

Chapter 6: Linear Model Selection and Regularization (Lecture 2)

Thiago G. Martins | NTNU & Verizon Media
Spring 2020

Previous lecture

Subset selection and shrinkage methods

- Subset selection and shrinkage methods have controlled variance in two ways:
 - Using a subset of the original predictors.
 - Shrinking their coefficients towards zero.
- Those methods use the original (possibly standardized) predictors X_1, \dots, X_p .

Dimension reduction methods

Dimension reduction methods

- Transform the original predictors

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

for $m = 1, \dots, M, j = 1, \dots, p$ and $M < p$

- If the constants ϕ_{jm} are chosen wisely, then such dimension reduction approaches can often outperform least squares regression.
- The dimension of the problem has been reduced from $p + 1$ to $M + 1$.

Constrained interpretation

- It can be shown that

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

- Such as large p in relation to n

Outline

- We will cover two approaches to dimensionality reduction:
 - Principal Components
 - Partial Least Squares

Principal Component Analysis (PCA)

Principal Component Analysis (PCA)

- Discussed in greater detail in Chapter 10 about unsupervised learning
- Focus in this lecture is how it can be applied for regression.
 - That is, in a supervised setting.

Principal Component Analysis (PCA)

- Discussed in greater detail in Chapter 10 about unsupervised learning
- Focus in this lecture is how it can be applied for regression.
 - That is, in a supervised setting.
- PCA is a (unsupervised) technique for reducing the dimension of a $n \times p$ data matrix X .

Principal Component Analysis (PCA)

- We want to create a $n \times M$ matrix Z , with $M < p$.
- The column Z_m of Z is the m -th principal component.

$$Z_m = \sum_{j=1}^p \phi_{jm} X_j \quad \text{subject to} \quad \sum_{j=1}^p \phi_{jm}^2 = 1$$

Principal Component Analysis (PCA)

- We want to create a $n \times M$ matrix Z , with $M < p$.
- The column Z_m of Z is the m -th principal component.

$$Z_m = \sum_{j=1}^p \phi_{jm} X_j \quad \text{subject to} \quad \sum_{j=1}^p \phi_{jm}^2 = 1$$

- We want Z_1 to have the highest possible variance.
 - That is, take the direction of the data where the observations vary the most.
 - Without the constrain we could get higher variance by increasing ϕ_j

Principal Component Analysis (PCA)

- Z_2 should be uncorrelated to Z_1 , and have the highest variance, subject to this constrain.
 - The direction of Z_1 must be perpendicular (or orthogonal) to the direction of Z_2

Principal Component Analysis (PCA)

- Z_2 should be uncorrelated to Z_1 , and have the highest variance, subject to this constrain.
 - The direction of Z_1 must be perpendicular (or orthogonal) to the direction of Z_2
- And so on ...
- We can construct up to p PCs that way.
 - In which case we have captured all the variability contained in the data
 - We have created a set of orthogonal predictors
 - But have **not** accomplished dimensionality reduction

- 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.
- The blue dashed line indicates the second principal component.

- A subset of the advertising data.
- Left: 1st PC
 - The dimension along which the data vary the most
 - The line that is closest to all n of the observations.
- Right: Rotated so that the 1st PC direction coincides with the x-axis.

- Plots 1st PC scores versus population and ad spending. The relationships are strong.
- Strong relationship: the 1st PC appears to capture most of the information contained in the pop and ad predictors.

- There is little relationship between 2nd PC and predictors
- Suggesting one only needs the first principal component in order to accurately represent the pop and ad budgets.

PCA - Overview

- Principal component analysis (PCA) is a dimensionality reduction technique
 - Our ability to visualize data is limited to 2 or 3 dimensions.

PCA - Overview

- Principal component analysis (PCA) is a dimensionality reduction technique
 - Our ability to visualize data is limited to 2 or 3 dimensions.
 - Lower dimension can reduce numerical algorithms computational time.

PCA - Overview

- Principal component analysis (PCA) is a dimensionality reduction technique
 - Our ability to visualize data is limited to 2 or 3 dimensions.
 - Lower dimension can reduce numerical algorithms computational time.
 - Many statistical models suffer from high correlation between covariates

PCA - Overview

- Principal component analysis (PCA) is a dimensionality reduction technique
 - Our ability to visualize data is limited to 2 or 3 dimensions.
 - Lower dimension can reduce numerical algorithms computational time.
 - Many statistical models suffer from high correlation between covariates
- PCA is not scale invariant,
 - standardize all the p variables before applying PCA.

PCA - Overview

- Principal component analysis (PCA) is a dimensionality reduction technique
 - Our ability to visualize data is limited to 2 or 3 dimensions.
 - Lower dimension can reduce numerical algorithms computational time.
 - Many statistical models suffer from high correlation between covariates
- PCA is not scale invariant,
 - standardize all the p variables before applying PCA.
- Assume Σ to be the covariance matrix associated with X .
 - The fraction of the original variance kept by the M principal component

$$R^2 = \frac{\sum_{i=1}^M \lambda_i}{\sum_{j=1}^p \lambda_j}, \quad \lambda_i' \text{'s eigenvalues of } \Sigma$$

Recommended exercise 7

How many principal components should we use for the Credit Dataset? Justify?

Principal Components Regression (PCR)

Principal Components Regression (PCR)

- Principal Components Regression involves:
 - Constructing the first M principal components $\mathbf{Z}_1, \dots, \mathbf{Z}_M$
 - Using these components as the predictors in a standard linear regression model

- In other words, we assume that:
 - The directions in which the predictors show the most variation are the directions that are associated with Y .

Principal Components Regression (PCR)

- Principal Components Regression involves:
 - Constructing the first M principal components $\mathbf{Z}_1, \dots, \mathbf{Z}_M$
 - Using these components as the predictors in a standard linear regression model
- Key assumptions: A small number of principal components suffice to explain:
 1. Most of the variability in the data.
 2. The relationship with the response.
- The assumptions above are not guaranteed to hold in every case.
 - This is true specially for assumption 2 above.
 - Since the PCs are selected via unsupervised learning.

- Simulated dataset in which the first five principal components of X contain all the information about the response Y .
- Dashed line: Irreducible error (Simulated data)
- Left: Results for PCR. Right: Results for lasso (solid) and ridge regression (dotted).

Example: PCR vs. Lasso and Ridge (Simulated data)

- PCR performed well on simulated data, recovering the need for $M = 5$
 - However, results are only slightly better than lasso and very similar to Ridge.

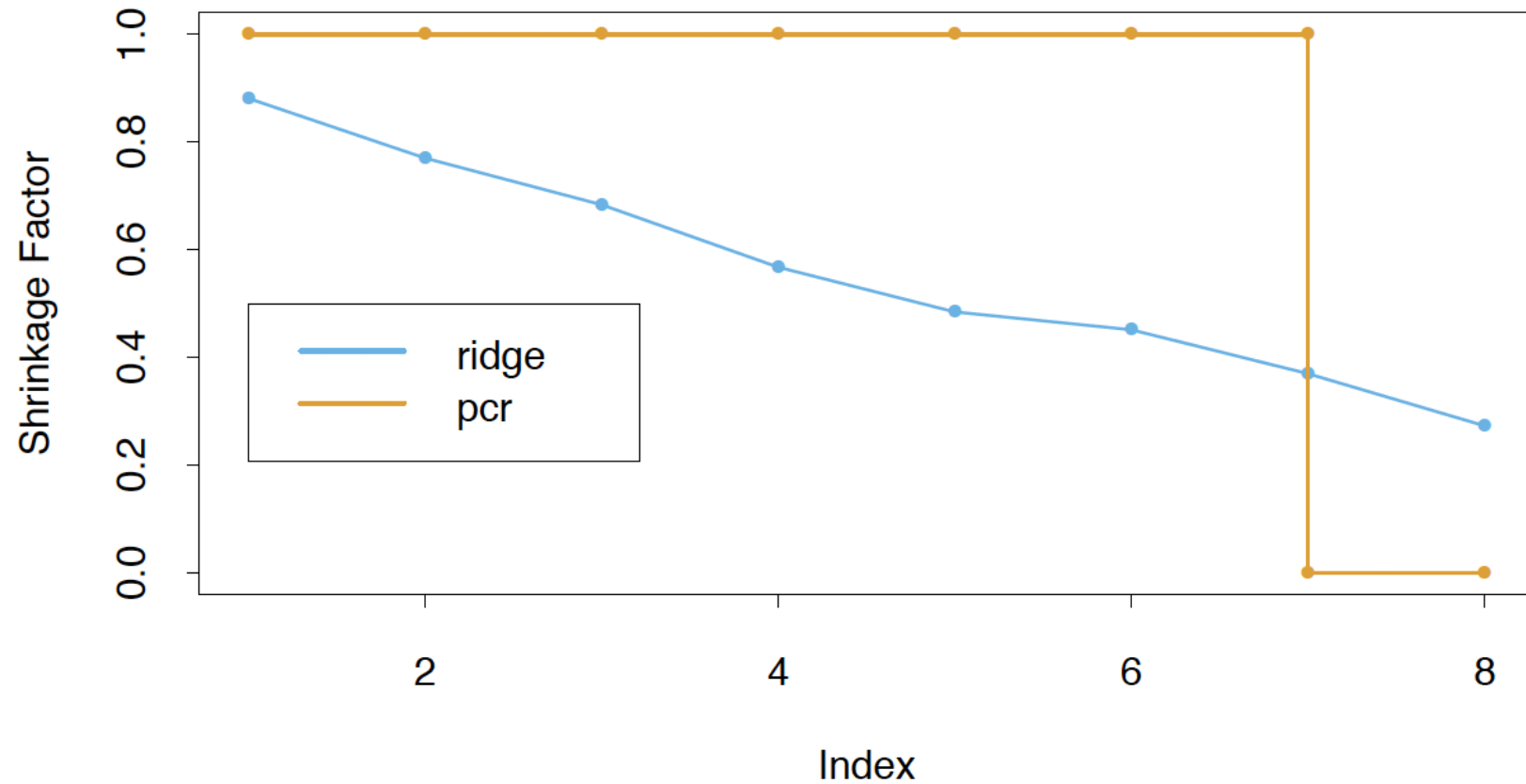
Example: PCR vs. Lasso and Ridge (Simulated data)

- PCR performed well on simulated data, recovering the need for $M = 5$
 - However, results are only slightly better than lasso and very similar to Ridge.
- Similar to Ridge, PCR does not perform feature selection
 - PCs are linear combination of all predictors

Example: PCR vs. Lasso and Ridge (Simulated data)

- PCR performed well on simulated data, recovering the need for $M = 5$
 - However, results are only slightly better than lasso and very similar to Ridge.
- Similar to Ridge, PCR does not perform feature selection
 - PCs are linear combination of all predictors
- PCR can be seen as discretized version of Ridge regression.
 - Ridge shrinks coefs. of the PCs by $\lambda_j^2 / (\lambda_j^2 + \lambda)$
 - Higher pressure on less important PCs
 - PCR discards the $p - M$ smallest eigenvalue components.

Example: Shrinkage Factor



- The lowest cross-validation error occurs when there are $M = 10$ components
 - almost no dimension reduction at all
- Discretized version of Ridge

Recommended exercise 11

Apply PCR on the Credit dataset and compare the results with the methods covered in Lecture 1.

PCR (Drawback)

- Dimensionality reduction is done via an unsupervised method (PCA)
 - No guarantee that the directions that best explain the predictors will also be the best directions to use for predicting the response.

Partial Least Squares (PLS)

Partial Least Squares (PLS)

- PLS works similar to PCR
 - Dimension reduction: $Z_1, \dots, Z_M, M < p$
 - Z_i linear combination of original predictors.
 - Apply standard linear model using Z_1, \dots, Z_M as predictors.

Partial Least Squares (PLS)

- PLS works similar to PCR
 - Dimension reduction: $Z_1, \dots, Z_M, M < p$
 - Z_i linear combination of original predictors.
 - Apply standard linear model using Z_1, \dots, Z_M as predictors.
- But it uses the response Y in order to identify new features
 - attempts to find directions that help explain both the response and the predictors.

Partial Least Squares (Algorithm)

- $Z_1 = \sum_{j=1}^p \phi_{j1} X_j$
 - ϕ_{j1} is the coefficient from the simple linear regression of Y onto X_j .
 - this coefficient is proportional to the correlation between Y and X_j .
 - PLS puts highest weight on the variables that are most strongly related to the response.

Partial Least Squares (Algorithm)

- $Z_1 = \sum_{j=1}^p \phi_{j1} X_j$
 - ϕ_{j1} is the coefficient from the simple linear regression of Y onto X_j .
 - this coefficient is proportional to the correlation between Y and X_j .
 - PLS puts highest weight on the variables that are most strongly related to the response.
- To obtain the second PLS direction, Z_2 :
 - We regress each variable on Z_1 and take the residuals
 - The residuals are remained info not explained by Z_1
 - We the compute Z_2 using this orthogonalized data, similarly to Z_1 .

Partial Least Squares (Algorithm)

- $Z_1 = \sum_{j=1}^p \phi_{j1} X_j$
 - ϕ_{j1} is the coefficient from the simple linear regression of Y onto X_j .
 - this coefficient is proportional to the correlation between Y and X_j .
 - PLS puts highest weight on the variables that are most strongly related to the response.
- To obtain the second PLS direction, Z_2 :
 - We regress each variable on Z_1 and take the residuals
 - The residuals are remained info not explained by Z_1
 - We the compute Z_2 using this orthogonalized data, similarly to Z_1 .
- We can repeat this iteration process M times to get Z_1, \dots, Z_M .

Recommended exercise 12

Apply PLS on the Credit dataset and compare the results with the methods covered in Lecture 1 and PCR.

Partial Least Squares (Performance)

- In practice, PLS often performs no better than ridge regression or PCR.
 - Supervised dimension reduction of PLS can reduce bias.
 - It also has the potential to increase variance.

In summary

- PLS, PCR and ridge regression tend to behave similarly.

In summary

- PLS, PCR and ridge regression tend to behave similarly.
- Ridge regression may be preferred because it shrinks smoothly, rather than in discrete steps.

- I would say:
 - If you only concerned with prediction accuracy, either ridge or lasso.
 - If model interpretability is desirable, lasso is preferred.

Considerations in high dimensions

High dimension

- High dimension problems: $n < p$
- More common nowadays

High dimension issues (Example)

- Standard linear regression cannot be applied.
 - Perfect fit to the data, regardless of relationship
 - Unfortunately, the C_p , AIC, and BIC approaches are problematic (hard to estimate σ^2)

- Simulated example with $n = 20$ training observations.
- features that are completely unrelated to the outcome are added to the model.

- The lasso was performed with $n = 100$ observations and three values of p , the number of features.
- Of the p features, 20 were associated with the response.
- When $p = 2000$ the lasso performed poorly regardless of the amount of regularization.

The danger of too many features

- In general, adding additional signal features helps (smaller test set errors)

The danger of too many features

- In general, adding additional signal features helps (smaller test set errors)
- However, adding noise features that are not truly associated with the response increases test set error.
 - Noise features exacerbating the risk of overfitting
 - Previous example shows that regularizations does not eliminate the problem

The danger of too many features

- In general, adding additional signal features helps (smaller test set errors)
- However, adding noise features that are not truly associated with the response increases test set error.
 - Noise features exacerbating the risk of overfitting
 - Previous example shows that regularizations does not eliminate the problem
- New technologies that allow for the collection of measurements for thousands or millions of features are a double-edged sword

- Multicollinearity: any variable in the model can be written as a linear combination of all of the other variables in the model.

Interpreting results in high dimension

- In the high-dimensional setting, the multicollinearity problem is extreme
- Essentially, this means:
 - We can never know exactly which variables (if any) truly are predictive of the outcome.
 - We can never identify the best coefficients for use in the regression.

Interpreting results in high dimension

- In the high-dimensional setting, the multicollinearity problem is extreme
- Essentially, this means:
 - We can never know exactly which variables (if any) truly are predictive of the outcome.
 - We can never identify the best coefficients for use in the regression.
 - At most, we can hope to assign large regression coefficients to variables that are correlated with the variables that truly are predictive of the outcome.
 - We will find one of possibly many suitable predictive models.

The end

Thank you for showing up