

# Chapter 6: Linear Model Selection and Regularization

*Thiago G. Martins / NTNU & AIA Science*

*Spring 2018*

## Recap

### Statistical Learning

- We have the predictors  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$  associated with the responses  $y_1, y_2, \dots, y_n$ .

### Statistical Learning

- We have the predictors  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$  associated with the responses  $y_1, y_2, \dots, y_n$ .
- We assume there is a function  $f$  such that

$$y_i = f(\mathbf{x}_i) + \epsilon$$

where the noise  $\epsilon$  has zero mean and variance  $\sigma^2$

- $\epsilon$  represents the irreducible error that cannot be explained by  $f$ .

### Statistical Learning

- We have the predictors  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$  associated with the responses  $y_1, y_2, \dots, y_n$ .
- We assume there is a function  $f$  such that

$$y_i = f(\mathbf{x}_i) + \epsilon$$

where the noise  $\epsilon$  has zero mean and variance  $\sigma^2$

- We want to find a model  $\hat{f}$ , such that  $(y_i - \hat{y}_i)$  is minimal for both in-sample and out-sample  $\mathbf{x}_i$ , where  $\hat{y}_i = \hat{f}(\mathbf{x}_i)$
- This is a general description, as we have not imposed assumptions on  $f$ .

### Bias-variance trade-off

- Independent of the model we choose,  $\hat{f}$ , we can decompose its expected error on an out-sample  $\mathbf{x}_i$  as follows:

$$E[(y_i - \hat{y}_i)^2] = \text{Bias}[\hat{y}_i]^2 + \text{Var}[\hat{y}_i] + \sigma^2$$

where

$$\text{Bias}[\hat{y}_i] = E[\hat{y}_i - y_i]$$

and

$$\text{Var}[\hat{y}_i] = E[\hat{y}_i^2] - E[\hat{y}_i]^2$$

- The expectation ranges over different choices of the training set all sampled from the same joint distribution  $P(x, y)$ .
- A very complex model tend to have low bias but high variance

- Simpler models will tend to have higher bias but lower variance
- Principle of parsimony: We want to have the simplest possible method to obtain low bias while minimizing the variance.

## Recommended exercise 1

Show that

$$E[(y_i - \hat{y}_i)^2] = \text{Bias}[\hat{y}_i]^2 + \text{Var}[\hat{y}_i] + \sigma^2$$

assuming the notation used in the previous two slides.

- This was already properly covered in previous modules, but it is a good exercise.

## RSS, $R^2$ , MSE

- Residual Sum of Squares (RSS)

$$RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

## RSS, $R^2$ , MSE

- Residual Sum of Squares (RSS)

$$RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

- Mean Square Error (MSE)

$$MSE = \frac{RSS}{n} = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}$$

## RSS, $R^2$ , MSE

- Residual Sum of Squares (RSS)

$$RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

- Mean Square Error (MSE)

$$MSE = \frac{RSS}{n} = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}$$

- Coefficient of determination,  $R^2$

$$R^2 = \frac{TSS - RSS}{TSS} = 1 - \frac{RSS}{TSS} = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y}_i)^2}$$

- The  $R^2$  statistic is the fraction of variance explained by the model.
- The TSS is proportional to the estimated variance of the response variable.
- We want  $R^2$  to be as high as possible.

## Training and test error

- Training error: Average loss over the training sample

$$\text{MSE}_{\text{train}} = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}$$

## Training and test error

- Training error: Average loss over the training sample

$$\text{MSE}_{\text{train}} = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}$$

- Test error: training sample  $\mathcal{T}$ , expectation wrt  $P(\mathbf{X}, Y)$

$$\text{MSE}_{\text{test}} = E[(Y - \hat{f}(\mathbf{X}))^2 | \mathcal{T}]$$

## Training and test error

- Training error: Average loss over the training sample

$$\text{MSE}_{\text{train}} = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}$$

- Test error: training sample  $\mathcal{T}$ , expectation wrt  $P(\mathbf{X}, Y)$

$$\text{MSE}_{\text{test}} = E[(Y - \hat{f}(\mathbf{X}))^2 | \mathcal{T}]$$

- Training error is not a good estimate of the test error.
  - over-fitting leads to poor generalization
- The generalization performance of a learning method relates to its prediction capability on independent test data.

## Model Selection and Model Assessment

- Two different goals
  - Model selection
  - Model assessment
- Model selection: estimating the performance of different models in order to choose the best one.
- Model assessment: having chosen a final model, estimating its prediction error (generalization error) on new data.

## Model Selection and Model Assessment

- Two different goals
  - Model selection
  - Model assessment
- Data-rich situation: a training set, a validation set, and a test set.
  - The validation error will likely underestimate the test error.
  - For this reason, it is important to keep the test set in a “vault”.
  - No golden rule about the % in each set.

## Penalties and CV

- If there is insufficient data to split it into three parts:
- We have at least two options.

## Penalties and CV

- If there is insufficient data to split it into three parts:
  - Indirectly estimate test error by making an adjustment to the training error to account for the bias due to overfitting.

## Penalties and CV

- If there is insufficient data to split it into three parts:
  - Indirectly estimate test error by making an adjustment to the training error to account for the bias due to overfitting.
  - Directly estimate the test error, using cross-validation approach.

## Standard Linear Models

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \epsilon$$

or in matrix form:

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

- So far, we have not made assumptions about  $f$ .
  - But from now on we assume  $f$  to be linear on the coeffs.
- Assumptions:
  - $(x_i^T, y_i)$  is independent from  $(x_j^T, y_j)$ ,  $\forall i \neq j$ .
  - The design matrix has full rank,  $\text{rank}(\mathbf{X}) = p + 1$ ,  $n \gg (p + 1)$
  - $\boldsymbol{\epsilon} \sim \mathcal{N}_n(\mathbf{0}, \sigma^2 \mathbf{I})$

## Standard Linear Models

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \epsilon$$

or in matrix form:

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

- Least Square Fitting: Minimize the RSS

$$RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 = (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}})$$

## Standard Linear Models

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \epsilon$$

or in matrix form:

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

- Least Square Fitting: Minimize the RSS

$$RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 = (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}})$$

- Least squares and maximum likelihood estimator for  $\boldsymbol{\beta}$ :

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

## Recommended exercise 2

1. Show that the least square estimator of a standard linear model is given by

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

2. Show that the maximum likelihood estimator is equal to the least square estimator for the standard linear model.

## Penalties to training error

- Indirectly estimate test error by penalizing the training error
  - $C_p$ , AIC, BIC and Adjusted  $R^2$
  - All of the options above add a penalty to the training error that increase with the number of predictors in the model
  - All have rigorous theoretical justifications that are beyond the scope of this course.

## Penalties to training error ( $C_p$ and AIC)

- $C_p$

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

- $C_p$  is an unbiased estimator of the test MSE, if  $\hat{\sigma}^2$  is an unbiased estimator of  $\sigma^2$
- lower  $C_p$  is better
- $d$  is the number of predictors in the model.

## Penalties to training error ( $C_p$ and AIC)

- $C_p$

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

- $C_p$  is an unbiased estimator of the test MSE, if  $\hat{\sigma}^2$  is an unbiased estimator of  $\sigma^2$
- lower  $C_p$  is better

- AIC

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

- AIC is proportional to  $C_p$  in the case of the standard linear model.

## Penalties to training error (BIC and Adjusted $R^2$ )

- BIC

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

- Derived from a Bayesian point of view.
- the lower the better

- BIC

- Since  $\log(n) > 2$  for any  $n > 7$ , the BIC statistic generally places a heavier penalty on models with many variables, and hence results in the selection of smaller models than  $C_p$ .

## Penalties to training error (BIC and Adjusted $R^2$ )

- BIC

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

- Derived from a Bayesian point of view.
- the lower the better

- Adjusted  $R^2$

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

- the higher the better

- Adjusted  $R^2$

- The intuition behind the adjusted  $R^2$  is that once all of the correct variables have been included in the model, adding additional noise variables will lead to only a very small decrease in RSS.

## Cross-validation

- Directly estimate the test error, using cross-validation approach.

## Cross-validation

- Directly estimate the test error, using cross-validation approach.
  - More computationally expensive than the quantities above

## Cross-validation

- Directly estimate the test error, using cross-validation approach.
  - More computationally expensive than the quantities above
  - Makes fewer assumptions about the true underlying model.
- For AIC, for example, we assume that the approximation model  $f(x)$  is correct, in the sense that exist a  $\hat{\beta}$  that recover the true model of the data.

## Cross-validation

- Directly estimate the test error, using cross-validation approach.
  - More computationally expensive than the quantities above
  - Makes fewer assumptions about the true underlying model.
  - Can be used in cases where it is hard to pinpoint the model degrees of freedom
- For standard linear models, the number of effective parameters is the number of predictors, but this does not hold for more complex models, such as dynamic models used in signal processing.

## Credit Dataset

- Balance: average credit card debt

## Recommended exercise 3

Write R code to create a similar representation of the Credit data figure of the previous slide. That is, try to recreate a similar plot in R.

## Introduction

### Objective of the module

Improve linear models **prediction accuracy** and/or **model interpretability** by replacing least square fitting with some alternative fitting procedures.

### Prediction accuracy ...

... when using standard linear models

Assuming true relationship is approx. linear: **low bias**.

- $n \gg p$ : **low variance**
- $n$  not much larger than  $p$ : **high variance**
- $n < p$ : multiple solutions available, **infinite variance**, model cannot be used.

### Prediction accuracy ...

... when using standard linear models

Assuming true relationship is approx. linear: **low bias**.

- $n \gg p$ : **low variance**
- $n$  not much larger than  $p$ : **high variance**
- $n < p$ : multiple solutions available, **infinite variance**, model cannot be used.

By constraining or shrinking the estimated coefficients:

- often substantially reduce the variance at the cost of a negligible increase in bias.
- better generalization for out of sample prediction

## Recommended exercise 4

Show why fitting a standard linear regression model when  $n < p$  is not an option.

## Model Interpretability

- Some or many of the variables might be irrelevant wrt the response variable
- Some of the discussed approaches lead to automatically performing feature/variable selection.

## Outline

We will cover the following alternatives to using least squares to fit linear models

- **Subset Selection:** Identifying a subset of the  $p$  predictors that we believe to be related to the response.

## Outline

We will cover the following alternatives to using least squares to fit linear models

- **Subset Selection:** Identifying a subset of the  $p$  predictors that we believe to be related to the response.
- **Shrinkage:** fitting a model involving all  $p$  predictors with the estimated coefficients shrunken towards zero relative to the least squares estimates.

## Outline

We will cover the following alternatives to using least squares to fit linear models

- **Subset Selection:** Identifying a subset of the  $p$  predictors that we believe to be related to the response.
- **Shrinkage:** fitting a model involving all  $p$  predictors with the estimated coefficients shrunken towards zero relative to the least squares estimates.
- **Dimension Reduction:** This approach involves projecting the  $p$  predictors into a  $M$ -dimensional subspace, where  $M < p$ .

## Subset Selection

### Subset Selection

Identifying a subset of the  $p$  predictors that we believe to be related to the response.



## Subset Selection

Identifying a subset of the  $p$  predictors that we believe to be related to the response.

Outline:

- Best subset selection
- Stepwise model selection

## Best Subset Selection

1. Fit a least square regression for each possible combination of the  $p$  predictors.
2. Look at all the resulting models and pick the best.

## Best Subset Selection

1. Fit a least square regression for each possible combination of the  $p$  predictors.
2. Look at all the resulting models and pick the best.

Number of models considered:

$$\binom{p}{1} + \binom{p}{2} + \dots + \binom{p}{p} = 2^p$$

## Best Subset Selection (Algorithm)

- Step 2 identifies the best model (on the training data) for each subset size
  - Reduces the problem from  $2^p$  to  $p + 1$  models to select from
- Step 3 choose among the  $p + 1$  using the test error
  - Otherwise we would always choose the model with all parameters

## Best Subset Selection (Credit Data Example)

- The red frontier tracks the best model for a given number of predictors, according to RSS and  $R^2$ .

## Recommended exercise 5

1. For the Credit Dataset, pick the best model using Best Subset Selection according to  $C_p$ ,  $BIC$  and Adjusted  $R^2$ 
  - Hint: Use the `regsubsets()` of the `leaps` library, similar to what was done in Lab 1 of the book.
2. For the Credit Dataset, pick the best model using Best Subset Selection according to a 10-fold CV
  - Hint: Use the output obtained in the previous step and build your own CV function to pick the best model.
3. Compare the result obtained in Step 1 and Step 2.

## Best Subset Selection (Drawbacks)

- Does not scale well -> the number of models to consider explode as  $p$  increases
  - $p = 10$  leads to approx. 1000 possibilities
  - $p = 20$  leads to over 1 million possibilities

## Best Subset Selection (Drawbacks)

- Does not scale well -> the number of models to consider explode as  $p$  increases
  - $p = 10$  leads to approx. 1000 possibilities
  - $p = 20$  leads to over 1 million possibilities
- Large search space might lead to overfitting on training data

## Stepwise selection

Add and/or remove one predictor at a time.

## Stepwise selection

Add and/or remove one predictor at a time.

Methods outline:

- Forward Stepwise Selection
- Backward Stepwise Selection
- Hybrid approaches

## Forward Stepwise Selection

- Starts with a model containing no predictors,  $\mathcal{M}_0$
- Adds predictors to the model, one at a time, until all of the predictors are in the model
  - $\mathcal{M}_1, \mathcal{M}_2, \dots, \mathcal{M}_p$
- Select the best model among  $\mathcal{M}_0, \mathcal{M}_1, \dots, \mathcal{M}_p$

## Forward Stepwise Selection (Algorithm)

### Forward Stepwise Selection (About the Algorithm)

- Goes from fitting  $2^p$  models to  $1 + \sum_{k=0}^{p-1} (p - k) = 1 + p(p + 1)/2$  models

### Forward Stepwise Selection (About the Algorithm)

- Goes from fitting  $2^p$  models to  $1 + \sum_{k=0}^{p-1} (p - k) = 1 + p(p + 1)/2$  models
- It is a guided search, we don't choose  $1 + p(p + 1)/2$  models to consider at random.

### Forward Stepwise Selection (About the Algorithm)

- Goes from fitting  $2^p$  models to  $1 + \sum_{k=0}^{p-1} (p - k) = 1 + p(p + 1)/2$  models
- It is a guided search, we don't choose  $1 + p(p + 1)/2$  models to consider at random.
- Not guaranteed to yield the best model containing a subset of the  $p$  predictors.

## Forward Stepwise Selection (About the Algorithm)

- Goes from fitting  $2^p$  models to  $1 + \sum_{k=0}^{p-1} (p-k) = 1 + p(p+1)/2$  models
- It is a guided search, we don't choose  $1 + p(p+1)/2$  models to consider at random.
- Not guaranteed to yield the best model containing a subset of the  $p$  predictors.
- Forward stepwise selection can be applied even in the high-dimensional setting where  $n < p$ 
  - By limiting the algorithm to submodels  $\mathcal{M}_0, \dots, \mathcal{M}_{n-1}$  only

## Forward Stepwise Selection (Credit Data Example)

- The first three models are identical but the fourth models differ.

## Backward Stepwise Selection

- Starts with a model containing all predictors,  $\mathcal{M}_p$ .
- Iteratively removes the least useful predictor, one-at-a-time, until all the predictors have been removed.
  - $\mathcal{M}_{p-1}, \mathcal{M}_{p-2}, \dots, \mathcal{M}_0$
- Select the best model among  $\mathcal{M}_0, \mathcal{M}_1, \dots, \mathcal{M}_p$

## Backward Stepwise Selection (Algorithm)

## Backward Stepwise Selection (About the Algorithm)

- Similar properties to the Forward algorithm
  - Search  $1 + p(p+1)/2$  models instead of  $2^p$  models
  - It is a guided search, we don't choose  $1 + p(p+1)/2$  models to consider at random.
  - Not guaranteed to yield the best model containing a subset of the  $p$  predictors.

## Backward Stepwise Selection (About the Algorithm)

- Similar properties to the Forward algorithm
  - Search  $1 + p(p+1)/2$  models instead of  $2^p$  models
  - It is a guided search, we don't choose  $1 + p(p+1)/2$  models to consider at random.
  - Not guaranteed to yield the best model containing a subset of the  $p$  predictors.
- However, Backward selection requires that the number of samples  $n$  is larger than the number of variables  $p$ 
  - So that the full model can be fit.

## Hybrid Approach

- Similarly to forward selection, variables are added to the model sequentially.
- However, after adding each new variable, the method may also remove any variables that no longer provide an improvement in the model fit.
- Better model space exploration while retaining computational advantages of stepwise selection.

## Recommended exercise 6

1. Select the best model for the Credit Data using Forward, Backward and Hybrid (sequential replacement) Stepwise Selection.
  - Hint: Use the `regsubsets()` of the `leaps` library
2. Compare with the results obtained with Best Subset Selection.

## Shrinkage Methods

### Shrinkage Methods

- 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.

### Shrinkage Methods

- 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.
- Reduce the number of effective parameters
  - While retaining the ability to capture the most interesting aspects of the problem.

### Shrinkage Methods

- 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.
- Reduce the number of effective parameters
  - While retaining the ability to capture the most interesting aspects of the problem.
- The two best-known techniques for shrinking the regression coefficients towards zero are:
  - the ridge regression.
  - the lasso.

## Ridge regression

The ridge regression coeffs  $\beta^R$  are the ones that minimize

$$RSS + \lambda \sum_{j=1}^p \beta_j^2$$

with  $\lambda > 0$  being a tuning parameter.

## Ridge regression

The ridge regression coeffs  $\beta^R$  are the ones that minimize

$$RSS + \lambda \sum_{j=1}^p \beta_j^2$$

with  $\lambda > 0$  being a tuning parameter.

- Note that the penalty is not applied to the intercept,  $\beta_0$ .
  - If we included the intercept,  $\beta^R$  would depend on the average of the response.
  - We want to shrink the estimated association of each feature with the response.
- Intercept is the mean value of the response when the covariates are set to zero

## Ridge regression

- Ridge regression are not scale-invariant
  - The standard least square are scale-invariant
- multiplying  $X_j$  by a constant  $c$  simply leads to a scaling of the least squares coefficient estimates by a factor of  $1/c$ .

## Ridge regression

- Ridge regression are not scale-invariant
  - The standard least square are scale-invariant
  - $\beta^R$  will not only depend on  $\lambda$  but also on the scaling of the  $j$ th predictor

## Ridge regression

- Ridge regression are not scale-invariant
  - The standard least square are scale-invariant
  - $\beta^R$  will not only depend on  $\lambda$  but also on the scaling of the  $j$ th predictor
  - Apply Ridge regression after standardizing the predictors

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

## Ridge regression (Credit Data Example)

The standardized ridge regression coefficients are displayed for the Credit data set.

## Ridge regression (Effectiveness)

- Why does it work?
  - As  $\lambda$  increase, the flexibility of the fit decreases.
  - Leading to a decrease variance but increased bias

## Ridge regression (Effectiveness)

- Why does it work?
  - As  $\lambda$  increase, the flexibility of the fit decreases.
  - Leading to a decrease variance but increased bias
- MSE is a function of the variance and the squared bias
  - Need to find sweet spot (see next Fig.)

## Ridge regression (Effectiveness)

- Why does it work?
  - As  $\lambda$  increase, the flexibility of the fit decreases.
  - Leading to a decrease variance but increased bias
- MSE is a function of the variance and the squared bias
  - Need to find sweet spot (see next Fig.)
- Therefore, ridge regression works best for the cases where
  - The relationship between covariates and response is close to linear (low bias)
  - And the least square estimates have high variance (high  $p$  in relation to  $n$ )

## Ridge regression (MSE)

- Squared bias (black), variance (green), and test mean squared error (purple) for the ridge regression predictions on a simulated data set.
- The horizontal dashed lines indicate the minimum possible MSE. The purple crosses indicate the ridge regression models for which the MSE is smallest.

## Ridge regression (Computationally efficient)

- The computations required to solve  $\beta_\lambda^R$ , simultaneously for all values of  $\lambda$ , are almost identical to those for fitting a model using least squares.
  - See (Friedman, Hastie, and Tibshirani 2010) and the references therein.

## Ridge regression (Disadvantages)

- Unlike previous methods, ridge regression will include all  $p$  predictors in the final model.
  - The penalty  $\lambda$  will shrink all of the coefficients towards zero.
  - But it will not set any of them exactly to zero (unless  $\lambda = \infty$ ).

## Ridge regression (Disadvantages)

- Unlike previous methods, ridge regression will include all  $p$  predictors in the final model.
  - The penalty  $\lambda$  will shrink all of the coefficients towards zero.
  - But it will not set any of them exactly to zero (unless  $\lambda = \infty$ ).
- This may not be a problem for prediction accuracy, but makes model interpretation hard for large  $p$ .

## Recommended exercise 7

1. Apply Ridge regression to the Credit Dataset.
2. Compare the results with the standard linear regression.

## Lasso

- The Lasso regression coeffs  $\beta^L$  are the ones that minimize

$$RSS + \lambda \sum_{j=1}^p |\beta_j|$$

with  $\lambda > 0$  being a tuning parameter.

## Lasso

- The Lasso regression coeffs  $\beta^L$  are the ones that minimize

$$RSS + \lambda \sum_{j=1}^p |\beta_j|$$

with  $\lambda > 0$  being a tuning parameter.

- Lasso also shrinks the coefficients towards zero
- In addition, the  $\uparrow_1$  penalty has the effect of forcing some of the coefficients to be exactly zero when  $\lambda$  is large enough

## Lasso

- The Lasso regression coeffs  $\beta^L$  are the ones that minimize

$$RSS + \lambda \sum_{j=1}^p |\beta_j|$$

with  $\lambda > 0$  being a tuning parameter.

- Lasso also shrinks the coefficients towards zero
- In addition, the  $\uparrow_1$  penalty has the effect of forcing some of the coefficients to be exactly zero when  $\lambda$  is large enough
- A geometric explanation will be presented in a future slide.

## Lasso regression (Credit Data Example)

The standardized lasso coefficients are displayed for the Credit data set.

## Lasso regression (Simulated Data Example)

- $p = 45$ ,  $n = 50$  and 2 out of 45 predictors related to the response.
- grey lines represent unrelated predictors
- Minimum CV error points to only the two real predictors have coeffs != zero
- least-square estimate assign high value to one of the two predictors
  - Many unrelated predictors have non-zero values

## Ridge and Lasso (Different formulations)

- Lasso

$$\text{minimize}_{\beta} \left\{ \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \right\} \text{ subject to } \sum_{j=1}^p |\beta_j| \leq s$$

- Ridge

$$\text{minimize}_{\beta} \left\{ \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \right\} \text{ subject to } \sum_{j=1}^p \beta_j^2 \leq s$$

## Ridge and Lasso (Geometric intuition)

- The red ellipses are the contours of the RSS.
- The solid blue areas are the constraint regions,  $|\beta_1| + |\beta_2| \leq s$  and  $\beta_1^2 + \beta_2^2 \leq s$
- the explanation holds for  $p > 2$ , just harder to visualize

## Comparison between Ridge and Lasso

- Neither is universally better than the other
- One expects lasso to perform better for cases where a relatively small number of predictors have coeffs that are very small or zero

## Comparison between Ridge and Lasso

- Neither is universally better than the other
- One expects lasso to perform better for cases where a relatively small number of predictors have coeffs that are very small or zero
- One expects ridge to be better when the response is a function of many predictors, all with roughly equal size



## Comparison between Ridge and Lasso

- Neither is universally better than the other
- One expects lasso to perform better for cases where a relatively small number of predictors have coeffs that are very small or zero
- One expects ridge to be better when the response is a function of many predictors, all with roughly equal size
- Hard to know a priori, techniques such as CV required

## Recommended exercise 8

1. Apply Lasso regression to the Credit Dataset.
2. Compare the results with the standard linear regression and the Ridge regression.

## Bayesian interpretation

- Gaussian prior with zero mean and std. dev. as function of lambda
  - posterior mode is the ridge regression solution
- Laplace prior with zero mean and scale parameter as a function of lambda
  - posterior mode is the lasso solution

Left: Ridge regression is the posterior mode for  $\beta$  under a Gaussian prior. Right: The lasso is the posterior mode for  $\beta$  under a double-exponential prior.

## Selecting $\lambda$

- Pick  $\lambda$  for which the cross-validation error is smallest.
- re-fit using all of the available observations and the selected value of  $\lambda$ .

## Recommended exercise 9

Through out the recommended exercises, you have applied the following techniques to the Credit Dataset:

1. Standard linear regression
2. Best subset selection
3. Stepwise selection (forward, backward and hybrid)
4. Ridge regression
5. Lasso regression

Which method worked best for this particular dataset? Elaborate.

## References

Friedman, Jerome, Trevor Hastie, and Rob Tibshirani. 2010. "Regularization Paths for Generalized Linear Models via Coordinate Descent." *Journal of Statistical Software* 33 (1). NIH Public Access: 1.