THE UNIVERSITY OF CHICAGO

Introduction

Climate science has recently and rapidly become data-rich [1], due to a combination of

- New satellites and ground-based sensors collecting and storing massive amounts of high-resolution climate and ecosystem data
- Increased computational power to generate and archive physics-based climate model simulations

Advances in statistics and machine learning have only recently been applied to the Earth sciences, due to

- Lack of labeled training data
- Small sample sizes and high-dimensional settings
- Strong spatiotemporal dependencies among features

Seasonal Forecasting

Goal: Predict winter precipitation across the southwestern US using summer sea-surface temperatures over the Pacific





Linear model: For year *i* and climate division *r*

- $\hat{y}_r^{(i)} = X_1^{(i)} \hat{\beta}_1 + \dots + X_p^{(i)} \hat{\beta}_p$
- $y_r^{(i)}$: winter precipitation
- $X_i^{(\prime)}$: sea-surface temperatures at (location, time) j
- $y \in \mathbb{R}^n$, $X \in \mathbb{R}^{n \times p}$, $p \gg n$

Graph Total Variation

GTV: Estimate sparse coefficients so that the higher the correlation of X_i and X_k (σ_{ik}) the more similar $\hat{\beta}_i$ and $\hat{\beta}_k$ [2].



Affiliations

¹Department of Statistics, University of Chicago, ² Department of CS, University of Chicago, ³Department of Civil and Environmental Engineering, University of California, Irvine, ⁴Department of Earth System Science, University of California, Irvine, ⁵Department of CS, University of California, Irvine, ⁶CS Department, University of Wisconsin, Madison

Graph-guided regularization for SWUS seasonal forecasting

Abby Stevens¹, Rebecca Willett^{1,2}, Antonios Mamalakis³, Efi Foufoula-Georgiou^{3,4}, James Randerson⁴, Padhraic Smyth⁵, and Stephen Wright⁶

Large Ensemble Climate Simulations



www.climate.gov

CESM Large Ensemble Project (LENS)

- 40-member ensemble of physics-based climate conditions
- Combining observations and simulations shows potential for improving performance of ML methods in climate science
- LENS simulations are leveraged to estimate Σ , the covariance matrix of X. The simulated SSTs are interpolated onto the same spatial grid as SST observations and stacked to form a matrix $X_L \in \mathbb{R}^{40n \times p}$. Let $\hat{\Sigma}_L$ be the sample covariance of X_L . $\hat{\Sigma}_L$ is used to form the covariance graph for the GTV estimator.

MultiGTV

Responses from *m* different regions across the southwestern US. We assume that

- The subset of relevant features is preserved across regions
- Coefficients in each region are aligned with the covariance graph

MultiGTV: Extend GTV to multi-task setting $\cdot \parallel \mathbf{V} \quad \mathbf{V} \mathbf{D} \parallel^2$ Ô

$$= \arg \min ||Y - XB||_{F}^{-}$$

$$+ \lambda_{1} \sum_{r=1}^{B} \sum_{j,k \in E} \sigma_{j,k} |\beta_{j}^{(r)} - \beta_{k}^{(r)}|$$

$$+ \lambda_{2} \sum_{i=1}^{p} ||\beta_{i}^{(i)}||_{2}$$

$$V = [V_{i}, V_{i}] = R = R$$

Covariance Graph

- Assume $X^{(i)} \sim N(0, \Sigma)$, $X \in \mathbb{R}^{n \times p}$. Define graph G = (V, E, W)• $V = \{X_1, ..., X_p\}, \quad W_{jk} = \sigma_{jk}$
- $(j,k) \in E \iff |\Sigma_{j,k}| > 0$
- $p \gg n \implies$ sample covariance not consistent estimator for Σ



simulations, each subject to perturbed initial

Data fit

Within-region GTV

Group sparsity

 $Y = [y_1, y_2, \dots, y_m]$ $B = |\beta^{(1)}, \beta^{(2)}, \dots, \beta^{(m)}|$

Experiments

Data: Precipitation and sea-surface temperatures from 1940-2014 ($y \in \mathbb{R}^{75}$, $X \in \mathbb{R}^{75 imes 904}$)

- on the held out last 25 years of data

Predictive performance



Model selection

- October for different methods
- aggregate estimates across regions



References

- CRC Press, 2017, ch. 2, pp. 13–32.

• Precipitation over 16 climate divisions across California, Nevada, Utah, and Arizona, aggregated over November-March • Sea-surface temperatures aggregated over $10^{\circ} \times 10^{\circ}$ regions across the Pacific in July, August, September, and October • Models trained on the first 50 years of data and errors reported

• Markers indicate location and relative magnitude of selected coefficients in

• The area-weighted average precipitation across all SWUS regions is used as the response for Lasso, Fused Lasso, and GTV, and for MultiGTV we

[1] A. R. Goncalves, A. Banerjee, V. Sivakumar, , and S. Chatterjee, "Structured estimation in high dimensions: applications in climate," in *Large-scale machine learning in the earth sciences*.

[2] Y. Li, B. Mark, G. Rasutti, and R. Willett, "Graph-based regularization for regression problems with highly-correlated designs," arXiv preprint arXiv:1410.5093, 2018.