# Linear Methods for Regression

## Outline

- The simple linear regression model

- Multiple linear regression

- Model selection and shrinkage—the state of the art

## Preliminaries

Data $(x_1, y_1), \ldots (x_N, y_N)$.

$x_i$ is the predictor (regressor, covariate, feature, independent variable)

$y_i$ is the response (dependent variable, outcome)

We denote the *regression function* by

$$\eta(x) = \mathrm{E}\,(Y|x)$$

This is the conditional expectation of $Y$ given $x$.

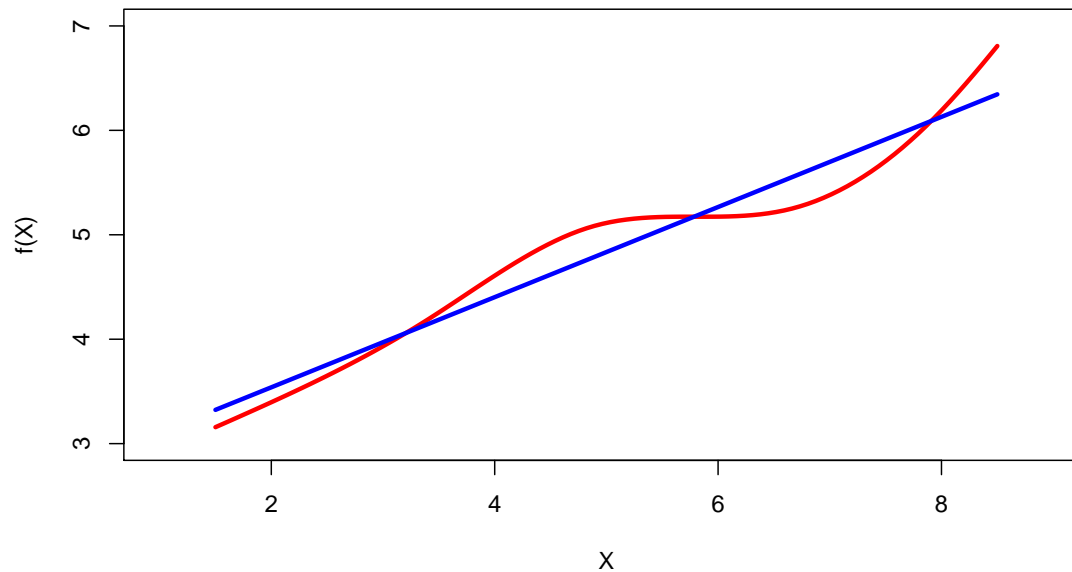The linear regression model assumes a specific linear form for $\eta$

$$\eta(x) = \alpha + \beta x$$

which is usually thought of as an approximation to the truth.

## Linearity assumption?

$$\eta(x) = \alpha + \beta x$$

Almost always thought of as an approximation to the truth. Functions in nature are rarely linear.

## Fitting by least squares

Minimize:

$$\hat{\beta}_0, \hat{\beta} = \mathrm{argmin}_{\beta_0,\beta} \sum_{i=1}^{N} (y_i - \beta_0 - \beta x_i)^2$$
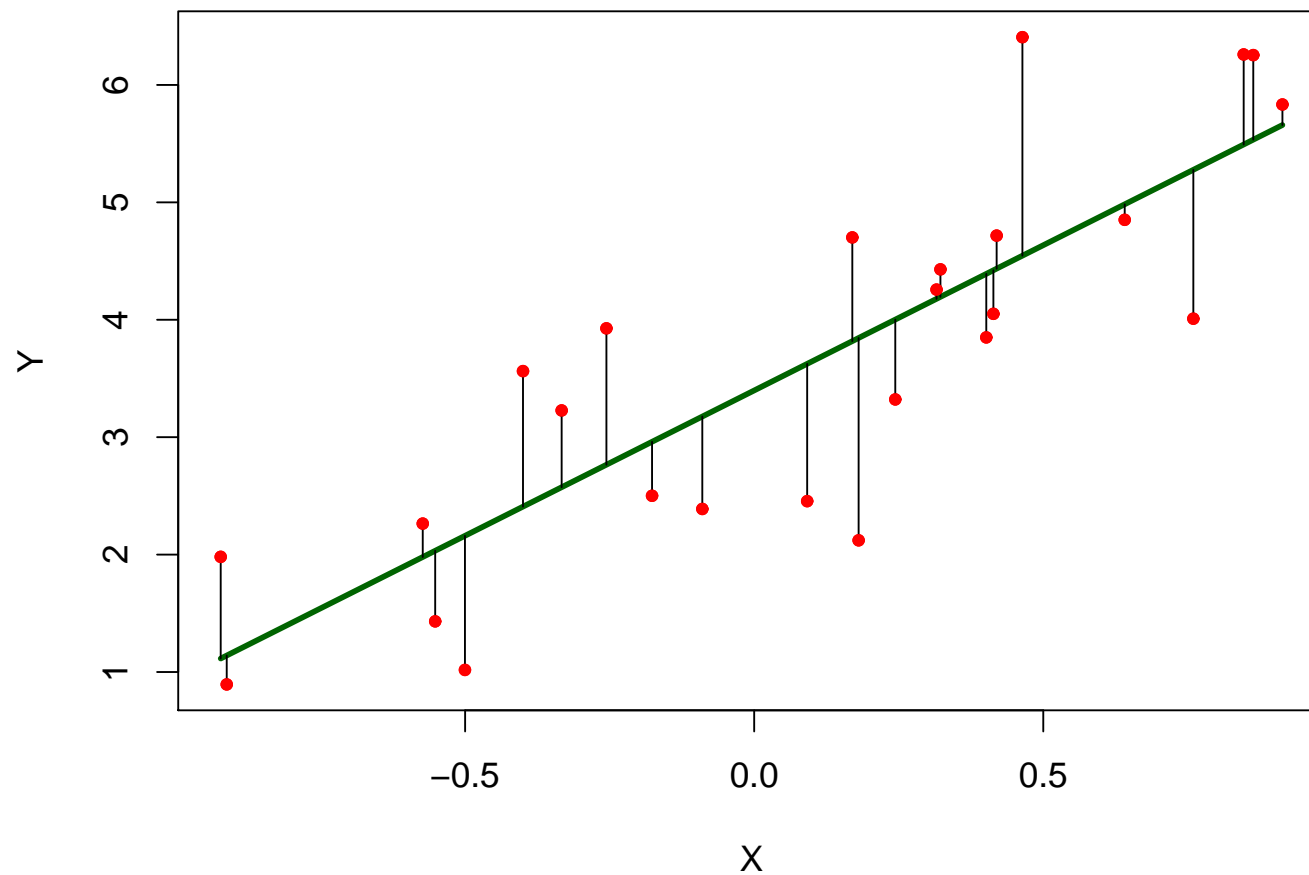
Solutions are

$$\hat{\beta} = \frac{\sum_{j=1}^{N}(x_i - \bar{x})y_i}{\sum_{j=1}^{N}(x_i - \bar{x})^2}$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}\bar{x}$$

$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}x_i$ are called the fitted or predicted values

$r_i = y_i - \hat{\beta}_0 - \hat{\beta}x_i$ are called the residuals

Least squares estimation for linear regression model.

# Standard errors & confidence intervals

We often assume further that

$$y_i = \beta_0 + \beta x_i + \epsilon_i$$

where $\mathrm{E}\,(\epsilon_i) = 0$ and $\mathrm{Var}\,(\epsilon_i) = \sigma^2$. Then

$$\mathrm{se}\,(\hat{\beta}) = \left[ \frac{\sigma^2}{\sum(x_i - \bar{x})^2} \right]^{\frac{1}{2}}$$

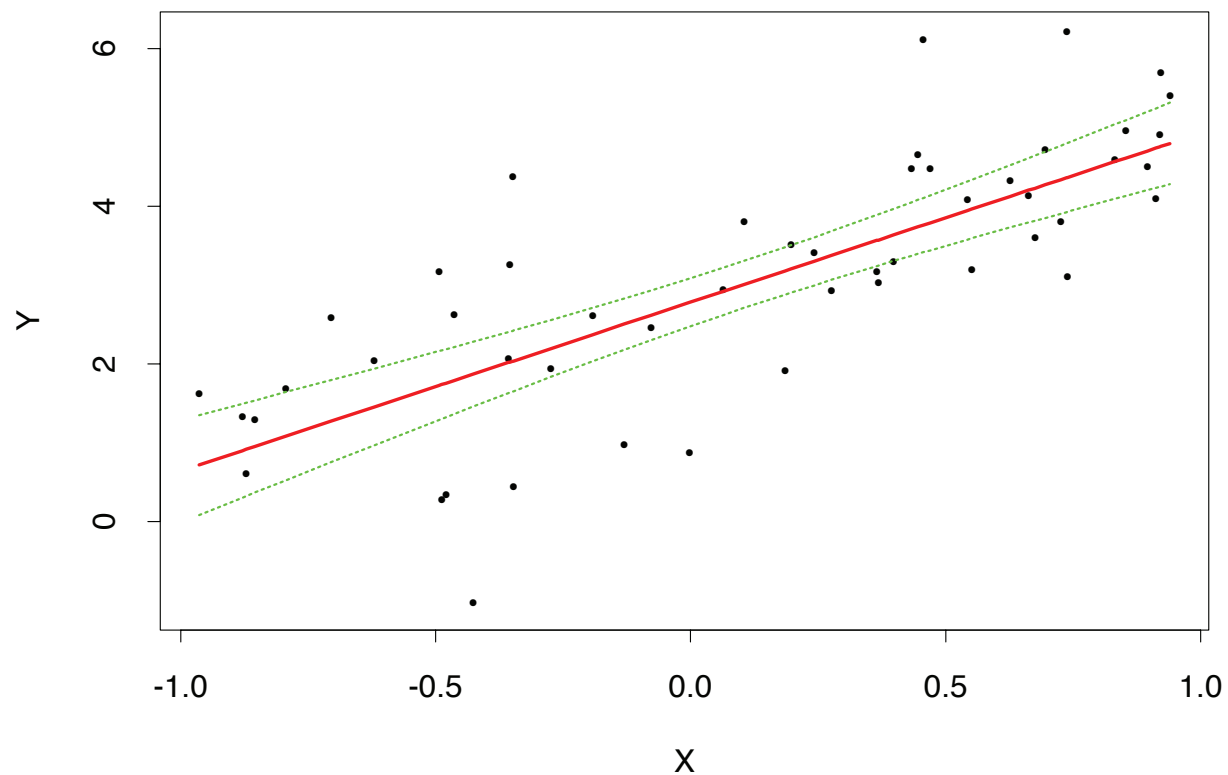Estimate $\sigma^2$ by $\hat{\sigma}^2 = \sum(y_i - \hat{y}_i)^2/(N-2)$.

Under additional assumption of normality for the $\epsilon_i$s, a 95% confidence interval for $\beta$ is: $\hat{\beta} \pm 1.96\hat{\mathrm{se}}(\hat{\beta})$

$$\hat{\mathrm{se}}\,(\hat{\beta}) = \left[ \frac{\hat{\sigma}^2}{\sum(x_i - \bar{x})^2} \right]^{\frac{1}{2}}$$

# Fitted Line and Standard Errors

$$\hat{\eta}(x) \; = \; \hat{\beta}_0 + \hat{\beta}x$$

$$= \; \bar{y} + \hat{\beta}(x - \bar{x})$$

$$\text{se}[\hat{\eta}(x)] \; = \; \left[ \text{var}(\bar{y}) + \text{var}(\hat{\beta})(x - \bar{x})^2 \right]^{\frac{1}{2}}$$

$$= \; \left[ \frac{\sigma^2}{n} + \frac{\sigma^2 (x - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right]^{\frac{1}{2}}$$

Fitted regression line with pointwise standard errors: $\hat{\eta}(x) \pm 2\hat{se}[\hat{\eta}(x)]$.

# Multiple linear regression

Model is

$$f(x_i) = \beta_0 + \sum_{j=1}^{p} x_{ij}\beta_j$$

Equivalently in matrix notation:

$$\mathbf{f} = \mathbf{X}\beta$$

$\mathbf{f}$ is $N$-vector of predicted values

$\mathbf{X}$ is $N \times (p+1)$ matrix of regressors, with ones in the first column

$\beta$ is a $p + 1$-vector of parameters
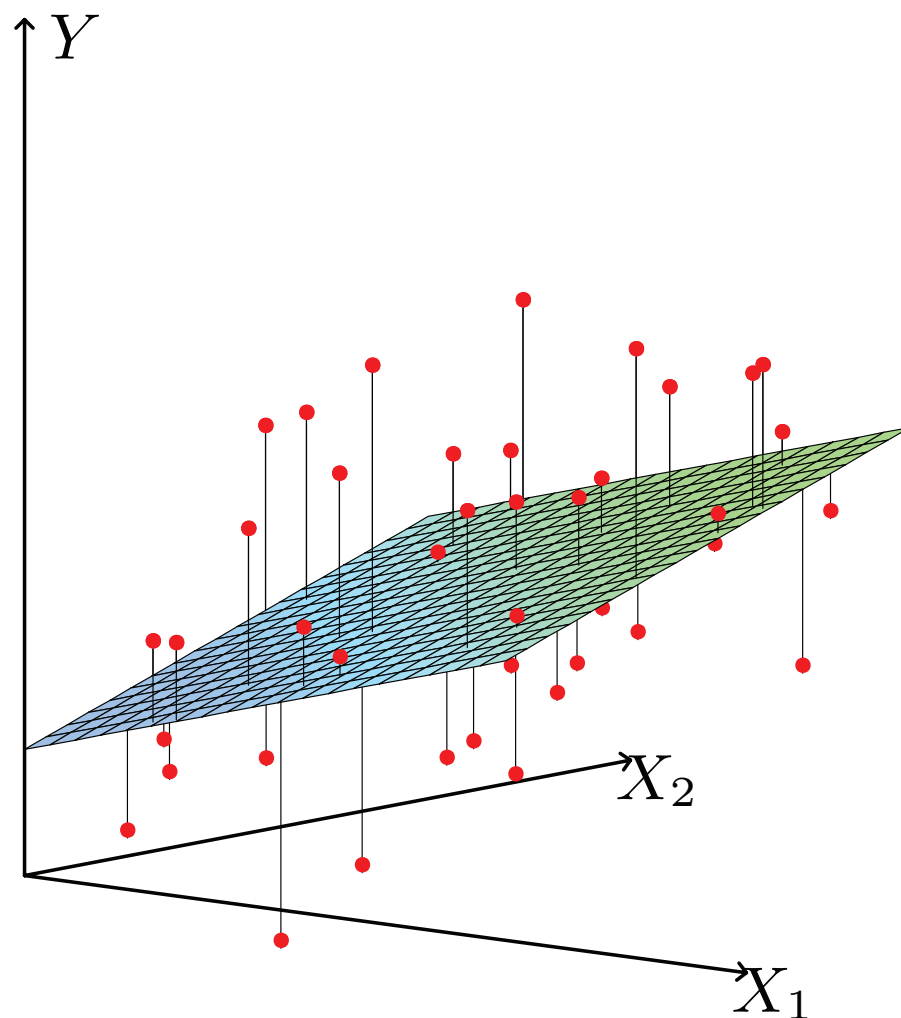
## Estimation by least squares

$$\hat{\beta} \;=\; \text{argmin} \sum_i (y_i - \beta_0 - \sum_{j=1}^{p} x_{ij}\beta_j)^2$$

$$\;=\; \text{argmin}(\mathbf{y} - \mathbf{X}\beta)^T(\mathbf{y} - \mathbf{X}\beta)$$

Solution satisfies *normal equations:* $\mathbf{X}^T(\mathbf{y} - \mathbf{X}\beta) = 0$.

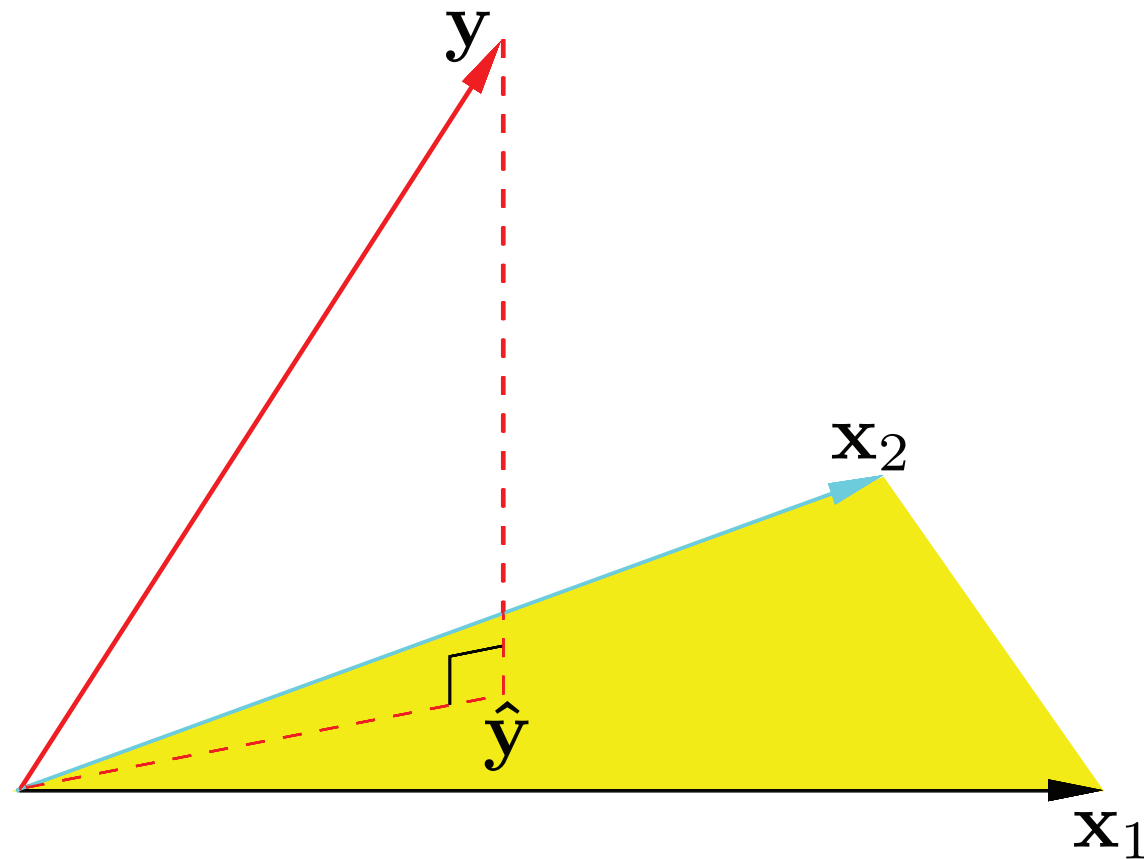If $\mathbf{X}$ full column rank,

$$\hat{\beta} \;=\; (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$$

$$\hat{\mathbf{y}} \;=\; \mathbf{X}\hat{\beta}$$

Also $\text{Var}\,(\hat{\beta}) = (\mathbf{X}^T\mathbf{X})^{-1}\sigma^2$

The $(p + 1)$-dimensional geometry of least squares estimation.

The $N$-dimensional geometry of least squares estimation

$$(\mathbf{y} - \hat{\mathbf{y}}) \perp \mathbf{x}_j, \; j = 1, \ldots, p$$

# Properties of OLS

- If $\mathbf{X}_1$ and $\mathbf{X}_2$ are mutually orthogonal matrices ($\mathbf{X}_1^T \mathbf{X}_2 = 0$), then the joint regression coefficients for $\mathbf{X} = (\mathbf{X}_1, \mathbf{X}_2)$ on $\mathbf{y}$ can be found from the separate regressions.
  Proof: $\mathbf{X}_1^T(\mathbf{y} - \mathbf{X}\hat{\beta}) = \mathbf{X}_1^T(\mathbf{y} - \mathbf{X}_1\hat{\beta}_1) = 0$. Same for $\hat{\beta}_2$.

- OLS is equivariant under non-singular linear transformations of $\mathbf{X}$. i.e. if $\hat{\beta}$ is OLS solution for $\mathbf{X}$, then $\beta^* = A^{-1}\hat{\beta}$ is OLS solution for $\mathbf{X}^* = \mathbf{X}A$ for $A_{p \times p}$ nonsingular.
  Proof: OLS is defined by orthogonal projection onto column space of $\mathbf{X}$. So $\hat{\mathbf{y}} = \mathbf{X}\hat{\beta} = \mathbf{X}^*\beta^*$.

- Let $\mathbf{X}_{(p)}$ be the submatrix of $\mathbf{X}$ excluding the last column $\mathbf{x}_p$. Let $\mathbf{z}_p = \mathbf{x}_p - \mathbf{X}_{(p)}\gamma$ (for any $\gamma$). Then OLS coefficient of $\mathbf{x}_p$ is the same as OLS coeffcient of $\mathbf{z}_p$ if we replace $\mathbf{x}_p$ by $\mathbf{z}_p$. Proof: previous point.

- Let $\gamma$ be the OLS coefficient of $\mathbf{x}_p$ on $\mathbf{X}_{(p)}$. Hence $\mathbf{z}_p$ is the residual obtained by *adjusting* $\mathbf{x}_p$ for all the other variables in the model.

- $\mathbf{X}_{(p)}^T \mathbf{z}_p = 0$ so the regression of $\mathbf{y}$ on $(\mathbf{X}_{(p)}, \mathbf{z}_p)$ decouples.

- The multiple regression coefficient of $\mathbf{x}_p$ is the same as the univariate coefficient in the regression of $\mathbf{y}$ on $\mathbf{z}_p$ i.e. $\mathbf{x}_p$ adjusted for the rest!

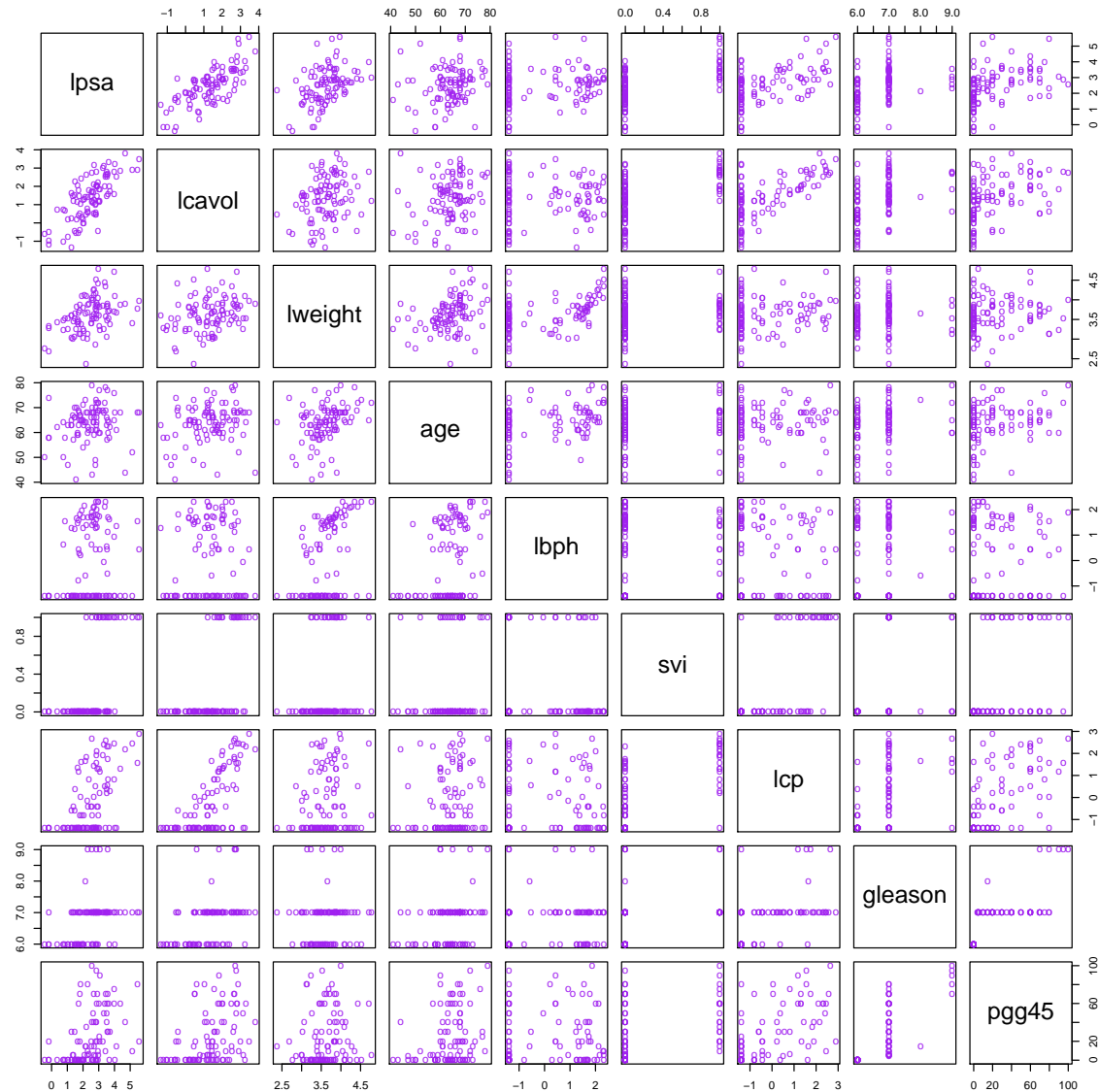$$\hat{\beta}_p = \frac{\langle \mathbf{z}_p, \mathbf{y} \rangle}{||\mathbf{z}_p||^2}$$

-

$$\mathrm{Var}\,(\hat{\beta}_p) = \frac{\sigma^2}{||\mathbf{z}_p||^2}$$

- Last statements true for all $j$, not just the last term $p$.

The course website has some additional more technical notes (linearR.pdf) on multiple linear regression, with an emphasis on computations.

# Example: Prostate cancer

- 97 observations on 9 variables (Stamey et al, 1989)

- Goal to predict log(PSA) from 8 clinical measurements/ demographics on men who were about to have their prostate removed.

- Next page shows a scatterplot matrix of all the date. This is created using the R expression

  ```
  pairs(lpsa ~., data=lprostate).
  ```

- Notice that several variables are correlated, and that `svi` is binary.

Scatterplot matrix of the *Prostate Cancer* data. The top variable is the response `lpsa`, and the top row shows its relationship with each of the others.

```
> lmfit=lm(lpsa~.,data=lprostate)
> summary(lmfit)


Call:
lm(formula = lpsa ~ ., data = lprostate)


Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)   0.180899   1.320601   0.137  0.89136
lcavol        0.564355   0.087831   6.425 6.54e-09 ***
lweight       0.622081   0.200892   3.097  0.00263 **
age          -0.021248   0.011084  -1.917  0.05848 .
lbph          0.096676   0.057915   1.669  0.09862 .
svi           0.761652   0.241173   3.158  0.00218 **
lcp          -0.106055   0.089866  -1.180  0.24112
gleason       0.049287   0.155340   0.317  0.75178
pgg45         0.004458   0.004365   1.021  0.30999


Residual standard error: 0.6995 on 88 degrees of freedom
Multiple R-squared: 0.6634,     Adjusted R-squared: 0.6328
F-statistic: 21.68 on 8 and 88 DF,  p-value: < 2.2e-16
```

## The Bias-variance tradeoff

A good measure of the quality of an estimator $\hat{\mathbf{f}}(x)$ is the mean squared error. Let $\mathbf{f}_0(x)$ be the true value of $\mathbf{f}(x)$ at the point $x$. Then

$$\text{Mse}\,[\hat{\mathbf{f}}(x)] = \text{E}\,[\hat{\mathbf{f}}(x) - \mathbf{f}_0(x)]^2$$

This can be written as

$$\text{Mse}\,[\hat{\mathbf{f}}(x)] = \text{Var}\,[\hat{\mathbf{f}}(x)] + [\text{E}\,\hat{\mathbf{f}}(x) - \mathbf{f}_0(x)]^2$$
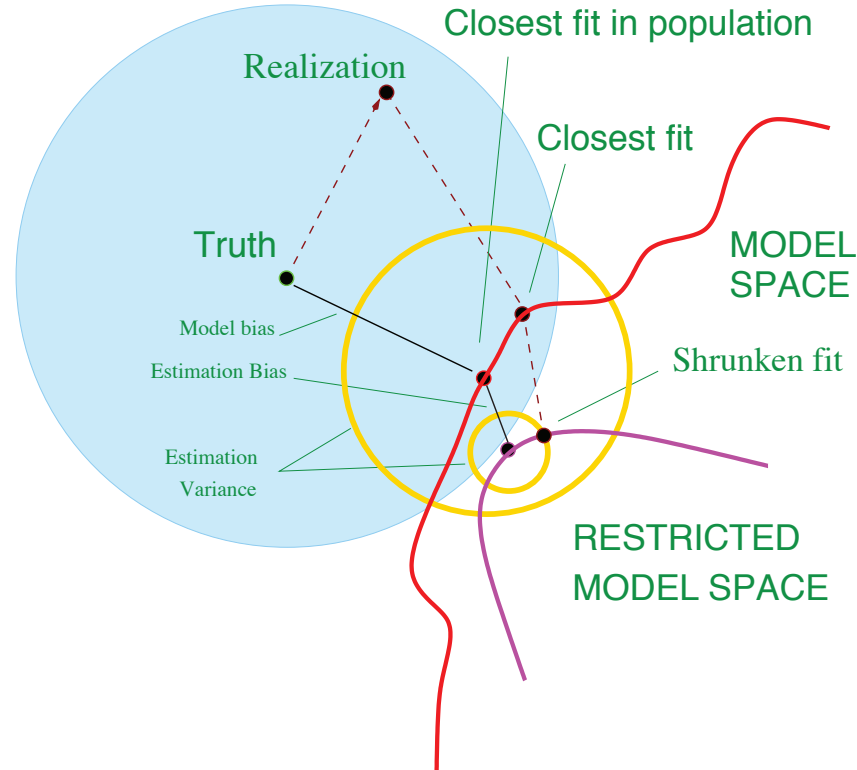
This is *variance* plus squared *bias*.

Typically, when bias is low, variance is high and vice-versa. Choosing estimators often involves a tradeoff between bias and variance.

- If the linear model is correct for a given problem, then the least squares prediction $\hat{\mathbf{f}}$ is unbiased, and has the lowest variance among all unbiased estimators that are linear functions of $\mathbf{y}$

- But there can be (and often exist) biased estimators with smaller Mse .

- Generally, by *regularizing* (shrinking, dampening, controlling) the estimator in some way, its variance will be reduced; if the corresponding increase in bias is small, this will be worthwhile.

- Examples of regularization: subset selection (forward, backward, all subsets); ridge regression, the lasso.

- In reality models are almost never correct, so there is an additional *model bias* between the closest member of the linear model class and the truth.

# Model Selection

Often we prefer a restricted estimate because of its reduced estimation variance.

## Analysis of time series data

Two approaches: *frequency domain* (fourier)—see discussion of wavelet smoothing.

*Time domain*. Main tool is auto-regressive (AR) model of order $k$:

$$y_t = \beta_1 y_{t-1} + \beta_2 y_{t-2} \cdots + \beta_k y_{t-k} + \epsilon_t$$

Fit by linear least squares regression on lagged data

$$
\begin{aligned}
y_t &= \beta_1 y_{t-1} + \beta_2 y_{t-2} \cdots \beta_k y_{t-k} \\
y_{t-1} &= \beta_1 y_{t-2} + \beta_2 y_{t-3} \cdots \beta_k y_{t-k-1} \\
\vdots &= \vdots \\
y_{k+1} &= \beta_1 y_k + \beta_2 y_{k-1} \cdots \beta_k y_1
\end{aligned}
$$

## Example: NYSE data

Time series of 6200 daily measurements, 1962-1987

`volume` — log(trading volume) — *outcome*

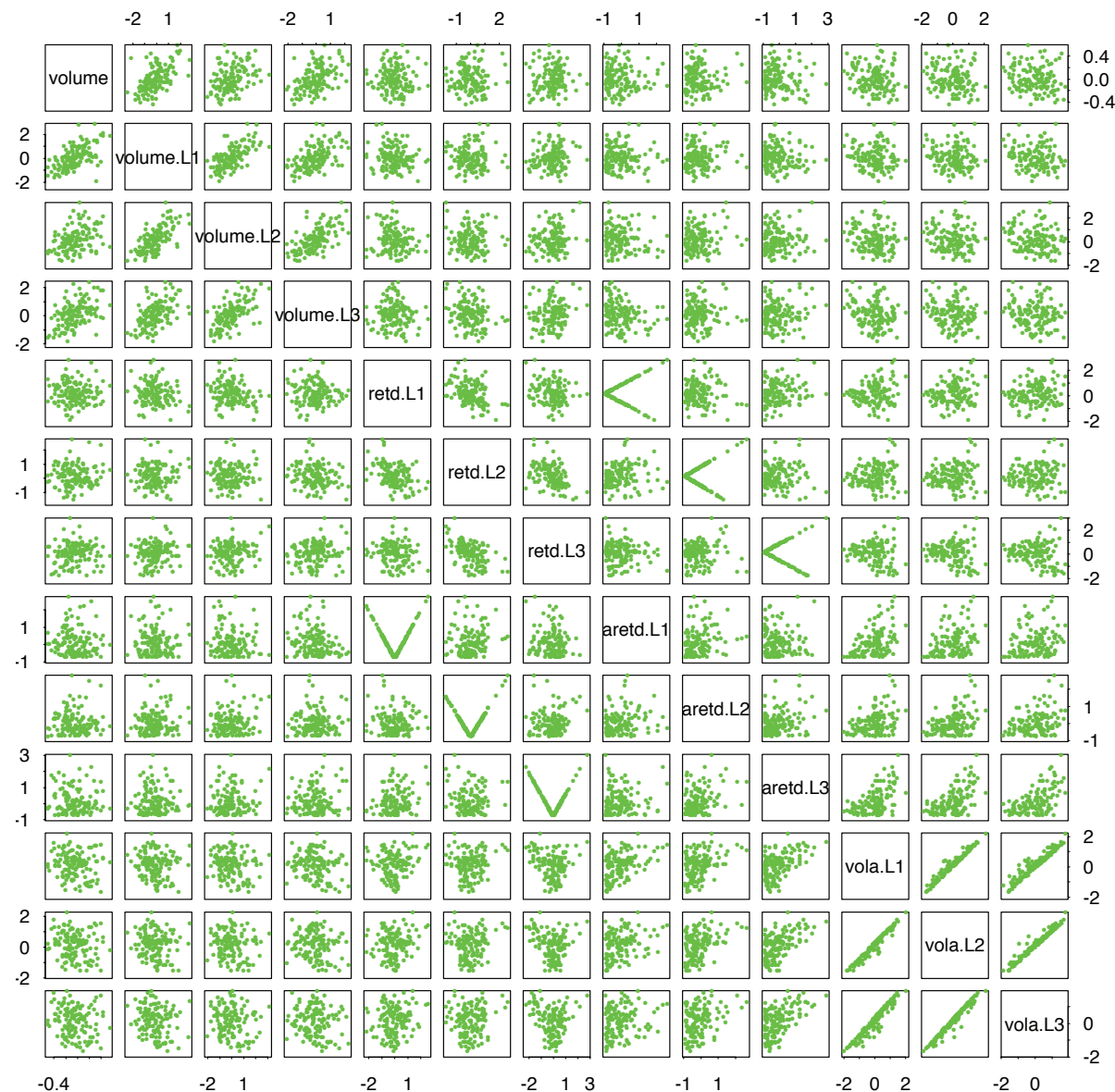`volume.Lj` — log(trading volume)$_{\text{day}-j}$, $j = 1, 2, 3$

`ret.Lj` — $\Delta$ log(Dow Jones)$_{\text{day}-j}$, $j = 1, 2, 3$

`aret.Lj` — $|\Delta\log(\text{Dow Jones})|_{\text{day}-j}$, $j = 1, 2, 3$

`vola.Lj` — volatility$_{\text{day}-j}$, $j = 1, 2, 3$

Source—Weigend and LeBaron (1994)

We randomly selected a training set of size 50 and a test set of size 500, from the first 600 observations.

# OLS Fit

Results of ordinary least squares analysis of NYSE data

| Term | Coefficient | Std. Error | t-Statistic |
|---|---|---|---|
| Intercept | -0.02 | 0.04 | -0.64 |
| volume.L1 | 0.09 | 0.05 | 1.80 |
| volume.L2 | 0.06 | 0.05 | 1.19 |
| volume.L3 | 0.04 | 0.05 | 0.81 |
| retd.L1 | 0.00 | 0.04 | 0.11 |
| retd.L2 | -0.02 | 0.05 | -0.46 |
| retd.L3 | -0.03 | 0.04 | -0.65 |
| aretd.L1 | 0.08 | 0.07 | 1.12 |
| aretd.L2 | -0.02 | 0.05 | -0.45 |
| aretd.L3 | 0.03 | 0.04 | 0.77 |
| vola.L1 | 0.20 | 0.30 | 0.66 |
| vola.L2 | -0.50 | 0.40 | -1.25 |
| vola.L3 | 0.27 | 0.34 | 0.78 |

# Variable subset selection

We retain only a subset of the coefficients and set to zero the coefficients of the rest.

There are different strategies:

- *All subsets regression* finds for each $s \in 0, 1, 2, \ldots p$ the subset of size $s$ that gives smallest residual sum of squares. The question of how to choose $s$ involves the tradeoff between bias and variance: can use cross-validation (see below)

- Rather than search through all possible subsets, we can seek a good path through them. *Forward stepwise selection* starts with the intercept and then sequentially adds into the model the variable that most improves the fit. The improvement in fit is usually based on the

$F$ ratio

$$F = \frac{RSS(\hat{\beta}^{old}) - RSS(\hat{\beta}^{new})}{RSS(\hat{\beta}^{new})/(N - s)}$$

- *Backward stepwise selection* starts with the full OLS model, and sequentially deletes variables.

- There are also hybrid *stepwise selection* strategies which add in the best variable and delete the least important variable, in a sequential manner.

- Each procedure has one or more *tuning parameters*:

  - subset size

  - P-values for adding or dropping terms

# Model Assessment

Objectives:

1. Choose a value of a tuning parameter for a technique

2. Estimate the prediction performance of a given model

For both of these purposes, the best approach is to run the procedure on an independent test set, if one is available

If possible one should use different test data for (1) and (2) above: a *validation set* for (1) and a *test set* for (2)

Often there is insufficient data to create a separate validation or test set. In this instance *Cross-Validation* is useful.

# $K$-Fold Cross-Validation

Primary method for estimating a tuning parameter $\lambda$ (such as subset size)

Divide the data into $K$ roughly equal parts (typically $K$=5 or 10)

| 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|
| Train | Train | Test | Train | Train |

- for each $k = 1, 2, \ldots K$, fit the model with parameter $\lambda$ to the other $K - 1$ parts, giving $\hat{\beta}^{-k}(\lambda)$ and compute its error in predicting the $k$th part:
$E_k(\lambda) = \sum_{i \in kth\ part}(y_i - \mathbf{x}_i \hat{\beta}^{-k}(\lambda))^2.$
This gives the cross-validation error

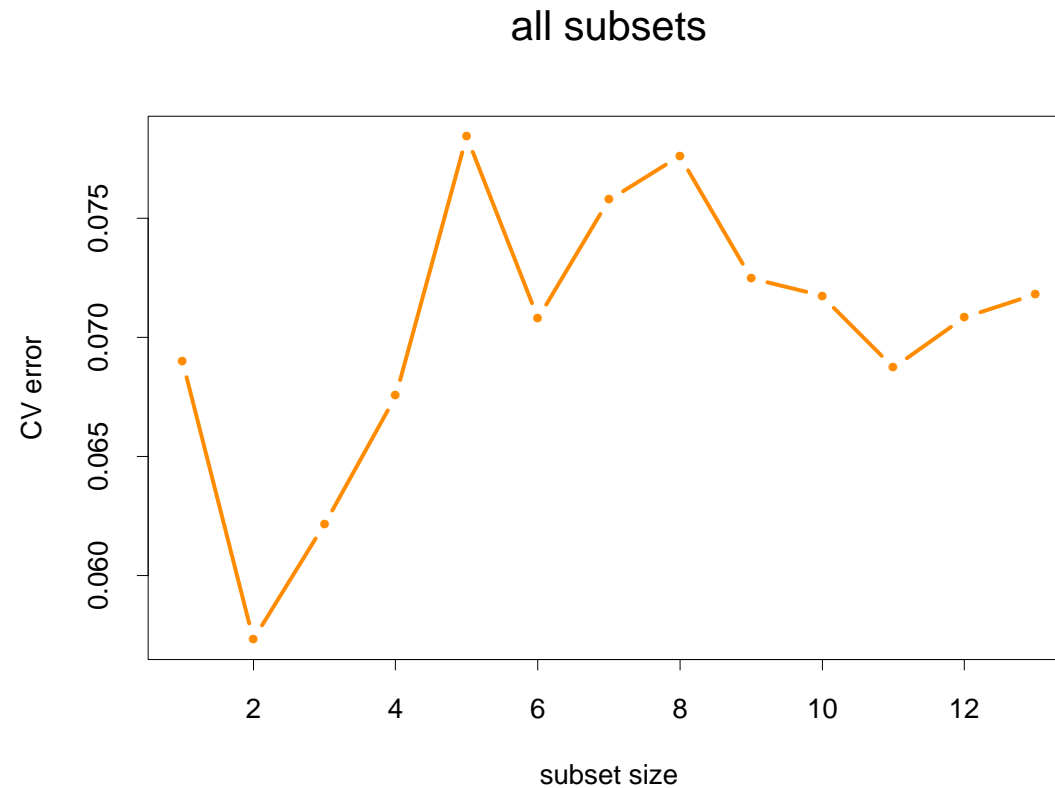$$CV(\lambda) = \frac{1}{K} \sum_{k=1}^{K} E_k(\lambda)$$

- do this for many values of $\lambda$ and choose the value of $\lambda$ that makes $CV(\lambda)$ smallest.

- In our variable subsets example, $\lambda$ is the subset size

- $\hat{\beta}^{-k}(\lambda)$ are the coefficients for the best subset of size $\lambda$, found from the training set that leaves out the $k$th part of the data

- $E_k(\lambda)$ is the estimated test error for this best subset.

- from the $K$ cross-validation training sets, the $K$ test error estimates are averaged to give
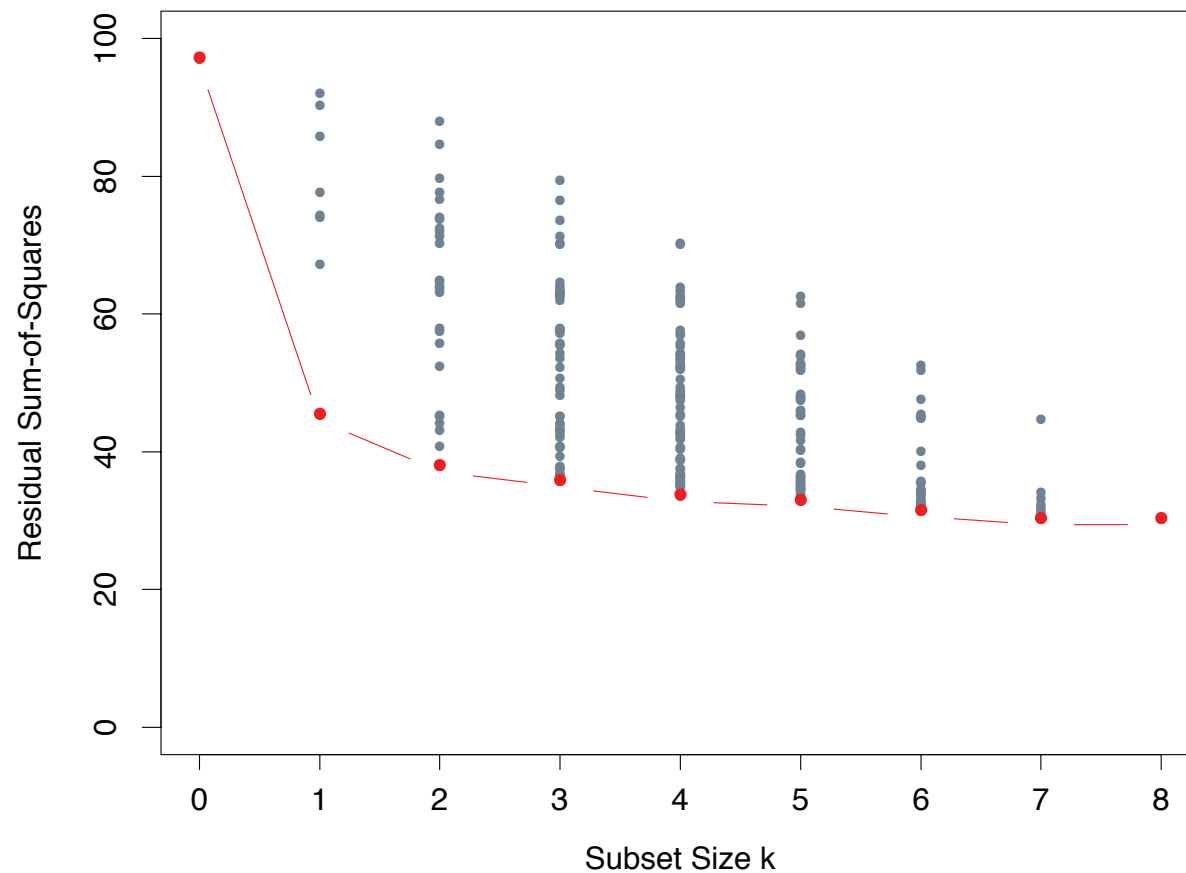
$$CV(\lambda) = (1/K) \sum_{k=1}^{K} E_k(\lambda).$$

- Note that different subsets of size $\lambda$ will (probably) be found from each of the $K$ cross-validation training sets. Doesn't matter: focus is on subset size, not the actual subset.

all subsets



**CV curve for NYSE data**

- The focus is on *subset size*—not which variables are in the model.

- Variance increases slowly—typically $\sigma^2/N$ per variable.

All possible subset models for the prostate cancer example. At each subset size is shown the residual sum-of-squares for each model of that size.

# The Bootstrap approach

- Bootstrap works by sampling $N$ times with replacement from training set to form a "bootstrap" data set. Then model is estimated on bootstrap data set, and predictions are made for original training set.

- This process is repeated many times and the results are averaged.

- Bootstrap most useful for estimating standard errors of predictions.

- Can also use modified versions of the bootstrap to estimate prediction error. Sometimes produces better estimates than cross-validation

  Eg — for each bootstrap sample, estimate errors using only observations excluded from bootstrap sample.

# Cross-validation- revisited

Consider a simple classifier for wide data:

1.  Starting with 5000 predictors, find the 200 predictors having the largest correlation with the class labels

2.  Carry about nearest-centroid classification using only these 200 genes

How do we estimate the test set performance of this classifier?

✘  *Wrong:* Apply cross-validation in step 2.

✔  *Right:* Apply cross-validation to steps 1 and 2.

It is easy to simulate realistic data with the class labels independent of the outcome — so that true test error $=50\%$ — but *Wrong* CV error estimate is zero! We have seen this error made in 4 high profile microarray papers in the last couple of years. See Ambroise and McLachlan PNAS 2002.
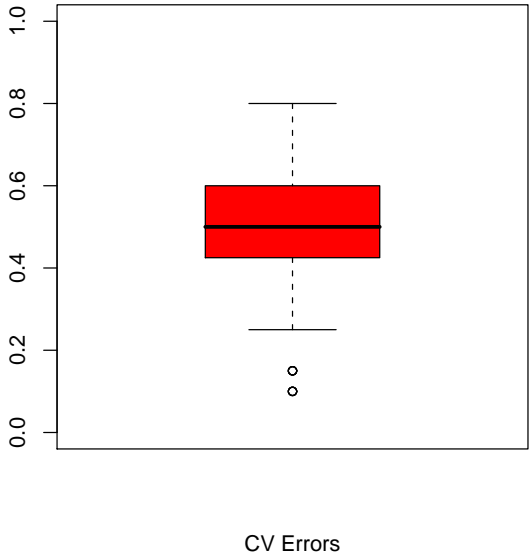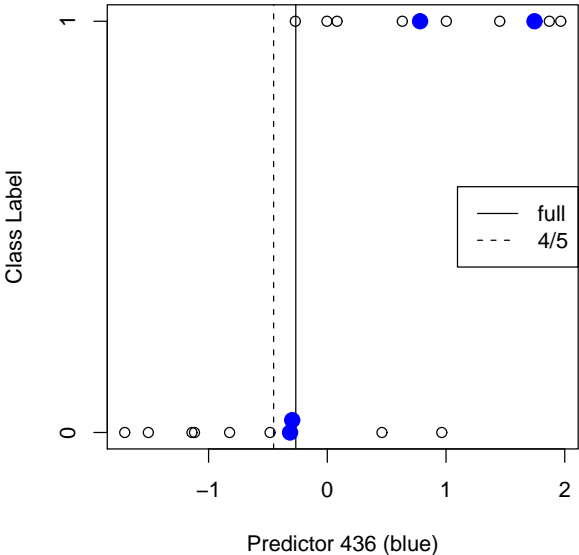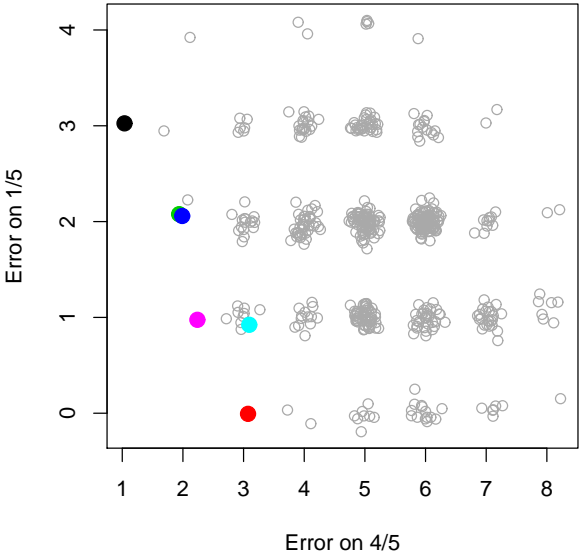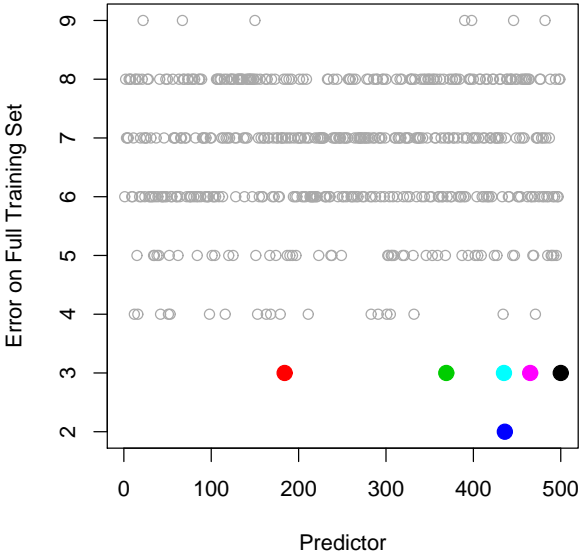
*A little cheating goes a long way!*

# Validation and test set issues

- Important to have both cross-validation and test sets, since we often run CV many times, fiddling with different parameters. This can bias the CV results

- A separate test set provides a convincing, independent assessment of a model's performance

- Test-set results might still overestimate actual performance, as a real future test set may differ in many ways from today's data

# Does cross-validation really work?

- Consider a scenario with $N = 20$ samples in two equal-sized classes, and $p = 500$ quantitative features that are independent of the class labels. The true error rate of any classifier is 50%.

- Consider a simple univariate classifier — a single split that minimizes the misclassification error (a "stump").

- Fitting to the entire training set, we will find a feature that splits the data very well

- If we do 5-fold CV, this same feature should split any $4/5$ths and $1/5$th of the data well too, and hence its CV error will be small (much less than 50%)

- Thus CV does not give an accurate estimate of error.

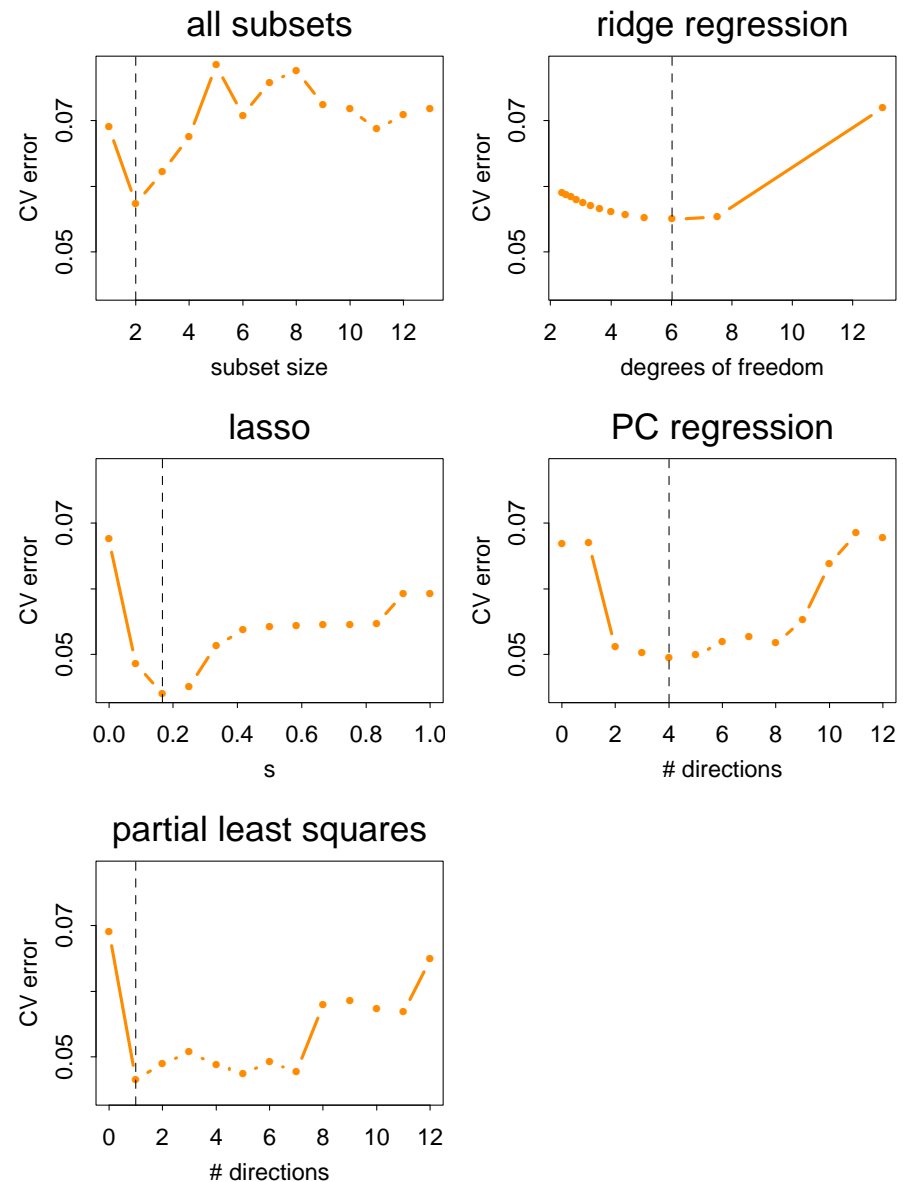- Is this argument correct? (Details in Section 7.10.3).

# NYSE example continued

Table shows the coefficients from a number of different selection and shrinkage methods, applied to the NYSE data.

| Term | OLS | VSS | Ridge | Lasso | PCR | PLS |
|---|---|---|---|---|---|---|
| Intercept | -0.02 | 0.00 | -0.01 | -0.02 | -0.02 | -0.04 |
| volume.L1 | 0.09 | 0.16 | 0.06 | 0.09 | 0.05 | 0.06 |
| volume.L2 | 0.06 | 0.00 | 0.04 | 0.02 | 0.06 | 0.06 |
| volume.L3 | 0.04 | 0.00 | 0.04 | 0.03 | 0.04 | 0.05 |
| retd.L1 | 0.00 | 0.00 | 0.01 | 0.01 | 0.02 | 0.01 |
| retd.L2 | -0.02 | 0.00 | -0.01 | 0.00 | -0.01 | -0.02 |
| retd.L3 | -0.03 | 0.00 | -0.01 | 0.00 | -0.02 | 0.00 |
| aretd.L1 | 0.08 | 0.00 | 0.03 | 0.02 | -0.02 | 0.00 |
| aretd.L2 | -0.02 | -0.05 | -0.03 | -0.03 | -0.01 | -0.01 |
| aretd.L3 | 0.03 | 0.00 | 0.01 | 0.00 | 0.02 | 0.01 |
| vola.L1 | 0.20 | 0.00 | 0.00 | 0.00 | -0.01 | -0.01 |
| vola.L2 | -0.50 | 0.00 | -0.01 | 0.00 | -0.01 | -0.01 |
| vola.L3 | 0.27 | 0.00 | -0.01 | 0.00 | -0.01 | -0.01 |
| Test err | 0.050 | 0.041 | 0.042 | 0.039 | 0.045 | 0.044 |
| SE | 0.007 | 0.005 | 0.005 | 0.005 | 0.006 | 0.006 |

CV was used on the 50 training observations (except for OLS). Test error for constant: 0.061.
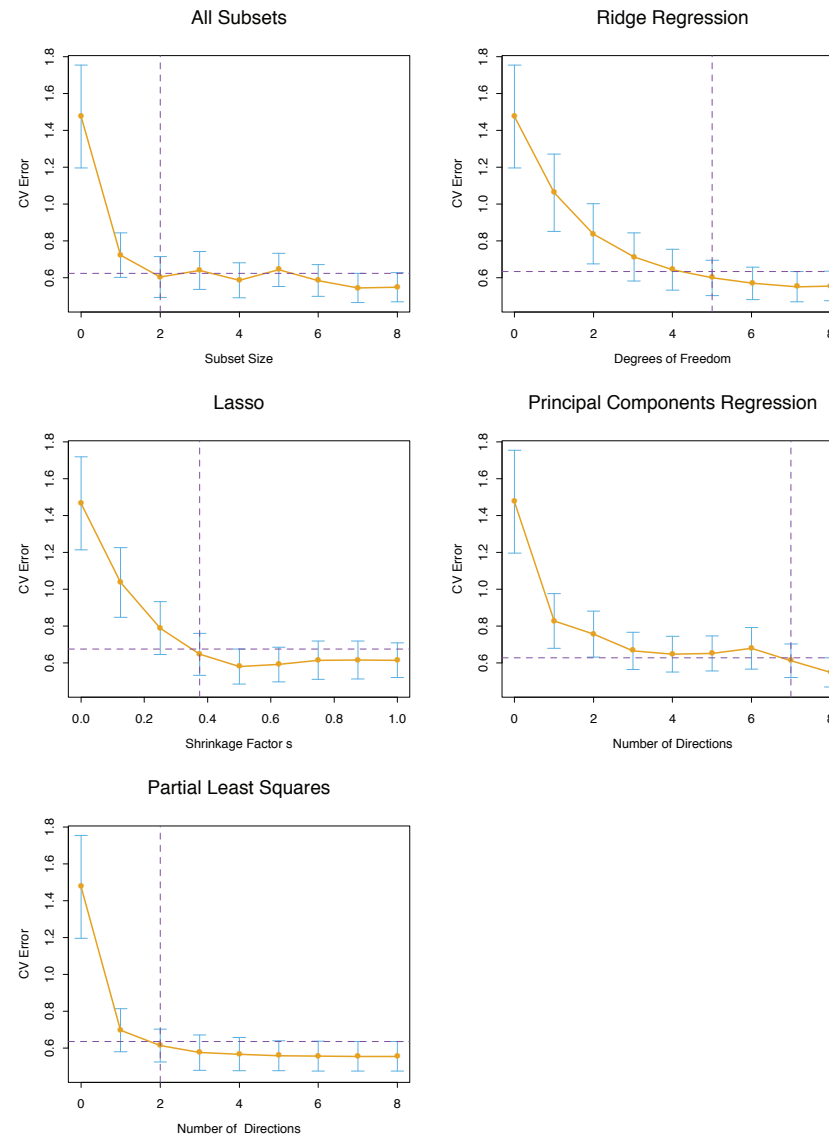
# Crossvalidation results for NYSE data

Estimated prediction error curves for the various selection and shrinkage methods. The arrow indicates the estimated minimizing value of the complexity parameter. Training sample size = 50.

# Crossvalidation results for Prostate Cancer data

Estimated prediction error curves for the various selection and shrinkage methods. The arrow indicates the estimated minimizing value of the complexity parameter. Training sample size = 67.

# Ridge Regression

The ridge estimator is defined by

$$\hat{\beta}^{\text{ridge}} = \text{argmin}(\mathbf{y} - \mathbf{X}\beta)^T(\mathbf{y} - \mathbf{X}\beta) + \lambda\beta^T\beta$$
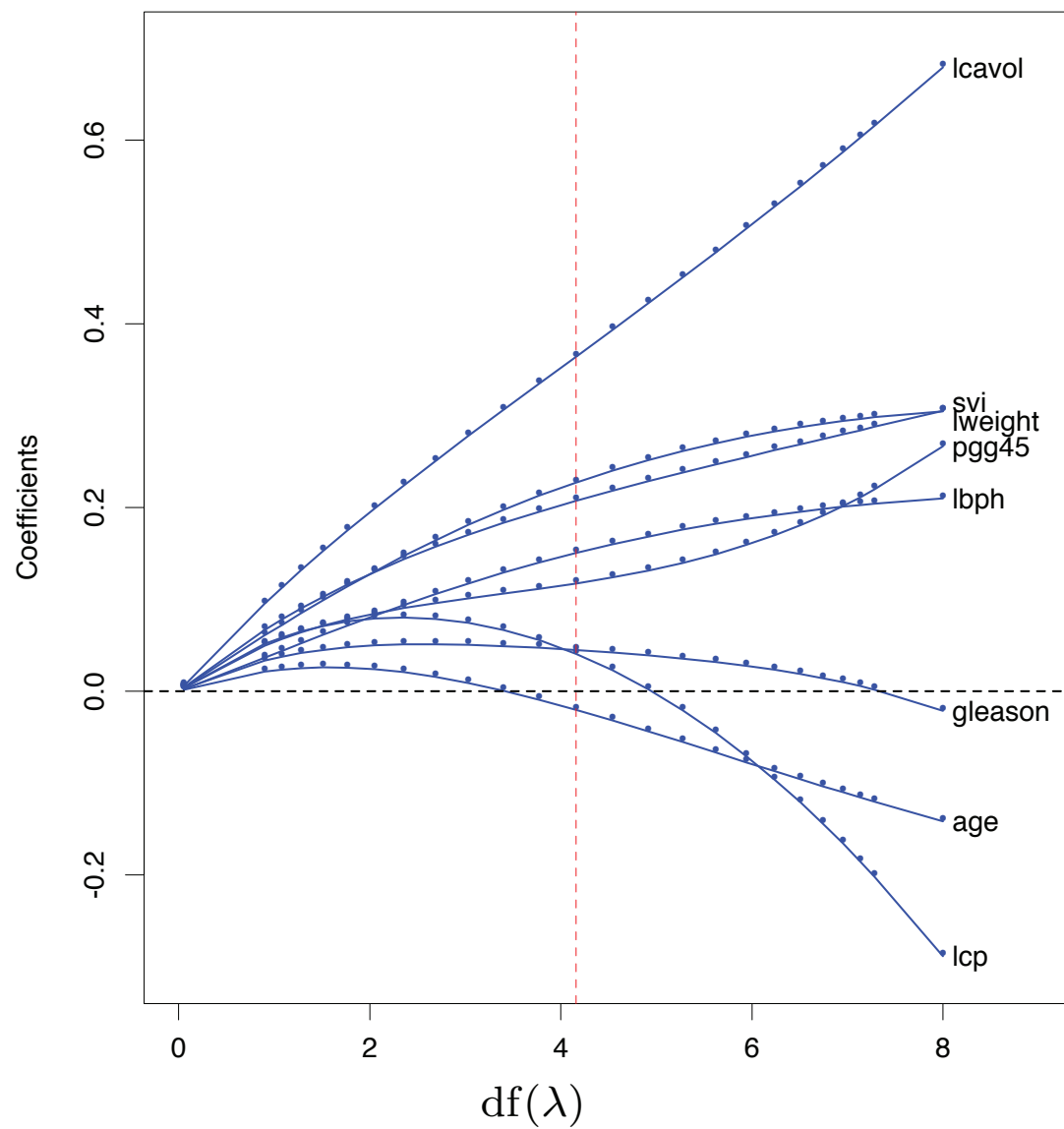
Equivalently,

$$\hat{\beta}^{\text{ridge}} = \text{argmin } (\mathbf{y} - \mathbf{X}\beta)^T(\mathbf{y} - \mathbf{X}\beta)$$
$$\text{subject to } \sum \beta_j^2 \leq s.$$

The parameter $\lambda > 0$ penalizes $\beta_j$ proportional to its size $\beta_j^2$. Solution is

$$\hat{\beta}_\lambda = (\mathbf{X}^T\mathbf{X} + \lambda I)^{-1}\mathbf{X}^T\mathbf{y}$$

where $I$ is the identity matrix. This is a biased estimator that for some value of $\lambda > 0$ may have smaller mean squared error than the least squares estimator.

Note $\lambda = 0$ gives the least squares estimator; if $\lambda \to \infty$, then $\hat{\beta} \to 0$.
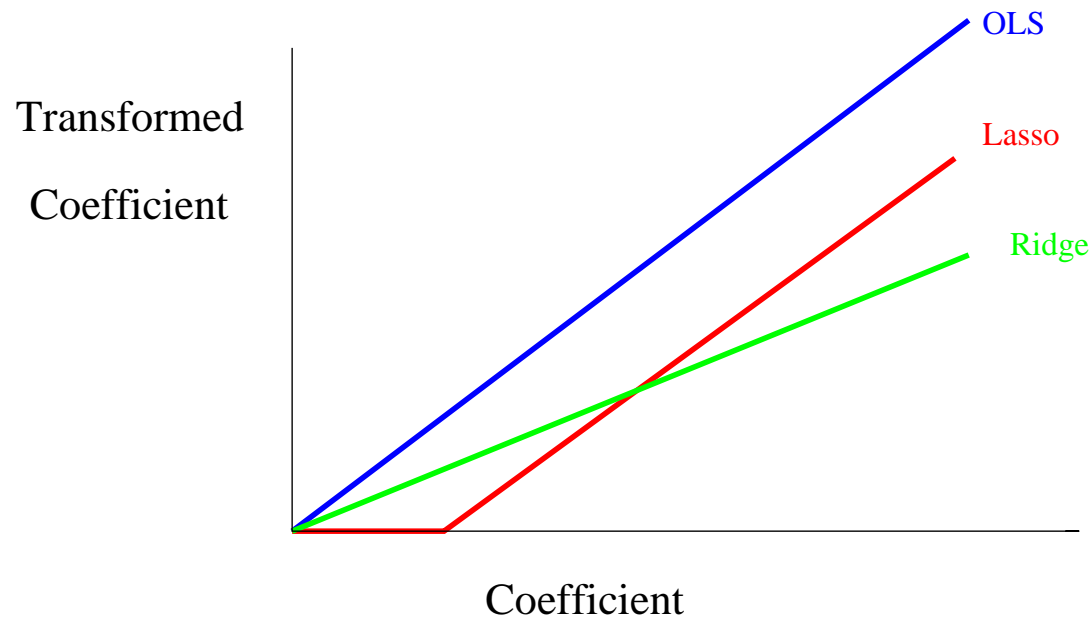
# The Lasso

The lasso is a shrinkage method like ridge, but acts in a nonlinear manner on the outcome $y$.

The lasso is defined by

$$
\begin{aligned}
\hat{\beta}^{\text{lasso}} \quad &= \quad \text{argmin} \; (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) \\
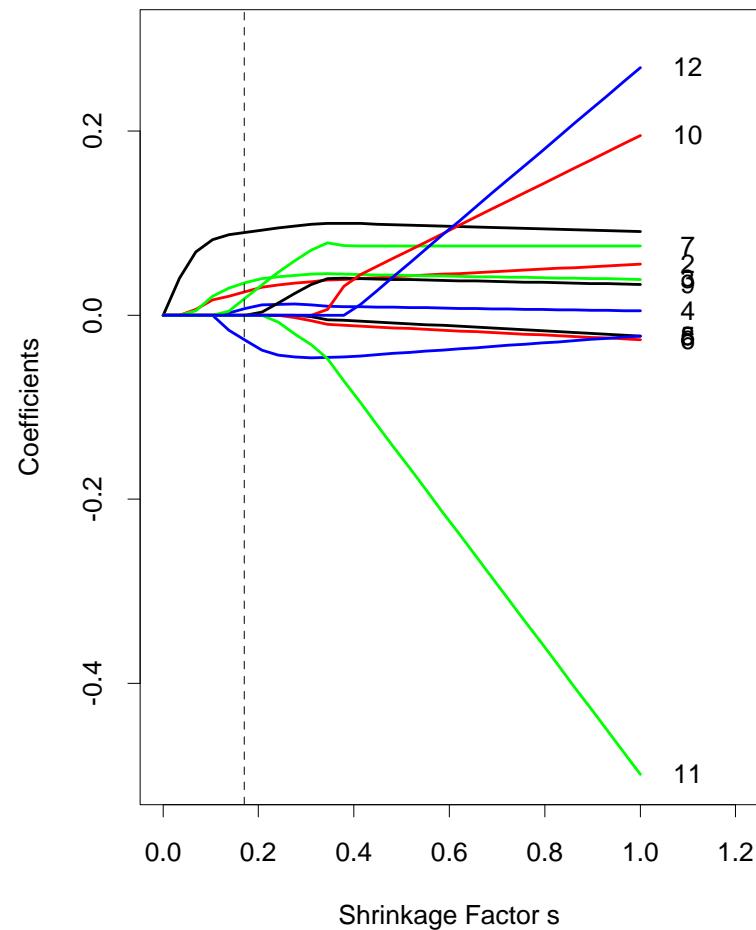&\qquad \text{subject to } \sum |\beta_j| \leq t
\end{aligned}
$$

- Notice that ridge penalty $\sum \beta_j^2$ is replaced by $\sum |\beta_j|$.

- this makes the solutions nonlinear in $\mathbf{y}$, and a quadratic programming algorithm is used to compute them.

- because of the nature of the constraint, if $t$ is chosen small enough then the lasso will set some coefficients exactly to zero. Thus the lasso does a kind of continuous model selection.
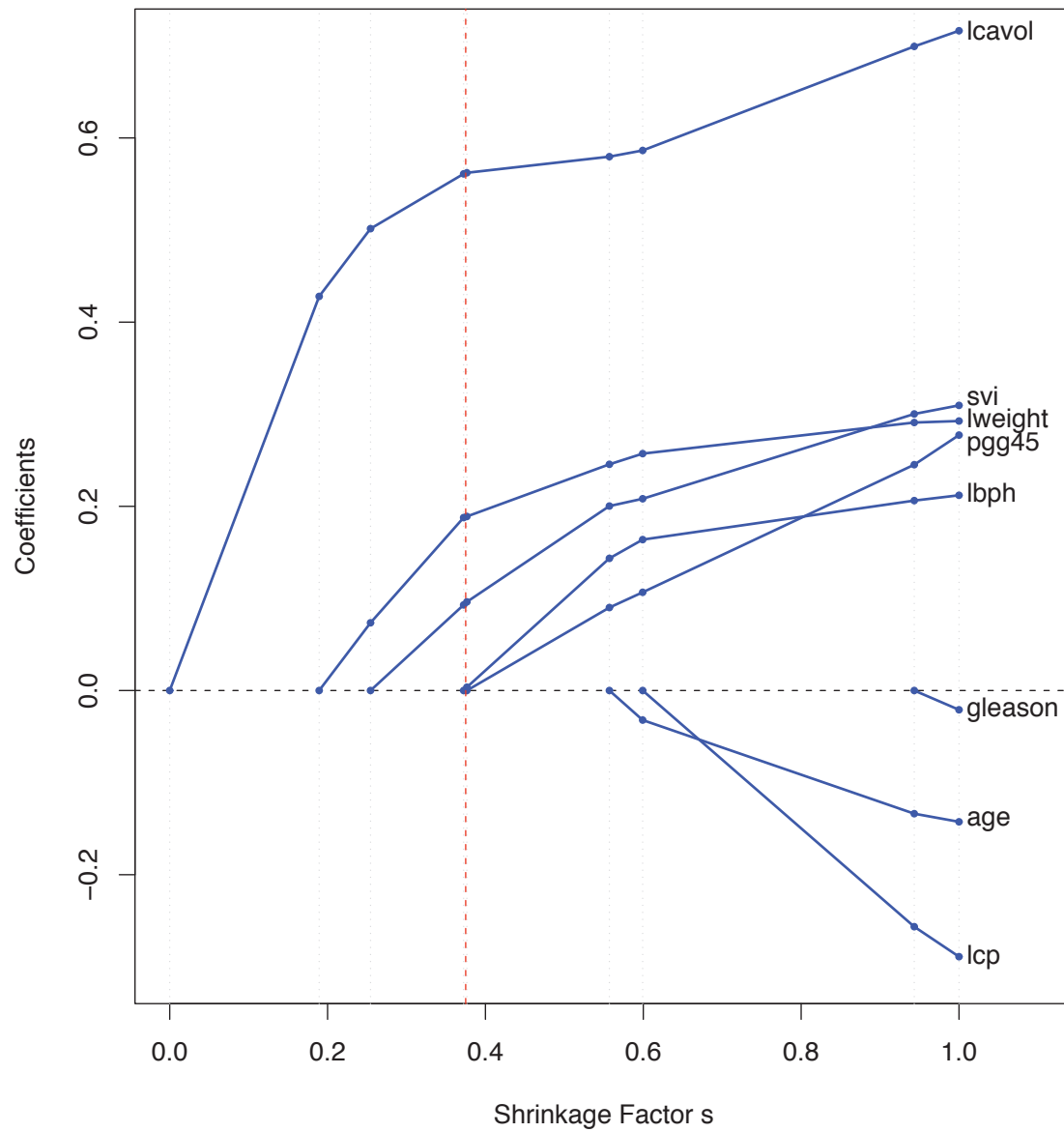
- The parameter $t$ should be adaptively chosen to minimize an estimate of expected, using say cross-validation

- *Ridge vs Lasso:* if inputs are orthogonal, ridge *multiplies* least squares coefficients by a constant $< 1$, lasso *translates* them towards zero by a constant, truncating at zero.

Transformed

Coefficient
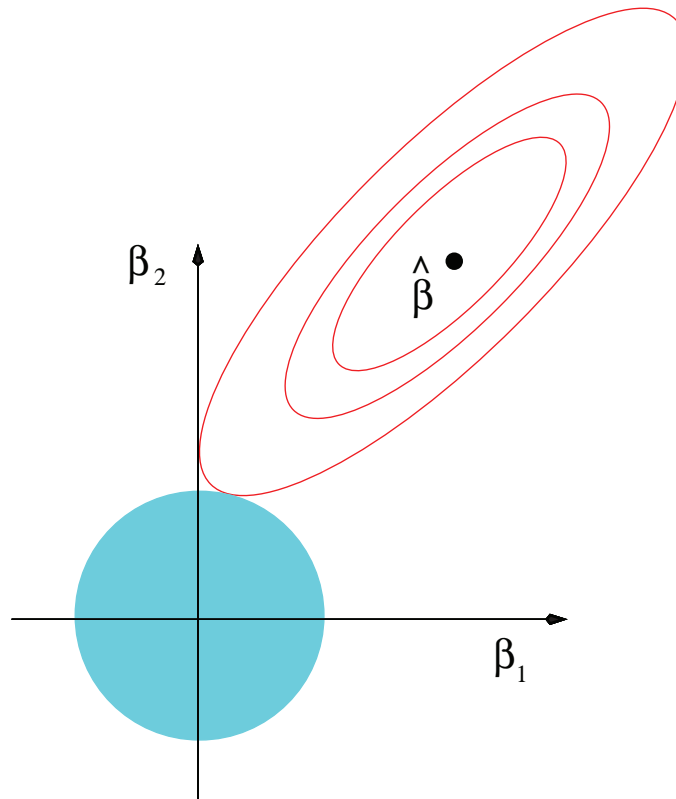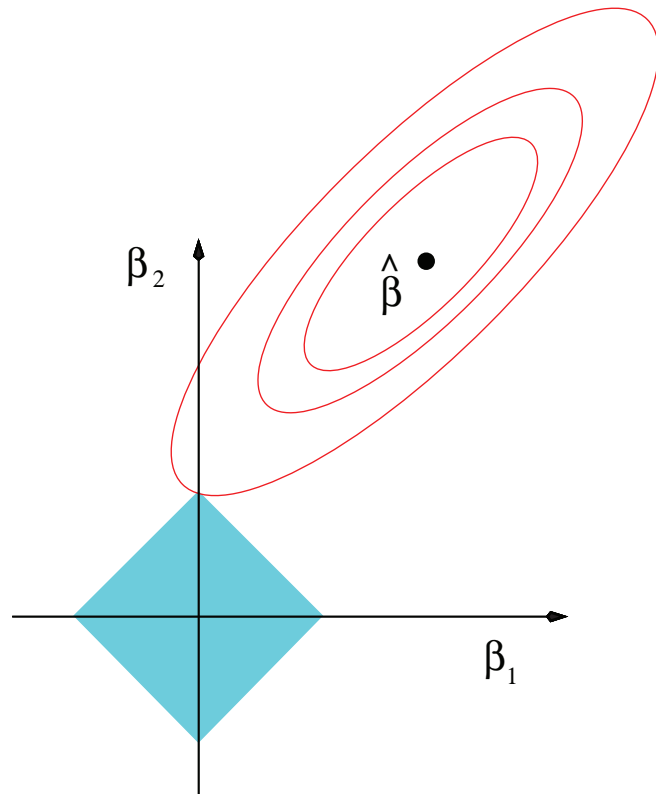
OLS

Lasso

Ridge

Coefficient

# Lasso in Action

Profiles of coefficients for NYSE data as lasso shrinkage is varied. $s = t/t_0 \in [0, 1]$, where $t_0 = \sum |\hat{\beta}_{OLS}|$.
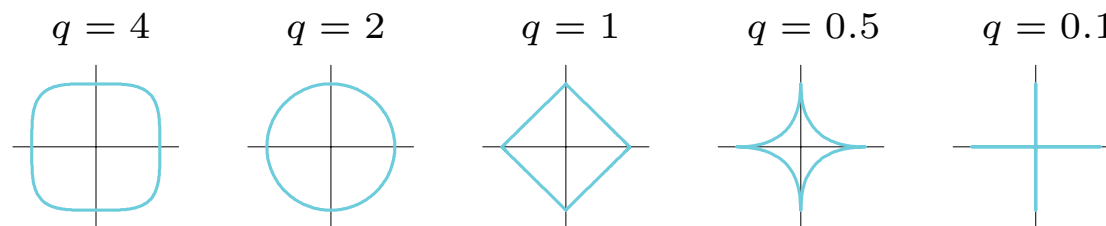
## A family of shrinkage estimators

Consider the criterion

$$\tilde{\beta} \;\; = \;\; \text{argmin}_{\beta} \sum_{i=1}^{N} (y_i - x_i^T \beta)^2$$

$$\text{subject to } \sum |\beta_j|^q \leq s$$

for $q \geq 0$. The contours of constant value of $\sum_j |\beta_j|^q$ are shown for the case of two inputs.



$q = 4 \qquad q = 2 \qquad q = 1 \qquad q = 0.5 \qquad q = 0.1$

*Contours of constant value of $\sum_j |\beta_j|^q$ for given values of q.*

Thinking of $|\beta_j|^q$ as the log-prior density for $\beta_j$, these are also the equi-contours of the prior.

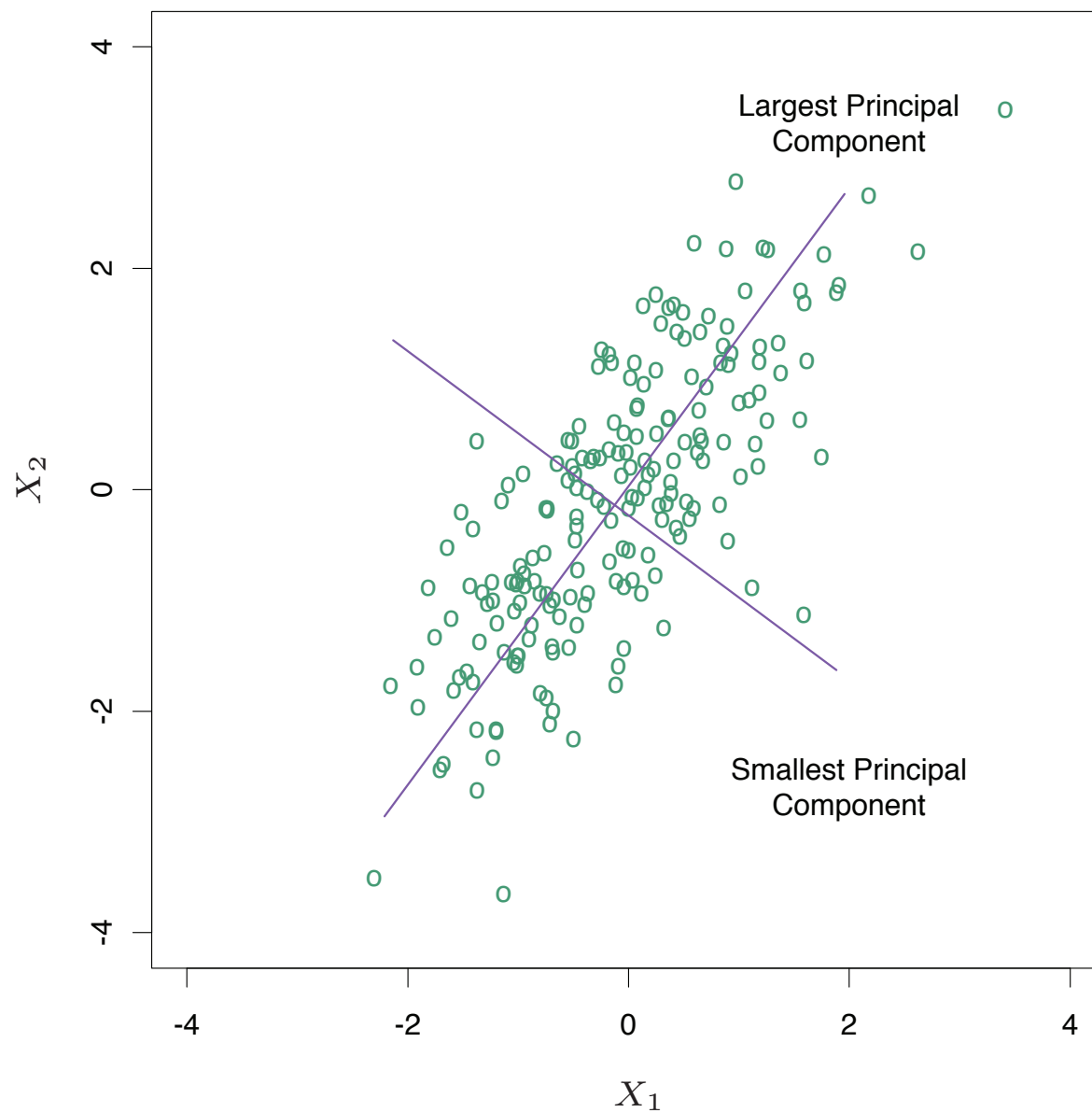## Use of derived input directions — Principal Component Regression

We choose a set of linear combinations of the $X_j$s, and then regress the outcome on these linear combinations.

The largest principal component $Z_1$ is the standardized linear combination of the $X_j$s with largest variance. Subsequent principal components $Z_2$, $Z_3$, ... maximize variance subject top being uncorrelated with the preceding components.

If $\mathbf{S}$ is the sample covariance matrix of $\mathbf{x}_1, \ldots, \mathbf{x}_p$, then the eigenvector equations

$$\mathbf{S}\mathbf{v}_j = d_j^2 \mathbf{v}_j, \quad j = 1, \ldots, p$$

define the principal components of $\mathbf{S}$.

Largest Principal Component

Smallest Principal Component

$X_2$

$X_1$

# PCA regression continued

- The $\mathbf{v}_j$ are the (ordered) principal component directions; the derived principal component variables are given by $\mathbf{z}_j = \mathbf{X}\mathbf{v}_j$. The variance of $\mathbf{z}_j$ is $d_j^2$, and determines the ordering.

- Then principal components regression regresses $\mathbf{y}$ on $\mathbf{z}_1, \mathbf{z}_2, \dots \mathbf{z}_J$ for some $J \leq p$.

- Since the $\mathbf{z}_j$s are orthogonal, this regression is just a sum of univariate regressions:

$$\hat{\mathbf{y}}^{\mathrm{pcr}} = \bar{y} + \sum_{j=1}^{J} \hat{\gamma}_j \mathbf{z}_j$$

where $\hat{\gamma}_j$ is the univariate regression coefficient of $\mathbf{y}$ on $\mathbf{z}_j$.
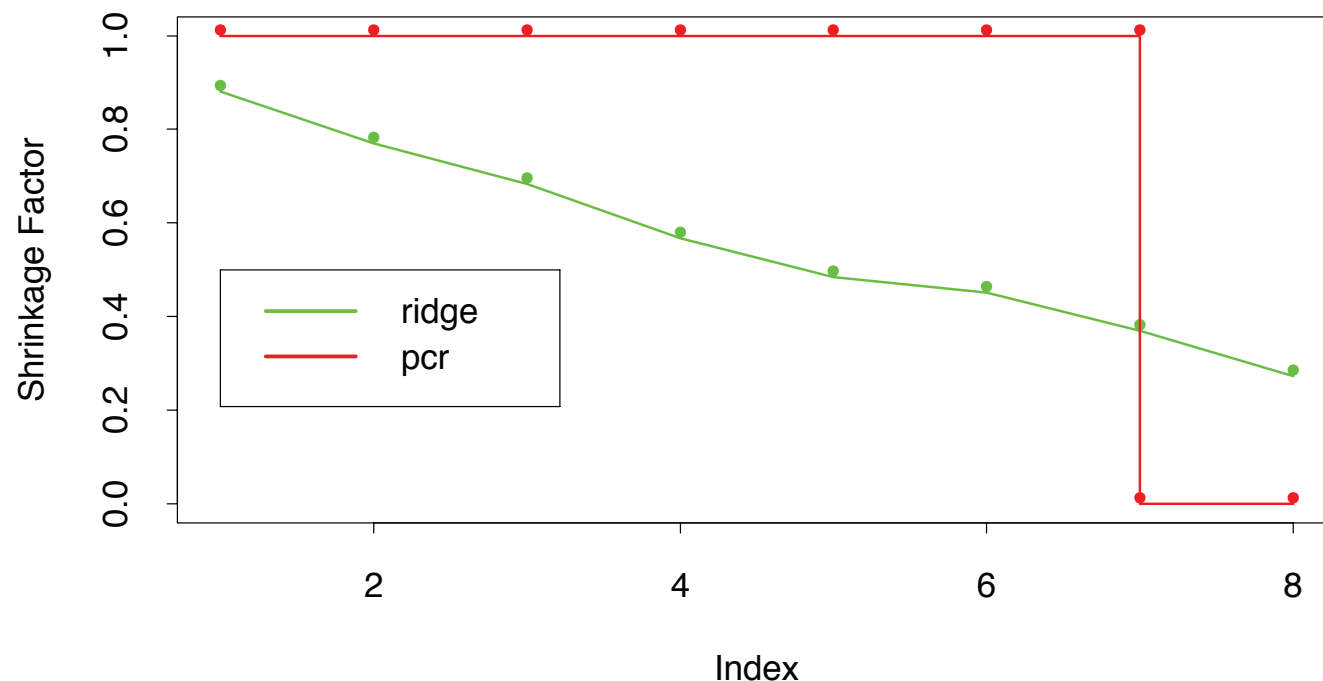
# Principal components and the SVD

Let $\tilde{\mathbf{X}} = \mathbf{U}\mathbf{D}\mathbf{V}^T$ be the SVD of the *centered* version $\tilde{\mathbf{X}}$ of the model matrix $\mathbf{X}$. This SVD provides the principal components of $\mathbf{X}$.

Proof:

$$\mathbf{S} = \frac{1}{N}\tilde{\mathbf{X}}^T\tilde{\mathbf{X}} = \frac{1}{N}\mathbf{V}\mathbf{D}^2\mathbf{V}^T$$

Digression: there are some notes on the coursework materials webpage on the SVD and PCA: SVD.pdf

- Principal components regression is very similar to ridge regression: both operate on the principal components of the input matrix.

- Ridge regression shrinks the coefficients of the principal components, with relatively more shrinkage applied to the smaller components than the larger; principal components regression discards the $p - J + 1$ smallest eigenvalue components.

Ridge fit vector: $\mathbf{X}\hat{\beta} = \mathbf{X}(\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})^{-1}\mathbf{X}^T\mathbf{y}$

SVD of $\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T$

$$
\begin{aligned}
\mathbf{X}\hat{\beta} &= \mathbf{U}\mathbf{D}\mathbf{V}^T(\mathbf{V}\mathbf{D}^2\mathbf{V}^T + \lambda\mathbf{I})^{-1}\mathbf{V}\mathbf{D}\mathbf{U}^T\mathbf{y} \\
&= \mathbf{U}\mathbf{D}(\mathbf{D}^2 + \lambda\mathbf{I})^{-1}\mathbf{D}\mathbf{U}^T\mathbf{y} \\
&= \sum_{j=1}^{p} \mathbf{u}_j \frac{d_j^2}{d_j^2 + \lambda} \langle \mathbf{u}_j, \mathbf{y} \rangle
\end{aligned}
$$

$\mathbf{u}_j$ is the $j$th (standardized) principal component ($\mathbf{z}_j = \mathbf{X}\mathbf{v}_j = \mathbf{u}_j d_j$), so $\langle \mathbf{u}_j, \mathbf{y} \rangle$ is the regression coefficient of $\mathbf{y}$ on $\mathbf{u}_j$. If $\lambda = 0$, this is the OLS fit—a projection onto $\mathbf{U}$; with $\lambda > 0$, the fit is shrunk, increasingly for smaller principal components.

## Partial least squares

This technique also constructs a set of linear combinations of the $x_j$s for regression, but unlike principal components regression, it uses $\mathbf{y}$ (in addition to $\mathbf{X}$) for this construction.

- We assume that $\mathbf{y}$ is centered and begin by computing the univariate regression coefficients $\hat{\gamma}_\ell$ of $\mathbf{y}$ on each $\mathbf{x}_\ell$

- From these we construct the derived input $\mathbf{z}_1 = \sum \hat{\gamma}_\ell \mathbf{x}_\ell$, which is the first partial-least-squares direction.

- The outcome $\mathbf{y}$ is regressed on $\mathbf{z}_1$, giving coefficient $\hat{\beta}_1$, and then we orthogonalize $\mathbf{y}, \mathbf{x}_1, \ldots \mathbf{x}_p$ with respect to $\mathbf{z}_1$: $\mathbf{r}_1 = \mathbf{y} - \hat{\beta}_1 \mathbf{z}_1$, and $\mathbf{x}_\ell^* = \mathbf{x}_\ell - \hat{\theta}_\ell \mathbf{z}_1$

- We repeat this process, until $J$ directions have been obtained.

- In this manner, partial least squares produces a sequence of derived inputs or directions $\mathbf{z}_1, \mathbf{z}_2, \ldots \mathbf{z}_J$.

- As with principal components regression, if we continue on to construct $J = p$ new directions we get back the ordinary least squares estimates; use of $J < p$ directions produces a reduced regression

- Notice that in the construction of each $\mathbf{z}_j$, the inputs are weighted by the strength of their univariate effect on $\mathbf{y}$.

- It can also be shown that the sequence $\mathbf{z}_1, \mathbf{z}_2, \ldots \mathbf{z}_p$ represents the conjugate gradient sequence for computing the ordinary least squares solutions.

# Ridge vs PCR vs PLS vs Lasso

Frank and Friedman (1993) show that ridge and PCR outperform PLS in prediction, and that they are simpler to understand.

Lasso outperforms ridge when there are a moderate number of sizable effects, rather than many small effects. It also produces more interpretable models.

These are all topics for ongoing research, and have become extremely relevant with massively wide datasets.