

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, \dots, 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, \dots, z_M represent $M < p$ linear combinations of our original predictors.

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

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

- ② Fit the linear regression model using least squares

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

↙ ↘
regression coefficients.

If $\left\{ \phi_{jm} \right\}_{\substack{j=1, \dots, p \\ m=1, \dots, M}}$ 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 $M < p$.

NOTE:

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

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

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

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

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)

- ↳ it can bias coefficient estimate
- ↳ 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 (get $\{\phi_{jm}\}_{j=1, \dots, p}^{m=1, \dots, M}$)
 - ② ^{linear} Model is fit using M transformed predictors from ①.
- selection of $\{\phi_{jm}\}_{j=1, \dots, p}^{m=1, \dots, M}$ can be done in multiple ways.
We will talk about 2.

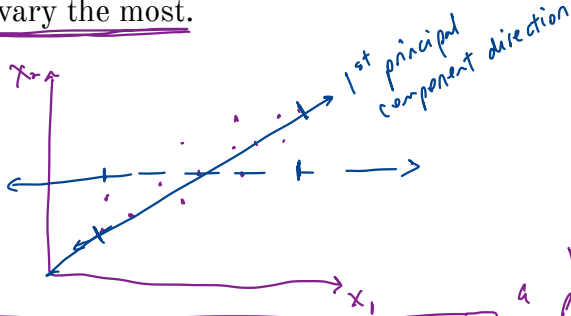
One way to choose $Z_{(1) \dots} Z_M$

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.

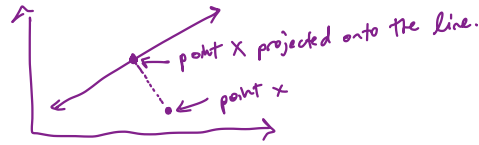
PCA is an unsupervised approach for reducing the dimension of an $n \times p$ data matrix X .

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 principal component direction.

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



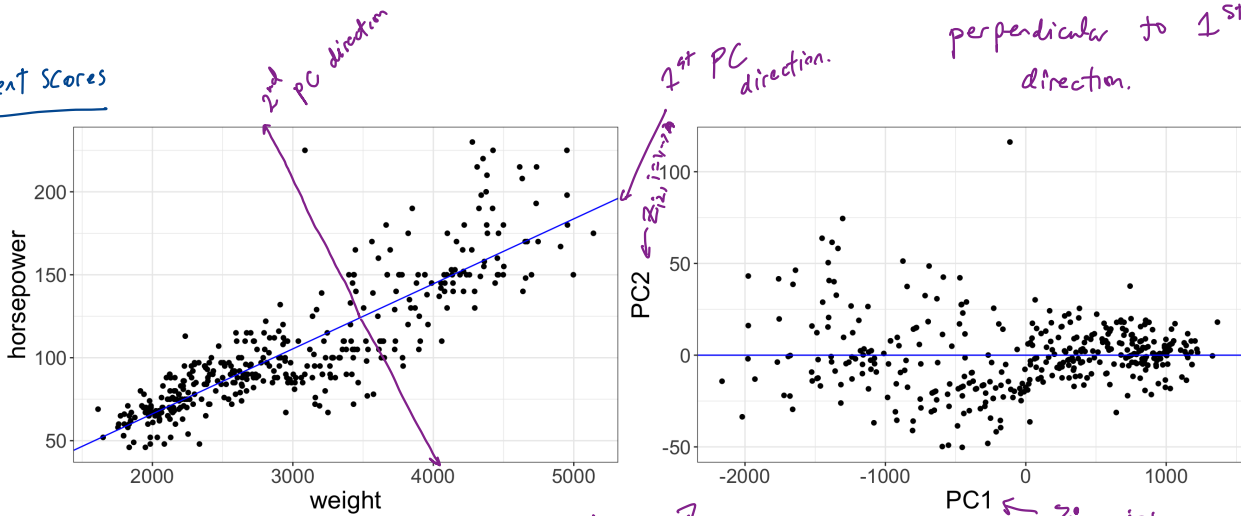
out of every possible linear combination of x_1 and x_2 such that $\phi_{11}^2 + \phi_{21}^2 = 1$, choose ϕ_{i1} s.t. $\text{Var}(\phi_{11}(x_1 - \bar{x}_1) + \phi_{21}(x_2 - \bar{x}_2))$ is maximized.

$Z_{i1} = \phi_{11}(x_{i1} - \bar{x}_1) + \phi_{21}(x_{i2} - \bar{x}_2)$ 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 to 1st PC direction.

for $i=1, \dots, n$ are called 1st principal component scores

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



projected onto principal component directions.

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 score vectors. Z_1, \dots, Z_M ← a choice we are making.
2. Fit a linear regression model w/ Z_1, \dots, 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, \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 model on X_1, \dots, 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 $M \uparrow$, 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, \dots, X_p to each have st. dev. = 1 before getting the PCs.

3.2 Partial Least Squares

can do this using pls function in pls package.

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

We identified these directions in an unsupervised way (response Y is not used to help us determine the directions).

Consequently, PCR suffers from a drawback

There is no guarantee that the direction that best explains the predictors will also be the best directions to explain the relationship between predictors and response.

Alternatively, partial least squares (PLS) is a supervised version of dimension reduction

① identify new features Z_1, \dots, Z_m linear combination of original predictors

② fit linear model (least squares) using transformed predictors.

PLS also going to use Y (not just X) to find linear combinations of X_1, \dots, X_p i.e. uses $Y \text{ \& } X$ to find $\phi_{1m}, \phi_{2m}, \dots, \phi_{pm}$ $m=1, \dots, M$.

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

The first PLS direction is computed,

① standardize the p predictors (all have std. dev. = 1).

② set each ϕ_{j1} equal to the coefficient from a simple linear regression $Y \sim X_j$

Since the coefficient from SLR of $Y \sim X_j$ is $\text{Cor}(Y, X_j)$ PLS places highest weight on variables that are strongly related (linearly) to the response.

To identify the second PLS direction, $X_j \sim Z_1$

① regress each variable X_1, \dots, X_p on Z_1 and get residuals ($r_{ji} = X_{ji} - \hat{X}_{ji}$ $i=1, \dots, n$ $j=1, \dots, p$)

② compute Z_2 by setting each ϕ_{j2} equal to coefficient from SLR $Y \sim r_{ji}$ ← residuals from step ①

The residuals $r_{11}, \dots, r_{p1} \approx$ remaining information not explained by 1st PLS direction.

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

Generally standardize the predictors and response before performing PLS.

In practice, PLS usually performs no better than ridge or PCR.

\hookrightarrow supervised nature of problem does reduce, but often increases variance. \Rightarrow not always better.

4 Considerations in High Dimensions

Most traditional statistical techniques for regression and classification are intended for the low-dimensional setting. $n \gg p$

This is because throughout the history of the field, the bulk of scientific problems requiring statistics have been low dimensional.

e.g. ag field trials.

In the past [?]25 years, new technologies have changed the way that data are collected in many fields. It is now commonplace to collect an almost unlimited number of feature measurements. (p very large).

But n can still be limited due to cost, sampling availability, etc.

e.g. Could predict BP on age, gender, BMI, and also collect half million SNPs

now $p \approx 500,000$ but SNPs are expensive to collect, maybe only get ~ 200 of them.

e.g. Consider trying to model online shopping patterns. We could treat all search terms in a person's month-long browsing history as features in a "bag-of-words" model.

But we may only have a few thousand people who have consented to use their history.

For a given user the features would be absence (0) or presence (1) of each potential search term. $\Rightarrow p$ large but $n \approx 2000$

Data sets containing more features than observations are often referred to as high-dimensional.
 $p > n$

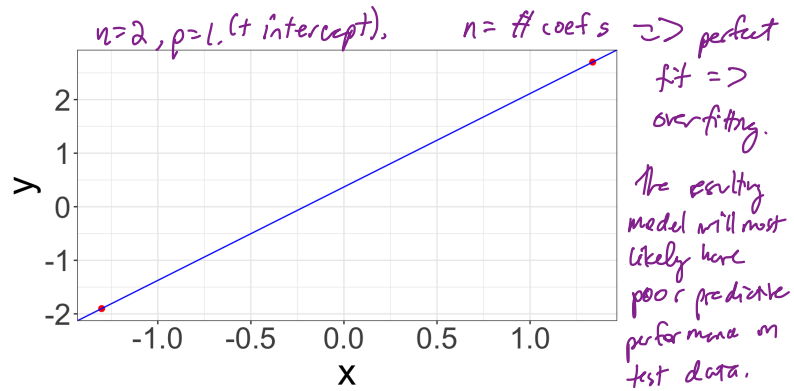
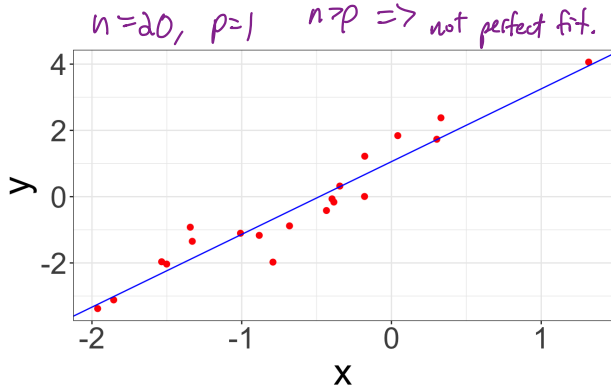
classical approaches (like least squares) are not appropriate in this setting.
(why? think bias-variance trade-off and overfitting).

\Rightarrow we need to be careful when $n \approx p$ or $n < p$.

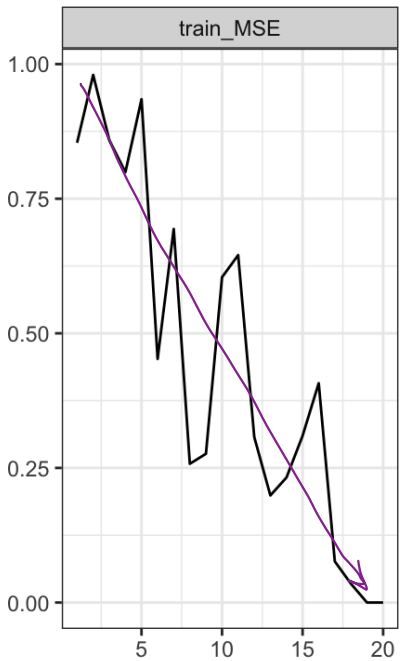
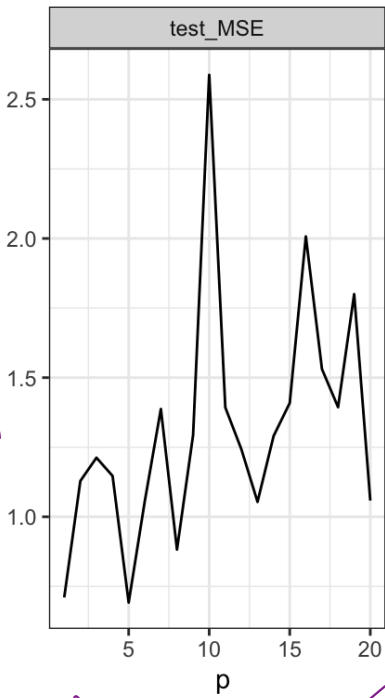
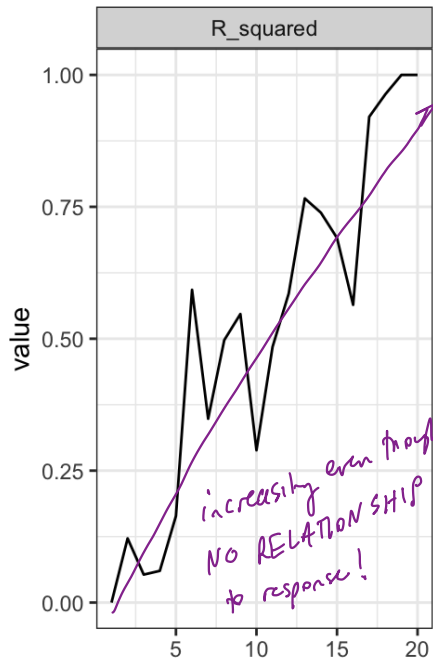
What can go wrong in high dimensions?

Going the talk about least squares but same issues arise w/ logistic regression or LDA

If $p \approx n$ or $p > n$ regardless of if there is a relationship w/ response least squares will yield a set of coefficients that is an (almost if $n < p$ but $n \leq p$) perfect fit. $p \geq n \Rightarrow$ residuals = 0.



Simulated data w/ $n=20$ and regression performed with between 1 and 20 features. Features generated w/ no relationship to response



test MSE never very good b/c not a good predictive fit (no relationship)

Note: we didn't see methods for adjusting R^2 to better estimate test MSE: C_p , AIC, BIC, adj R^2
BUT: in high dim settings, we can't compute. (linear models?)

\Rightarrow we need to be very careful when analyzing data w/ many predictors. Always need to evaluate model performance on indep test set

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

Key points

1. regularization or shrinkage plays a key role in high dimensional problems.
2. appropriate tuning parameter selection is critical for good predictive performance.
3. The test error tends to \uparrow as $p \uparrow$ UNLESS the additional features are truly associated w/ response.

This is related to the curse of dimensionality

adding noise will deteriorate our fitted model $\Rightarrow \uparrow$ test error
vs. adding signal features will improve our model (fitted).

(\uparrow dimension $\Rightarrow \uparrow$ risk of overfitting due to noise.)

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 predictors will be highly correlated
 \Rightarrow any variable in the model could be written (almost) as a linear combination of other variables in the model.

This means we can never really know if any are truly predictive of the response.
 \Rightarrow we can never identify which are the best 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.

★ \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 test data sets.

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