# 3 Dimension Reduction Methods

So far we have controlled variance in two ways:

① Used a subset of original variables
- best subset, forward/backward selection, lasso

② shrinking coefficients towards zero
- ridge regression, lasso.

These methods all defined using original predictor variables $x_1, \ldots, x_p$.

We now explore a class of approaches that

① transform the predictors

② then perform least squares using transformed variables.

We refer to these techniques as *dimension reduction* methods.

① Let $z_1, \ldots, z_M$ represent $M < p$ linear combinations of our original predictors.

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

for constants $\phi_{1m}, \ldots, \phi_{pm} \quad m = 1, \ldots, M$

② Fit the linear regression model using least squares

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

↖ ↗ regression coefficients.

If $\left\{ \phi_{jm} \right\}_{\substack{j=1, \ldots, p \\ m=1, \ldots, M}}$ chosen well, this can <u>outperform</u> 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 $M < p$.

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

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

NOTE:

from least squares → have to come from somewhere else (hopefully we've picked well).

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

$$= \sum_{j=1}^{p} \beta_j x_{ij}$$

Dimension reduction serves to constrain $\beta_j$, 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 $\beta_j$ constrained)
  ↳ if can bias coefficient estimate
  ↳ if $p > n$ ($\overset{or}{p} \approx n$) selecting $M \ll p$ can reduce variance.

All dimension reduction methods work in two steps.

① transformed predictors are obtained ( get $\{ \phi_{jm} \}_{\substack{j=1,\dots,p \\ m=1,\dots,M}}$ )

② Model is fit using $M$ transformed predictors from ①.
   linear

selection of $\{ \phi_{jm} \}_{\substack{j=1,\dots,p \\ m=1,\dots,M}}$   can be done in multiple ways.
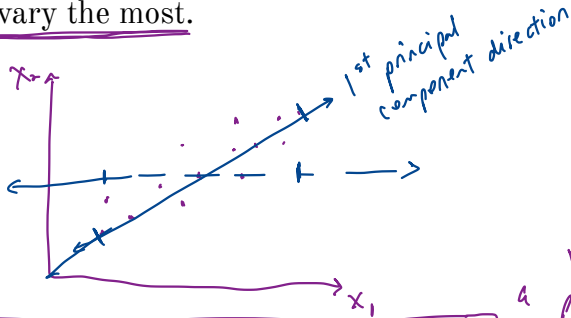
We will talk about ②.

One way to choose $Z_{(1)}...Z_m$

# 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.
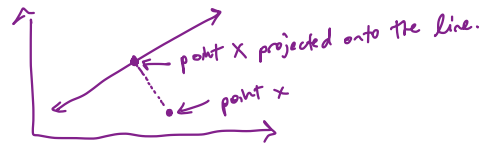
PCA is an <u>unsupervised approach</u> for reducing the dimension of an $n \times p$ data matrix $X$.

The *first principal component* directions of the data is that along which the obervations <u>vary the most.</u>



1st principal component direction

The 1st principal components are obtained by <u>projecting</u> the data onto the 1st principal component direction.

A <u>point is projected</u> onto a line by finding the point on the line closest to the original point.
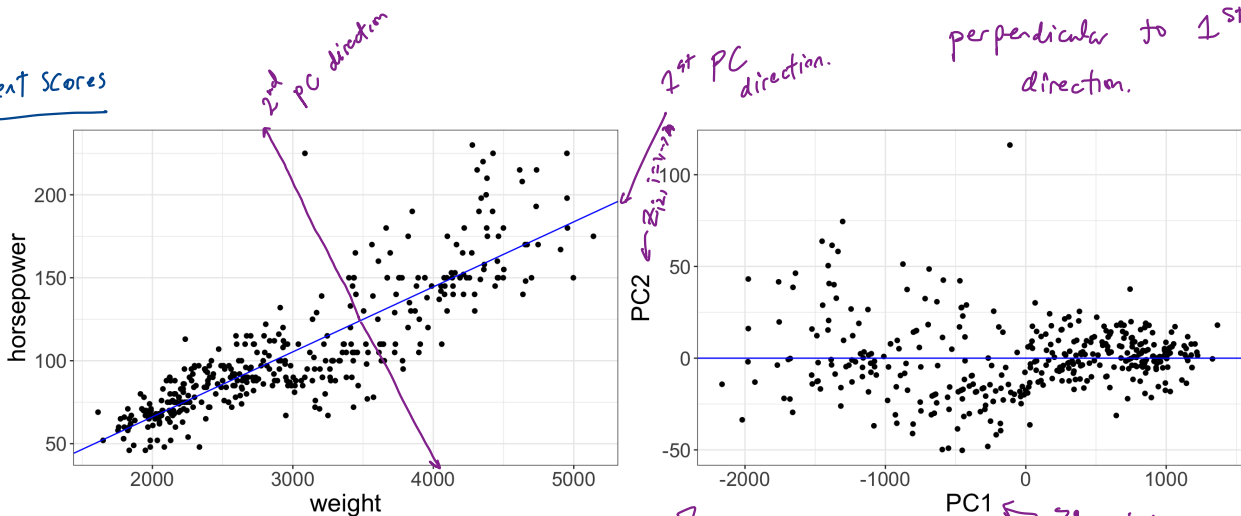
out of every possible linear combination of $X_1$ and $X_2$ such that $\phi_{1r}^2 + \phi_{2r}^2 = 1$, choose $\phi_{11}$ s.t. $Var\left(\phi_{11}(X_1-\bar{X_1}) + \phi_{21}(X_2-\bar{X_2})\right)$ is maximized.

point X projected onto the line.

point x

$Z_{i1} = \phi_{11}(x_{i1}-\bar{x_1}) + \phi_{21}(x_{2i}-\bar{x_2})$
for $i=1,...,n$
are called 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 <u>uncorrelated to the first principal component</u> and has the largest variance subject to this constraint.

$\Rightarrow$ 2nd PC direction is perpendicular to 1st PC direction.

1st PC direction = dimension along which data vary the most.

2nd PC direction

1st PC direction.

$\leftarrow Z_{i1}, i=1,...,n$



projected onto principal component directions.

1st PC contains the most information $\longrightarrow$ $p^{th}$ PC contains the least.

The Principal Components Regression approach (PCR) involves

1. Construct first $M$ principal components $Z_1, ..., Z_M$ *a choice we are making.*
   *score vectors.*

2. Fit a linear regression model w/ $Z_1, ..., Z_M$ as predictors using least squares.

Key idea: Often a small # of PC suffice to explain most of the variability in the data, as well as the relationship w/ predictor.

In other words, we assume that the directions in which $X_1, ..., X_p$ show the most varia-
tion 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 model on $X_1, ..., X_p$ because we can mitigate overfitting.

How to choose $M$, the number of components?

$M$ can be thought of as a tuning parameter
$\Rightarrow$ use CV method to choose!

as $M \uparrow p$, PCR $\rightarrow$ least squares. $\Rightarrow$ bias $\downarrow$ but variance $\uparrow$, will see bias-variance trade-off in the form of a U-shape in the test MSE.

Note: PCR is not feature selection!

each of the $M$ principal components used in the linear regression is a linear combination of all $p$ of the original predictors!

$\Rightarrow$ while PCR works well to reduce variance, it doesn't give us a sparse model.

PCR more like ridge regression than the lasso. (not going to help w/ interpretation)

NOTE: recommended standardizing predictors $X_1, ..., X_p$ to each have st. dev. = 1 before getting the PCs.

## 3.2 Partial Least Squares

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

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 reponse 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.
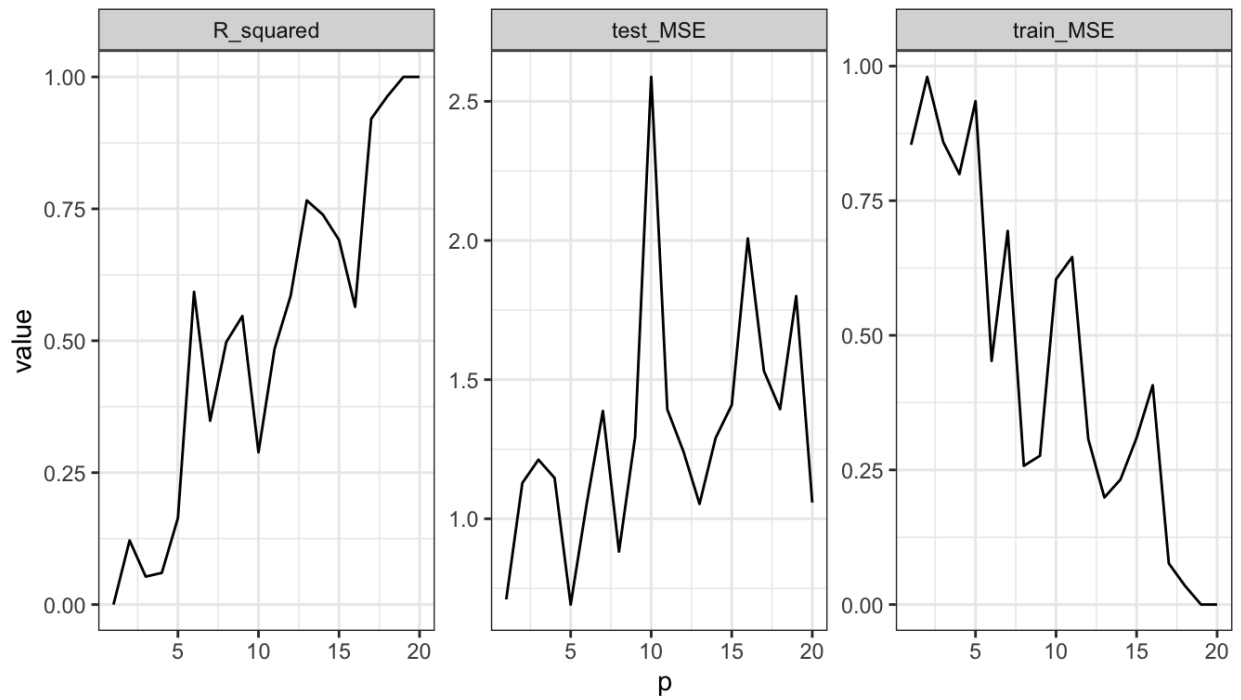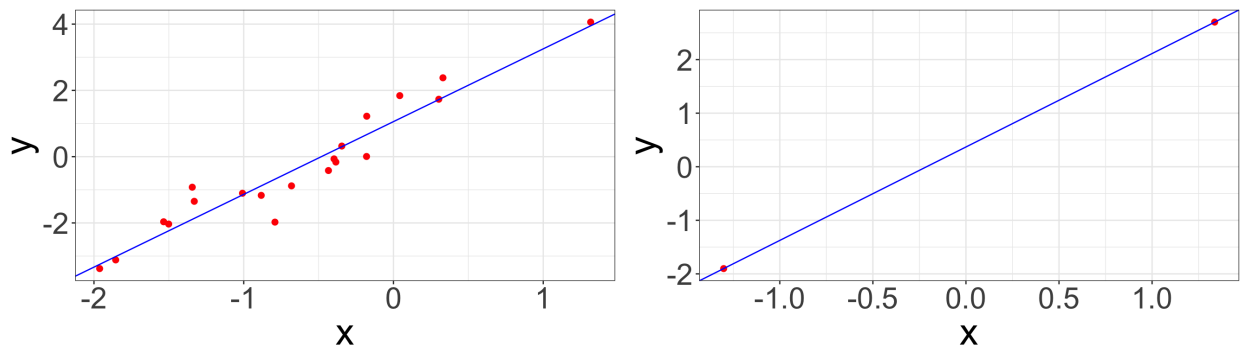
# 4 Considerations in High Dimensions

Most traditional statistical techniques for regression and classification are intendend 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 methds 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.