## 6.2.2: The Lasso

- A disadvantage of ridge regression is that it will include all $p$ predictors in the final model.

- The lasso overcomes this. It minimizes

$$\sum_{i=1}^{n}\left(y_i - \beta_0 - \sum_{j=1}^{p}\beta_j x_{ij}\right) + \lambda \sum_{j=1}^{p}|\beta_j|$$

This forces some coefficients to be exactly $0$ when $\lambda$ is large enough. So: the lasso performs variable selection yielding sparse models

### The Variable Selection Property of the Lasso

Has to do with Lagrange multipliers

Use Lasso if you expect that some coefficients are $0$.

## 6.2.3: Selecting the Tuning Parameter
Need to select $\lambda$.

Cross-validation approach:
1) Choose a grid of $\lambda$ values
2) Compute CV error for each $\lambda$
3) Select $\lambda$ where CV error is smallest, call it $\lambda^*$
4) Refit the model on all the data using $\lambda^*$

## 2/19/23 6.3: Dimension Reduction Methods

Day 14  This is a class of approaches that transforms the predictors and then fits least squares using the transformed variables.

Let $Z_1, \ldots, Z_M$ represent $M < p$ <u>linear combinations</u> of our original $p$ predictors,

$$Z_m = \sum_{j=1}^{p} \phi_{jm} X_j, \quad \text{for constants } \phi_{1m}, \ldots, \phi_{pm}, \ m \in [M].$$

Then we can fit the linear regression model (w/ least square)

$$y_i = \theta_0 + \sum_{m=1}^{M} \theta_m z_{im} + \varepsilon_i, \quad i = 1, \ldots, n$$

If $\phi_{1m}, \ldots, \phi_{pm}$ are chosen wisely, then dimension reduction approaches can outperform least squares regression.

This works well when ~~$M < p$~~ $p > n$, especially if you choose $M \ll p$.

All dimension reduction methods follow two steps:
1) Get $Z_1, \ldots, Z_m$
2) Estimate $\theta_1, \ldots, \theta_m$

We review two approaches: <u>Principal Components</u> and <u>Partial Least Squares</u>

6.3.1: <u>Principal Components Regression (PCA)</u>
<u>An Overview of PCA</u>
• PCA is a technique for reducing the dimension of an $n \times p$ data matrix $X$.
• The <u>first principal component</u> direction of the data is that along which the observations vary the most.
   - i.e. Projecting the 100 observations onto this line gives the largest variance possible.

So, (1) center the data, (2) fit ~~original~~ orthogonal components to the data that maximize variance.

This is about simplifying the predictors.

# Principal Component Regression

- Find the first M components, then use these components in a linear regression.

- We assume that the component directions are associated with Y.

- Can do better than least squares linear regression by using most of the information while avoiding overfitting.

- PCR does better if fewer components are required.

- PCR is ~~not~~ a feature extraction method

- "One can think of ridge regression as continuous PCR." Cool

- In PCR, M is ^often chosen by cross validation.

- ~~In PCR~~ Standardize the predictors prior to doing PCR, otherwise a variable having high variance can have an outsized impact on the result.

## 6.3.2: Partial Least Squares (PLS)

- In PCR, the principal components are identified in an unsupervised way, as Y is not used to find the principal components.

- Hence, there is no guarantee that the principal components are associated with the response.

- PLS is a supervised alternative to PCR.

## How do you compute the PLS directions?

1) The first direction $Z_1$ is computed by setting each $\phi_{j1}$ equal to $\beta_j$ from a simple linear regression of $Y$ onto $X_j$. So

$$Z_1 = \sum_{j=1}^{p} \phi_{j1} X_j$$

is ~~the~~ strongly correlated with the response.

Okay, I don't understand this