

Analyzing High-Dimensional Data



Challenges in analysis of 'omics data

- New technologies now make it possible to capture huge amounts of data on experimental units.
 - **Genomics**
 - Can sequence and entire genome
 - Sequence all RNA in a tissues
 - **Phenomics**
 - Collect multi-spectral or hyper spectral data on plants throughout the growing season
- These technologies provide opportunities to increase understanding of complex phenotypes but can be challenging to model.
 - The 'curse' of dimensionality



R Exercise – $N \ll P$ Challenges



Methods for dealing with high-dimensional data

- **Penalized Methods** apply some penalty to the solutions to avoid overfitting the data when there are a large number of explanatory variables relative to the number of independent observations.
 - Various penalized methods differ in the type of penalty applied
 - Penalize methods yield solutions that are not unbiased – at least in the same sense that OLS solutions are unbiased
- **Dimension Reduction Methods** solve the issue of dimensionality by reducing the number of explanatory variables used to fit the regression.
 - Reducing the dimensions of the data inevitably results in the loss of some information



Penalized methods – Ridge Regression

- Ridge Regression simply adds some constant value to the diagonal of $\mathbf{X}'\mathbf{X}$

OLS solutions $\hat{\mathbf{b}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$

Ridge Regression
Solutions solutions

$$\hat{\mathbf{b}}^* = (\mathbf{X}'\mathbf{X} + \alpha\mathbf{I})^{-1}\mathbf{X}'\mathbf{y}$$



Genomic Prediction using Mixed Models

$$u = \sum w_i \beta_i \quad W = M - P \quad \alpha = \frac{\sigma_e^2}{\sigma_\beta^2}$$

$$\begin{bmatrix} X'X & X'W \\ W'X & W'W + \alpha I \end{bmatrix} \begin{bmatrix} b \\ \beta \end{bmatrix} = \begin{bmatrix} X'y \\ W'y \end{bmatrix}$$



R – Exercise – Ridge Regression



Penalized Methods - LASSO

OLS minimizes the function: $(y - \mathbf{X}b)'(y - \mathbf{X}b)$

LASSO minimizes the function: $(y - \mathbf{X}b)'(y - \mathbf{X}b) + \lambda ||b||$

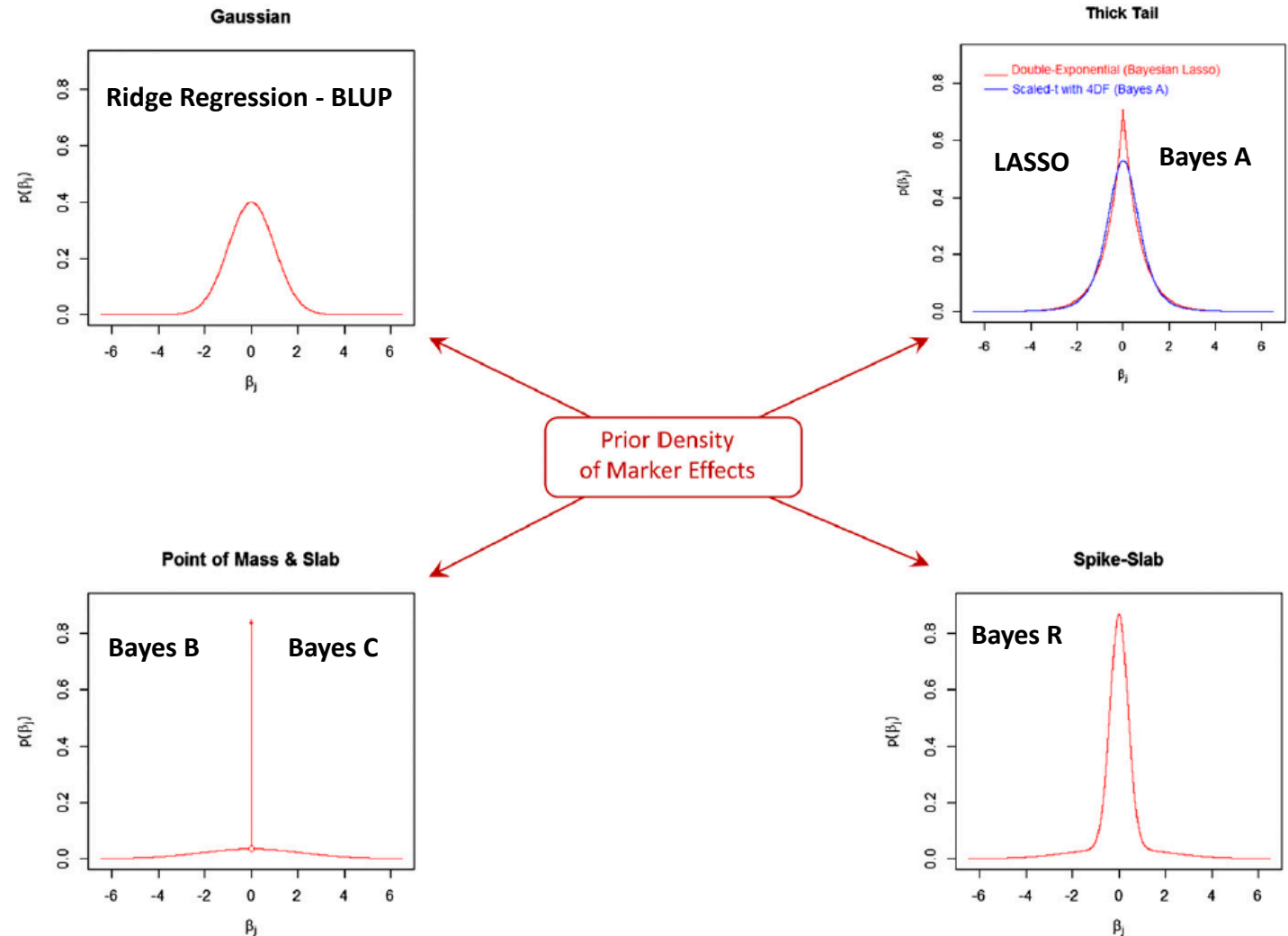
- LASSO optimizes the absolute value of elements of b with the variation explained.
- Only elements of b that explain a large proportion of variance are allowed to have a large absolute value



Penalized Methods

There are many models for dealing with over-parameterized data using penalized methods.

These methods vary in the underlying assumptions of the data distribution, which determines the types of penalties applied.



Dimension Reduction – Feature Selection

- Filter methods
 - Markers scored and ranked
- Wrapper methods
 - Often deploy optimized searching algorithms (Genetic Algorithm, Ant Colony Algorithm, Simulated Annealing)
 - Initialized based on some prior information
 - Select features, perform cross validation, and update based on predictive performance
- Work best when there are a small number of markers that explain a large proportion of variation.
 - Commonly used in disease diagnostics
 - Not as effective for complex, additive traits traits often modeling in plant and animal breeding.



Wrapper Approach using Swarm Intelligence



Real Ant Colony

m ants searching for best rout to food source

Communicate through a chemical pheromone trail

Pheromone level changes each trip based on the time it takes to reach food source

Artificial Ant colony

m ants searching for best subset of genes

Communicate through a PDF:

$$P_{mc}(t) = \frac{(\tau_{mc}(t))^{\alpha} \eta_{mc}^{\beta}}{\sum_{m=1}^{nf} (\tau_{mc}(t))^{\alpha} \eta_{mc}^{\beta}}$$

Pheromone level changes each iteration based on the prediction accuracy of selected genes

$$\tau_m(t+1) = (1-\rho) * \tau_m(t) + \Delta\tau_m(t)$$



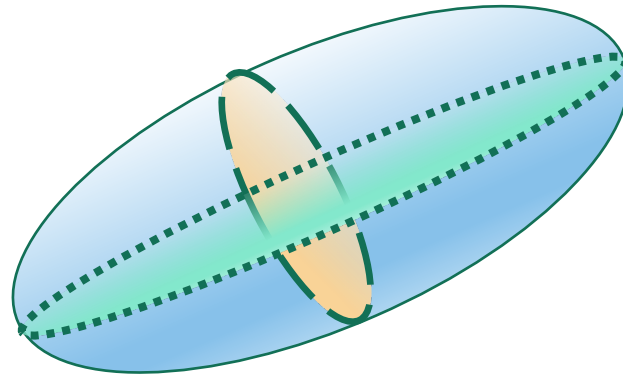
Dimension Reduction - Principle Component Regression.

- The goal of PCR is to decompose \mathbf{X} in m orthogonal vectors (we will call this matrix \mathbf{T}). If there are correlations in the original vectors (columns of \mathbf{X}) we can often find a small subset of vectors in \mathbf{T} that explain most of the variance in \mathbf{X} .
- We then use this reduced full rank matrix \mathbf{T} to calculate OLS
- We can do this decomposition on \mathbf{X} or $\mathbf{X}'\mathbf{X}$ – here will focus on the spectral decomposition (S.D.) of $\mathbf{X}'\mathbf{X}$



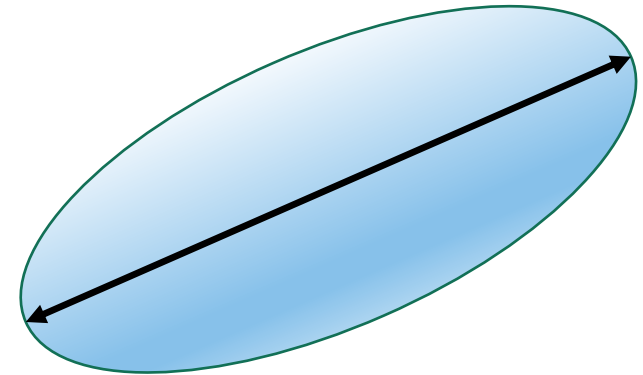
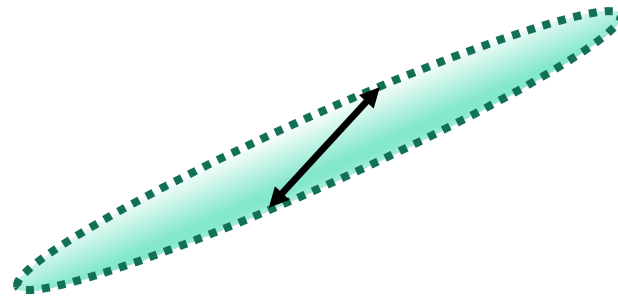
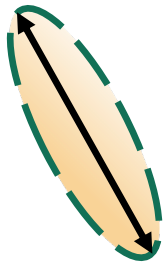
Principle Component Regression.

Orthogonalization – A process by which you take correlated vectors (columns in a matrix) and decompose them into orthogonal components.



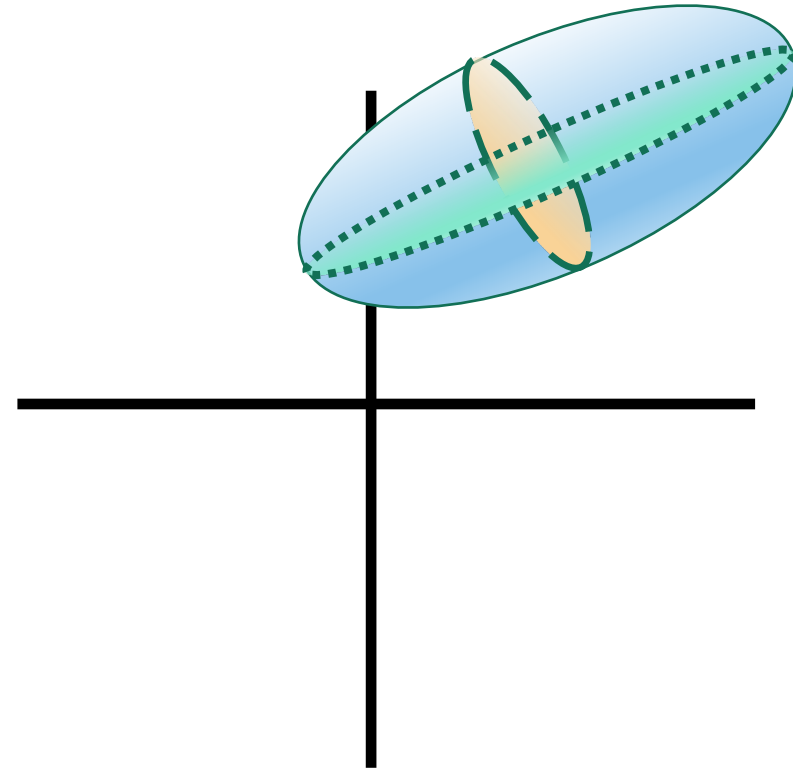
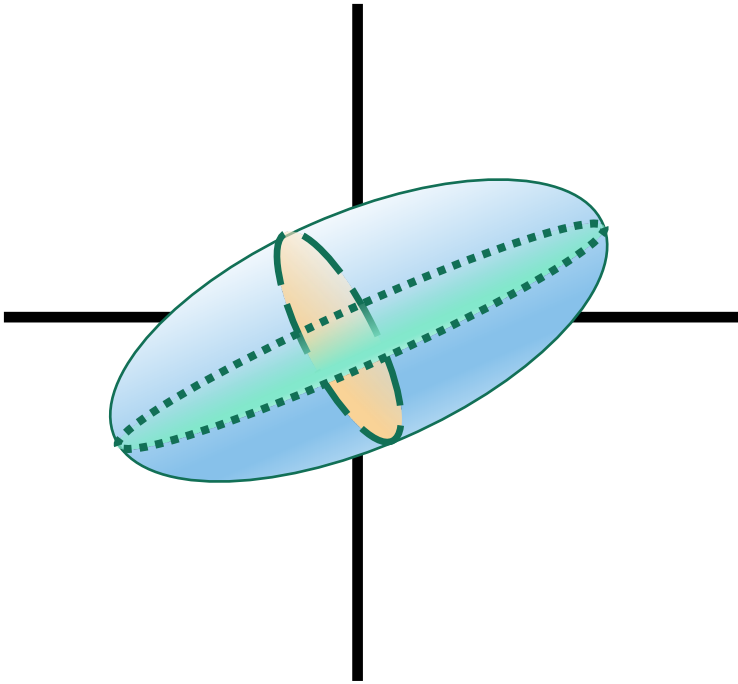
Principle Component Regression.

There are many ways to do orthogonalization – PCR uses an eigenvalue decomposition or spectral decomposition.



Principle Component Regression.

The first step in PCR is to center the matrix.



Spectral Decomposition

\mathbf{X} is a $n \times m$ incidence matrix – n observations and m covariates

$\mathbf{X}'\mathbf{X}$ is a $m \times m$ S.D. $(\mathbf{X}'\mathbf{X}) = \mathbf{U}\mathbf{D}\mathbf{U}' = \mathbf{X}'\mathbf{X}$

\mathbf{U} is a symmetric $m \times m$ matrix of eigen vectors $\mathbf{U}\mathbf{U}' = \mathbf{U}'\mathbf{U} = \mathbf{I}$

\mathbf{D} is a $m \times m$ diagonal matrix of eigen values



Principle Component Regression.

To get \mathbf{T} we need to multiple the matrix \mathbf{X} by some matrix \mathbf{P} to generate a matrix with orthogonal columns such that:

$$\mathbf{T} = \mathbf{XP} \quad \text{and}$$

$\mathbf{T}'\mathbf{T}$ is a diagonal matrix with non-zero values on the diagonal



Some Matrix Algebra ...

$$\mathbf{T}'\mathbf{T} = \mathbf{P}'\mathbf{X}'\mathbf{X}\mathbf{P}$$

$$\text{Given } \rightarrow \mathbf{U}\mathbf{D}\mathbf{U}' = \mathbf{X}'\mathbf{X}$$

$$\mathbf{T}'\mathbf{T} = \mathbf{P}'\mathbf{U}\mathbf{D}\mathbf{U}'\mathbf{P}$$

By setting $\mathbf{P} = \mathbf{U}$

$$\mathbf{T}'\mathbf{T} = \mathbf{P}'\mathbf{U}\mathbf{D}\mathbf{U}'\mathbf{P} = \mathbf{U}'\mathbf{U}\mathbf{D}\mathbf{U}'\mathbf{U} = \mathbf{D}$$



Using T for OLS

OLS using \mathbf{X}

$$\hat{\mathbf{b}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$$

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\mathbf{b}}$$

$$\mathbf{T} = \mathbf{X}\mathbf{U}$$

OLS using \mathbf{T}

$$\hat{\mathbf{b}}^* = (\mathbf{T}'\mathbf{T})^{-1}\mathbf{T}'\mathbf{y} \quad \longrightarrow \quad \widehat{\mathbf{b}}^* = (\mathbf{D})^{-1}\mathbf{T}'\mathbf{y}$$

$$\hat{\mathbf{y}} = \mathbf{T}\widehat{\mathbf{b}}^*$$



R Exercise - PCR



Dimension Reduction – Partial Least Squares

- In PCR we orthogonalize \mathbf{X} to get a matrix \mathbf{T} which contain most of the variation in \mathbf{X} but is full rank.
- This approach is superior to feature selection when there are several highly correlated covariates all explaining some variation in \mathbf{y}
- The draw back is that \mathbf{y} is never considered when decomposing \mathbf{X}
- Ideally, we would decompose \mathbf{X} into orthogonal components that explain the most variation in \mathbf{y} not \mathbf{X} .
- **PLS differs from PCR in that we account for \mathbf{y} when decomposing \mathbf{X} .**



Nonlinear Iterative Partial Least Squares

- NIPALS is an iterative algorithm to decompose a matrix \mathbf{X}
- Going back to the PCR approach we could decompose \mathbf{X} as follows:

Initialize $\mathbf{t} := \mathbf{x}_j$ for some column j in \mathbf{X}

Loop until \mathbf{t} converges

$$\mathbf{p} := \frac{\mathbf{X}'\mathbf{t}}{\|\mathbf{X}'\mathbf{t}\|}$$

$$\mathbf{t} := \mathbf{X}\mathbf{p}$$

Then set $\mathbf{X} := \mathbf{X} - \mathbf{t}\mathbf{p}'$

Repeat until you have all columns in \mathbf{T}



Partial Least Squares

- For PLS we substitute in \mathbf{y} to give the following algorithm:

Loop until \mathbf{t} converges

$$\mathbf{p} := \frac{\mathbf{X}'\mathbf{y}}{\|\mathbf{X}'\mathbf{y}\|}$$

$$\mathbf{t} := \mathbf{X}\mathbf{p}$$

Then set $\mathbf{X} := \mathbf{X} - \mathbf{t}\mathbf{p}'$

Repeat until you have all columns in \mathbf{T}

The above algorithm is for a single response variable – this can be generalized to multiple response variables – see notes on Canvas



Questions



Course Review

