



Modeling compositional data

Background

NAPAP, 1980' s

**Workshop on biological monitoring,
1986**

Dirichlet process: Gary Grunwald, 1987

**Current framework: Dean Billheimer,
1995**

**Other co-workers: Adrian Raftery,
Mariabeth Silkey, Eun-Sug Park**

Compositional data

Vector of proportions

$$\mathbf{z} = (z_1, \dots, z_k)^T \quad z_i > 0 \quad \sum_{i=1}^k z_i = 1 \quad \mathbf{z} \in \nabla^{k-1}$$

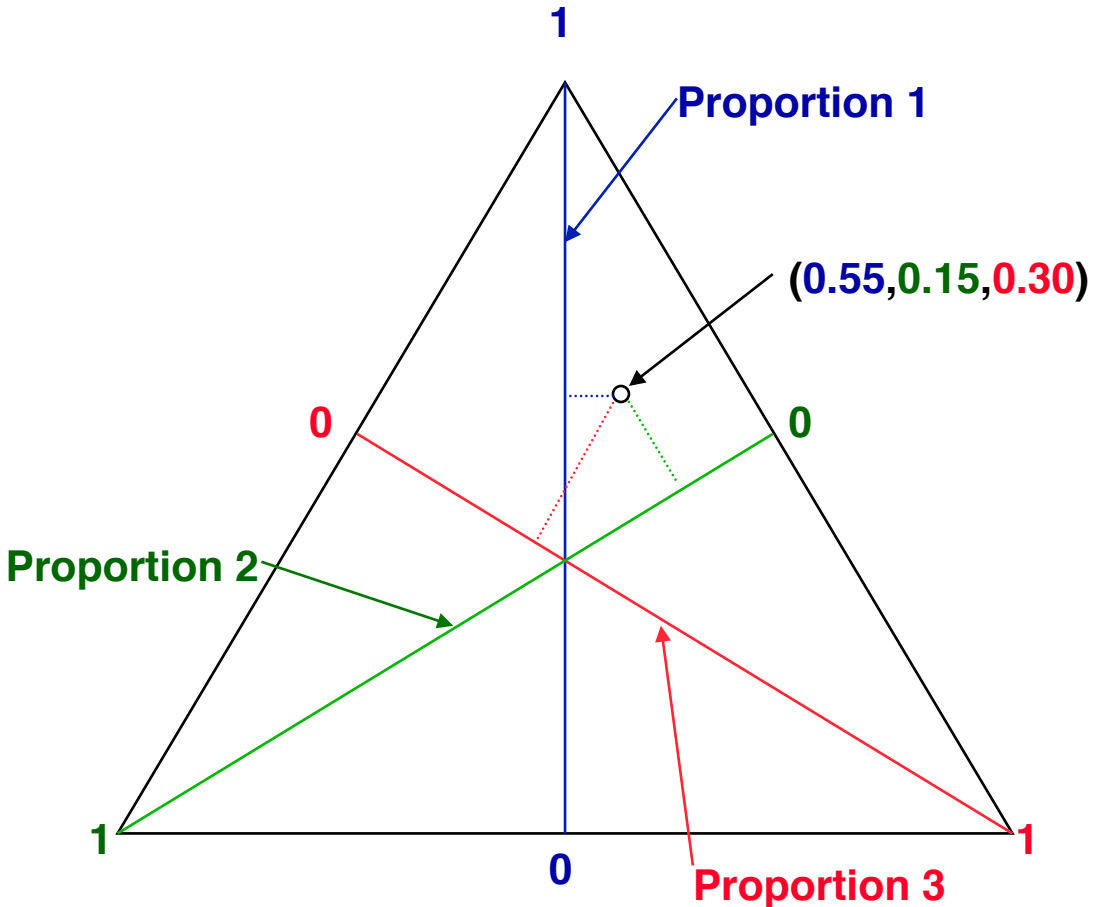
Proportion of taxes in different categories

Composition of rock samples

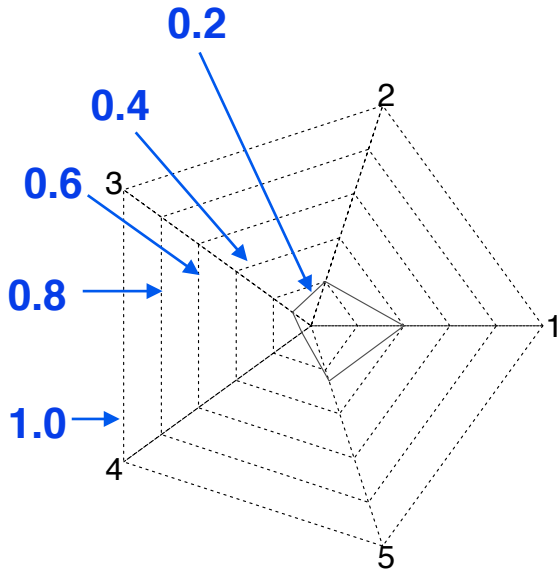
Composition of biological populations

Composition of air pollution

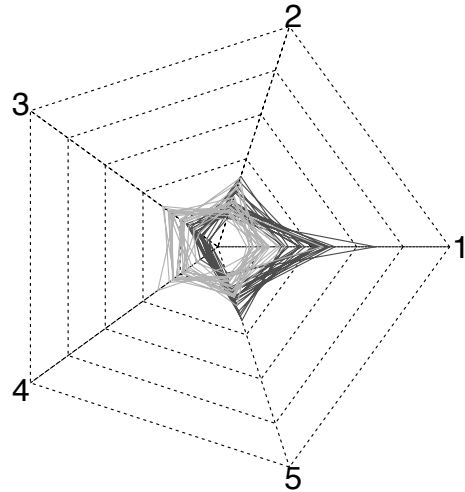
The triangle plot



The spider plot



(0.40,0.20,0.10,0.05,0.25)



An algebra for compositions

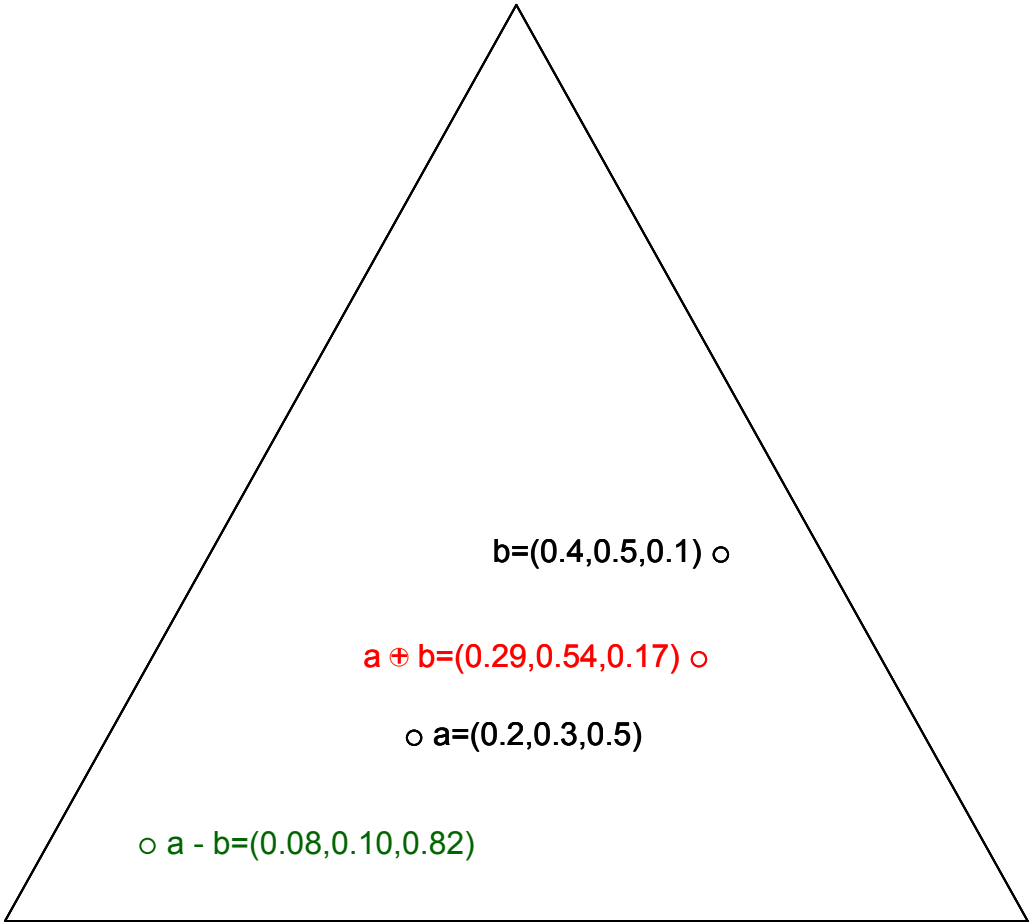
Perturbation: For $\xi, \alpha \in \nabla^{k-1}$ define

$$\xi \oplus \alpha = \left(\frac{\xi_1 \alpha_1}{\sum_1^k \xi_i \alpha_i}, \dots, \frac{\xi_k \alpha_k}{\sum_1^k \xi_i \alpha_i} \right) \in \nabla^{k-1}$$

The composition $\iota = \left(\frac{1}{k}, \dots, \frac{1}{k} \right)$ acts as a zero, so $\xi \oplus \iota = \xi$.

Set $\xi^{-1} = \left(\frac{1}{\xi_1}, \dots, \frac{1}{\xi_k} \right)$ so $\xi \oplus \xi^{-1} = \iota$.

Finally define $\xi - \eta = \xi \oplus \eta^{-1}$.



The logistic normal

$$\text{If } \text{alr}(\mathbf{z}) = \left(\log \frac{z_1}{z_k}, \dots, \log \frac{z_{k-1}}{z_k} \right)^T \sim \text{MVN}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

we say that \mathbf{z} is *logistic normal*, in short $\mathbf{Z} \sim \text{LN}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$.

alr has a unique inverse.

Other distributions on the simplex:

Dirichlet — ratios of independent gammas

“Danish” — ratios of independent inverse Gaussian

Both have very limited correlation structure.

Scalar multiplication

Let a be a scalar. Define

$$\xi \otimes a = \left(\frac{\xi_1^a}{\sum \xi_i^a}, \dots, \frac{\xi_k^a}{\sum \xi_i^a} \right)$$

$(\nabla^{k-1}, \oplus, \otimes)$ is a complete inner product space, with inner product given, e.g., by

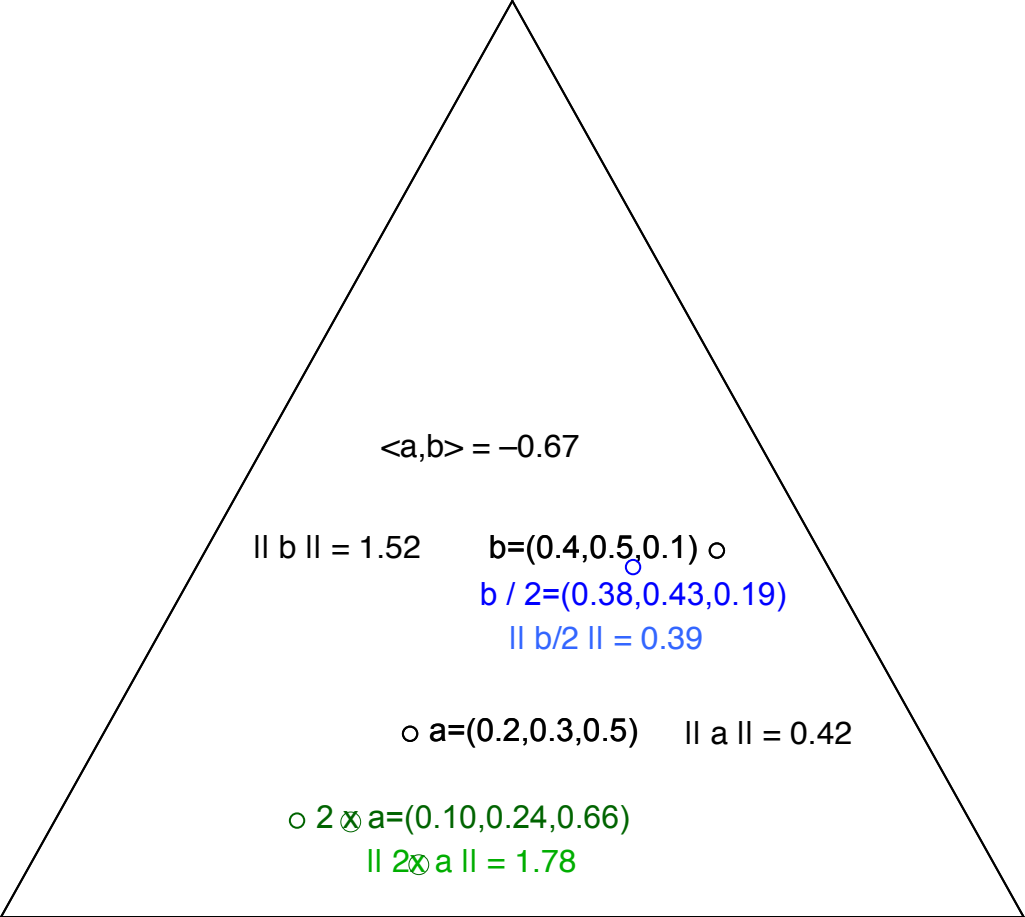
$$\langle \xi, \eta \rangle = \text{alr}(\xi)^T N^{-1} \text{alr}(\eta)$$

N is the precision matrix $N = I + jj^T$

j is a vector of $k-1$ ones.

$\|\xi\| = \langle \xi, \xi \rangle$ is a norm on the simplex.

The inner product and norm are invariant to permutations of the components of the composition.


$$\langle a, b \rangle = -0.67$$

$$\| b \| = 1.52$$

$$b = (0.4, 0.5, 0.1)$$

$$b / 2 = (0.38, 0.43, 0.19)$$

$$\| b / 2 \| = 0.39$$

$$a = (0.2, 0.3, 0.5) \quad \| a \| = 0.42$$

$$2 \otimes a = (0.10, 0.24, 0.66)$$

$$\| 2 \otimes a \| = 1.78$$

Some models

Measurement error:

$$\mathbf{z}_j = \boldsymbol{\xi} \oplus \boldsymbol{\varepsilon}_j \quad \text{where } \boldsymbol{\varepsilon}_j \sim \text{LN}(0, \boldsymbol{\Sigma}) .$$

Regression:

$$\boldsymbol{\xi}_j = \boldsymbol{\xi} \oplus \boldsymbol{\gamma} \otimes \mathbf{u}_j \quad \leftarrow \begin{array}{l} \text{centered} \\ \text{covariate} \end{array}$$

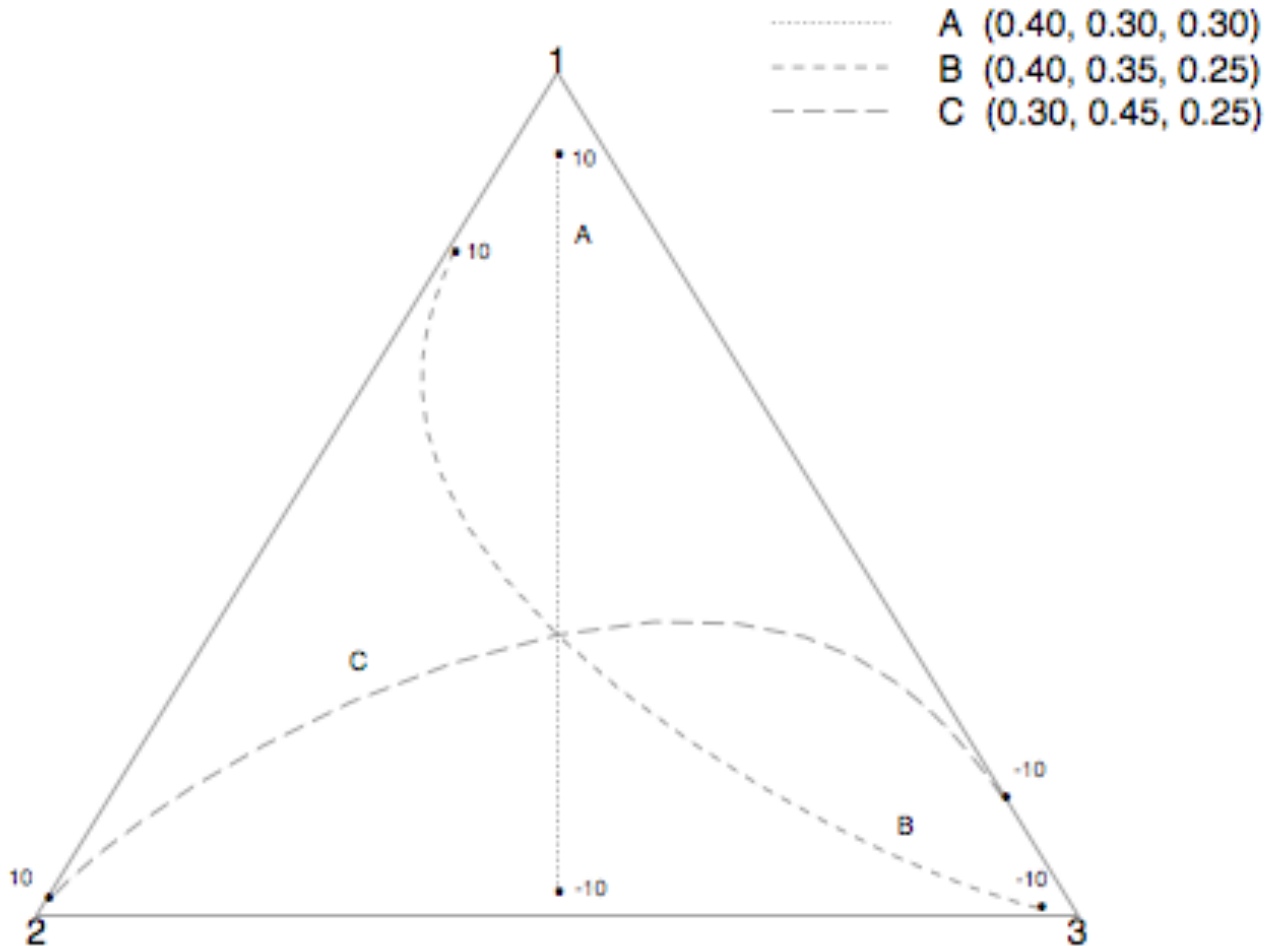
↑
compositions

Correspondence in Euclidean space:

$$\boldsymbol{\mu}_j = \boldsymbol{\beta}_0 + \boldsymbol{\beta}_1 (\mathbf{x}_j - \bar{\mathbf{x}})$$

$$\underset{\boldsymbol{\xi}_j}{\text{alr}^{-1}(\boldsymbol{\mu}_j)} = \underset{\boldsymbol{\xi}}{\text{alr}^{-1}(\boldsymbol{\beta}_0)} \oplus \underset{\boldsymbol{\gamma}}{\text{alr}^{-1}(\boldsymbol{\beta}_1)} \otimes \underset{\mathbf{u}_j}{(\mathbf{x}_j - \bar{\mathbf{x}})}$$

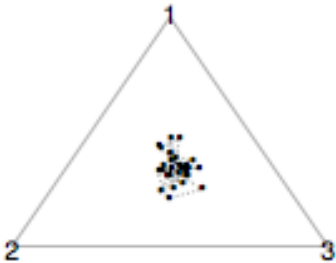
Some regression lines



Time series (AR 1)

$$\mathbf{z}_{k+1} = \phi \otimes \mathbf{z}_k \oplus \varepsilon_k$$

AR parameter = 0.2



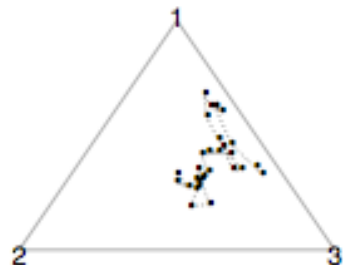
AR parameter = 0.6



AR parameter = 0.95



AR parameter = 1



A source receptor model

Observe relative concentration Y_i of k species at a location over time.

Consider p sources with chemical profiles θ_j . Let α_i be the vector of mixing proportions of the different sources at the receptor on day i .

$$EY_i = \sum_{j=1}^p \alpha_{ij} \theta_j = \Theta \alpha_i$$

$$Y_i = \Theta \alpha_i \oplus \varepsilon_i$$

$\Theta \sim \text{LN}$, $\alpha_i \sim \text{indep LN}$, $\varepsilon_i \sim \text{zero mean LN}$

Juneau air quality

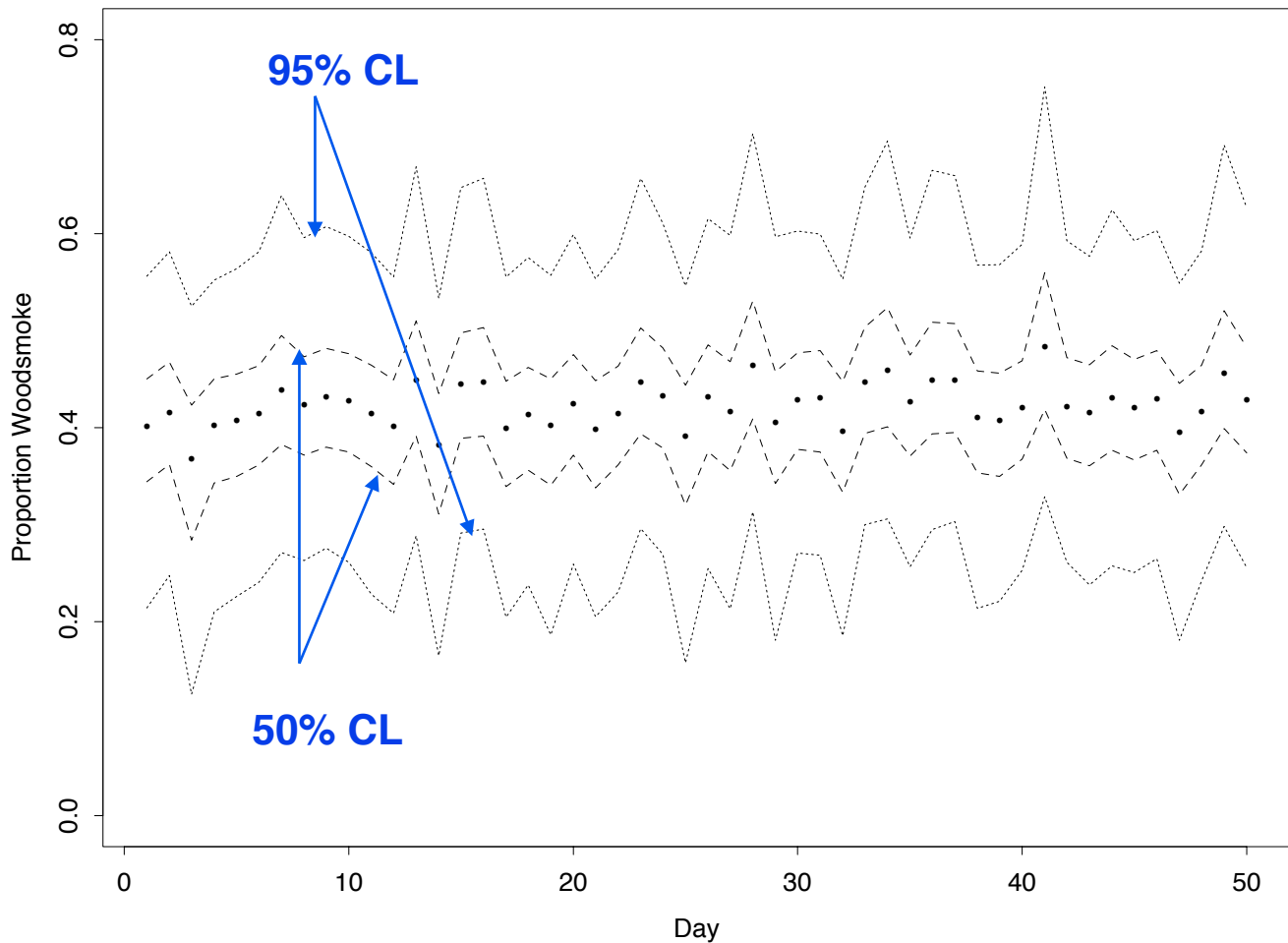
50 observations of relative mass of 5 chemical species. Goal: determine the contribution of wood smoke to local pollution load.

Prior specification:

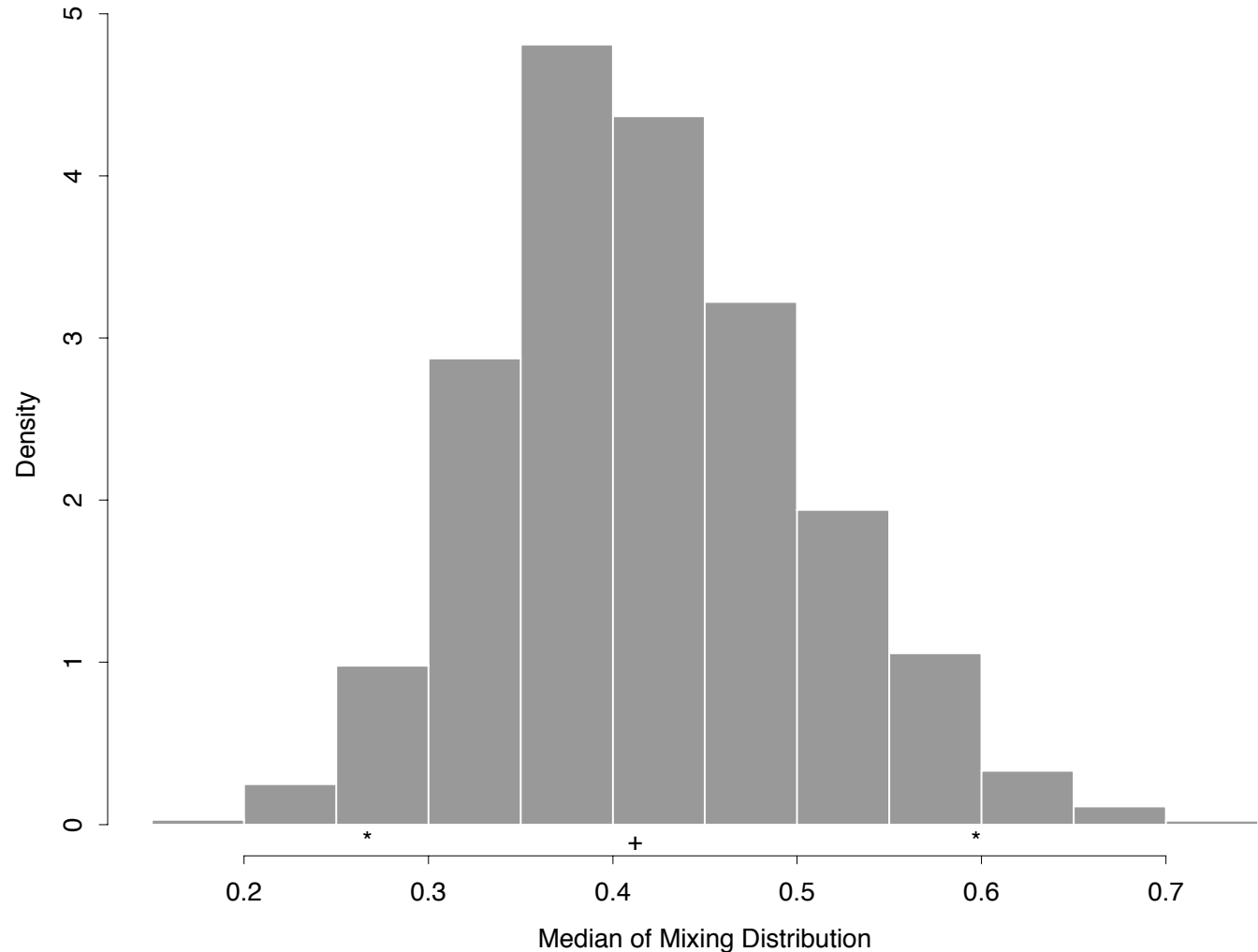
$$\begin{aligned} f(\Theta, \alpha_i, \varepsilon_i, \mu_\alpha, \Gamma, \Sigma_\varepsilon) = \\ f(\alpha_i | \mu_\alpha, \Gamma) f(\varepsilon_i | \Sigma_\varepsilon) f(\mu_\alpha) f(\Gamma) f(\Sigma_\varepsilon) \end{aligned}$$

Inference by MCMC.

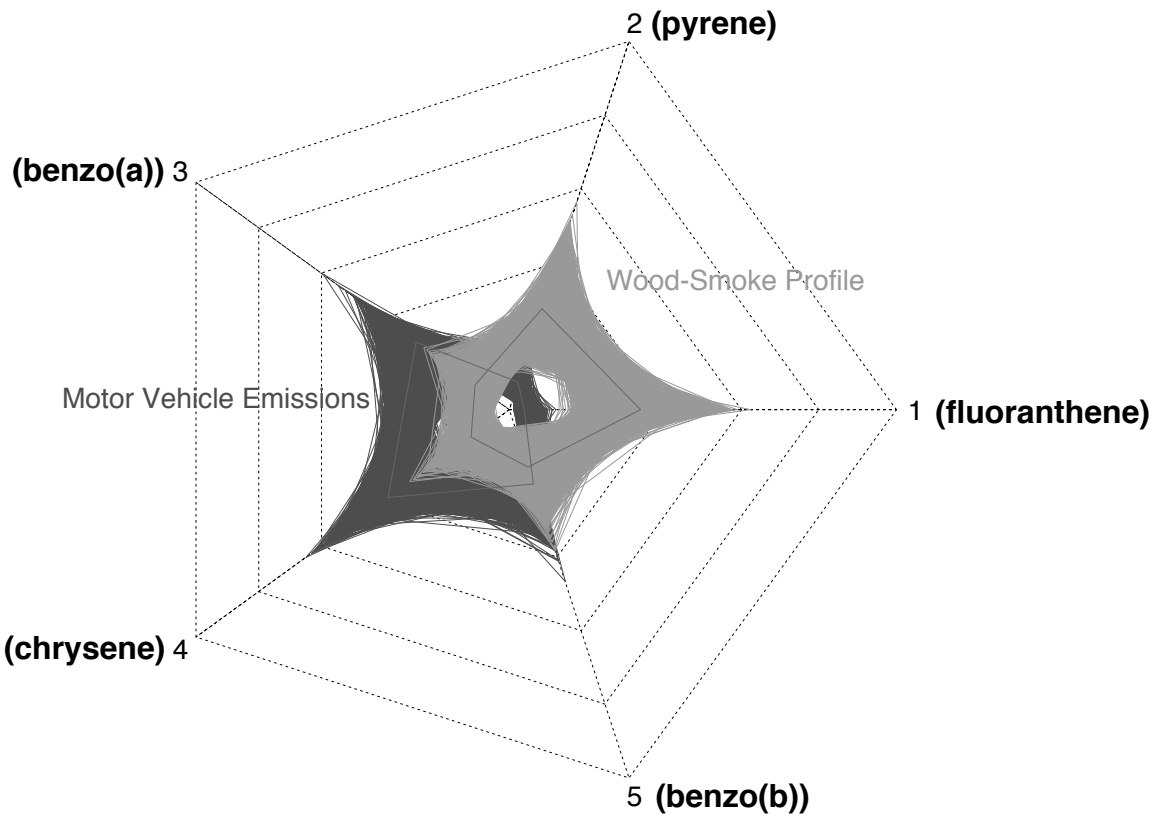
Wood smoke contribution



Wood smoke proportion



Posterior source profiles



State-space model

Space-time model of proportions

State-space model:

z_j unobservable composition $\sim \text{LN}(\mu_j, \Sigma_j)$
 y_j k-vector of counts $\sim \text{Mult}(\sum_{i=1}^k [y_j]_i, z_j)$

Inference using MCMC again

Stability of arthropod food webs

Omnivory thought to destabilize ecological communities

Stability: Capacity to recover from shock (relative abundance in trophic classes)

Mount St. Helens experiment: 6 treatments in 2-way factorial design; 5 reps.

- **Predator manipulation (more omnivores, more specialists, control)**
- **Vegetation disturbance (50% reduction, control)**

Count arthropods, 6 wks after treatment.

Divide into specialized herbivores, general herbivores, predators.

Manipulated species

**Omnivore:
Wolf spider**



**Specialist predator
Big-eyed bug**



Vegetation

fireweed

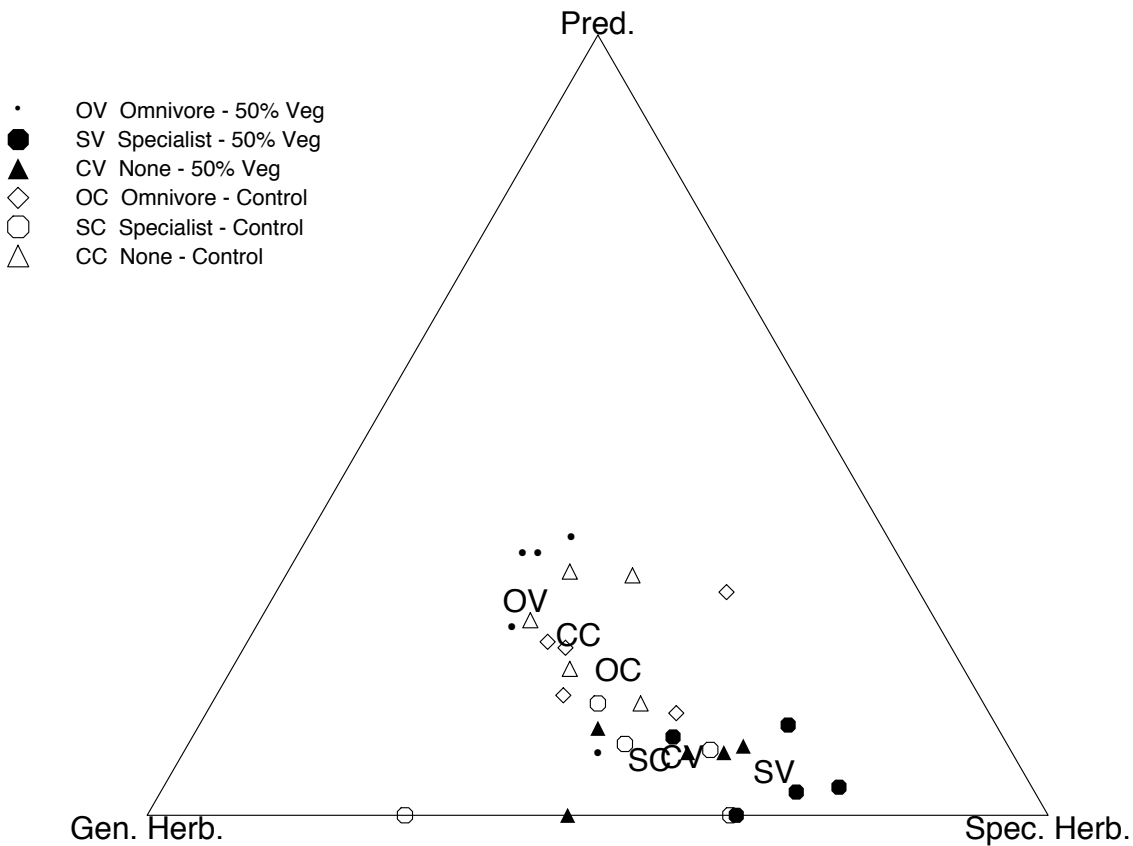


**pearly-
everlasting**

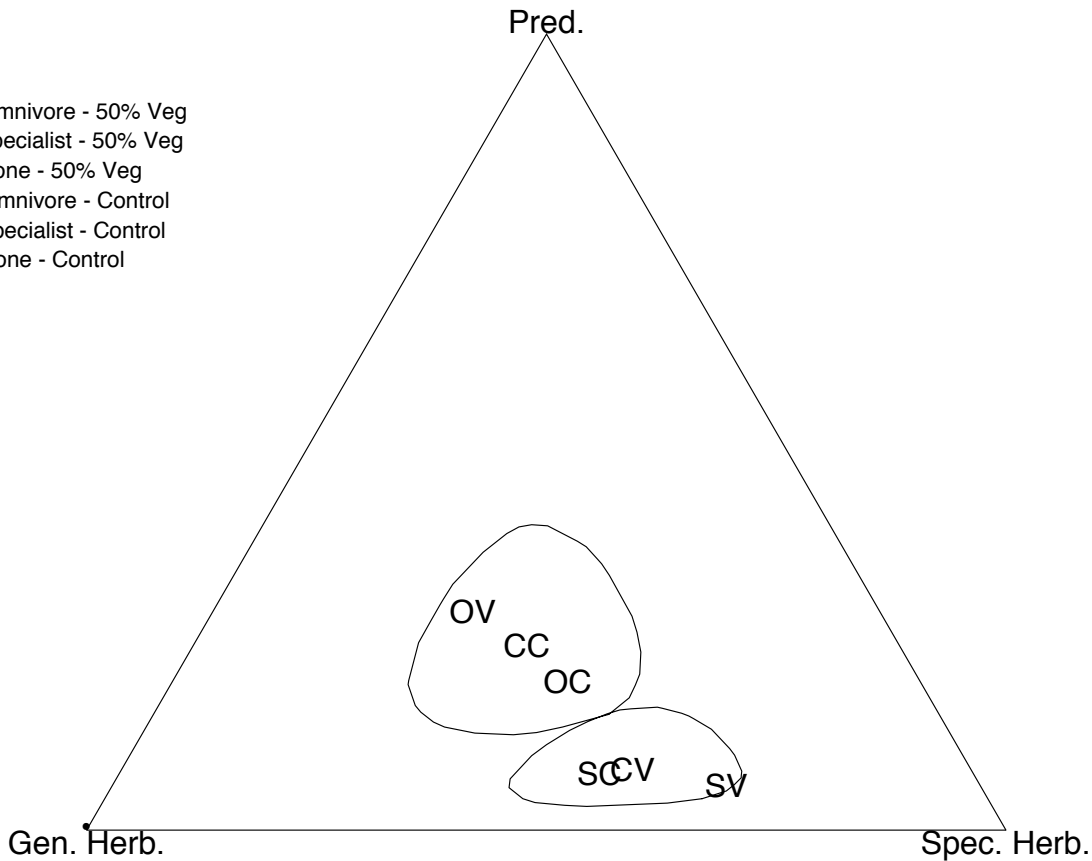


Specification of structure

Σ is generated from independent observations at each treatment mean depends only on treatment



OV Omnivore - 50% Veg
SV Specialist - 50% Veg
CV None - 50% Veg
OC Omnivore - Control
SC Specialist - Control
CC None - Control



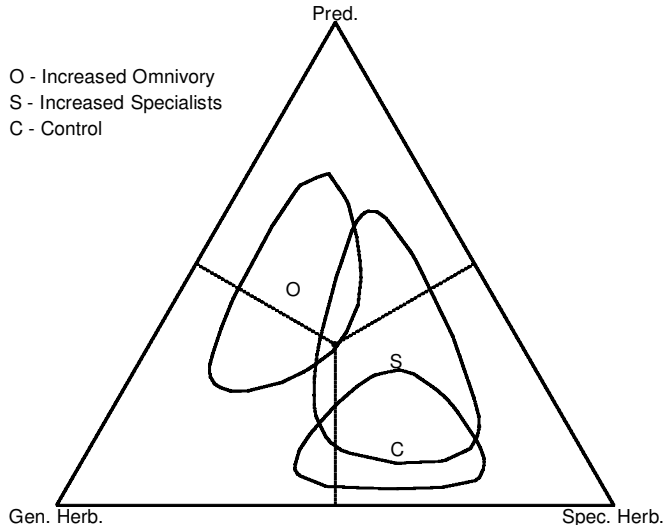
Interaction effect

ANOVA interaction effect

$$z_{ij} - \bar{z}_{i.} - \bar{z}_{.j} + \bar{z}_{..}$$

alr inverse to get

$$\xi_{ij} - \bar{\xi}_{i.} - \bar{\xi}_{.j} \oplus \bar{\xi}_{..}$$



Benthic invertebrates in estuary

**EMAP estuaries monitoring program:
Delaware Bay 1990. 25 locations, 3 grab
samples of bottom sediment during
summer**

Invertebrates in samples classified into

–pollution tolerant



–pollution intolerant



–suspension feeders (control group)



Site j, subsample t

$$z_{jt} \sim \text{LN}(\theta_j + \beta x_j, \Psi)$$

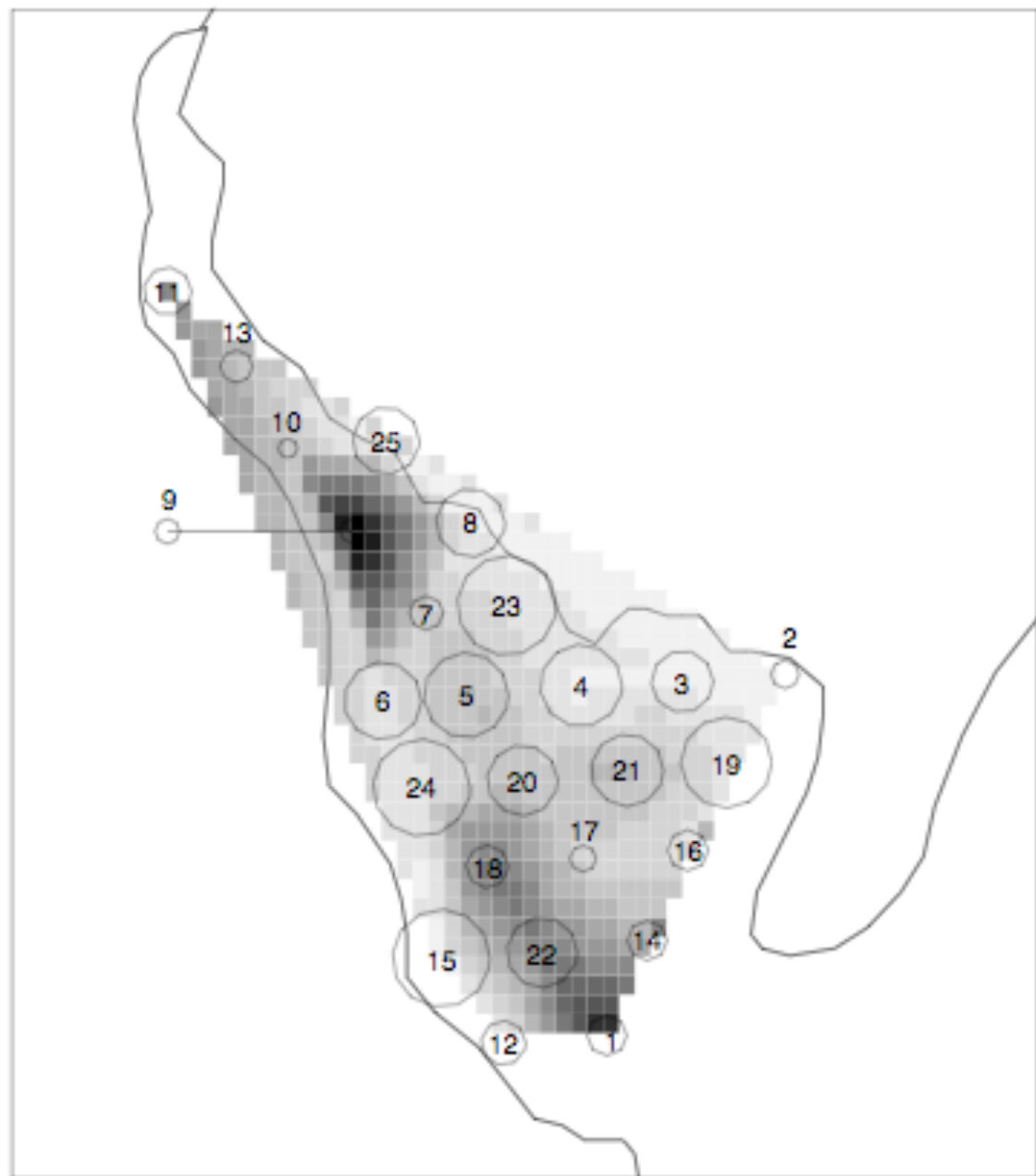
standardized
covariate



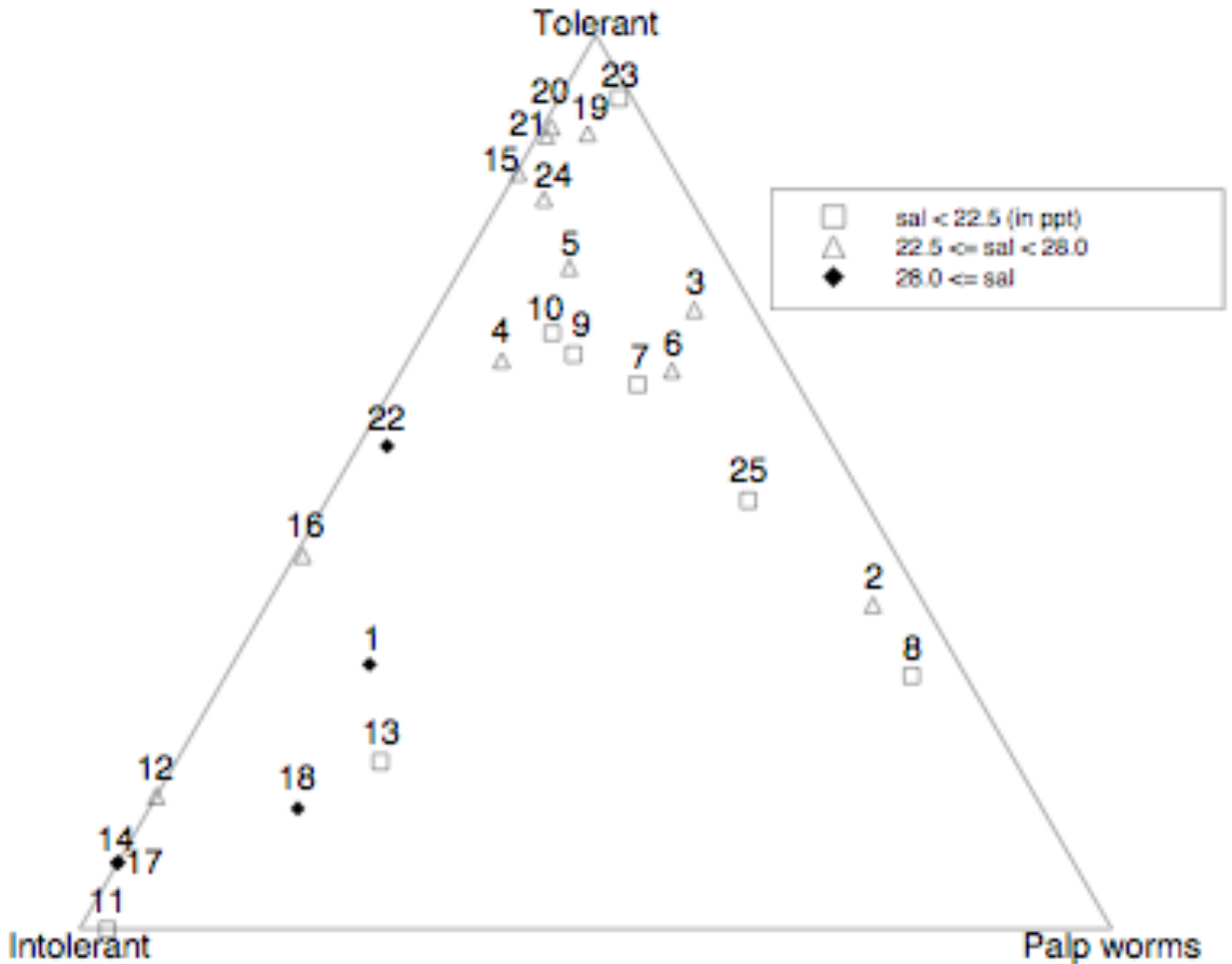
$\theta_j \sim$ **CAR process**

$$E(\theta_j | \theta_{-j}) = \mu + \sum_{k \in N(j)} \frac{\lambda}{n_j} (\theta_k - \mu)$$

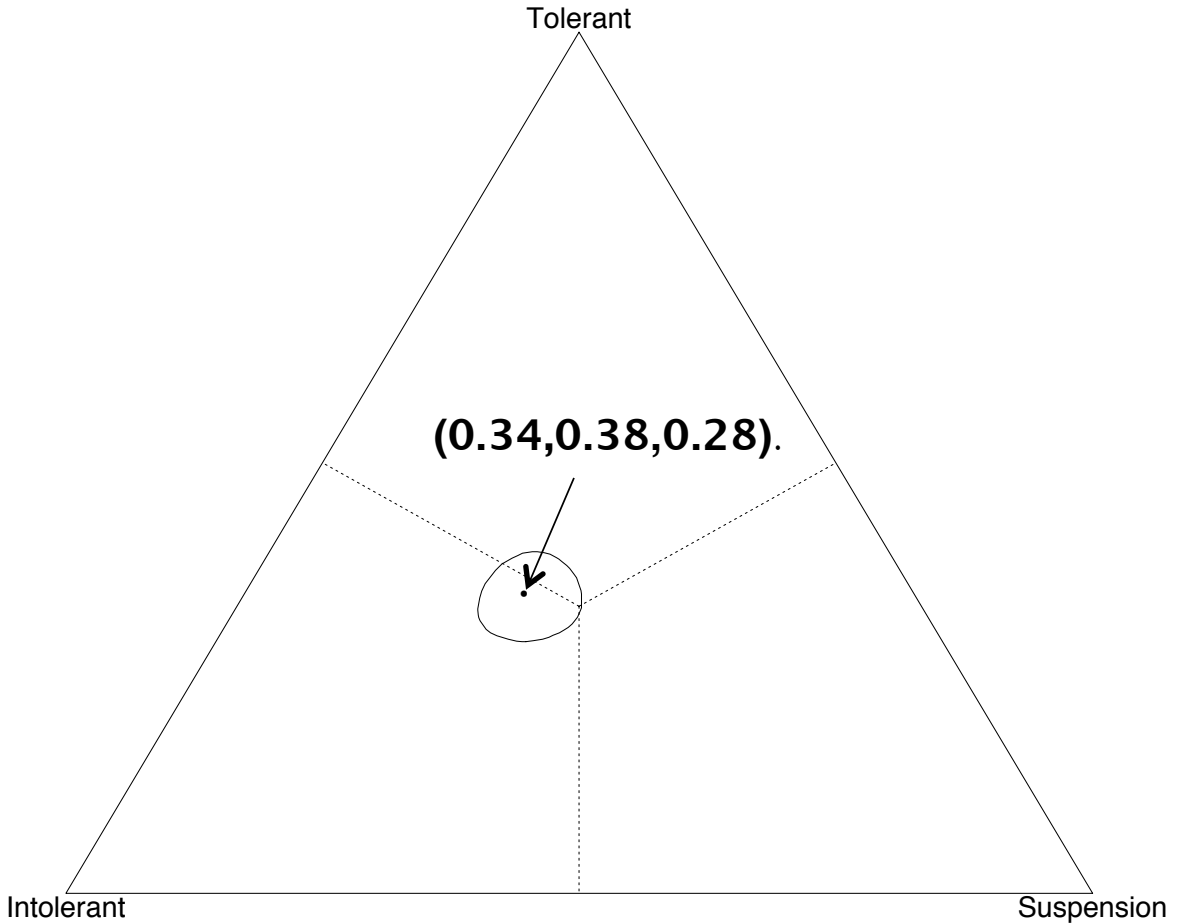
$$\text{Var}(\theta_j | \theta_{-j}) = \frac{\Gamma}{n_j}$$



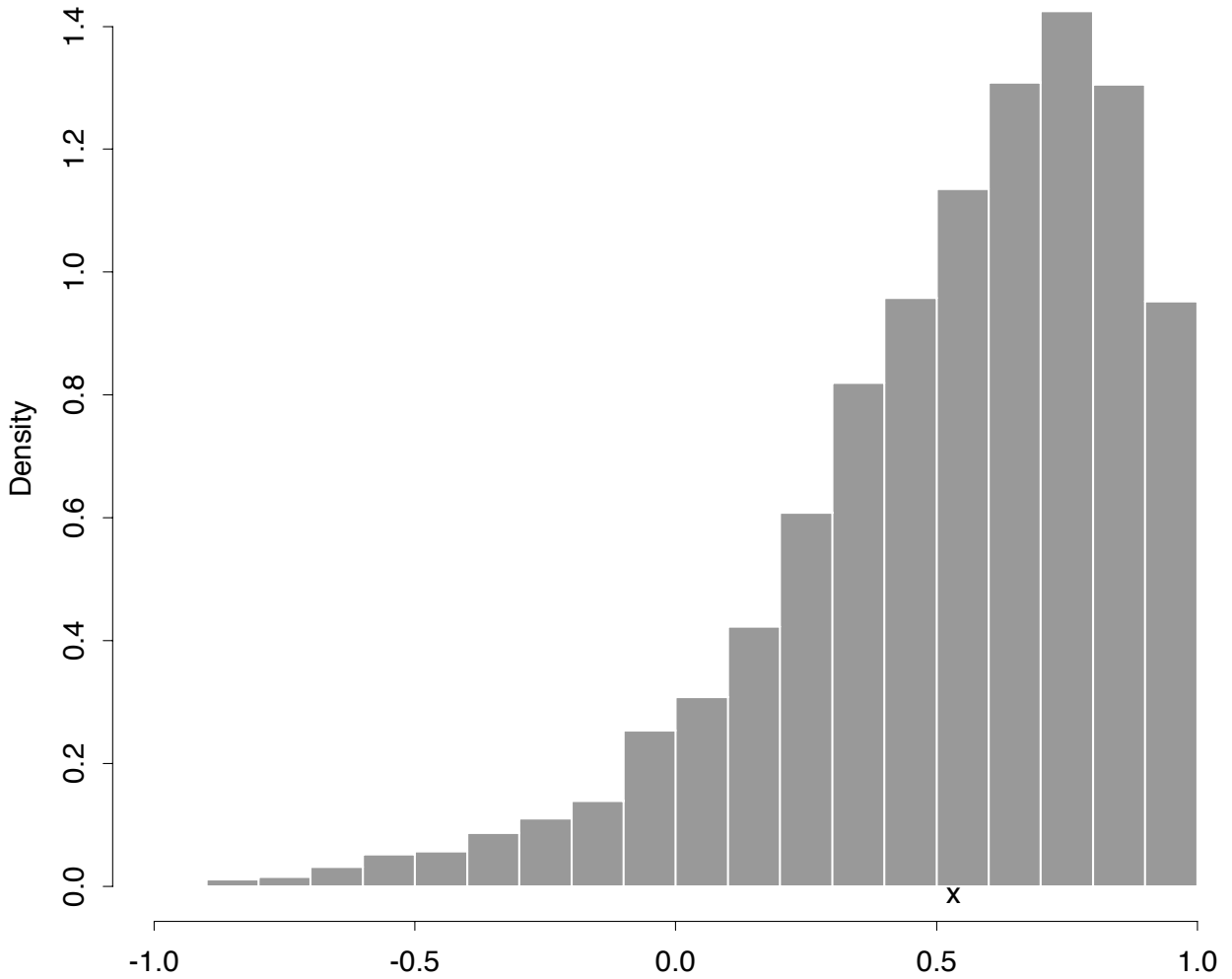
Effect of salinity



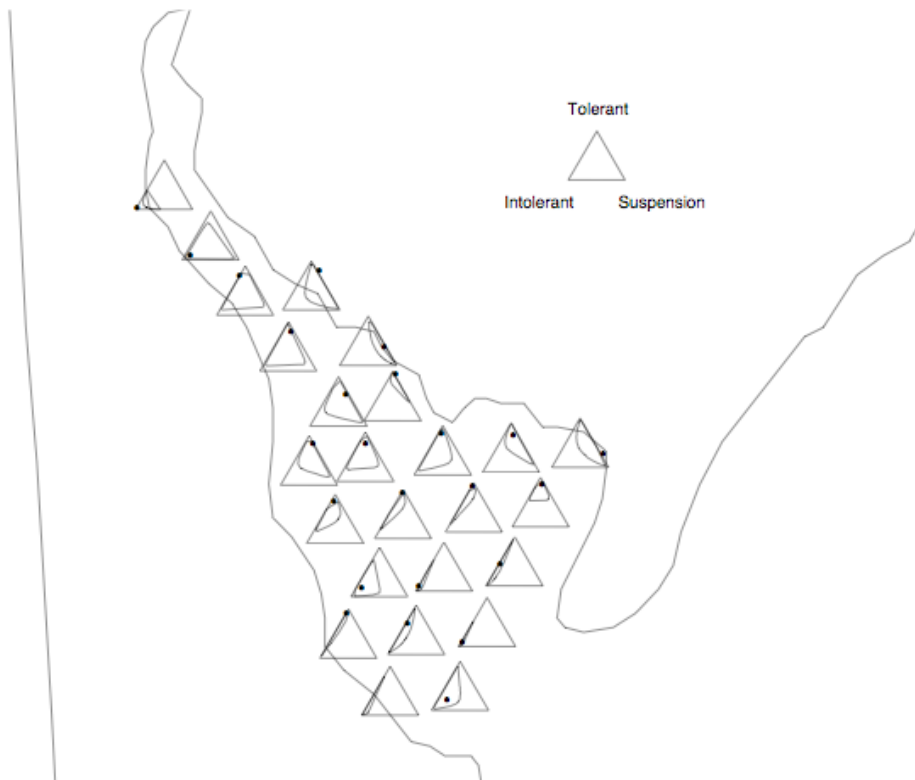
95% Credible Region for Salinity Regression Composition



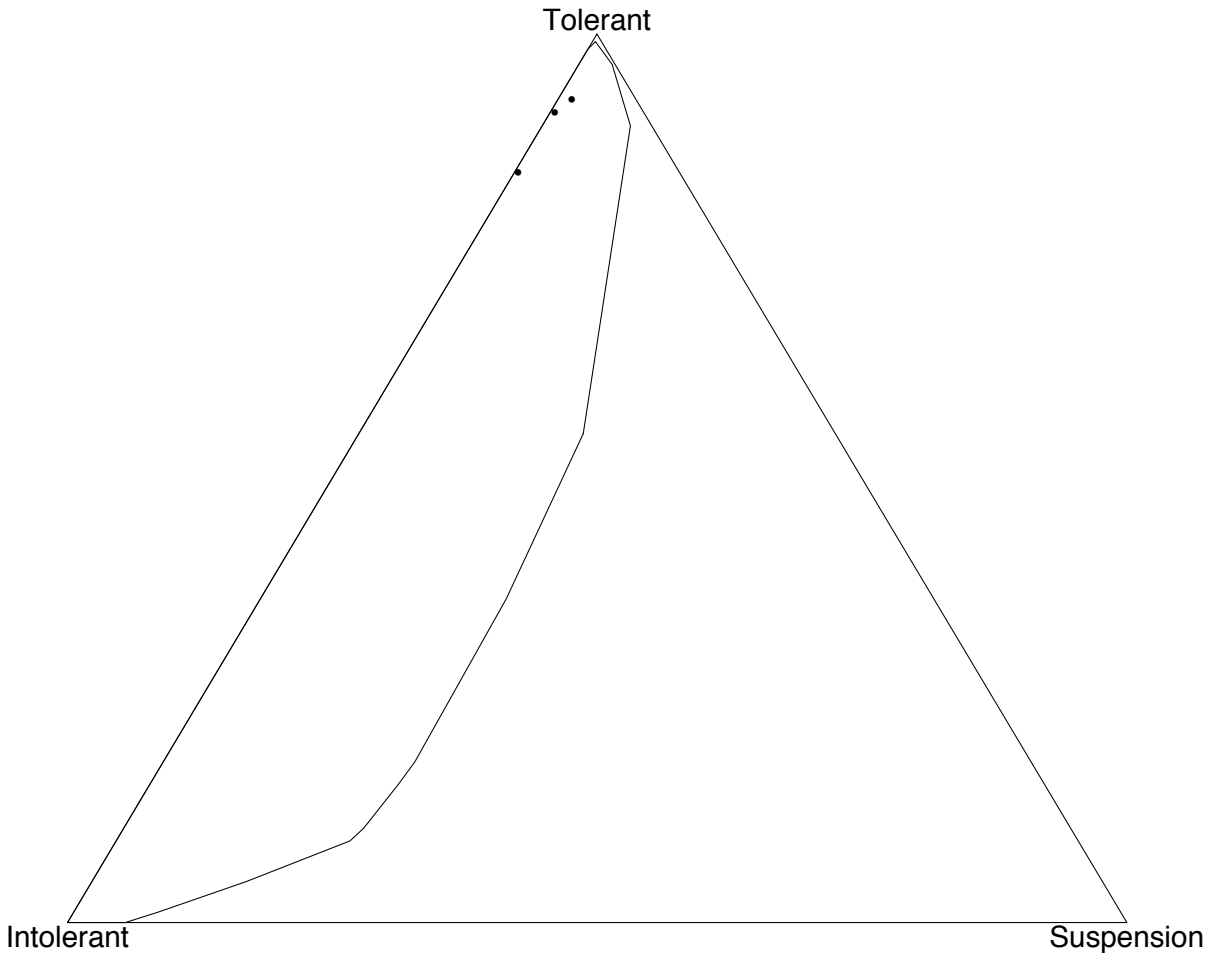
Spatial Dependence Parameter



95% Prediction Regions for Hold-out Sub-Sample Compositions



95% Prediction Region Site 20



95% Prediction Region Site 23

