

Some research in multivariate analysis

Stat 542

Peter Hoff

University of Washington

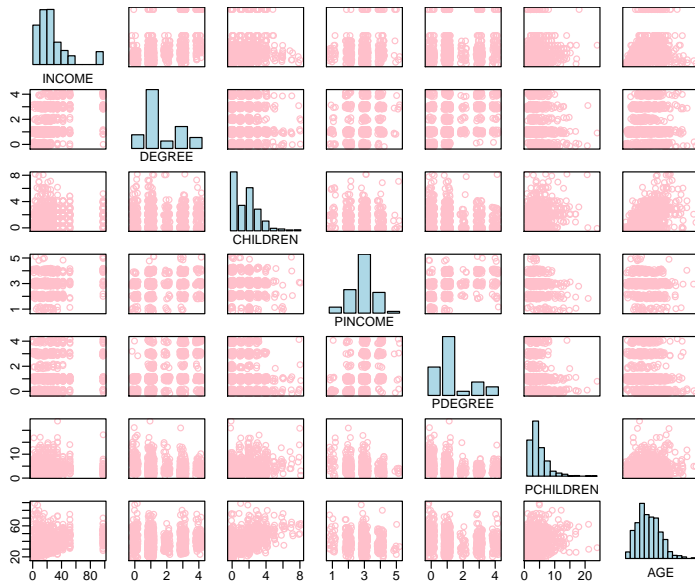
Outline

Copula modeling

Covariance regression

Hierarchical models for matrix-variate data

GSS data



Conditional models

Interest is typically in the **conditional** relationship between pairs of variables, accounting for heterogeneity in other variables of less interest. Standard bivariate rank-based methods are inappropriate.

Model 1

$$\text{INC}_i = \beta_0 + \beta_1 \text{CHILD}_i + \beta_2 \text{DEG}_i + \beta_3 \text{AGE}_i + \beta_4 \text{PCHILD}_i + \beta_5 \text{PINC}_i + \beta_6 \text{PDEG}_i + \epsilon_i$$

p-value for β_1 is 0.11: "little evidence" that $\beta_1 \neq 0$

Model 2

$$\text{CHILD}_i \sim \text{Pois}(\exp\{\beta_0 + \beta_1 \text{INC}_i + \beta_2 \text{DEG}_i + \beta_3 \text{AGE}_i + \beta_4 \text{PCHILD}_i + \beta_5 \text{PINC}_i + \beta_6 \text{PDEG}_i\})$$

p-value for β_1 is 0.01: "strong evidence" that $\beta_1 \neq 0$.

| Response | | | | Predictor AGE | | | | |
|----------|------------------|-------------------|-------------|------------------|-------------|-----------|------------|--|
| | INC | CHILD | DEG | | PCHILD | PINC | PDEG | |
| INC | NA | 1.10 (.11) | 7.03 (<.01) | .34 (<.01) | 4.07 (<.01) | .28 (.41) | 1.40 (.12) | |
| CHILD | .01 (.01) | NA | -.07 (.06) | .04 (<.01) | -.06 (.20) | .02 (.08) | -.05 (.20) | |

Multivariate normal copula model

This idea motivates the following “latent variable” model:

$$\begin{aligned}(z_1, \dots, z_p) &\sim \text{mvn}(\mathbf{0}, \Sigma) \\ (y_1, \dots, y_p) &= (F_1^{-1}(z_1), \dots, F_p^{-1}(z_p))\end{aligned}$$

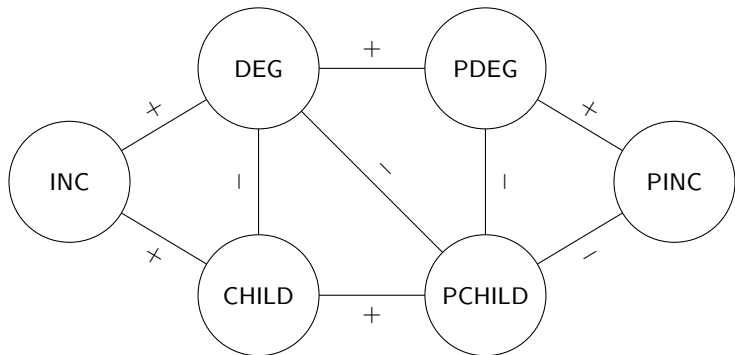
Σ parameterizes the dependence, $F_1^{-1}, \dots, F_p^{-1}$ the marginal distributions.

This model

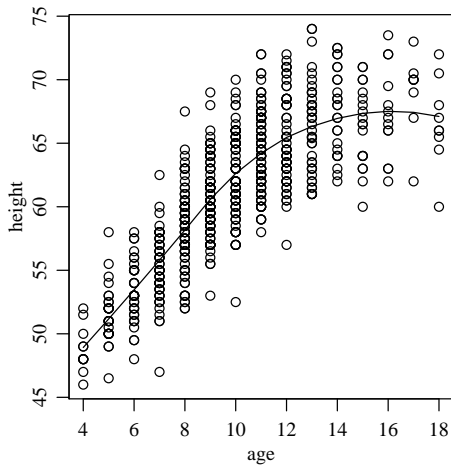
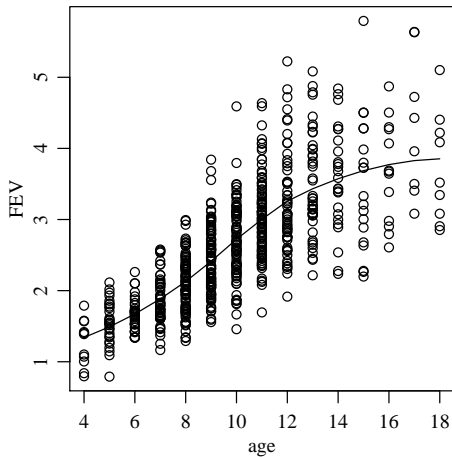
- is scale free
- is appropriate for discrete and continuous data
- gives compatible full conditional distributions

$$E[z_j | z_{-j}] = \Sigma_{[j, -j]} \Sigma_{[-j, -j]}^{-1} z_{-j}$$

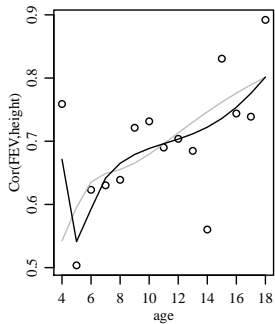
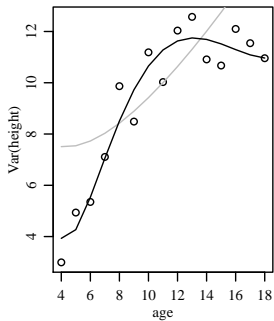
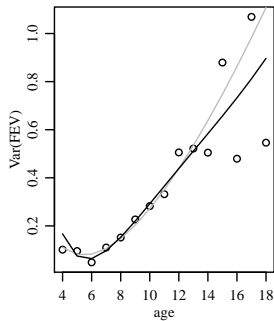
Compatible regression coefficients



FEV data

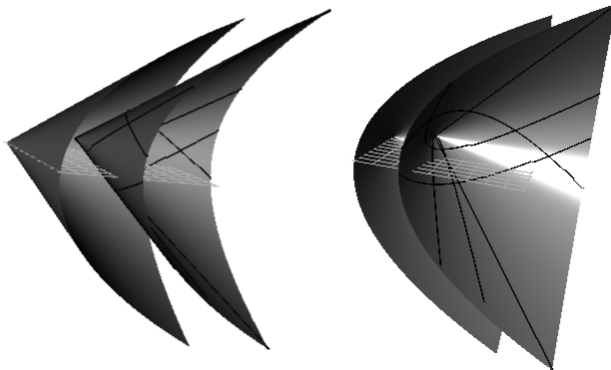


FEV data

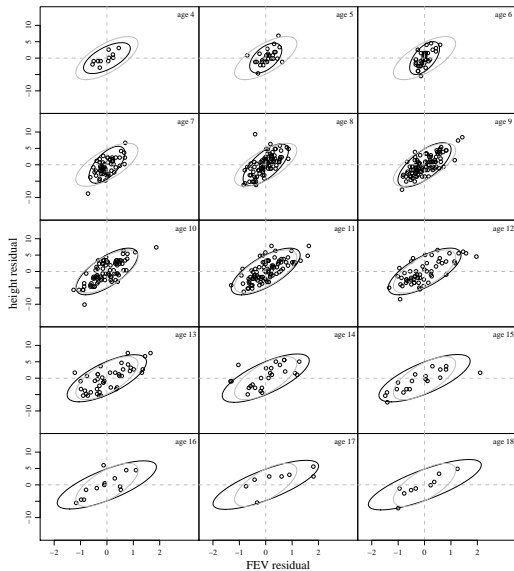


A covariance regression model

$$\Sigma_x = \mathbf{A} + \mathbf{B}\mathbf{x}\mathbf{x}^T\mathbf{B}^T \quad (1)$$



FEV data by age

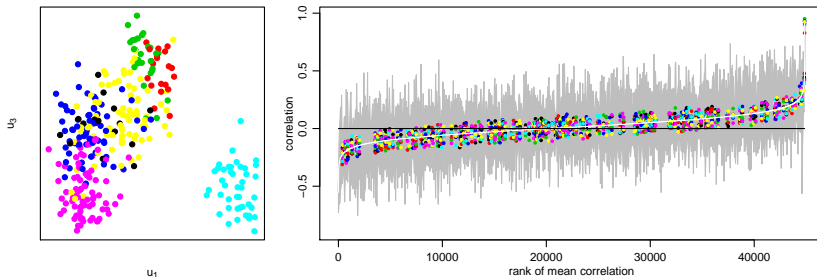


Leukemia data

Gene expression data on 327 cancer patients, each in one of seven groups:

| group | BCR | E2A | Hyperdip50 | MLL | T | TEL | other |
|-------------|-----|-----|------------|-----|----|-----|-------|
| sample size | 15 | 27 | 64 | 20 | 43 | 79 | 79 |

We look at the 300 genes with highest rank variation across subjects.



$\mathbf{Y} = \mathbf{U}\mathbf{D}\mathbf{V}^T$ Left-singular vectors of \mathbf{U} separate the groups.

$\mathbf{Y}_k = \mathbf{U}_k\mathbf{D}_k\mathbf{V}_k^T$ How do covariances $\mathbf{V}_k\mathbf{D}_k^2\mathbf{V}_k^T$ differ across groups?

Reduced rank matrix approximation

For high-dimensional data, low rank approximations are useful for describing the main patterns in row/column variability:

Symmetric matrices: $\mathbf{Y} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T + \mathbf{E}$, $y_{i,j} = \mathbf{u}_i^T \mathbf{\Lambda} \mathbf{u}_j + e_{i,j}$

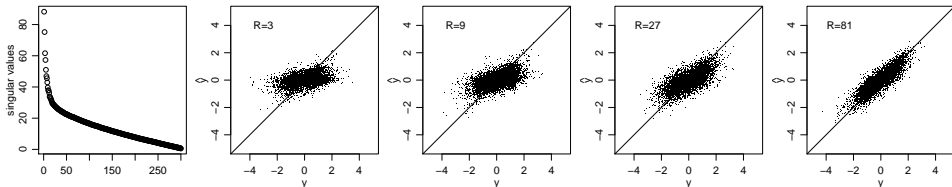
Rectangular matrices: $\mathbf{Y} = \mathbf{U}\mathbf{D}\mathbf{V}^T + \mathbf{E}$, $y_{i,j} = \mathbf{u}_i^T \mathbf{D} \mathbf{v}_j + e_{i,j}$

The column dimension R of \mathbf{U} is generally much smaller than that of \mathbf{Y} ,

$$R \ll \min(m, n)$$

so that $\mathbf{U}\mathbf{U}^T$, $\mathbf{U}\mathbf{D}\mathbf{V}^T$ provide low-rank approximations to \mathbf{Y} .

$$\min_{\mathbf{M}: \text{rank}(\mathbf{M})=R} \|\mathbf{Y} - \mathbf{M}\|^2 = \|\mathbf{Y} - \hat{\mathbf{U}}_{[:,1:R]} \hat{\mathbf{D}}_{[1:R,1:R]} \hat{\mathbf{V}}_{[:,1:R]}^T\|^2$$



A hierarchical eigenmodel

$$\begin{aligned}\mathbf{Y}_1 &= \mathbf{U}_1 \mathbf{D}_1 \mathbf{V}_1^T + \mathbf{E}_1 \\ &\vdots \quad \vdots \quad \vdots \\ \mathbf{Y}_K &= \mathbf{U}_K \mathbf{D}_K \mathbf{V}_K^T + \mathbf{E}_K\end{aligned}$$

$$\begin{aligned}\mathbf{U}_1 &\sim \text{uniform}(\mathcal{V}_{n_1, R}) & \text{diag}(\mathbf{D}_1) &\sim \text{normal}(\mathbf{0}, \tau^2 I) & \mathbf{V}_1 &\sim \text{Bingham}(\mathbf{A}, \mathbf{B}, \mathbf{V}) \\ & & \vdots & & \\ \mathbf{U}_K &\sim \text{uniform}(\mathcal{V}_{n_K, R}) & \text{diag}(\mathbf{D}_K) &\sim \text{normal}(\mathbf{0}, \tau^2 I) & \mathbf{V}_K &\sim \text{Bingham}(\mathbf{A}, \mathbf{B}, \mathbf{V})\end{aligned}$$

$$\mathbf{V} \sim \text{uniform}(\mathcal{O}_R)$$