Probabilistic Graphical Models for Climate Data Analysis

Arindam Banerjee <u>banerjee@cs.umn.edu</u>

Dept of Computer Science & Engineering University of Minnesota, Twin Cities

Aug 15, 2013

Climate Data Analysis

- Key Challenges
 - High-dimensional dependent data, small sample size
 - Spatial and temporal dependencies, temporal lags
 - Oscillations with frequency and phase variations
 - Important variables are unreliable, e.g., precipitation
 - Several others: Nonlinearity, heavy tails, ...
- Potential portunities
 - Multi-model ensembles: Regional skills vs global performance
 - Statistical Downscaling: Coarse to fine scale, capture dependencies
 - Understanding tails: Extreme precipitation, mega-droughts, heat waves, etc.
 - Understanding dependencies: Statistical dependencies, not correlation



Source: Overpeck et al., *Science*, (2011)

Graphical Models

- Graphical models
 - Dependencies between (random) variables, avoid I.I.D. assumptions
 - Closer to reality, learning/inference is much more difficult
- Basic nomenclature
 - Node = Random Variable, Edge = Statistical Dependency
- Directed Graphs
 - A *directed* graph between random variables
 - Example: Bayesian networks, Hidden Markov Models
 - Joint distribution is a product of P(child|parents)
- Undirected Graphs
 - An *undirected* graph between random variables
 - Example: Markov/Conditional random fields
 - Joint distribution in terms of potential functions



Graphical Models: Key Problems

- Structure Learning
 - Given: Samples
 - Problem: Learn the Structure
- Parameter Estimation
 - Given: Samples and Structure
 - Problem: Estimate Parameters
- P(E)P(B)Burglary Earthquake .001 .002 P(A|B,E)В Ε Т Т .95 Alarm Т F .94 F Т .29 F F .001 P(J|A) А A P(M|A) Т JohnCalls .90 MaryCalls Т .70 F .05 F .01

- Inference
 - Given: Structure, Parameters, and some variables (part of a Sample)
 - Problem: Find other variables (part of a Sample)

Global Climate Models (GCMs)



Atmospheric GCM

Atmospheric Model

Ocean Model

Layers

Layers

Ocean OCM

Combining GCM Outputs

- Several ways to combine the model outputs
 - Average: Equal weightage to all models (IPCC AR4 2007, Reifen and Toumi 2009)
 - Superensemble: Least Squares (Krishnamurti et al., 2002)
 - REA: Reliability based ensemble averaging (Giorgi et al., 2002)
 - Bayesian: Probabilistic estimates of climate variables (Tebaldi et al., 2005, Smith et al., 2011)
 - Online Learning: Tracking climate models (Monteleoni et al., 2011)
- Our work
 - Hypothesis: Certain models do well in certain climatic conditions
 - Goal: Climate model combination
 - Different weighs at different locations
 - Similar climatic conditions should get similar weights
 - Builds on superensemble and probabilistic approaches

Smooth Model Combination (SMC)



SMC: Error Term



SMC: Smoothness Term



SMC: Graphical Model Perspective

Generative Model

Prior on rows: $\theta_j \sim N(0, A^{-1})$ Conditional: $y_i \sim N(X_i \ \theta_i, \sigma^2)$

- Precision matrix specification
 - Gaussian Markov random field (GMRF)
 - Precision A = L = D W, the discrete graph Laplacian
 - Intrinsic Conditionally Autoregressive Model (ICAR)
 - Spatial statistics literature (Diggle et al., 1998, Besag et al., 1995, Banerjee et al., 2004, Rue et al., 2005)
- Estimation of precision matrix
 - Estimate which locations are 'similar'
 - Estimated precision is full rank but sparse



Data Set and Methodology

GCM	Origin
BCCR_BCM2	Bjerknes Centre for Climate Research, Norway
CCMA_CGCM3	Canadian Centre for Climate Modelling and Analysis, Canada
MICRO3-2-Hires	Center for Climate System Research, Univ. of Tokyo, Japan
CNRM_CM3	Center for National Weather Research, France
GFDL_CM2_1	Geophysical Fluid Dynamics Laboratory, USA
GISS_E_H	Goddard Institute for Space Studies, USA
INGV_ECHAM4	European Center for Medium-Range Weather Forecasts, UK
IPSL	Institut Pierre Simon Laplace, France

- GCM output: Monthly average surface temperature
- Target variable: Temperature from Climatic Research Unit (CRU)
- Error/accuracy measures: RMSE and MAE
- Smoothness measures:
 - Kendall τ
 - Spearman ρ



Spatial Error Profile: AVE vs SMC



- SMC has lower compared to AVE: lower by $\sim 0.5^{\circ}$ C
- Errors (visibly) reduced in many regions
 - Africa, Greenland, Southeast Asia, Siberia
- High errors in some regions
 - Northern Europe/Russia, China/Tibet, West South America, North America

Error vs Smoothness



13

Mega-Droughts

- Mega-Droughts
 - Persistent over space and time
 - Catastrophic consequences
- Examples
 - Late 1906s Sahel drought
 - 1930s North American Dust Bowl
- Discrete Markov Random Field (MRF)
 - Each node x_i is "wet" or "dry"
 - Observations: Precipitation
 - Smoothness in space and time
 - Most likely state assignments
 - Each (lat,long,time) gets "wet" or "dry"
 - Advanced analysis
 - Soil moisture, hydrology/watershed models
 - Multiple states based on severity, e.g., lower quantiles





Results: Droughts starting in 1920-30s



Results: Droughts starting in 1960-70s



Major Droughts: 1901-2006



Learning dependencies



Sparse Regression, Structure Learning



Gaussian Copula: $(f_1(x_1), \dots, f_i(x_i), \dots, f_p(x_p)) \sim N(0, \Sigma)$, precision $A = \Sigma^{-1}$

- <u>Not</u> $(x_1,...,x_i,...,x_p) \sim N(0, \Sigma)$, f_i are monotonic transformations
- Sparse A can be consistently estimated (H. Liu et al., 2012)
- Linear Programming (LP) based estimator (CLIME) (T. Cai et al., 2011)
- LP estimator scales to high-dimensional copulas (Our work)
 - Millions of variables, trillions of edges

References

- K. Subbian and A. Banerjee, Climate Multi-model Regression Using Spatial Smoothing, *SIAM Data Mining (SDM)*, 2013. [Best Application Paper Award]
- Q. Fu, A. Banerjee, S. Liess, and P. Snyder. Drought Detection for the Last Century: A MRF-based Approach, *SIAM Data Mining (SDM)*, 2012.
- Q. Fu, H. Wang, and A. Banerjee, Bethe-ADMM for Tree Decomposition based Parallel MAP inference, *Uncertainty in Artificial Intelligence (UAI)*, 2013.
- C. Jin, Q. Fu, H. Wang, A. Agrawal, W. Hendrix, W.-K. Liao, M. A. Patwary, A. Banerjee, and A. Choudhary, Solving Combinatorial Optimization Problems using Relaxed Linear Programming: A High Performance Computing Perspective, *BigMine workshop (KDD)*, 2013. [Best Paper Award]
- C. Hsieh, I. Dhillon, P. Ravikumar, A. Banerjee, A Divide-and-Conquer Method for Sparse Inverse Covariance Estimation, *Advances in Neural Information Processing Systems (NIPS)*, 2012.
- S. Chatterjee, K. Steinhaeuser, A. Banerjee, S. Chatterjee, and A. Ganguly, Sparse Group Lasso: Consistency and Climate Applications, *SIAM Data Mining (SDM)*, 2012. [Best Student Paper Award]