

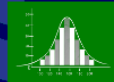
Introduction to Data Science

Winter Semester 2018/19

Oliver Ernst

TU Chemnitz, Fakultät für Mathematik, Professur Numerische Mathematik

Lecture Slides



Contents I

① What is Data Science?

② Learning Theory

2.1 What is Statistical Learning?

2.2 Assessing Model Accuracy

③ Linear Regression

3.1 Simple Linear Regression

3.2 Multiple Linear Regression

3.3 Other Considerations in the Regression Model

3.4 Revisiting the Marketing Data Questions

3.5 Linear Regression vs. K -Nearest Neighbors

④ Classification

4.1 Overview of Classification

4.2 Why Not Linear Regression?

4.3 Logistic Regression

4.4 Linear Discriminant Analysis

4.5 A Comparison of Classification Methods

⑤ Resampling Methods

Contents II

5.1 Cross Validation

5.2 The Bootstrap

6 Linear Model Selection and Regularization

6.1 Subset Selection

6.2 Shrinkage Methods

6.3 Dimension Reduction Methods

6.4 Considerations in High Dimensions

6.5 Miscellanea

7 Nonlinear Regression Models

7.1 Polynomial Regression

7.2 Step Functions

7.3 Regression Splines

7.4 Smoothing Splines

7.5 Generalized Additive Models

8 Tree-Based Methods

8.1 Decision Tree Fundamentals

8.2 Bagging, Random Forests and Boosting

Contents III

9 Support Vector Machines

- 9.1 Maximal Margin Classifier
- 9.2 Support Vector Classifiers
- 9.3 Support Vector Machines
- 9.4 SVMs with More than Two Classes
- 9.5 Relationship to Logistic Regression

10 Unsupervised Learning

- 10.1 Principal Components Analysis
- 10.2 Clustering Methods

⑥ Linear Model Selection and Regularization

6.1 Subset Selection

6.2 Shrinkage Methods

6.3 Dimension Reduction Methods

6.4 Considerations in High Dimensions

6.5 Miscellanea

Linear Model Selection and Regularization

Chapter overview

- Alternative fitting procedures to least squares (LS) for standard linear model

$$Y = \beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p + \varepsilon \quad (6.1)$$

to improve **prediction accuracy** and **model interpretability**.

- Prediction accuracy: for approximately linear (true) model, LS has low bias and, if $n \gg p$, also low variance. More variability if $n \gtrsim p$, no unique minimizer if $n < p$.

Idea: constraining or **shrinking** estimated coefficients reduces variability in these cases at negligible increase in bias, improving prediction accuracy.

- Model interpretability: some predictor variables may be irrelevant for response; LS will not remove these, hence consider other methods for **feature selection** or **variable selection** to exclude irrelevant variables from multiple regression model (by producing zero coefficients for these).

Linear Model Selection and Regularization

Alternative fitting procedures

We consider three classes of fitting alternatives to LS:

- **Subset selection**: Find subset of initial p predictor variables which are relevant, fit model using LS for reduced set of variables.
- **Shrinkage**: fit all p variables, shrink coefficients towards zero relative to LS estimate. Shrinkage (also known as **regularization**) reduces variance, some coefficients shrunk to zero, can be viewed as variable selection.
- **Dimension reduction**: project p predictors into subspace of dimension $M < p$, i.e., construct M linearly independent *pseudo-variables* which depend linearly on original p predictor variables. Use these as new predictors for LS fit.
- Same concepts apply to other methods (e.g. classification).

⑥ Linear Model Selection and Regularization

6.1 Subset Selection

6.2 Shrinkage Methods

6.3 Dimension Reduction Methods

6.4 Considerations in High Dimensions

6.5 Miscellanea

Linear Model Selection and Regularization

Best subset selection

Idea: Perform separate LS regression for *all possible subsets* of given p predictor variables.

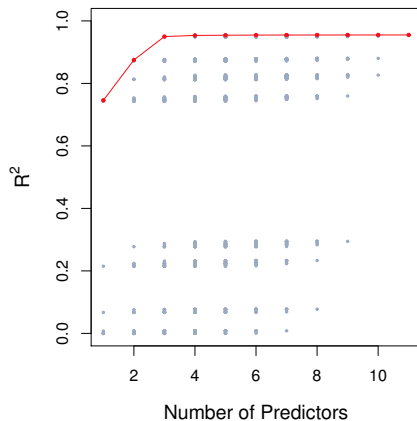
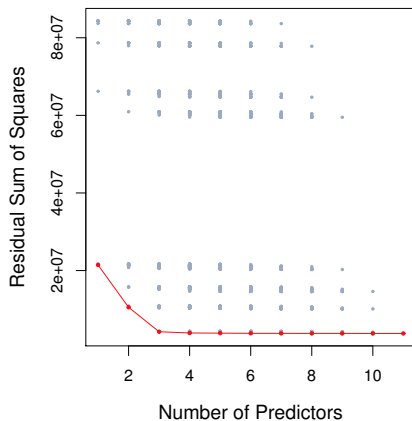
Algorithm 1: Best subset selection.

- 1 Set \mathcal{M}_0 to be the **null model**, i.e., containing only constant term β_0 .
 - 2 **for** $k = 1, 2, \dots, p$
 - a Fit all $\binom{p}{k}$ models containing exactly k predictors.
 - b Pick best (smallest RSS, i.e., largest R^2) among these, call it \mathcal{M}_k .
 - 3 Select single best model among $\mathcal{M}_0, \dots, \mathcal{M}_p$ using model selection criterion (later).
-

- Step 2 reduces # model candidates from 2^p to $p + 1$.
- Models in Step 3 display monotone decreasing RSS (increasing R^2) as # variables increases.
- Want low test error rather than low training error.

Linear Model Selection and Regularization

Best subset selection



Best subset selection for **Credit** data set: 10 predictors (three-valued variable **ethnicity** coded using two dummy variables selected separately).

Red line indicates model with smallest RSS (largest R^2).

Linear Model Selection and Regularization

Best subset selection

- Can apply to classification problems using **deviance** in place of RSS ($-2 \cdot$ maximized log-likelihood).
- Best subset selection simple, but $\#$ regression fits to compare grows exponentially with p (e.g. 1024 for $p = 10$, over 1 million for $p = 20$).
- Also, statistical problems for large p : the larger the search space, the higher the chance of finding models performing well on training set, but badly for test set.

Linear Model Selection and Regularization

Forward stepwise selection

Idea: Add predictors to model one at a time, at each step adding variable leading to greatest additional improvement.

Algorithm 2: Forward stepwise selection.

- 1 Set \mathcal{M}_0 to be the **null model**, i.e., containing only constant term β_0 .
 - 2 **for** $k = 0, 1, \dots, p - 1$
 - a Consider all $p - k$ models augmenting \mathcal{M}_k by one additional predictor.
 - b Pick best (smallest RSS, i.e., largest R^2) among these, call it \mathcal{M}_{k+1} .
 - 3 Select single best model among $\mathcal{M}_0, \dots, \mathcal{M}_p$ using model selection criterion (later).
-

- Rather than 2^p models considered by best subset selection, forward stepwise selection requires only $1 + p(p + 1)/2$ LS fits.

E.g. $p = 20$: 1,048,576 models for best subset selection, 211 models for forward stepwise selection.

Linear Model Selection and Regularization

Forward stepwise selection

- Forward stepwise selection not guaranteed to find best model out of 2^p possible. E.g. for $p = 3$, best single-variable model could consist of X_1 , while best two-variable model consists of X_2, X_3 .
- First 4 selected models for best subset selection and forward stepwise selection on [Credit](#) data set:

# variables	Best subset	Forward stepwise
1	rating	rating
2	rating, income	rating, income
3	rating, income, student	rating, income, student
4	cards, income student, limit	rating, income student, limit

- Can use forward stepwise selection in high-dimensional case when $n < p$. However, can only construct submodels $\mathcal{M}_0, \dots, \mathcal{M}_{n-1}$, since LS can uniquely fit at most $n - 1$ variables.

Linear Model Selection and Regularization

Backward stepwise selection

Idea: Begin with full LS model, successively remove least useful predictor.

Algorithm 3: Backward stepwise selection.

- 1 Set \mathcal{M}_p to be the **full model**, containing all p predictors.
 - 2 **for** $k = p, p - 1, \dots, 1$
 - a Consider all k models containing all but one of the predictors in \mathcal{M}_k .
 - b Pick best (smallest RSS, i.e., largest R^2) among these k models, call it \mathcal{M}_{k-1} .
 - 3 Select single best model among $\mathcal{M}_0, \dots, \mathcal{M}_p$ using model selection criterion (later).
-

- Again only $1 + p(p + 1)/2$ model fits.
- No guarantee of finding best model.
- Requires $n > p$.
- **Hybrid** approaches possible, where addition step followed by removal step.

Linear Model Selection and Regularization

Optimal model selection

- In best subset selection, forward selection and backward selection, need to choose best among models containing different # variables.
- RSS and R^2 measures will always select model with all p variables.
- Goal: select best model with respect to *test* error.
- Two basic approaches:
 - ① Indirectly estimate test error by making an adjustment to training error to account for bias due to overfitting.
 - ② Directly estimate test error using either validation set approach or cross-validation approach.

Linear Model Selection and Regularization

C_p , AIC, BIC, adjusted R^2

- Training set MSE generally underestimates test MSE (recall $\text{MSE} = \text{RSS} / n$)
- For LS regression: coefficients determined by minimization of RSS.
- Therefore training error decreases as variables added to model; not so for test error.
- For fitted LS model containing d predictors, C_p **estimate** defined by

$$C_p := \frac{1}{n}(\text{RSS} + 2d \hat{\sigma}^2), \quad (6.2)$$

where $\hat{\sigma}^2$ is an estimate of $\text{Var } \epsilon$, typically computed using full model. Adds penalty term $2d \hat{\sigma}^2$ to training RSS to compensate for underestimating test error.

Can show: C_p unbiased estimate of test MSE if $\hat{\sigma}^2$ unbiased estimate of σ^2 .

Hence C_p small for models with small test MSE.

Linear Model Selection and Regularization

AIC

- **Akaike information criterion** (AIC) defined for models fit by maximum likelihood.
- For standard linear model (6.1) with Gaussian noise maximum likelihood fit coincides with LS fit.

- In this case

$$AIC = \frac{1}{n\hat{\sigma}^2} (RSS + 2d\hat{\sigma}^2)$$

(have omitted additive constant).

- Hence, for LS models C_p and AIC proportional.

Linear Model Selection and Regularization

BIC

- **Bayesian information criterion** (BIC), derived from Bayes point of view, is given by (up to irrelevant constants)

$$\text{BIC} = \frac{1}{n\hat{\sigma}^2}(\text{RSS} + d\hat{\sigma}^2 \log n) \quad (6.3)$$

- Also tends to be small for models with small test error.
- Replaces $2d\hat{\sigma}^2$ used by C_p with $d\hat{\sigma}^2 \log n$, hence places heavier penalty on models with many variables, results in selection of smaller models than C_p .

Linear Model Selection and Regularization

Adjusted R^2

- Recall $R^2 = 1 - \text{TSS} / \text{RSS}$,
 $\text{TSS} = \sum (y_i - \bar{y})^2$ total sum of squares for response.
- R^2 increases as variables added to LS model.
- For LS model with d variables, **adjusted R^2** statistic given by

$$\text{Adjusted } R^2 := 1 - \frac{\text{RSS} / (n - d - 1)}{\text{TSS} / (n - 1)} = 1 - \frac{\text{RSS}}{\text{TSS}} \cdot \frac{n - 1}{n - d - 1}. \quad (6.4)$$

- Unlike C_p , AIC and BIC, where small value indicates model with low test error, here a *large* value of the adjusted R^2 statistic indicates a model with a small test error.
- Maximizing adjusted R^2 equivalent to minimizing $\text{RSS} / (n - d - 1)$.
- Intuition: once all relevant variables have been included, adding additional noise variables will only lead to small decrease in RSS.
- Compared to R^2 , adjusted R^2 pays a price for adding irrelevant variables.

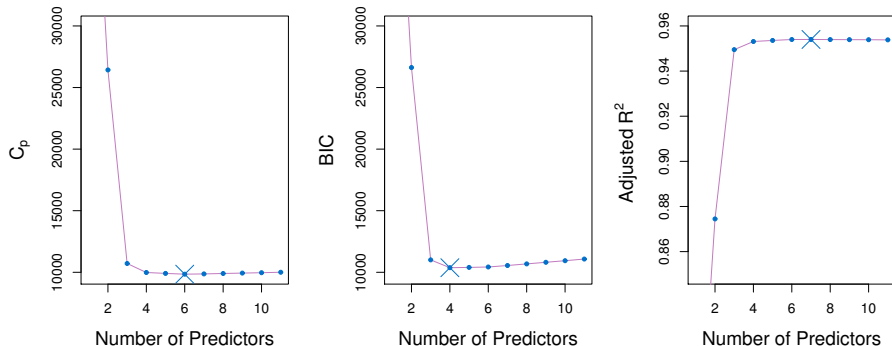
Linear Model Selection and Regularization

C_p , AIC, BIC, adjusted R^2

- Rigorous justifications of C_p , AIC, BIC rely on asymptotic arguments (large n limit).
- Adjusted R^2 popular, intuitive, but not as well motivated statistically.
- All measures simple to use and compute.
- Modified formulas for more general models.

Linear Model Selection and Regularization

C_p , AIC, BIC, adjusted R^2



C_p , BIC and adjusted R^2 for best models of each size for Credit data set (red curve in previous plot).

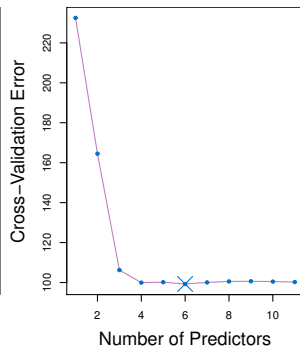
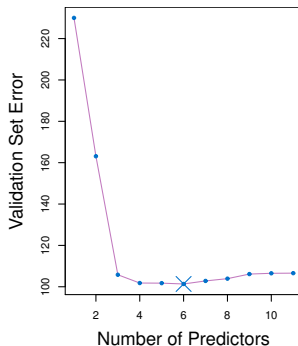
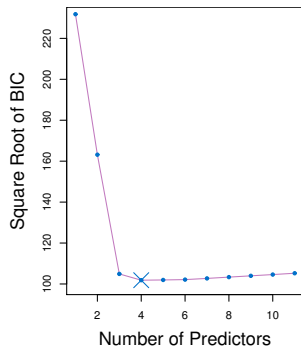
Linear Model Selection and Regularization

Cross-validation

- Can apply validation and cross-validation to each model and select that with lowest estimate.
- Advantage over C_p , AIC, BIC, adjusted R^2 : direct estimate of test error, fewer assumptions about underlying model.
- More widely useable, e.g., when noise variance estimates difficult to obtain.
- CV initially less popular than C_p , AIC, BIC, adjusted R^2 due to computational cost; this is less and less an issue.
- Apply to **Credit** data set: display BIC, validation set errors, cross-validation errors as function of $d = \#$ variables in model.
Validation: randomly choose 3/4 of observations as training set, remainder as validation set.
Cross-validation using $k = 10$ folds.

Linear Model Selection and Regularization

Cross-validation



Credit data: 3 model error estimates for best model containing 1 to 11 predictors. Both validation set and CV result in 6-variable models. All approaches agree: not much difference in test error for 4, 5, 6-variable models.

Left: $\sqrt{\text{BIC}}$; center: validation set errors; right: CV errors.

Linear Model Selection and Regularization

Cross-validation

- Observation: all 3 error estimates quite flat from 4 variables onward.
- Error estimate-minimizing model likely to change for different partitions of observations or different choice of CV folds.
- **One-standard-error rule:** calculate standard error of estimated test MSE for each model size, then select smallest model for which estimated test error is within one standard error of lowest point on curve.
Rationale: if several models appear equally good, may as well choose simplest.
Here: rule leads to 3-variable model.

⑥ Linear Model Selection and Regularization

6.1 Subset Selection

6.2 Shrinkage Methods

6.3 Dimension Reduction Methods

6.4 Considerations in High Dimensions

6.5 Miscellanea

Linear Model Selection and Regularization

Shrinkage

- **Inverse problems**: branch of applied mathematics for solving problems where solution extremely sensitive to data and/or solution not unique (e.g.: X-ray tomography, image deblurring).
- Prevalent strategy: instead of original problem, solve *nearby* problem with better stability properties: **regularization**.
- In LS methods: modify objective function by minimizing different norm or adding **penalty term**, thus imposing “a priori information” on the coefficients.
- In statistics, particularly in LS regression, regularization is known as **shrinkage**, as certain coefficients are “shrunk” in magnitude relative to their values under LS estimation.
- Here we introduce two popular shrinkage techniques: **ridge regression** and the **LASSO**.

Linear Model Selection and Regularization

Ridge regression

Least-squares fitting determines coefficients β_0, \dots, β_p by minimizing

$$\text{RSS} = \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j X_j \right)^2 = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2.$$

In **ridge regression**, one minimizes instead the objective function

$$\sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j X_j \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 = \text{RSS} + \lambda \|\tilde{\boldsymbol{\beta}}\|_2^2, \quad (6.5)$$

where λ is a tuning parameter to be suitably chosen and $\tilde{\boldsymbol{\beta}} := (\beta_1, \dots, \beta_p)^\top \in \mathbb{R}^p$. From now on $\boldsymbol{\beta} \in \mathbb{R}^p$ and tilde omitted.

In the inverse problems community, this general approach is known as **Tikhonov regularization** and λ is called the **regularization parameter**.

Linear Model Selection and Regularization

Ridge regression

- Tuning λ constitutes tradeoff between two objectives: minimizing RSS (good fit to data) and minimizing **shrinkage penalty** $\lambda \|\boldsymbol{\beta}\|_2^2$, which shrinks β_1, \dots, β_p to zero.
- $\lambda = 0$: recover standard LS estimate.
- $\lambda \rightarrow \infty$: $\boldsymbol{\beta} \rightarrow \mathbf{0}$.
- Different estimate for each value of λ , choice critical.
- Intercept omitted from shrinkage: this is just the mean value of response when all predictor variables are zero.

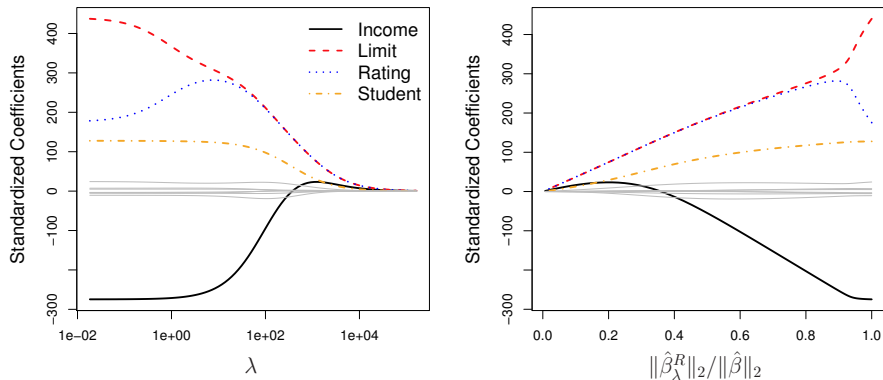
Under assumption that all columns of data matrix \mathbf{X} have been **centered** to have mean zero, then

$$\hat{\beta}_0 = \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i.$$

- In the following, for the standard linear model, we tacitly assume \mathbf{X} to be centered, the coefficient β_0 to be set to its optimal value \bar{y} and the coefficient vector to be estimated to consist of the components β_1, \dots, β_p .

Linear Model Selection and Regularization

Ridge regression



Ridge regression applied to **Credit** data set: values of coefficients of the 10 predictor variables against λ . Lines for largest coefficients **income**, **limit**, **rating** and **student** displayed in distinct colors. Right: x-axis is $\|\hat{\beta}_\lambda^R\|_2 / \|\hat{\beta}\|_2$ in place of λ .

Predictor variables standardized before carrying out ridge regression.

Linear Model Selection and Regularization

Ridge regression: standardizing the predictors

- For LS estimation of standard linear model, rescaling a predictor variable $X_j \leftarrow cX_j$ simply results in reciprocal rescaling of estimate as $\hat{\beta}_j \leftarrow \hat{\beta}_j/c$. Consequence: $\hat{\beta}_j X_j$, hence data fit, remains the same. This property is called **scale equivariance**.
- This is no longer the case for ridge regression: value of $\hat{\beta}_{j,\lambda}^R X_j$ depends on λ as well as the scaling of X_j (possibly even the scaling of other predictors).
- Therefore, best to **standardize** predictor variables by transformation

$$x_{i,j} \leftarrow \tilde{x}_{i,j} := \frac{x_{i,j}}{s_j}, \quad s_j := \sqrt{\frac{1}{n} \sum_{i=1}^n (x_{i,j} - \bar{x}_j)^2}. \quad (6.6)$$

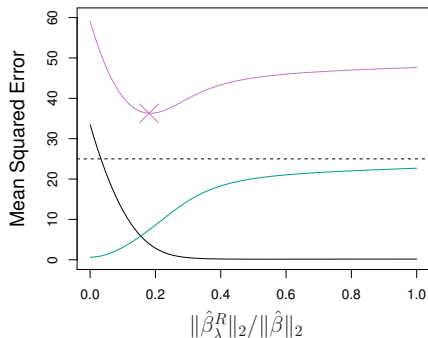
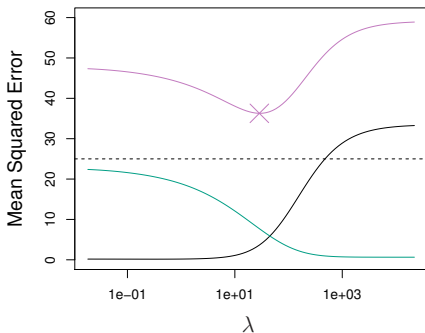
Denominator s_j estimates variance of j -th predictor.

- Standardized predictor observations have unit variance estimate.

Linear Model Selection and Regularization

Ridge regression: improvement over LS

Bias-variance tradeoff: as λ increases, model flexibility decreased, reducing variance, increasing bias.



Simulated data, $p = 45$ predictors, $n = 50$ observations. Test MSE (purple), squared bias (black) and variance (green) of ridge regression predictions. Cross: minimal MSE. Dashed line: minimal possible MSE.

Linear Model Selection and Regularization

Ridge regression: improvement over LS

- In general: for almost linear (true) model, LS estimate has low bias, but possibly high variance, particularly when p large relative to n .
- For $p > n$ LS fit not unique, but ridge regression still works, trading off slight bias for much reduced variance.
- Computational advantage over best subset selection: ridge regression for many values of λ can be computed at cost of essentially one LS fit, compared to comparing 2^p models.

Linear Model Selection and Regularization

The LASSO

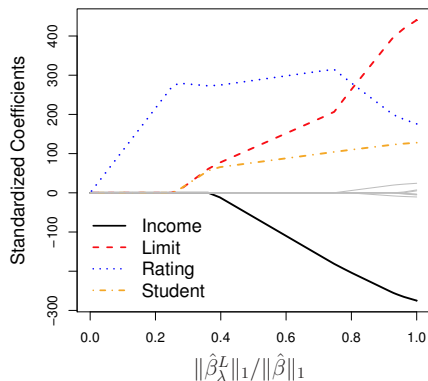
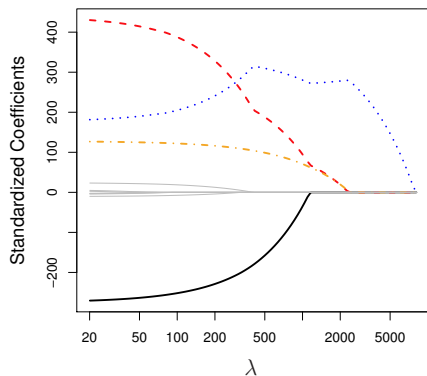
- Disadvantage of ridge regression: will generally include all p predictors in the model, in contrast with subset selection methods.
- OK for prediction, challenging for interpretation.
- Example: **Credit** data set; most important variables are **income**, **limit**, **rating** and **student**. Model including just these desirable, ridge regression will generally include all 10 predictors.
- LASSO** (least absolute shrinkage and selection operator): choose coefficients β_j to minimize

$$\sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{i,j} \right)^2 + \lambda \sum_{j=1}^p |\beta_j| = \text{RSS} + \lambda \|\boldsymbol{\beta}\|_1. \quad (6.7)$$

- ℓ^2 -penalty in ridge regression replaced by ℓ^1 -penalty, $\|\boldsymbol{\beta}\|_1 = |\beta_1| + \dots + |\beta_p|$.
- ℓ^1 -penalty: for λ sufficiently large, results in some estimates $\hat{\beta}_{j,\lambda}^L$ being exactly zero, effecting an implicit **variable selection**, yielding in **sparse models**, which are easier to interpret.

Linear Model Selection and Regularization

The LASSO



LASSO applied to [Credit](#) data set. Note difference to ridge regression for intermediate values of λ : as λ increases, coefficients are successively set to zero, thereby removed from model.

Linear Model Selection and Regularization

Equivalent constrained minimization problem

Can show: ridge regression and LASSO estimates solve **constrained minimization problems**

$$\hat{\beta}_{\lambda}^L = \arg \min_{\beta} \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{i,j} \right)^2 \quad \text{subject to } \|\beta\|_1 \leq s \quad (6.8)$$

and

$$\hat{\beta}_{\lambda}^R = \arg \min_{\beta} \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{i,j} \right)^2 \quad \text{subject to } \|\beta\|_2^2 \leq s, \quad (6.9)$$

respectively.

In other words: for each value of λ , there is a corresponding value of s , such that both problems give the same estimates.

Linear Model Selection and Regularization

LASSO: relation to best subset selection

- Consider constrained minimization problem

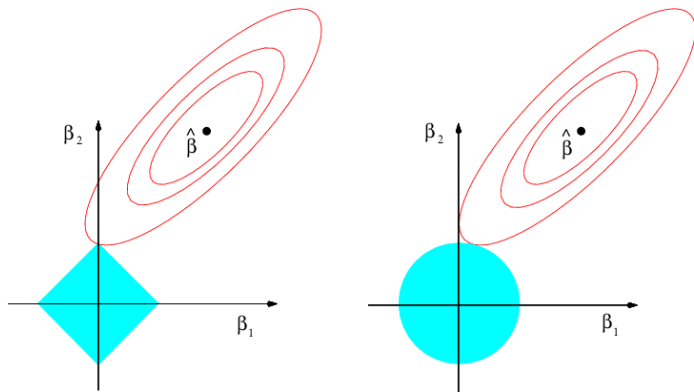
$$\hat{\beta} = \arg \min_{\beta} \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{i,j} \right)^2 \quad \text{subject to} \quad \sum_{j=1}^p \mathbf{1}_{\{\beta_j \neq 0\}} \leq s \quad (6.10)$$

- Minimizes RSS subject to constraint that no more than s coefficients are nonzero.
- This is equivalent to best subset selection.
- Computationally infeasible for large p , since it involves considering all $\binom{p}{s}$ models containing s predictors.
- Hence ridge regression / LASSO computationally feasible alternatives to best subset selection replacing intractable form of budget in (6.10).

Linear Model Selection and Regularization

LASSO: variable selection property

- Formulations (6.8) and (6.9) key to understanding variable selection property of LASSO:

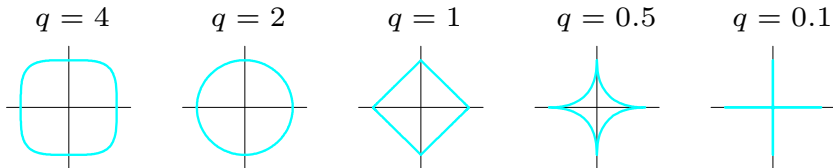


Red: RSS contours, blue: constraints $|\beta_1| + |\beta_2| \leq s$ (left) and $\beta_1^2 + \beta_2^2 = s$ (right).

Linear Model Selection and Regularization

LASSO: variable selection property

- Unit spheres of $\sum_{j=1}^p |\beta_j|^q$ for $q < 2$ progressively sharper (no longer a norm for $q < 1$).

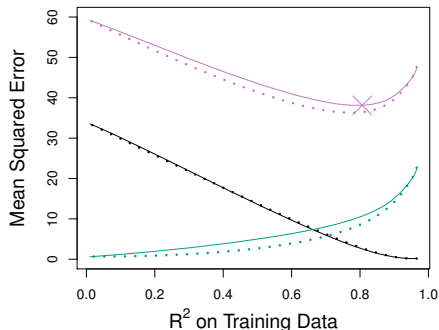
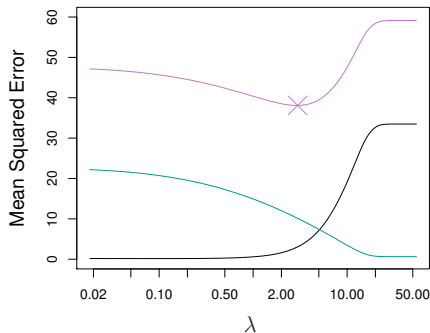


- Limiting case: $q = 0$ counts $\#$ nonzero components.

Linear Model Selection and Regularization

Comparison of ridge regression with LASSO

Simulated data using all $p = 45$ predictors: ($\beta_j \neq 0 \forall j$ in true model)



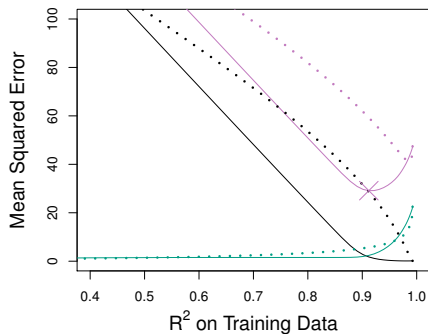
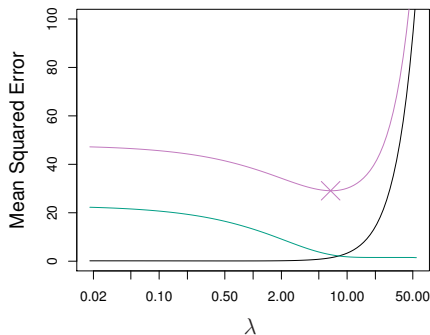
Left: Test MSE (purple), squared bias (black) and variance (green) of LASSO for different values of λ .

Right: Comparison of test MSE (purple), squared bias (black) and variance (green) against training R^2 ; dotted lines denote corresponding quantities for ridge regression.

Linear Model Selection and Regularization

Comparison of ridge regression with LASSO

Simulated data using only 2 out of $p = 45$ predictors: (only two $\beta_j \neq 0$ in true model)



Left: Test MSE (purple), squared bias (black) and variance (green) of LASSO for different values of λ .

Right: Comparison of test MSE (purple), squared bias (black) and variance (green) against training R^2 ; dotted lines denote corresponding quantities for ridge regression.

Linear Model Selection and Regularization

Simple special case for ridge regression and the lasso

Assume data matrix $\mathbf{X} = \mathbf{I}$ ($p = n$) and $\bar{y} = 0$.

LS problem reduces to minimizing

$$\sum_{j=1}^p (y_j - \beta_j)^2, \quad \text{hence } \beta_j = y_j, \quad j = 1, \dots, p. \quad (6.11)$$

Ridge regression and lasso estimation result from minimizing

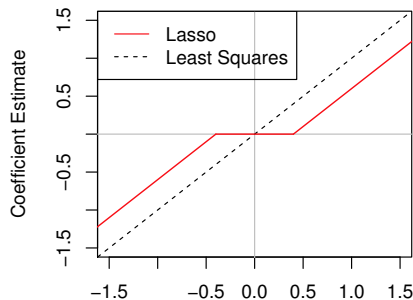
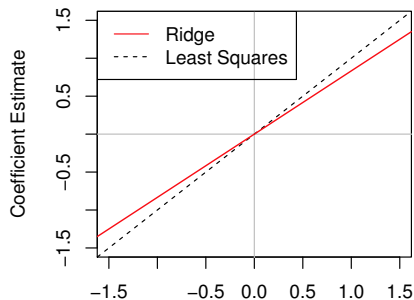
$$\sum_{j=1}^p (y_j - \beta_j)^2 + \lambda \sum_{j=1}^p \beta_j^2 \quad \text{and} \quad \sum_{j=1}^p (y_j - \beta_j)^2 + \lambda \sum_{j=1}^p |\beta_j|,$$

respectively, with solutions

$$\hat{\beta}^R = \frac{y_j}{1 + \lambda}, \quad \hat{\beta}^L = \begin{cases} y_j - \lambda/2 & \text{if } y_j > \lambda/2, \\ y_j + \lambda/2 & \text{if } y_j < -\lambda/2, \\ 0 & \text{if } |y_j| \leq \lambda/2. \end{cases}$$

Linear Model Selection and Regularization

Simple special case for ridge regression and the lasso



Ridge regression (left) and lasso (right) estimates for one variable of special case $\mathbf{X} = \mathbf{I}$ and $p = n$.

General case: more complicated (of course), but basic mechanism still holds:

- Ridge regression: shrinks every dimension roughly by same proportion.
- Lasso: shrinks all components to zero by similar amount, sufficiently small coefficients damped to zero.

Linear Model Selection and Regularization

Bayesian interpretation for ridge regression and the lasso

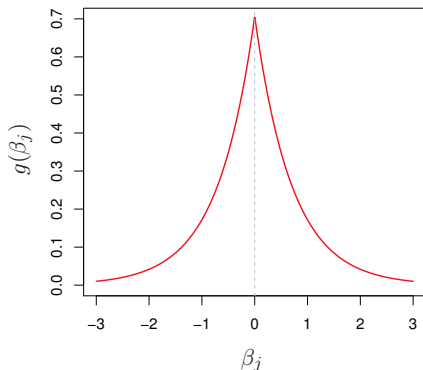
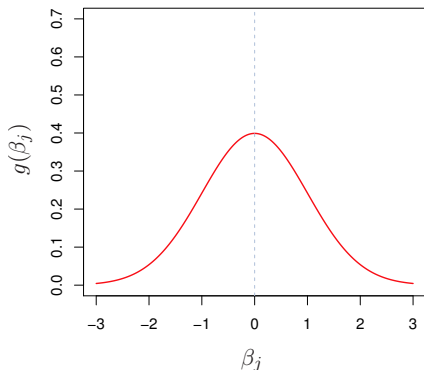
- Assume prior distribution on $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^\top$, with density $p(\boldsymbol{\beta})$.
- Likelihood of data: $f(Y|X, \boldsymbol{\beta})$, $X = (X_1, \dots, X_p)$.
- Bayes' rule then says (noting X is fixed)

$$p(\boldsymbol{\beta}|X, Y) \propto f(Y|X, \boldsymbol{\beta}) \cdot p(\boldsymbol{\beta}|X) = f(Y|X, \boldsymbol{\beta}) \cdot p(\boldsymbol{\beta}).$$

- Assume standard linear model $Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \varepsilon$, with independent Gaussian noise and $p(\boldsymbol{\beta}) = \prod_{j=1}^p g(\beta_j)$ for pdf g .
- Ridge regression/lasso results from two special cases for g :
 - g centered Gaussian, λ -dependent variance, then ridge regression estimate is **posterior mode** (and posterior mean) of $\boldsymbol{\beta}$.
 - g centered Laplace distribution with λ -dependent scale parameter, then posterior mode for $\boldsymbol{\beta}$ given by lasso estimate. (Not posterior mean; posterior mean itself not sparse.)

Linear Model Selection and Regularization

Bayesian interpretation for ridge regression and the lasso



Prior densities for Bayesian interpretation of shrinkage methods.

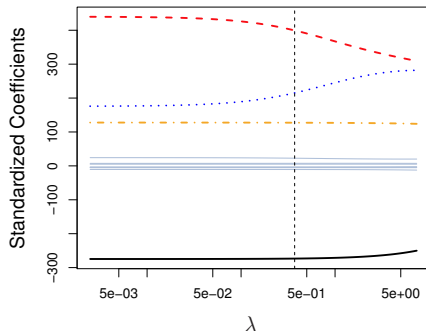
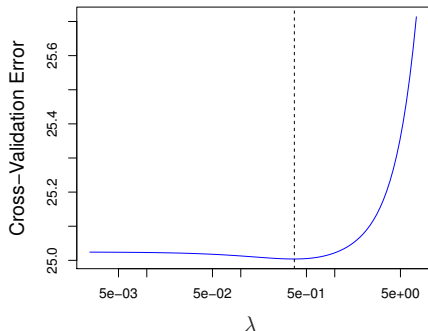
Left: centered Gaussian prior density, results in posterior distribution with ridge regression solution as posterior mode.

Right: centered Laplace (double-exponential) prior density, results on lasso solution as posterior mode.

Linear Model Selection and Regularization

Selection of λ

- Model selection methods required measure of goodness to compare models.
- Shrinkage methods require selection of shrinkage parameter λ .
- **Cross-validation** approach: fix a grid of λ values; compute cross-validation error for each λ ; select λ with smallest error; refit this model with all available observations.

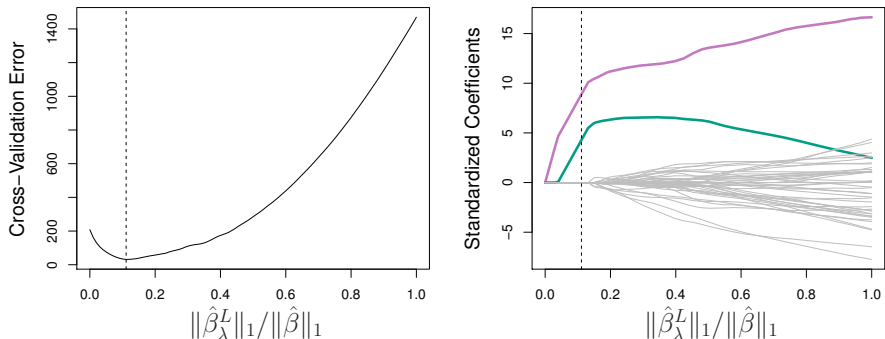


Left: LOOCV errors vs. λ for ridge regression applied to **Credit** data set.

Right: Coefficient estimates vs. λ . Vertical dashed line indicates selected λ .

Linear Model Selection and Regularization

Selection of λ



10-fold CV applied to data set from Slide 281.

Left: CV error. Right: coefficient estimates. Vertical dashed line indicates CV error-minimizing λ . Colored lines represent 2 predictors related to response, grey lines unrelated predictors (**signal** vs. **noise**).

Lasso assigns relevant predictors much larger estimates; CV chooses λ for which irrelevant predictors set to zero. Compare LS estimate (far right).

⑥ Linear Model Selection and Regularization

6.1 Subset Selection

6.2 Shrinkage Methods

6.3 Dimension Reduction Methods

6.4 Considerations in High Dimensions

6.5 Miscellanea

Linear Model Selection and Regularization

Dimension reduction methods

- Up to now: control variance by removing predictor variables or shrinking coefficients.
- Now: reduce variance by projecting into subspace of dimension $M < p$.
- Set

$$Z_m := \sum_{j=1}^M \phi_{j,m} X_j, \quad m = 1, \dots, M, \quad \text{i.e., } \mathbf{Z} = \mathbf{X}\Phi, \quad \Phi \in \mathbb{R}^{p \times M}. \quad (6.12)$$

- Fit standard linear regression model

$$Y = \theta_0 + \theta_1 Z_1 + \dots + \theta_M Z_M + \varepsilon. \quad (6.13)$$

- **Dimension reduction**: fit $M + 1 < p + 1$ coefficients.
- For well-chosen Φ , this reduced-dimension approach can outperform LS.

Linear Model Selection and Regularization

Dimension reduction methods

- Note:

$$\sum_{m=1}^M z_{i,m} \theta_m = \sum_{m=1}^M \sum_{j=1}^p \phi_{j,m} x_{i,j} \theta_m = \sum_{j=1}^p x_{i,j} \sum_{m=1}^M \phi_{j,m} \theta_m = \sum_{j=1}^p x_{i,j} \beta_j$$

with $\beta_j := \sum_{m=1}^M \phi_{j,m} \theta_m$.

- In matrix terms:

$$\mathbf{Z}\boldsymbol{\theta} = \mathbf{X}\boldsymbol{\Phi}\boldsymbol{\theta} =: \mathbf{X}\boldsymbol{\beta}, \quad \boldsymbol{\beta} := \boldsymbol{\Phi}\boldsymbol{\theta}$$

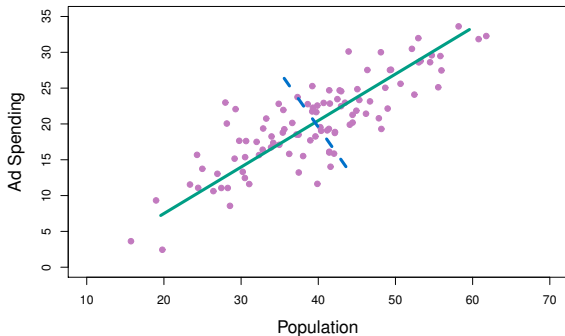
Hence can view (6.13) as special case of original linear model (6.1).

- Dimension reduction constrains $\boldsymbol{\beta}$ by making it a linear function of $M < p$ variables $\{\theta_m\}_{m=1}^M$.
- May introduce bias, but when $p \gg n$ this is outweighed by resulting variance reduction.
- Next: 2 ways of choosing $\boldsymbol{\Phi}$.

Linear Model Selection and Regularization

Principal components regression

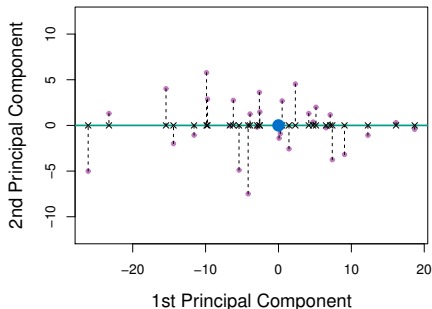
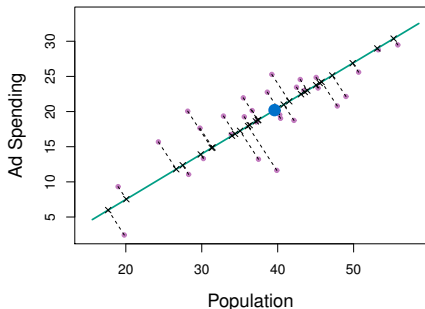
- **Principal components analysis** (PCA): approach for deriving a low-dimensional feature set from a large set of variables.
- First principal component: direction in \mathbb{R}^p in which observations vary the most.



Population size **pop** vs. ad spending **ad** for 100 cities (purple dots). Green solid line: first principal component; blue dashed line: second principal component.

Linear Model Selection and Regularization

Principal components regression



- Project data on direction (line) along which it varies most.
- For `pop` / `ad` data: $\phi_{1,1} = 0.839$, $\phi_{2,1} = 0.544$, giving

$$Z_1 = 0.839 \times (\text{pop} - \overline{\text{pop}}) + 0.544 \times (\text{ad} - \overline{\text{ad}})$$

- Out of every (normalized) linear combination of the (centered) variables, Z_1 has maximal variance.

Linear Model Selection and Regularization

Principal components regression

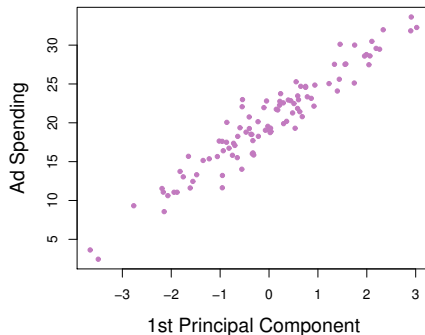
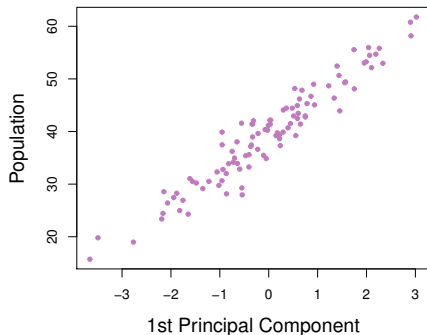
- Principal component data vector (“**scores**”) has same length n , e.g.

$$z_{i,1} = 0.839 \times (\text{pop}_i - \overline{\text{pop}}) + 0.544 \times (\text{ad}_i - \overline{\text{ad}}), \quad i = 1, \dots, n.$$

- Alternative interpretation of PCA: 1st principal component vector defines line as close as possible to data in sense of minimizing sum of squared perpendicular distances between each data point and this line.
-

Linear Model Selection and Regularization

Principal components regression



First principal components scores $z_{i,1}$ for [pop](#) and [ad](#). Strong relationship in both cases, i.e., principal component captures most of the information contained in the two predictors.

Linear Model Selection and Regularization

Principal components regression

- Second principal component Z_2 : direction of largest variance among all linear combinations of predictor variables which is *orthogonal* to (uncorrelated with) Z_1 .
- Here:

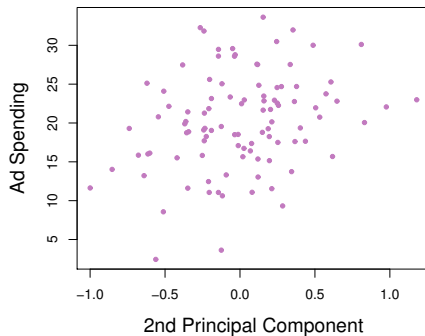
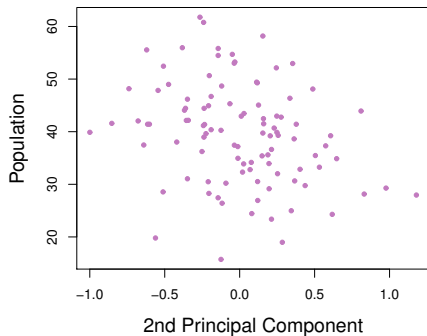
$$Z_2 = 0.544 \times (\text{pop} - \overline{\text{pop}}) - 0.839 \times (\text{ad} - \overline{\text{ad}}).$$

Since $p = 2$, this covers all of remaining variance.

- Of these, Z_1 contains most of the information, cf. much larger variation in Z_1 -coordinate than Z_2 -coordinate in right panel of figure on Slide 292.
- Plot on Slide 296 displays $z_{i,2}$ against **pop** and **ad** predictors: much less relationship than with Z_1 .
Thus, Z_1 sufficient to explain most of variability in data set.
- For p predictor variables, can construct up to p principal components.

Linear Model Selection and Regularization

Principal components regression



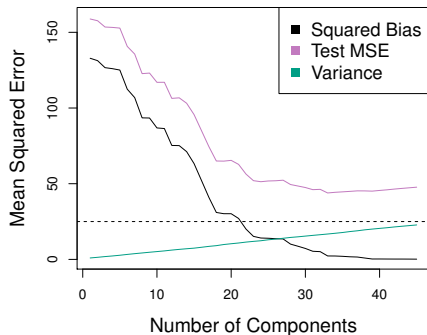
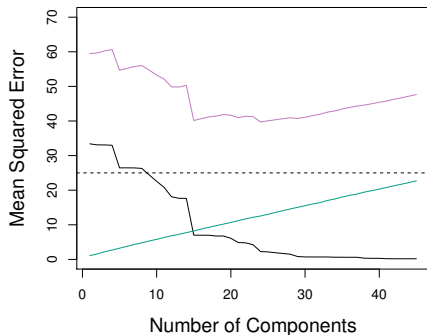
Linear Model Selection and Regularization

Principal components regression

- **Principal components regression** (PCR): construct first M principal components Z_1, \dots, Z_M , use these in a linear regression model fit by LS.
- Guiding principle: directions in span of X_1, \dots, X_p with most variance are the directions associated with response Y .
- Under this assumption, fitting LS model to Z_1, \dots, Z_M will yield better predictions than fitting X_1, \dots, X_p , since most information related to response Y contained in Z_1, \dots, Z_M , and estimating $M \ll p$ coefficients avoids overfitting.

Linear Model Selection and Regularization

Principal components regression

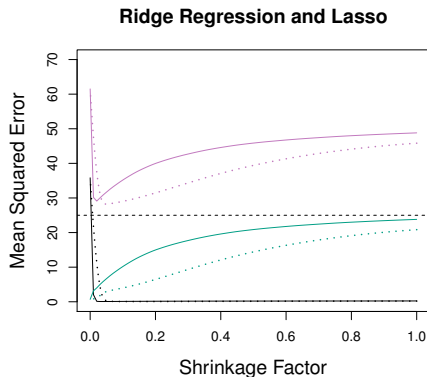
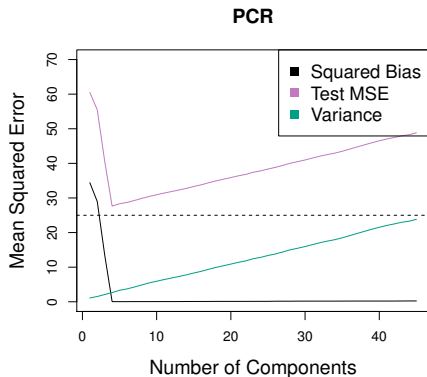


PCR fits to data sets from Slide 280 (left) and Slide 281 (right): MSE against # principal components M . More components reduces bias, increases variance (U-shape). $p = M$ coincides with LS fit of original predictors. Compared with ridge regression and lasso results in figures on Slides 272, 280 and 281, PCR seen to underperform shrinkage.

Linear Model Selection and Regularization

Principal components regression

Worse performance of PCR in previous example due to fact that many principal components needed needed to explain response.



Data generated in such a way that response depends exclusively on first 5 principal components. Left: PCR, MSE has clear minimum at $M = 5$. Right: ridge regression (dotted) and lasso (solid) results.

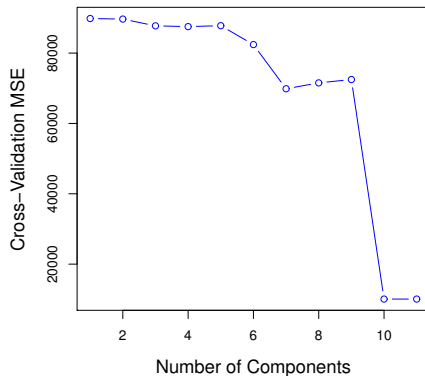
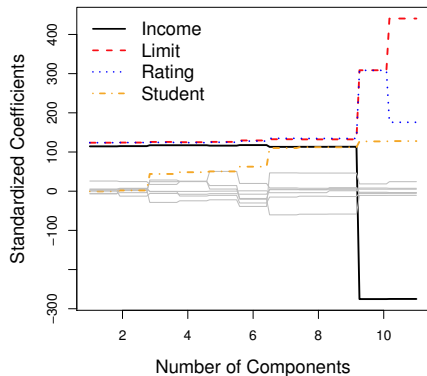
Linear Model Selection and Regularization

Principal components regression

- PCR uses $M < p$ new variables, but these all still depend on original predictors.
- Hence, PCR not a feature selection method.
- In this aspect, PCR closer to ridge regression than lasso.
- Ridge regression can be viewed as a continuous version of PCR.
- # principal components M can be chosen by CV.
- Recommended: first standardize data.

Linear Model Selection and Regularization

Principal components regression



PCR applied to [Credit](#) data set.

Left: standardized coefficients. Right: CV MSE against M .

Lowest error for 10 components (only one less than full model).

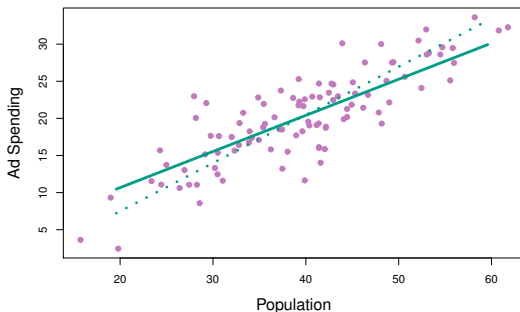
Linear Model Selection and Regularization

Partial least squares

- PCR only looks at predictor variability, not at response.
- In this sense, it is **unsupervised**.
- **Partial least squares** (PLS): **supervised** variant of PCR: find linear combination of predictors containing most variability and best explain response.
- To construct Z_1 , set each coefficient in $Z_1 = \sum_{j=1}^p \phi_{j,1} X_j$ to coefficient of simple linear regression of Y onto X_j . Results in coefficient proportional to $\text{Cor}(X_j, Y)$.
This places highest weight on variables most strongly related to response Y .
- To identify Z_2 , first adjust all predictors for Z_1 by regressing these on Z_1 and taking residuals. Interpretation: remaining information not explained by first PLS direction. Compute Z_2 using this orthogonalized data just as Z_1 was computed using original data.
- In the same way, compute further PLS directions Z_3, \dots, Z_M .

Linear Model Selection and Regularization

Partial least squares



PLS on synthetic data set giving **Sales** data in each of 100 regions as response to two predictors **Population Size** and **Advertising Spending**. Solid line: first PLS direction, dotted: first principal components direction.

⑥ Linear Model Selection and Regularization

6.1 Subset Selection

6.2 Shrinkage Methods

6.3 Dimension Reduction Methods

6.4 Considerations in High Dimensions

6.5 Miscellanea

Linear Model Selection and Regularization

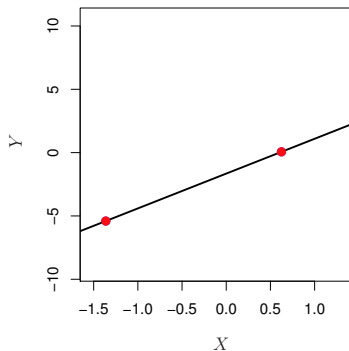
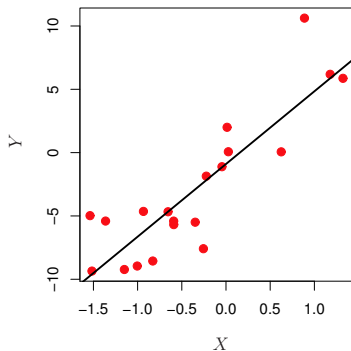
The high-dimensional setting

- Most traditional statistical techniques: $n \gg p$ (**low-dimensional setting**).
- Typical example: Predict patient's blood pressure based on age, gender, body mass index (BMI). Three predictors, and typically thousands of patients' data.
- More recently, in many fields such as medicine, finance, marketing, trend towards collecting almost unlimited number of feature measurements (p large), while cost of obtaining sufficiently many samples prohibitive.
- Example: in place of age, gender, BMI, collect measurements of half million **single nucleotide polymorphisms**, i.e., common individual DNA mutations. Results in $p \approx 500,000$, $n \approx 200$.
- Example: 'Bag-of-words' model to understand customers' online shopping patterns, using as features all search terms entered in search engine (binary feature vector). Only few hundred users consented to their data being used. Results in $n \approx 100$, p much larger.

Linear Model Selection and Regularization

The high-dimensional setting: what goes wrong?

- When $p \geq n$ LS cannot (should not) be used, since data will be fit perfectly.
- Example: $p = 1$ $n = 20$ vs. $n = 2$:

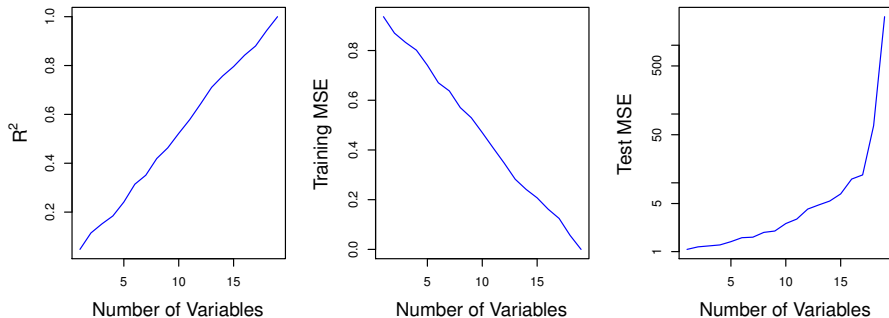


Right model will not generalize well (overfitting), model too flexible..

Linear Model Selection and Regularization

The high-dimensional setting: what goes wrong?

Another example: $n = 20$ observations for $1 \leq p \leq 20$ features, each *completely unrelated* to response.



As p increases, $R^2 \rightarrow 1$, training MSE $\rightarrow 0$ despite no relation of predictors to response. At the same time, test MSE sharply increases as model increasingly flexible.

Casual observer may find large model superior if only first two quantities monitored.

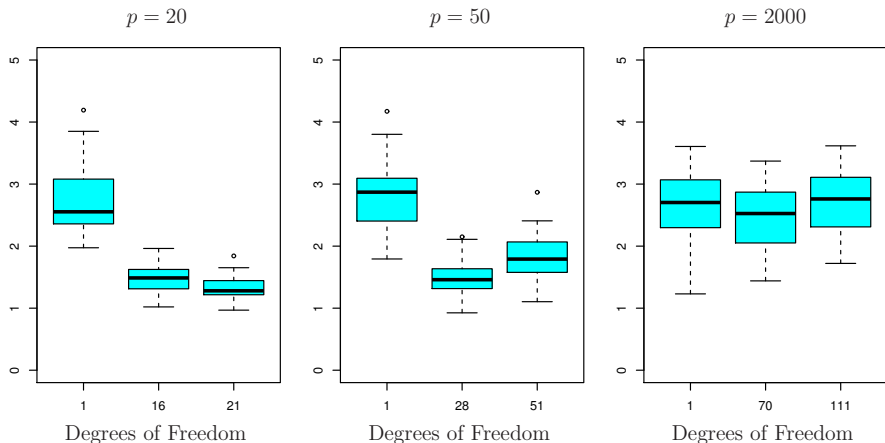
Linear Model Selection and Regularization

The high-dimensional setting: what goes wrong?

- Model selection techniques based on C_p , AIC, BIC not appropriate for high-dimensional setting, as estimating $\hat{\sigma}^2$ problematic.
- Adjusted R^2 may easily yield value of 1 in high-dimensional setting.
- Less flexible regression models (stepwise selection, shrinkage, PCR) particularly useful in high dimensions. Avoid overfitting by constraining flexibility.
- Next figure: Lasso on $n = 100$ simulated training observations using $p = 20, 50$ and $2,000$ features, of which 20 related to response.
Then MSE evaluated on independent test set.

Linear Model Selection and Regularization

The high-dimensional setting: what goes wrong?



- For $p = 20$, lowest test MSE for low value of λ . For larger p , best model obtained for larger λ . When $p = 2000$, lasso performs badly for all values of λ .
- Rather than λ , plot shows **degrees of freedom** of model, i.e., # nonzero coefficients of lasso estimate.

Linear Model Selection and Regularization

The high-dimensional setting: what goes wrong?

Summary:

- ➊ Shrinkage plays key role in high dimensions.
 - ➋ Correct value of tuning parameter essential.
 - ➌ Test error increases with dimension, unless additional features informative.
- Third observation related to **curse of dimensionality**: quality of model need not increase as features added.
 - Compare left and right panel in figure: test MSE almost doubles as p increased from 20 to 2000.
 - **Noise features** (not related to response) increase dimension, exacerbate overfitting danger.
 - Adding features truly related to response will generally improve model.
 - New sensor technology allowing for millions of observations can lead to worse results if features not relevant. Even if relevant, variance incurred by fitting their coefficients may outweigh reduction in bias from additional features.

Linear Model Selection and Regularization

The high-dimensional setting: what goes wrong?

- In high dimensions: **collinearity** problem extreme. (Why?)
- Never know which variables truly predictive, can never obtain best coefficients.
- At best: assign large coefficients to variables correlated with variables truly predictive for response.
- For $p > n$ can easily obtain useless model with zero residual.
- Traditional measures of model quality based on training data often highly misleading in high dimensions.
- Reporting MSE on independent test data particularly important here.

⑥ Linear Model Selection and Regularization

6.1 Subset Selection

6.2 Shrinkage Methods

6.3 Dimension Reduction Methods

6.4 Considerations in High Dimensions

6.5 Miscellanea

Optimality of LS Estimate

The Gauss-Markov⁷ theorem

Theorem 6.1 (Gauss-Markov theorem)

Given observations $y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + \varepsilon_i$, $i = 1, \dots, n$, for which the uncorrelated random noise variables ε_i have mean zero and constant variance $\sigma^2 > 0$, and assuming that the observation vectors $\mathbf{x}_1, \dots, \mathbf{x}_p \in \mathbb{R}^n$ are linearly independent, then the least squares estimate

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}, \quad \mathbf{X} = [\mathbf{x}_1 | \dots | \mathbf{x}_p] \in \mathbb{R}^{n \times p}, \quad \mathbf{y} \in \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

has minimal variance among all linear unbiased estimators of $\boldsymbol{\beta}$.

⁷C.F. Gauss, 1777–1855, A.A. Markov, 1856–1922

Optimality of LS Estimate

Remarks

- No assumption is made on the distribution of the errors, only on their first two moments.
- The theorem also holds if $\mathbf{Var} \varepsilon = \mathbf{\Sigma}$ is a (nonsingular) covariance matrix. In this case the best linear unbiased estimator solves the **weighted least squares** problem

$$\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_{\mathbf{\Sigma}} \rightarrow \min_{\boldsymbol{\beta}}, \quad \|\mathbf{x}\|_{\mathbf{\Sigma}}^2 = \mathbf{x}^T \mathbf{\Sigma}^{-1} \mathbf{x}.$$

- The theorem does not say there are no better estimators than LS, only that any better estimators are either nonlinear or biased.
- Examples of biased estimators are ridge regression and the lasso.
- ESL:

"Most models are distortions of the truth, and hence are biased; picking the right model amounts to creating the right balance between bias and variance."

Singular Value Decomposition

Definition

Theorem 6.2 (Singular value decomposition)

For any matrix $\mathbf{A} \in \mathbb{R}^{n \times p}$ of rank r , there exist orthogonal matrices $\mathbf{U} \in \mathbb{R}^{n \times n}$ and $\mathbf{V} \in \mathbb{R}^{p \times p}$ as well as a “diagonal” matrix

$$\mathbf{\Sigma} = \begin{bmatrix} \mathbf{\Sigma}_r & \mathbf{O} \\ \mathbf{O} & \mathbf{O} \end{bmatrix} \in \mathbb{R}^{n \times p} \quad \text{where} \quad \mathbf{\Sigma}_r = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_r) \in \mathbb{R}^{r \times r}$$

and $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$, such that

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top. \quad (\text{SVD})$$

- The positive numbers $\sigma_1, \dots, \sigma_r$ are called the **singular values** of \mathbf{A} .
- The columns of $\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n]$ and $\mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_p]$ are the **left and right singular vectors**, respectively.

Singular Value Decomposition

Properties

- 1 Representation of \mathbf{A} as sum of rank-1 matrices:

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r] \mathbf{\Sigma}_r [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r]^\top = \sum_{k=1}^r \sigma_k \mathbf{u}_k \mathbf{v}_k^\top$$

- 2 Singular vector mapping properties:

$$\mathbf{A}\mathbf{v}_k = \begin{cases} \sigma_k \mathbf{u}_k & k = 1, 2, \dots, r, \\ \mathbf{0} & k = r+1, \dots, p \end{cases}$$

and

$$\mathbf{A}^\top \mathbf{u}_k = \begin{cases} \sigma_k \mathbf{v}_k & k = 1, 2, \dots, r, \\ \mathbf{0} & k = r+1, \dots, n. \end{cases}$$

3

$\{\mathbf{u}_1, \dots, \mathbf{u}_r\}$	is an ON-basis of	$\mathcal{R}(\mathbf{A})$.
$\{\mathbf{u}_{r+1}, \dots, \mathbf{u}_m\}$	is an ON-basis of	$\mathcal{N}(\mathbf{A}^\top) = \mathcal{R}(\mathbf{A})^\perp$.
$\{\mathbf{v}_1, \dots, \mathbf{v}_r\}$	is an ON-basis of	$\mathcal{R}(\mathbf{A}^\top) = \mathcal{N}(\mathbf{A})^\perp$.
$\{\mathbf{v}_{r+1}, \dots, \mathbf{v}_n\}$	is an ON-basis of	$\mathcal{N}(\mathbf{A})$.

Singular Value Decomposition

Properties

4 Eigenspaces of $\mathbf{A}\mathbf{A}^\top$ and $\mathbf{A}^\top\mathbf{A}$:

- $\sigma_1^2, \dots, \sigma_r^2$ are the non-zero eigenvalues of $\mathbf{A}^\top\mathbf{A}$ and $\mathbf{A}\mathbf{A}^\top$, respectively:

$$\mathbf{A}^\top\mathbf{A} = \mathbf{V}\mathbf{\Sigma}^\top\mathbf{\Sigma}\mathbf{V}^\top = \mathbf{V} \begin{bmatrix} \mathbf{\Sigma}_r^2 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{V}^\top,$$

$$\mathbf{A}\mathbf{A}^\top = \mathbf{U}\mathbf{\Sigma}\mathbf{\Sigma}^\top\mathbf{U}^\top = \mathbf{U} \begin{bmatrix} \mathbf{\Sigma}_r^2 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{U}^\top.$$

- In particular, the singular values $\sigma_1, \dots, \sigma_r$ are uniquely determined by \mathbf{A} .
- The right singular vectors $\mathbf{v}_1, \dots, \mathbf{v}_p$ form an ON-basis of \mathbb{R}^p of eigenvectors of $\mathbf{A}^\top\mathbf{A}$:

$$\mathbf{A}^\top\mathbf{A}\mathbf{v}_k = \begin{cases} \sigma_k^2\mathbf{v}_k & k = 1, 2, \dots, r, \\ \mathbf{0} & k = r + 1, \dots, p. \end{cases}$$

The left singular vectors $\mathbf{u}_1, \dots, \mathbf{u}_n$ form an ON-basis of \mathbb{R}^n of eigenvectors of $\mathbf{A}\mathbf{A}^\top$:

$$\mathbf{A}\mathbf{A}^\top\mathbf{u}_k = \begin{cases} \sigma_k^2\mathbf{u}_k & k = 1, 2, \dots, r, \\ \mathbf{0} & k = r + 1, \dots, n. \end{cases}$$

Singular Value Decomposition

Properties

- 5 If $\mathbf{A} = \mathbf{A}^\top \in \mathbb{R}^{n \times n}$ with non-zero eigenvalues

$$\lambda_1, \dots, \lambda_r, \quad |\lambda_1| \geq \dots \geq |\lambda_r| > 0,$$

then the singular values of \mathbf{A} are given by $\sigma_k = |\lambda_k|$.

- 6 The (p -dimensional) unit sphere is mapped by \mathbf{A} to an ellipsoid (in \mathbb{R}^n) with center $\mathbf{0}$ and semi-axes $\sigma_k \mathbf{u}_k$ ($\sigma_k := 0$ für $k > r$).
- 7 For $\mathbf{A} \in \mathbb{R}^{n \times p}$ there holds $\|\mathbf{A}\|_2 = \sigma_1$ and $\|\mathbf{A}\|_F = \sqrt{\sigma_1^2 + \dots + \sigma_r^2}$.
For $\mathbf{A} \in \mathbb{R}^{n \times n}$ invertible, there holds in addition that $\|\mathbf{A}^{-1}\|_2 = \sigma_n^{-1}$.
- 8 Analogous statements hold for complex-valued matrices $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^H$ (\mathbf{U}, \mathbf{V} unitary). (In (5) replace $\mathbf{A} = \mathbf{A}^\top$ by ' \mathbf{A} normal'.)

Singular Value Decomposition

Best rank- k approximation

- E. Schmidt. Zur Theorie der linearen und nichtlinearen Integralgleichungen. 1. Teil: Entwicklung willkürlicher Funktionen nach Systemen vorgeschriebener. Math. Ann., 63 (1907), pp. 433–476
- C. Eckart, G. Young. The approximation of one matrix by another of lower rank. Psychometrika, 1 (1936), pp. 211–218
- L. Mirsky. Symmetric gauge functions and unitarily invariant norms. Quart. J. Math. Oxford, 11 (1960), pp. 50–59

Theorem 6.3 (Best approximation by matrices of lower rank)

For a matrix $\mathbf{A} \in \mathbb{R}^{n \times p}$ of rank r with SVD $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$ the best approximation problem

$$\min\{\|\mathbf{A} - \mathbf{B}\|_2 : \mathbf{B} \in \mathbb{R}^{n \times p} \text{ and } \text{rank}(\mathbf{B}) \leq k\}$$

for $k < r$ is solved by

$$\mathbf{A}_k := \sum_{i=1}^k \sigma_i \mathbf{u}_i \mathbf{v}_i^\top \quad \text{with} \quad \|\mathbf{A} - \mathbf{A}_k\|_2 = \sigma_{k+1}.$$

- \mathbf{A}_k as above is also the closest rank- k matrix to \mathbf{A} in the Frobenius-norm $\|\cdot\|_F$, with distance $\|\mathbf{A} - \mathbf{A}_k\|_F = \sqrt{\sigma_{k+1}^2 + \cdots + \sigma_r^2}$.

Singular Value Decomposition

Ridge regression

- Recall the ridge regression estimate $\hat{\boldsymbol{\beta}}_R$ for the LS problem $\mathbf{X}\boldsymbol{\beta} \approx \mathbf{y}$ with data matrix $\mathbf{X} \in \mathbb{R}^{n \times p}$ and observation vector $\mathbf{y} \in \mathbb{R}^n$: for a given value of the tuning (or regularization) parameter $\lambda \geq 0$ it was defined by

$$\hat{\boldsymbol{\beta}}_R = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^p} Q_\lambda(\boldsymbol{\beta}), \quad Q_\lambda(\boldsymbol{\beta}) := \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_2^2.$$

- Rewriting the objective function $Q_\lambda(\boldsymbol{\beta})$ as

$$\begin{aligned} Q_\lambda(\boldsymbol{\beta}) &= (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \boldsymbol{\beta}^\top \boldsymbol{\beta} = \begin{bmatrix} \mathbf{y} - \mathbf{X}\boldsymbol{\beta} \\ \sqrt{\lambda} \boldsymbol{\beta} \end{bmatrix}^\top \begin{bmatrix} \mathbf{y} - \mathbf{X}\boldsymbol{\beta} \\ \sqrt{\lambda} \boldsymbol{\beta} \end{bmatrix} \\ &= \left\| \begin{bmatrix} \mathbf{y} \\ \mathbf{0} \end{bmatrix} - \begin{bmatrix} \mathbf{X} \\ \sqrt{\lambda} \mathbf{I} \end{bmatrix} \boldsymbol{\beta} \right\|_2^2, \end{aligned}$$

we observe that ridge regression can be viewed as a standard LS formulation for the augmented problem

$$\begin{bmatrix} \mathbf{X} \\ \sqrt{\lambda} \mathbf{I} \end{bmatrix} \boldsymbol{\beta} \approx \begin{bmatrix} \mathbf{y} \\ \mathbf{0} \end{bmatrix}.$$

Singular Value Decomposition

Ridge regression

- The associated normal equations of ridge regression

$$(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}) \boldsymbol{\beta} = \mathbf{X}^\top \mathbf{y} \quad (6.14)$$

are obtained from those of original LS problem by adding $\lambda \mathbf{I}$ to the coefficient matrix, guaranteeing positive definiteness for $\lambda > 0$.

- Given an SVD $\mathbf{X} = \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^\top$ of the data matrix \mathbf{X} with orthogonal matrices $\mathbf{U} = [\mathbf{u}_1 | \dots | \mathbf{u}_n] \in \mathbb{R}^{n \times n}$, $\mathbf{V} = [\mathbf{v}_1 | \dots | \mathbf{v}_p] \in \mathbb{R}^{p \times p}$ and, and assuming it has full rank $p \leq n$, $\boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_p \\ \mathbf{O} \end{bmatrix}$, $\boldsymbol{\Sigma}_p = \text{diag}(\sigma_1, \dots, \sigma_p)$, $\sigma_1 \geq \dots \geq \sigma_p > 0$, this implies

$$\mathbf{X}^\top \mathbf{X} = \mathbf{V} \boldsymbol{\Sigma}^\top \boldsymbol{\Sigma} \mathbf{V}^\top, \quad \boldsymbol{\Sigma}^\top \boldsymbol{\Sigma} = \text{diag}(\sigma_1^2, \dots, \sigma_p^2),$$

$$\mathbf{X}^\top \mathbf{y} = \mathbf{V} \boldsymbol{\Sigma}^\top \mathbf{U}^\top \mathbf{y} = \sum_{j=1}^p \sigma_p (u_j^\top \mathbf{y}) \mathbf{v}_j$$

Singular Value Decomposition

Ridge regression

- Inserting these expressions into the normal equations (6.14) yields

$$\mathbf{V}(\mathbf{\Sigma}^\top \mathbf{\Sigma} + \lambda \mathbf{I})\mathbf{V}^\top \boldsymbol{\beta} = \mathbf{V}\mathbf{\Sigma}^\top \mathbf{U}^\top \mathbf{y}$$

or, setting $\boldsymbol{\gamma} := \mathbf{V}^\top \boldsymbol{\beta}$,

$$(\mathbf{\Sigma}^\top \mathbf{\Sigma} + \lambda \mathbf{I})\boldsymbol{\gamma} = \mathbf{\Sigma}^\top \mathbf{U}^\top \mathbf{y}, \quad \text{giving} \quad \gamma_j = \frac{\sigma_j}{\sigma_j^2 + \lambda} \mathbf{u}_j^\top \mathbf{y}, \quad j = 1, \dots, p,$$

and finally, with $\boldsymbol{\beta} = \mathbf{V}\boldsymbol{\gamma}$, the ridge regression estimate

$$\hat{\boldsymbol{\beta}}_R = \sum_{j=1}^p \frac{\sigma_j}{\sigma_j^2 + \lambda} (\mathbf{u}_j^\top \mathbf{y}) \mathbf{v}_j.$$

- Observe that $\hat{\boldsymbol{\beta}}_R$ is obtained from the standard LS estimate $\hat{\boldsymbol{\beta}} = \sum_{j=1}^p \frac{\mathbf{u}_j^\top \mathbf{y}}{\sigma_j} \mathbf{v}_j$ by multiplying each coefficient with the **filter factor**

$$\frac{\sigma_j^2}{\sigma_j^2 + \lambda}, \quad j = 1, \dots, p.$$

- Given SVD, ridge regression estimates for additional λ essentially for free.

Principal Components

Covariance matrix of a random vector

- Recall: the variance of a random variable X with expectation $\mu := \mathbf{E}[X]$ is given by

$$\sigma^2 = \mathbf{Var} X = \mathbf{E}[(X - \mu)^2].$$

- For a **random vector** $X = (X_1, \dots, X_p)^\top \in \mathbb{R}^p$ with expectation $\mu := \mathbf{E}[X]$, the **variance** or **covariance matrix** is given by

$$\mathbf{C} := \mathbf{Var} X = \mathbf{E}[(X - \mu)(X - \mu)^\top] = \mathbf{C}^\top \in \mathbb{R}^{p \times p},$$

with matrix entries

$$C_{i,j} = \mathbf{E}[(X_i - \mu_i)(X_j - \mu_j)] = \mathbf{Cov}(X_i, X_j), \quad i, j = 1, \dots, p.$$

Principal Components

Total variance a random vector

- A scalar measure of the total variance contained in a random vector $X \in \mathbb{R}^p$ is provided by the **trace** of its covariance matrix

$$\text{tr } \mathbf{C} = \sum_{j=1}^p C_{j,j} = \sum_{j=1}^p \mathbf{Cov}(X_j, X_j) = \sum_{j=1}^p \mathbf{Var } X_j.$$

- Justification:

$$\begin{aligned} \mathbf{E} [\|X - \mathbf{E}[X]\|_2^2] &= \mathbf{E} [\|X - \boldsymbol{\mu}\|_2^2] = \mathbf{E} [(X - \boldsymbol{\mu})^\top (X - \boldsymbol{\mu})] \\ &= \mathbf{E} \left[\sum_{j=1}^p (X_j - \mu_j)^2 \right] = \sum_{j=1}^p \mathbf{E} [(X_j - \mu_j)^2] = \sum_{j=1}^p \mathbf{Var } X_j. \end{aligned}$$

- By a well-known result from linear algebra, if $\lambda_j(\mathbf{C})$ denotes the j -th eigenvalue (in descending order) of \mathbf{C} ⁸, there also holds

$$\text{tr } \mathbf{C} = \sum_{j=1}^p \lambda_j(\mathbf{C}).$$

⁸Note that these are real and positive as \mathbf{C} is symmetric and positive-definite.

Principal Components

Total variance a random vector

- Given a spectral decomposition

$$\mathbf{C} = \mathbf{W}\mathbf{\Lambda}\mathbf{W}^\top, \quad \mathbf{W}^\top\mathbf{W} = \mathbf{I}, \quad \mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_p),$$

of \mathbf{C} and the fact that the **Frobenius norm** $\|\cdot\|_F$ is unitarily invariant, we also have

$$\text{tr } \mathbf{C} = \|\mathbf{\Lambda}^{1/2}\|_F^2 = \|\mathbf{W}\mathbf{\Lambda}^{1/2}\mathbf{W}^\top\|_F^2 = \|\mathbf{C}^{1/2}\|_F^2.$$

- In view of the fact that $|\lambda_j(\mathbf{C})| = \lambda_j(\mathbf{C})$ for covariance matrices, the spectral decomposition $\mathbf{W}\mathbf{\Lambda}\mathbf{W}^\top$ is also a singular value decomposition.
- Combining with Theorem 6.3, we conclude that for any $k \in \{1, \dots, p\}$ the matrix

$$\mathbf{C}_k = \sum_{j=1}^k \lambda_j \mathbf{w}_j \mathbf{w}_j^\top,$$

where $\mathbf{W} = [\mathbf{w}_1 | \dots | \mathbf{w}_p]$, is the best approximation of the covariance matrix \mathbf{C} in the spectral and Frobenius norms among all matrices of rank $\leq k$.

Principal Components

Linear combinations of random vector components

- Given a random vector $X = (X_1, \dots, X_p)^\top \in \mathbb{R}^p$ and \mathbf{w}_j a normalized eigenvector of its covariance matrix \mathbf{C} with associated eigenvalue λ_j , define the scalar random variable $Z_j := \mathbf{w}_j^\top X$. Then

$$\begin{aligned}\mathbf{Var} Z_j &= \mathbf{E} [(\mathbf{w}_j^\top X - \mathbf{E} [\mathbf{w}_j^\top X])^2] = \mathbf{E} [(\mathbf{w}_j^\top (X - \boldsymbol{\mu}))^2] \\ &= \mathbf{E} [(\mathbf{w}_j^\top (X - \boldsymbol{\mu}))(X - \boldsymbol{\mu})^\top \mathbf{w}_j] = \mathbf{w}_j^\top \mathbf{E} [(X - \boldsymbol{\mu})(X - \boldsymbol{\mu})^\top] \mathbf{w}_j \\ &= \mathbf{w}_j^\top \mathbf{C} \mathbf{w}_j = \lambda_j.\end{aligned}$$

- More generally, for any linear combination $Z = \boldsymbol{\phi}^\top X$, $\boldsymbol{\phi} = (\phi_1, \dots, \phi_p)^\top$, we have

$$\begin{aligned}\mathbf{Var} Z &= \mathbf{E} [(\boldsymbol{\phi}^\top X - \mathbf{E} [\boldsymbol{\phi}^\top X])^2] = \mathbf{E} [(\boldsymbol{\phi}^\top (X - \boldsymbol{\mu}))^2] \\ &= \mathbf{E} \left[\left(\sum_{j=1}^p \phi_j (X_j - \mu_j) \right)^2 \right] = \sum_{j,k=1}^p \phi_j \phi_k \mathbf{E} [(X_j - \mu_j)(X_k - \mu_k)] \\ &= \boldsymbol{\phi}^\top \mathbf{C} \boldsymbol{\phi}.\end{aligned}$$

Principal Components

Linear combinations of random vector components

- For two general linear combinations $Z_1 = \phi_1^\top X$, $Z_2 = \phi_2^\top X$, we conclude by an analogous calculation that

$$\mathbf{Cov}(Z_1, Z_2) = \phi_2^\top \mathbf{C} \phi_1$$

and therefore that Z_1 and Z_2 are **uncorrelated** if and only if $\phi_2^\top \mathbf{C} \phi_1 = 0$, i.e., if the coefficient vectors ϕ_1 and ϕ_2 are orthogonal in the inner product generated by the (symmetric and positive definite) matrix \mathbf{C} .

- If we seek a change of variables $Z = \Phi^\top X$ with a nonsingular $\Phi \in \mathbb{R}^{p \times p}$ such that the components of Z are uncorrelated with unit variance, then it is necessary that

$$\mathbf{I} = \mathbf{E}[(Z - \mathbf{E}[Z])(Z - \mathbf{E}[Z])^\top] = \mathbf{E}[\Phi^\top (X - \mu)(X - \mu)^\top \Phi] = \Phi^\top \mathbf{C} \Phi.$$

The set of all matrices $\Phi \in \mathbb{R}^{p \times p}$ which achieve this is precisely the set of all **congruences** taking \mathbf{C} to \mathbf{I} .

Principal Components

Linear combinations of random vector components

- **Example 1:** given **Cholesky factorization** $C = LL^T$, choosing $\Phi := L^{-T}$ gives

$$\Phi^T C \Phi = L^{-1}(LL^T)L^{-T} = I.$$

- **Example 2:** given **spectral decomposition** $C = W\Lambda W^T$, choosing $\Phi := W\Lambda^{-1/2}$ gives

$$\Phi^T C \Phi = \Lambda^{-1/2}W^T(W\Lambda W^T)W\Lambda^{-1/2} = I.$$

- **Example 3:** given **square-root-free Cholesky factorization** $C = LDL^T$, where L is lower triangular with a unit diagonal and D is diagonal, choosing $\Phi := L^{-T}$ gives

$$\Phi^T C \Phi = L^{-1}(LDL^T)L^{-T} = D.$$

- **Example 4:** given **spectral decomposition** $C = W\Lambda W^T$, choosing $\Phi := W$ gives

$$\Phi^T C \Phi = W^T(W\Lambda W^T)W = \Lambda.$$

Principal Components

Courant-Fischer min-max-characterization

For a square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ the expression $\frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}}$, $\mathbf{0} \neq \mathbf{x} \in \mathbb{R}^n$, is called a **Rayleigh quotient**.

Theorem 6.4 (Fischer, 1905; Courant, 1920)

Let $\mathbf{A} \in \mathbb{R}^{n \times n}$ be a symmetric matrix with eigenvalues $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ and $k \in \{1, 2, \dots, n\}$. Then

$$\lambda_k = \min_{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_{n-k} \in \mathbb{R}^n} \max_{\substack{\mathbf{0} \neq \mathbf{x} \in \mathbb{R}^n \\ \mathbf{x} \perp \mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_{n-k}}} \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}}, \quad (6.15)$$

$$\lambda_k = \max_{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_{k-1} \in \mathbb{R}^n} \min_{\substack{\mathbf{0} \neq \mathbf{x} \in \mathbb{R}^n \\ \mathbf{x} \perp \mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_{k-1}}} \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}} \quad (6.16)$$

- The extremal values of the Rayleigh quotient are attained when \mathbf{x} is an eigenvector associated with λ_1 or λ_n , respectively.

Principal Components

Courant-Fischer min-max-characterization

Consequences of Theorem 6.4:

- Linear combination $\boldsymbol{\phi}^\top X$ where $\|\boldsymbol{\phi}\|_2 = 1$ with maximal variance obtained for $\boldsymbol{\phi} = \boldsymbol{\phi}_1 = \mathbf{w}_1$. This is the first principal component.
- Linear combination $\boldsymbol{\phi}^\top X$ where $\|\boldsymbol{\phi}\|_2 = 1$ with maximal variance subject to $\boldsymbol{\phi} \perp \mathbf{w}_1$ obtained for $\boldsymbol{\phi} = \boldsymbol{\phi}_2 = \mathbf{w}_2$ (second principal component).
- Linear combination $\boldsymbol{\phi}^\top X$ where $\|\boldsymbol{\phi}\|_2 = 1$ with maximal variance subject to $\boldsymbol{\phi} \perp \mathbf{w}_1, \dots, \mathbf{w}_{j-1}$ obtained for $\boldsymbol{\phi} = \boldsymbol{\phi}_j = \mathbf{w}_j$ (j -th principal component).
- The change of variables afforded by replacing the original random variables X_1, \dots, X_p by the principal components $Z = \mathbf{W}^\top X$ is the (unscaled) congruence obtained from the spectral decomposition.

The total variance contained in Z is given by

$$\mathbf{E} [\|Z - \mathbf{E}[Z]\|_2^2] = \sum_{j=1}^p \mathbf{Var} Z_j = \sum_{j=1}^p \lambda_j = \text{tr} \sum_{j=1}^p \lambda_j \mathbf{w}_j \mathbf{w}_j^\top = \text{tr} \mathbf{C},$$

which coincides with the total variance contained in X .

Principal Components

PCR

- Performing regression of a data vector \mathbf{y} on $M < p$ principal components results in **principal components regression** (PCR).
- The total variance contained in random vector $(Z_1, \dots, Z_M)^\top$ is

$$\mathbf{E} [\|Z - \mathbf{E}[Z]\|_2^2] = \sum_{j=1}^M \mathbf{Var} Z_j = \sum_{j=1}^M \lambda_j = \text{tr} \sum_{j=1}^M \lambda_j \mathbf{w}_j \mathbf{w}_j^\top = \text{tr} \mathbf{C}_k.$$

- The fraction of neglected variance in PCR using M principal components is

$$\frac{\sum_{j=k+1}^p \lambda_j}{\sum_{j=1}^p \lambda_j} = 1 - \frac{\sum_{j=1}^k \lambda_j}{\sum_{j=1}^p \lambda_j}.$$

Principal Components

Data

- The covariance matrix \mathbf{C} and expectation vector $\boldsymbol{\mu}$ are theoretical constructs and typically unavailable hence estimated from data.
- As usual, we denote the data matrix (design matrix) by

$$\mathbf{X} = \begin{bmatrix} x_{1,1} & \dots & x_{1,p} \\ \vdots & & \vdots \\ x_{n,1} & \dots & x_{n,p} \end{bmatrix} = [\mathbf{x}_1 | \dots | \mathbf{x}_p] \in \mathbb{R}^{n \times p},$$

each column corresponding to one of p predictor variables (features) and each row to one of n observations (samples, realizations).

- We denote the vector of **sample means** by $\bar{\mathbf{X}} := \frac{1}{n} \mathbf{e}^\top \mathbf{X} = [\bar{x}_1, \dots, \bar{x}_p]$ and obtain the **centered data matrix** as

$$\tilde{\mathbf{X}} := [\mathbf{x}_1 - \bar{x}_1 \mathbf{e} | \dots | \mathbf{x}_p - \bar{x}_p \mathbf{e}] = \mathbf{X} - \mathbf{e} \bar{\mathbf{X}} = \mathbf{X} - \mathbf{e} \frac{1}{n} \mathbf{e}^\top \mathbf{X} = (\mathbf{I} - \frac{1}{n} \mathbf{e} \mathbf{e}^\top) \mathbf{X}.$$

- Finally, the **unbiased sample covariance matrix** is

$$\mathbf{S}_n := \frac{1}{n-1} \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} = \frac{1}{n-1} \mathbf{X} (\mathbf{I} - \frac{1}{n} \mathbf{e} \mathbf{e}^\top)^2 \mathbf{X} = \frac{1}{n-1} \mathbf{X} (\mathbf{I} - \frac{1}{n} \mathbf{e} \mathbf{e}^\top) \mathbf{X}.$$

Principal Components

Data

- In practice the sample covariance matrix \mathbf{S}_n takes the place of the covariance matrix \mathbf{C} .
- For PCA/PCR, one can compute a spectral decomposition of \mathbf{S}_n .
- Alternatively, given an SVD $\tilde{\mathbf{X}} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$, a spectral decomposition of \mathbf{S}_n is obtained as

$$\mathbf{S}_n = \frac{1}{n-1} \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} = \frac{1}{n-1} \mathbf{V} \mathbf{\Sigma}^\top \mathbf{\Sigma} \mathbf{V}^\top.$$

- The SVD approach is generally numerically stabler, in particular if $\tilde{\mathbf{X}}$ is ill-conditioned. The spectral decomposition may be cheaper, as $\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}}$ is smaller than $\tilde{\mathbf{X}}$.