

Generalizing the Matrix Normal Distribution

– An application to spatio-temporal data

Ekvall, Karl Oskar [University of Minnesota] and Gray, Brian [US Geological Survey]
Contact: ekvall@umn.edu

Goal

Develop model suitable for estimating trends in and fitting splines to time series of environmental and ecological outcomes from multiple spatial locations that takes into account both spatial and temporal dependence.

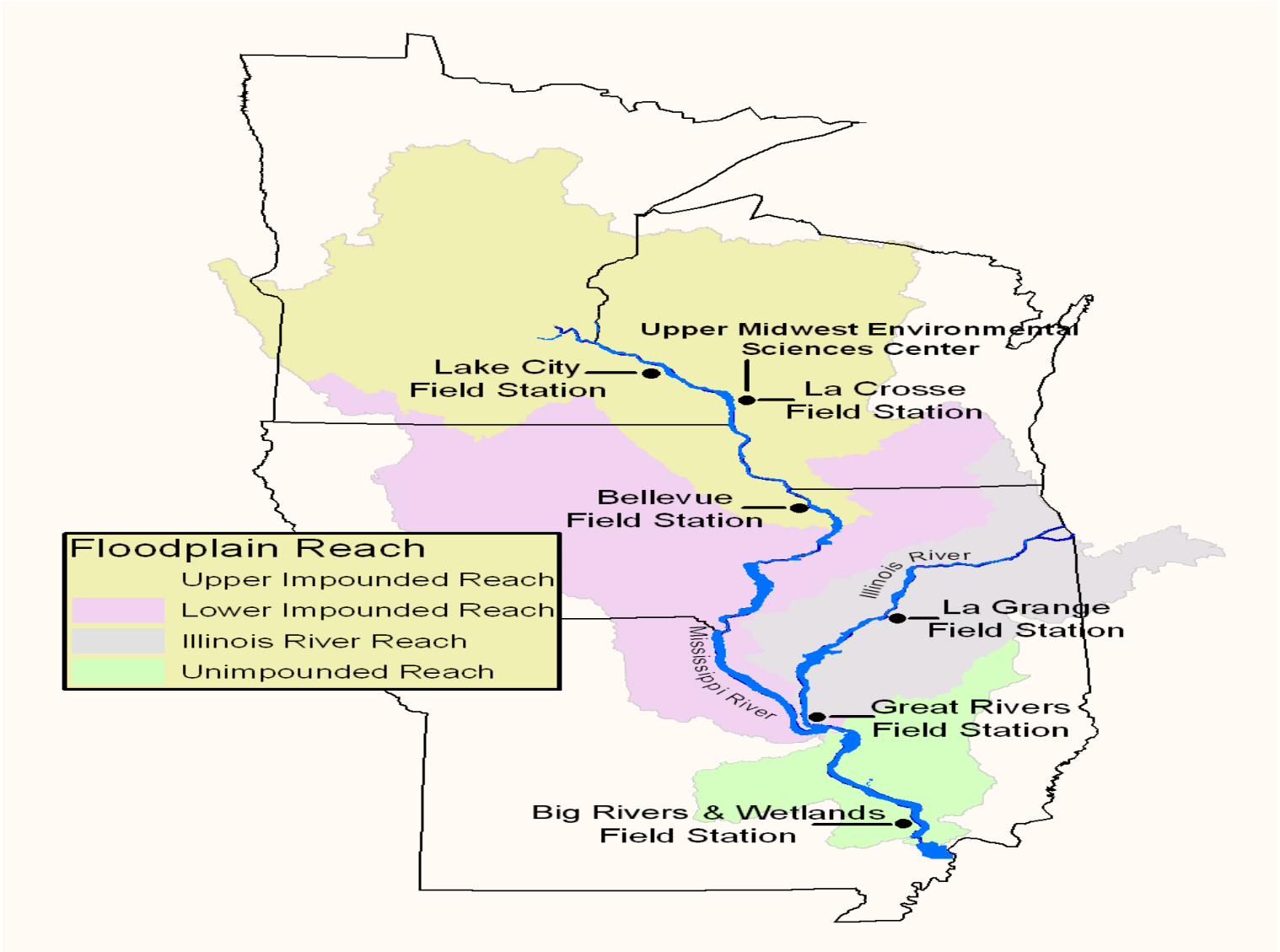
Data

- Water temperature measurements from ~ 20 locations on the Mississippi river
- Sampled quarterly for 20 years
- Each response component is the sample mean of $\sim 5 - 150$ measurements

Spatial Structure

Locations represent an incomplete crossing of reaches and sampling strata.

Reaches on the Mississippi [Image source: www.umesc.usgs.gov]



Temporal Structure

- May assume data from different years are independent
- There may or may not be dependence between data from different seasons within the same year

$$\mathbf{Y}_t = \begin{pmatrix} y_{1,1} & y_{1,2} & y_{1,3} \\ \vdots & \vdots & \vdots \\ y_{18,1} & y_{18,2} & y_{18,3} \end{pmatrix}, t = 1, \dots, n$$

A Data Excerpt

```
## # A tibble: 1,080 x 5
##   location year season   temp     n
##   <fctr> <int> <fctr>   <dbl> <int>
## 1      1:1 1994     SP 10.12000 25
## 2      1:2 1994     SP 10.62333 30
## 3      1:3 1994     SP 12.32600 50
## 4      1:4 1994     SP 11.08333 30
## 5      2:1 1994     SP 12.18000 25
## 6      2:2 1994     SP 11.82333 30
## 7      2:3 1994     SP 12.95500 60
## 8      2:5 1994     SP 11.69600 25
## 9      2:6 1994     SP 15.41000 10
## 10     3:1 1994     SP 15.85000 30
## # ... with 1,070 more rows
```

The Matrix Normal Distribution

$$\mathbf{y}_t := \text{vec}(\mathbf{Y}_t) \sim \text{N}_{rc}(\text{vec}(\mathbf{X}_t \boldsymbol{\beta}), \boldsymbol{\Sigma})$$

$$\boldsymbol{\Sigma} = \mathbf{V} \otimes \mathbf{U}$$

- Restriction: $\text{cov}(y_{i,j}, y_{i',j'}) = \mathbf{U}_{i,i'} \mathbf{V}_{j,j'}$ [separability of covariance]
- Gain: $\mathcal{O}(r^2 + c^2)$ covariance parameters instead of $\mathcal{O}(r^2 c^2)$
- Note: Independence would imply $\boldsymbol{\Sigma} = \text{diag}(\sigma_1^2, \dots, \sigma_{rc}^2) \in \mathcal{O}(rc)$

Theory: MLE of $\boldsymbol{\Sigma}$ exists and is unique when $n > r/c + c/r + p$, which is proved in [1]. The key observation is that the negative log-likelihood is geodesically convex in (\mathbf{V}, \mathbf{U}) with the geodesics $d_t(\mathbf{A}, \mathbf{B}) := \mathbf{A}^{1/2}(\mathbf{A}^{-1/2} \mathbf{B} \mathbf{A}^{-1/2})^t \mathbf{A}^{1/2}$

Our Generalization

$$\mathbf{y}_t := \text{vec}(\mathbf{Y}_t) \sim \text{N}_{rc}(\text{vec}(\mathbf{X}_t \boldsymbol{\beta}), \boldsymbol{\Sigma})$$

$$\boldsymbol{\Sigma} = \mathbf{C}(\mathbf{A} \otimes \mathbf{B})\mathbf{C},$$

- where $\mathbf{C} = \text{diag}(\sigma_1, \dots, \sigma_{rc})$.
- Restriction: $\text{corr}(y_{i,j}, y_{i',j'}) = \mathbf{B}_{i,i'} \mathbf{A}_{j,j'}$ [separability of correlation]
 - Gain: Completely general variance structure
 - Still with $\mathcal{O}(r^2 + c^2) \ni \mathcal{O}(r^2 + c^2 + rc)$ parameters

Theory: The negative log-likelihood is convex in each of \mathbf{V} , \mathbf{U} and $\mathbf{C}_{ii}, i = 1, \dots, rc$, but not jointly. We do not know whether the negative log-likelihood is geodesically convex under some appropriate choice of geodesics.

Optimization

We propose a blockwise coordinate descent algorithm for maximizing the likelihood.

- In each update, we set an unrestricted gradient to zero and then rescale variables to satisfy the constraint
- All updates are closed form

For example, treating \mathbf{A} as a *general* matrix: update \mathbf{A}^{k+1} to the solution of $\nabla_{\mathbf{A}} \ell(\mathbf{A}, \mathbf{B}^k, \mathbf{C}^k) = 0$. Then, rescale:

- $\mathbf{C}^k \leftarrow \mathbf{C}^k (\text{diag}(\mathbf{A}^{k+1})^{1/2} \otimes \mathbf{I})$
- $\mathbf{A}^{k+1} \leftarrow \text{diag}(\mathbf{A}^{k+1})^{-1/2} \mathbf{A}^{k+1} \text{diag}(\mathbf{A}^{k+1})^{-1/2}$

Rescaling after updating both \mathbf{A} and \mathbf{B} and then updating \mathbf{C} component by component gives the following algorithm.

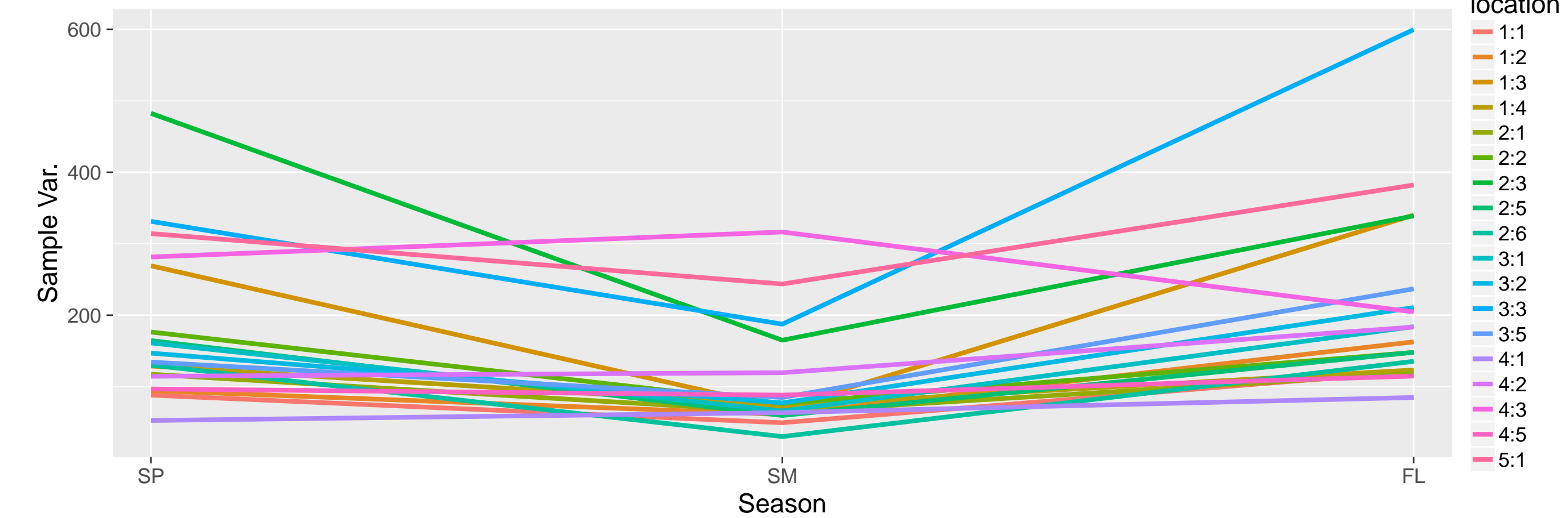
Algorithm: Blockwise coordinate descent

- 1: Initialize $\mathbf{A}^0, \mathbf{B}^0, \mathbf{C}^0, k = 0$
- 2: **repeat**
- 3: Set \mathbf{A}^{k+1} to the solution of $\nabla_{\mathbf{A}} \ell(\mathbf{A}, \mathbf{B}^k, \mathbf{C}^k) = 0$
- 4: Set \mathbf{B}^{k+1} to the solution of $\nabla_{\mathbf{B}} \ell(\mathbf{A}^{k+1}, \mathbf{B}, \mathbf{C}^k) = 0$
- 5: Rescale $\mathbf{A}^{k+1}, \mathbf{B}^{k+1}$ and \mathbf{C}^k to satisfy constraints
- 6: **for** $j = 1, \dots, rc$ **do**
- 7: Set θ_j to the solution of $\nabla_{\theta_j} \ell(\mathbf{A}^{k+1}, \mathbf{B}^{k+1}, \theta_1^{k+1}, \dots, \theta_{j-1}^{k+1}, \theta_j, \theta_{j+1}^k, \dots, \theta_{rc}^k) = 0$
- 8: **end for**
- 9: $k \leftarrow k + 1$
- 10: **until** $|\ell^k - \ell^{k-1}| \leq \epsilon$

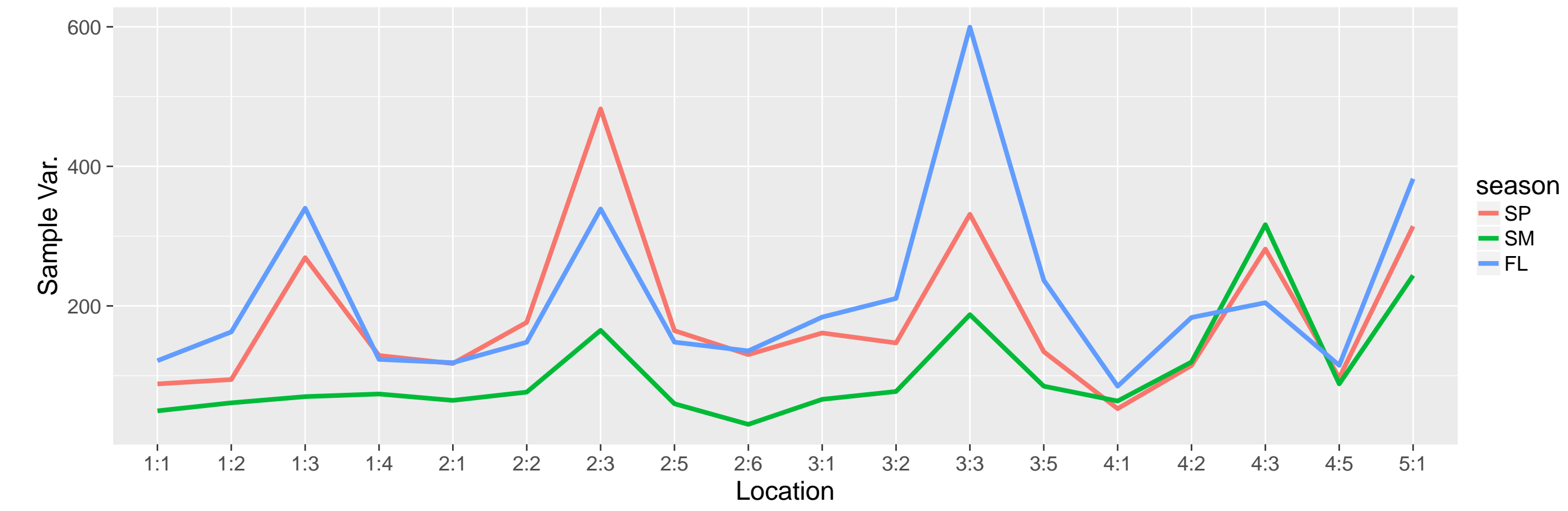
Notation: ℓ^k is the log-likelihood evaluated at the k th iterate of all optimization variables. ∇_Z denotes gradient w.r.t. variable Z . ϵ is a tolerance parameter for terminating the algorithm. We parameterize $\mathbf{C}_{i,i} = 1/\theta_i$.

Results

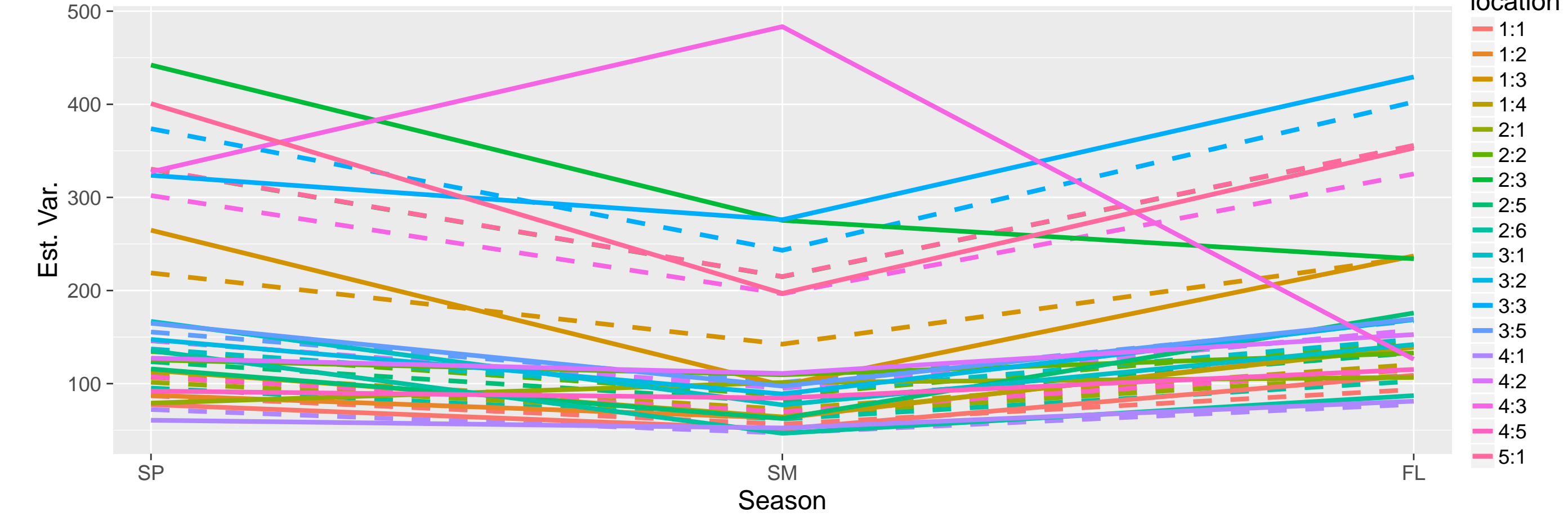
Sample Variance v. Season



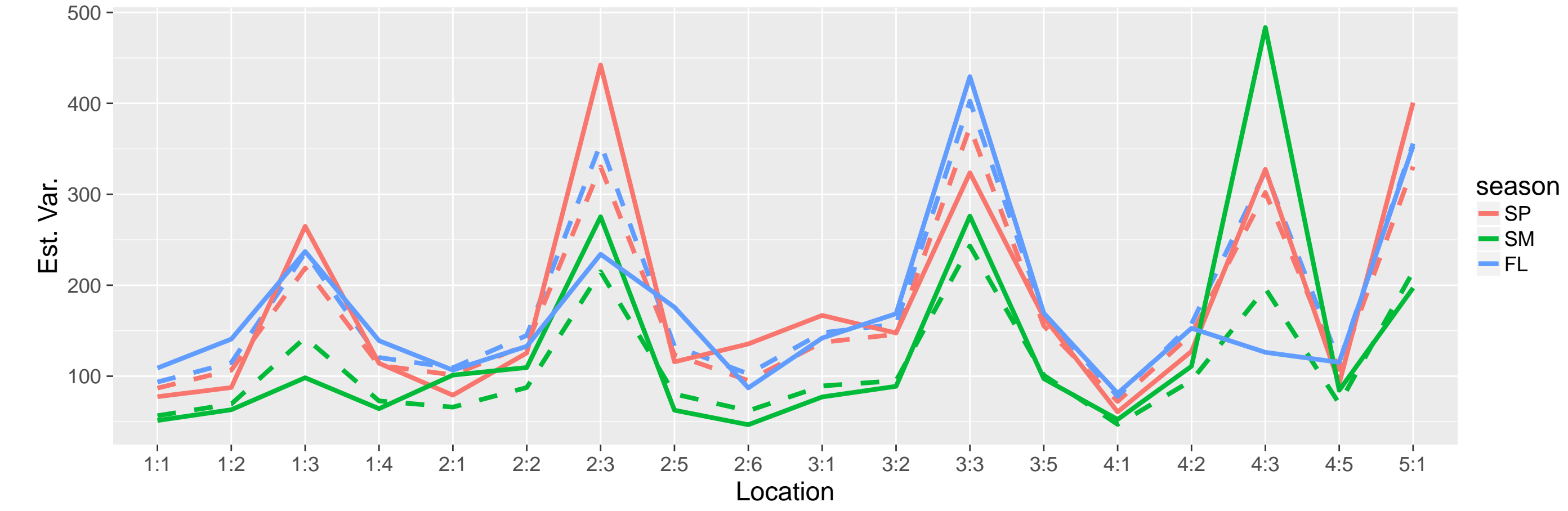
Sample Variance v. Location



Estimated Var v. Season: Solid = Generalization, Dashed = Matrix Normal



Estimated Var v. Location: Solid = Generalization, Dashed = Matrix Normal



References

[1] I. Solovveychik and D. Trushin. Gaussian and robust Kronecker product covariance estimation: Existence and uniqueness. *Journal of Multivariate Analysis*, 149:92 – 113, 2016.