

# Day 7: Resampling methods and model selection

ME314: Introduction to Data Science and Machine Learning

LSE Methods Summer Programme 2019

7 August 2019

# Day 7 Outline

## Cross-validation

- Validation-set approach

- K-fold Cross-validation

## Bootstrap

- Bootstrap mechanism

- Bootstrap and prediction

## Linear Model Selection and Regularization

## Shrinkage Methods

## Dimensionality Reduction Methods

## **Cross-validation**

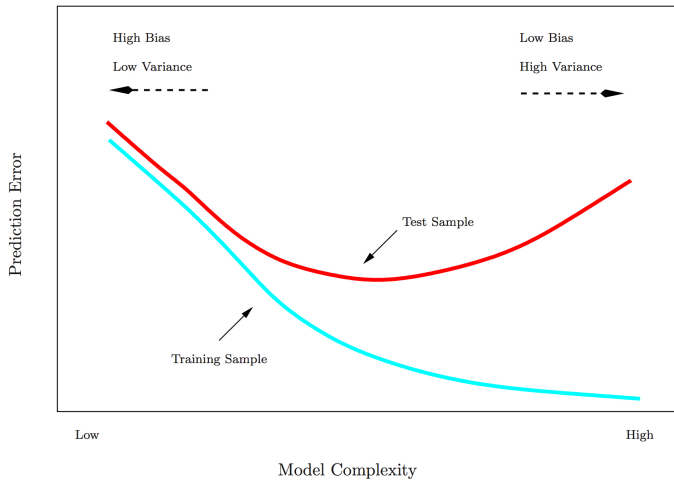
# Resampling

- ▶ Today we discuss two resampling methods: cross-validation and the bootstrap.
- ▶ These methods refit a model of interest to samples formed from the training set, in order to obtain additional information about the fitted model.
- ▶ E.g., they provide estimates of test-set prediction error, and the standard deviation and bias of our parameter estimates.

# Training Error versus Test error

- ▶ The **test error** is the average error that results from using a statistical learning method to predict the response on a new observation, one that was not used in training the method.
- ▶ In contrast, the **training error** can be easily calculated by applying the statistical learning method to the observations used in its training.
- ▶ Training error rate often is quite different from the test error rate, and in particular the former can **dramatically underestimate** the latter.

# Training- versus Test-Set Performance



## More on prediction-error estimates

- ▶ Best solution: a large designated test set. Often not available.
- ▶ Some methods make a **mathematical adjustment** to the training error rate in order to estimate the test error rate. These include the **Cp statistic**, **AIC** and **BIC**.
- ▶ First, however, we consider a class of methods that estimate the test error by holding out a subset of the training observations from the fitting process, and then applying the statistical learning method to those held out observations.

# Validation-set approach

- ▶ We randomly divide the available set of samples into two parts: a **training set** and a **validation** or **hold-out set**.
- ▶ The model is fit on the training set, and the fitted model is used to predict the responses for the observations in the validation set.
- ▶ The resulting validation-set error provides an estimate of the test error. This is typically assessed using MSE in the case of a quantitative response and misclassification rate for qualitative response models.



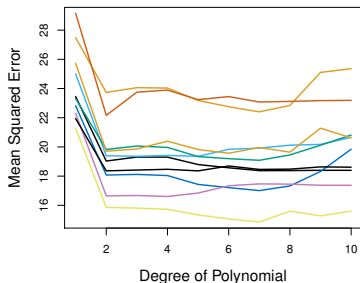
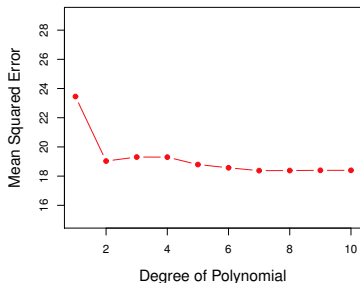
# The Validation process



A random splitting into two halves: left part is training set, right part is validation set.

## Example

- ▶ We want to compare linear vs higher-order polynomial terms in a linear regression with our Auto dataset.
- ▶ We randomly split the 392 observations into two sets, a training set containing 196 of the data points, and a validation set containing the remaining 196 observations.



Left panel shows single split; right panel shows multiple splits.

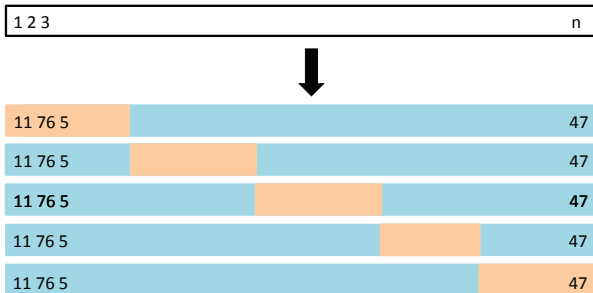
## Drawbacks of validation set approach

- ▶ The validation estimate of the test error can be highly variable, depending on precisely which observations are included in the training set and which observations are included in the validation set.
- ▶ In the validation approach, only a subset of the observations – those that are included in the training set rather than in the validation set – are used to fit the model.
- ▶ This suggests that the validation set error may tend to **overestimate** the test error for the model fit on the entire dataset.

# K-fold Cross-validation

- ▶ Very popular approach for estimating test error.
- ▶ Estimates can be used to select best model, and to give an idea of the test error of the final chosen model.
- ▶ Idea is to randomly divide the data into  $K$  equal-sized parts. We leave out part  $k$ , fit the model to the other  $K - 1$  parts (combined), and then obtain predictions for the left-out  $k$ th part.
- ▶ This is done in turn for each part  $k = 1, 2, \dots, K$ , and then the results are combined.

## 5-fold CV



# Mechanism

- ▶ Let the  $K$  parts be  $C_1, C_2, \dots, C_K$ , where  $C_K$  denotes the indices of the observations in part  $k$ . There are  $n_k$  observations in part  $k$ : if  $N$  is a multiple of  $K$ , then  $n_k = n/K$ .
- ▶ Compute

$$CV_{(K)} = \sum_{k=1}^K \frac{n_k}{n} \text{MSE}_k$$

where  $\text{MSE}_k = \sum_{i \in C_k} (y_i - \hat{y}_i)^2 / n_k$ , and  $\hat{y}_i$  is the fit for observation  $i$ , obtained from the data with part  $k$  removed.

- ▶ Setting  $K = n$  yields  $n$ -fold or **leave-one out cross-validation** (LOOCV).

# LOOCV

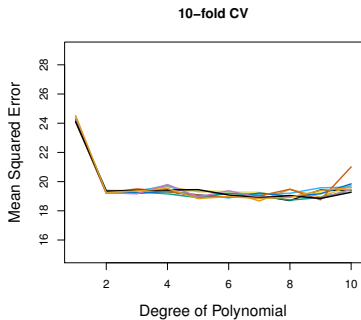
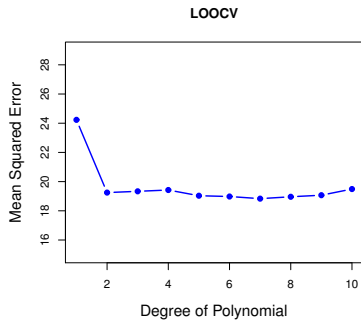


# LOOCV

- ▶ The estimates from each fold are highly correlated and hence their average can have high variance.
- ▶ A better choice is  $K = 5$  or  $10$ .



# Auto data example



## Additional issues with CV

- ▶ Since each training set is only  $(K - 1)/K$  as big as the original training set, the estimates of prediction error will typically be biased upward.
- ▶ This bias is minimized when  $K = n$  (LOOCV), but this estimate has high variance, as noted earlier.
- ▶  $K = 5$  or  $10$  provides a good balance for this bias-variance tradeoff.

## CV for classification

- ▶ We divide the data into  $K$  roughly equal-sized parts  $C_1, C_2, \dots, C_K$ .  $C_k$  denotes the indices of the observations in part  $k$ . There are  $n_k$  observations in part  $k$ : if  $n$  is a multiple of  $K$ , then  $n_k = n/K$ .
- ▶ Compute

$$CV_K = \sum_{k=1}^K \frac{n_k}{n} \text{Err}_k$$

where  $\text{Err}_k = \sum_{i \in C_k} I(y_i \neq \hat{y}_i) / n_k$ .

- ▶ The estimated standard deviation of  $CV_K$  is

$$\widehat{\text{SE}}(CV_K) = \sqrt{\sum_{k=1}^K (\text{Err}_k - \bar{\text{Err}})^2 / (K - 1)}$$

## CV: right and wrong

- ▶ Consider a simple classifier applied to some two-class data:
  1. Starting with 5000 predictors and 50 samples, find the 100 predictors having the largest correlation with the class labels.
  2. We then apply a classifier such as logistic regression, using only these 100 predictors.
- ▶ How do we estimate the test set performance of this classifier?
- ▶ Can we apply cross-validation in step 2, ignoring step 1?

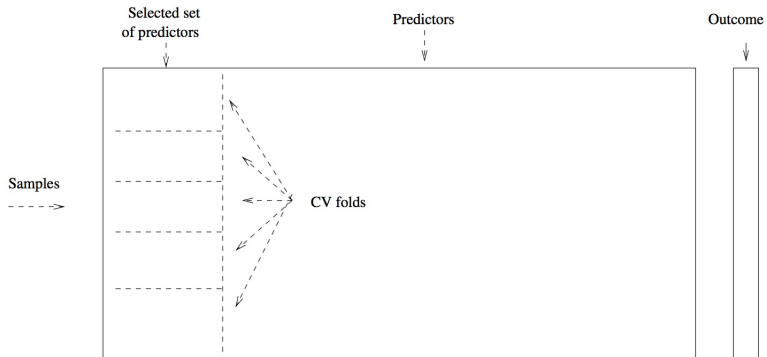
## CV: right and wrong

- ▶ This would ignore the fact that in Step 1, the procedure has already seen the labels of the training data, and made use of them. This is a form of training and must be included in the validation process.
- ▶ It is easy to simulate realistic data with the class labels independent of the outcome, so that true test error = 50%, but the CV error estimate that ignores Step 1 is zero.

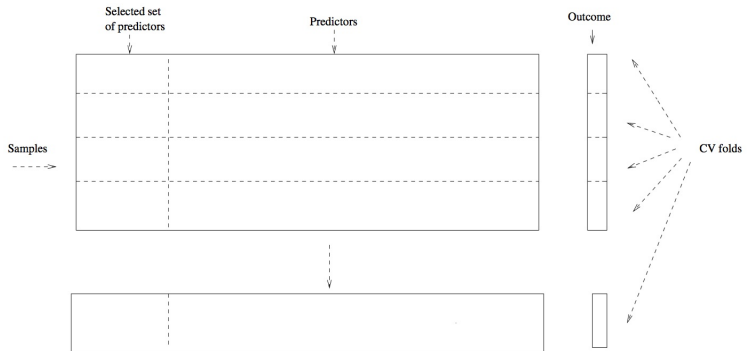
## CV: right and wrong

- ▶ **Incorrect:** Apply cross-validation in step 2.
- ▶ **Correct:** Apply cross-validation to steps 1 and 2.

# Wrong



# Right





**Bootstrap**

# Bootstrap

- ▶ The **bootstrap** is a flexible and powerful statistical tool that can be used to quantify the uncertainty associated with a given estimator or statistical learning method.
- ▶ E.g., it can provide an estimate of the standard error of a coefficient, or a confidence interval for that coefficient.

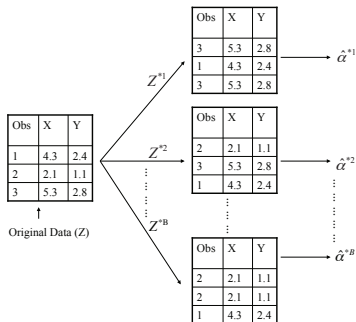
# Where does the name come from?

- ▶ The use of the term bootstrap derives from the phrase to pull oneself up by one's bootstraps, widely thought to be based on one of the eighteenth century "The Surprising Adventures of Baron Munchausen" by Rudolph Erich Raspe:  
The Baron had fallen to the bottom of a deep lake. Just when it looked like all was lost, he thought to pick himself up by his own bootstraps.
- ▶ It is not the same as the term "bootstrap" used in computer science meaning to "boot" a computer from a set of core instructions, though the derivation is similar.

# Intuition

- ▶ We usually care about uncertainty around our estimates from a sample. One way would be to repeatedly sample the population. But we cannot do that in real world.
- ▶ The bootstrap approach allows us to replicate this process of obtaining new data sets, so that we can estimate the variability of our estimate without generating additional samples.
- ▶ Rather than repeatedly obtaining independent data sets from the population, we instead obtain distinct data sets by repeatedly sampling observations from the original data set with replacement.
- ▶ Each of these “bootstrap data sets” is created by sampling with replacement, and is the same size as our original dataset. As a result some observations may appear more than once in a given bootstrap data set and some not at all.

# Example with three observations



- ▶ A graphical illustration of the bootstrap approach on a small sample containing  $n = 3$  observations.
- ▶ Each bootstrap dataset contains  $n$  observations, sampled with replacement from the original data set.
- ▶ Each bootstrap data set is used to obtain an estimate of  $\alpha$ .

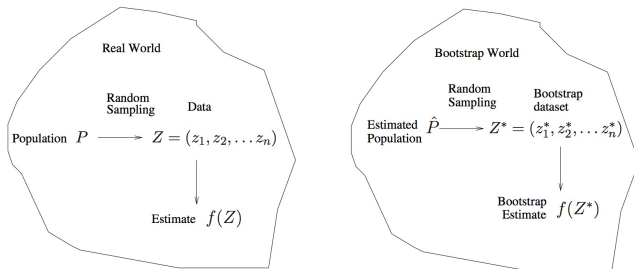
# Mechanism

- ▶ Denoting the first bootstrap dataset by  $Z^{*1}$ , we use  $Z^{*1}$  to produce a new bootstrap estimate for  $\alpha$ , which we call  $\hat{\alpha}^{*1}$ .
- ▶ This procedure is repeated  $B$  times for some large value of  $B$  (say 1000 or 10,000), in order to produce  $B$  different bootstrap datasets,  $Z^{*1}, Z^{*2}, \dots, Z^{*B}$ , and  $B$  corresponding  $\alpha$  estimates,  $\hat{\alpha}^{*1}, \hat{\alpha}^{*2}, \dots, \hat{\alpha}^{*B}$ .
- ▶ We estimate the standard error of these bootstrap estimates using the formula

$$SE_B(\hat{\alpha}) = \sqrt{\frac{1}{B-1} \sum_{r=1}^B (\hat{\alpha}^{*r} - \bar{\hat{\alpha}}^*)^2}$$

- ▶ This serves as an estimate of the standard error of  $\hat{\alpha}$  estimated from the original data set.

# A general picture for the bootstrap



# The bootstrap in general

- ▶ In more complex data situations, figuring out the appropriate way to generate bootstrap samples can require some thought.
- ▶ For example, if the data is a time series, we can't simply sample the observations with replacement.
- ▶ We can instead create blocks of consecutive observations, and sample those with replacements. Then we paste together sampled blocks to obtain a bootstrap dataset.



## Other uses of the bootstrap

- ▶ Primarily used to obtain standard errors of an estimate (e.g. median).
- ▶ Also provides approximate confidence intervals for a population parameter. Also called a Bootstrap Percentile confidence interval. It is the simplest method (among many approaches) for obtaining a confidence interval from the bootstrap.

# Bootstrap and prediction error

- ▶ In cross-validation, each of the  $K$  validation folds is distinct from the other  $K - 1$  folds used for training: **there is no overlap**. This is crucial for its success.
- ▶ To estimate prediction error using the bootstrap, we could think about using each bootstrap dataset as our training sample, and the original sample as our validation sample.
- ▶ But each bootstrap sample has significant overlap with the original data. About two-thirds of the original data points appear in each bootstrap sample.
- ▶ This will cause the bootstrap to seriously underestimate the true prediction error.
- ▶ The other way around – with original sample = training sample, bootstrap dataset = validation sample – is even worse.

## Removing the overlap

- ▶ We can partly fix this problem by only using predictions for those observations that did not (by chance) occur in the current bootstrap sample.
- ▶ But the method gets complicated, and in the end, cross-validation provides a simpler, more attractive approach for estimating prediction error.

## **Linear Model Selection and Regularization**

# Linear Model Selection and Regularization

- ▶ Recall the linear model

$$Y = \beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p + \epsilon.$$

- ▶ Tomorrow, we consider some approaches for extending the linear model framework.
- ▶ We generalize the linear model in order to accommodate **non-linear**, but still **additive**, relationships.
- ▶ We also consider even more general **non-linear** models.

# Linear models

- ▶ Despite its simplicity, the linear model has distinct advantages in terms of its **interpretability** and often shows good **predictive performance**.
- ▶ Hence we discuss some ways in which the simple linear model can be improved, by replacing ordinary least squares fitting with some alternative fitting procedures.

# Why alternatives to least squares?

- ▶ **Prediction Accuracy:** especially when  $p > n$ , to control the variance.
- ▶ **Model Interpretability:** By removing irrelevant features – that is, by setting the corresponding coefficient estimates to zero – we can obtain a model that is more easily interpreted.
- ▶ Some approaches for automatically performing **feature selection**.

# Three classes of methods

- ▶ **Subset Selection.** We identify a subset of the  $p$  predictors that we believe to be related to the response. We then fit a model using least squares on the reduced set of variables.
- ▶ **Shrinkage.** We fit a model involving all  $p$  predictors, but the estimated coefficients are shrunk towards zero relative to the least squares estimates. This shrinkage (also known as regularization) has the effect of reducing variance and can also perform variable selection.
- ▶ **Dimension Reduction.** We project the  $p$  predictors into a  $M$ -dimensional subspace, where  $M < p$ . This is achieved by computing  $M$  different **linear combinations**, or **projections**, of the variables. Then these  $M$  projections are used as predictors to fit a linear regression model by least squares.



# Shrinkage Methods

## Ridge regression and Lasso

- ▶ The subset selection methods use least squares to fit a linear model that contains a subset of the predictors.
- ▶ As an alternative, we can fit a model containing all  $p$  predictors using a technique that **constrains** or **regularizes** the coefficient estimates, or equivalently, that **shrinks** the coefficient estimates towards zero.
- ▶ It may not be immediately obvious why such a constraint should improve the fit, but it turns out that shrinking the coefficient estimates can significantly reduce their variance.

# Ridge regression

- Recall that the least squares fitting procedure estimates  $\beta_0, \beta_1, \dots, \beta_p$  using the values that minimize

$$RSS = \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2.$$

- In contrast, the ridge regression coefficient estimates  $\hat{\beta}^R$  are the values that minimize

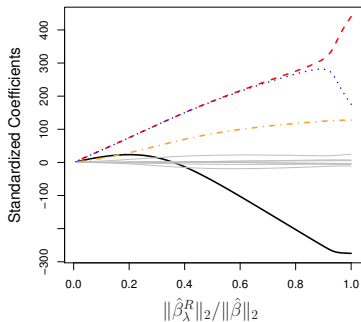
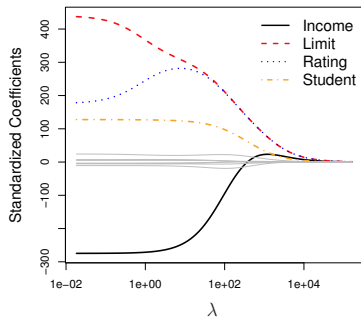
$$\sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 = RSS + \lambda \sum_{j=1}^p \beta_j^2,$$

where  $\lambda \geq 0$  is a **tuning parameter**, to be determined separately.

## Ridge regression: continued

- ▶ As with least squares, ridge regression seeks coefficient estimates that fit the data well, by making the RSS small.
- ▶ However, the second term,  $\lambda \sum_{j=1} \beta_j^2$ , called a **shrinkage penalty**, is small when  $\beta_1, \dots, \beta_p$  are close to zero, and so it has the effect of **shrinking** the estimates of  $\beta_j$  towards zero.
- ▶ The tuning parameter  $\lambda$  serves to control the relative impact of these two terms on the regression coefficient estimates.
- ▶ Selecting a good value for  $\lambda$  is critical; cross-validation is used for this.

## Example: Credit data



## Details of Previous Figure

- ▶ In the left-hand panel, each curve corresponds to the ridge regression coefficient estimate for one of the ten variables, plotted as a function of  $\lambda$ .
- ▶ The right-hand panel displays the same ridge coefficient estimates as the left-hand panel, but instead of displaying  $\lambda$  on the x-axis, we now display  $\|\hat{\beta}_\lambda^R\|_2 / \|\hat{\beta}\|_2$ , where  $\hat{\beta}$  denotes the vector of least squares coefficient estimates.
- ▶ The notation  $\|\beta\|_2$  denotes the  $\ell_2$  norm (pronounced “ell 2”) of a vector, and is defined as  $\|\hat{\beta}\|_2 = \sqrt{\sum_{j=1}^p \beta_j^2}$ .

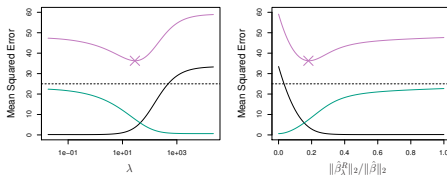
## Ridge regression: scaling of predictors

- ▶ The standard least squares coefficient estimates are **scale equivariant**: multiplying  $X_j$  by a constant  $c$  simply leads to a scaling of the least squares coefficient estimates by a factor of  $1/c$ . In other words, regardless of how the  $j$ th predictor is scaled,  $X_j \hat{\beta}_j$  will remain the same.
- ▶ In contrast, the ridge regression coefficient estimates can change **substantially** when multiplying a given predictor by a constant, due to the sum of squared coefficients term in the penalty part of the ridge regression objective function.
- ▶ Therefore, it is best to apply ridge regression after **standardizing the predictors**, using the formula

$$\tilde{x}_{ij} = \frac{x_{ij}}{\sqrt{\frac{1}{n} \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2}}$$

# Why Does Ridge Regression Improve Over Least Squares?

## The Bias-Variance tradeoff



- ▶ Simulated data with  $n = 50$  observations,  $p = 45$  predictors, all having nonzero coefficients.
- ▶ Squared bias (black), variance (green), and test mean squared error (purple) for the ridge regression predictions on a simulated data set, as a function of  $\lambda$  and  $\|\hat{\beta}_\lambda^R\|_2^2 / \|\hat{\beta}\|_2^2$ .
- ▶ The horizontal dashed lines indicate the minimum possible MSE.
- ▶ The purple crosses indicate the ridge regression models for which the MSE is smallest.

# The Lasso

- ▶ Ridge regression does have one obvious disadvantage: unlike subset selection, which will generally select models that involve just a subset of the variables, ridge regression will include all  $p$  predictors in the final model.
- ▶ The **Lasso** is a relatively recent alternative to ridge regression that overcomes this disadvantage. The lasso coefficients,  $\hat{\beta}_{\lambda}^L$ , minimize the quantity

$$\sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p |\beta_j| = RSS + \lambda \sum_{j=1}^p |\beta_j|$$

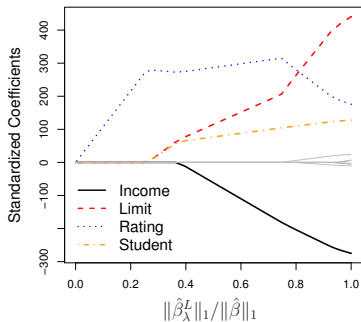
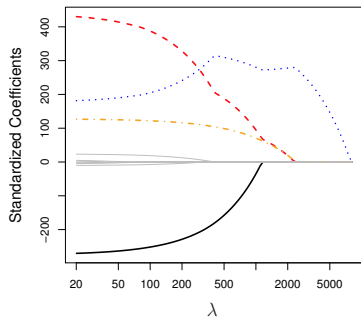
- ▶ In statistical parlance, the lasso uses an  $\ell_1$  (pronounced “ell 1”) penalty instead of an  $\ell_2$  penalty. The  $\ell_1$  norm of a coefficient vector  $\beta$  is given by  $\|\beta\|_1 = \sum |\beta_j|$ .



## The Lasso: continued

- ▶ As with ridge regression, the lasso shrinks the coefficient estimates towards zero.
- ▶ However, in the case of the lasso, the  $\ell_1$  penalty has the effect of forcing some of the coefficient estimates to be exactly equal to zero when the tuning parameter  $\lambda$  is sufficiently large.
- ▶ Hence, much like best subset selection, the lasso performs **variable selection**.
- ▶ We say that the lasso yields **sparse** models – that is, models that involve only a subset of the variables.
- ▶ As in ridge regression, selecting a good value of  $\lambda$  for the lasso is critical; cross-validation is again the method of choice.

## Example: Credit data



# Comparing the Lasso and Ridge Regression

- ▶ Neither ridge regression nor the lasso will universally dominate the other.
- ▶ In general, one might expect the lasso to perform better when the response is a function of only a relatively small number of predictors.
- ▶ However, the number of predictors that is related to the response is never known *a priori* for real data sets.
- ▶ A technique such as cross-validation can be used in order to determine which approach is better on a particular data set.

# Selecting the Tuning Parameter for Ridge Regression and Lasso

- ▶ As for subset selection, for ridge regression and lasso we require a method to determine which of the models under consideration is best.
- ▶ That is, we require a method selecting a value for the tuning parameter  $\lambda$  or equivalently, the value of the constraint  $s$ .
- ▶ **Cross-validation** provides a simple way to tackle this problem. We choose a grid of  $\lambda$  values, and compute the cross-validation error rate for each value of  $\lambda$ .
- ▶ We then select the tuning parameter value for which the cross-validation error is smallest.
- ▶ Finally, the model is re-fit using all of the available observations and the selected value of the tuning parameter.

## **Dimensionality Reduction Methods**

# Dimensionality Reduction Methods

- ▶ The methods that we have discussed so far today have involved fitting linear regression models, via least squares or a shrunk approach, using the original predictors,  $X_1, X_2, \dots, X_p$ .
- ▶ We now explore a class of approaches that **transform** the predictors and then fit a least squares model using the transformed variables.
- ▶ We will refer to these techniques as **dimension reduction methods**.

# Dimensionality Reduction Methods: details

- ▶ Let  $Z_1, Z_2, \dots, Z_M$  represent  $M < p$  linear combinations of our original  $p$  predictors. That is,

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

for some constants  $\phi_{m1}, \dots, \phi_{mp}$ .

- ▶ We can then fit the linear regression model,

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

using ordinary least squares.

- ▶ Note that in the linear regression model above, the regression coefficients are given  $\theta_0, \theta_1, \dots, \theta_M$ . If the constants  $\phi_{m1}, \dots, \phi_{mp}$  are chosen wisely, then such dimension reduction approaches can often outperform OLS regression.

- Notice that from the definition of  $Z_m$  above,

$$\sum_{m=1}^M \theta_m Z_{im} = \sum_{m=1}^M \theta_m \sum_{j=1}^p \phi_{mj} x_{ij} = \sum_{j=1}^p \sum_{m=1}^M \theta_m \phi_{mj} x_{ij} = \sum_{j=1}^p \beta_j x_{ij},$$

where

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

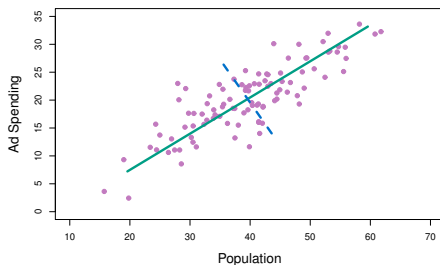
- Hence linear regression model above can be thought of as a special case of the original linear regression model.
- Dimension reduction serves to constrain the estimated  $\beta_j$  coefficients, since now they must take the above form.
- Doing this we can win in the bias-variance tradeoff.



# Principal Components Regression

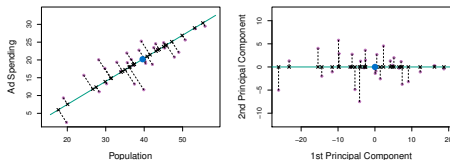
- ▶ Here we apply principal components analysis (PCA) (discussed later in the course) to define the linear combinations of the predictors, for use in our regression.
- ▶ The first principal component is that (normalized) linear combination of the variables with the largest variance.
- ▶ The second principal component has largest variance, subject to being uncorrelated with the first.
- ▶ And so on.
- ▶ Hence with many correlated original variables, we replace them with a small set of principal components that capture their joint variation.

# Illustration of PCA



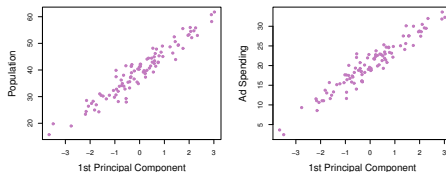
- ▶ The population size ( $pop$ ) and ad spending ( $ad$ ) for 100 different cities are shown as purple circles.
- ▶ The green solid line indicates the first principal component, and the blue dashed line indicates the second principal component.

# Illustration of PCA: continued



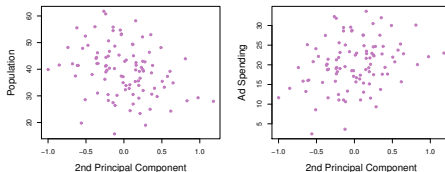
- ▶ A subset of the advertising data.
- ▶ **Left:** The first principal component, chosen to minimize the sum of the squared perpendicular distances to each point, is shown in green. These distances are represented using the black dashed line segments.
- ▶ **Right:** The left-hand panel has been rotated so that the first principal component lies on the x-axis.

# Illustration of PCA: continued



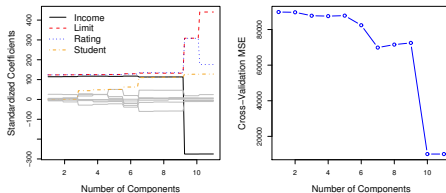
- ▶ Plots of the first principal component scores  $z_{i1}$  versus  $pop$  and  $ad$ .
- ▶ The relationships are strong.

## Illustration of PCA: continued



- ▶ Plots of the second principal component scores  $z_{i2}$  versus *pop* and *ad*.
- ▶ The relationships are weak.

# Choosing the number of components $M$



- ▶ **Left:** PCR standardized coefficient estimates on the *Credit* data set for different values of  $M$ .
- ▶ **Right:** The 10-fold cross validation MSE obtained using PCR, as a function of  $M$ .

# Summary

- ▶ Model selection methods are an essential tool for data analysis, especially for big datasets involving many predictors.
- ▶ Research into methods that give **sparsity**, such as the lasso is an especially hot area.