# Spatial Modeling of Zero-Inflated Data with Copula Models
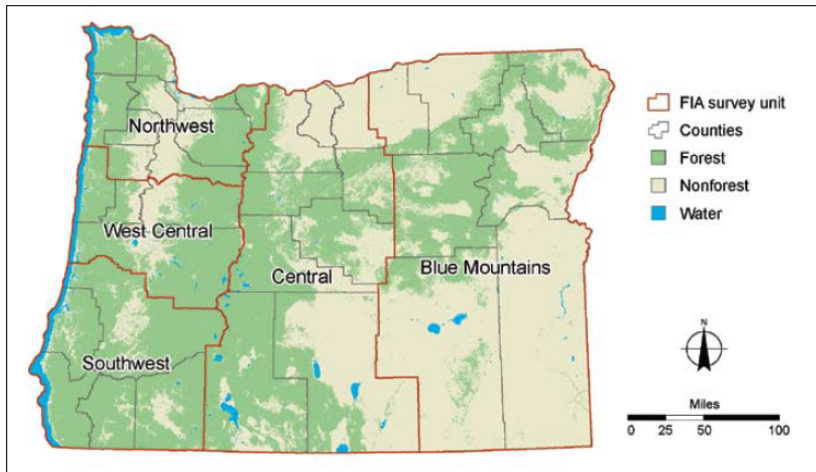
Lisa Madsen [1]    Vicente Monleon [2]

[1] Oregon State University
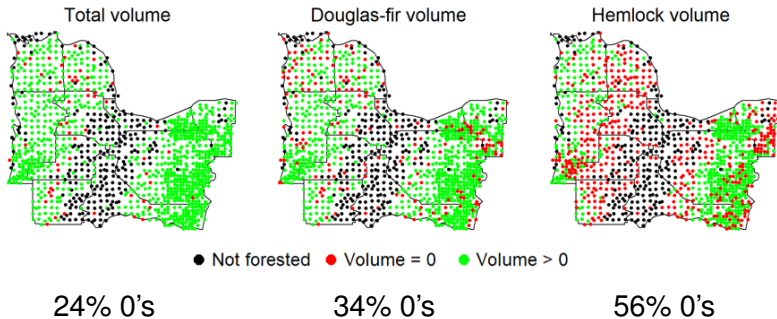
[2] USDA Forest Service, PNW Research Station

WNAR 2021

# Forest Inventory and Analysis (FIA)

# Variables of Interest



● Not forested  ● Volume = 0  ● Volume > 0

24% 0's         34% 0's         56% 0's

## Challenges

- Spatial dependence
- Non-normality
- Zero-inflation

## Challenges

- Spatial dependence
- Non-normality
- Zero-inflation
- Big data

## Outline

Motivation

Model Components
    Gaussian Copula
    Marginal Distributions
    Spatial Dependence Model

Zero-inflated Spatial Model

Model Fitting

Spatial Prediction

Conclusions

Motivation   Model Components   Zero-inflated Spatial Model   Model Fitting   Spatial Prediction   Conclusions
0000         ●0000              00                            00000000000   0000000000000       000
             000000
             00000000000

Gaussian Copula

## Outline

Motivation

### Model Components
#### Gaussian Copula
Marginal Distributions
Spatial Dependence Model

Zero-inflated Spatial Model

Model Fitting

Spatial Prediction

Conclusions

Motivation 0000    Model Components ○●○○○ ○○○○○○ ○○○○○○○○○○○    Zero-inflated Spatial Model ○○    Model Fitting ○○○○○○○○○○○    Spatial Prediction ○○○○○○○○○○○○○    Conclusions ○○○

Gaussian Copula

## Copula Models

Definition: A **copula** is a multivariate distribution with uniform marginals.

Motivation   Model Components   Zero-inflated Spatial Model   Model Fitting   Spatial Prediction   Conclusions
0000         00000                00                            00000000000    0000000000000        000
             000000
             00000000000

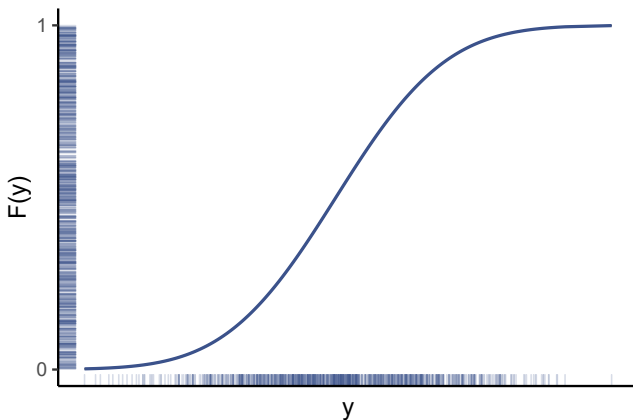Gaussian Copula

# Copula Models

Definition: A **copula** is a multivariate distribution with uniform marginals.

**Copula models** are useful for modeling multivariate data with arbitrary marginal distributions.

| Motivation | Model Components | Zero-inflated Spatial Model | Model Fitting | Spatial Prediction | Conclusions |
| 0000 | 0●000 | 00 | 0000000000 | 000000000000 | 000 |
| | 000000 | | | | |
| | 00000000000 | | | | |

Gaussian Copula

## Copula Models

Definition: A **copula** is a multivariate distribution with uniform marginals.

**Copula models** are useful for modeling multivariate data with arbitrary marginal distributions.

The **Gaussian copula** allows us to adapt methodology based on the multivariate normal distribution to non-normal data.
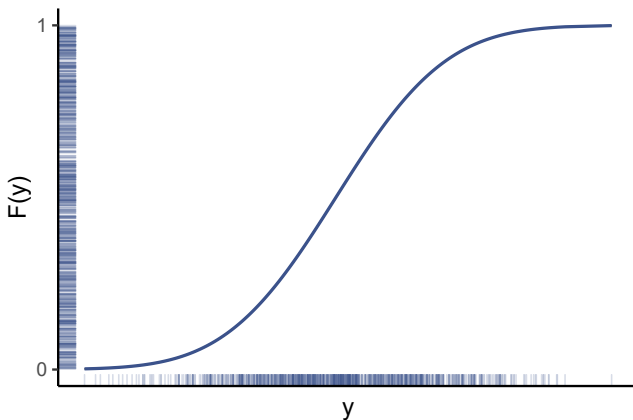
Motivation    Model Components    Zero-inflated Spatial Model    Model Fitting    Spatial Prediction    Conclusions
0000          00●00                00                            00000000000      0000000000000        000
              000000
              00000000000

Gaussian Copula

## Copula Construction

If $Y$ has continuous cdf $F(y)$, then $F(Y) \sim U(0,1)$.

Motivation    Model Components    Zero-inflated Spatial Model    Model Fitting    Spatial Prediction    Conclusions
0000          00●00               00                            00000000000      00000000000000        000
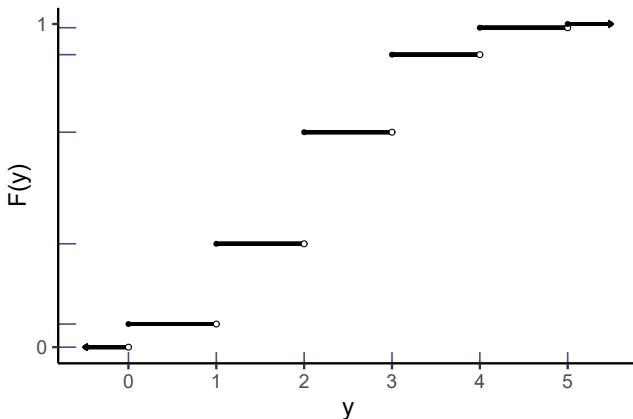              000000
              00000000000

Gaussian Copula

## Copula Construction

If $Y$ has continuous cdf $F(y)$, then $F(Y) \sim U(0, 1)$.



Conversely, if $U \sim U(0, 1)$, then $F^{-1}(U) \sim F$.

Motivation  Model Components  Zero-inflated Spatial Model  Model Fitting  Spatial Prediction  Conclusions
0000        00000                00                          00000000000    0000000000000       000
            000000
            00000000000

Gaussian Copula

## Discrete Marginals

If $Y$ is discrete, then $F^{-1}(U) \sim F$, but $F(Y)$ is not uniform.

Motivation    Model Components    Zero-inflated Spatial Model    Model Fitting    Spatial Prediction    Conclusions
0000          0000●               00                             00000000000      00000000000000       000
              000000
              0000000000000

Gaussian Copula

## Gaussian Copula

Notation:

$\Phi_{\Sigma}$ denotes the *n*-dimensional standard normal cdf with correlation matrix $\Sigma$.

Gaussian Copula

## Gaussian Copula

Notation:

$\Phi_\Sigma$ denotes the *n*-dimensional standard normal cdf with correlation matrix $\Sigma$.

$\Phi$ denotes the univariate standard normal cdf.

Motivation   Model Components   Zero-inflated Spatial Model   Model Fitting   Spatial Prediction   Conclusions
0000         0000●                00                           00000000000     00000000000000     000
             000000
             00000000000

Gaussian Copula

## Gaussian Copula

Notation:

$\Phi_{\Sigma}$ denotes the *n*-dimensional standard normal cdf with correlation matrix $\Sigma$.

$\Phi$ denotes the univariate standard normal cdf.

$F_1, \ldots, F_n$ are (preferably continuous) marginal cdfs.

## Gaussian Copula

Notation:

$\Phi_{\Sigma}$ denotes the *n*-dimensional standard normal cdf with correlation matrix $\Sigma$.

$\Phi$ denotes the univariate standard normal cdf.

$F_1, \ldots, F_n$ are (preferably continuous) marginal cdfs.

Copula CDF:

$$C(\boldsymbol{y}; \Sigma) = \Phi_{\Sigma}[\Phi^{-1}\{F_1(y_1)\}, \ldots, \Phi^{-1}\{F_n(y_n)\}]$$

Motivation  Model Components  Zero-inflated Spatial Model  Model Fitting  Spatial Prediction  Conclusions
0000        00000             00                            00000000000    0000000000000     000
            000000
            0000000000000

Gaussian Copula

## Gaussian Copula

Notation:

$\Phi_{\Sigma}$ denotes the $n$-dimensional standard normal cdf with correlation matrix $\Sigma$.

$\Phi$ denotes the univariate standard normal cdf.

$F_1, \ldots, F_n$ are (preferably continuous) marginal cdfs.

Copula CDF:

$$C(\boldsymbol{y}; \Sigma) = \Phi_{\Sigma}[\Phi^{-1}\{F_1(y_1)\}, \ldots, \Phi^{-1}\{F_n(y_n)\}]$$

Gaussian copula's normalizing transformation for continuous $F_i$:

$Y_i \sim F_i$

Motivation 0000 | Model Components 00000 000000 00000000000 | Zero-inflated Spatial Model 00 | Model Fitting 00000000000 | Spatial Prediction 0000000000000 | Conclusions 000

Gaussian Copula

## Gaussian Copula

Notation:

$\Phi_{\Sigma}$ denotes the *n*-dimensional standard normal cdf with correlation matrix $\Sigma$.

$\Phi$ denotes the univariate standard normal cdf.

$F_1, \ldots, F_n$ are (preferably continuous) marginal cdfs.

Copula CDF:

$$C(\boldsymbol{y}; \Sigma) = \Phi_{\Sigma}[\Phi^{-1}\{F_1(y_1)\}, \ldots, \Phi^{-1}\{F_n(y_n)\}]$$

Gaussian copula's normalizing transformation for continuous $F_i$:

$Y_i \sim F_i \implies F_i(Y_i) \sim \text{Uniform}(0, 1)$

Motivation 0000 | Model Components 0000● 000000 00000000000 | Zero-inflated Spatial Model 00 | Model Fitting 00000000000 | Spatial Prediction 0000000000000 | Conclusions 000

Gaussian Copula

## Gaussian Copula

Notation:

$\Phi_\Sigma$ denotes the $n$-dimensional standard normal cdf with correlation matrix $\Sigma$.

$\Phi$ denotes the univariate standard normal cdf.

$F_1, \ldots, F_n$ are (preferably continuous) marginal cdfs.

Copula CDF:

$$C(\boldsymbol{y}; \Sigma) = \Phi_\Sigma[\Phi^{-1}\{F_1(y_1)\}, \ldots, \Phi^{-1}\{F_n(y_n)\}]$$

Gaussian copula's normalizing transformation for continuous $F_i$:

$$Y_i \sim F_i \;\Rightarrow\; F_i(Y_i) \sim \text{Uniform}(0, 1) \;\Rightarrow\; \Phi^{-1}\{F_i(Y_i)\} \sim N(0, 1)$$

# Outline

Motivation

Model Components
Gaussian Copula
Marginal Distributions
Spatial Dependence Model

Zero-inflated Spatial Model

Model Fitting

Spatial Prediction

Conclusions

Motivation    **Model Components**    Zero-inflated Spatial Model    Model Fitting    Spatial Prediction    Conclusions
0000          00000                   00                            00000000000     0000000000000       000
              ○●0000
              00000000000

Marginal Distributions

Motivation    **Model Components**    Zero-inflated Spatial Model    Model Fitting    Spatial Prediction    Conclusions
0000          00000                   00                            0000000000       0000000000000        000
              ○●0000
              00000000000

Marginal Distributions

24% of plots have zero total volume.

Marginal Distributions

# Zero-inflated Lognormal Distribution

$$B \sim \text{Bernoulli}(1 - p)$$

Motivation   **Model Components**   Zero-inflated Spatial Model   Model Fitting   Spatial Prediction   Conclusions
0000         00000                  00                           00000000000    0000000000000      000
             000●000
             00000000000

Marginal Distributions

# Zero-inflated Lognormal Distribution

$$
\begin{aligned}
B &\sim \text{Bernoulli}(1 - p) \\
W &\sim \text{Lognormal}(\mu, \sigma^2)
\end{aligned}
$$

## Zero-inflated Lognormal Distribution

$$
\begin{aligned}
B &\sim \text{Bernoulli}(1 - p) \\
W &\sim \text{Lognormal}(\mu, \sigma^2) \\
Y &= \begin{cases} 0 & B = 1 \\ W & B = 0 \end{cases}
\end{aligned}
$$

Motivation  Model Components  Zero-inflated Spatial Model  Model Fitting  Spatial Prediction  Conclusions
oooo  ooooo  oo  ooooooooooo  oooooooooooo  ooo
       oo●ooo
       ooooooooooo

Marginal Distributions

## Zero-inflated Lognormal Distribution

$$
\begin{aligned}
B &\sim \text{Bernoulli}(1-p) \\
W &\sim \text{Lognormal}(\mu, \sigma^2) \\
Y &= \left\{ \begin{array}{ll} 0 & B = 1 \\ W & B = 0 \end{array} \right. \\
P(Y = 0) &= p
\end{aligned}
$$

## Zero-inflated Lognormal Distribution

$$
\begin{aligned}
B &\sim \text{Bernoulli}(1 - p) \\
W &\sim \text{Lognormal}(\mu, \sigma^2) \\
Y &= \left\{ \begin{array}{ll} 0 & B = 1 \\ W & B = 0 \end{array} \right.
\end{aligned}
$$

$$
\begin{aligned}
P(Y = 0) &= p \\
\log(W) &\sim N(\mu, \sigma^2)
\end{aligned}
$$

# Zero-inflated Lognormal CDF

Motivation    Model Components    Zero-inflated Spatial Model    Model Fitting    Spatial Prediction    Conclusions
0000          00000               00                             00000000000      0000000000000         000
              000000
              00000000000

Marginal Distributions

# Zero-inflated Lognormal CDF



CDF discontinuous at 0.

Motivation  Model Components  Zero-inflated Spatial Model  Model Fitting  Spatial Prediction  Conclusions
0000        00000                00                          00000000000    00000000000000      000
            0000●0
            0000000000000

Marginal Distributions

## Continuous Zero-inflated Lognormal Distribution

$$B \sim \text{Bernoulli}(1 - p)$$

Motivation  **Model Components**  Zero-inflated Spatial Model  Model Fitting  Spatial Prediction  Conclusions
0000          00000                  00                          00000000000     0000000000000     000
              000000
              00000000000

Marginal Distributions

## Continuous Zero-inflated Lognormal Distribution

$$
\begin{aligned}
B &\sim \text{Bernoulli}(1 - p) \\
W - \epsilon &\sim \text{Lognormal}(\mu, \sigma^2)
\end{aligned}
$$

## Continuous Zero-inflated Lognormal Distribution

$$
\begin{aligned}
B &\sim \text{Bernoulli}(1 - p) \\
W - \epsilon &\sim \text{Lognormal}(\mu, \sigma^2) \\
Y &= \begin{cases} \text{Uniform}(0, \epsilon) & B = 1 \\ W & B = 0 \end{cases}
\end{aligned}
$$

Motivation 0000 | Model Components 00000 0000•0 00000000000 | Zero-inflated Spatial Model 00 | Model Fitting 00000000000 | Spatial Prediction 0000000000000 | Conclusions 000

Marginal Distributions

## Continuous Zero-inflated Lognormal Distribution

$$
\begin{aligned}
B &\sim \text{Bernoulli}(1 - p) \\
W - \epsilon &\sim \text{Lognormal}(\mu, \sigma^2) \\
Y &= \begin{cases} \text{Uniform}(0, \epsilon) & B = 1 \\ W & B = 0 \end{cases}
\end{aligned}
$$

$$
F(y) = \begin{cases}
0, & y < 0 \\
y \cdot p/\epsilon, & 0 \leq y < \epsilon \\
p + (1 - p)F_{\text{lnorm}}(y - \epsilon; \mu, \sigma^2), & y \geq \epsilon
\end{cases}
$$

| Motivation | Model Components | Zero-inflated Spatial Model | Model Fitting | Spatial Prediction | Conclusions |
|------------|------------------|-----------------------------|---------------|-------------------|-------------|
| oooo | ooooo | oo | ooooooooooo | ooooooooooooo | ooo |
|  | ooooo● |  |  |  |  |
|  | ooooooooooo |  |  |  |  |

Marginal Distributions

# Outline

Motivation

Model Components
  Gaussian Copula
  Marginal Distributions
  Spatial Dependence Model

Zero-inflated Spatial Model

Model Fitting

Spatial Prediction

Conclusions

The Gaussian copula induces dependence via the copula association matrix $\Sigma$.

$$C(\boldsymbol{y}; \Sigma) = \Phi_{\Sigma}[\Phi^{-1}\{F_1(y_1)\}, \ldots, \Phi^{-1}\{F_n(y_n)\}]$$

Motivation 0000 | Model Components 00000 000000 0●000000000 | Zero-inflated Spatial Model 00 | Model Fitting 00000000000 | Spatial Prediction 0000000000000 | Conclusions 000

Spatial Dependence Model

The Gaussian copula induces dependence via the copula association matrix $\Sigma$.

$$C(\mathbf{y}; \Sigma) = \Phi_\Sigma[\Phi^{-1}\{F_1(y_1)\}, \ldots, \Phi^{-1}\{F_n(y_n)\}]$$

$$\Sigma_{ij} = \text{corr}\left(\Phi^{-1}\{F_i(Y_i)\}, \Phi^{-1}\{F_j(Y_j)\}\right)$$

The Gaussian copula induces dependence via the copula association matrix $\Sigma$.

$$C(\boldsymbol{y}; \Sigma) = \Phi_{\Sigma}[\Phi^{-1}\{F_1(y_1)\}, \ldots, \Phi^{-1}\{F_n(y_n)\}]$$

$$\Sigma_{ij} = \text{corr}\left(\Phi^{-1}\{F_i(Y_i)\}, \Phi^{-1}\{F_j(Y_j)\}\right) \neq \text{corr}(Y_i, Y_j)$$

The Gaussian copula induces dependence via the copula association matrix $\Sigma$.

$$C(\boldsymbol{y}; \Sigma) = \Phi_\Sigma[\Phi^{-1}\{F_1(y_1)\}, \ldots, \Phi^{-1}\{F_n(y_n)\}]$$

$$\Sigma_{ij} = \text{corr}\left(\Phi^{-1}\{F_i(Y_i)\}, \Phi^{-1}\{F_j(Y_j)\}\right) \neq \text{corr}(Y_i, Y_j)$$

If $F_i$ and $F_j$ are continuous, $\Sigma_{ij}$ is the **rank** correlation between $Y_i$ and $Y_j$.

The Gaussian copula induces dependence via the copula association matrix $\Sigma$.

$$C(\mathbf{y}; \Sigma) = \Phi_{\Sigma}[\Phi^{-1}\{F_1(y_1)\}, \ldots, \Phi^{-1}\{F_n(y_n)\}]$$

$$\Sigma_{ij} = \text{corr}\left(\Phi^{-1}\{F_i(Y_i)\}, \Phi^{-1}\{F_j(Y_j)\}\right) \neq \text{corr}(Y_i, Y_j)$$

If $F_i$ and $F_j$ are continuous, $\Sigma_{ij}$ is the **rank** correlation between $Y_i$ and $Y_j$.

Model the elements of $\Sigma$ as a decreasing function of distance.

Motivation  Model Components  Zero-inflated Spatial Model  Model Fitting  Spatial Prediction  Conclusions
0000  00000  00  0000000000  00000000000000  000
      00000
      000●000000000

Spatial Dependence Model

## Variogram Models

Definition: When $h_{ij}$ denotes the distance between locations of observations $i$ and $j$, the **isotropic variogram** is

$$2\gamma(h_{ij}) = \mathrm{var}(Y_i - Y_j)$$

Motivation   Model Components   Zero-inflated Spatial Model   Model Fitting   Spatial Prediction   Conclusions
0000         00000              00                           00000000000     0000000000000     000
             000000
             0000000000000

Spatial Dependence Model

## Variogram Models

Definition: When $h_{ij}$ denotes the distance between locations of observations $i$ and $j$, the **isotropic variogram** is

$$2\gamma(h_{ij}) = \text{var}(Y_i - Y_j)$$

Exponential Semivariogram:

$$\gamma(h) = \begin{cases} 0, & h = 0 \\ \theta_n + \theta_s[1 - \exp(-h/\theta_r)], & h > 0 \end{cases}$$

Motivation  Model Components  Zero-inflated Spatial Model  Model Fitting  Spatial Prediction  Conclusions
0000  00000  00  00000000000  0000000000000  000
      000000
      0000000000

Spatial Dependence Model

## Variogram Models

Definition: When $h_{ij}$ denotes the distance between locations of observations $i$ and $j$, the **isotropic variogram** is

$$2\gamma(h_{ij}) = \text{var}(Y_i - Y_j)$$

Exponential Semivariogram:

$$\gamma(h) = \begin{cases} 0, & h = 0 \\ \theta_n + \theta_s[1 - \exp(-h/\theta_r)], & h > 0 \end{cases}$$

$\theta_n$ is the **nugget**

Motivation   Model Components   Zero-inflated Spatial Model   Model Fitting   Spatial Prediction   Conclusions
0000          00000                00                            0000000000      0000000000000        000
              000000
              0000●000000

Spatial Dependence Model

## Variogram Models

Definition: When $h_{ij}$ denotes the distance between locations of observations $i$ and $j$, the **isotropic variogram** is

$$2\gamma(h_{ij}) = \mathrm{var}(Y_i - Y_j)$$

Exponential Semivariogram:

$$\gamma(h) = \begin{cases} 0, & h = 0 \\ \theta_n + \theta_s[1 - \exp(-h/\theta_r)], & h > 0 \end{cases}$$

$\theta_n$ is the **nugget**

$\theta_s$ is the **partial sill**

Motivation  Model Components  Zero-inflated Spatial Model  Model Fitting  Spatial Prediction  Conclusions
0000  00000  00  00000000000  0000000000000  000
       000000
       00000●00000

Spatial Dependence Model

## Variogram Models

Definition: When $h_{ij}$ denotes the distance between locations of observations $i$ and $j$, the **isotropic variogram** is

$$2\gamma(h_{ij}) = \text{var}(Y_i - Y_j)$$

Exponential Semivariogram:

$$\gamma(h) = \begin{cases} 0, & h = 0 \\ \theta_n + \theta_s[1 - \exp(-h/\theta_r)], & h > 0 \end{cases}$$

$\theta_n$ is the **nugget**

$\theta_s$ is the **partial sill**

$\theta_r$ is the **range**

# Exponential Semivariogram

## Covariance

If $\mathrm{var}(Y_i) = \mathrm{var}(Y_j)$, then $2\gamma(h_{ij}) = 2\,\mathrm{var}(Y_i) - 2\,\mathrm{cov}(Y_i, Y_j)$

## Covariance

If $\text{var}(Y_i) = \text{var}(Y_j)$, then $2\gamma(h_{ij}) = 2\,\text{var}(Y_i) - 2\,\text{cov}(Y_i, Y_j)$

Exponential covariance function:

$$C(h) = \left\{ \begin{array}{ll} \theta_n + \theta_s, & h = 0 \\ \theta_s \exp(-h/\theta_r), & h > 0 \end{array} \right.$$

Motivation   Model Components   Zero-inflated Spatial Model   Model Fitting   Spatial Prediction   Conclusions
0000          00000               00                            00000000000     0000000000000        000
              000000
              00000000●00

Spatial Dependence Model

# Exponential Covariance Function

## Correlation

Exponential correlation function:

$$\rho(h) = \begin{cases} 1, & h = 0 \\ \frac{\theta_s}{\theta_s + \theta_n} \exp(-h/\theta_r), & h > 0 \end{cases}$$

# Exponential Correlation Function

## Spatial Gaussian Copula Model

Let $\boldsymbol{Y} = [Y_1, \ldots Y_n]$ have cdf

$$F(\boldsymbol{y}; \Sigma) = \Phi_\Sigma[\Phi^{-1}\{F_1(y_1)\}, \ldots, \Phi^{-1}\{F_n(y_n)\}].$$

# Spatial Gaussian Copula Model

Let $\boldsymbol{Y} = [Y_1, \ldots Y_n]$ have cdf

$$F(\boldsymbol{y}; \Sigma) = \Phi_{\Sigma}[\Phi^{-1}\{F_1(y_1)\}, \ldots, \Phi^{-1}\{F_n(y_n)\}].$$

Association matrix $\Sigma$ has $ij$th element

$$\Sigma_{ij} = \rho(h_{ij}) = \left\{ \begin{array}{ll} 1, & h_{ij} = 0 \\ \alpha_N \exp(-h_{ij}/\alpha_R), & h_{ij} > 0 \end{array} \right. ,$$

where $h_{ij}$ denotes the Euclidean distance between locations of observations $i$ and $j$.

## Zero Inflated Continuous Marginals

$Y_i \sim F_i$ with

$$
F_i(y) = \begin{cases} 0, & y < 0 \\ y \cdot p_i/\epsilon, & 0 \le y < \epsilon \\ p_i + (1-p_i)F_{\mathsf{lnorm}}(y - \epsilon; \mu_i, \sigma^2), & y \ge \epsilon \end{cases}
$$

where $p_i$ and $\mu_i$ may depend on covariates.

# Survey Unit 1 Volume Map



Total Volume (m$^3$ha$^{-1}$)

$n = 1224$ plots
298 with zero total volume

## Logistic Submodel

Logistic regression model for the Bernoulli process:

$$
\begin{aligned}
B_i &= \begin{cases} 1, & \text{plot } i \text{ total volume} = 0 \\ 0, & \text{otherwise} \end{cases} \\
B_i &\sim \text{Bernoulli}(p_i) \\
\text{logit}(p_i) &= \boldsymbol{X}_i \boldsymbol{\beta} \\
\boldsymbol{X}_i &= \text{row vector of covariates}
\end{aligned}
$$

## Lognormal Submodel

Log-linear regression model for the 926 non-0 volume
observations:

$$
\begin{aligned}
Y_i &= \text{total volume in } i\text{th plot, if positive} \\
\log(\boldsymbol{Y}) &\sim N(\boldsymbol{\mu}, \sigma^2) \\
\boldsymbol{\mu} &= \boldsymbol{X}_y \gamma
\end{aligned}
$$

where $\boldsymbol{X}_y$ is a design matrix of covariates.

Potential covariates:

| | |
|---|---|
| forind | indicator of forest |
| annpre | mean annual precipitation |
| anntmp | mean annual temperature |
| smrtp | moisture stress during growing season |
| ndvi | vegetation greenness |
| tc1 | brightness |
| tc2 | greenness |
| tc3 | wetness |

## Logistic Model Fit

From R's `glm` function.

```
Coefficients:
            Estimate Std. Error z value Pr(>|z|)
(Intercept) -2.7570     0.2355  -11.709  < 2e-16 ***
forind      -0.5957     0.1425   -4.181 2.90e-05 ***
tc1          1.2641     0.2315    5.460 4.77e-08 ***
tc2         -0.6842     0.2146   -3.188 0.001435 **
tc3         -0.9236     0.1861   -4.963 6.95e-07 ***
annpre      -1.2260     0.3363   -3.646 0.000267 ***
anntmp       2.1818     0.3405    6.407 1.49e-10 ***
smrtp       -1.3573     0.4446   -3.053 0.002268 **
ndvi        -0.7923     0.2051   -3.862 0.000112 ***
forind:tc2  -0.2492     0.1154   -2.160 0.030756 *
forind:tc3  -0.4386     0.1351   -3.247 0.001167 **
```

## Lognormal Model Fit

From R's `lm` function.

```
Coefficients:
             Estimate Std. Error t value Pr(>|t|)
(Intercept)  5.14220    0.05453  94.308  < 2e-16 ***
tc1         -0.57893    0.10593  -5.465 5.96e-08 ***
tc2         -0.13481    0.07386  -1.825  0.06828 .
tc3          0.82963    0.08732   9.501  < 2e-16 ***
annpre       0.17192    0.04515   3.808  0.00015 ***
anntmp       0.11513    0.04663   2.469  0.01373 *
tc1:ndvi     0.08206    0.06750   1.216  0.22437
tc2:ndvi     0.09974    0.05807   1.718  0.08617 .
anntmp:ndvi -0.31491    0.05041  -6.247 6.38e-10 ***
tc3:smrtp    0.27859    0.08820   3.159  0.00164 **
tc3:annpre   0.02789    0.08764   0.318  0.75040
```

## Fitting the Spatial Model

Copula association matrix $\Sigma$ is the spatial correlation matrix of $\Phi^{-1}\{F_1(Y_1)\}, \ldots, \Phi^{-1}\{F_n(Y_n)\}$.

## Fitting the Spatial Model

Copula association matrix $\Sigma$ is the spatial correlation matrix of
$\Phi^{-1}\{F_1(Y_1)\}, \ldots, \Phi^{-1}\{F_n(Y_n)\}$.

Estimate variogram of $\Phi^{-1}\{\widehat{F}_1(Y_1)\}, \ldots, \Phi^{-1}\{\widehat{F}_n(Y_n)\}$, where

$$
\widehat{F}_i(y) = \begin{cases} 0, & y < 0 \\ y \cdot \widehat{p}_i/\epsilon, & 0 \le y < \epsilon \\ \widehat{p}_i + (1 - \widehat{p}_i)F_{\text{lnorm}}(y - \epsilon; \widehat{\mu}_i, \widehat{\sigma}^2), & y \ge \epsilon \end{cases}
$$

## Fitting the Spatial Model

Copula association matrix $\Sigma$ is the spatial correlation matrix of $\Phi^{-1}\{F_1(Y_1)\}, \ldots, \Phi^{-1}\{F_n(Y_n)\}$.

Estimate variogram of $\Phi^{-1}\{\widehat{F}_1(Y_1)\}, \ldots, \Phi^{-1}\{\widehat{F}_n(Y_n)\}$, where

$$
\widehat{F}_i(y) = \begin{cases} 0, & y < 0 \\ y \cdot \widehat{p}_i/\epsilon, & 0 \leq y < \epsilon \\ \widehat{p}_i + (1 - \widehat{p}_i)F_{\mathsf{lnorm}}(y - \epsilon; \widehat{\mu}_i, \widehat{\sigma}^2), & y \geq \epsilon \end{cases}
$$

Choose $\epsilon < \min\{Y_i | Y_i > 0\}$

# Empirical Variogram

Calculate empirical variogram of $\Phi^{-1}\{\widehat{F_1}(y_1)\}, \ldots, \Phi^{-1}\{\widehat{F_n}(y_n)\}$ using R `gstat` package.

## Fitted Exponential Variogram

|   | model | psill | range |
|---|-------|-------|-------|
| 1 | Nug | 0.7201172 | 0.00 |
| 2 | Exp | 0.1581205 | 28559.81 |

# Fitted Exponential Variogram

```
    model      psill       range
1    Nug 0.7201172       0.00
2    Exp 0.1581205 28559.81
```

## Estimated Exponential Correlation Function

$$
\begin{aligned}
\widehat{\alpha_N} &= \frac{0.1581205}{0.1581205 + 0.7201172} &\approx& \quad 0.18 \\
\widehat{\alpha_R} &\approx 28560
\end{aligned}
$$

For plot $i$ and plot $j$ separated by $h_{ij}$ meters,
$\widehat{\text{corr}}\left[\Phi^{-1}\{F_i(y_i)\}, \Phi^{-1}\{F_j(y_j)\}\right] = 0.18\exp(-h_{ij}/28560).$

## Estimated Exponential Correlation Function

$$\widehat{\alpha_N} = \frac{0.1581205}{0.1581205 + 0.7201172} \approx 0.18$$
$$\widehat{\alpha_R} \approx 28560$$

For plot $i$ and plot $j$ separated by $h_{ij}$ meters,
$\widehat{corr}\left[\Phi^{-1}\{F_i(y_i)\}, \Phi^{-1}\{F_j(y_j)\}\right] = 0.18 \exp(-h_{ij}/28560)$.

## Fitted Variogram for Intercept-only Models

$$\text{logit}(p_i) = \beta_0$$
$$\mu_i = \mu$$

# Fitted Variogram for Intercept-only Models

$$\text{logit}(p_i) = \beta_0$$
$$\mu_i = \mu$$

# Survey Unit 0

# Block Kriging

# Kriging With Normal Data

Suppose

$$\begin{bmatrix} \boldsymbol{Z}_o \\ \boldsymbol{Z}_u \end{bmatrix} \sim N\left( \begin{bmatrix} \mu_o \\ \mu_u \end{bmatrix}, \begin{bmatrix} \Sigma_o & \Sigma_{ou} \\ \Sigma'_{ou} & \Sigma_u \end{bmatrix} \right)$$

where

$$\begin{aligned} \boldsymbol{Z}_o &\sim N(\mu_o, \Sigma_o) \text{ are the observed data} \\ \boldsymbol{Z}_u &\sim N(\mu_u, \Sigma_u) \text{ are the unobserved data} \\ \Sigma_{ou} &= \text{cov}(\boldsymbol{Z}_o, \boldsymbol{Z}_u) \text{ is the cross-covariance matrix} \end{aligned}$$

## Predicting $Z_u$

Then

$$
\begin{aligned}
E(Z_u|Z_o) &= \mu_u + \Sigma'_{ou}\Sigma_o^{-1}(Z_o - \mu_o) \\
\text{var}(Z_u|Z_o) &= \Sigma_u - \Sigma'_{ou}\Sigma_o^{-1}\Sigma_{ou}
\end{aligned}
$$

## Predicting $Z_u$

Then

$$
\begin{aligned}
E(Z_u|Z_o) &= \mu_u + \Sigma'_{ou}\Sigma_o^{-1}(Z_o - \mu_o) \\
\mathrm{var}(Z_u|Z_o) &= \Sigma_u - \Sigma'_{ou}\Sigma_o^{-1}\Sigma_{ou}
\end{aligned}
$$

Empirical Best Linear Unbiased Predictor is

$$
\widehat{Z}_u = \widehat{\mu}_u + \widehat{\Sigma}'_{ou}\widehat{\Sigma}_o^{-1}(Z_o - \widehat{\mu}_o),
$$

where $\widehat{\mu}_u$ and $\widehat{\mu}_o$ are weighted least squares estimates.

## Predicting $Z_u$

Then

$$
\begin{aligned}
E(Z_u|Z_o) &= \mu_u + \Sigma'_{ou}\Sigma_o^{-1}(Z_o - \mu_o) \\
\mathrm{var}(Z_u|Z_o) &= \Sigma_u - \Sigma'_{ou}\Sigma_o^{-1}\Sigma_{ou}
\end{aligned}
$$

Empirical Best Linear Unbiased Predictor is

$$
\widehat{Z}_u = \widehat{\mu}_u + \widehat{\Sigma}'_{ou}\widehat{\Sigma}_o^{-1}(Z_o - \widehat{\mu}_o),
$$

where $\widehat{\mu}_u$ and $\widehat{\mu}_o$ are weighted least squares estimates.

Block kriging estimate of total is $\widehat{T} = \mathbf{1}'\widehat{Z}_u + \mathbf{1}'Z_o$.

## Predicting $Z_{\mathsf{u}}$

Then

$$
\begin{aligned}
E(Z_{\mathsf{u}}|Z_{\mathsf{o}}) &= \mu_{\mathsf{u}} + \Sigma'_{\mathsf{ou}}\Sigma_{\mathsf{o}}^{-1}(Z_{\mathsf{o}} - \mu_{\mathsf{o}}) \\
\operatorname{var}(Z_{\mathsf{u}}|Z_{\mathsf{o}}) &= \Sigma_{\mathsf{u}} - \Sigma'_{\mathsf{ou}}\Sigma_{\mathsf{o}}^{-1}\Sigma_{\mathsf{ou}}
\end{aligned}
$$

Empirical Best Linear Unbiased Predictor is

$$
\widehat{Z}_{\mathsf{u}} = \widehat{\mu}_{\mathsf{u}} + \widehat{\Sigma}'_{\mathsf{ou}}\widehat{\Sigma}_{\mathsf{o}}^{-1}(Z_{\mathsf{o}} - \widehat{\mu}_{\mathsf{o}}),
$$

where $\widehat{\mu}_{\mathsf{u}}$ and $\widehat{\mu}_{\mathsf{o}}$ are weighted least squares estimates.

Block kriging estimate of total is $\widehat{T} = \mathbf{1}'\widehat{Z}_{\mathsf{u}} + \mathbf{1}'Z_{\mathsf{o}}$.

$\operatorname{var}(\widehat{T} - T)$ depends on $\Sigma$.

## Block Kriging with the *Y*'s

Let $\boldsymbol{Y}_O = [Y_1, \ldots, Y_n]$ be the vector of responses on the observed plots and $\boldsymbol{Y}_U = [Y_{n+1}, \ldots, Y_{n+m}]$ be the vector of responses on the unobserved plots.

Definitions:

$$
\begin{aligned}
Z_i &= \Phi^{-1}\{F_i(Y_i)\}, i = 1, \ldots, n+m \\
\boldsymbol{Z}_O &= Z_1, \ldots, Z_n \\
\boldsymbol{Z}_U &= Z_{n+1}, \ldots, Z_{n+m}
\end{aligned}
$$

## Block Kriging with the $Y$'s

Let $\boldsymbol{Y}_O = [Y_1, \ldots, Y_n]$ be the vector of responses on the observed plots and $\boldsymbol{Y}_U = [Y_{n+1}, \ldots, Y_{n+m}]$ be the vector of responses on the unobserved plots.

Definitions:

$$
\begin{aligned}
Z_i &= \Phi^{-1}\{F_i(Y_i)\}, i = 1, \ldots, n+m \\
\boldsymbol{Z}_O &= Z_1, \ldots, Z_n \\
\boldsymbol{Z}_U &= Z_{n+1}, \ldots, Z_{n+m} \\
\widehat{Z}_i &= \Phi^{-1}\{\widehat{F}_i(Y_i)\}, i = 1, \ldots, n \\
\widehat{\boldsymbol{z}}_O &= \widehat{Z}_1, \ldots, \widehat{Z}_n
\end{aligned}
$$

## Exploit the Copula Normalizing Transformation

Copula model:

$$\begin{bmatrix} \boldsymbol{Z}_\text{o} \\ \boldsymbol{Z}_\text{u} \end{bmatrix} \sim N\left( \begin{bmatrix} \boldsymbol{0}_n \\ \boldsymbol{0}_m \end{bmatrix}, \begin{bmatrix} \Sigma_\text{o} & \Sigma_\text{ou} \\ \Sigma'_\text{ou} & \Sigma_\text{u} \end{bmatrix} \right)$$

## Exploit the Copula Normalizing Transformation

Copula model:

$$\begin{bmatrix} \mathbf{Z}_\mathsf{o} \\ \mathbf{Z}_\mathsf{u} \end{bmatrix} \sim N \left( \begin{bmatrix} \mathbf{0}_n \\ \mathbf{0}_m \end{bmatrix}, \begin{bmatrix} \Sigma_\mathsf{o} & \Sigma_\mathsf{ou} \\ \Sigma'_\mathsf{ou} & \Sigma_\mathsf{u} \end{bmatrix} \right)$$

Predict $\mathbf{Z}_\mathsf{u}$ as

$$\widehat{\mathbf{Z}}_\mathsf{u} = \widehat{\Sigma}'_\mathsf{ou} \widehat{\Sigma}_\mathsf{o}^{-1} \widehat{\mathbf{Z}}_\mathsf{o}.$$

## Exploit the Copula Normalizing Transformation

Copula model:

$$\begin{bmatrix} \boldsymbol{Z}_\mathsf{O} \\ \boldsymbol{Z}_\mathsf{U} \end{bmatrix} \sim N\left( \begin{bmatrix} \boldsymbol{0}_n \\ \boldsymbol{0}_m \end{bmatrix}, \begin{bmatrix} \Sigma_\mathsf{O} & \Sigma_\mathsf{OU} \\ \Sigma'_\mathsf{OU} & \Sigma_\mathsf{U} \end{bmatrix} \right)$$

Predict $\boldsymbol{Z}_\mathsf{U}$ as

$$\widehat{\boldsymbol{Z}}_\mathsf{U} = \widehat{\Sigma}'_\mathsf{OU} \widehat{\Sigma}_\mathsf{O}^{-1} \widehat{\boldsymbol{Z}}_\mathsf{O}.$$

Then let $\widehat{Y}_i = \widehat{F}_i^{-1}\{\Phi(\widehat{Z}_i)\}, i = 1, \ldots, n + m$

## Exploit the Copula Normalizing Transformation

Copula model:

$$\begin{bmatrix} \mathbf{Z}_O \\ \mathbf{Z}_U \end{bmatrix} \sim N \left( \begin{bmatrix} \mathbf{0}_n \\ \mathbf{0}_m \end{bmatrix}, \begin{bmatrix} \Sigma_O & \Sigma_{OU} \\ \Sigma'_{OU} & \Sigma_U \end{bmatrix} \right)$$

Predict $\mathbf{Z}_U$ as

$$\widehat{\mathbf{Z}}_U = \widehat{\Sigma}'_{OU} \widehat{\Sigma}_O^{-1} \widehat{\mathbf{Z}}_O.$$

Then let $\widehat{Y}_i = \widehat{F}_i^{-1} \{ \Phi(\widehat{Z}_i) \}, i = 1, \ldots, n + m$

Estimate total $T$ as $\sum_{i=1}^{n+m} \widehat{Y}_i$.

## Bad News

This doesn't work.

## Bad News

This doesn't work.

$\widehat{\boldsymbol{Z}}_{\text{u}}$ estimates $E(\boldsymbol{Z}_{\text{u}}|\boldsymbol{Z}_{\text{o}})$, and $F_i^{-1}\{\Phi(E(Z_i|\boldsymbol{Z}_{\text{o}}))\} \neq E(Y_i|\boldsymbol{Y}_{\text{o}})$.

## Bad News

This doesn't work.

$\widehat{\boldsymbol{Z}}_{\mathsf{u}}$ estimates $E(\boldsymbol{Z}_{\mathsf{u}}|\boldsymbol{Z}_{\mathsf{o}})$, and $F_i^{-1}\{\Phi(E(Z_i|\boldsymbol{Z}_{\mathsf{o}}))\} \neq E(Y_i|\boldsymbol{Y}_{\mathsf{o}})$.

Prediction does not account for variance of parameter estimates.

## Proposed Solution–Parametric Bootstrap

- Estimate parameters $\theta = [\alpha_N, \alpha_R, \theta, \gamma, \sigma^2]$ by maximizing copula likelihood.

## Proposed Solution–Parametric Bootstrap

- Estimate parameters $\boldsymbol{\theta} = [\alpha_N, \alpha_R, \boldsymbol{\theta}, \boldsymbol{\gamma}, \sigma^2]$ by maximizing copula likelihood.
- Generate $k = 1, \ldots, N_b$ realizations of $\widehat{\boldsymbol{\theta}}$ from estimated asymptotic distribution.

## Proposed Solution–Parametric Bootstrap

- Estimate parameters $\boldsymbol{\theta} = [\alpha_N, \alpha_R, \boldsymbol{\theta}, \boldsymbol{\gamma}, \sigma^2]$ by maximizing copula likelihood.
- Generate $k = 1, \ldots, N_b$ realizations of $\widehat{\boldsymbol{\theta}}$ from estimated asymptotic distribution.
- For each realization of $\widehat{\boldsymbol{\theta}}$,
  - Generate $\widehat{\boldsymbol{Z}}_\mathsf{u} \sim N(\boldsymbol{0}, \widehat{\Sigma}_\mathsf{u} - \widehat{\Sigma}'_\mathsf{ou} \widehat{\Sigma}_\mathsf{o}^{-1} \widehat{\Sigma}_\mathsf{ou})$.

## Proposed Solution–Parametric Bootstrap

- Estimate parameters $\boldsymbol{\theta} = [\alpha_N, \alpha_R, \boldsymbol{\theta}, \boldsymbol{\gamma}, \sigma^2]$ by maximizing copula likelihood.
- Generate $k = 1, \ldots, N_b$ realizations of $\widehat{\boldsymbol{\theta}}$ from estimated asymptotic distribution.
- For each realization of $\widehat{\boldsymbol{\theta}}$,
  - Generate $\widehat{\boldsymbol{Z}}_U \sim N(\boldsymbol{0}, \widehat{\Sigma}_U - \widehat{\Sigma}'_{OU} \widehat{\Sigma}_O^{-1} \widehat{\Sigma}_{OU})$.
  - Transform elements of $\widehat{\boldsymbol{Z}}_U$ to data scale: $\widehat{Y}_i = \widehat{F}_i^{-1} \{\Phi(\widehat{Z}_i)\}$, where $\widehat{F}_i$ is based on the realization of $\widehat{\boldsymbol{\theta}}$.

## Proposed Solution–Parametric Bootstrap

- Estimate parameters $\boldsymbol{\theta} = [\alpha_N, \alpha_R, \boldsymbol{\theta}, \boldsymbol{\gamma}, \sigma^2]$ by maximizing copula likelihood.
- Generate $k = 1, \ldots, N_b$ realizations of $\widehat{\boldsymbol{\theta}}$ from estimated asymptotic distribution.
- For each realization of $\widehat{\boldsymbol{\theta}}$,
  - Generate $\widehat{\boldsymbol{Z}}_{\mathsf{u}} \sim N(\boldsymbol{0}, \widehat{\Sigma}_{\mathsf{u}} - \widehat{\Sigma}'_{\mathsf{ou}} \widehat{\Sigma}_{\mathsf{o}}^{-1} \widehat{\Sigma}_{\mathsf{ou}})$.
  - Transform elements of $\widehat{\boldsymbol{Z}}_{\mathsf{u}}$ to data scale: $\widehat{Y}_i = \widehat{F}_i^{-1}\{\Phi(\widehat{Z}_i)\}$, where $\widehat{F}_i$ is based on the realization of $\widehat{\boldsymbol{\theta}}$.
  - Calculate $\widehat{T} = \boldsymbol{1}'\widehat{\boldsymbol{Y}}_{\mathsf{u}} + \boldsymbol{1}'\boldsymbol{Y}_{\mathsf{o}}$

## Proposed Solution–Parametric Bootstrap

- Estimate parameters $\theta = [\alpha_N, \alpha_R, \theta, \gamma, \sigma^2]$ by maximizing copula likelihood.
- Generate $k = 1, \ldots, N_b$ realizations of $\widehat{\theta}$ from estimated asymptotic distribution.
- For each realization of $\widehat{\theta}$,
  - Generate $\widehat{Z}_U \sim N(\mathbf{0}, \widehat{\Sigma}_U - \widehat{\Sigma}'_{OU}\widehat{\Sigma}_O^{-1}\widehat{\Sigma}_{OU})$.
  - Transform elements of $\widehat{Z}_U$ to data scale: $\widehat{Y}_i = \widehat{F}_i^{-1}\{\Phi(\widehat{Z}_i)\}$, where $\widehat{F}_i$ is based on the realization of $\widehat{\theta}$.
  - Calculate $\widehat{T} = \mathbf{1}'\widehat{Y}_U + \mathbf{1}'Y_O$
- This yields a bootstrapped distribution of $\widehat{T}$'s. Take the median as the point prediction of the block total, and take $\alpha/2$ and $1 - \alpha/2$ quantiles as lower and upper $1 - \alpha$ prediction limits.

## Preliminary Simulations

Covariate vectors $\boldsymbol{X}_1$ and $\boldsymbol{X}_2$ iid Uniform$(0, 1)$.

Bernoulli model: $\text{logit}(p_i) = 1 - 3X_{1\,i}$

Lognormal model: $\log(Y_i) \sim N(1 + 5X_{2\,i}, 1)$
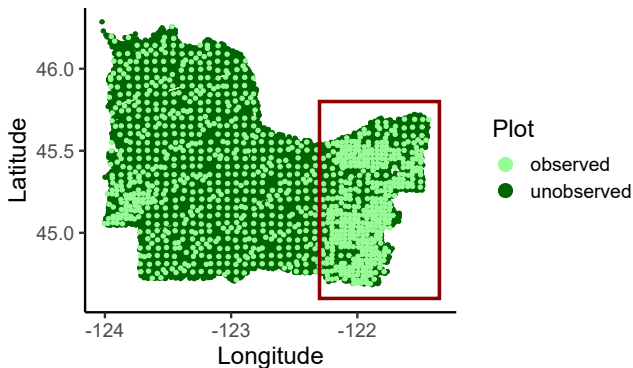
## Preliminary Simulations

Covariate vectors $\boldsymbol{X}_1$ and $\boldsymbol{X}_2$ iid Uniform$(0, 1)$.

Bernoulli model: $\text{logit}(p_i) = 1 - 3X_{1\,i}$

Lognormal model: $\log(Y_i) \sim N(1 + 5X_{2\,i}, 1)$

This gives approximately 40% 0's.

## Locations



481 observed plots
3047 unobserved plots

## Exponential Spatial Dependence

Range: $\alpha_R = 25000$

Nugget: $\alpha_N = 0.3, 0.5, 0.8$ (weak, medium, strong dependence)

## Exponential Spatial Dependence

Range: $\alpha_R = 25000$

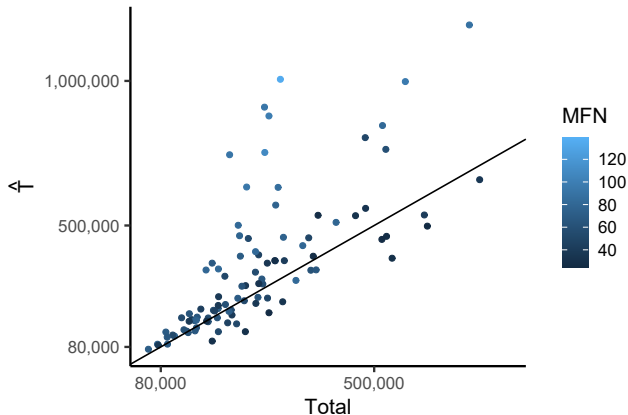Nugget: $\alpha_N = 0.3, 0.5, 0.8$ (weak, medium, strong dependence)
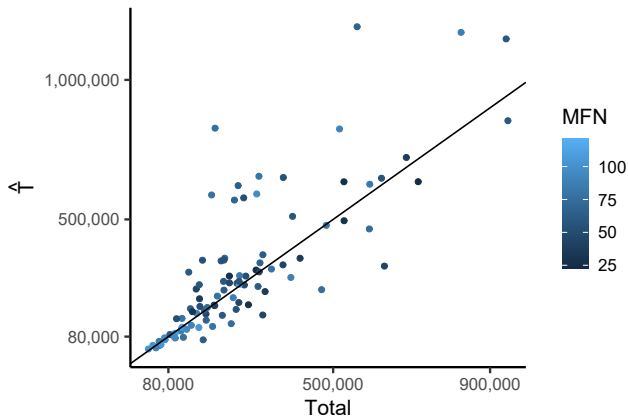
## Results for $\alpha_N = 0.3$



MFN is the median Frobenius norm of $\Sigma - \widehat{\Sigma}$ over the bootstrapped sample.

# Results for $\alpha_N = 0.5$

Motivation
oooo

Model Components
ooooo
oooooo
ooooooooooo

Zero-inflated Spatial Model
oo

Model Fitting
ooooooooooo

Spatial Prediction
oooooooooooooo●

Conclusions
ooo

# Results for $\alpha_N = 0.8$

## Summary

- Gaussian copula models spatial dependence.
- Continuous zero-inflated lognormal marginal models accommodate large percentage of zeros.
- For FIA data, marginal models account for most of the spatial pattern, so spatial prediction not needed.

## Next Steps

- Simulations
- Model assumptions
- Big data

## Thanks