

3 Dimension Reduction Methods

So far we have controlled variance^{of estimates} in two ways:

- ① Use a subset of original variable
 - best subset, forward/backward selection, lasso
- ② shrinking coefficients towards zero
 - ridge regression, lasso

These methods all defined using original predictor variables X_1, \dots, X_p .

We now explore a class of approaches that

- ① transform predictors
- ② then fit least squares regression model using transformed variables

We refer to these techniques as dimension reduction methods.

- ① let Z_1, \dots, Z_M represent $M < p$ linear combinations of our original predictors.

$$Z_m = \sum_{j=1}^p \phi_{jm} X_j$$

for constants ϕ_{jm} , $m=1, \dots, M$.

- ② Fit linear regression model using least squares

$$y_i = \theta_0 + \left(\sum_{m=1}^M \theta_m Z_{im} \right) + \varepsilon_i, \quad i=1, \dots, n$$

↑
regression coefficients.

If ϕ_{jm} chosen well, this can outperform least squares.

The term *dimension reduction* comes from the fact that this approach reduces the problem of estimating $p + 1$ coefficients to the problem of estimating $M + 1$ coefficients where

$$M < p.$$

$$\text{Note: } \sum_{m=1}^M \theta_m z_{im} = \sum_{m=1}^M \theta_m \sum_{j=1}^p \phi_{jm} x_{ij} = \sum_{j=1}^p \left[\sum_{m=1}^M \theta_m \phi_{jm} \right] x_{ij} = \sum_{j=1}^p [B_j] x_{ij}$$

β_j is linear combination of x_i 's

$$\theta_0, \theta_1, \dots, \theta_M$$

Dimension reduction serves to constrain β_j since now they must take a particular form.

$$\beta_j = \sum_{m=1}^M \theta_m \phi_{jm}$$

\Rightarrow special case of original linear regression problem

with β_j constrained \rightarrow can bias our coefficient estimates
 \rightarrow if $p \geq n$ (or $p \approx n$), selecting $M < p$ can reduce variance.

All dimension reduction methods work in two steps.

- ① transformed predictors are obtained.
- ② model is fit using M transformed predictors.

The selection of ϕ_{jm} 's can be done in multiple ways.

\hookrightarrow We will talk about 2 ways.

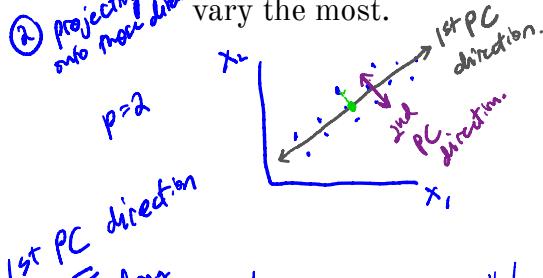
First way to choose obj's $\Rightarrow z_{11} \dots z_m$

3.1 Principle Component Regression

Principal Components Analysis (PCA) is a popular approach for deriving a low-dimensional set of features from a large set of variables.

PCA is an unsupervised approach for reducing the dimension of a $n \times p$ data matrix X .

- ① Finding PC directions
 - ② projecting data onto more directions
- The first principal component direction of the data is that along which the observations vary the most.



The 1st principal components are obtained by projecting the data onto the 1st PC direction.

↳ A point is projected onto a line by finding the point on the line closest to the original point.

1st PC direction = direction along which data varies most
2nd PC direction = direction along which data varies least
line closest to all directions! (least squares line).

out of every possible linear combination of x_1 and x_2 such that $\phi_{11}^2 + \phi_{21}^2 = 1$, choose linear combinations such that

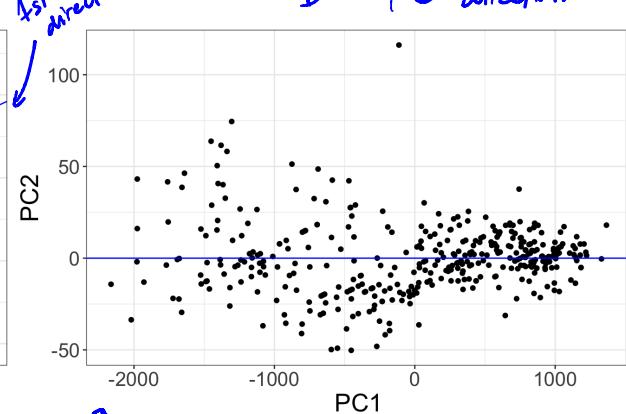
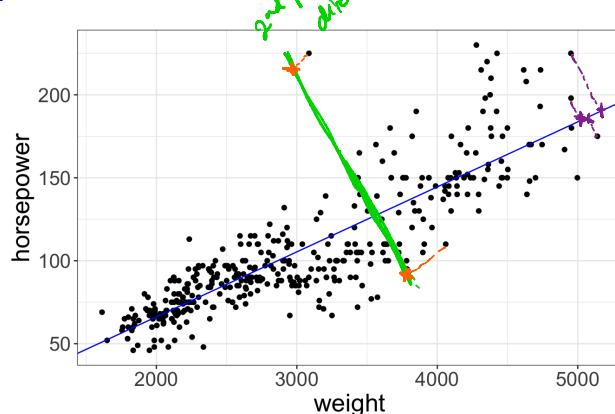
$$\text{Var} [\phi_{11}(x_1 - \bar{x}_1) + \phi_{21}(x_2 - \bar{x}_2)] \text{ is maximized.}$$

$\Rightarrow z_{ii} = \phi_{11}(x_{1i} - \bar{x}_1) + \phi_{21}(x_{2i} - \bar{x}_2)$ for $i=1, \dots, n$ are "principal component scores"
We can construct up to p principal components, where the 2nd principal component is a linear combination of the variables that are uncorrelated to the first principal component and has the largest variance subject to this constraint.

1st.

⇒ 2nd PC direction is perpendicular (orthogonal) to 1st PC direction.

From Auto data in ISLR package.



projected into 1st + 2nd PC directions

The 1st PC contains the most information \rightarrow p^{th} PC contains the least.

The Principal Components Regression approach (PCR) involves

1. Construct first M principal components Z_1, \dots, Z_M
2. fit linear regression model predicting Y using Z_1, \dots, Z_M by least squares.

Key idea: often a small # of principal component(s) will suffice to explain most of variability in the data X , as well as the relationship with the response.

In other words, we assume that the directions in which X_1, \dots, X_p show the most variation are the directions that are associated with Y .

This is not guaranteed to be true, but often works well in practice.

If this assumption holds, fitting PCR will lead to better results than fitting least squares on X_1, \dots, X_p .

(We can mitigate overfitting).

How to choose M , the number of components?

Note: PCR is not feature selection!

3.2 Partial Least Squares

The PCR approach involved identifying linear combinations that best represent the predictors X_1, \dots, X_p .

Consequently, PCR suffers from a drawback

Alternatively, *partial least squares (PLS)* is a supervised version.

Roughly speaking, the PLS approach attempts to find directions that help explain both the response and the predictors.

The first PLS direction is computed,

To identify the second PLS direction,

As with PCR, the number of partial least squares directions is chosen as a tuning parameter.

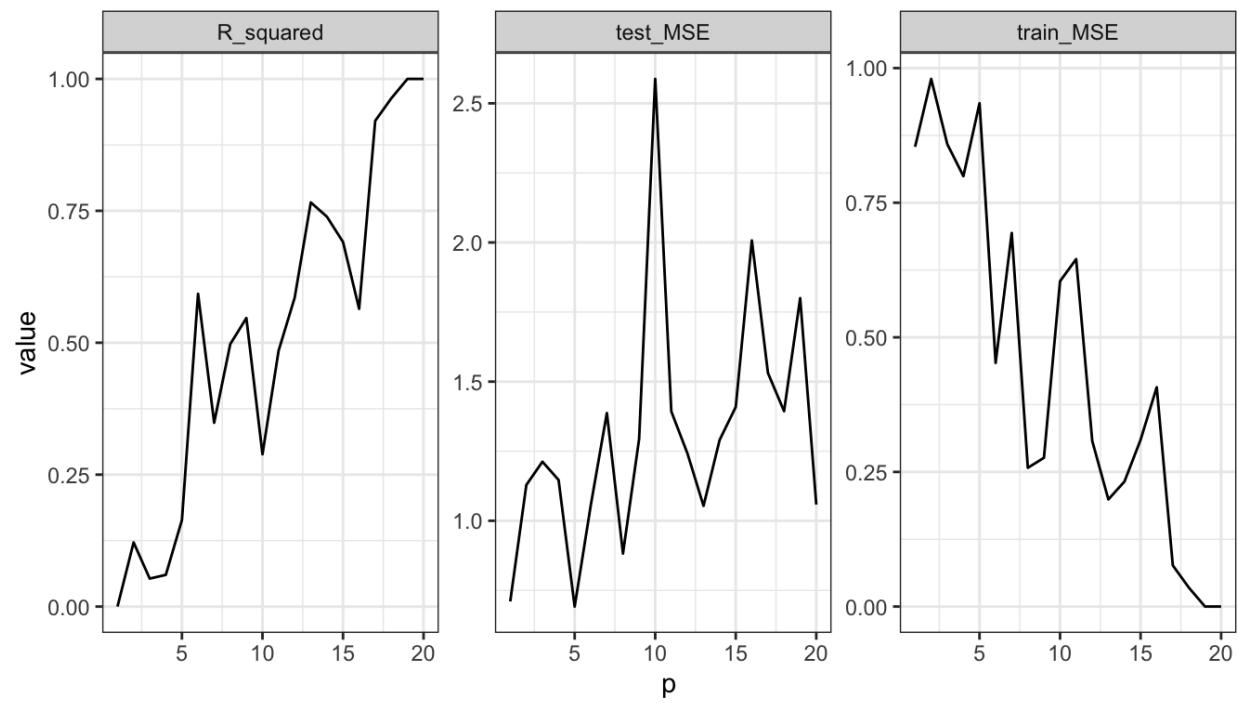
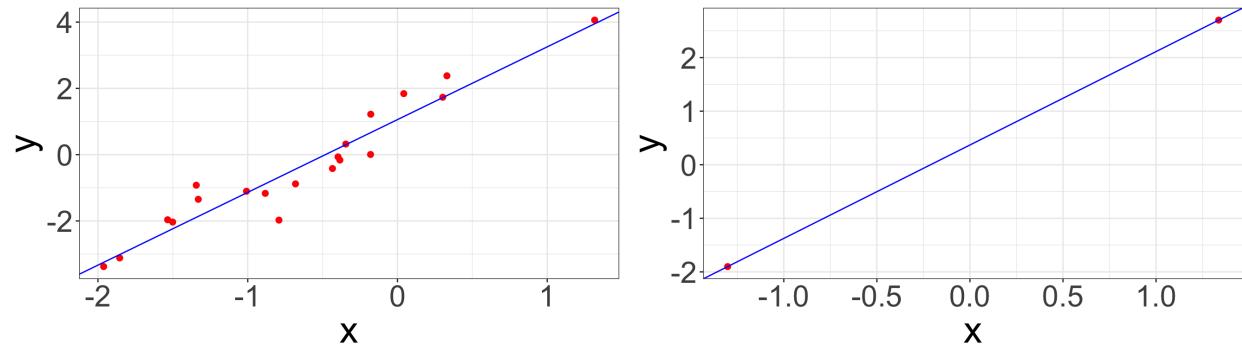
4 Considerations in High Dimensions

Most traditional statistical techniques for regression and classification are intended for the low-dimensional setting.

In the past 25 years, new technologies have changed the way that data are collected in many fields. It is not commonplace to collect an almost unlimited number of feature measurements.

Data sets containing more features than observations are often referred to as *high-dimensional*.

What can go wrong in high dimensions?



Many of the methods that we've seen for fitting *less flexible* models work well in the high-dimension setting.

1.

2.

3.

When we perform the lasso, ridge regression, or other regression procedures in the high-dimensional setting, we must be careful how we report our results.