

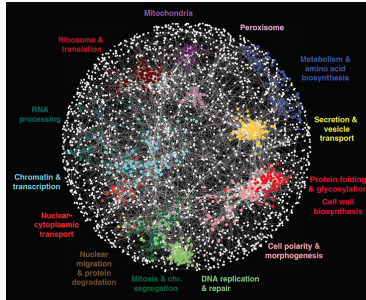
Computational and Statistical Methodology for Highly-Structured Data

Ph.D.Thesis Defense: 2020-09-10

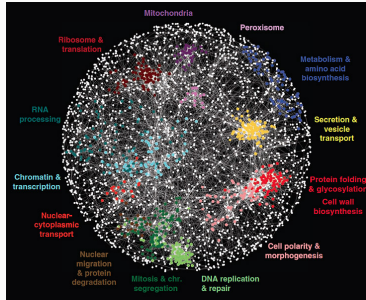
Michael Weylandt

Slides Available Online at <https://tinyurl.com/WeylandtDefense>

Rice University, Department of Statistics

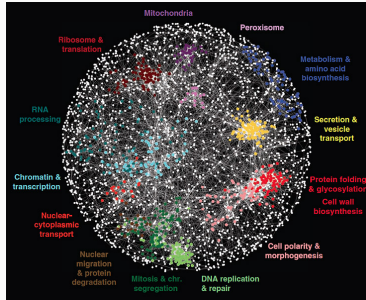


“Big Data” enables “Big Models”



“Big Data” enables “Big Models”

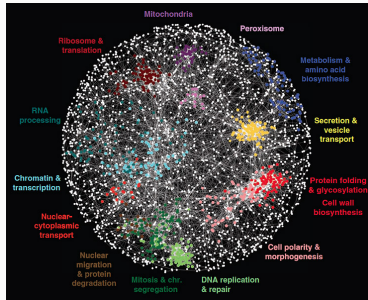
“Big Data” is never IID



“Big Data” enables “Big Models”

“Big Data” is never IID

“Big Data” **requires** “Big Models”

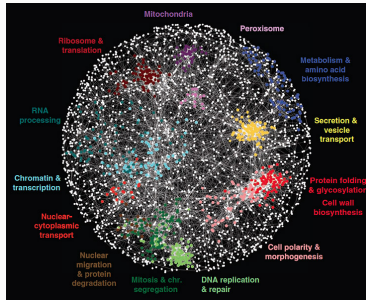


“Big Data” enables “Big Models”

“Big Data” is never IID

“Big Data” **requires** “Big Models”

Highly-structured data requires flexible but powerful models to reflect and capture dependencies in data



“Big Data” enables “Big Models”

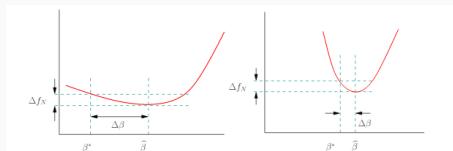
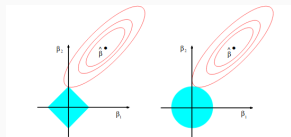
“Big Data” is never IID

“Big Data” **requires** “Big Models”

Highly-structured data requires flexible but powerful models to reflect and capture dependencies in data

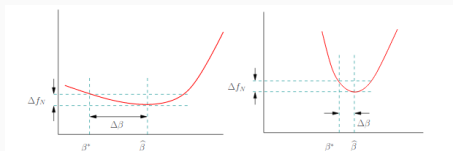
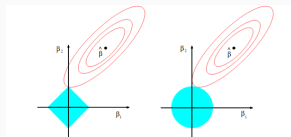
Big data allows us to fit such models

Convex Revolution in Statistical Machine Learning



Biggest advance in 21st c. statistics – convex analysis and optimization:

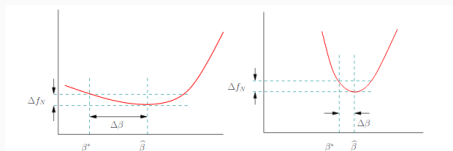
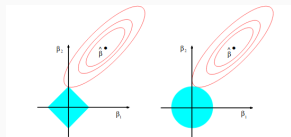
Convex Revolution in Statistical Machine Learning



Biggest advance in 21st c. statistics – convex analysis and optimization:

- Development of novel regularized estimation schemes

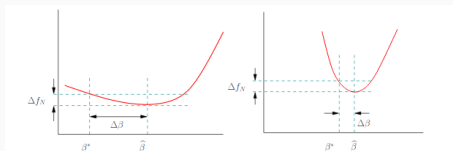
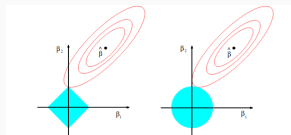
Convex Revolution in Statistical Machine Learning



Biggest advance in 21st c. statistics – convex analysis and optimization:

- Development of novel regularized estimation schemes
- Algorithms that efficiently scale to enormous data sets

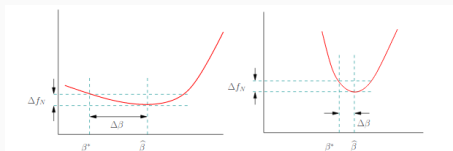
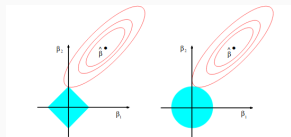
Convex Revolution in Statistical Machine Learning



Biggest advance in 21st c. statistics – convex analysis and optimization:

- Development of novel regularized estimation schemes
- Algorithms that efficiently scale to enormous data sets
- Theoretical advances based on powerful convex analysis

Convex Revolution in Statistical Machine Learning



Biggest advance in 21st c. statistics – convex analysis and optimization:

- Development of novel regularized estimation schemes
- Algorithms that efficiently scale to enormous data sets
- Theoretical advances based on powerful convex analysis

**Develop Methodology for Big Highly-Structured Data
Built on Powerful Convex Analysis and Optimization**

Agenda

Splitting Methods for Clustering

Multi-Rank Regularized PCA

Multivariate Models for Gas Markets

Complex Convex Analysis

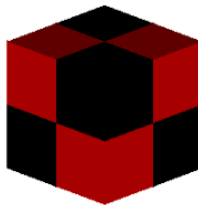
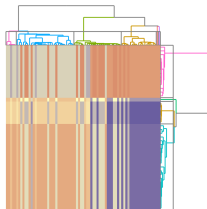
Conclusion & Discussion

Splitting Methods for Clustering

Tensor Co-Clustering

Co-Clustering:

- Simultaneous clustering along all faces of a tensor
- Discover “checkerboard” patterns in data
- “Cluster Heatmap” for 2-tensors
- Manifold learning for K -tensors (Mishne *et al.*, 2019)



Convex Bi-Clustering

Convex formulation of co-clustering: Chi *et al.* (2017) and Chi *et al.* (2018)

- Frobenius norm loss \implies approximate observed data
- Convex fusion penalty \implies encourages clustering

Convex Bi-Clustering

Convex formulation of co-clustering: Chi *et al.* (2017) and Chi *et al.* (2018)

- Frobenius norm loss \implies approximate observed data
- Convex fusion penalty \implies encourages clustering

Matrix (2-tensor) case:

$$\hat{\mathbf{U}} = \arg \min_{\mathbf{U} \in \mathbb{R}^{n \times p}} \frac{1}{2} \|\mathbf{X} - \mathbf{U}\|_F^2 + \lambda \left(\sum_{\substack{i,j=1 \\ i \neq j}}^n w_{ij} \|\mathbf{U}_i - \mathbf{U}_j\|_q + \sum_{\substack{k,l=1 \\ k \neq l}}^p \tilde{w}_{kl} \|\mathbf{U}_k - \mathbf{U}_l\|_q \right)$$

Convex Bi-Clustering

Convex formulation of co-clustering: Chi *et al.* (2017) and Chi *et al.* (2018)

- Frobenius norm loss \implies approximate observed data
- Convex fusion penalty \implies encourages clustering

Matrix (2-tensor) case:

$$\hat{\mathbf{U}} = \arg \min_{\mathbf{U} \in \mathbb{R}^{n \times p}} \frac{1}{2} \|\mathbf{X} - \mathbf{U}\|_F^2 + \lambda \left(\sum_{\substack{i,j=1 \\ i \neq j}}^n w_{ij} \|\mathbf{U}_{\cdot i} - \mathbf{U}_{\cdot j}\|_q + \sum_{\substack{k,l=1 \\ k \neq l}}^p \tilde{w}_{kl} \|\mathbf{U}_{\cdot k} - \mathbf{U}_{\cdot l}\|_q \right)$$

Simultaneous clustering of rows and columns:

- Rows are clustered together if $\hat{\mathbf{U}}_{\cdot i} = \hat{\mathbf{U}}_{\cdot j}$.
- Columns are clustered together if $\hat{\mathbf{U}}_{\cdot k} = \hat{\mathbf{U}}_{\cdot l}$.
- Each element of \mathbf{X} is assigned to a single bi-cluster

Convex Bi-Clustering

Convex formulation of co-clustering: Chi *et al.* (2017) and Chi *et al.* (2018)

- Frobenius norm loss \implies approximate observed data
- Convex fusion penalty \implies encourages clustering

Matrix (2-tensor) case:

$$\hat{\mathbf{U}} = \arg \min_{\mathbf{U} \in \mathbb{R}^{n \times p}} \frac{1}{2} \|\mathbf{X} - \mathbf{U}\|_F^2 + \lambda \left(\sum_{\substack{i,j=1 \\ i \neq j}}^n w_{ij} \|\mathbf{U}_{\cdot i} - \mathbf{U}_{\cdot j}\|_q + \sum_{\substack{k,l=1 \\ k \neq l}}^p \tilde{w}_{kl} \|\mathbf{U}_{\cdot k} - \mathbf{U}_{\cdot l}\|_q \right)$$

Simultaneous clustering of **rows** and **columns**:

- Rows are clustered together if $\hat{\mathbf{U}}_{\cdot i} = \hat{\mathbf{U}}_{\cdot j}$.
- Columns are clustered together if $\hat{\mathbf{U}}_{\cdot k} = \hat{\mathbf{U}}_{\cdot l}$.
- Each element of \mathbf{X} is assigned to a single bi-cluster

λ controls the number of co-clusters smoothly

Splitting Methods for Convex Bi-Clustering

Simplified form:

$$\hat{\mathbf{U}} = \arg \min_{\mathbf{U} \in \mathbb{R}^{n \times p}} \frac{1}{2} \|\mathbf{X} - \mathbf{U}\|_F^2 + \lambda \left(\underbrace{\|\mathbf{D}_{\text{row}} \mathbf{U}\|_{\text{row},q}}_{P_{\text{row}}(\mathbf{D}_{\text{row}} \mathbf{U})} + \underbrace{\|\mathbf{U} \mathbf{D}_{\text{col}}\|_{\text{col},q}}_{P_{\text{col}}(\mathbf{U} \mathbf{D}_{\text{col}})} \right)$$

Splitting Methods for Convex Bi-Clustering

Simplified form:

$$\hat{\mathbf{U}} = \arg \min_{\mathbf{U} \in \mathbb{R}^{n \times p}} \frac{1}{2} \|\mathbf{X} - \mathbf{U}\|_F^2 + \lambda \left(\underbrace{\|\mathbf{D}_{\text{row}} \mathbf{U}\|_{\text{row},q}}_{P_{\text{row}}(\mathbf{D}_{\text{row}} \mathbf{U})} + \underbrace{\|\mathbf{U} \mathbf{D}_{\text{col}}\|_{\text{col},q}}_{P_{\text{col}}(\mathbf{U} \mathbf{D}_{\text{col}})} \right)$$

Current state of the art:

COBRA - Dykstra-Like Proximal Algorithm

(Bauschke and Combettes, 2008; Chi and Lange, 2015)

- Alternating row- and column-wise convex clustering

Convex clustering subproblems are still slow, so COBRA doesn't scale

Splitting Methods for Convex Bi-Clustering

$$\hat{\mathbf{U}} = \arg \min_{\mathbf{U} \in \mathbb{R}^{n \times p}} \frac{1}{2} \|\mathbf{X} - \mathbf{U}\|_F^2 + \lambda \left(\underbrace{\|\mathbf{D}_{\text{row}} \mathbf{U}\|_{\text{row},q}}_{P_{\text{row}}(\mathbf{D}_{\text{row}} \mathbf{U})} + \underbrace{\|\mathbf{U} \mathbf{D}_{\text{col}}\|_{\text{col},q}}_{P_{\text{col}}(\mathbf{U} \mathbf{D}_{\text{col}})} \right)$$

Can we develop a fast splitting approach?

Splitting Methods for Convex Bi-Clustering

$$\hat{\mathbf{U}} = \arg \min_{\mathbf{U} \in \mathbb{R}^{n \times p}} \frac{1}{2} \|\mathbf{X} - \mathbf{U}\|_F^2 + \lambda \left(\underbrace{\|\mathbf{D}_{\text{row}} \mathbf{U}\|_{\text{row},q}}_{P_{\text{row}}(\mathbf{D}_{\text{row}} \mathbf{U})} + \underbrace{\|\mathbf{U} \mathbf{D}_{\text{col}}\|_{\text{col},q}}_{P_{\text{col}}(\mathbf{U} \mathbf{D}_{\text{col}})} \right)$$

Can we develop a fast splitting approach?

Davis and Yin (2017) three-block ADMM:

1. $\mathbf{U}^{(k+1)} = \mathbf{X} - \mathbf{D}_{\text{row}}^T \mathbf{Z}_{\text{row}}^{(k)} - \mathbf{Z}_{\text{col}}^{(k)} \mathbf{D}_{\text{col}}^T$
- 2(a). $\mathbf{V}_{\text{row}}^{(k+1)} = \text{prox}_{\lambda/\rho, P_{\text{row}}(\cdot)}(\mathbf{D}_{\text{row}} \mathbf{U}^{(k+1)} + \mathbf{Z}_{\text{row}}^{(k)})$
- 2(b). $\mathbf{V}_{\text{col}}^{(k+1)} = \text{prox}_{\lambda/\rho, P_{\text{col}}(\cdot)}(\mathbf{U}^{(k+1)} \mathbf{D}_{\text{col}} + \mathbf{Z}_{\text{col}}^{(k)})$
- 3(a). $\mathbf{Z}_{\text{row}}^{(k+1)} = \mathbf{Z}_{\text{row}}^{(k)} + \rho(\mathbf{D}_{\text{row}} \mathbf{U}^{(k+1)} - \mathbf{V}_{\text{row}}^{(k+1)})$
- 3(b). $\mathbf{Z}_{\text{col}}^{(k+1)} = \mathbf{Z}_{\text{col}}^{(k)} + \rho(\mathbf{U}^{(k+1)} \mathbf{D}_{\text{col}} - \mathbf{V}_{\text{col}}^{(k+1)})$

Equivalent to AMA and to prox-gradient on the dual - very slow!

(Tseng, 1991)

Splitting Methods for Convex Bi-Clustering

Why not apply ADMM directly?

Splitting Methods for Convex Bi-Clustering

Why not apply ADMM directly? Lifted problem:

$$\arg \min_{\substack{\mathbf{U} \in \mathbb{R}^{n \times p} \\ (\mathbf{V}_{\text{row}}, \mathbf{V}_{\text{col}}) \in \mathbb{R}^{\# \text{row} \times p} \times \mathbb{R}^{n \times \# \text{col}}}} \frac{1}{2} \|\mathbf{X} - \mathbf{U}\|_F^2 + \lambda (P_{\text{row}}(\mathbf{V}_{\text{row}}) + P_{\text{col}}(\mathbf{V}_{\text{col}}))$$

subject to

$$\mathfrak{L}_1 \mathbf{U} - (\mathbf{V}_{\text{row}}, \mathbf{V}_{\text{col}}) = 0 \text{ where } \mathfrak{L}_1 \mathbf{U} = (\mathbf{D}_{\text{row}} \mathbf{U}, \mathbf{U} \mathbf{D}_{\text{col}})$$

Isomorphic, but much better computationally!

\mathbf{V}, \mathbf{Z} updates as before (separable penalties + Cartesian structure)

\mathbf{U} more complicated

Splitting Methods for Convex Bi-Clustering

U -subproblem:

$$\arg \min_{\mathbf{U} \in \mathbb{R}^{n \times p}} \frac{1}{2} \|\mathbf{X} - \mathbf{U}\|_F^2 + \frac{\rho}{2} \|\mathbf{D}_{\text{row}} \mathbf{U} - \mathbf{V}_{\text{row}}^{(k)} + \rho^{-1} \mathbf{Z}_{\text{row}}^{(k)}\|_F^2 + \frac{\rho}{2} \|\mathbf{U} \mathbf{D}_{\text{col}} - \mathbf{V}_{\text{col}}^{(k)} + \rho^{-1} \mathbf{Z}_{\text{col}}^{(k)}\|_F^2$$

Stationary condition - Sylvester equation:

$$\mathbf{X} + \mathbf{D}_{\text{row}}^T (\mathbf{V}_{\text{row}}^{(k)} - \rho^{-1} \mathbf{Z}_{\text{row}}^{(k)}) + (\mathbf{V}_{\text{col}}^{(k)} - \rho^{-1} \mathbf{Z}_{\text{col}}^{(k)}) \mathbf{D}_{\text{col}}^T = \mathbf{U} + \rho \mathbf{D}_{\text{row}}^T \mathbf{D}_{\text{row}} \mathbf{U} + \rho \mathbf{U} \mathbf{D}_{\text{col}} \mathbf{D}_{\text{col}}^T$$

Splitting Methods for Convex Bi-Clustering

\mathbf{U} -subproblem:

$$\arg \min_{\mathbf{U} \in \mathbb{R}^{n \times p}} \frac{1}{2} \|\mathbf{X} - \mathbf{U}\|_F^2 + \frac{\rho}{2} \|\mathbf{D}_{\text{row}} \mathbf{U} - \mathbf{V}_{\text{row}}^{(k)} + \rho^{-1} \mathbf{Z}_{\text{row}}^{(k)}\|_F^2 + \frac{\rho}{2} \|\mathbf{U} \mathbf{D}_{\text{col}} - \mathbf{V}_{\text{col}}^{(k)} + \rho^{-1} \mathbf{Z}_{\text{col}}^{(k)}\|_F^2$$

Stationary condition - Sylvester equation:

$$\mathbf{X} + \mathbf{D}_{\text{row}}^T (\mathbf{V}_{\text{row}}^{(k)} - \rho^{-1} \mathbf{Z}_{\text{row}}^{(k)}) + (\mathbf{V}_{\text{col}}^{(k)} - \rho^{-1} \mathbf{Z}_{\text{col}}^{(k)}) \mathbf{D}_{\text{col}}^T = \mathbf{U} + \rho \mathbf{D}_{\text{row}}^T \mathbf{D}_{\text{row}} \mathbf{U} + \rho \mathbf{U} \mathbf{D}_{\text{col}} \mathbf{D}_{\text{col}}^T$$

Alternative: add quadratic term to make \mathbf{U} -subproblem easier to solve
(Deng and Yin, 2016)

$$\arg \min_{\mathbf{U} \in \mathbb{R}^{n \times p}} \cdots + \alpha \|\mathbf{U}\|_F^2 - \rho \|\mathfrak{L}_1 \mathbf{U}\|^2 \text{ where } \mathfrak{L}_1 \mathbf{U} = (\mathbf{D}_{\text{row}} \mathbf{U}, \mathbf{U} \mathbf{D}_{\text{col}})$$

$$\begin{aligned} \mathbf{U}^{(k+1)} = & \left(\alpha \mathbf{U}^{(k)} + \mathbf{X} + \rho \mathbf{D}_{\text{row}}^T (\mathbf{V}_{\text{row}}^{(k)} - \rho^{-1} \mathbf{Z}_{\text{row}}^{(k)}) - \mathbf{D}_{\text{row}} \mathbf{U}^{(k)} \right) \\ & + \rho (\mathbf{V}_{\text{col}}^{(k)} - \rho^{-1} \mathbf{Z}_{\text{col}}^{(k)} - \mathbf{U}^{(k)} \mathbf{D}_{\text{col}}) \mathbf{D}_{\text{col}}^T / (1 + \alpha) \end{aligned}$$

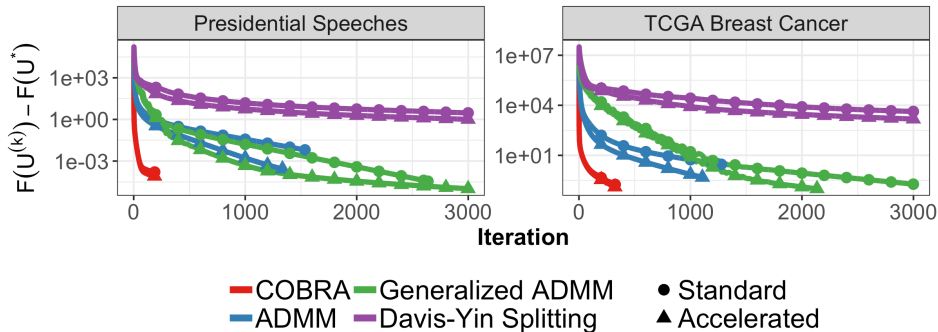
Compare:

- ADMM
- Generalized ADMM
- Davis-Yin Three-Block ADMM
- COBRA \implies alternating row- and column-clustering sub-problems

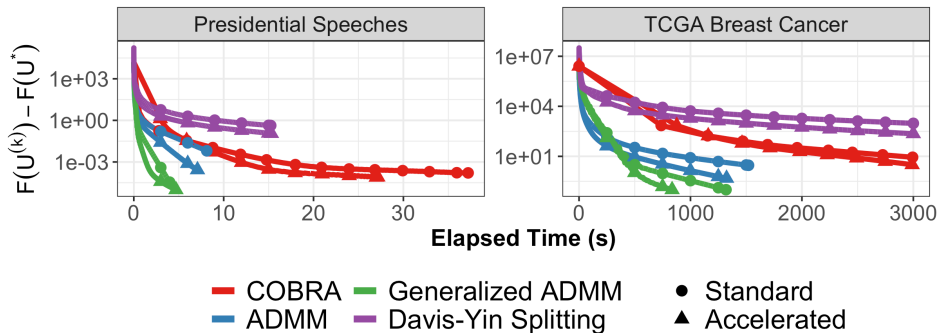
Data:

- Presidents $\in \mathbb{R}^{44 \times 75}$
- TCGA Breast Cancer $\in \mathbb{R}^{438 \times 353}$

Results: Iteration Count



Results: Elapsed Time



$$\hat{\mathcal{U}} = \arg \min_{\mathcal{U}} \frac{1}{2} \|\mathcal{X} - \mathcal{U}\|_F^2 + \lambda \sum_{j=1}^J \|\mathcal{U} \times_j \mathcal{D}_j\|_{j,q}$$

Higher-Order Extensions

$$\hat{\mathcal{U}} = \arg \min_{\mathcal{U}} \frac{1}{2} \|\mathcal{X} - \mathcal{U}\|_F^2 + \lambda \sum_{j=1}^J \|\mathcal{U} \times_j \mathcal{D}_j\|_{j,q}$$

Same “lifting” approach works for Generalized ADMM and Davis-Yin:

$$\mathcal{U}_{\text{Gen-ADMM}}^{(k+1)} = \frac{\alpha}{1+\alpha} \mathcal{U}^{(k)} + \frac{\mathcal{X}}{1+\alpha} + \frac{\rho}{1+\alpha} \sum_{j=1}^J (\mathcal{V}_j^{(k)} - \rho^{-1} \mathcal{Z}_j^{(k)} - \mathcal{U}^{(k)} \times_j \mathcal{D}_j) \times_j \mathcal{D}_j^T$$

$$\mathcal{U}_{\text{DY/AMA}}^{(k+1)} = \mathcal{X} - \sum_{j=1}^J \mathcal{Z}_j^{(k)} \times_j (\mathcal{D}_j)^T$$

$$\mathcal{V}_j^{(k+1)} = \underset{\lambda/\rho \|\cdot\|_{j,q}}{\text{prox}} \left(\mathcal{U}^{(k+1)} \times_j \mathcal{D}_j + \rho^{-1} \mathcal{Z}_j^{(k)} \right) \quad \forall j \in \{1, \dots, J\}$$

$$\mathcal{Z}_j^{(k+1)} = \mathcal{Z}_j^{(k)} + \rho (\mathcal{U}^{(k+1)} \times_j \mathcal{D}_j - \mathcal{V}_j^{(k+1)}) \quad \forall j \in \{1, \dots, J\}$$

Higher-Order Extensions

$$\hat{\mathcal{U}} = \arg \min_{\mathcal{U}} \frac{1}{2} \|\mathcal{X} - \mathcal{U}\|_F^2 + \lambda \sum_{j=1}^J \|\mathcal{U} \times_j \mathcal{D}_j\|_{j,q}$$

Same “lifting” approach works for Generalized ADMM and Davis-Yin:

$$\mathcal{U}_{\text{Gen-ADMM}}^{(k+1)} = \frac{\alpha}{1+\alpha} \mathcal{U}^{(k)} + \frac{\mathcal{X}}{1+\alpha} + \frac{\rho}{1+\alpha} \sum_{j=1}^J (\mathcal{V}_j^{(k)} - \rho^{-1} \mathcal{Z}_j^{(k)} - \mathcal{U}^{(k)} \times_j \mathcal{D}_j) \times_j \mathcal{D}_j^T$$

$$\mathcal{U}_{\text{DY/AMA}}^{(k+1)} = \mathcal{X} - \sum_{j=1}^J \mathcal{Z}_j^{(k)} \times_j (\mathcal{D}_j)^T$$

$$\mathcal{V}_j^{(k+1)} = \underset{\lambda/\rho \|\cdot\|_{j,q}}{\text{prox}} \left(\mathcal{U}^{(k+1)} \times_j \mathcal{D}_j + \rho^{-1} \mathcal{Z}_j^{(k)} \right) \quad \forall j \in \{1, \dots, J\}$$

$$\mathcal{Z}_j^{(k+1)} = \mathcal{Z}_j^{(k)} + \rho (\mathcal{U}^{(k+1)} \times_j \mathcal{D}_j - \mathcal{V}_j^{(k+1)}) \quad \forall j \in \{1, \dots, J\}$$

Standard ADMM \implies *tensor Sylvester equation*:

$$\mathcal{X} + \rho \sum_{j=1}^J (\mathcal{V}_j^{(k)} - \rho^{-1} \mathcal{Z}_j^{(k)}) \times_j \mathcal{D}_j^T = \mathcal{U}_{\text{ADMM}} + \rho \sum_{j=1}^J \mathcal{U}_{\text{ADMM}} \times_j \mathcal{D}_j \times_j \mathcal{D}_j^T.$$

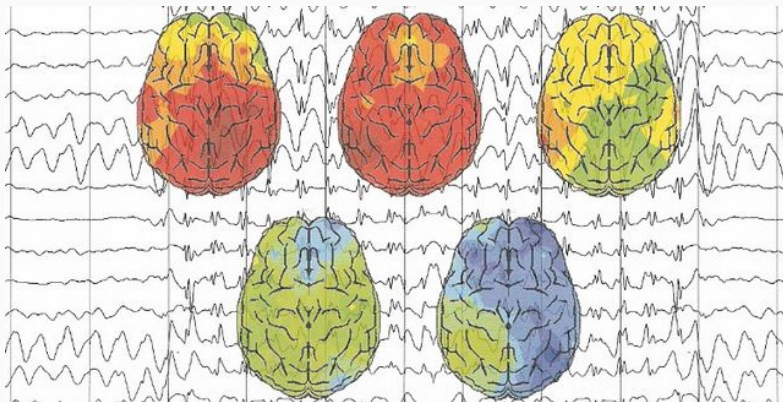
Efficient Convex Clustering Algorithm

Embed in More Complex Schemes (e.g., 2D trend filtering)

“Lifting” Trick Useful for Multiply-Regularized Problems

Multi-Rank Regularized PCA

Motivation



Principal Components Analysis:

- Exploratory Data Analysis
- Dimension Reduction
- Pattern Recognition
- Data Visualization

Regularization in PCA

Low-rank model for PCA - estimate low-rank mean of \mathbf{X} :

$$\mathbf{X} = \mathbf{u}\mathbf{v}^T + \mathbf{E} \text{ where } \mathbf{E} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2)$$

$$\arg \min_{\mathbf{u}, \mathbf{v}, d} \|\mathbf{X} - d\mathbf{u}\mathbf{v}^T\|_F^2 \Leftrightarrow \arg \max_{\mathbf{u}, \mathbf{v}} \mathbf{u}^T \mathbf{X} \mathbf{v} \quad \text{subject to } \|\mathbf{u}\| = \|\mathbf{v}\| = 1$$

Regularization in PCA

Low-rank model for PCA - estimate low-rank mean of \mathbf{X} :

$$\mathbf{X} = \mathbf{u}\mathbf{v}^T + \mathbf{E} \text{ where } \mathbf{E} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2)$$

$$\arg \min_{\mathbf{u}, \mathbf{v}, d} \|\mathbf{X} - d\mathbf{u}\mathbf{v}^T\|_F^2 \Leftrightarrow \arg \max_{\mathbf{u}, \mathbf{v}} \mathbf{u}^T \mathbf{X} \mathbf{v} \quad \text{subject to } \|\mathbf{u}\| = \|\mathbf{v}\| = 1$$

Advantages:

- Identify patterns in rows and columns of \mathbf{X}
- $\mathbf{u}, \mathbf{v}, d$ calculated using SVD of \mathbf{X}

Regularization in PCA

Low-rank model for PCA - estimate low-rank mean of \mathbf{X} :

$$\mathbf{X} = \mathbf{u}\mathbf{v}^T + \mathbf{E} \text{ where } \mathbf{E} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2)$$

$$\arg \min_{\mathbf{u}, \mathbf{v}, d} \|\mathbf{X} - d\mathbf{u}\mathbf{v}^T\|_F^2 \Leftrightarrow \arg \max_{\mathbf{u}, \mathbf{v}} \mathbf{u}^T \mathbf{X} \mathbf{v} \quad \text{subject to } \|\mathbf{u}\| = \|\mathbf{v}\| = 1$$

Advantages:

- Identify patterns in rows and columns of \mathbf{X}
- $\mathbf{u}, \mathbf{v}, d$ calculated using SVD of \mathbf{X}

PCA is consistent under standard ($n \rightarrow \infty$) asymptotics (Anderson, 1963)

Regularization in PCA

Low-rank model for PCA - estimate low-rank mean of \mathbf{X} :

$$\mathbf{X} = \mathbf{u}\mathbf{v}^T + \mathbf{E} \text{ where } \mathbf{E} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2)$$

$$\arg \min_{\mathbf{u}, \mathbf{v}, d} \|\mathbf{X} - d\mathbf{u}\mathbf{v}^T\|_F^2 \Leftrightarrow \arg \max_{\mathbf{u}, \mathbf{v}} \mathbf{u}^T \mathbf{X} \mathbf{v} \quad \text{subject to } \|\mathbf{u}\| = \|\mathbf{v}\| = 1$$

Advantages:

- Identify patterns in rows and columns of \mathbf{X}
- $\mathbf{u}, \mathbf{v}, d$ calculated using SVD of \mathbf{X}

PCA is consistent under standard ($n \rightarrow \infty$) asymptotics (Anderson, 1963)

Convergence is slow: RMT asymptotics ($p/n \rightarrow c$) more relevant

Regularization in PCA

Low-rank model for PCA - estimate low-rank mean of \mathbf{X} :

$$\mathbf{X} = \mathbf{u}\mathbf{v}^T + \mathbf{E} \text{ where } \mathbf{E} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2)$$

$$\arg \min_{\mathbf{u}, \mathbf{v}, d} \|\mathbf{X} - d\mathbf{u}\mathbf{v}^T\|_F^2 \Leftrightarrow \arg \max_{\mathbf{u}, \mathbf{v}} \mathbf{u}^T \mathbf{X} \mathbf{v} \quad \text{subject to } \|\mathbf{u}\| = \|\mathbf{v}\| = 1$$

Advantages:

- Identify patterns in rows and columns of \mathbf{X}
- $\mathbf{u}, \mathbf{v}, d$ calculated using SVD of \mathbf{X}

PCA is consistent under standard ($n \rightarrow \infty$) asymptotics (Anderson, 1963)

Convergence is slow: RMT asymptotics ($p/n \rightarrow c$) more relevant

PCA in high-dimensions is inconsistent (Johnstone and Lu, 2009)

Regularization in PCA

Low-rank model for PCA - estimate low-rank mean of \mathbf{X} :

$$\mathbf{X} = \mathbf{u}\mathbf{v}^T + \mathbf{E} \text{ where } \mathbf{E} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2)$$

$$\arg \min_{\mathbf{u}, \mathbf{v}, d} \|\mathbf{X} - d\mathbf{u}\mathbf{v}^T\|_F^2 \Leftrightarrow \arg \max_{\mathbf{u}, \mathbf{v}} \mathbf{u}^T \mathbf{X} \mathbf{v} \quad \text{subject to } \|\mathbf{u}\| = \|\mathbf{v}\| = 1$$

Advantages:

- Identify patterns in rows and columns of \mathbf{X}
- $\mathbf{u}, \mathbf{v}, d$ calculated using SVD of \mathbf{X}

PCA is consistent under standard ($n \rightarrow \infty$) asymptotics (Anderson, 1963)

Convergence is slow: RMT asymptotics ($p/n \rightarrow c$) more relevant

PCA in high-dimensions is inconsistent (Johnstone and Lu, 2009)

Low-rank model is *always* high-dimensional along one axis

Regularization in PCA

Low-rank model for PCA - estimate low-rank mean of \mathbf{X} :

$$\mathbf{X} = \mathbf{u}\mathbf{v}^T + \mathbf{E} \text{ where } \mathbf{E} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2)$$

$$\arg \min_{\mathbf{u}, \mathbf{v}, d} \|\mathbf{X} - d\mathbf{u}\mathbf{v}^T\|_F^2 \Leftrightarrow \arg \max_{\mathbf{u}, \mathbf{v}} \mathbf{u}^T \mathbf{X} \mathbf{v} \quad \text{subject to } \|\mathbf{u}\| = \|\mathbf{v}\| = 1$$

Advantages:

- Identify patterns in rows and columns of \mathbf{X}
- $\mathbf{u}, \mathbf{v}, d$ calculated using SVD of \mathbf{X}

PCA is consistent under standard ($n \rightarrow \infty$) asymptotics (Anderson, 1963)

Convergence is slow: RMT asymptotics ($p/n \rightarrow c$) more relevant

PCA in high-dimensions is inconsistent (Johnstone and Lu, 2009)

Low-rank model is *always* high-dimensional along one axis

Regularization Needed

Sparse and Functional PCA

Sparse and Functional PCA: Allen and W., (DSW 2019)

$$\arg \max_{\mathbf{u} \in \bar{\mathbb{B}}_{\mathbf{S}_u}^n, \mathbf{v} \in \bar{\mathbb{B}}_{\mathbf{S}_v}^p} \mathbf{u}^T \mathbf{X} \mathbf{v} - \lambda_u P_u(\mathbf{u}) - \lambda_v P_v(\mathbf{v})$$

where

$$\bar{\mathbb{B}}_{\mathbf{S}_u}^n = \{ \mathbf{x} \in \mathbb{R}^n : \mathbf{x}^T \mathbf{S}_u \mathbf{x} = \mathbf{x}^T (\mathbf{I} + \alpha_u \mathbf{\Omega}_u) \mathbf{x} \leq 1 \}$$

Sparse and Functional PCA

Sparse and Functional PCA: Allen and W., (DSW 2019)

$$\arg \max_{\mathbf{u} \in \bar{\mathbb{B}}_{\mathbf{S}_u}^n, \mathbf{v} \in \bar{\mathbb{B}}_{\mathbf{S}_v}^p} \mathbf{u}^T \mathbf{X} \mathbf{v} - \lambda_u P_u(\mathbf{u}) - \lambda_v P_v(\mathbf{v})$$

where

$$\bar{\mathbb{B}}_{\mathbf{S}_u}^n = \{ \mathbf{x} \in \mathbb{R}^n : \mathbf{x}^T \mathbf{S}_u \mathbf{x} = \mathbf{x}^T (\mathbf{I} + \alpha_u \mathbf{\Omega}_u) \mathbf{x} \leq 1 \}$$

Sparse and Functional PCA:

- Smoothness in \mathbf{u} - structure \mathbf{S}_u + strength α_u
- Sparsity in \mathbf{u} - structure P_u + strength λ_u
- Smoothness in \mathbf{v} - structure \mathbf{S}_v + strength α_v
- Sparsity in \mathbf{v} - structure P_v + strength λ_v

Sparse and Functional PCA

Sparse and Functional PCA: Allen and W., (DSW 2019)

$$\arg \max_{\mathbf{u} \in \bar{\mathbb{B}}_{\mathbf{S}_u}^n, \mathbf{v} \in \bar{\mathbb{B}}_{\mathbf{S}_v}^p} \mathbf{u}^T \mathbf{X} \mathbf{v} - \lambda_u P_u(\mathbf{u}) - \lambda_v P_v(\mathbf{v})$$

where

$$\bar{\mathbb{B}}_{\mathbf{S}_u}^n = \{ \mathbf{x} \in \mathbb{R}^n : \mathbf{x}^T \mathbf{S}_u \mathbf{x} = \mathbf{x}^T (\mathbf{I} + \alpha_u \mathbf{\Omega}_u) \mathbf{x} \leq 1 \}$$

Sparse and Functional PCA:

- Smoothness in \mathbf{u} - structure \mathbf{S}_u + strength α_u
- Sparsity in \mathbf{u} - structure P_u + strength λ_u
- Smoothness in \mathbf{v} - structure \mathbf{S}_v + strength α_v
- Sparsity in \mathbf{v} - structure P_v + strength λ_v

Well-posed and non-degenerate

$$\text{SFPCA: } \arg \max_{\mathbf{u} \in \bar{\mathbb{B}}_{S_u}^n, \mathbf{v} \in \bar{\mathbb{B}}_{S_v}^p} \mathbf{u}^T \mathbf{X} \mathbf{v} - \lambda_u P_u(\mathbf{u}) - \lambda_v P_v(\mathbf{v})$$

Bi-concave in \mathbf{u} and $\mathbf{v} \implies$ alternating maximization strategy

Projection + (accelerated) proximal gradient to solve sub-problems

$$\text{SFPCA: } \arg \max_{\mathbf{u} \in \overline{\mathbb{B}}_{S_u}^n, \mathbf{v} \in \overline{\mathbb{B}}_{S_v}^p} \mathbf{u}^T \mathbf{X} \mathbf{v} - \lambda_u P_u(\mathbf{u}) - \lambda_v P_v(\mathbf{v})$$

Bi-concave in \mathbf{u} and $\mathbf{v} \implies$ alternating maximization strategy

Projection + (accelerated) proximal gradient to solve sub-problems

Theorem

1. The \mathbf{u} -update converges to a solution of

$$\arg \min_{\mathbf{u} \in \overline{\mathbb{B}}_{S_u}^n} \frac{1}{2} \|\mathbf{X} \mathbf{v} - \mathbf{u}\|_2^2 + \lambda_u P_u(\mathbf{u}) + \frac{\alpha_u}{2} \mathbf{u}^T \Omega_u \mathbf{u}$$

2. \mathbf{u} -update finds global optimum for fixed \mathbf{v}

3. Converges to block-coordinate-wise global optima (Nash points)

Fast!

Orthogonality of PCs: interpretation and statistical independence

Can we do the same for SFPCA?

Multi-Rank SFPCA

Orthogonality of PCs: interpretation and statistical independence

Can we do the same for SFPCA?

Multi-rank extension of SFPCA:

$$\text{MR-SFPCA: } \arg \max_{\mathbf{U} \in \mathcal{V}_{S_u}^{n \times k}, \mathbf{V} \in \mathcal{V}_{S_v}^{p \times k}} \text{Tr}(\mathbf{U}^T \mathbf{X} \mathbf{V}) - \lambda_u P_u(\mathbf{U}) - \lambda_v P_v(\mathbf{V})$$

where $\mathcal{V}_{S_u}^{n \times k}$ is the k^{th} order generalized Stiefel manifold in \mathbb{R}^n :

$$\mathbf{U} \in \mathcal{V}_{S_u}^{n \times k} \Leftrightarrow \mathbf{U}^T S_u \mathbf{U} = \mathbf{I}_k$$

Multi-Rank SFPCA

Orthogonality of PCs: interpretation and statistical independence

Can we do the same for SFPCA?

Multi-rank extension of SFPCA:

$$\text{MR-SFPCA: } \arg \max_{\mathbf{U} \in \mathcal{V}_{\mathbf{S}_u}^{n \times k}, \mathbf{V} \in \mathcal{V}_{\mathbf{S}_v}^{p \times k}} \text{Tr}(\mathbf{U}^T \mathbf{X} \mathbf{V}) - \lambda_u P_u(\mathbf{U}) - \lambda_v P_v(\mathbf{V})$$

where $\mathcal{V}_{\mathbf{S}_u}^{n \times k}$ is the k^{th} order generalized Stiefel manifold in \mathbb{R}^n :

$$\mathbf{U} \in \mathcal{V}_{\mathbf{S}_u}^{n \times k} \Leftrightarrow \mathbf{U}^T \mathbf{S}_u \mathbf{U} = \mathbf{I}_k$$

Generalized Stiefel manifold constraint \implies manifold optimization (Absil *et al.*, 2007)

As with R1-SFPCA, alternating (partial) maximization

Manifold Proximal Gradient

Standard SFPCA \mathbf{u} -subproblem updates:

$$\mathbf{u} := \underset{\frac{\lambda \mathbf{u}}{L_{\mathbf{u}}} P_{\mathbf{u}}(\cdot)}{\text{prox}} \left(\mathbf{u} + L_{\mathbf{u}}^{-1} (\mathbf{X}\hat{\mathbf{v}} - \mathbf{S}_{\mathbf{u}}\mathbf{u}) \right) \quad \hat{\mathbf{u}} := \begin{cases} \mathbf{u} & \|\mathbf{u}\|_{\mathbf{S}_{\mathbf{u}}} \leq 1 \\ \mathbf{u} / \|\mathbf{u}\|_{\mathbf{S}_{\mathbf{u}}} & \text{otherwise} \end{cases}$$

Proximal + projected gradient descent

Manifold Proximal Gradient

Standard SFPCA \mathbf{u} -subproblem updates:

$$\mathbf{u} := \underset{\frac{\lambda_{\mathbf{u}}}{L_{\mathbf{u}}} P_{\mathbf{u}}(\cdot)}{\text{prox}} \left(\mathbf{u} + L_{\mathbf{u}}^{-1} (\mathbf{X}\hat{\mathbf{v}} - \mathbf{S}_{\mathbf{u}}\mathbf{u}) \right) \quad \hat{\mathbf{u}} := \begin{cases} \mathbf{u} & \|\mathbf{u}\|_{\mathbf{S}_{\mathbf{u}}} \leq 1 \\ \mathbf{u}/\|\mathbf{u}\|_{\mathbf{S}_{\mathbf{u}}} & \text{otherwise} \end{cases}$$

Proximal + projected gradient descent

Multi-Rank SFPCA \mathbf{U} -subproblem updates - Manifold Prox Gradient:
(Chen *et al.*, 2020a)

$$\begin{aligned} \hat{\mathbf{D}} &= \arg \min_{\mathbf{D} \in \mathbb{R}^{n \times k}} -\langle \mathbf{X}\hat{\mathbf{V}}, \mathbf{D} \rangle_F + \lambda_{\mathbf{U}} P_{\mathbf{U}}(\mathbf{U}^{(k)} + \mathbf{D}) \\ &\text{s.t. } \mathbf{D}^T \mathbf{S}_{\mathbf{U}} \mathbf{U}^{(k)} + (\mathbf{U}^{(k)})^T \mathbf{S}_{\mathbf{U}} \mathbf{D} = \mathbf{0} \\ \mathbf{U}^{(k+1)} &= \text{Retr}_{\mathbf{U}^{(k)}}(\eta \hat{\mathbf{D}}) \end{aligned}$$

Manifold Proximal Gradient

Standard SFPCA \mathbf{u} -subproblem updates:

$$\mathbf{u} := \underset{\frac{\lambda_{\mathbf{u}}}{L_{\mathbf{u}}} P_{\mathbf{u}}(\cdot)}{\text{prox}} \left(\mathbf{u} + L_{\mathbf{u}}^{-1} (\mathbf{X}\hat{\mathbf{v}} - \mathbf{S}_{\mathbf{u}}\mathbf{u}) \right) \quad \hat{\mathbf{u}} := \begin{cases} \mathbf{u} & \|\mathbf{u}\|_{\mathbf{S}_{\mathbf{u}}} \leq 1 \\ \mathbf{u}/\|\mathbf{u}\|_{\mathbf{S}_{\mathbf{u}}} & \text{otherwise} \end{cases}$$

Proximal + projected gradient descent

Multi-Rank SFPCA \mathbf{U} -subproblem updates - Manifold Prox Gradient:
(Chen *et al.*, 2020a)

$$\begin{aligned} \hat{\mathbf{D}} &= \arg \min_{\mathbf{D} \in \mathbb{R}^{n \times k}} -\langle \mathbf{X}\hat{\mathbf{V}}, \mathbf{D} \rangle_F + \lambda_{\mathbf{U}} P_{\mathbf{U}}(\mathbf{U}^{(k)} + \mathbf{D}) \\ &\text{s.t. } \mathbf{D}^T \mathbf{S}_{\mathbf{U}} \mathbf{U}^{(k)} + (\mathbf{U}^{(k)})^T \mathbf{S}_{\mathbf{U}} \mathbf{D} = \mathbf{0} \\ \mathbf{U}^{(k+1)} &= \text{Retr}_{\mathbf{U}^{(k)}}(\eta \hat{\mathbf{D}}) \end{aligned}$$

One step of each subproblem \implies convergence to stationary point:

- Constraint set smooth \implies Stationary points isolated
- Guaranteed descent at each iteration (Chen *et al.*, 2020b)

Easier \mathbf{U} -updates from Manifold ADMM: (Kovnatsky *et al.*, 2016)

$$\mathbf{U}^{(k+1)} = \arg \min_{\mathbf{U} \in \mathcal{V}_{n \times k}^{S_{\mathbf{U}}}} -\text{Tr}(\mathbf{U}^T \mathbf{X} \mathbf{V}) + \frac{\rho}{2} \|\mathbf{U} - \mathbf{W}^{(k)} + \mathbf{Z}^{(k)}\|_F^2$$

$$\mathbf{W}^{(k+1)} = \arg \min_{\mathbf{W} \in \mathbb{R}^{n \times k}} \lambda_{\mathbf{U}} P_{\mathbf{U}}(\mathbf{W}) + \frac{\rho}{2} \|\mathbf{U}^{(k+1)} - \mathbf{W} + \mathbf{Z}^{(k)}\|_F^2$$

$$= \text{prox}_{\lambda_{\mathbf{U}}/\rho P_{\mathbf{U}}(\cdot)} \left(\mathbf{U}^{(k+1)} + \mathbf{Z}^{(k)} \right)$$

$$\mathbf{Z}^{(k+1)} = \mathbf{Z}^{(k)} + \mathbf{U}^{(k+1)} - \mathbf{W}^{(k+1)}$$

Manifold ADMM

Easier \mathbf{U} -updates from Manifold ADMM: (Kovnatsky *et al.*, 2016)

$$\mathbf{U}^{(k+1)} = \arg \min_{\mathbf{U} \in \mathcal{V}_{n \times k}^{\mathcal{S}_U}} -\text{Tr}(\mathbf{U}^T \mathbf{X} \mathbf{V}) + \frac{\rho}{2} \|\mathbf{U} - \mathbf{W}^{(k)} + \mathbf{Z}^{(k)}\|_F^2$$

$$\mathbf{W}^{(k+1)} = \arg \min_{\mathbf{W} \in \mathbb{R}^{n \times k}} \lambda_U P_U(\mathbf{W}) + \frac{\rho}{2} \|\mathbf{U}^{(k+1)} - \mathbf{W} + \mathbf{Z}^{(k)}\|_F^2$$

$$= \text{prox}_{\lambda_U / \rho P_U(\cdot)} \left(\mathbf{U}^{(k+1)} + \mathbf{Z}^{(k)} \right)$$

$$\mathbf{Z}^{(k+1)} = \mathbf{Z}^{(k)} + \mathbf{U}^{(k+1)} - \mathbf{W}^{(k+1)}$$

First step is a *generalized unbalanced Procrustes problem* - analytical solution via SVD of $\mathbf{S}_U^{-1/2} \mathbf{X} \hat{\mathbf{V}} + \rho \mathbf{S}_U^{1/2} (\mathbf{W}^{(k)} - \mathbf{Z}^{(k)})$

Typically converges quickly and to a good solution (no theory)

$$\text{Rank-1 SFPCA: } \arg \max_{\mathbf{u} \in \bar{\mathbb{B}}_{S_u}^n, \mathbf{v} \in \bar{\mathbb{B}}_{S_v}^p} \mathbf{u}^T \mathbf{X} \mathbf{v} - \lambda_u P_u(\mathbf{u}) - \lambda_v P_v(\mathbf{v})$$

How to get additional *nested* SFPCA components?

Deflation Techniques

$$\text{Rank-1 SFPCA: } \arg \max_{\mathbf{u} \in \overline{\mathbb{B}}_{S_u}^n, \mathbf{v} \in \overline{\mathbb{B}}_{S_v}^p} \mathbf{u}^T \mathbf{X} \mathbf{v} - \lambda_u P_u(\mathbf{u}) - \lambda_v P_v(\mathbf{v})$$

How to get additional *nested* SFPCA components?

Deflation:

- Hotelling: $\mathbf{X}_t^{\text{HD}} := \mathbf{X}_{t-1} - \mathbf{U}_t(\mathbf{U}_t^T \mathbf{U}_t)^{-1} \mathbf{U}_t^T \mathbf{X}_{t-1} \mathbf{V}_t(\mathbf{V}_t^T \mathbf{V}_t)^{-1} \mathbf{V}_t^T$
- Projection: $\mathbf{X}_t^{\text{PD}} := (\mathbf{I}_n - \mathbf{U}_t(\mathbf{U}_t^T \mathbf{U}_t)^{-1} \mathbf{U}_t^T) \mathbf{X}_{t-1} (\mathbf{I}_p - \mathbf{V}_t(\mathbf{V}_t^T \mathbf{V}_t)^{-1} \mathbf{V}_t^T)$
- Schur Complement: $\mathbf{X}_t^{\text{SD}} := \mathbf{X}_{t-1} - \mathbf{X}_{t-1} \mathbf{V}_t(\mathbf{U}_t^T \mathbf{X}_{t-1} \mathbf{V}_t)^{-1} \mathbf{U}_t^T \mathbf{X}_{t-1}$

HD doesn't fully remove signal and may re-introduce if estimated PCs non-orthogonal

Deflation Techniques

$$\text{Rank-1 SFPCA: } \arg \max_{\mathbf{u} \in \overline{\mathbb{B}}_{S_u}^n, \mathbf{v} \in \overline{\mathbb{B}}_{S_v}^p} \mathbf{u}^T \mathbf{X} \mathbf{v} - \lambda_u P_u(\mathbf{u}) - \lambda_v P_v(\mathbf{v})$$

How to get additional *nested* SFPCA components?

Deflation:

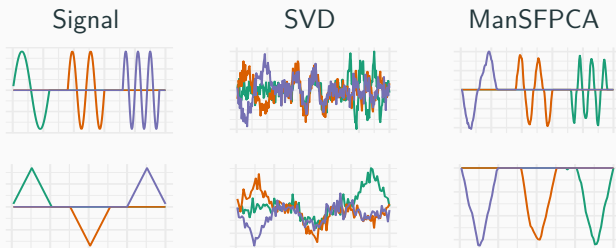
- Hotelling: $\mathbf{X}_t^{\text{HD}} := \mathbf{X}_{t-1} - \mathbf{U}_t (\mathbf{U}_t^T \mathbf{U}_t)^{-1} \mathbf{U}_t^T \mathbf{X}_{t-1} \mathbf{V}_t (\mathbf{V}_t^T \mathbf{V}_t)^{-1} \mathbf{V}_t^T$
- Projection: $\mathbf{X}_t^{\text{PD}} := (\mathbf{I}_n - \mathbf{U}_t (\mathbf{U}_t^T \mathbf{U}_t)^{-1} \mathbf{U}_t^T) \mathbf{X}_{t-1} (\mathbf{I}_p - \mathbf{V}_t (\mathbf{V}_t^T \mathbf{V}_t)^{-1} \mathbf{V}_t^T)$
- Schur Complement: $\mathbf{X}_t^{\text{SD}} := \mathbf{X}_{t-1} - \mathbf{X}_{t-1} \mathbf{V}_t (\mathbf{U}_t^T \mathbf{X}_{t-1} \mathbf{V}_t)^{-1} \mathbf{U}_t^T \mathbf{X}_{t-1}$

HD doesn't fully remove signal and may re-introduce if estimated PCs non-orthogonal

Method	Two-Way 0-ing $\mathbf{u}_t^T \mathbf{X}_t \mathbf{v}_t = 0$	One-Way 0-ing $\mathbf{u}_t^T \mathbf{X}_t, \mathbf{X}_t \mathbf{v}_t = 0$	Subsequent 0-ing ($\forall s \geq 0$) $\mathbf{u}_t^T \mathbf{X}_{t+s}, \mathbf{X}_{t+s} \mathbf{v}_t = 0$	Robust to Scale of $\mathbf{u}_t, \mathbf{v}_t$
Hotelling	✓	✗	✗	✗
Projection	✓	✓	✗	✗
Schur	✓	✓	✓	✓

Simulation: "On Model" Signal Recovery

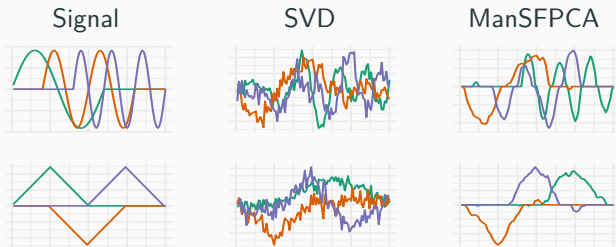
Scenario 1: \mathbf{U}^* and \mathbf{V}^* Orthogonal - SNR ≈ 1.2



		HD	PD	SD	ManSFPCA
CPVE	PC1	15.92%	21.05%	21.87%	
	PC2	22.21%	29.42%	30.59%	37.12%
	PC3	26.80%	35.57%	37.09%	
rSS-Error	\mathbf{U}	129.54%	129.55%	128.35%	68.66%
	\mathbf{V}	143.01%	143.72%	141.15%	36.98%

Simulation: ‘Off-Model’ Signal Recovery

Scenario 2: \mathbf{U}^* and \mathbf{V}^* **Not Orthogonal** - SNR ≈ 1.7



		HD	PD	SD	ManSFPCA
CPVE	PC1	8.85%	19.74%	29.80%	
	PC2	13.03%	28.30%	39.87%	50.85%
	PC3	16.16%	34.22%	46.48%	
rSS-Error	\mathbf{U}	215.73%	206.30%	205.74%	97.77%
	\mathbf{V}	211.15%	207.77%	204.38%	78.26%

Principled Approach to Multi-Rank PCA:

- Unifies **many** regularized PCA variants
- Extension to multiple PCs without losing orthogonality
- Deflation techniques applicable to all PCA techniques

Principled Approach to Multi-Rank PCA:

- Unifies **many** regularized PCA variants
- Extension to multiple PCs without losing orthogonality
- Deflation techniques applicable to all PCA techniques

Tensor extensions:

- Rank-1 SFPCA yields regularized CP
- MR-SFPCA yields regularized Tucker

Principled Approach to Multi-Rank PCA:

- Unifies **many** regularized PCA variants
- Extension to multiple PCs without losing orthogonality
- Deflation techniques applicable to all PCA techniques

Tensor extensions:

- Rank-1 SFPCA yields regularized CP
- MR-SFPCA yields regularized Tucker

Additional Multivariate Methods:

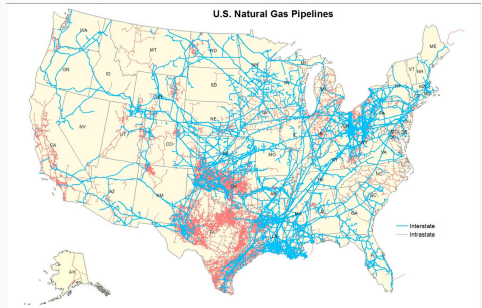
- CCA, LDA, PLS, *etc.* all SVD - can all be similarly treated

Multivariate Models for Gas Markets

Natural Gas Markets

LNG Markets:

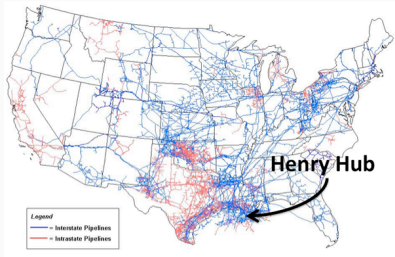
- 32% of all US Electricity (1.273 PWh in 2017)
- 3M Miles of NG Pipelines
- 150+ Trading Spots



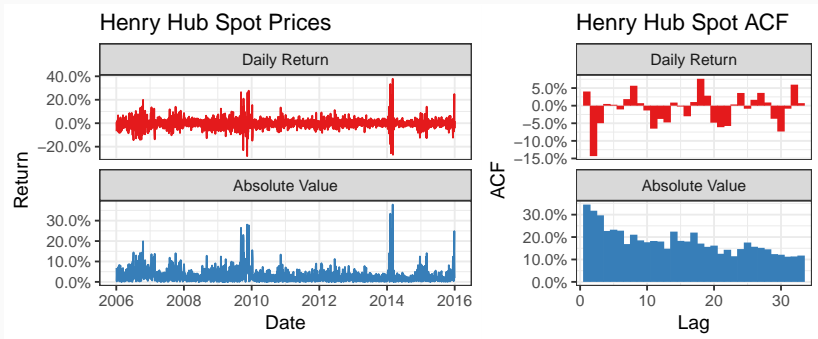
Natural Gas Markets

Henry Hub:

- \$14B+ Futures Volume Daily
- Common Proxy for Domestic NG Markets Broadly

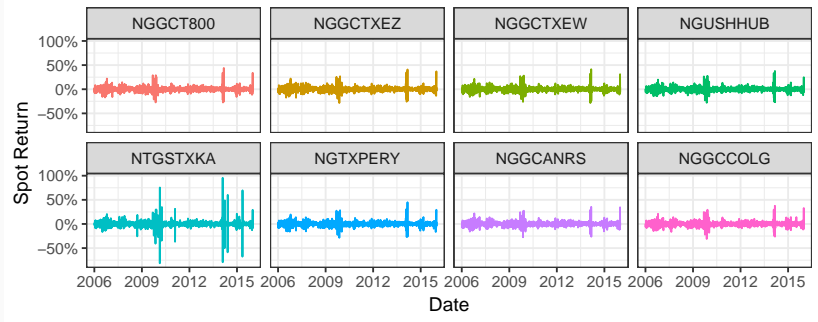


Natural Gas Markets



Equity-like dynamics: vol clustering, heavy-tails, 2nd moment autocorrelation

Natural Gas Markets



High inter-spot correlation: PC1 (74%) suggests single-factor model

We want to capture:

- High-Dimensional Multivariate Time Series
- Irregular Data Availability
 - NG Futures Priced Near-Continuously on Lit Markets
 - NG Spots Traded Over-the-Counter
- GARCH Type Behavior + Single-Factor Structure

We want to capture:

- High-Dimensional Multivariate Time Series
- Irregular Data Availability
 - NG Futures Priced Near-Continuously on Lit Markets
 - NG Spots Traded Over-the-Counter
- GARCH Type Behavior + Single-Factor Structure

💡 Realized Beta GARCH Model (Hansen *et al.*, 2014) combining:

- Intra-Day Futures Realized Volatility
- End-of-Day Spot Volatility
- 2nd Moment Single-Factor Dynamics

- Bi-Variate GARCH Model (Multivariate Skew Normal Specification):
 - “Realized” (High-Frequency) Volatility: improved estimate of σ_t^2
 - “Beta” Volatility linkage: Volatility at Henry \implies volatility in spots

- Bi-Variate GARCH Model (Multivariate Skew Normal Specification):
 - “Realized” (High-Frequency) Volatility: improved estimate of σ_t^2
 - “Beta” Volatility linkage: Volatility at Henry \implies volatility in spots
- Bayesian Estimation:
 - Priors calibrated to S&P 500 (equity) markets: improve estimation
 - Coherent uncertainty propagation
 - Improved out of sample forecasts

Does it work?

Fitting strategy:

- Fit to 250 window: refit every 50 days
- One-day rolling predictions for out-of-sample

Does it work?

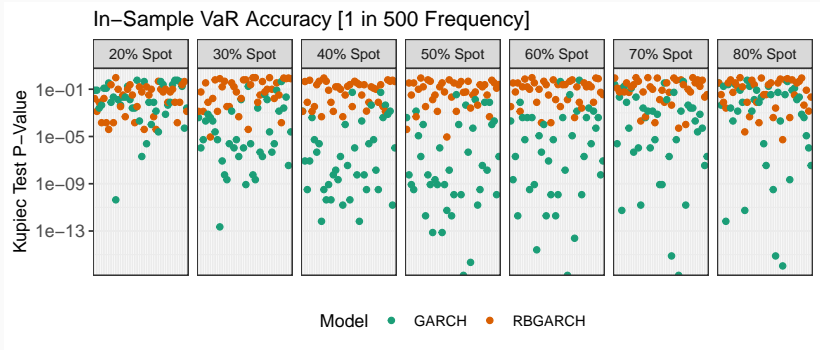
Fitting strategy:

- Fit to 250 window: refit every 50 days
- One-day rolling predictions for out-of-sample

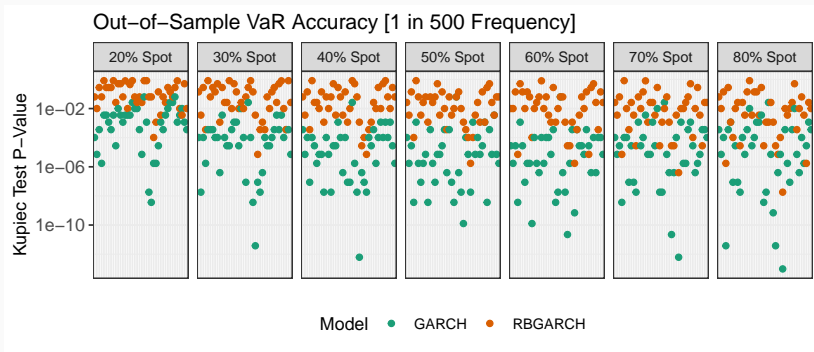
Measures of Fit:

- In-sample VaR Test (Kupiec, 1995)
 - Binomial test for number of VaR exceedances
- Out-of-sample VaR Test (Kupiec, 1995)
 - Estimated out-of-sample log-likelihood

Application to Tail Forecasting

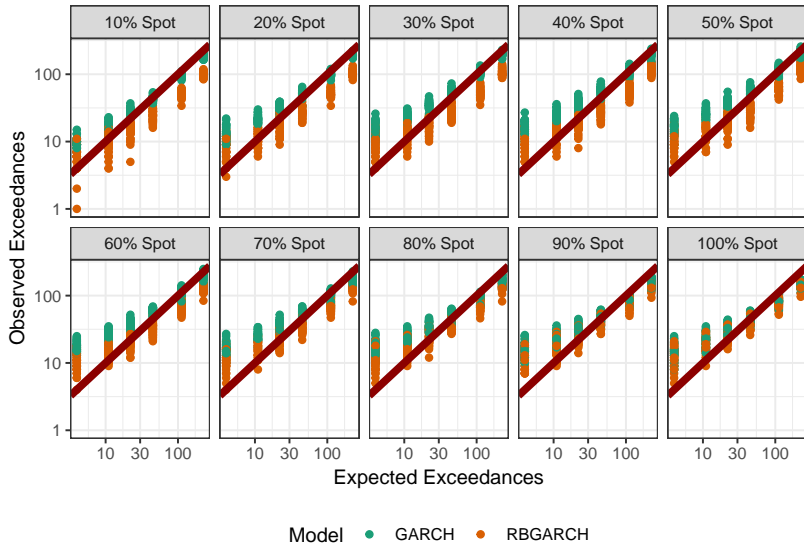


Application to Tail Forecasting

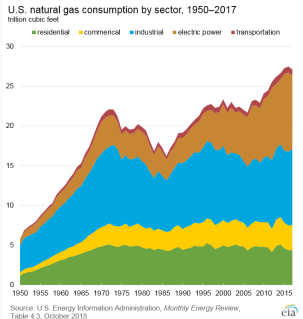


Application to Tail Forecasting

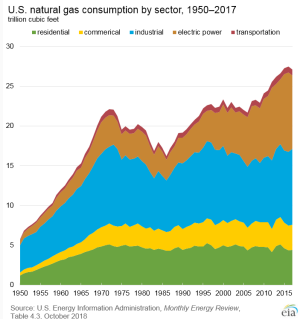
Out-of-Sample VaR Performance



- Multivariate and Multi-Resolution Model for NG Volatility
 - Daily and intra-day volatility measures
 - Multivariate Treatment of 50+ NG Trading Hubs

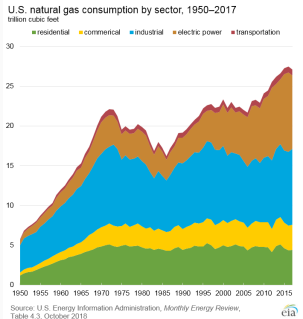


- Multivariate and Multi-Resolution Model for NG Volatility
 - Daily and intra-day volatility measures
 - Multivariate Treatment of 50+ NG Trading Hubs
- Bayesian Approach
 - Market Calibrated Priors



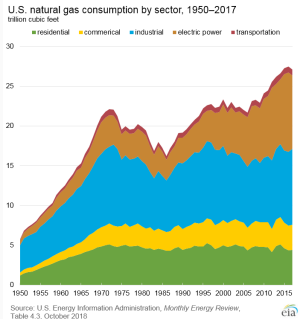
Implications

- Multivariate and Multi-Resolution Model for NG Volatility
 - Daily and intra-day volatility measures
 - Multivariate Treatment of 50+ NG Trading Hubs
- Bayesian Approach
 - Market Calibrated Priors
- Improved Out-of-Sample Prediction
 - VaR Estimates: more accurate + more conservative



Implications

- Multivariate and Multi-Resolution Model for NG Volatility
 - Daily and intra-day volatility measures
 - Multivariate Treatment of 50+ NG Trading Hubs
- Bayesian Approach
 - Market Calibrated Priors
- Improved Out-of-Sample Prediction
 - VaR Estimates: more accurate + more conservative
- Amenable to all commodities markets with irregular data availability



Complex Convex Analysis

Complex-data arise in many domains:

- Signal and radar processing (Schreier and Scharf, 2010; Candès *et al.*, 2015; Mechlenbrauker *et al.*, 2017)
- Neuroscience (Yu *et al.*, 2018; Adrian *et al.*, 2018)
- Geostatistics (de Iaco *et al.*, 2003; Mandic *et al.*, 2009)
- Astronomy (Zechmeister and Kürster, 2009)
- Econometrics (Granger and Engle, 1983)

Complex-data arise in many domains:

- Signal and radar processing (Schreier and Scharf, 2010; Candès *et al.*, 2015; Mechlenbrauker *et al.*, 2017)
- Neuroscience (Yu *et al.*, 2018; Adrian *et al.*, 2018)
- Geostatistics (de Iaco *et al.*, 2003; Mandic *et al.*, 2009)
- Astronomy (Zechmeister and Kürster, 2009)
- Econometrics (Granger and Engle, 1983)

Major sources:

- Spectral analysis (Fourier transforms)
- 2D directional data

Why Complex?

Why not treat complex data as real?

Why Complex?

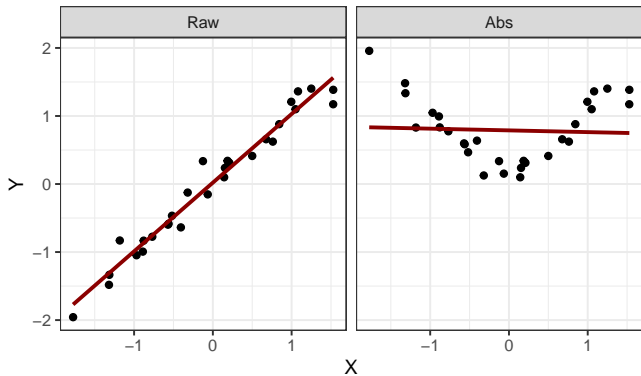
Why not treat complex data as real?

- Discards phase data (Canolty *et al.*, 2006)

Why Complex?

Why not treat complex data as real?

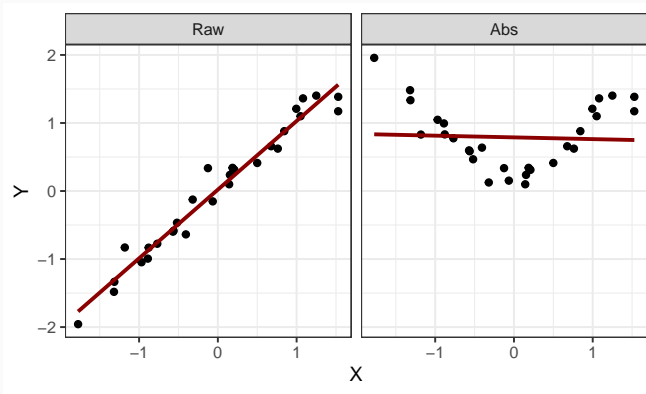
- Discards phase data (Canolty *et al.*, 2006)
- Corrupts statistical relationships:



Why Complex?

Why not treat complex data as real?

- Discards phase data (Canolty *et al.*, 2006)
- Corrupts statistical relationships:



- Natural domain for “spectral” phenomena

Proper and Improper RVs

Let Z be a *univariate* complex random variable:

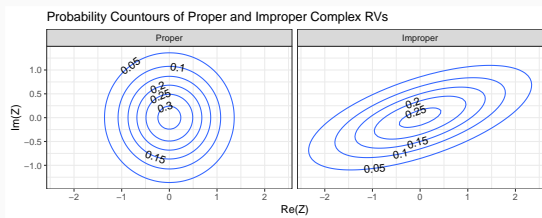
- Secretly “multivariate”: correlation between $\Re(Z)$ and $\Im(Z)$

Proper and Improper RVs

Let Z be a *univariate* complex random variable:

- Secretly “multivariate”: correlation between $\Re(Z)$ and $\Im(Z)$

Important case: *proper* (circular) RVs are those where $\Re(Z) \perp\!\!\!\perp \Im(Z)$

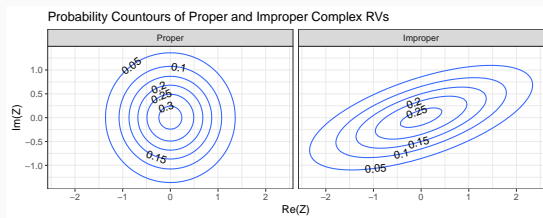


Proper and Improper RVs

Let Z be a *univariate* complex random variable:

- Secretly “multivariate”: correlation between $\Re(Z)$ and $\Im(Z)$

Important case: *proper* (circular) RVs are those where $\Re(Z) \perp\!\!\!\perp \Im(Z)$



Complex variables often arise as Fourier transform of *stationary* processes

- Only *relative* phase of multivariate signal matters
- Absolute phase is meaningless

$$\text{Law}[Z] = \text{Law}[e^{i\theta} Z] \implies Z \text{ proper}$$

Complex Gaussian Distribution

The *complex* Gaussian is a **three** parameter distribution:

- Mean: $\boldsymbol{\mu} = \mathbb{E}[\mathbf{Z}]$
- Covariance: $\boldsymbol{\Sigma} = \mathbb{E}[(\mathbf{Z} - \boldsymbol{\mu})(\mathbf{Z} - \boldsymbol{\mu})^H] = \mathbb{E}[(\mathbf{Z} - \boldsymbol{\mu})(\overline{\mathbf{Z} - \boldsymbol{\mu}})^T]$
- Pseudo-Covariance: $\boldsymbol{\Gamma} = \mathbb{E}[(\mathbf{Z} - \boldsymbol{\mu})(\mathbf{Z} - \boldsymbol{\mu})^T]$

Complex Gaussian Distribution

The *complex* Gaussian is a **three** parameter distribution:

- Mean: $\boldsymbol{\mu} = \mathbb{E}[\mathbf{Z}]$
- Covariance: $\boldsymbol{\Sigma} = \mathbb{E}[(\mathbf{Z} - \boldsymbol{\mu})(\mathbf{Z} - \boldsymbol{\mu})^H] = \mathbb{E}[(\mathbf{Z} - \boldsymbol{\mu})(\mathbf{Z} - \boldsymbol{\mu})^T]$
- Pseudo-Covariance: $\boldsymbol{\Gamma} = \mathbb{E}[(\mathbf{Z} - \boldsymbol{\mu})(\mathbf{Z} - \boldsymbol{\mu})^T]$

Requirements:

- $\boldsymbol{\Sigma}$ is positive-definite: $\Sigma_{ij} = \mathbb{E}[z_i \bar{z}_j]$ - non-negative when $i = j$
- $\boldsymbol{\Gamma}$ is indefinite: $\Gamma_{ij} = \mathbb{E}[z_i z_j]$ - possibly negative when $i = j$
- Additionally $\bar{\boldsymbol{\Sigma}} - \boldsymbol{\Gamma}^H \boldsymbol{\Sigma}^{-1} \boldsymbol{\Gamma} \succeq 0$

Complex Gaussian Distribution

The *complex* Gaussian is a **three** parameter distribution:

- Mean: $\boldsymbol{\mu} = \mathbb{E}[\mathbf{Z}]$
- Covariance: $\boldsymbol{\Sigma} = \mathbb{E}[(\mathbf{Z} - \boldsymbol{\mu})(\mathbf{Z} - \boldsymbol{\mu})^H] = \mathbb{E}[(\mathbf{Z} - \boldsymbol{\mu})(\mathbf{Z} - \boldsymbol{\mu})^T]$
- Pseudo-Covariance: $\boldsymbol{\Gamma} = \mathbb{E}[(\mathbf{Z} - \boldsymbol{\mu})(\mathbf{Z} - \boldsymbol{\mu})^T]$

Requirements:

- $\boldsymbol{\Sigma}$ is positive-definite: $\Sigma_{ij} = \mathbb{E}[z_i \bar{z}_j]$ - non-negative when $i = j$
- $\boldsymbol{\Gamma}$ is indefinite: $\Gamma_{ij} = \mathbb{E}[z_i z_j]$ - possibly negative when $i = j$
- Additionally $\bar{\boldsymbol{\Sigma}} - \boldsymbol{\Gamma}^H \boldsymbol{\Sigma}^{-1} \boldsymbol{\Gamma} \succeq 0$

Proper complex normal ($\boldsymbol{\Gamma} = \mathbf{0}$) almost universally assumed in statistics (Wooding, 1956; Goodman, 1963; Graczyk *et al.*, 2003)

Non-stationary DSP sometimes uses general case (van den Bos, 1995; Schreier and Scharf, 2010; Adali *et al.*, 2011)

Complex Convex Analysis: Subgradient Analysis

Penalized M -Estimation Paradigm:

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \mathcal{L}(\beta; \mathbf{X}, \mathbf{y}) + \lambda \mathcal{P}(\beta)$$

What if $\mathbf{X}, \mathbf{y}, \beta$ complex?

Complex Convex Analysis: Subgradient Analysis

Penalized M -Estimation Paradigm:

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \mathcal{L}(\beta; \mathbf{X}, \mathbf{y}) + \lambda \mathcal{P}(\beta)$$

What if $\mathbf{X}, \mathbf{y}, \beta$ complex? Well defined if $\mathcal{L}, \mathcal{P} \rightarrow \mathbb{R}$ (e.g. norms!)

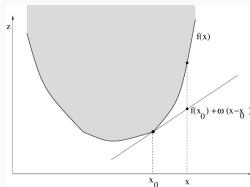
Complex Convex Analysis: Subgradient Analysis

Penalized M -Estimation Paradigm:

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \mathcal{L}(\beta; \mathbf{X}, \mathbf{y}) + \lambda \mathcal{P}(\beta)$$

What if $\mathbf{X}, \mathbf{y}, \beta$ complex? Well defined if $\mathcal{L}, \mathcal{P} \rightarrow \mathbb{R}$ (e.g. norms!)

Analysis: first order sub-gradient conditions



γ is a sub-gradient of f at x if

$$f(y) \geq f(x) + \gamma(y - x) \quad \text{for all } y$$

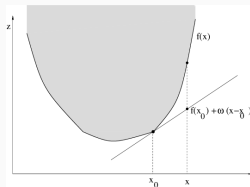
Complex Convex Analysis: Subgradient Analysis

Penalized M -Estimation Paradigm:

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \mathcal{L}(\beta; \mathbf{X}, \mathbf{y}) + \lambda \mathcal{P}(\beta)$$

What if $\mathbf{X}, \mathbf{y}, \beta$ complex? Well defined if $\mathcal{L}, \mathcal{P} \rightarrow \mathbb{R}$ (e.g. norms!)

Analysis: first order sub-gradient conditions



γ is a sub-gradient of f at x if

$$f(y) \geq f(x) + \gamma(y - x) \quad \text{for all } y$$

If x, y are *complex*, this is ill-defined!



The Wirtinger Fix

Wirtinger (1927) studied differentiability of functions $\mathbb{C} \rightarrow \mathbb{R}$

Never Cauchy-Riemann differentiable (holomorphic) \implies traditional complex analysis does not apply

The Wirtinger Fix

Wirtinger (1927) studied differentiability of functions $\mathbb{C} \rightarrow \mathbb{R}$

Never Cauchy-Riemann differentiable (holomorphic) \implies traditional complex analysis does not apply

Wirtinger's idea: write $f(z) = f(z, \bar{z})$ and differentiate with respect to z while holding \bar{z} fixed

The Wirtinger Fix

Wirtinger (1927) studied differentiability of functions $\mathbb{C} \rightarrow \mathbb{R}$

Never Cauchy-Riemann differentiable (holomorphic) \implies traditional complex analysis does not apply

Wirtinger's idea: write $f(z) = f(z, \bar{z})$ and differentiate with respect to z while holding \bar{z} fixed

Occasionally used as $\mathbb{C}\mathbb{R}$ -calculus (Kreutz-Delgado, 2009)

The Wirtinger Fix

Wirtinger (1927) studied differentiability of functions $\mathbb{C} \rightarrow \mathbb{R}$

Never Cauchy-Riemann differentiable (holomorphic) \implies traditional complex analysis does not apply

Wirtinger's idea: write $f(z) = f(z, \bar{z})$ and differentiate with respect to z while holding \bar{z} fixed

Occasionally used as $\mathbb{C}\mathbb{R}$ -calculus (Kreutz-Delgado, 2009)

GOAL: rigorously define this derivative and connect it to optimization

The Wirtinger Fix

\mathbb{C}^p is an *inner product space* over a field \mathbb{F} :

- Can add vectors (elements of \mathbb{C}^p) and multiply by \mathbb{F}
- Inner product $\langle \cdot, \cdot \rangle : \mathbb{C}^p \times \mathbb{C}^p \rightarrow \mathbb{F}$

Typically $\mathbb{F} = \mathbb{C}$

The Wirtinger Fix

\mathbb{C}^p is an *inner product space* over a field \mathbb{F} :

- Can add vectors (elements of \mathbb{C}^p) and multiply by \mathbb{F}
- Inner product $\langle \cdot, \cdot \rangle : \mathbb{C}^p \times \mathbb{C}^p \rightarrow \mathbb{F}$

Typically $\mathbb{F} = \mathbb{C}$

In this work, $\mathbb{F} = \mathbb{R}$!

$$\langle \mathbf{a}, \mathbf{b} \rangle = \frac{\mathbf{a}^H \mathbf{b} + \mathbf{a}^T \bar{\mathbf{b}}}{2}$$

Sub-gradient inequality becomes

$$f(\mathbf{w}) \geq f(\mathbf{z}) + \langle \boldsymbol{\gamma}, \mathbf{w} - \mathbf{z} \rangle \quad \text{for all } \mathbf{w} \in \mathbb{C}^p$$

All terms real \implies well-defined!

The Wirtinger Fix

Is this valid?

The Wirtinger Fix

Is this valid?

Identity and existence of minimizer are *topological* properties

$$\langle \mathbf{a}, \mathbf{a} \rangle = \sum_{i=1}^p |a_i|^2$$

- Real inner product \implies the same norm
- Same norm \implies same topology
- Same topology \implies minimizer unchanged

The Wirtinger Fix

Is this valid?

Identity and existence of minimizer are *topological* properties

$$\langle \mathbf{a}, \mathbf{a} \rangle = \sum_{i=1}^p |a_i|^2$$

- Real inner product \implies the same norm
- Same norm \implies same topology
- Same topology \implies minimizer unchanged

We have freedom to change *algebraic* structure used to analyze problem

Likelihood still defined with “regular” complex multiplication

The Wirtinger Fix

Convex Analysis for Wirtinger ($\mathbb{C}^p \rightarrow \mathbb{R}$) Functions:

The Wirtinger Fix

Convex Analysis for Wirtinger ($\mathbb{C}^p \rightarrow \mathbb{R}$) Functions:

- **W. Theorem 2.2:** If Wirtinger f is convex, it has sub-gradients.

The Wirtinger Fix

Convex Analysis for Wirtinger ($\mathbb{C}^p \rightarrow \mathbb{R}$) Functions:

- **W. Theorem 2.2:** If Wirtinger f is convex, it has sub-gradients.
- **W. Theorem 2.6:** Sub-gradients given by the Wirtinger (formal) derivative.

The Wirtinger Fix

Convex Analysis for Wirtinger ($\mathbb{C}^P \rightarrow \mathbb{R}$) Functions:

- **W. Theorem 2.2:** If Wirtinger f is convex, it has sub-gradients.
- **W. Theorem 2.6:** Sub-gradients given by the Wirtinger (formal) derivative.
- **W. Theorem 2.1:** If the Wirtinger derivative is 0, global minimum.

The Wirtinger Fix

Convex Analysis for Wirtinger ($\mathbb{C}^P \rightarrow \mathbb{R}$) Functions:

- **W. Theorem 2.2:** If Wirtinger f is convex, it has sub-gradients.
- **W. Theorem 2.6:** Sub-gradients given by the Wirtinger (formal) derivative.
- **W. Theorem 2.1:** If the Wirtinger derivative is 0, global minimum.
- **W. Theorem 2.8:** Key result for sparse models:

$$\partial|z| = \begin{cases} \angle z = z/|z| & z \neq 0 \\ \{w \in \mathbb{C} : |w| \leq 1\} & z = 0 \end{cases}$$

The Wirtinger Fix

Convex Analysis for Wirtinger ($\mathbb{C}^P \rightarrow \mathbb{R}$) Functions:

- **W. Theorem 2.2:** If Wirtinger f is convex, it has sub-gradients.
- **W. Theorem 2.6:** Sub-gradients given by the Wirtinger (formal) derivative.
- **W. Theorem 2.1:** If the Wirtinger derivative is 0, global minimum.
- **W. Theorem 2.8:** Key result for sparse models:

$$\partial|z| = \begin{cases} \angle z = z/|z| & z \neq 0 \\ \{w \in \mathbb{C} : |w| \leq 1\} & z = 0 \end{cases}$$

Change definition of “multiplication”
convex analysis still “works!”

Example: Complex OLS

$$\arg \min_{\beta \in \mathbb{C}^p} \|\mathbf{y} - \mathbf{X}\beta\|_2^2$$

Set Wirtinger derivative to $\mathbf{0}$:

$$\begin{aligned} f &= \|\mathbf{y} - \mathbf{X}\beta\|_2^2 \\ &= (\mathbf{y} - \mathbf{X}\beta)^H (\mathbf{y} - \mathbf{X}\beta) \\ &= (\mathbf{y}^H - \bar{\beta} \mathbf{X}^H) (\mathbf{y} - \mathbf{X}\beta) \\ \mathbf{0} &= \frac{\partial f}{\partial \beta} = -\mathbf{y}^H \mathbf{X} + \bar{\beta}^T \mathbf{X}^H \mathbf{X} \\ \implies \bar{\beta} &= (\mathbf{X}^H \mathbf{X})^{-1} \mathbf{X}^T \bar{\mathbf{y}} \\ \beta &= (\mathbf{X}^H \mathbf{X})^{-1} \mathbf{X}^H \mathbf{y} \end{aligned}$$

Intuition from real-domain translates to complex-domain!

Concentration Inequalities: Regression Noise (W., Lemma 3.9)

Suppose Z is mean-zero sub-Gaussian with variance proxy σ^2 :

(a) If Z is real:

$$P[|Z| \geq t] \leq 2 \exp\{-t^2/2\sigma^2\}$$

Concentration Inequalities: Regression Noise (W., Lemma 3.9)

Suppose Z is mean-zero sub-Gaussian with variance proxy σ^2 :

(a) If Z is real:

$$P[|Z| \geq t] \leq 2 \exp\{-t^2/2\sigma^2\}$$

(b) If Z is proper and complex:

$$P[|Z| \geq t] \leq 4 \exp\{-t^2/2\sigma^2\}$$

Proper and complex Z concentrates like real 2-vector.

Concentration Inequalities: Regression Noise (W., Lemma 3.9)

Suppose Z is mean-zero sub-Gaussian with variance proxy σ^2 :

(a) If Z is real:

$$P[|Z| \geq t] \leq 2 \exp\{-t^2/2\sigma^2\}$$

(b) If Z is proper and complex:

$$P[|Z| \geq t] \leq 4 \exp\{-t^2/2\sigma^2\}$$

(c) If Z is complex:

$$P[|Z| \geq t] \leq 2 \exp\{-t^2/16\sigma^2\}$$

Proper and complex Z concentrates like real 2-vector.

Penalty for unknown dependence in $\Re(Z)$ and $\Im(Z)$

Concentration Inequalities: Regression Noise (W., Lemma 3.9)

Suppose Z is mean-zero sub-Gaussian with variance proxy σ^2 :

(a) If Z is real:

$$P[|Z| \geq t] \leq 2 \exp\{-t^2/2\sigma^2\}$$

(b) If Z is proper and complex:

$$P[|Z| \geq t] \leq 4 \exp\{-t^2/2\sigma^2\}$$

(c) If Z is complex:

$$P[|Z| \geq t] \leq 2 \exp\{-t^2/16\sigma^2\}$$

Proper and complex Z concentrates like real 2-vector.

Penalty for unknown dependence in $\Re(Z)$ and $\Im(Z)$

Similar rates for effective noise $\|\mathbf{X}^H \epsilon\|_\infty$ depending on ϵ

Concentration Inequalities: Sample Covariance (W., Lemma 4.3)

Suppose \mathbf{Z} is mean-zero Gaussian vector with variance Σ^* :

$$[\sigma^2 = \max(\Sigma_{ii}^*)]$$

(a) If Z is real:

$$P[\|\hat{\Sigma} - \Sigma^*\|_{\max} \geq t\sigma^2] \leq 3p^2 e^{-nt^2/8}$$

Concentration Inequalities: Sample Covariance (W., Lemma 4.3)

Suppose \mathbf{Z} is mean-zero Gaussian vector with variance Σ^* :

$$[\sigma^2 = \max(\Sigma_{ii}^*)]$$

(a) If Z is real:

$$P[\|\hat{\Sigma} - \Sigma^*\|_{\max} \geq t\sigma^2] \leq 3p^2 e^{-nt^2/8}$$

(b) If Z is proper and complex :

$$P[\|\hat{\Sigma} - \Sigma^*\|_{\max} \geq t\sigma^2] \leq 3p^2 e^{-nt^2/4}$$

Better concentration for proper complex \mathbf{Z} than real \mathbf{Z} !

Intuition: $\Re(\mathbf{Z}) \perp \Im(\mathbf{Z}) \implies$ double sample size

Concentration Inequalities: Sample Covariance (W., Lemma 4.3)

Suppose \mathbf{Z} is mean-zero Gaussian vector with variance Σ^* :
[$\sigma^2 = \max(\Sigma_{ii}^*)$]

(a) If Z is real:

$$P[\|\hat{\Sigma} - \Sigma^*\|_{\max} \geq t\sigma^2] \leq 3p^2 e^{-nt^2/8}$$

(b) If Z is proper and complex :

$$P[\|\hat{\Sigma} - \Sigma^*\|_{\max} \geq t\sigma^2] \leq 3p^2 e^{-nt^2/4}$$

(c) If Z is complex:

$$P[\|\hat{\Sigma} - \Sigma^*\| \geq t\sigma^2] \leq 3p^2 e^{-nt^2/64}$$

Better concentration for proper complex \mathbf{Z} than real \mathbf{Z} !

Intuition: $\Re(\mathbf{Z}) \perp \Im(\mathbf{Z}) \implies$ double sample size

The Complex Lasso

Real LASSO: (Tibshirani, 1996; Chen *et al.*, 1998)

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_1$$

ℓ_1 -norm penalty $\implies \hat{\beta}$ will be *sparse* (have exact zeros):

- Compressed Sensing: can estimate β^* well even with $p \ll n$ elements
- Automatic Variable Selection: can guess the exact zeros in β^* so long as \mathbf{X} is not “too correlated”

The Complex Lasso

Real LASSO: (Tibshirani, 1996; Chen *et al.*, 1998)

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_1$$

ℓ_1 -norm penalty $\implies \hat{\beta}$ will be *sparse* (have exact zeros):

- Compressed Sensing: can estimate β^* well even with $p \ll n$ elements
- Automatic Variable Selection: can guess the exact zeros in β^* so long as \mathbf{X} is not “too correlated”

Rich theoretical literature Fu and Knight (2000), Greenshtein and Ritov (2004), Zhao and Yu (2006), Bickel *et al.* (2009), Zhang and Huang (2008), Bunea *et al.* (2007), Meinshausen and Yu (2009), and van de Geer and Bühlmann (2009) *etc.*

Results all translate to the complex lasso (CLASSO)!

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{C}^p} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_1$$

W., Theorem 3.3: Under standard assumptions (Wainwright, 2009), the Complex-Lasso (CLASSO) is model selection consistent with probability

(a) $\geq 1 - 2 \exp\{-(\tau - 2)/2 \log(p - s)\}$ for real ϵ (real \mathbf{X}, \mathbf{y})

$$\epsilon \stackrel{\text{iid}}{\sim} \text{subG}(0, \sigma^2) \quad \lambda_{\min}(\mathbf{X}_S^H \mathbf{X}_S / n) \geq c > 0 \quad \max_{j \in S^c} \|(\mathbf{X}_S^H \mathbf{X}_S)^{-1} \mathbf{X}_S^H \mathbf{x}_j\|_1 \leq 1 - \gamma$$

The Complex Lasso

W., Theorem 3.3: Under standard assumptions (Wainwright, 2009), the Complex-Lasso (CLASSO) is model selection consistent with probability

(a) $\geq 1 - 2 \exp\{-(\tau - 2)/2 \log(p - s)\}$ for real ϵ (real \mathbf{X}, \mathbf{y})

(b) $\geq 1 - 4 \exp\{-(\tau - 2)/2 \log(p - s)\}$ for real ϵ (complex \mathbf{X}, \mathbf{y})

$$\epsilon \stackrel{\text{i.i.d.}}{\sim} \text{subG}(0, \sigma^2) \quad \lambda_{\min}(\mathbf{X}_S^H \mathbf{X}_S / n) \geq c > 0 \quad \max_{j \in S^c} \|(\mathbf{X}_S^H \mathbf{X}_S)^{-1} \mathbf{X}_S^H \mathbf{x}_j\|_1 \leq 1 - \gamma$$

W., Theorem 3.3: Under standard assumptions (Wainwright, 2009), the Complex-Lasso (CLASSO) is model selection consistent with probability

(a) $\geq 1 - 2 \exp\{-(\tau - 2)/2 \log(p - s)\}$ for real ϵ (real \mathbf{X}, \mathbf{y})

(b) $\geq 1 - 4 \exp\{-(\tau - 2)/2 \log(p - s)\}$ for real ϵ (complex \mathbf{X}, \mathbf{y})

(c) $\geq 1 - 4 \exp\{-(\tau - 16)/16 \log(p - s)\}$ for complex ϵ

$$\epsilon \stackrel{\text{iid}}{\sim} \text{subG}(0, \sigma^2) \quad \lambda_{\min}(\mathbf{X}_S^H \mathbf{X}_S / n) \geq c > 0 \quad \max_{j \in S^c} \|(\mathbf{X}_S^H \mathbf{X}_S)^{-1} \mathbf{X}_S^H \mathbf{x}_j\|_1 \leq 1 - \gamma$$

W., Theorem 3.3: Under standard assumptions (Wainwright, 2009), the Complex-Lasso (CLASSO) is model selection consistent with probability

(a) $\geq 1 - 2 \exp\{-(\tau - 2)/2 \log(p - s)\}$ for real ϵ (real \mathbf{X}, \mathbf{y})

(b) $\geq 1 - 4 \exp\{-(\tau - 2)/2 \log(p - s)\}$ for real ϵ (complex \mathbf{X}, \mathbf{y})

(c) $\geq 1 - 4 \exp\{-(\tau - 16)/16 \log(p - s)\}$ for complex ϵ

(d) $\geq 1 - 8 \exp\{-(\tau - 2)/2 \log(p - s)\}$ for proper complex ϵ

$$\epsilon \stackrel{\text{iid}}{\sim} \text{subG}(0, \sigma^2) \quad \lambda_{\min}(\mathbf{X}_S^H \mathbf{X}_S / n) \geq c > 0 \quad \max_{j \in S^c} \|(\mathbf{X}_S^H \mathbf{X}_S)^{-1} \mathbf{X}_S^H \mathbf{x}_j\|_1 \leq 1 - \gamma$$

W., Theorem 3.3: Under standard assumptions (Wainwright, 2009), the Complex-Lasso (CLASSO) is model selection consistent with probability

(a) $\geq 1 - 2 \exp\{-(\tau - 2)/2 \log(p - s)\}$ for real ϵ (real \mathbf{X}, \mathbf{y})

(b) $\geq 1 - 4 \exp\{-(\tau - 2)/2 \log(p - s)\}$ for real ϵ (complex \mathbf{X}, \mathbf{y})

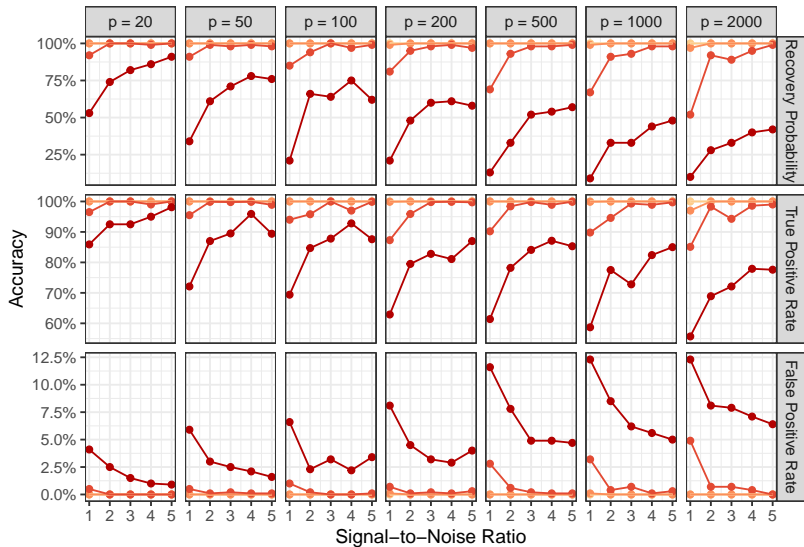
(c) $\geq 1 - 4 \exp\{-(\tau - 16)/16 \log(p - s)\}$ for complex ϵ

(d) $\geq 1 - 8 \exp\{-(\tau - 2)/2 \log(p - s)\}$ for proper complex ϵ

$$\epsilon \stackrel{\text{iid}}{\sim} \text{subG}(0, \sigma^2) \quad \lambda_{\min}(\mathbf{X}_S^H \mathbf{X}_S / n) \geq c > 0 \quad \max_{j \in S^c} \|(\mathbf{X}_S^H \mathbf{X}_S)^{-1} \mathbf{X}_S^H \mathbf{X}_j\|_1 \leq 1 - \gamma$$

First precise finite sample results for CLASSO: previously studied by Yang and Zhang (2011), Maleki *et al.* (2013), and Mechlenbrauker *et al.* (2017)

Model Selection Consistency of C_Lasso



p 0 0.2 0.4 0.6 0.8

The Complex Graphical Lasso

Suppose \mathbf{Z} is drawn from a p -variate complex Gaussian with precision matrix $\Theta^* = (\Sigma^*)^{-1}$. CGLASSO gives a sparse estimate of Θ^* :

$$\hat{\Theta} = \arg \min_{\Theta \in \mathbb{C}_{\geq 0}^{p \times p}} -\log \det \Theta + \text{Tr}(\hat{\Sigma} \Theta) + \lambda \|\Theta\|_{1, \text{off-diag}}$$

where $\hat{\Sigma}$ is the sample covariance (Yuan and Lin, 2007; Friedman *et al.*, 2008; Ravikumar *et al.*, 2011).

The Complex Graphical Lasso

Suppose \mathbf{Z} is drawn from a p -variate complex Gaussian with precision matrix $\Theta^* = (\Sigma^*)^{-1}$. CGLASSO gives a sparse estimate of Θ^* :

$$\hat{\Theta} = \arg \min_{\Theta \in \mathbb{C}_{\geq 0}^{p \times p}} -\log \det \Theta + \text{Tr}(\hat{\Sigma} \Theta) + \lambda \|\Theta\|_{1, \text{off-diag}}$$

where $\hat{\Sigma}$ is the sample covariance (Yuan and Lin, 2007; Friedman *et al.*, 2008; Ravikumar *et al.*, 2011).

W. Theorem 4.1: If \mathbf{Z} is a proper complex Gaussian satisfying standard assumptions, then the CGLASSO recovers the true sparsity pattern of Θ^* with probability

$$\geq 1 - 3p^2 \exp \left\{ -\frac{\log p}{64\sigma^2} - \frac{\delta n}{256\sigma^2} \right\}$$

The Complex Graphical Lasso

Suppose \mathbf{Z} is drawn from a p -variate complex Gaussian with precision matrix $\Theta^* = (\Sigma^*)^{-1}$. CGLASSO gives a sparse estimate of Θ^* :

$$\hat{\Theta} = \arg \min_{\Theta \in \mathbb{C}_{\geq 0}^{p \times p}} -\log \det \Theta + \text{Tr}(\hat{\Sigma}\Theta) + \lambda \|\Theta\|_{1, \text{off-diag}}$$

where $\hat{\Sigma}$ is the sample covariance (Yuan and Lin, 2007; Friedman *et al.*, 2008; Ravikumar *et al.*, 2011).

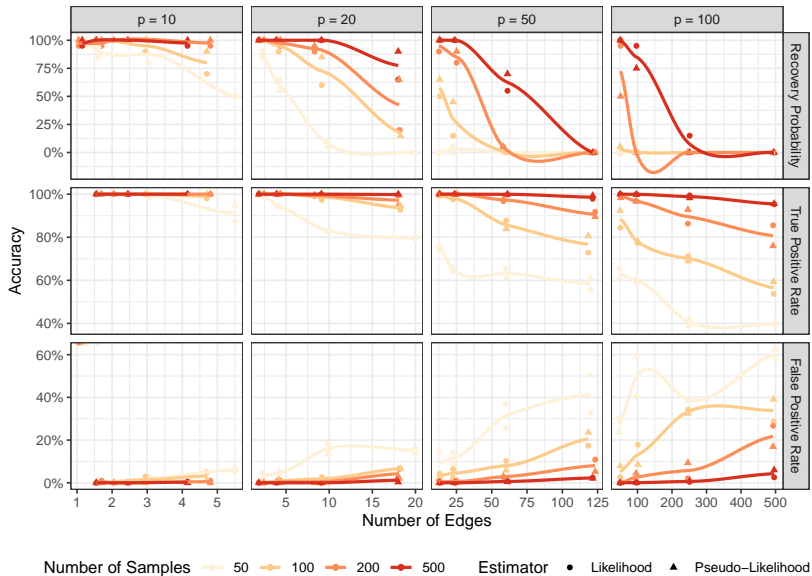
W. Theorem 4.1: If \mathbf{Z} is a proper complex Gaussian satisfying standard assumptions, then the CGLASSO recovers the true sparsity pattern of Θ^* with probability

$$\geq 1 - 3p^2 \exp \left\{ -\frac{\log p}{64\sigma^2} - \frac{\delta n}{256\sigma^2} \right\}$$

Similar results for neighborhood selection (regularized pseudo-likelihood) (Meinshausen and Bühlmann, 2006)

First theoretical results for CGLASSO: Tugnait (2018, 2019a, 2019b) gave experimental results and algorithms

Model Selection Consistency of CGLasso



Improper and Dependent Observations

In addition to incoherence of Θ^* , CGLASSO requires only that

$$\|\hat{\Sigma} - \Sigma^*\|_{\max} < C\sigma\sqrt{\frac{\log p}{n}} \quad \text{for fixed } C$$

Improper and Dependent Observations

In addition to incoherence of Θ^* , CGLASSO requires only that

$$\|\hat{\Sigma} - \Sigma^*\|_{\max} < C\sigma\sqrt{\frac{\log p}{n}} \quad \text{for fixed } C$$

This condition can be satisfied under weaker assumptions,

- Improper \mathbf{Z} : consistent with same rate, worse constants
- Time series spectrum (Fiecas *et al.*, 2019): $\mathcal{O}(p^2/n)$

Improper and Dependent Observations

In addition to incoherence of Θ^* , CGLASSO requires only that

$$\|\hat{\Sigma} - \Sigma^*\|_{\max} < C\sigma\sqrt{\frac{\log p}{n}} \quad \text{for fixed } C$$

This condition can be satisfied under weaker assumptions,

- Improper \mathbf{Z} : consistent with same rate, worse constants
- Time series spectrum (Fiecas *et al.*, 2019): $\mathcal{O}(p^2/n)$

W. Theorem 4.1 + Fiecas *et al.* (2019) + Dahlhaus (2000)

Suppose $\mathcal{Z} = \{\mathbf{Z}_t\}_{t=1}^T$ is a stationary Gaussian linear p -dimensional time series with spectrum $\Gamma(k) \in \mathbb{C}_{\geq 0}^{p \times p}$, such that $\Gamma^{-1}(k)$ satisfies the incoherence conditions at all k . Let $\hat{\Gamma}(k)$ be the sample averaged periodogram based on T observations. Then the graphical model $\hat{\mathcal{G}} = (\mathcal{V}, \hat{\mathcal{E}})$ with

$$(i, j) \notin \hat{\mathcal{E}} \Leftrightarrow \hat{\Theta}_{ij}^{(k)} = 0 \quad \text{for } \hat{\Theta}^{(k)} = \text{CGLASSO}(\hat{\Gamma}(k)) \quad \text{for all } k < T/2$$

correctly estimates the conditional independence structures in \mathcal{Z} at all lags with probability $\geq 1 - Cp^2/T$ for T sufficiently large.

Developed Wirtinger convex analysis and applied it to CLASSO and CGLASSO estimators:

Developed Wirtinger convex analysis and applied it to CLASSO and CGLASSO estimators:

- Foundational Optimization and Statistical Theory for *Complex Machine Learning*

Developed Wirtinger convex analysis and applied it to CLASSO and CGLASSO estimators:

- Foundational Optimization and Statistical Theory for *Complex Machine Learning*
- First Finite-Sample Results for Complex M -Estimation

Developed Wirtinger convex analysis and applied it to CLASSO and CGLASSO estimators:

- Foundational Optimization and Statistical Theory for *Complex Machine Learning*
- First Finite-Sample Results for Complex M -Estimation
- First Statistical Study of Improper Gaussian Distributions

Developed Wirtinger convex analysis and applied it to CLASSO and CGLASSO estimators:

- Foundational Optimization and Statistical Theory for *Complex Machine Learning*
- First Finite-Sample Results for Complex M -Estimation
- First Statistical Study of Improper Gaussian Distributions
- Exciting Implications for Multivariate Time Series

Conclusion & Discussion

Novel Methodologies with Sophisticated Structure: tensors, low-rank estimation, misaligned & high-frequency time series, complex-variates

Novel Methodologies with Sophisticated Structure: tensors, low-rank estimation, misaligned & high-frequency time series, complex-variates

Efficient algorithms + strong convergence guarantees built on robust convex analysis

Novel Methodologies with Sophisticated Structure: tensors, low-rank estimation, misaligned & high-frequency time series, complex-variates

Efficient algorithms + strong convergence guarantees built on robust convex analysis

Ability to flexibly incorporate problem structure into estimation methodology directly

Novel Methodologies with Sophisticated Structure: tensors, low-rank estimation, misaligned & high-frequency time series, complex-variates

Efficient algorithms + strong convergence guarantees built on robust convex analysis

Ability to flexibly incorporate problem structure into estimation methodology directly

Extensions of classical M -estimation and convex analysis to Wirtinger functions

Acknowledgements

M G Z L U K E M I N J I E F M
M A J F G G L I B B Y U X R I
A A R O A N D E R S E N N E T
K U T G H I S G E O R G E D C
E I G T A A K R I S T E N Y H
A I M U T R N J U L I A H G T
N U G A S O E M E W V T G U E
D U T A T T M T A T A K B R R
R G I U U T I R Z K N L E D E
E R A K Y T E N I L N U N C N
A F N K O V A O E O U I H O C
B P Y M E W N M N H C S A F E
N F I N B T A N E J C M A E G
I S E Y O O A L A J I I M S P
E G F O B C D R J O N E S S E

W. “Splitting Methods for Convex Bi-Clustering and Co-Clustering.”
DSW 2019

W. “Multi-Rank Sparse and Functional PCA: Manifold Optimization
and Iterative Deflation Techniques.” *CAMSAP 2019*

W, Y. Han, and K. B. Ensor. “Multivariate Modeling of Natural Gas
Spot Trading Hubs Incorporating Futures Market Realized Volatility.”
Under revision at JASA: ACS

2020 ASA Business & Economic Statistics Student Paper Award

W. “High-Dimensional Spectral Graph Estimation via Wirtinger
Graphical Models for Complex-Valued Data.” *In preparation for AoS*

Thank you!