## Treelet Covariance Smoothers Estimation of Genetic Parameters

#### B. Draves<sup>1</sup>

<sup>1</sup>Department of Mathematics Lafayette College

Advisor: T. Gaugler

Moravian College, 2017

B. Draves (Lafayette College)

Treelet Covariance Smoothers

Moravian College, 2017 1 / 11

#### • Principle Component Analysis

B. Draves (Lafayette College)

Treelet Covariance Smoothers

Moravian College, 2017 2 / 11

3

<ロ> (日) (日) (日) (日) (日)

- Principle Component Analysis
- Motivation in Statistical Genetics

3

∃ → ( ∃ →

< 67 ▶

- Principle Component Analysis
- Motivation in Statistical Genetics
- Existing Methods

3

∃ → ( ∃ →

< 一型

- Principle Component Analysis
- Motivation in Statistical Genetics
- Existing Methods
- New Methods & Results

3

∃ ► < ∃ ►</p>

• PCA looks to reduce the number of input variables  $X_1, X_2, ..., X_k$  (k large) to  $k_0 < k$  variables

3

(日) (周) (三) (三)

- PCA looks to reduce the number of input variables X<sub>1</sub>, X<sub>2</sub>,..., X<sub>k</sub> (k large) to k<sub>0</sub> < k variables</li>
- Reduce to fewer inputs while preserving as much variability in the data as possible

イロト 不得下 イヨト イヨト

- PCA looks to reduce the number of input variables X<sub>1</sub>, X<sub>2</sub>,..., X<sub>k</sub> (k large) to k<sub>0</sub> < k variables</li>
- Reduce to fewer inputs while preserving as much variability in the data as possible



- PCA looks to reduce the number of input variables X<sub>1</sub>, X<sub>2</sub>,..., X<sub>k</sub> (k large) to k<sub>0</sub> < k variables</li>
- Reduce to fewer inputs while preserving as much variability in the data as possible



- PCA looks to reduce the number of input variables X<sub>1</sub>, X<sub>2</sub>,..., X<sub>k</sub> (k large) to k<sub>0</sub> < k variables</li>
- Reduce to fewer inputs while preserving as much variability in the data as possible



 Individual's genetic material can be described by a panel of genotyped SNPs

3

∃ ► < ∃ ►</p>

- Individual's genetic material can be described by a panel of genotyped SNPs
- Using this genetic information, an estimate of the relationship matrix, *A*, can be calculated

∃ ► < ∃ ►</p>

- Individual's genetic material can be described by a panel of genotyped SNPs
- Using this genetic information, an estimate of the relationship matrix, *A*, can be calculated
- Genetically inferred relationship matrices are typically very noisy

E 6 4 E 6

- Individual's genetic material can be described by a panel of genotyped SNPs
- Using this genetic information, an estimate of the relationship matrix, *A*, can be calculated
- Genetically inferred relationship matrices are typically very noisy
- Idea: use PCA to set the data on top of a more "natural" basis

∃ → ( ∃ →

#### • Goal: preserve local structure of data

- Goal: preserve local structure of data
- Iteratively change basis via PCA to preserve the variability between two most closely related individuals

- Goal: preserve local structure of data
- Iteratively change basis via PCA to preserve the variability between two most closely related individuals
- Repeat this process until all individuals are processed

- Goal: preserve local structure of data
- Iteratively change basis via PCA to preserve the variability between two most closely related individuals
- Repeat this process until all individuals are processed
- Once all individuals are processed, enforce sparsity to improve the estimator





P3











◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?















◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへ⊙



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへ⊙











Moravian College, 2017 6 / 11

• Why do we require merging into one cluster?

3

- Why do we require merging into one cluster?
- Idea: Stop merging variables to utilize familiar blocks

4 A N

- Why do we require merging into one cluster?
- Idea: Stop merging variables to utilize familiar blocks
- Smooth by projecting data onto sum variables
  - Treelet Covariance Blocking (TCB)

A B F A B F

< A >

- Why do we require merging into one cluster?
- Idea: Stop merging variables to utilize familiar blocks
- Smooth by projecting data onto sum variables
  - Treelet Covariance Blocking (TCB)
- Further enforce sparsity by thresholding estimates
  - Treelet Covariance Blocked Smoothing (TCBS)

4 3 5 4 3 5 5

• At each level of the tree there are basis vectors  $\{v_1^{(\ell)},\ldots,v_N^{(\ell)}\}$ 

3

(日) (同) (三) (三)

At each level of the tree there are basis vectors {v<sub>1</sub><sup>(ℓ)</sup>,...,v<sub>N</sub><sup>(ℓ)</sup>}
We can then write our estimate of A, Σ̂, by

$$\widehat{\boldsymbol{\Sigma}} = \sum_{i \in \widehat{\boldsymbol{S}}_{\ell}} \widehat{\gamma}_{i,i}^{(\ell)} \widehat{\boldsymbol{v}}_{i}^{(\ell)} \left( \widehat{\boldsymbol{v}}_{i}^{(\ell)} \right)^{t} + \sum_{i,j \in \widehat{\boldsymbol{S}}_{\ell}, i \neq j} \widehat{\gamma}_{i,j}^{(\ell)} \widehat{\boldsymbol{v}}_{i}^{(\ell)} \left( \widehat{\boldsymbol{v}}_{j}^{(\ell)} \right)^{t}$$

3

( )

At each level of the tree there are basis vectors {v<sub>1</sub><sup>(ℓ)</sup>,...,v<sub>N</sub><sup>(ℓ)</sup>}
We can then write our estimate of A, Σ̂, by

$$\widehat{\boldsymbol{\Sigma}} = \sum_{i \in \widehat{S}_{\ell}} \widehat{\gamma}_{i,i}^{(\ell)} \widehat{\mathbf{v}}_{i}^{(\ell)} \left( \widehat{\mathbf{v}}_{i}^{(\ell)} \right)^{t} + \sum_{i,j \in \widehat{S}_{\ell}, i \neq j} \widehat{\gamma}_{i,j}^{(\ell)} \widehat{\mathbf{v}}_{i}^{(\ell)} \left( \widehat{\mathbf{v}}_{j}^{(\ell)} \right)^{t}$$

• Here  $\gamma_{i,j}^{(\ell)}$  are the variance - covariance estimates of transformed relationship variables

At each level of the tree there are basis vectors {v<sub>1</sub><sup>(ℓ)</sup>,...,v<sub>N</sub><sup>(ℓ)</sup>}
We can then write our estimate of A, Σ̂, by

$$\widehat{\boldsymbol{\Sigma}} = \sum_{i \in \widehat{S}_{\ell}} \widehat{\gamma}_{i,i}^{(\ell)} \widehat{\boldsymbol{v}}_{i}^{(\ell)} \left( \widehat{\boldsymbol{v}}_{i}^{(\ell)} \right)^{t} + \sum_{i,j \in \widehat{S}_{\ell}, i \neq j} \widehat{\gamma}_{i,j}^{(\ell)} \widehat{\boldsymbol{v}}_{i}^{(\ell)} \left( \widehat{\boldsymbol{v}}_{j}^{(\ell)} \right)^{t}$$

- Here  $\gamma_{i,j}^{(\ell)}$  are the variance covariance estimates of transformed relationship variables
- We can enforce sparsity by further thresholding these  $\gamma_{i,j}^{(\ell)}$  values

#### Simulation Results



B. Draves (Lafayette College)

Treelet Covariance Smoothers

Moravian College, 2017 9 / 11

#### Conclusion

• Treelet Covariance Smoothers is a class of methods which improve the estimation of distant relationships

э

## Conclusion

- Treelet Covariance Smoothers is a class of methods which improve the estimation of distant relationships
- Estimating relationships is the centerpiece of successful estimation of other genetic parameters such as heritability

## Conclusion

- Treelet Covariance Smoothers is a class of methods which improve the estimation of distant relationships
- Estimating relationships is the centerpiece of successful estimation of other genetic parameters such as heritability
- Understanding of genetic basis of heritability can lead to better treatment

A B K A B K

#### Thanks for listening

Questions? Comments?

3

イロト イポト イヨト イヨト