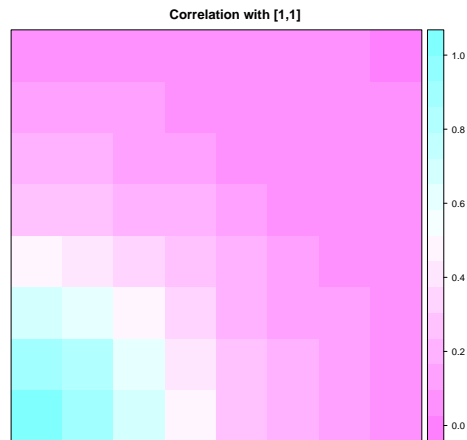
- $$\begin{matrix} & 1 & \mathbf{1} & 2 \\ \mathbf{3} & 10 & \mathbf{1} \\ 1 & \mathbf{2} & 3 \end{matrix}$$
 - $\sum b_{ij}\varepsilon(s_j) = 1.75$, $\nu(s_i) = 8.25$
 - much larger than neighbor errors, so independent contribution, $\nu(s_i)$ is large

- $$\begin{matrix} & -5 & -\mathbf{7} & -10 \\ -\mathbf{3} & 10 & -\mathbf{6} \\ -1 & \mathbf{2} & 0 \end{matrix}$$
 - $\sum b_{ij}\varepsilon(s_j) = -3.5$, $\nu(s_i) = 13.5$
 - neighbors suggest a negative error, so independent contribution, $\nu(s_i)$ is very large

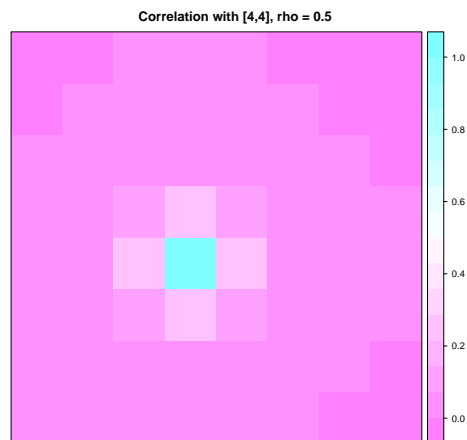
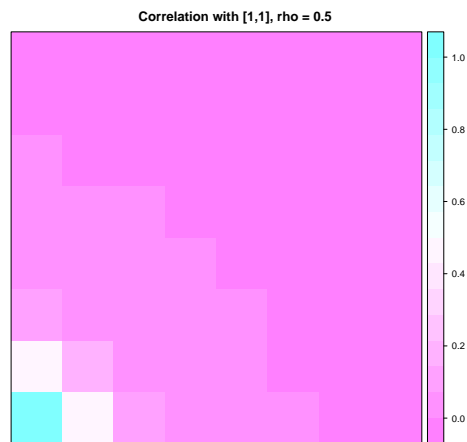
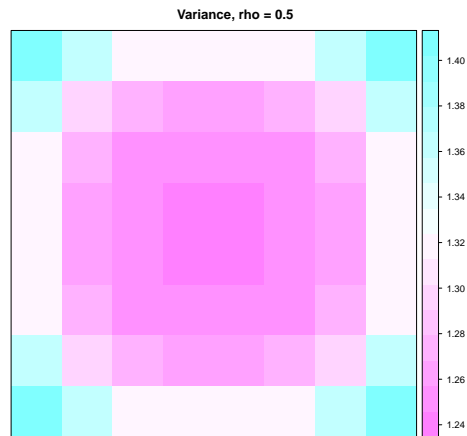
SAR models

- When part of the variation in error values can be “explained” by neighbors, the “explained” part is removed under the SAR model.
- Both SAR and CAR models do this, but they define “explained” differently
- SAR model: applies to all locations simultaneously
- For one location: $\varepsilon(s_i) = \sum_{j=1}^N b_{ij}\varepsilon_j + \nu(s_i)$
- For all locations: $\varepsilon = \mathbf{B}\varepsilon + \nu$, which means:
 - $\nu = \varepsilon - \mathbf{B}\varepsilon = (\mathbf{I} - \mathbf{B})\varepsilon$
 - $\varepsilon = (\mathbf{I} - \mathbf{B})^{-1}\nu$
- $\mathbf{Y} = \mathbf{X}\beta + (\mathbf{I} - \mathbf{B})^{-1}\nu$
 - where \mathbf{I} is the $N \times N$ identity matrix. \mathbf{B} is the matrix of “connectivity” coefficients
 - and ν is the vector of independent errors

- If consider $\mathbf{Y} = \mathbf{X}\beta + \varepsilon$, then the Var-Cov matrix of ε is $(\mathbf{I} - \mathbf{B})^{-1}\Sigma_{\nu}(\mathbf{I} - \mathbf{B}')^{-1}$
 - \mathbf{B} can be arbitrary, but it has $N(N-1)$ parameters. usually simplify using a model
 - usually write $\mathbf{B} = \rho\mathbf{W}$, where ρ is a coefficient of spatial dependence and \mathbf{W} is the spatial weight matrix
 - Note $\rho = 0 \Rightarrow$ locations are connected (elements of $\mathbf{W} > 0$), but not spatially dependent
- What does this model imply about VC matrix of the observations?
- Consider $\rho = 0.9$, \mathbf{W} is rook's neighbors, row standardized, $\Sigma_{\nu} = \sigma^2\mathbf{I}$
- Pictures on next few slides
 - Correlation declines with distance: good!
 - But Var \mathbf{Y} not constant - largest in corners, smallest in middle of region



- Notice modeling correlation between observations indirectly
 - B describes how one error value depends on the values of neighboring values
 - describing connections between observations, not correlation
 - Correlated errors is a consequence of those connections
 - As is unequal variances
- A geostatistical model directly describes correlations
- Preference is subject-matter dependent
 - Spatial econometrics, most areal data analysis, connections
 - Spatial statistics, geostatistical data, correlations
- Degree of non-constant variance depends on magnitude of ρ
- Similar plots for $\rho = 0.5$



- Reminder: model is $\mathbf{Y} = \mathbf{X}\beta + (\mathbf{I} - \mathbf{B})^{-1}\nu$,
- where $\mathbf{B} = \rho\mathbf{W}$
- \mathbf{W} is the known spatial weight matrix
- IF ρ is known, then \mathbf{B} is known, only need to estimate $\hat{\beta}$
- Easy:
 - 1) calculate $\Sigma_\varepsilon = (\mathbf{I} - \mathbf{B})^{-1}\Sigma_\nu(\mathbf{I} - \mathbf{B})^{-1}$ and use GLS, or
 - 2) Note that: $(\mathbf{I} - \mathbf{B})\mathbf{Y} = (\mathbf{I} - \mathbf{B})(\mathbf{X}\beta) + \nu$
transform \mathbf{Y} vector and \mathbf{X} matrix, and you have an OLS problem.
- Usual situation: ρ is unknown, need to estimate
 - use maximum likelihood
 - general alternative to LS for any statistical problem
 - Have already seen the VC matrix for ε : $\Sigma_\varepsilon = (\mathbf{I} - \mathbf{B})^{-1}\Sigma_\nu(\mathbf{I} - \mathbf{B})^{-1}$,
where $\mathbf{B} = \rho\mathbf{W}$
 - mvNormal lnL is

$$-\frac{k}{2} \log 2\pi - \frac{1}{2} \log |\Sigma_\varepsilon| - \frac{1}{2} (\mathbf{y} - \mathbf{X}\beta)' \Sigma_\varepsilon^{-1} (\mathbf{y} - \mathbf{X}\beta)$$

Maximizing the mvN lnL

- Iterative algorithm.
- Key insight is that given Σ_ε , mle of β is trivial (GLS)
- Assume $\rho = 0$, find OLS estimate of β
- Condition on β , use numerical maximization to find $\hat{\rho} | \beta$
- find GLS estimate of β for $\Sigma_\varepsilon(\rho)$
- repeat last two steps until convergence.
- Traditional frequentist approach is to estimate β and Var β conditional on $\hat{\rho}$
- Bayesian approach incorporates uncertainty in $\hat{\rho}$ into uncertainty about β
- Not an issue for simple problems (e.g. $\Sigma_\varepsilon = \sigma^2\mathbf{I}$) because in this case, $\hat{\beta}$ independent of σ^2
- Is in issue in these models because $\hat{\rho}$ and $\hat{\beta}$ are not independent.
- For sp diversity data, $\hat{\rho} = 0.914$.
That is strong positive spatial dependence.

Useful things about likelihood

- test hypotheses using Likelihood Ratio Test (LRT)
- e.g. test $H_0 : \rho = 0$
 - calculate lnL given $\rho = 0$ (need to maximize over β): $\ln L_0$
 - calculate lnL at mle's of all parameters: $\ln L_A$
 - $\ln L_A \geq \ln L_0$, because ρ probably not 0
 - but by how much? $\sim 0 \Rightarrow$ data consistent with $\rho = 0 \Rightarrow$ accept H_0
 - but how much is "too far" from 0?
 - General result: When H_0 true, $-2(\ln L_0 - \ln L_A) \sim \chi_k^2$ where k is the difference in # parameters between the two models
 - This is an asymptotic result, but surprisingly effective in small samples
 - Here, $k = 1$ because testing hypothesis about 1 parameter
 - $\ln L_A = -112.137$
 $\ln L_0 = -158.306$
 $\Delta = -2(\ln L_0 - \ln L_A) = 92.34$
 - $\chi_{1,0.95}^2 = 3.84$. Here $p < 0.0001$
 - LRT only works when one model is a simplification of another

Useful things about likelihood

- Model selection: e.g. compare different \mathbf{W} matrices
 - AIC = $-2 \ln L + 2k$
 - k is # of parameters. spdep counts β , σ^2 , and ρ
 - Can compare models with same number of parameters, or diff. # parameters
 - choose model with smallest AIC
 - Compare row-standardized (style='W') to binary (style='B') weights
 - Different models for connections between areas
 - \Rightarrow Different variance and covariance relationships
 - Here, 3 parameters (μ , σ^2 , and ρ)
 - Row std: AIC = $224.27 + 6 = 230.2742$
 - Binary: AIC = $235.54 + 6 = 241.54$
 - Choose model with Row std weights.

- Size of AIC irrelevant
 - doesn't matter whether best = 1000 or best = -250
 - Only comparisons among models fit to the same data
- Does not mean row std. is the "correct" model.
 - Only best among the set being considered.
 - model diagnostics still very important and useful
 - I don't know any that evaluate choice of weight matrix
- Smallest AIC is "best" model - Is anything almost as good?
 - Models with AIC within 2 units of the best are likely alternatives
 - Models with AIC > 10 units of the best are unlikely

Useful things about likelihood

- confidence interval for ρ by profile likelihood
 - Concept: Repeat LRT for many values of ρ
 - include inside 95% ci all values of ρ for which test has $p > 0.05$
 - here (0.82, 0.98)
- consequences of choice of ρ on inference about β
 - $\rho = 0$ (independence): $\hat{\mu} = 6.09$, se = 0.36
 - $\rho = 0.914$: $\hat{\mu} = 5.69$, se = 1.69
- Two points:
 - just demonstrated non-independence of $\hat{\rho}$ and $\hat{\beta}$ in a SAR model
 - se of $\hat{\mu}$ when account for correlation much larger
 - observations are positively correlated, so fewer "effective" observations

CAR models

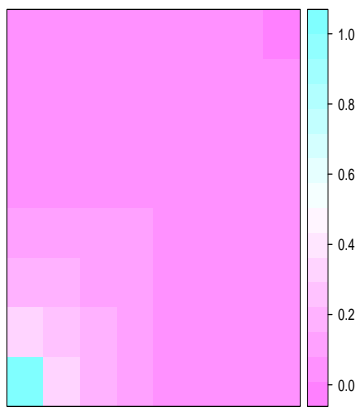
- SAR models all Y values simultaneously
- CAR models distrib of a each Y given the values of its neighbors

$$Y_i | Y_{-i} = \mathbf{X}_i \beta + \sum_{j=1}^N c_{ij} (Y_j - \mathbf{X}_j \beta) + \nu_i$$

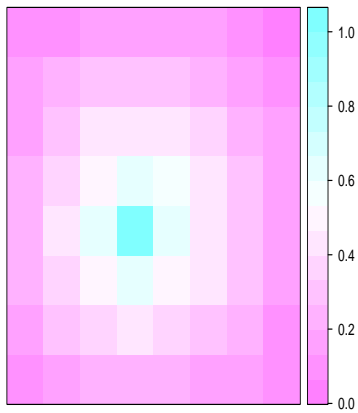
- RHS same as SAR model
- The difference is that Y_{-i} now treated as fixed values when specifying distrib of Y_i
- Diff. VC matrix for obs.: $\Sigma_\varepsilon = \text{Var } Y = (\mathbf{I} - \mathbf{C})^{-1} \Sigma_\nu$

- to be a valid VC matrix, requires some conditions:
 - ρ can not be too large
 - Definition depends on the weight matrix, \mathbf{W}
 - \mathbf{C} must be symmetric, $C_{ij} = C_{ji}$
 - so my example now uses binary weights
- Practical difference: correlation pattern similar
 - Pictures on next slide
 - correl between pairs in middle higher among than between middle and edge
 - similar to SAR pattern, but details slightly different
- But pattern of variance quite different
 - biggest var in the middle of the area (more connections to other points)
- Fitting a CAR model gives: $\hat{\mu} = 5.10$, se = 0.83, $\hat{\rho} = 0.253$.
- 0.253 doesn't seem large, but close to maximum possible value

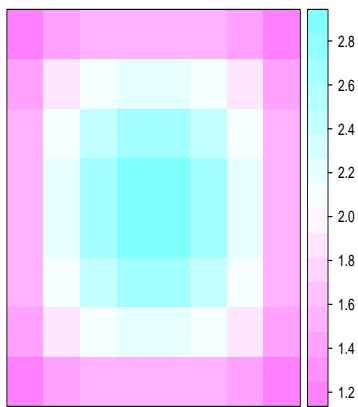
Correlation with [1,1], CAR



Correlation with [4,4], CAR



Variance, CAR

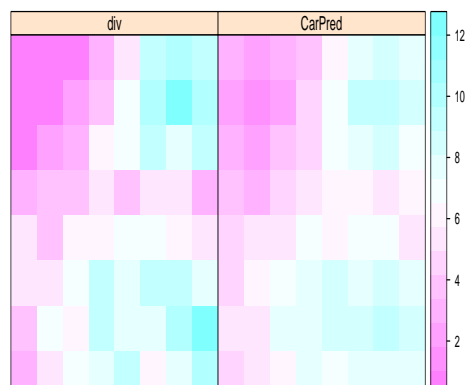


- Different variance and correlation pattern \Rightarrow CAR and SAR are not the same models
- Which fits the data better? How will you answer this Q?

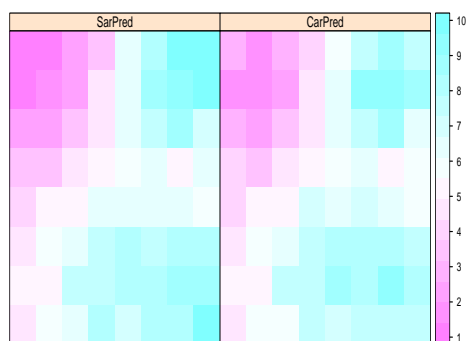
Model comparison

- I would use AIC
 - model AIC
 - SAR 241.54
 - CAR 262.38
 - Indep 320.61
 - Clear dominance of SAR model
- Traditional interpretation
 - AIC w/i 2 of the top: worse model is reasonable competitor
 - AIC more than 10 from the top: worse model is very unlikely
- For Bayesian analysis, CAR models very, very popular
 - They are easy to use in an MCMC chain, because they are conditional distributions

Smoothing noisy areal data



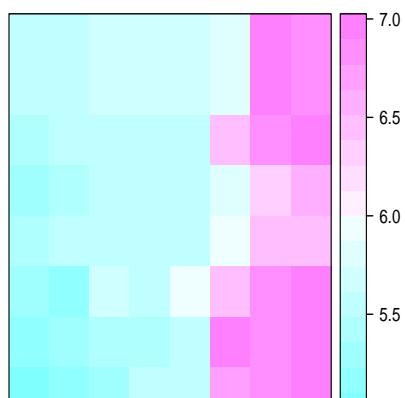
Comparison of smoothing using CAR and SAR



More complicated models

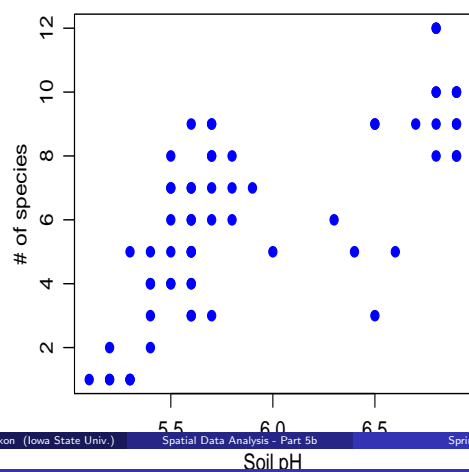
- What if you have X variables measured for each area?
- Easy to include in the model
- My example: soil pH measured in each area
- the small quadrats cross from igneous parent material (lower pH) to limestone parent material (higher pH)
- Residuals from OLS regression are spatially correlated
 - Moran's I: 0.62
 - less than 1 for observations = 0.79
- Pictures on next three slides

Soil pH



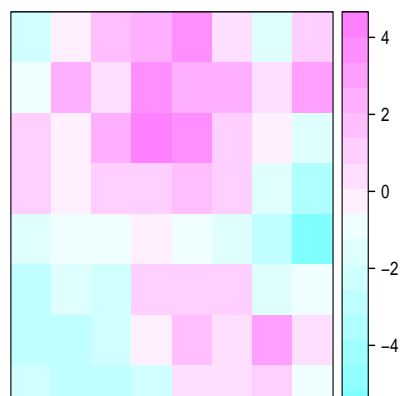
© Philip M. Dixon (Iowa State Univ.) Spatial Data Analysis - Part 5b Spring 2020 45 / 53

species vs. Soil pH



© Philip M. Dixon (Iowa State Univ.) Spatial Data Analysis - Part 5b Spring 2020 46 / 53

Residuals from linear regression

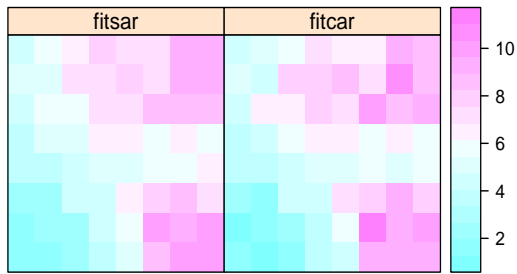


© Philip M. Dixon (Iowa State Univ.) Spatial Data Analysis - Part 5b Spring 2020 47 / 53

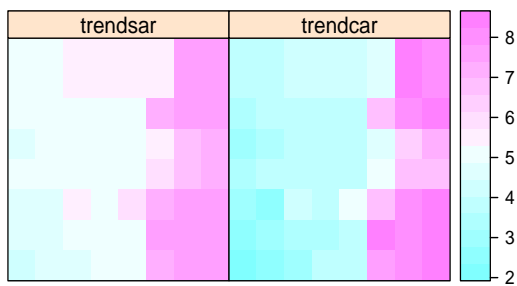
Results from pH models

Model	$\hat{\beta}_{pH}$	se	$\hat{\rho}$	AIC
Indep.	3.56	0.43	—	275.06
SAR	1.93	0.66	0.86	225.77
CAR	3.28	0.52	0.26	233.78

- Pictures of the two fits on next slide



SAR and CAR fits: trend component = X beta



Why is SAR slope for pH much lower?

- When X variable is spatially correlated
 - X and the spatial correlation "fight" to predict Y
 - analogous to two correlated X variables: pH and "space"
- Commonly seen when using spatial correlated errors
 - "The effect you love disappears"
- For these data:
 - SAR fit: smaller β_{pH} , higher correlation
 - CAR fit: larger β_{pH} , smaller correlation

Comparison of intercept only to pH SAR models

