

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

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

↑
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

$$M < p.$$

$$\beta_0, \beta_1, \dots, \beta_p$$

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

$z = \text{linear combination of } x\text{'s}$

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

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

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

⇒ special case of original linear regression problem

with β_j constrained → can bias our coefficient estimates
 → if $p > n$ (or $p \approx n$), selecting $M \ll p$ 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.

↳ We will talk about 2 ways.

First way to choose ϕ_{im} 's $\Rightarrow z_{11} \dots z_{1n}$

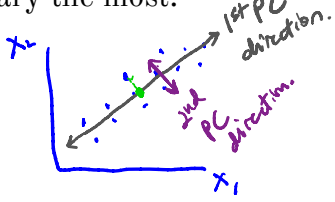
3.1 Principle Component Regression

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 into those directions

The first principal component directions 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.

\hookrightarrow 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 = line closest to all observations! (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 combination such that

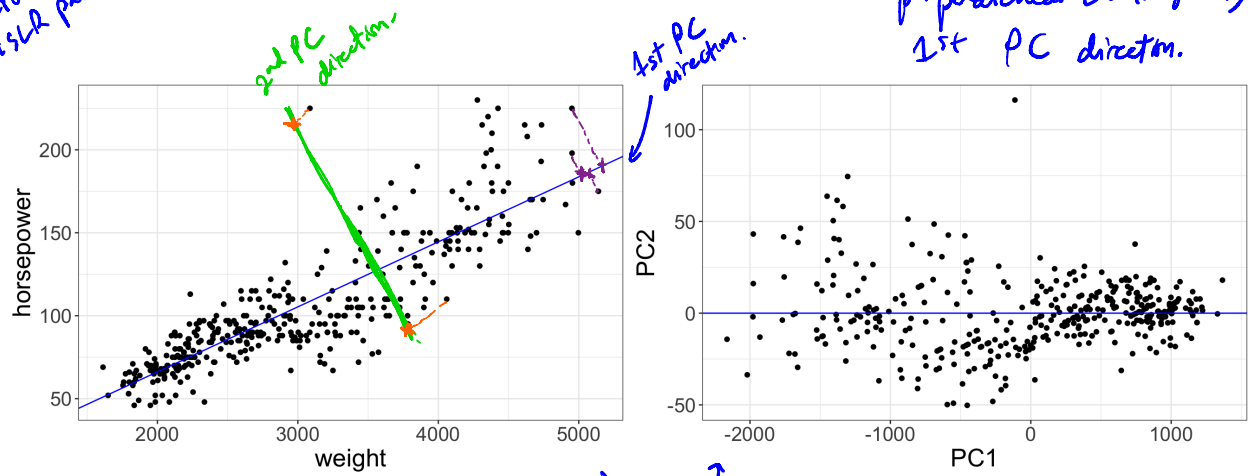
$$\text{Var} [\phi_{11}(x_1 - \bar{x}_1) + \phi_{21}(x_2 - \bar{x}_2)] \text{ is maximized.}$$

$$\Rightarrow z_{i1} = \phi_{11}(x_{1i} - \bar{x}_1) + \phi_{21}(x_{2i} - \bar{x}_2) \text{ for } i=1, \dots, n \text{ are "principal component scores" } 1^{st}.$$

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.

\Rightarrow 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 components 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?

M is tuning parameter,
use C.V.

Note: PCR is not feature selection!

Z_i depend on all X 's.

PCR is not sparse model

PCR more like ridge than lasso.

3.2 Partial Least Squares

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

$$Z_1, \dots, Z_m \quad Z_m = \sum_{j=1}^p a_{jm} X_j$$

Consequently, PCR suffers from a drawback

Explains var in X , not necessarily Y .

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

$$Z_m = \sum \phi_{jm} X_j$$

Determine ϕ with both X and y

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

ϕ , linear combination

The first PLS direction is computed,

① Standardize X $\text{var}(X_j) = 1$

② Set ϕ_{j1} according to simple regression $y \sim X_j$
 $\phi_{j1} = \beta$

Places more weight on Predictive X 's

To identify the second PLS direction,

① Regress each $X_j \sim Z_1$, take residuals

② Repeat above process, but with residuals in stead of X 's.

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

M is a tuning parameter
 Use C.V. to find M .

Neither PCR nor PLS is "better"

4 Considerations in High Dimensions

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

$n \gg p$. Throughout history, low dimension was most prevalent.

Ex: $y = \text{blood pressure}$, $X = \text{age, gen, bmi}$ $p = 3$
 $n = 100$

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.

Ex: Blood pressure. Now, we may have X 's for every SNP (gene sequencing)
 $p \gg n$

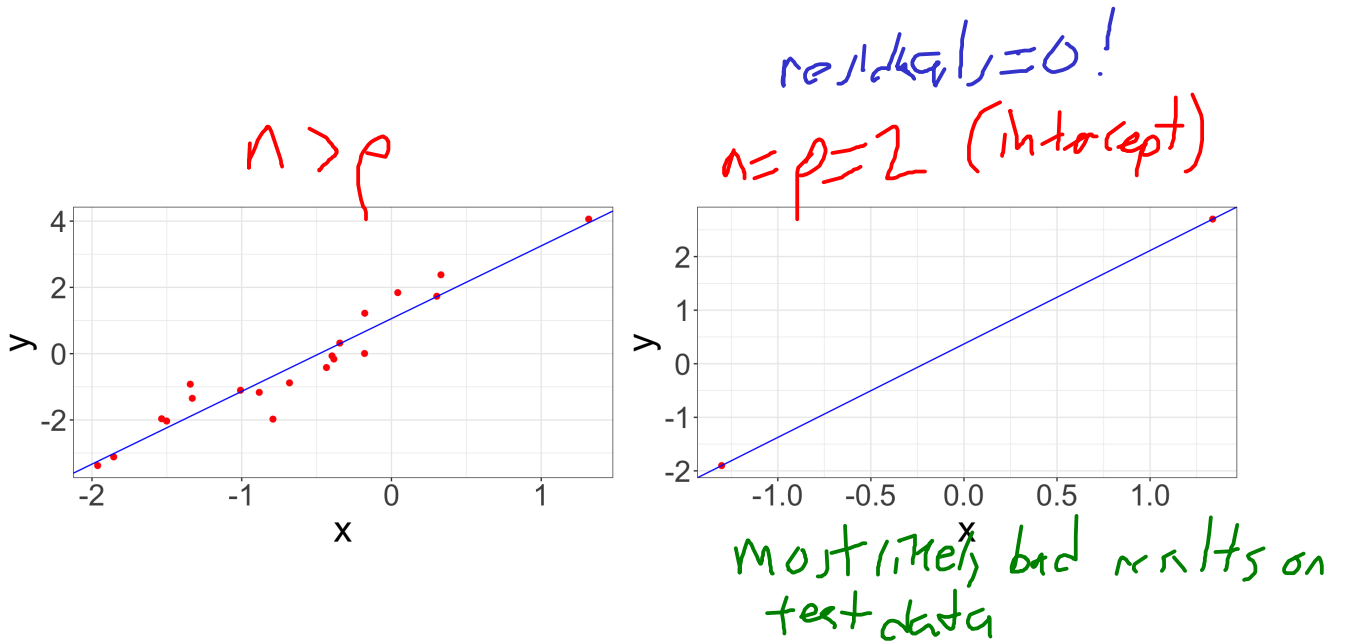
Ex: Shopping patterns
 X for every other purchase
 X for shopping carts, wish list
 $p \gg n$

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

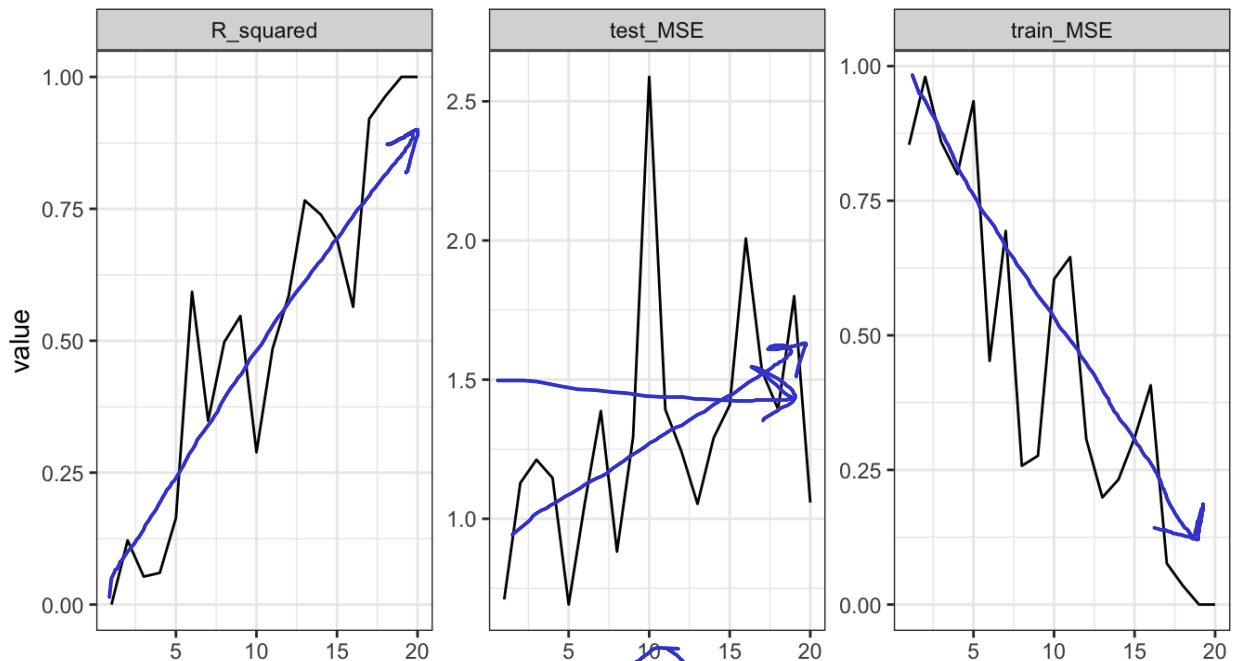
Least squares + other standard methods may not work here.

Need to be careful when $n \approx p$ or $n \ll p$

What can go wrong in high dimensions?



1. Be careful.
2. Always check test set performance



R^2 looks good!

Test set looks meh

Training set looks good!

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

1. Regularization or shrinkage (lasso, ridge)
2. Appropriate tuning parameter (r.v.)
3. Test error \uparrow $\rho \uparrow$ unless new predictors are totally associated w/ y . \leftarrow curse of dimensionality
 Adding predictors will improve training perf. at the cost of more variance. \Rightarrow test error \uparrow
 Risk of overfitting

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.

In high dimensional setting, it is more likely that variables will be highly correlated.
 \Rightarrow any variable in the model could be written as a linear combination of other variables in the model.

This means we can never really know if any variables are truly predictive of the response.

\Rightarrow we can never identify which are best variables to include.

at best, we can only hope to assign large regression coefficients to variables that are highly correlated to variables that are truly predictive of the response.

$\star \Rightarrow$ When we use lasso/feature selection, etc. we should be clear that we have identified one of many possible models for predicting the response and should be validated on many independent data sets. (reproducibility)

\star also important to report test errors (not R^2 , training errors) because we know $R^2 \uparrow$ as $p \uparrow$ but this doesn't mean we have a good model.