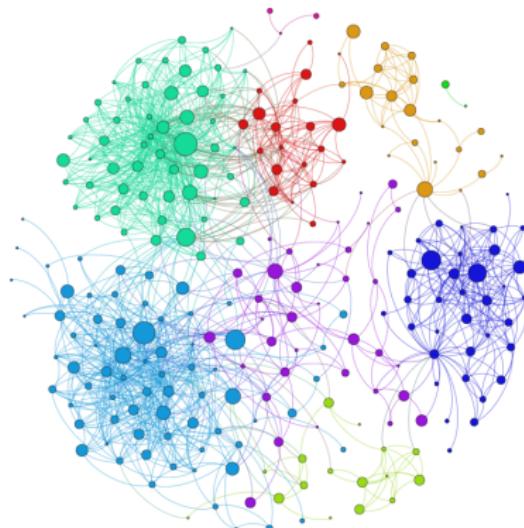


# Model Selection and Regularization

Jiaming Mao

Xiamen University



Copyright © 2017–2019, by Jiaming Mao

This version: Spring 2019

Contact: [jmao@xmu.edu.cn](mailto:jmao@xmu.edu.cn)

Course homepage: [jiamingmao.github.io/data-analysis](https://jiamingmao.github.io/data-analysis)



All materials are licensed under the [Creative Commons Attribution-NonCommercial 4.0 International License](#).

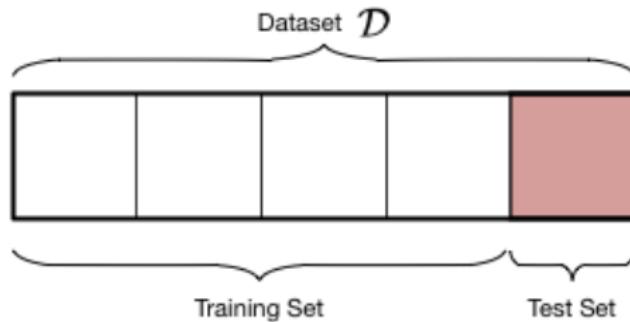
# Model Assessment and Model Selection

- **Model assessment:** evaluating a model's performance
- **Model selection:** choosing the model with the best performance
- Both model assessment and model selection require estimate of a model's **out-of-sample error**.

# Model Assessment

To evaluate the performance of a model  $\mathcal{H}$ , we can *randomly* split our data into two parts:

- **Training set:** used to fit the model (select  $g$  from  $\mathcal{H}$ )
- **Test set:** used to see how good the fit is (evaluate  $g$ 's performance)



## Training versus Testing

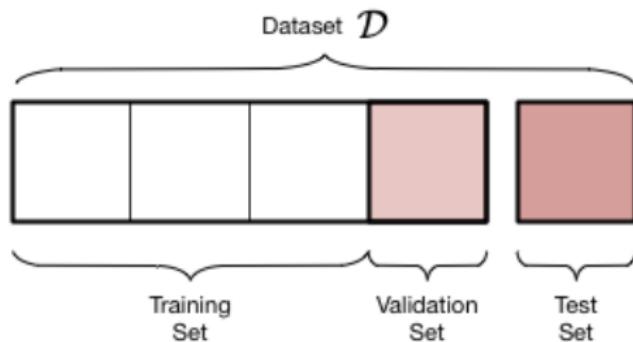
$$\text{Training: } \Pr(|E_{in}(g) - E_{out}(g)| > \epsilon) \leq 4m_{\mathcal{H}}(2N) e^{-\frac{1}{8}\epsilon^2 N} \quad (1)$$

$$\text{Testing: } \Pr(|E_{test}(g) - E_{out}(g)| > \epsilon) \leq 2e^{-2\epsilon^2 N} \quad (2)$$

- The generalization bound for test error  $E_{test}(g)$  is much tighter.
- The test set is not biased, whereas the training set has an **optimistic** bias, since it is used to choose a hypothesis that looks good *on it*.
- The price for a test set is fewer data for training.

# Model Selection

To choose the best model  $\mathcal{H}_*$  among a set of models  $\mathbb{H} = \{\mathcal{H}_0, \mathcal{H}_1, \dots, \mathcal{H}_n\}$ , we can *randomly* split our data into three parts: a training set, a test set, and a **validation (hold-out) set**:



# Model Selection

Then follow the following steps:

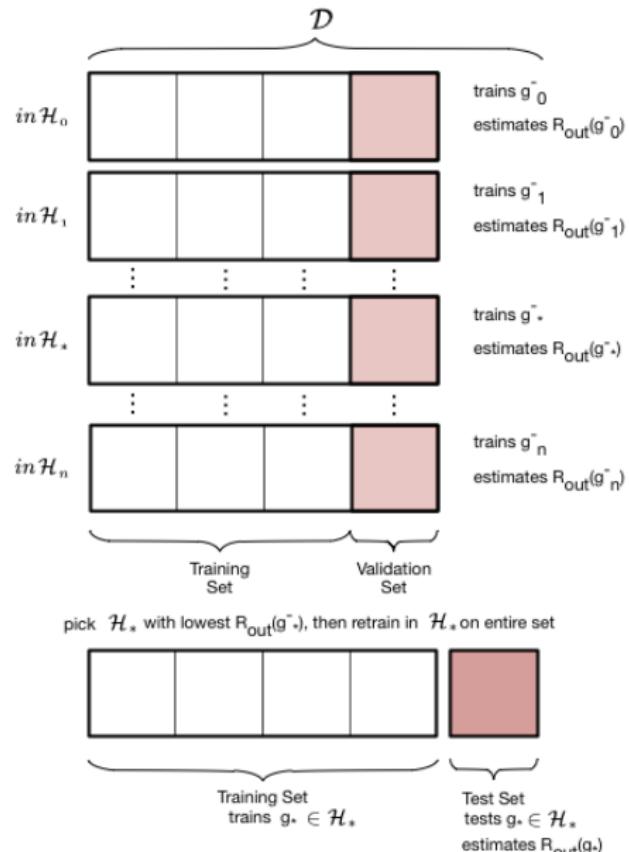
- ① Fit each model  $\mathcal{H}_i$  on the training set and obtain  $g_i$  for  $i = 0, 1, \dots, n$
- ② Use the validation set to evaluate the performance of each  $g_i$  and select the hypothesis  $g_* \in \mathcal{H}_*$  with the smallest validation-set error –  $\mathcal{H}_*$  is our selected model.
- ③ Combine the training set and the validation set. Refit  $\mathcal{H}_*$  on this larger set<sup>1</sup> to obtain  $g_{**}$  – this is our final selected hypothesis.
- ④ Assess the performance of  $g_{**}$  on the test set.

This is called the **validation set approach**.

---

<sup>1</sup>which corresponds to the training set in model assessment, since once we have picked a model – here  $\mathcal{H}_*$  – what is left to do is model assessment.

# Model Selection



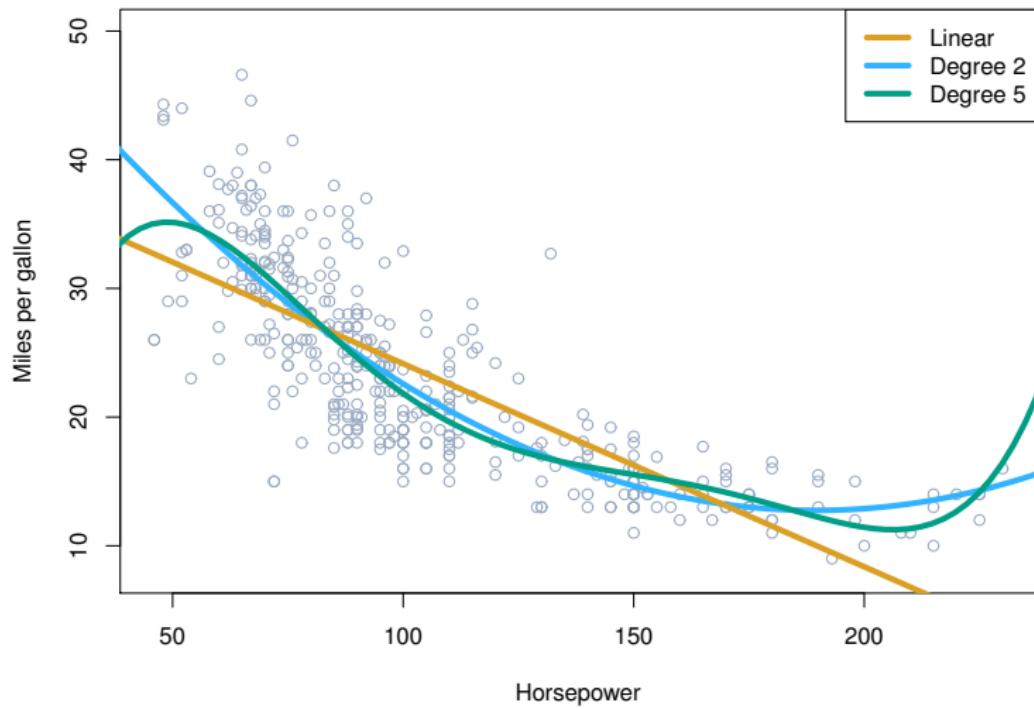
## Model Selection

- A validation set is needed because we cannot use the test set to both help select the best model and assess the performance of the final hypothesis.
- Once a data set has been used in the learning process, it is “contaminated” – it obtains an *optimistic* bias, and the error calculated on the data set no longer has the tight bound in (2).

## Model Selection

- Often times, the different models in  $\mathbb{H}$  can be indexed by a complexity parameter  $\lambda$ , which we call a **hyperparameter**. Choosing the best model amounts to finding the best value of  $\lambda$ .
  - For example, if we are choosing among polynomial models of varying degrees, then  $\lambda$  is the degree of the polynomial.
- Using the validation set approach, the training set is used to fit each model with a given  $\lambda$ , the validation set is the set on which  $\lambda$  itself is fit, while the test set is used to estimate the true out-of-sample performance of the final hypothesis.

# MPG and Horsepower



## MPG and Horsepower

- Want to compare linear vs higher-order polynomial terms in a linear regression of miles per gallon (MPG) on horsepower on a data set of 392 vehicles.
- Suppose there exists an independent test data set somewhere else, so we only need this data set for model selection (without assessment). Then we can randomly split the 392 observations into two sets, a training set containing 196 of the data points, and a validation set containing the remaining 196 observations.

# MPG and Horsepower

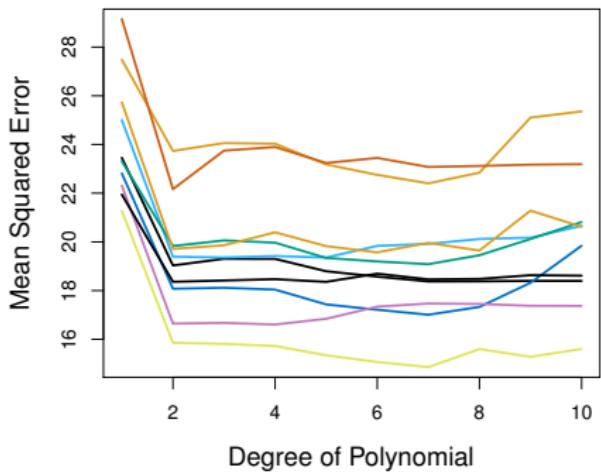
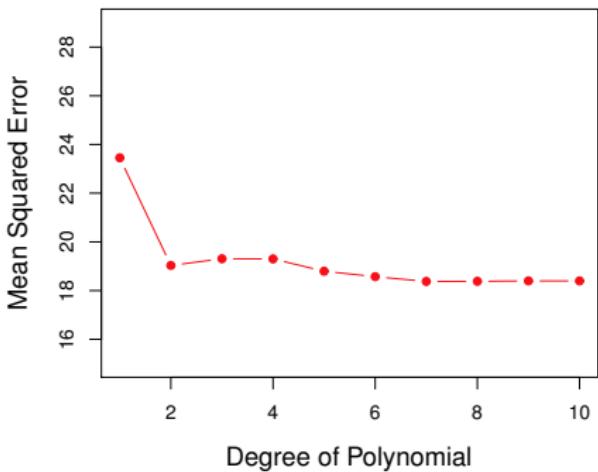
```
require(ISLR) # contains the data set 'Auto'  
attach(Auto)  
train <- sample(392,196) #draw a random subset from the sample  
fit <- lm(mpg ~ horsepower, subset=train)  
yhat <- predict(fit,Auto)  
V.e <- (mpg - yhat)[-train] #validation set error  
V.MSE <- mean(V.e^2) #validation set MSE  
V.MSE  
  
## [1] 23.29559
```

# MPG and Horsepower

```
V.MSE <- rep(0,5)
for (i in 1:5){
  fit <- lm(mpg ~ poly(horsepower,i), subset=train)
  V.e <- (mpg - predict(fit,Auto))[-train]
  V.MSE[i] <- mean(V.e^2)
}
V.MSE

## [1] 23.29559 18.90124 19.25740 20.38538 20.29775
```

# MPG and Horsepower



Left: Validation-set MSE for a single split into training and validation data sets.  
Right: The validation method was repeated ten times, each time using a different random split.

# The Validation Set Approach

## Problems:

- Results can be highly variable, depending on which observations are included in the training set and which are in the validation set.
- The larger the validation set is, the smaller the training set has to be, hence the better the validation set errors are as estimates of true out-of-sample errors, but the poorer the model fits produced by the training set are.
- Conversely, the smaller the validation set, the better the model fits on the training set, but the poorer the validation set errors are as estimates of true out-of-sample errors.

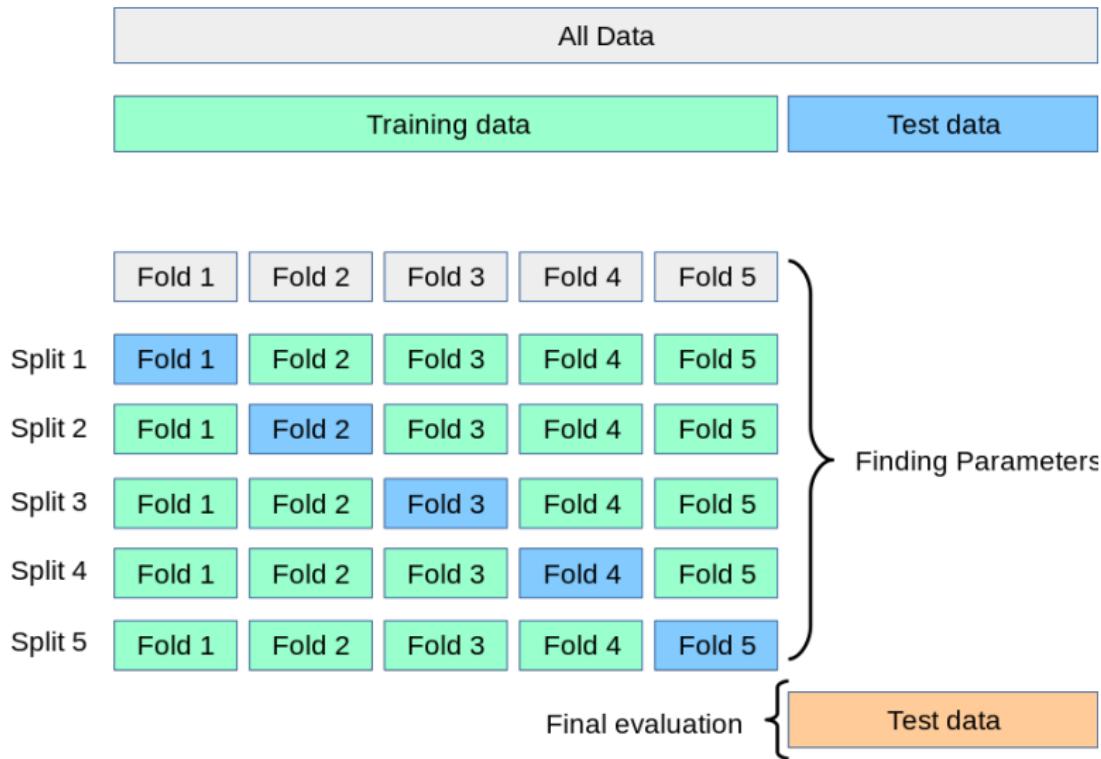
## Cross Validation

- The results of the validation set approach are variable due to the whims of a single random split. Thus, it might be better to randomly split the data multiple times and average the results – this is what **cross-validation** does.
- Like the bootstrap, cross-validation is a **resampling method**: these methods involve repeatedly drawing samples from a training set and refitting a model of interest on each sample in order to obtain additional information about the fitted model.

## Cross Validation

- The cross-validation approach splits the original data into two parts: a training set and a test set, with no separate validation set.
- The training set is then randomly divided into  $K$  equal-sized folds.
- In turn, training is performed on  $K - 1$  folds (combined), and the remainder fold is used for evaluation. The  $K$  evaluation results are then averaged to produce an estimate of a model's out-of-sample error.

# Cross Validation



source

## Cross Validation

The  $K$ -fold cross-validation error of a model  $\mathcal{H}$  is

$$CV_{(K)}(\mathcal{H}) = \frac{1}{K} \sum_{k=1}^K E_{in}^{(k)}\left(g^{(-k)}(\mathcal{H})\right)$$

, where  $g^{(-k)}(\mathcal{H})$  is the hypothesis selected by fitting  $\mathcal{H}$  on the training data with the  $k^{th}$  fold removed, and  $E_{in}^{(k)}(.)$  is the error calculated on the  $k^{th}$  fold data<sup>2</sup>.

---

<sup>2</sup>Note that  $g^{(-k)}(\mathcal{H})$  can be different for different  $k$ , hence the cross-validation error is obtained by averaging over the performance of different hypotheses.

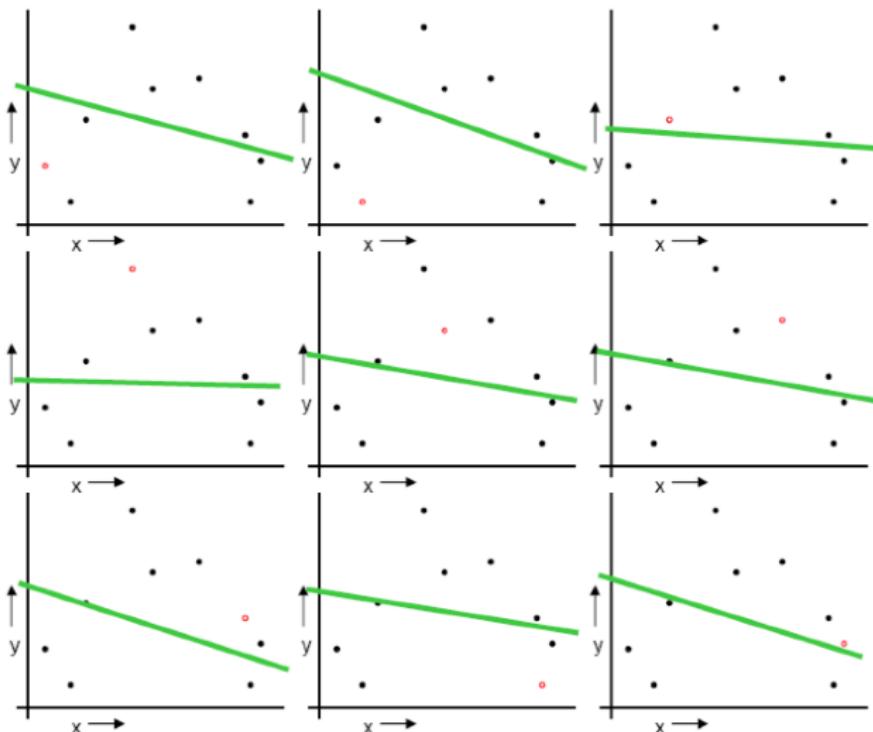
## Cross Validation

Setting  $K = N$  yields the **leave-one out cross-validation (LOOCV)** error:

$$CV_{(N)}(\mathcal{H}) = \frac{1}{N} \sum_{i=1}^N E_{in}^{(i)} \left( g^{(-i)}(\mathcal{H}) \right)$$

, in which case  $g^{(-i)}(\mathcal{H})$  is the result of fitting  $\mathcal{H}$  on the entire training set with only the  $i^{th}$  data point removed, and  $E_{in}^{(i)}(.)$  is the difference between the prediction made by  $g^{(-i)}(\mathcal{H})$  on the  $i^{th}$  data point and its true value, according to the loss function used.

# LOOCV



# MPG and Horsepower

```
## LOOCV
require(boot)
LOOCV.MSE <- rep(0,5)
for (i in 1:5){
  fit <- glm(mpg ~ poly(horsepower,i))
  LOOCV.MSE[i] <- cv.glm(Auto,fit)$delta[1]
}
LOOCV.MSE

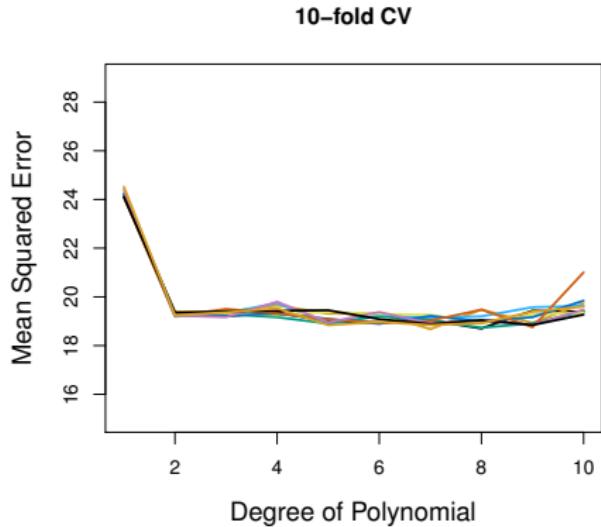
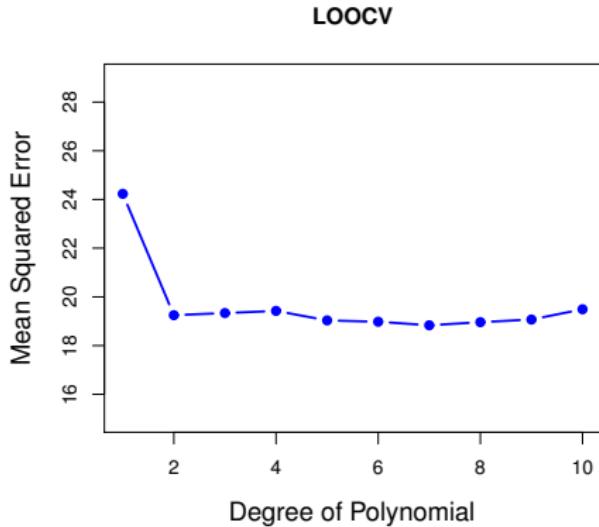
## [1] 24.23151 19.24821 19.33498 19.42443 19.03321
```

# MPG and Horsepower

```
## 10-fold CV
CV.MSE <- rep(0,5)
for (i in 1:5){
  fit <- glm(mpg ~ poly(horsepower,i))
  CV.MSE[i] <- cv.glm(Auto,fit,K=10)$delta[1]
}
CV.MSE

## [1] 24.16778 19.13634 19.15571 19.33862 19.30632
```

# MPG and Horsepower

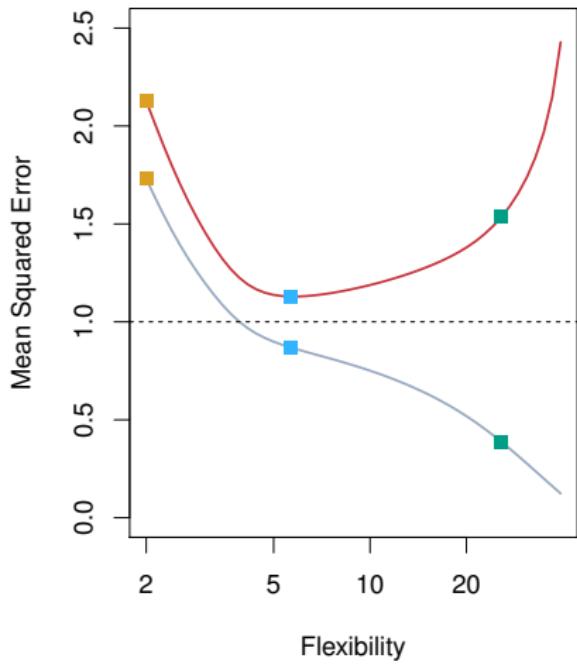
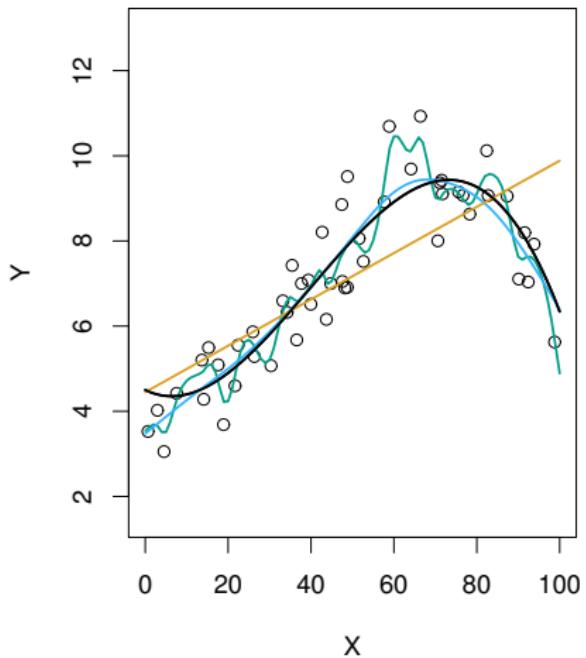


Left: The LOOCV error curve. Right: 10-fold CV was run nine separate times, each with a different random split of the data into ten folds. The figure shows the nine slightly different CV error curves.

## Choice of K

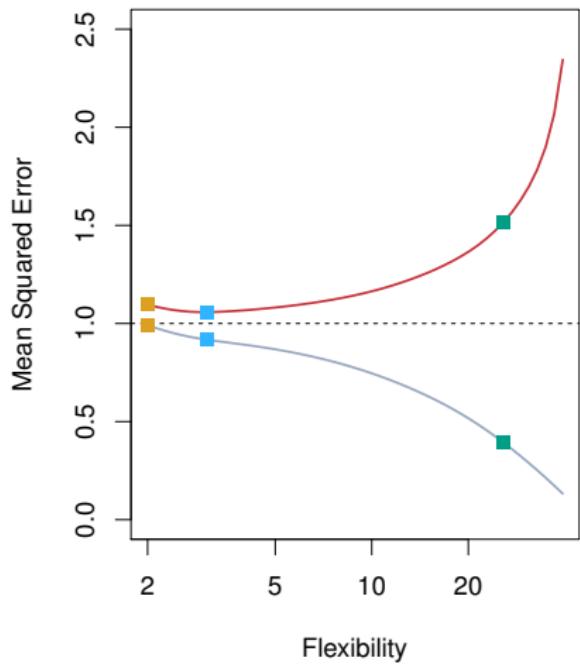
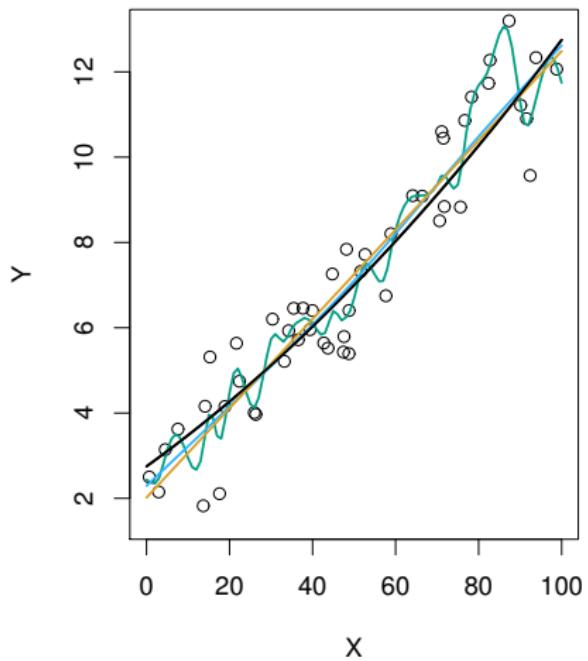
- When  $N$  is large, LOOCV can be computationally expensive: need to fit the model  $N$  times!
- LOOCV gives almost unbiased estimate of the out-of-sample error, since its training set contains almost the entire data set.
- However, the LOOCV estimate has high variance: when we perform LOOCV, we are in effect averaging the outputs of  $N$  fitted models, each of which is trained on an almost identical set of observations. Therefore, these outputs are highly positively correlated with each other.
- There exists a *bias-variance tradeoff*: when  $K$  is small, bias is high and variance is low. When  $K$  is large, bias is low and variance is high.
- Rule of thumb:  $K = 5$  or  $K = 10$ .

# Cross-Validation for Regression



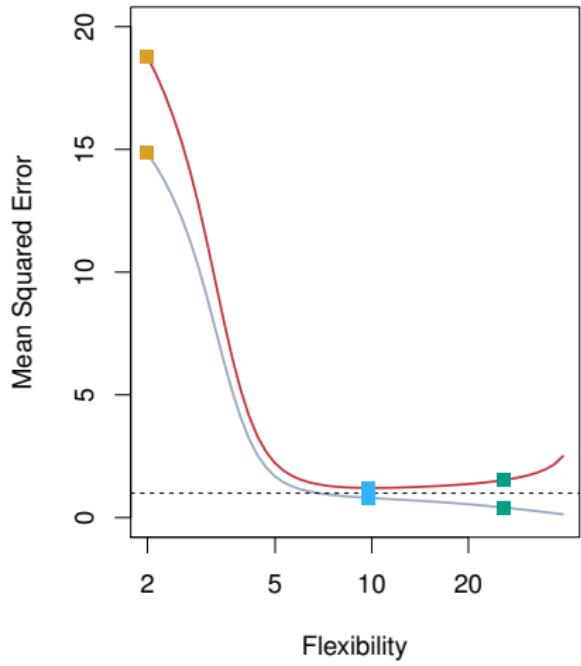
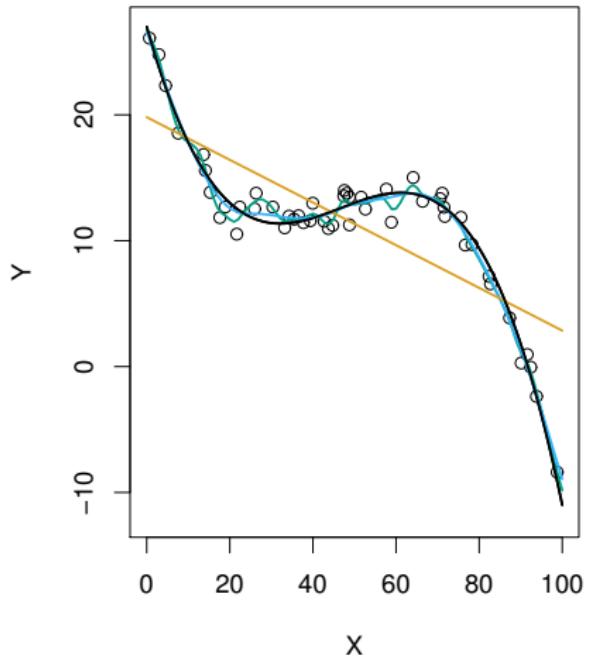
Left: target function (black), linear fit (orange), smoothing spline fits (blue & green). Right: training error (grey), prediction error (red),  $\text{Var}(e)$  (dashed).

# Cross-Validation for Regression



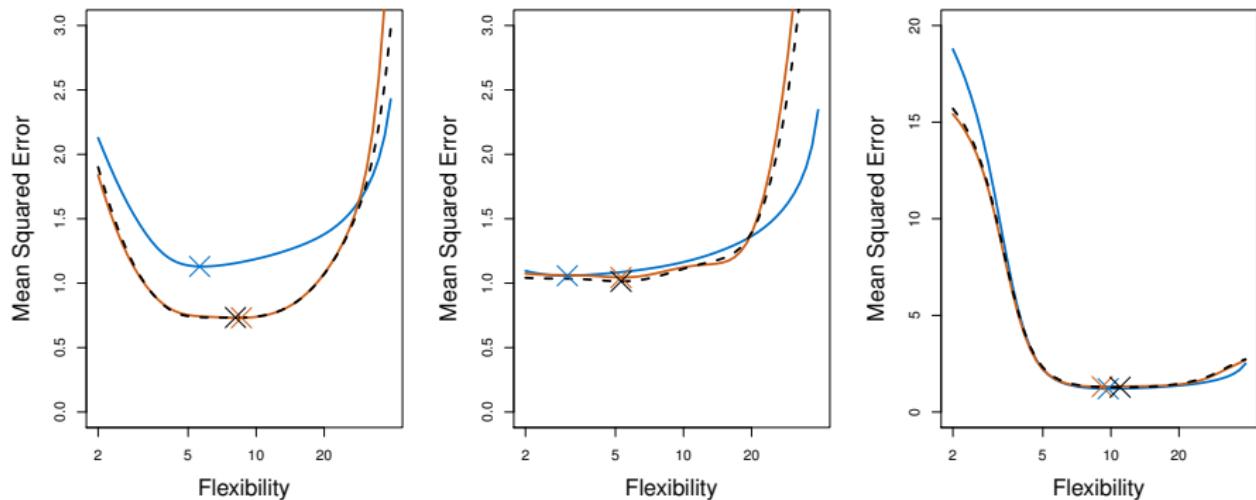
Left: target function (black), linear fit (orange), smoothing spline fits (blue & green). Right: training error (grey), prediction error (red),  $\text{Var}(e)$  (dashed).

# Cross-Validation for Regression



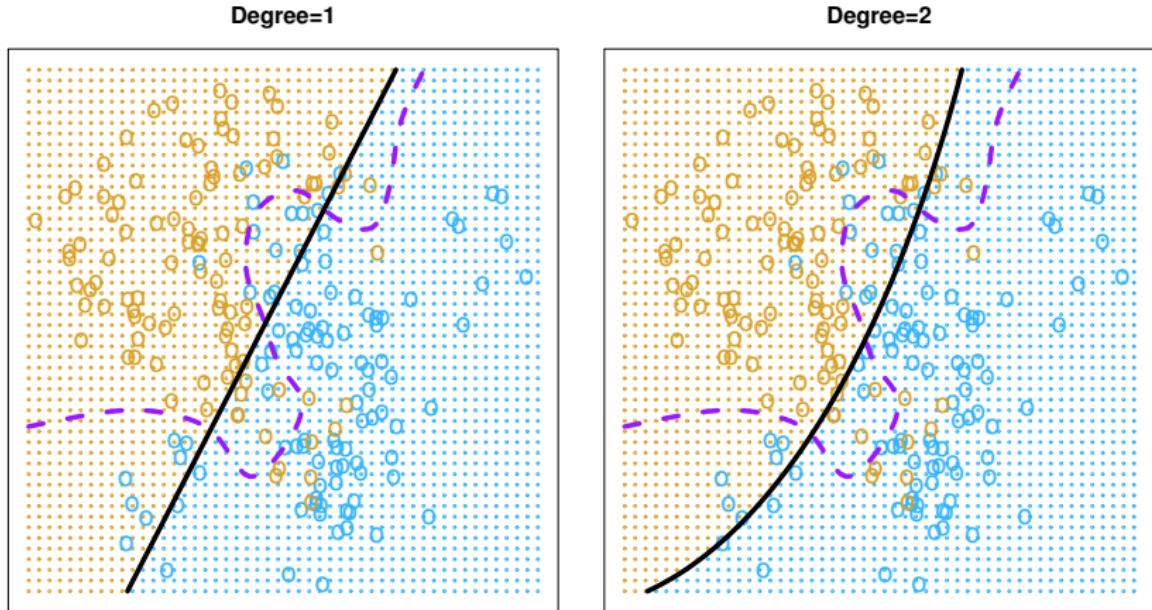
Left: target function (black), linear fit (orange), smoothing spline fits (blue & green). Right: training error (grey), prediction error (red),  $\text{Var}(e)$  (dashed).

# Cross-Validation for Regression



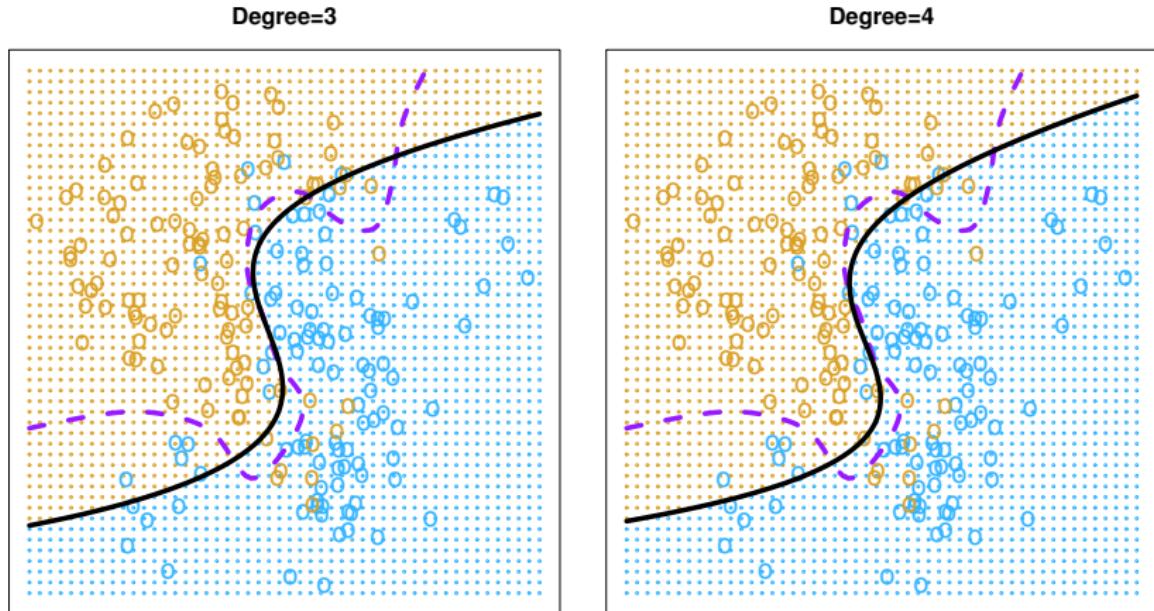
True and estimated prediction error for the simulated data sets. True prediction error (blue), LOOCV estimate (black dashed), 10-fold CV estimate (orange). Crosses indicate the minimum of each of the error curves.

# Cross-Validation for Classification



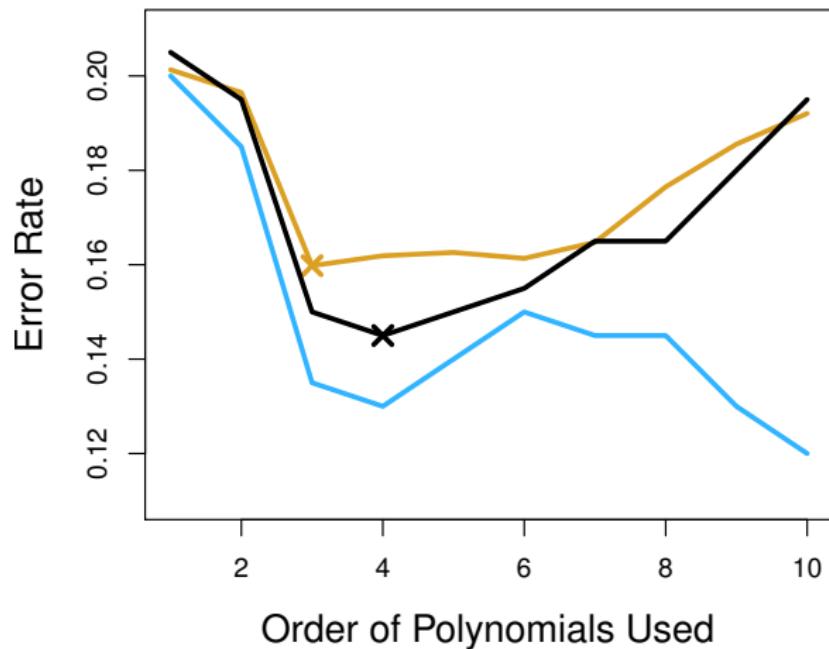
Two-dimensional classification data. Purple dashed: the Bayes decision boundary;  
Black: estimated decision boundaries from linear (left) and quadratic (right)  
logistic regressions.

# Cross-Validation for Classification



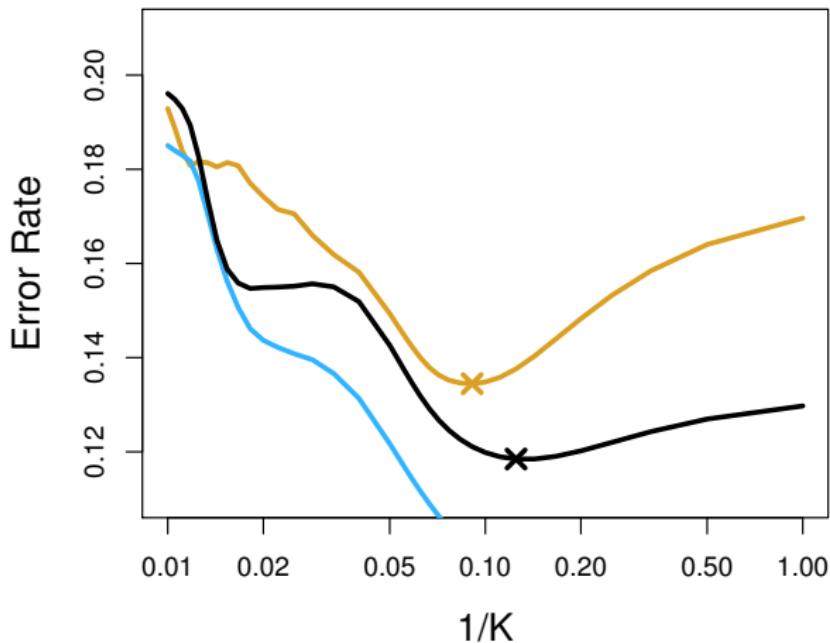
Two-dimensional classification data. Purple dashed: the Bayes decision boundary; Black: estimated decision boundaries from cubic (left) and quartic (right) logistic regressions.

# Cross-Validation for Classification



Logistic regression using polynomial functions of the predictors. True misclassification rate (brown), training error (blue), 10-fold CV error (black).

# Cross-Validation for Classification



KNN classifier. True misclassification rate (brown),  
training error (blue), 10-fold CV error (black).

# Note on Cross-Validation

## Note

In choosing the  $K$  folds at random, we are assuming that  $(x_i, y_i)$  are independently drawn from the underlying population  $p(x, y)$ , i.e., our data  $\mathcal{D} = \{(x_1, y_1), \dots, (x_N, y_N)\}$  constitutes a *random sample*.

This would not be true in *structured problems* where  $(x_i, y_i)$  are not independent – for example, when  $\mathcal{D}$  is time-series data<sup>a</sup>.

---

<sup>a</sup>See [here](#) for discussions on cross-validation for time series data.

## Alternatives to CV: Information Criteria

Suppose for dependent variable  $y$ , we have two potential predictors,  $x_1$  and  $x_2$ . If we limit our attention to linear models without interactions, we could fit the following four models:

$$\text{Model A: } y_i = \beta_0 + e_i$$

$$\text{Model B: } y_i = \beta_0 + \beta_1 x_{i1} + e_i$$

$$\text{Model C: } y_i = \beta_0 + \beta_2 x_{i2} + e_i$$

$$\text{Model D: } y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + e_i$$

- Model A is the **null model**, which contains no predictors.
- Model D is the **full model**, which contains all potential predictors.

## Alternatives to CV: Information Criteria

We could use CV to select the best model. Alternatively, we can use **information criteria** for model selection. These are metrics that make adjustment to the training error in order to account for the bias due to overfitting.

### Akaike's Information Criterion (AIC)

$$\text{AIC} = -2 \log \mathcal{L} + 2p = \text{Deviance} + 2p$$

, where  $\mathcal{L}$  is the value of the likelihood function<sup>3</sup>, and  $p$  is the number of parameters.

### Bayesian Information Criterion (BIC)

$$\text{BIC} = -2 \log \mathcal{L} + p \log N = \text{Deviance} + p \log N$$

, where  $N$  is the number of data points.

---

<sup>3</sup> achieved at the maximum-likelihood estimate.

## Alternatives to CV: Information Criteria

- For both AIC and BIC, smaller values are better. We can use them for model selection by choosing the model with the minimum AIC or BIC.
- The term  $2p$  in AIC and  $p \log N$  in BIC can be regarded as **penalty** terms that favor simpler models over more complicated ones.
- Both AIC and BIC have theoretical justifications that rely on asymptotic arguments.
- Since  $\log N > 2$  for  $N > 7$ , BIC generally increases with  $p$  faster than AIC. Thus BIC tends to lead to the selection of smaller models than AIC.

# MPG and Horsepower

```
require(foreach)
foreach (i = 1:5) %do% {
  model <- lm(mpg ~ poly(horsepower,i))
  c(AIC(model), BIC(model))
}

## [[1]]
## [1] 2363.324 2375.237
##
## [[2]]
## [1] 2274.354 2290.239
##
## [[3]]
## [1] 2275.531 2295.388
##
## [[4]]
## [1] 2276.108 2299.936
##
## [[5]]
## [1] 2268.663 2296.462
```

## Subset Selection

- In general, given  $p$  potential predictors in a linear model, we can use CV or information criteria to choose the best subset of predictors. This is also called **variable selection**.
- However, given  $p$  potential predictors, there are  $2^p$  possible models – each involves a subset of the  $p$  predictors. Comparing all of them becomes computationally infeasible when  $p$  is large.
- **Forward stepwise selection** is a computationally efficient alternative that begins with the null model, and then adds predictors one-at-a-time.

# Subset Selection

## Forward Stepwise Selection

- ① Let  $\mathcal{M}_0$  denote the null model.
- ② Fit all univariate models. Choose the one with the best in-sample fit (smallest RSS, highest  $R^2$ ) and add that variable – say  $x_{(1)}$  – to  $\mathcal{M}_0$ . Call the resulting model  $\mathcal{M}_1$ .
- ③ Fit all bivariate models that include  $x_{(1)}$ :  $y \sim \beta_0 + \beta_{(1)}x_{(1)} + \beta_j x_j$ , and add  $x_j$  from the one with the best in-sample fit to  $\mathcal{M}_1$ . Call the resulting model  $\mathcal{M}_2$ .
- ④ Continue until your model selection rule (cross-validation error, AIC, BIC) is lower for the current model than for any of the models that add one variable.

## Subset Selection

- Similarly, we can also do **backward stepwise selection**: start with the full model and remove predictors one-at-a-time.
- Forward is in general better than backward:
  - ▶ The full model can be expensive to fit, while the null model is usually available in closed form.
  - ▶ The full model fit has high variance (because it is usually overfit). The null model is always the same.
  - ▶ Backward selection requires  $N > p$  (so that the full model can be fit). Forward selection can be used even when  $p > N$ .
- Neither forward nor backward selection is guaranteed to find the best subset of predictors.

# Credit Card Balance

```
require(AER)
credit <- read.csv("Credit.csv")
fit <- lm(Balance~.,credit)
coeftest(fit)

##
## t test of coefficients:
##  

##           Estimate Std. Error t value Pr(>|t|)  

## (Intercept) -484.017314   26.076322 -18.5616 < 2.2e-16 ***  

## Income       -7.798884    0.233752 -33.3639 < 2.2e-16 ***  

## Limit        0.191288    0.032555  5.8759 9.011e-09 ***  

## Rating       1.129091    0.487398  2.3166  0.02104 *  

## Cards        17.919886   4.329445  4.1391 4.274e-05 ***  

## Age          -0.638650    0.293004 -2.1797  0.02988 *  

## GenderFemale -10.327215   9.896964 -1.0435  0.29737  

## StudentYes    425.498590  16.608864 25.6188 < 2.2e-16 ***  

## MarriedYes   -7.482805   10.249865 -0.7300  0.46580  

## ---  

## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

# Credit Card Balance

```
require(leaps)
all.fit <- regsubsets(Balance~., credit)
all <- summary(all.fit)
all$outmat
```

```
##           Income Limit Rating Cards Age GenderFemale StudentYes MarriedYe
## 1      ( 1 ) " "   " "   "*"   " "   " "   " "   " "
## 2      ( 1 ) "*"   " "   "*"   " "   " "   " "   " "
## 3      ( 1 ) "*"   " "   "*"   " "   " "   " "   "*"   " "
## 4      ( 1 ) "*"   "*"   " "   "*"   " "   " "   "*"   " "
## 5      ( 1 ) "*"   "*"   "*"   "*"   " "   " "   "*"   " "
## 6      ( 1 ) "*"   "*"   "*"   "*"   "*"   " "   "*"   " "
## 7      ( 1 ) "*"   "*"   "*"   "*"   "*"   "*"   "*"   " "
## 8      ( 1 ) "*"   "*"   "*"   "*"   "*"   "*"   "*"   "*"   "
```

# Credit Card Balance

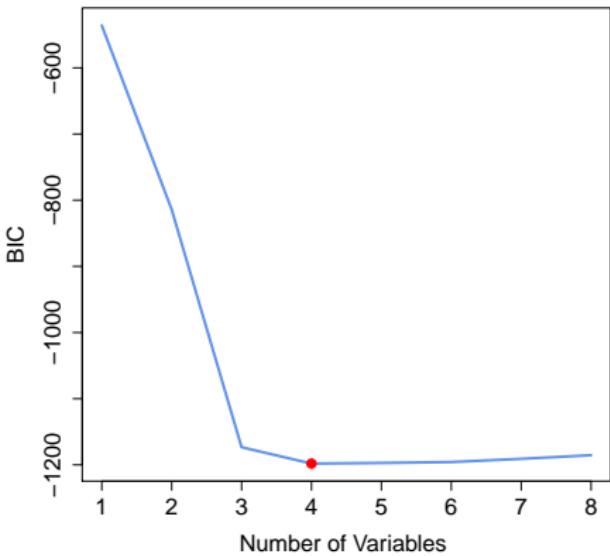
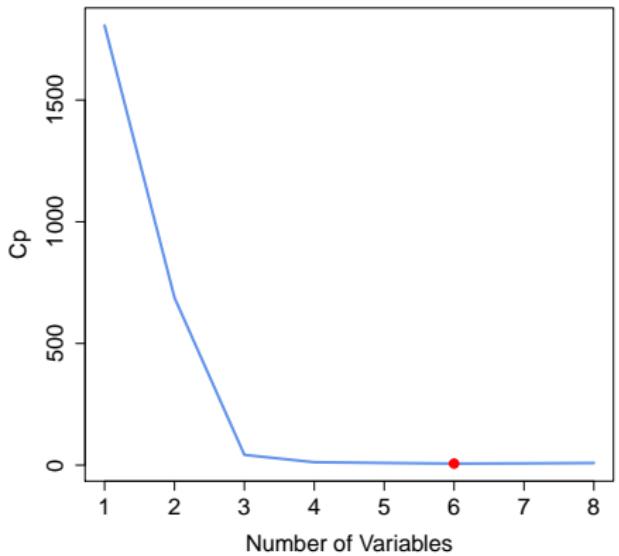
```
# Mallow's Cp statistic is equivalent to AIC in the case of
# linear models with Gaussian errors
all$cp #Mallow's Cp

## [1] 1806.332193 688.156412 42.321818 12.249136 9.218038 6.648
## [7] 7.532958 9.000000

all$bic

## [1] -535.9468 -814.1798 -1173.3585 -1198.0527 -1197.0957 -1195.7321
## [7] -1190.8790 -1185.4324
```

# Credit Card Balance



# Credit Card Balance

```
## best model according to AIC
```

```
coef(all.fit, 6)
```

```
## (Intercept) Income Limit Rating Cards
## -493.7341870 -7.7950824 0.1936914 1.0911874 18.2118976
## Age StudentYes
## -0.6240560 425.6099369
```

```
## best model according to BIC
```

```
coef(all.fit, 4)
```

```
## (Intercept) Income Limit Cards StudentYes
## -499.7272117 -7.8392288 0.2666445 23.1753794 429.6064203
```

# Credit Card Balance

```
## Forward stepwise selection using AIC
null <- lm(Balance~1,credit)
full <- lm(Balance~.,credit)
fwd <- step(null, scope = formula(full), direction="forward")
```

```
## Start:  AIC=4905.56
## Balance ~ 1
##
##           Df Sum of Sq      RSS      AIC
## + Rating   1  62904790  21435122  4359.6
## + Limit    1  62624255  21715657  4364.8
## + Income   1  18131167  66208745  4810.7
## + Student  1   5658372  78681540  4879.8
## + Cards    1   630416   83709496  4904.6
## <none>          84339912  4905.6
## + Gender   1    38892  84301020  4907.4
## + Married  1     2715  84337197  4907.5
## + Age      1      284  84339628  4907.6
##
```

# Credit Card Balance

```
##  
## Step: AIC=4359.63  
## Balance ~ Rating  
##  
##           Df Sum of Sq      RSS      AIC  
## + Income   1  10902581 10532541 4077.4  
## + Student  1   5735163 15699959 4237.1  
## + Age      1    649110 20786012 4349.3  
## + Cards    1    138580 21296542 4359.0  
## + Married  1    118209 21316913 4359.4  
## <none>          21435122 4359.6  
## + Gender   1     16065 21419057 4361.3  
## + Limit    1      7960 21427162 4361.5  
##
```

# Credit Card Balance

```
##  
## Step: AIC=4077.41  
## Balance ~ Rating + Income  
##  
##           Df Sum of Sq      RSS      AIC  
## + Student  1   6305322  4227219  3714.2  
## + Married  1     95068 10437473  4075.8  
## + Limit    1     94545 10437996  4075.8  
## + Age      1     90286 10442255  4076.0  
## <none>          10532541 4077.4  
## + Cards    1     2094 10530447  4079.3  
## + Gender   1      948 10531593  4079.4  
##
```

# Credit Card Balance

```
##  
## Step: AIC=3714.24  
## Balance ~ Rating + Income + Student  
##  
##          Df Sum of Sq    RSS     AIC  
## + Limit   1  194718 4032502 3697.4  
## + Age     1   44620 4182600 3712.0  
## <none>           4227219 3714.2  
## + Married 1   13083 4214137 3715.0  
## + Gender   1   12168 4215051 3715.1  
## + Cards   1   10608 4216611 3715.2  
##
```

# Credit Card Balance

```
##  
## Step: AIC=3697.37  
## Balance ~ Rating + Income + Student + Limit  
##  
##          Df Sum of Sq      RSS      AIC  
## + Cards     1   166410 3866091 3682.5  
## + Age       1    37952 3994549 3695.6  
## <none>           4032502 3697.4  
## + Gender    1    13345 4019157 3698.0  
## + Married   1     6660 4025842 3698.7  
##
```

# Credit Card Balance

```
##  
## Step: AIC=3682.52  
## Balance ~ Rating + Income + Student + Limit + Cards  
##  
##           Df Sum of Sq      RSS      AIC  
## + Age     1   44472 3821620 3679.9  
## <none>            3866091 3682.5  
## + Gender   1   11350 3854741 3683.3  
## + Married  1    3121 3862970 3684.2  
##
```

# Credit Card Balance

```
##  
## Step: AIC=3679.89  
## Balance ~ Rating + Income + Student + Limit + Cards + Age  
##  
##           Df Sum of Sq      RSS      AIC  
## <none>            3821620 3679.9  
## + Gender    1   10860.9 3810759 3680.7  
## + Married   1    5450.6 3816169 3681.3  
##
```

# Credit Card Balance

```
coeftest(fwd)

##
## t test of coefficients:
##
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) -493.734187   24.824765 -19.8888 < 2.2e-16 ***
## Rating       1.091187    0.484804   2.2508   0.02495 *
## Income      -7.795082    0.233417 -33.3955 < 2.2e-16 ***
## StudentYes  425.609937   16.509565  25.7796 < 2.2e-16 ***
## Limit        0.193691    0.032383   5.9813 4.980e-09 ***
## Cards        18.211898    4.318650   4.2170 3.075e-05 ***
## Age         -0.624056    0.291817  -2.1385   0.03309 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

# Credit Card Balance

```
## Forward stepwise selection using BIC
n <- nrow(credit)
fwd <- step(null, scope = formula(full), direction="forward", k=log(n))

coeftest(fwd)

##
## t test of coefficients:
##
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) -526.155523   19.746614 -26.6454 < 2.2e-16 ***
## Rating       1.087901    0.486995   2.2339   0.02605 *
## Income      -7.874924    0.231455 -34.0236 < 2.2e-16 ***
## StudentYes  426.850146   16.574025  25.7542 < 2.2e-16 ***
## Limit        0.194409    0.032527   5.9768 5.099e-09 ***
## Cards        17.851731    4.334889   4.1182 4.656e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

# The Bias-Variance Trade-off

The goal of model selection is to choose the model (out of the set of candidate models) that can best balance the bias-variance tradeoff to achieve the smallest out-of-sample error.

In the case of linear models, we have:

$$E \left[ (y - x' \hat{\beta})^2 \right] = Var(x' \hat{\beta}) + [bias(x' \hat{\beta})]^2 + Var(e)$$

- The OLS estimate  $\hat{\beta}^{LS}$  is unbiased:  $E[\hat{\beta}^{LS}] = \beta^*$ . Hence  $bias(x' \hat{\beta}^{LS}) = E[x' \beta^* - x' \hat{\beta}^{LS}] = 0$ . What about  $Var(x' \hat{\beta}^{LS})$ ?

# The Gauss-Markov Theorem

Note that  $\hat{\beta}^{LS}$  is a linear function of  $Y$ :  $\hat{\beta}^{LS} = (X'X)^{-1} X' Y = \sum_i w_i y_i$ , where  $w_i = (X'X)^{-1} x_i$ .

## Gauss-Markov Theorem

Under the assumption that

- ① The true CEF is linear:  $E[y|x] = x'\beta^*$
- ② Homoskedasticity:  $E[(e^*)^2 | x] = \sigma^2$

The OLS estimator  $\hat{\beta}^{LS}$  is **BLUE**: Best Linear Unbiased Estimator, in the sense that among all *unbiased* linear estimators (estimators that are linear functions of  $Y$ ),  $\hat{\beta}^{LS}$  has the *smallest* variance.

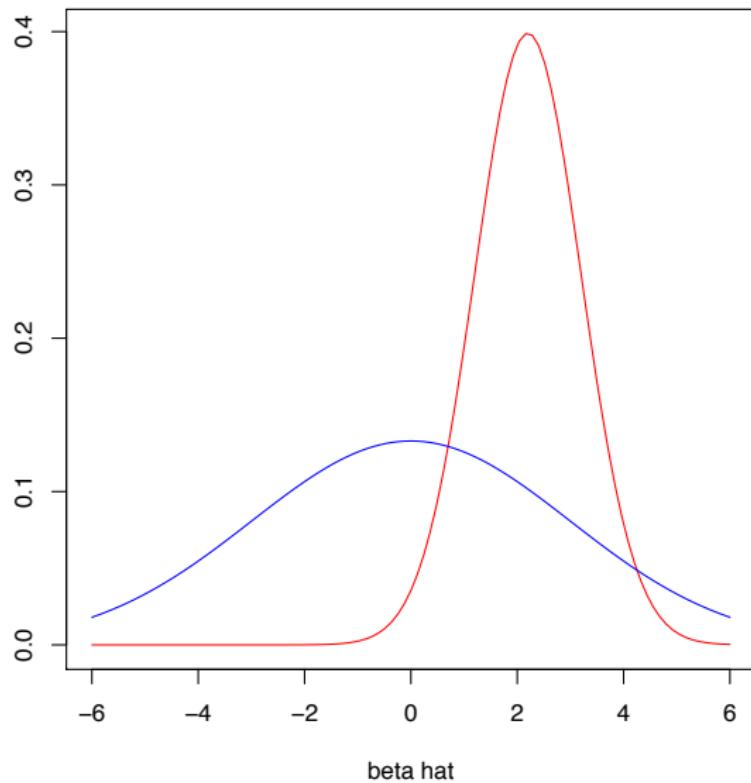
# The Gauss-Markov Theorem

Under the assumptions of Gauss-Markov,

$$E \left[ (y - x' \hat{\beta}^{LS})^2 \right] = \text{Var} (x' \hat{\beta}^{LS}) + \sigma^2 \quad (3)$$

- The theorem says that if we have any other *unbiased* linear estimator  $\hat{\beta}$ , then  $\text{Var} (x' \hat{\beta}^{LS}) \leq \text{Var} (x' \hat{\beta})$ .
- But can we do better – reduce the prediction error – if we are willing to allow a little bias in exchange for a larger reduction in variance?

# Exploring the Bias-Variance Trade-off



# Exploring the Bias-Variance Trade-off

(3)  $\Rightarrow$

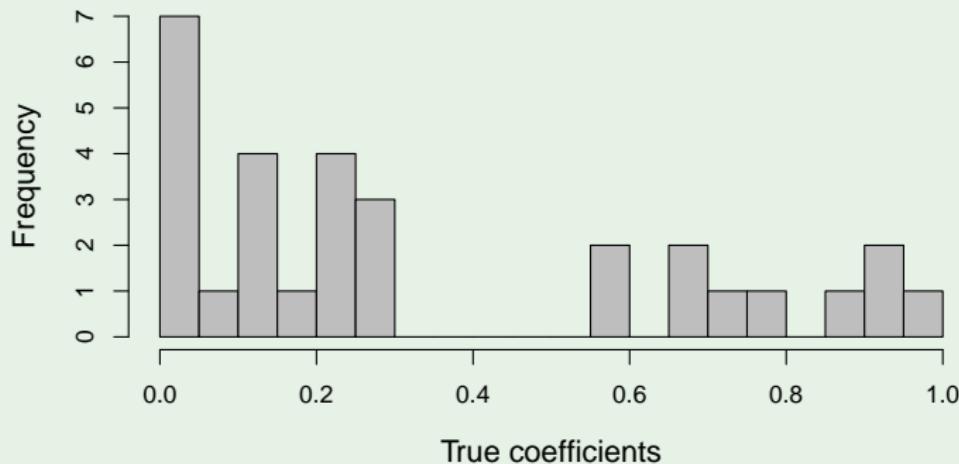
$$\begin{aligned} E_{\mathcal{D},x} \left[ (y - x' \hat{\beta}^{LS})^2 \right] &\approx \frac{1}{N} \sum_{i=1}^N E_{\mathcal{D}} \left[ (y_i - x_i' \hat{\beta}^{LS})^2 \right] \\ &= \frac{1}{N} \sum_{i=1}^N \text{Var} (x_i' \hat{\beta}^{LS}) + \sigma^2 \\ &= \frac{1}{N} \text{trace} (\text{Var} (X \hat{\beta}^{LS})) + \sigma^2 \\ &\approx \frac{1}{N} \text{trace} (X (X' X)^{-1} X') \sigma^2 + \sigma^2 \\ &= \frac{1}{N} (1 + p) \sigma^2 + \sigma^2 \end{aligned}$$

, i.e., the prediction error scales *linearly* with the number of predictors  $p$ .

# Exploring the Bias-Variance Trade-off

## Simulation 1

Simulation:  $N = 50, p = 30$ .  $y \sim \mathcal{N}(x'\beta^*, 1)$ , where  $\beta^* \in \mathbb{R}^{30}$  with 10 “large” coefficients (between 0.5 and 1) and 20 “small” ones (between 0 and 0.3).

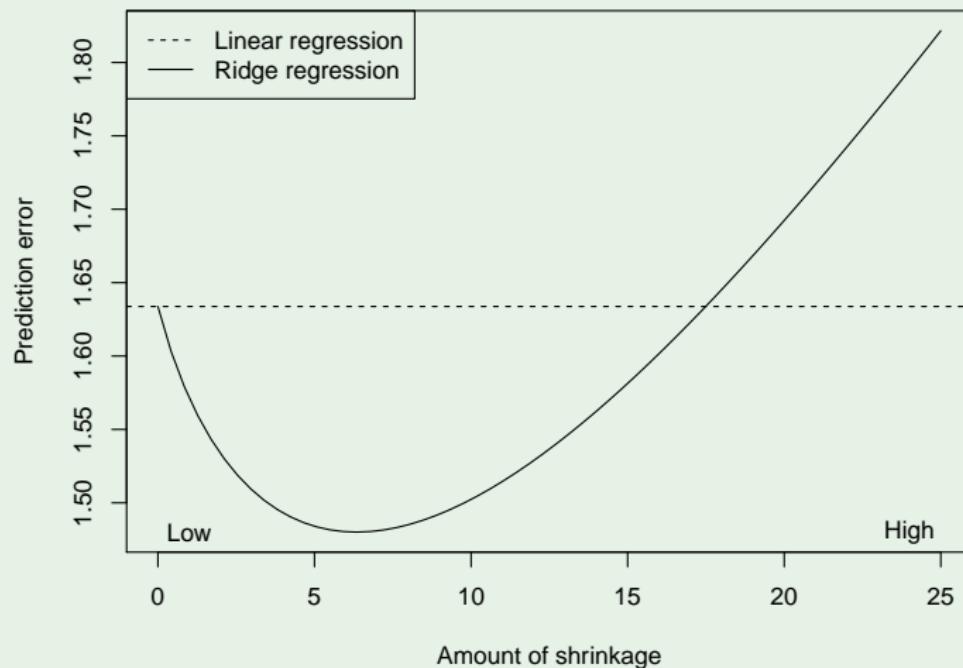


## Exploring the Bias-Variance Trade-off

- Each additional predictor adds the same amount of variance  $\frac{\sigma^2}{N}$ , regardless of whether its true coefficient is large or small (or zero).
- So in this example, we are “spending” variance in trying to fit truly small coefficients—20 of them, out of 30 total.
- We can do better by **shrinking** small coefficients towards zero, incurring some bias, so as to reduce the variance. Methods that do this are called **shrinkage methods**.

# Exploring the Bias-Variance Trade-off

## Simulation 1



# Exploring the Bias-Variance Trade-off

## Simulation 1

Linear regression:

$$\text{Square bias} \approx 0.006$$

$$\text{Variance} \approx 0.627$$

$$\text{prediction error} \approx 1 + 0.006 + 0.627 = 1.633$$

Ridge regression, at its best:

$$\text{Square bias} \approx 0.077$$

$$\text{Variance} \approx 0.403$$

$$\text{prediction error} \approx 1 + 0.077 + 0.403 = 1.48$$

# Ridge Regression

$$\hat{\beta}^R = \arg \min_{\beta} \left\{ \sum_{i=1}^N \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \underbrace{\sum_{j=1}^p \beta_j^2}_{\text{penalty}} \right\} \quad (4)$$

, where  $\lambda \geq 0$  is a **tuning parameter** that controls the strength of the penalty term, which has the effect of shrinking the estimates towards zero.

- $\lambda = 0 \Rightarrow \hat{\beta}^R = \hat{\beta}^{LS}$ <sup>4</sup>
- $\lambda = \infty \Rightarrow \hat{\beta}_1^R = \dots = \hat{\beta}_p^R = 0$ <sup>5</sup>

---

<sup>4</sup>  $\hat{\beta}^{LS}$  denotes the least square estimate.

<sup>5</sup>  $\hat{\beta}_0$  is not shrunk toward 0.

## Ridge Regression

An equivalent way to write the ridge problem is:

$$\hat{\beta}^R = \arg \min_{\beta} \left\{ \sum_{i=1}^N \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \right\} \quad (5)$$

subject to  $\sum_{j=1}^p \beta_j^2 \leq C$

(4) is the Lagrangian form of (5). There is a one-to-one correspondence between  $\lambda$  in (4) and  $C$  in (5).

# Ridge Regression

- Least square estimates are **scale equivariant**: multiplying  $x_j$  by a constant  $c$  leads to a scaling of  $\hat{\beta}_j$  by  $1/c$ .  $\hat{\beta}_j x_j$  remains the same.
- Ridge regression estimates are *not* scale equivariant. They can change substantially when multiplying a given predictor by a constant. Therefore,  $x_{ij}$  need to be *standardized* before applying ridge regression.
- In addition, note that  $\beta_0$  is not penalized: we do *not* shrink  $\hat{\beta}_0$ . Otherwise we could add some constant  $c$  to  $y$  and the solution would be different.

# Ridge Regression

Let  $\tilde{x}_{ij} = \frac{x_{ij} - \bar{x}_j}{\sqrt{\frac{1}{N} \sum_{i=1}^N (x_{ij} - \bar{x}_j)^2}}$ , so that  $x_{ij}$ s are both *centered* and *standardized*.

If we use  $\tilde{x}_{ij}$  as the regressors, then  $\hat{\beta}_0 = \bar{y}$ <sup>6</sup> and we can estimate the the rest of the  $\beta$ s by a ridge regression without intercept:

$$\hat{\beta}^R = \arg \min_{\beta} \left\{ \sum_{i=1}^N (\tilde{y}_i - \tilde{x}'_i \beta)^2 + \lambda \|\beta\|_2^2 \right\} \quad (6)$$

, where  $\|\cdot\|_2$  denotes the  $\ell_2$  norm,  $\tilde{y}_i = y_i - \bar{y}$ ,  $\tilde{x}_i = [\tilde{x}_{i1}, \dots, \tilde{x}_{ip}]'$ , and  $\beta = [\beta_1, \dots, \beta_p]'$ .

---

<sup>6</sup>This is true as long as  $\tilde{x}_{ij}$  are centered.

# Ridge Regression

(6)  $\Rightarrow$

$$\hat{\beta}^R = (\tilde{X}'\tilde{X} + \lambda I)^{-1} \tilde{X}'Y \quad (7)$$

, where  $I$  is the  $p \times p$  identity matrix, and  $\tilde{X}$  is centered and standardized.

In the case that  $\tilde{X}$  is *orthonormal*<sup>7</sup>, (7)  $\Rightarrow$

$$\hat{\beta}_j^R = \frac{\hat{\beta}_j^{LS}}{1 + \lambda} \quad (8)$$

---

<sup>7</sup>so that  $\tilde{X}'\tilde{X} = I$

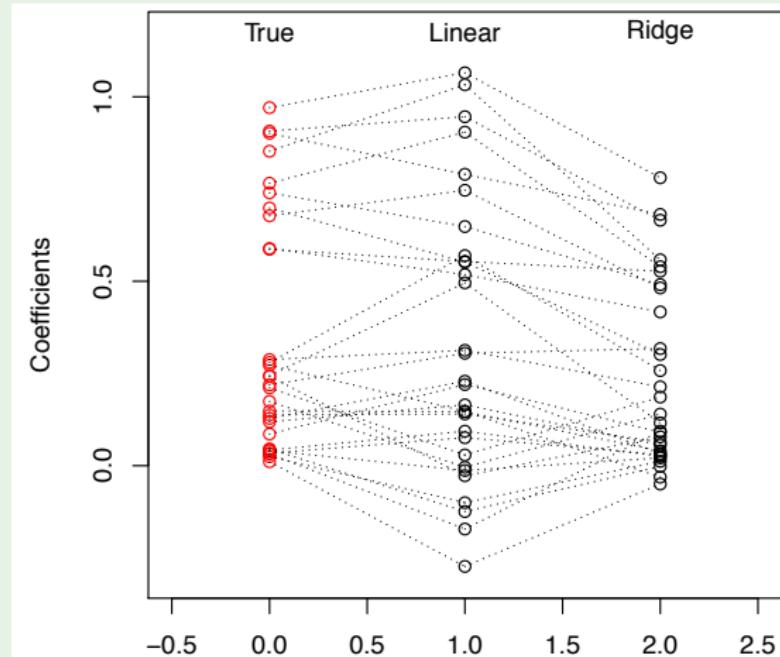
## Ridge Regression

To choose  $\lambda$ , we can treat  $\lambda$  as a **hyperparameter** and select its value via cross-validation:

- Given a grid of  $\lambda$  values, perform cross-validation on the training set to pick the  $\lambda$  for which the cross-validation error is smallest.
- Re-fit the model with the selected  $\lambda$  using all the training data.
- Benchmark the performance of the final fitted model using a test set.

# Ridge Regression

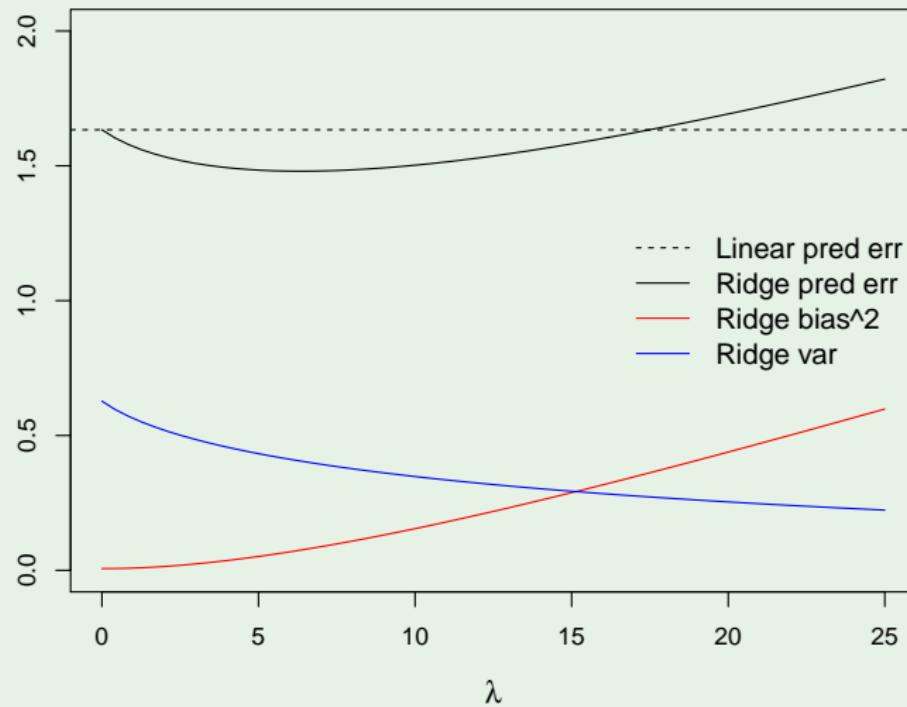
## Simulation 1



Ridge regression coefficients at  $\lambda = 25$

# Ridge Regression

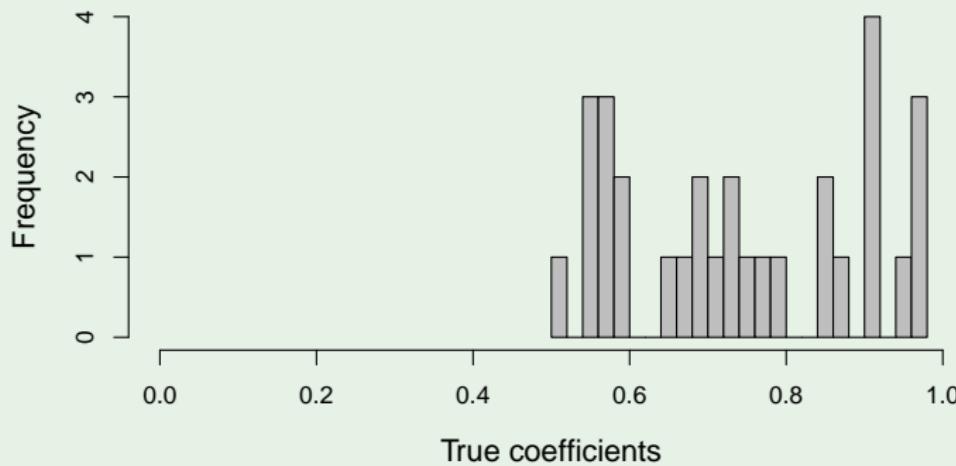
## Simulation 1



# Ridge Regression

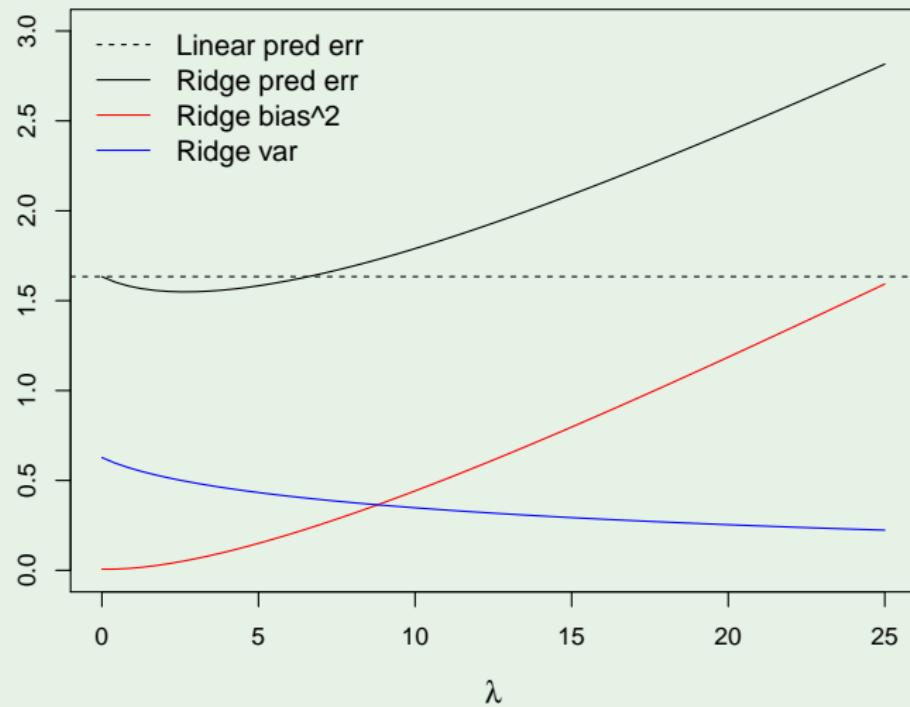
## Simulation 2

Simulation:  $N = 50, p = 30$ .  $y \sim \mathcal{N}(x'\beta^*, 1)$ , where  $\beta^* \in \mathbb{R}^{30}$  with all 30 coefficients between 0.5 and 1.



# Ridge Regression

## Simulation 2



# Credit Card Balance

```
require(glmnet)
x <- model.matrix(Balance~.,credit)[,-1] #no intercept
y <- credit$Balance
ridgefit <- glmnet(x,y,alpha=0,lambda=0.01)
coef(ridgefit)

## 9 x 1 sparse Matrix of class "dgCMatrix"
##                               s0
## (Intercept) -479.7851081
## Income       -7.8002927
## Limit        0.2049782
## Rating       0.9251803
## Cards        18.9116337
## Age          -0.6364639
## GenderFemale -10.3247237
## StudentYes   426.0874930
## MarriedYes   -7.0597648
```

# Credit Card Balance

```
ridgefit <- glmnet(x,y,alpha=0,lambda=1e10)
coef(ridgefit)

## 9 x 1 sparse Matrix of class "dgCMatrix"
##                               s0
## (Intercept) 5.200149e+02
## Income      2.777310e-07
## Limit       7.881306e-09
## Rating      1.178376e-07
## Cards       1.331034e-06
## Age         2.245918e-09
## GenderFemale 9.061130e-07
## StudentYes   1.820460e-05
## MarriedYes   -2.455471e-07
```

# Credit Card Balance

```
ridgefit <- glmnet(x,y,alpha=0) #automatically select a range of lambda
dim(coef(ridgefit))

## [1] 9 100

c(ridgefit$lambda[1],ridgefit$lambda[100])

## [1] 396562.69957      39.65627

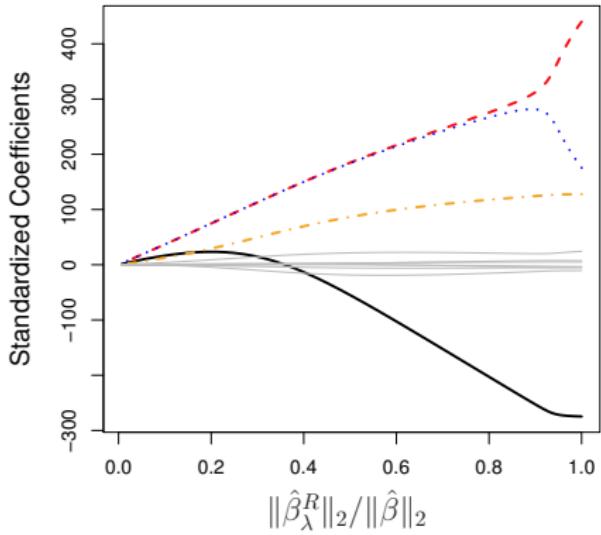
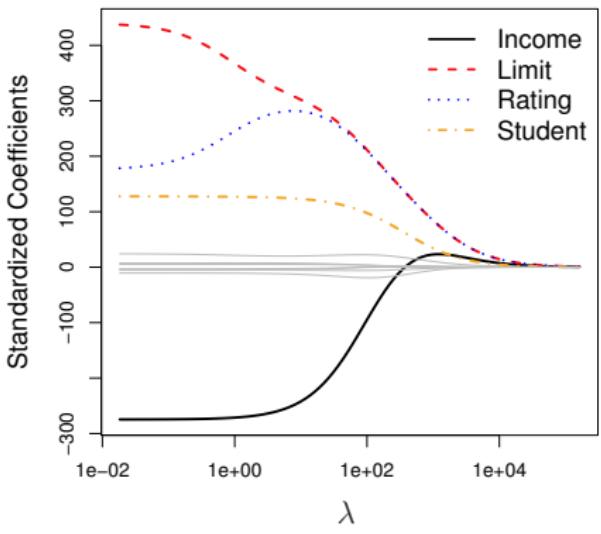
ridgefit$lambda[50]

## [1] 4154.453

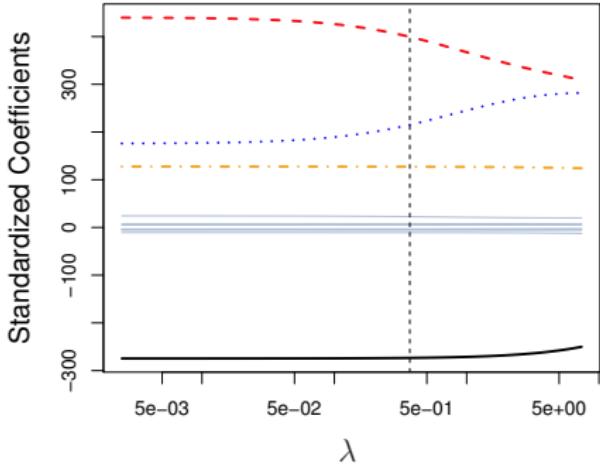
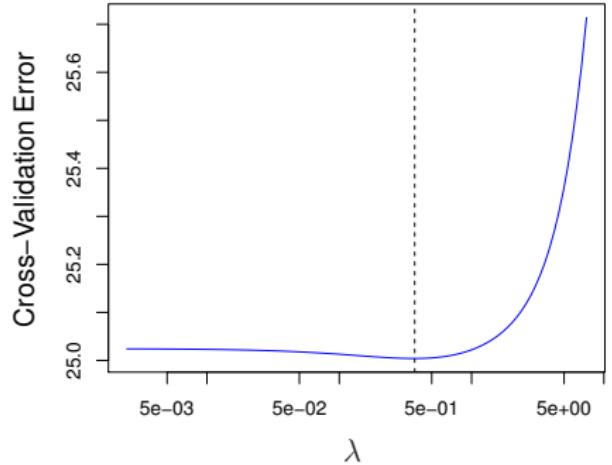
coef(ridgefit)[,50]

## (Intercept)      Income       Limit      Rating      Cards
## 339.02562071   0.44648135   0.01505380   0.22506908   2.77187416
##          Age GenderFemale StudentYes MarriedYes
## -0.05112386   1.76037503   39.44486798  -0.95049050
```

# Credit Card Balance



# Credit Card Balance



# Lasso

The lasso<sup>8</sup> is defined as:

$$\hat{\beta}^L = \arg \min_{\beta} \left\{ \sum_{i=1}^N \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p |\beta_j| \right\} \quad (9)$$

Equivalently,

$$\begin{aligned} \hat{\beta}^L &= \arg \min_{\beta} \left\{ \sum_{i=1}^N \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \right\} \\ &\text{subject to } \sum_{j=1}^p |\beta_j| \leq C \end{aligned} \quad (10)$$

---

<sup>8</sup> acronym for: Least Absolute Selection and Shrinkage Operator. Not to be confused with wonder woman's **weapon of choice**.

# Lasso

Centering and standardizing  $x_{ij} \Rightarrow$

$$\hat{\beta}^L = \arg \min_{\beta} \left\{ \sum_{i=1}^N (\tilde{y}_i - \tilde{x}'_i \beta)^2 + \lambda \|\beta\|_1 \right\} \quad (11)$$

, where  $\|\cdot\|_1$  denotes the  $\ell_1$  norm,  $\tilde{y}_i = y_i - \bar{y}$ ,  $\tilde{x}_{ij} = \frac{x_{ij} - \bar{x}_j}{\sqrt{\frac{1}{N} \sum_{i=1}^N (x_{ij} - \bar{x}_j)^2}}$ , and  $\beta = [\beta_1, \dots, \beta_p]'$ .

# Lasso

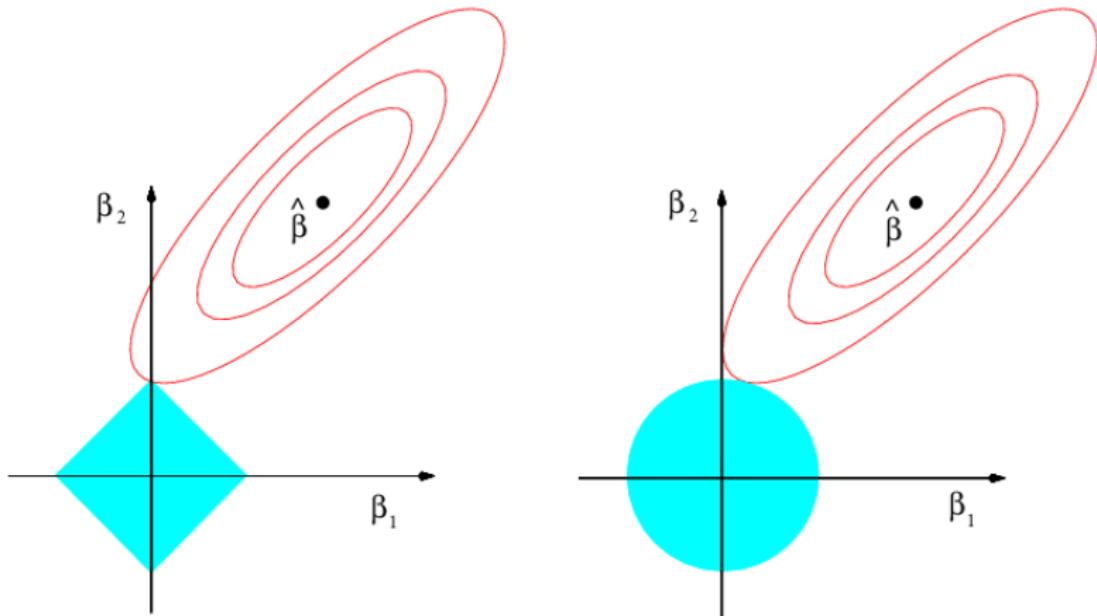
There is in general no closed-form solution to (11). In the case that  $X$  is *orthonormal*, we have:

$$\hat{\beta}^L = \begin{cases} \hat{\beta}^{LS} - \frac{\lambda}{2} & \text{if } \hat{\beta}^{LS} > \frac{\lambda}{2} \\ \hat{\beta}^{LS} + \frac{\lambda}{2} & \text{if } \hat{\beta}^{LS} < -\frac{\lambda}{2} \\ 0 & \text{if } |\hat{\beta}^{LS}| \leq \frac{\lambda}{2} \end{cases}$$

## Lasso

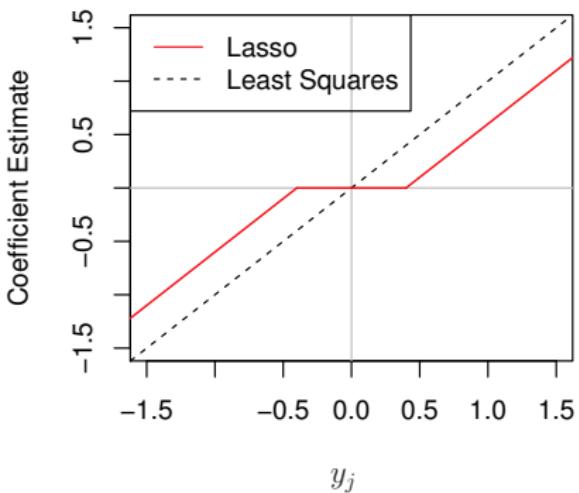
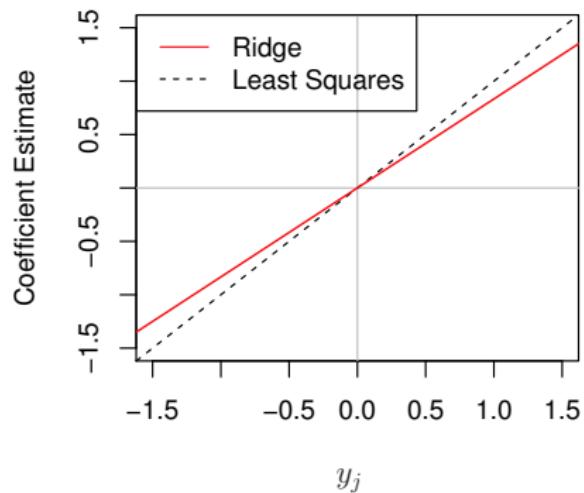
- As with ridge regression, the lasso shrinks the coefficient estimates towards zero.
- However, in the case of the lasso, the  $\ell_2$  penalty has the effect of forcing some of the coefficient estimates to be exactly equal to zero when the tuning parameter is sufficiently large.
- Hence, the lasso performs **variable selection**. We say that the lasso yields **sparse models** — that is, models that involve only a subset of the variables.
- As in ridge regression, selecting a good value of  $\lambda$  for the lasso is critical and is in practice done via cross-validation.

# Lasso



Contours of the error and constraint functions for the lasso (left) and ridge regression (right). The solid blue areas are the constraint regions,  $|\beta_1| + |\beta_2| \leq C$  and  $\beta_1^2 + \beta_2^2 \leq C$ , while the red ellipses are the contours of the RSS.

# Lasso



The ridge regression and the lasso coefficient estimates for a simple setting with  $N = p$  and  $X$  being orthonormal.

# Credit Card Balance

```
lassofit <- glmnet(x,y,alpha=1,lambda=0.01)
coef(lassofit)

## 9 x 1 sparse Matrix of class "dgCMatrix"
##                               s0
## (Intercept) -479.5814236
## Income        -7.7997555
## Limit         0.2056525
## Rating        0.9149775
## Cards         18.9546633
## Age           -0.6358133
## GenderFemale -10.3038053
## StudentYes    426.0922168
## MarriedYes    -7.0184389
```

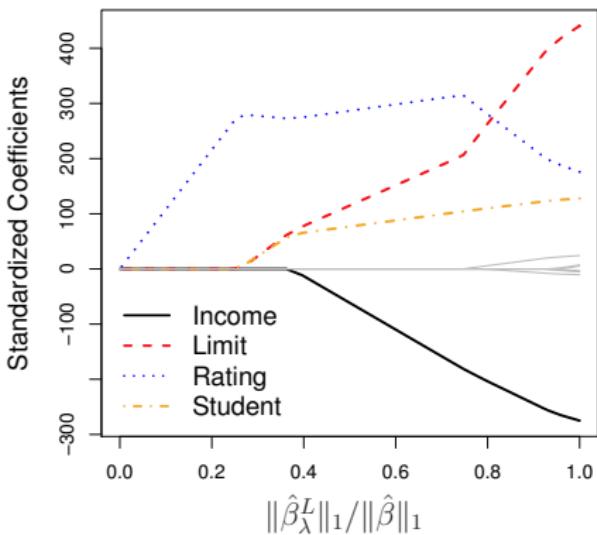
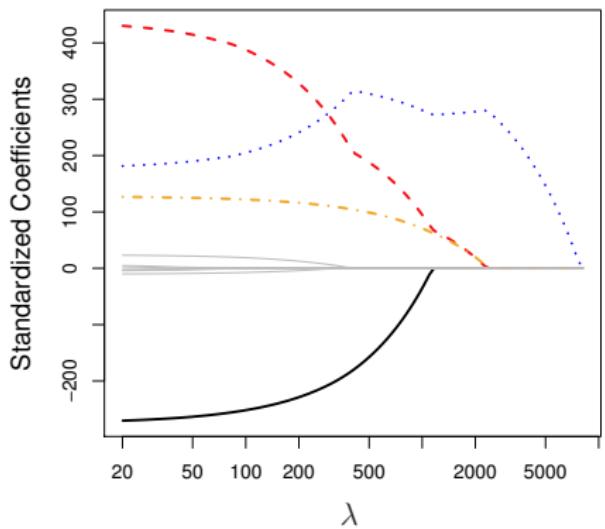
# Credit Card Balance

```
lassofit <- glmnet(x,y,alpha=1,lambda=100)
coef(lassofit)

## 9 x 1 sparse Matrix of class "dgCMatrix"
##                               s0
## (Intercept) -157.92073077
## Income        .
## Limit         0.01822877
## Rating        1.64828689
## Cards         .
## Age           .
## GenderFemale  .
## StudentYes    65.68634041
## MarriedYes    .
```

# Credit Card Balance

```
lassofit <- glmnet(x,y,alpha=1) #automatically select a range of lambda
```



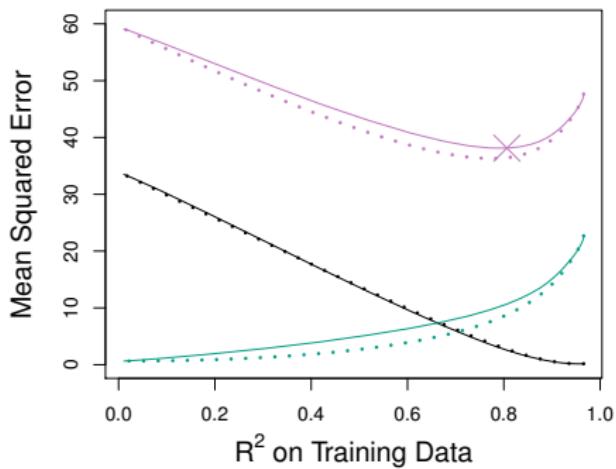
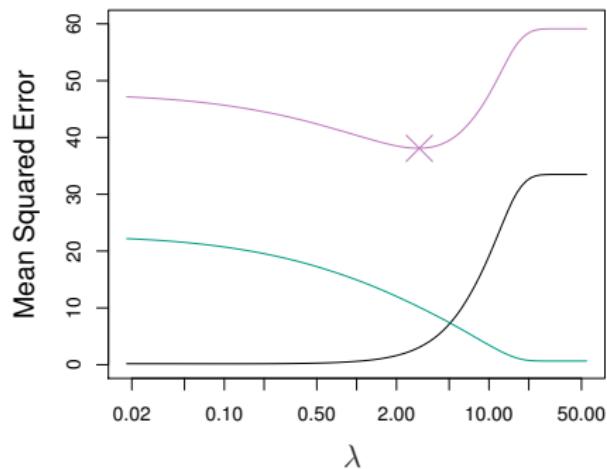
## Soft vs. Hard Order Constraint

- Consider a series of polynomial models:  $\mathcal{H}_0 \subset \mathcal{H}_1 \subset \mathcal{H}_2 \subset \dots$ , where  $\mathcal{H}_q = \{h(x) = \beta_0 + \beta_1 x + \dots + \beta_q x^q\}$ , then estimating, say  $\mathcal{H}_2$ , is equivalent to estimating, say  $\mathcal{H}_{10}$ , with the constraint that  $\beta_3 = \dots = \beta_{10} = 0$ <sup>9</sup>.
- Instead of imposing such “hard order constraints”, shrinkage methods impose “soft order constraints” by giving  $\beta$  a “budget” (e.g.,  $\|\beta\|_2^2 \leq C$  or  $\|\beta\|_1 \leq C$ ) that have the effect of shrinking  $\beta$  toward 0.
- In this sense, shrinkage can be thought of as selecting the final hypothesis from a continuous set of models rather than a discrete one (constant, linear, quadratic ...).

---

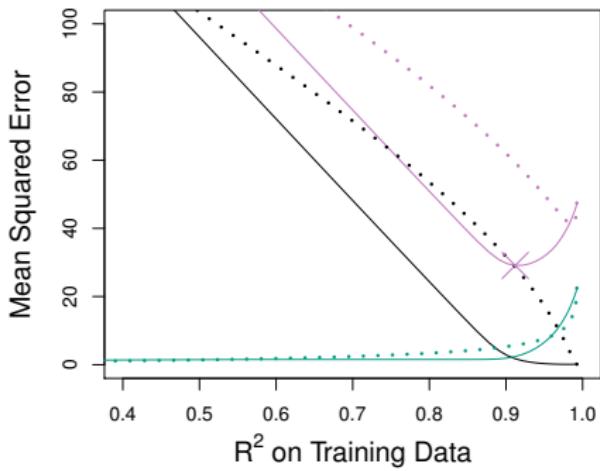
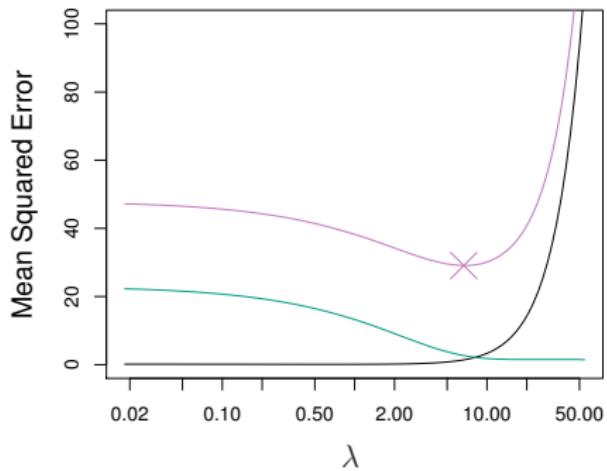
<sup>9</sup>Technically, this is only true if the polynomials are expressed as linear combinations of Legendre polynomials in  $x$  rather than consecutive powers of  $x$ .

# Comparing the Lasso and Ridge Regression



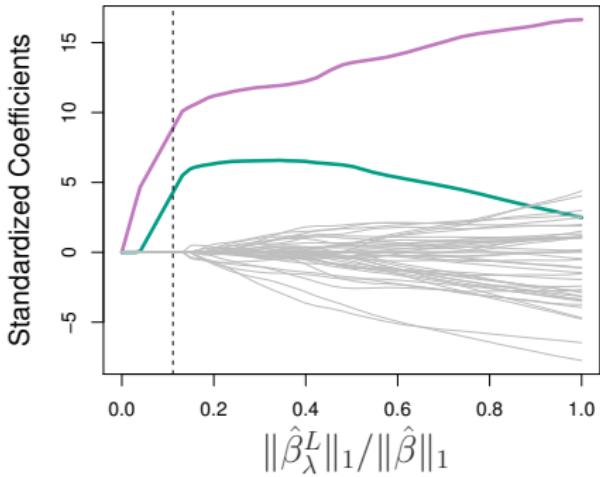
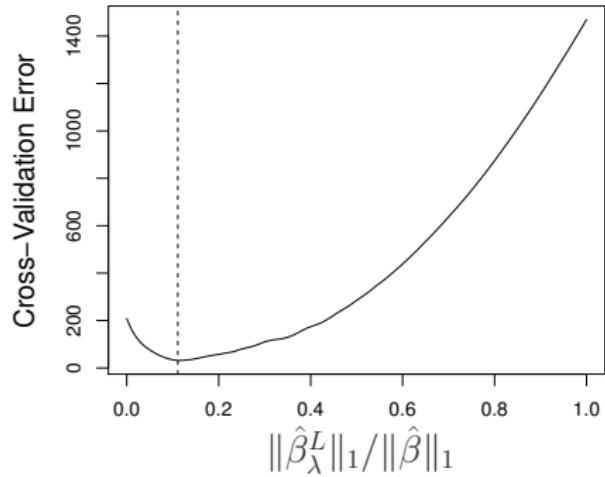
Plots of squared bias (black), variance (green), and prediction error (purple) for the lasso (solid) and ridge (dashed) on a simulated data set with  $p = 45$  and  $n = 50$ . All 45 predictors are related to the response.

# Comparing the Lasso and Ridge Regression



Comparison of squared bias, variance and prediction error between the lasso (solid) and ridge (dashed). Now only two predictors in the simulated data set are related to the response.

# Comparing the Lasso and Ridge Regression



Ten-fold cross-validation error for the lasso, applied to the simulated data set with 2 out of 45 predictors related to the response.

# Supermarket Entry

Supermarket entry in geographical markets. For each market, data include:

- total population
- population density
- average income
- percentage population with college degree
- percentage population between 25 and 64 years old
- percentage married
- housing price
- unemployment rate
- crime rate
- tax rate
- etc.

# Supermarket Entry

```
supermarket <- read.csv("Supermarket.csv")
supermarket$entry <- factor(supermarket$entry, levels=c(0,1),
                             labels=c("FALSE", "TRUE"))
summary(supermarket$entry)

## FALSE TRUE
## 1119 881

head(supermarket,2)

##   entry      pop density gas unemp  tax union housing crime publicTrans
## 1 FALSE 7064692     0.42 2.0 0.017 0.05      0    106  0.05          69
## 2  TRUE 2961394     0.39 4.9 0.039 0.07      1    121  0.07          47
##   cost_insurance young adult old married college sexratio meanincome
## 1                 2762  0.18  0.46 0.37     0.49    0.43    0.52        84603
## 2                 1258  0.28  0.39 0.33     0.51    0.30    0.53        67249
##   income25 income50 income75 cost_wage
## 1      55123    70779   116696     42930
## 2      47798    61374   87093     37225
```

# Supermarket Entry

```
# Split into training and test data sets
n <- nrow(supermarket)
train <- sample(n,n/2)
supermarket.tr <- supermarket[train,]
supermarket.te <- supermarket[-train,]
#
#####
# Logistic Regression: Full Model #
#####
# Fit on training data
logitfit <- glm(entry ~.,supermarket.tr,family='binomial')
```

# Supermarket Entry

```
require(AER)
coefest(logitfit)

##
## z test of coefficients:
##
##           Estimate   Std. Error z value Pr(>|z|)
## (Intercept) -0.7197401841 1.9375965904 -0.37  0.71029
## pop          0.0000000705 0.0000000239  2.94  0.00323 **
## density      1.8151478321 0.4272158129  4.25  0.000021 ***
## gas          0.0079951115 0.0721180090  0.11  0.91173
## unemp        -17.0951338342 13.6282037274 -1.25  0.20970
## tax          -5.7681705719 1.6654230134 -3.46  0.00053 ***
## union         -0.1975921113 0.1590731169 -1.24  0.21418
## housing      -0.0200346920 0.0071127694 -2.82  0.00485 **
## crime         3.9358630077 2.4066828117  1.64  0.10197
## publicTrans   0.0027691613 0.0049642106  0.56  0.57696
## cost_insurance 0.0001071340 0.0000823863  1.30  0.19347
## young         -0.2318216399 1.3439223932 -0.17  0.86305
## adult          0.6225533757 1.4354506835  0.43  0.66451
## old            NA          NA          NA          NA
## married        2.6820666759 0.7537144293  3.56  0.00037 ***
## college        5.9639748862 2.3073285487  2.58  0.00974 **
## sexratio       -0.9159391203 1.5161018930 -0.60  0.54575
## meanincome     -0.0001205415 0.0000536641 -2.25  0.02469 *
## income25       -0.0000492954 0.0001458362 -0.34  0.73535
## income50       0.0000293917 0.0000164330  1.79  0.07368 .
## income75       0.0000110611 0.0000097409  1.14  0.25616
## cost_wage      0.0001705942 0.0002028866  0.84  0.40044
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

# Supermarket Entry

```
# Predict on test data
p = predict(logitfit,supermarket.te,type="response")
logitpred = as.factor(p > 0.5)
table(logitpred,supermarket.te$entry)

##
## logitpred FALSE TRUE
##      FALSE    424   193
##      TRUE     146   237

logiterr <- 1-mean(logitpred==supermarket.te$entry)
logiterr

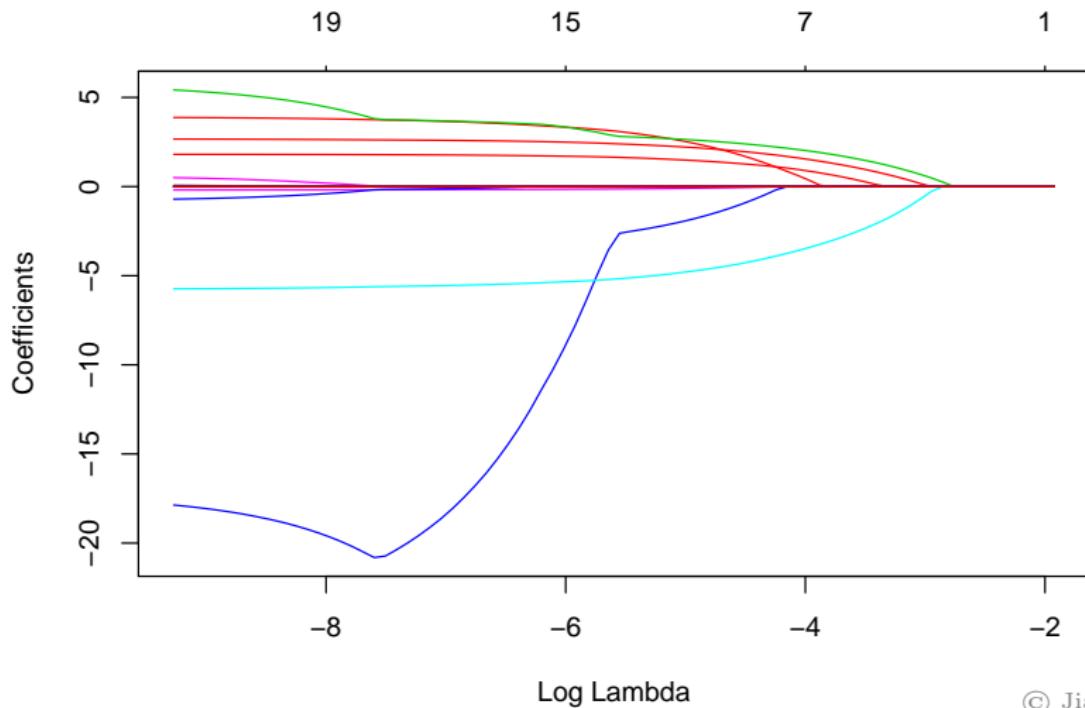
## [1] 0.34
```

# Supermarket Entry

```
#####
# Logistic Regression: Lasso #
#####
# Fit on training data
require(glmnet)
x <- model.matrix(entry ~.,supermarket.tr)[,-1] #no intercept
y <- supermarket.tr$entry
lassofit.all <- glmnet(x,y,alpha=1,family="binomial")
cv.lasso <- cv.glmnet(x,y,alpha=1,family="binomial") #cross-validation
lambda.star <- cv.lasso$lambda.min #lambda associated with min
#
#                                     cross-validated error
# Alternatively: lambda.star <- cv.lasso$lambda.1se
# now fit Lasso with optimal lambda on the entire training data set
lassofit.star <- glmnet(x,y,alpha=1,lambda=lambda.star,family="binomial")
```

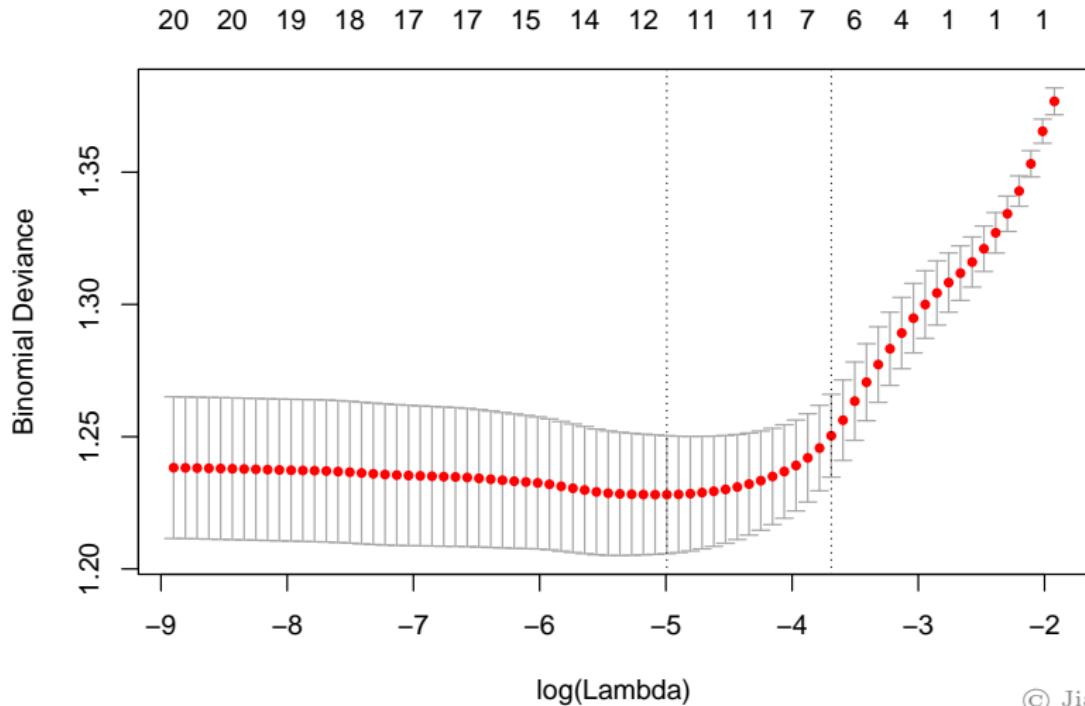
# Supermarket Entry

```
plot(lassofit.all,xvar="lambda")
```



# Supermarket Entry

```
plot(cv.lasso)
```



# Supermarket Entry

```
coef(lassofit.star)

## 22 x 1 sparse Matrix of class "dgCMatrix"
##           s0
## (Intercept) -0.913802643
## pop          0.000000057
## density      1.484797053
## gas          .
## unemp         -1.909544384
## tax           -4.789096976
## union         -0.098232058
## housing       -0.023352708
## crime         2.515947491
## publicTrans   0.001273307
## cost_insurance 0.000054071
## young         .
## adult         .
## old           .
## married        2.212484918
## college        2.638237373
## sexratio       .
## meanincome    .
## income25      .
## income50      .
## income75      .
## cost_wage     .
```

# Supermarket Entry

```
# Predict on test data
newx <- model.matrix(entry ~.,supermarket.te)[,-1]
lassopred <- predict(lassofit.star,newx,type="class")
table(lassopred,supermarket.te$entry)

##
## lassopred FALSE TRUE
##      FALSE    440   188
##      TRUE     130   242

lassoerr <- 1-mean(lassopred==supermarket.te$entry)
lassoerr

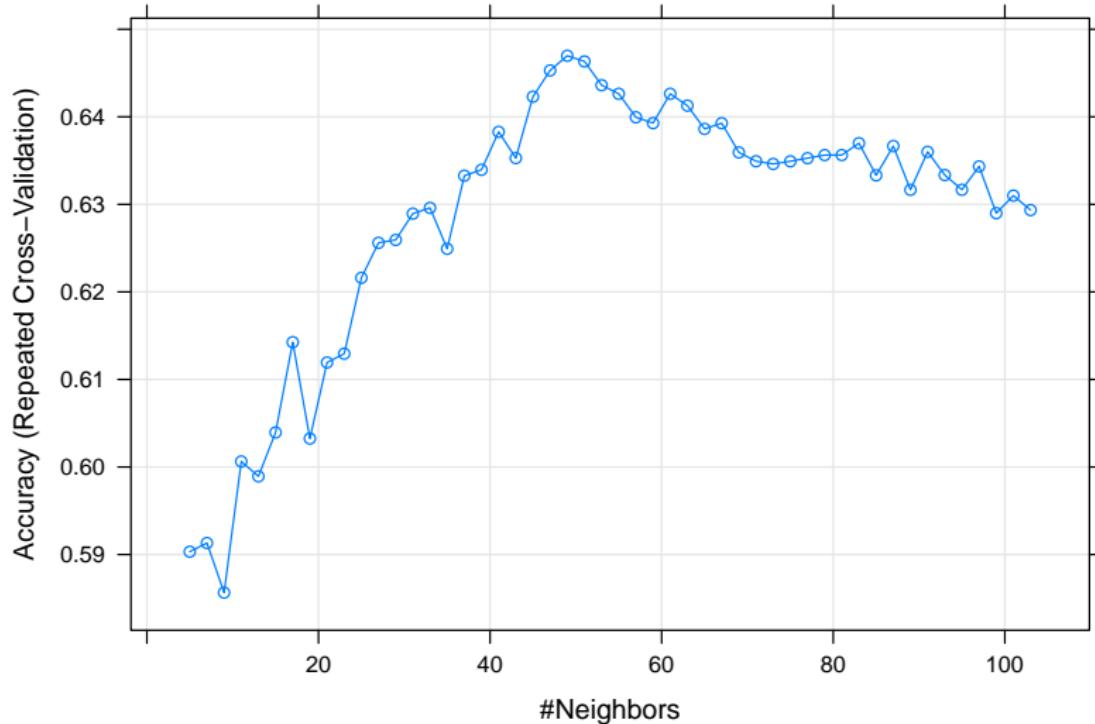
## [1] 0.32
```

# Supermarket Entry

```
#####
# KNN #
#####

# Fit on training data
# use (repeated) cross-validation to choose optimal K
require("caret")
knnfit <- train(entry ~., data=supermarket.tr, method="knn",
                 trControl=trainControl(method="repeatedcv",repeats=3),
                 preProcess=c("center","scale"),tuneLength = 50)
plot(knnfit)
```

# Supermarket Entry



# Supermarket Entry

```
# Predict on test data
knnpred <- predict(knnfit,supermarket.te)
table(knnpred,supermarket.te$entry)

##
## knnpred FALSE TRUE
##   FALSE    426   227
##   TRUE     144   203

knnerr <- 1-mean(knnpred==supermarket.te$entry)
knnerr

## [1] 0.37
```

## The one standard error rule for choosing $\lambda$

Let  $CV(\lambda)$  denote the cross-validation error at  $\lambda$ . Let  $SE(\lambda)$  be its standard error<sup>10</sup>. Let  $\lambda_{min} = \arg \min_{\lambda} CV(\lambda)$ .

Instead of choosing  $\lambda_{min}$ , the **one standard error rule** chooses the largest  $\lambda$  whose  $CV(\lambda)$  is within one standard error of  $CV(\lambda_{min})$ . i.e.,  $\lambda^{1se}$  is the largest  $\lambda$  such that

$$CV(\lambda^{1se}) \leq CV(\lambda_{min}) + SE(\lambda_{min})$$

Idea: “All else equal (up to one standard error), go for the simpler (more regularized) model”

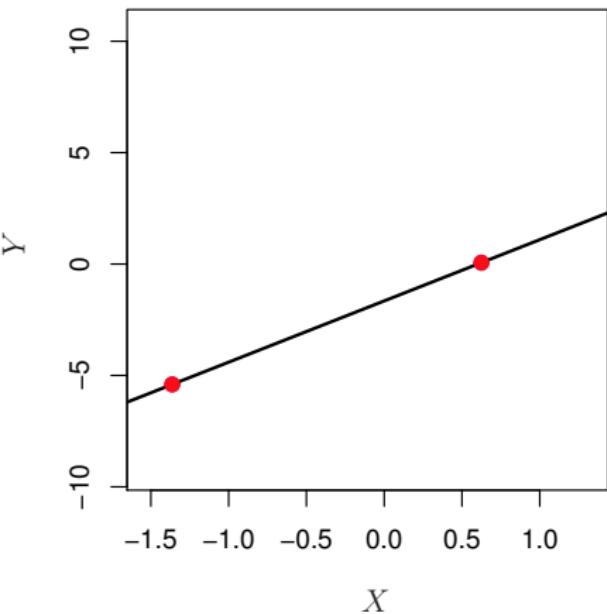
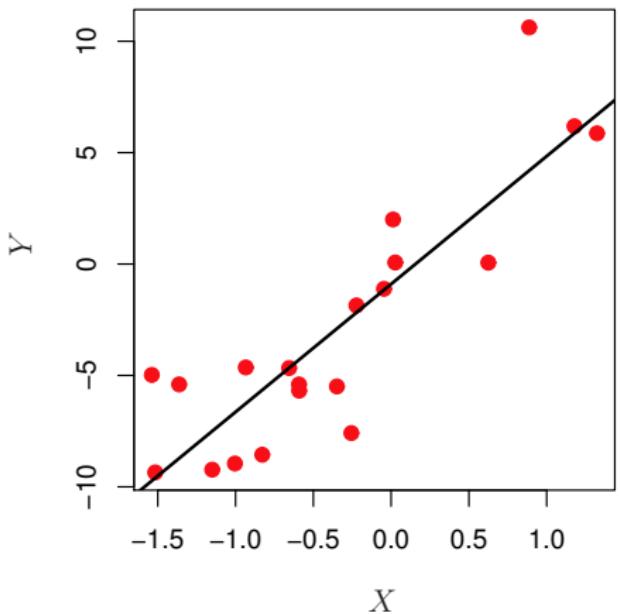
---

<sup>10</sup>i.e.,  $SE(\lambda)$  is the standard deviation of  $CV(\lambda)$  generated by repeated cross-validations.

# Considerations in High Dimensions

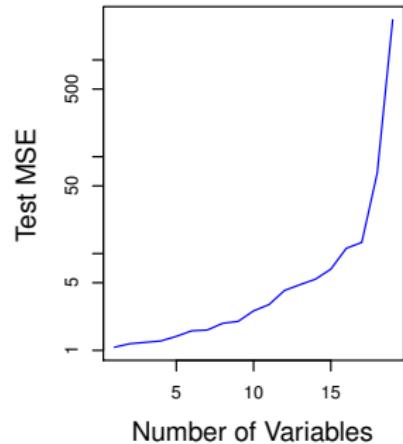
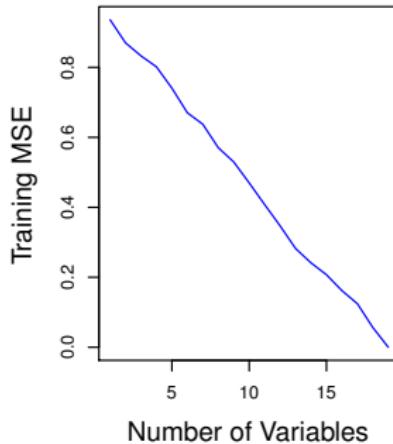
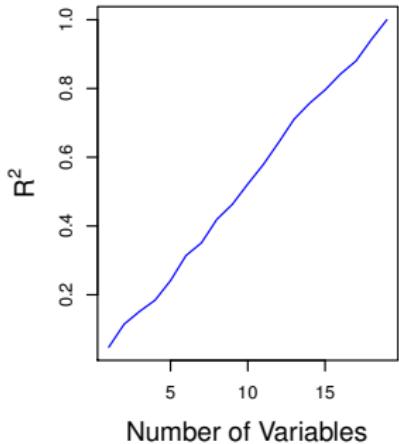
- Most traditional statistical techniques for regression and classification are intended for the **low-dimensional** settings in which  $N \gg p$ .
- Problems with  $p > N$  are often referred to as **high-dimensional**. Classical approaches such as least squares linear regression are not appropriate in this setting.

# Considerations in High Dimensions



Left: Least squares regression in the low-dimensional setting. Right: Least squares regression with  $N = 2$  observations and two parameters to be estimated.

# Considerations in High Dimensions



$R^2$ , training MSE, and test MSE on a sample with  $N = 20$  observations, as more (unrelated) features are added to the model.

# Considerations in High Dimensions

- In high dimensions, even linear models are *too flexible*.
- Information criteria such as AIC and BIC are *not* appropriate.
- When  $p$  is large, multicollinearity problem can be extreme.
  - ▶ This means that we can never know exactly which variables (if any) truly are predictive of the outcome, and 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.

# Regularization

- Shrinkage is a form of **regularization**: methods that **constrain** the complexity of learning algorithms in order to avoid overfitting and improve out-of-sample error.
- In least squares regression, given a model  $\mathcal{H}$ , we choose  $g \in \mathcal{H}$  to minimize the in-sample error  $E_{in}$ . This is called **empirical risk minimization (ERM)**.
- Regularization methods choose  $g$  by solving the following problem:

$$\min_{h \in \mathcal{H}} E_{in}(h) \quad (12)$$

subject to  $\Omega(h) \leq C$

, where  $\Omega(h)$  is a measure of the complexity of  $h$  and is called a **regularizer**<sup>11</sup>.

---

<sup>11</sup>When  $\mathcal{H}$  is parametrized by  $\beta$ ,  $\Omega(h)$  can be written as  $\Omega(\beta)$ .

# Regularization

- Equivalently, (12) can be written as

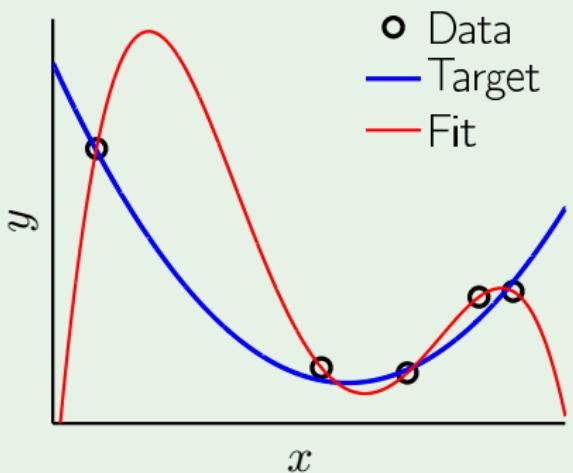
$$\begin{aligned} & \min_{h \in \mathcal{H}} \{E_{in}(h) + \lambda \Omega(h)\} \\ &= \min_{h \in \mathcal{H}} E_{aug}(h) \end{aligned} \tag{13}$$

, where we define **augmented error**  $E_{aug}(h) \equiv E_{in}(h) + \lambda \Omega(h)$ .

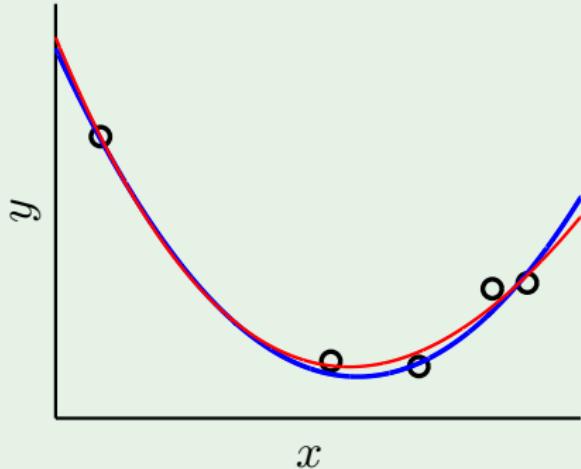
- (12) is called the **Ivanov form** and expresses regularization as constrained ERM.
- (13) is called the **Tikhonov form** and expresses regularization as penalized ERM.

# Regularization

## Regularization as Constrained ERM



free fit

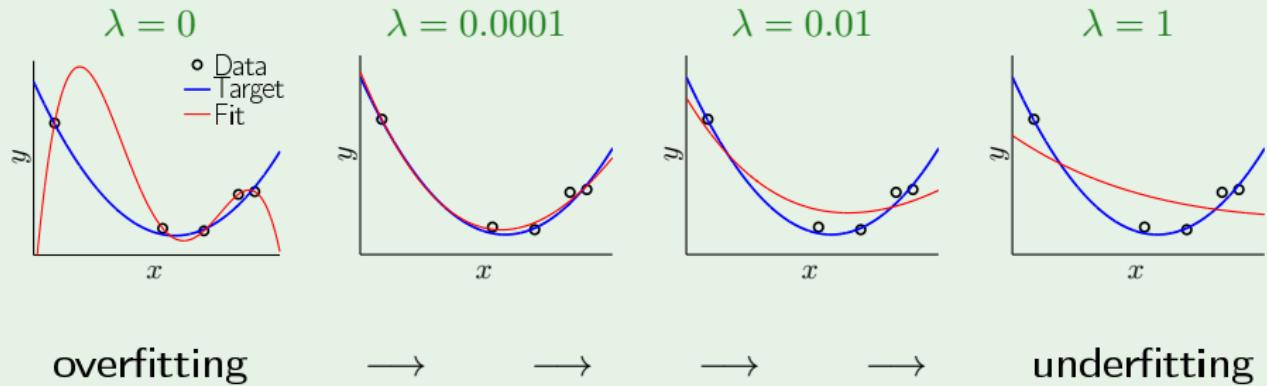


restrained fit

# Regularization

## Regularization as Constrained ERM

Minimizing  $E_{in}(\beta) + \lambda \|\beta\|_2^2$  for different  $\lambda$ s



# Regularization

## Bias-Variance Decomposition

$$f : [-1, 1] \rightarrow \mathbb{R} \quad f(x) = \sin(\pi x)$$

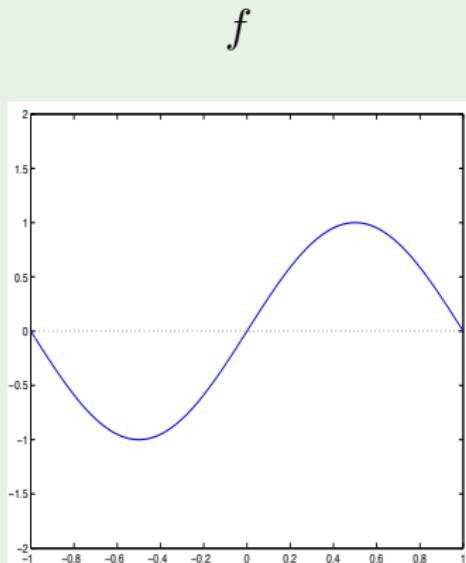
Only two training examples!  $N = 2$

Two models used for learning:

$$\mathcal{H}_0: h(x) = b$$

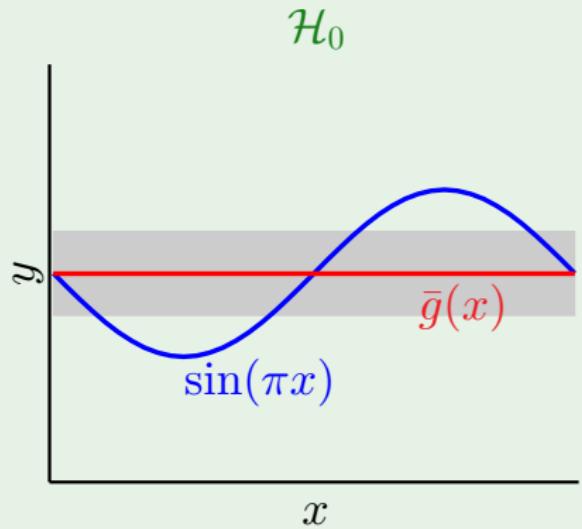
$$\mathcal{H}_1: h(x) = ax + b$$

Which is better,  $\mathcal{H}_0$  or  $\mathcal{H}_1$ ?



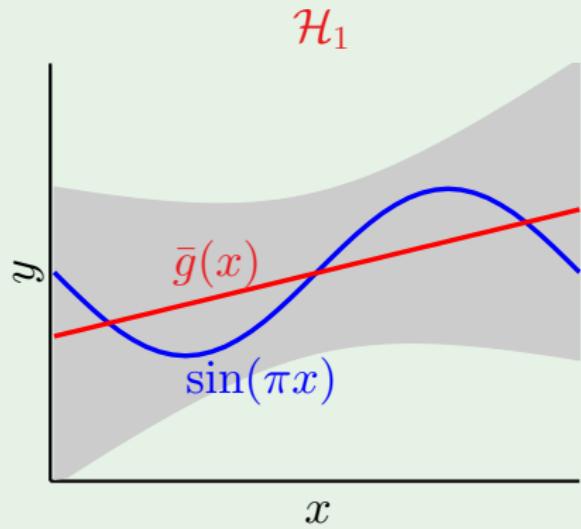
# Regularization

## Bias-Variance Decomposition



bias = **0.50**

var = **0.25**

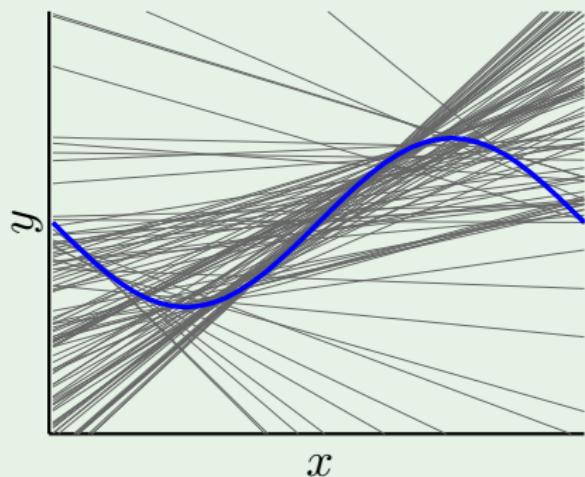


bias = **0.21**

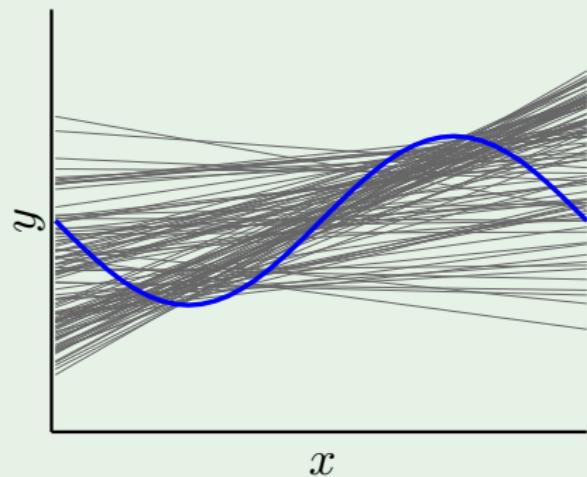
var = **1.69**

# Regularization

## Bias-Variance Decomposition



without regularization

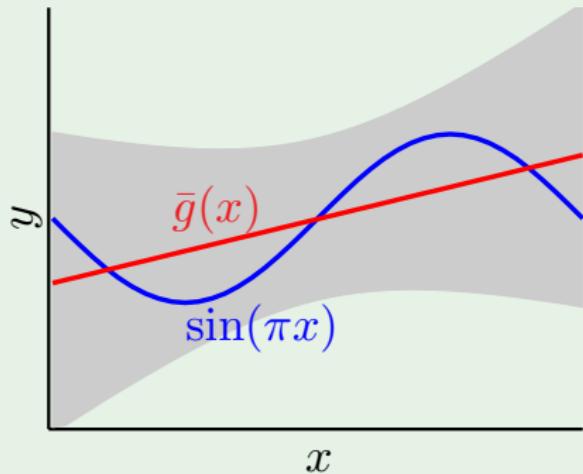


with regularization

# Regularization

## Bias-Variance Decomposition

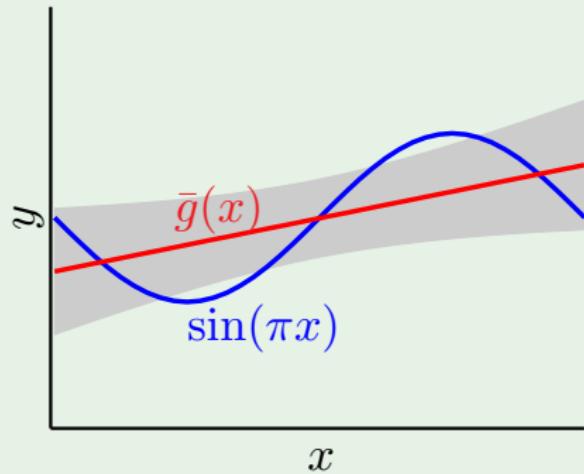
without regularization



bias = **0.21**

var = **1.69**

with regularization



bias = **0.23**

var = **0.33**

# Regularization

$E_{aug}$  is a better proxy for  $E_{out}$  than  $E_{in}$ .

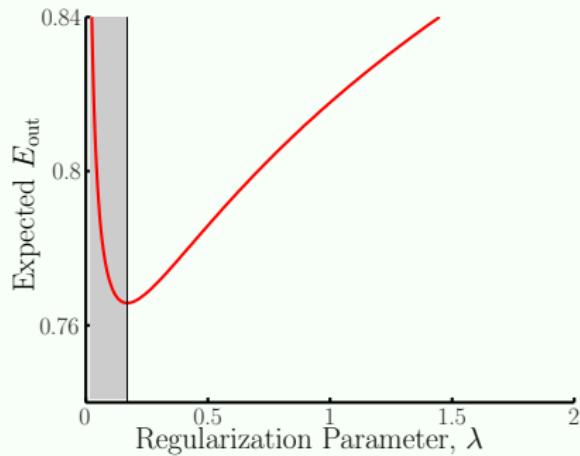
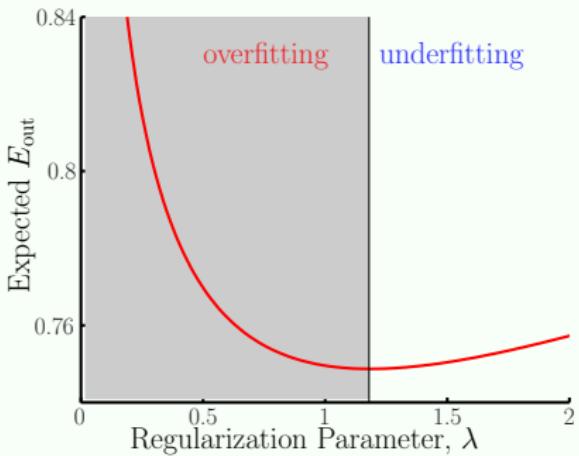
VC generalization bound:

$$E_{out}(h) \leq E_{in}(h) + \mathcal{O}\left(\sqrt{d_{VC}(\mathcal{H}) \frac{\ln N}{N}}\right)$$

Augmented error:

$$E_{aug}(h) = E_{in}(h) + \lambda \Omega(h)$$

# Choosing a Regularizer



Simulated data set with the target function being a 15<sup>th</sup> degree polynomial.

$$\text{Left: } \Omega(\beta) = \sum_{j=0}^{15} \beta_j^2; \text{ Right: } \Omega(\beta) = \sum_{j=0}^{15} j\beta_j^2$$

# Choosing a Regularizer

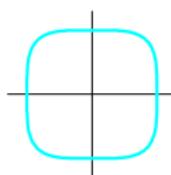
$L_0$  regularizer:  $\Omega(\beta) = \|\beta\|_0 = \sum_j \mathcal{I}\{\beta_j \neq 0\}$

$L_1$  regularizer:  $\Omega(\beta) = \|\beta\|_1 = \sum_j |\beta_j|$

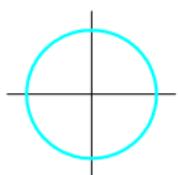
$L_2$  regularizer:  $\Omega(\beta) = \|\beta\|_2^2 = \sum_j \beta_j^2$

$L_q$  regularizer:  $\Omega(\beta) = \|\beta\|_q^q = \sum_j |\beta_j|^q$

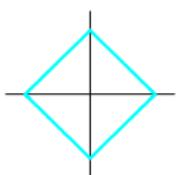
$$q = 4$$



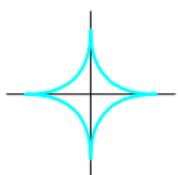
$$q = 2$$



$$q = 1$$



$$q = 0.5$$



$$q = 0.1$$



Contours of constant value of  $L_q$  regularizer for different values of  $q$

# Elastic Net

In the presence of highly correlated features, the results of the lasso can be arbitrary and unstable.

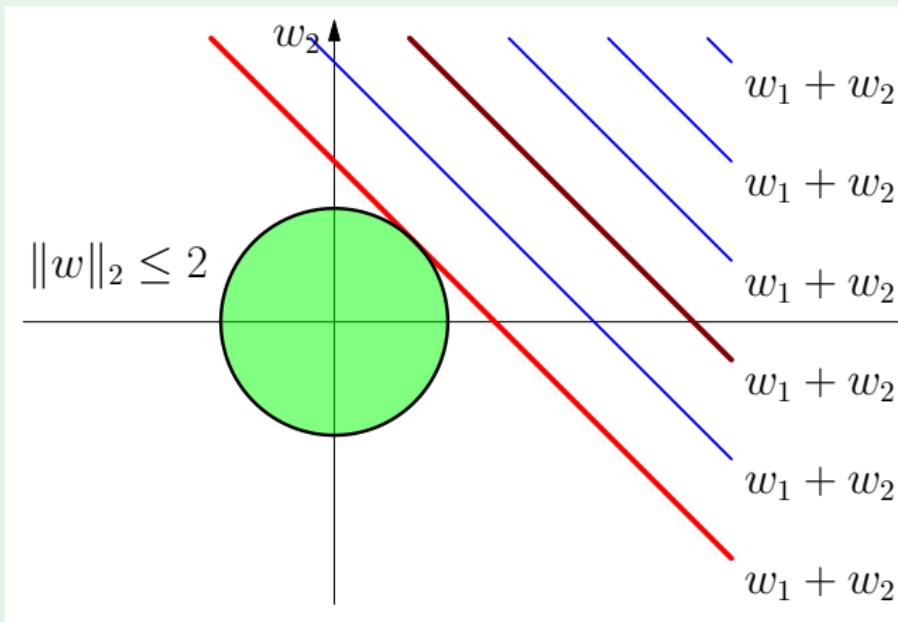
- Target function:  $y = f(x) = 4x_1$ ; Model:  $h(x) = w_1x_1 + w_2x_2$
- Suppose  $x_1 = x_2$ , then any  $h(x)$  with  $w_1 + w_2 = 4$  minimizes  $E_{in}$ , but different regularizers penalize them differently:

$w_1$	$w_2$	$\ w\ _1$	$\ w\ _2^2$
4	0	4	16
2	2	4	8
1	3	4	10
-1	5	6	26

- $\|w\|_1$  does not discriminate (as long as all have same sign)
- $\|w\|_2^2$  minimized when  $w$  is spread equally

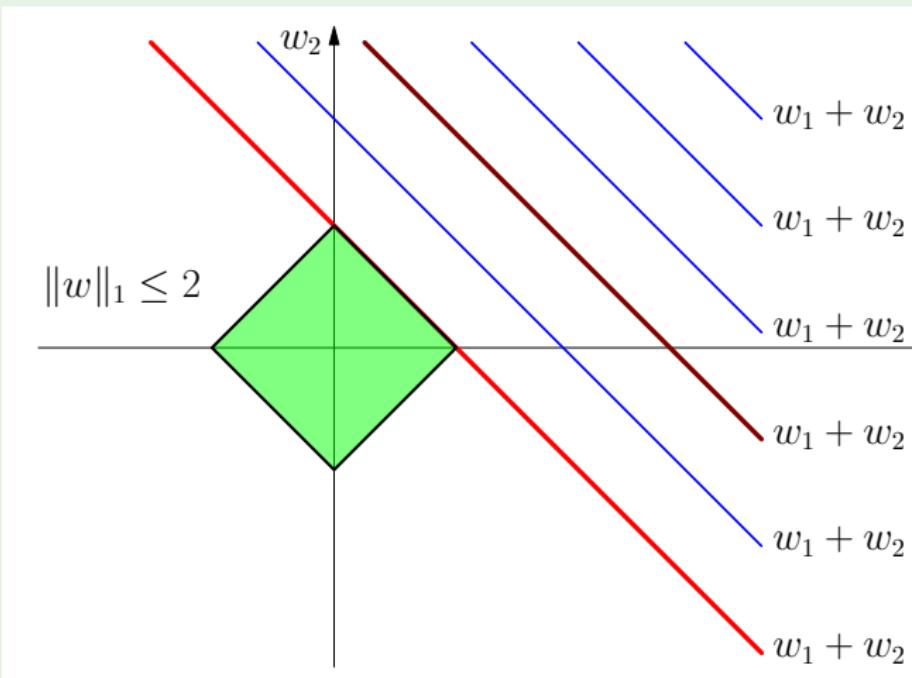
# Elastic Net

## Correlated Features (cont.)



# Elastic Net

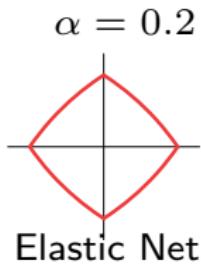
## Correlated Features (cont.)



# Elastic Net

The **elastic net** combines the lasso and ridge penalties:

$$\Omega(\beta) = \alpha \|\beta\|_1 + (1 - \alpha) \|\beta\|_2^2$$



- The elastic net shrinks together the coefficients of correlated features like ridge regression, while encouraging sparse solutions like the lasso.

# Elastic Net

Simulation:

$$z_1, z_2, e, \epsilon_1, \dots, \epsilon_6 \sim^{i.i.d.} \mathcal{N}(0, 1)$$

$$y \sim 3z_1 - 1.5z_2 + 2e$$

$$x_j = \begin{cases} z_1 + 0.2\epsilon_j & j = 1, 2, 3 \\ z_2 + 0.2\epsilon_j & j = 4, 5, 6 \end{cases}$$

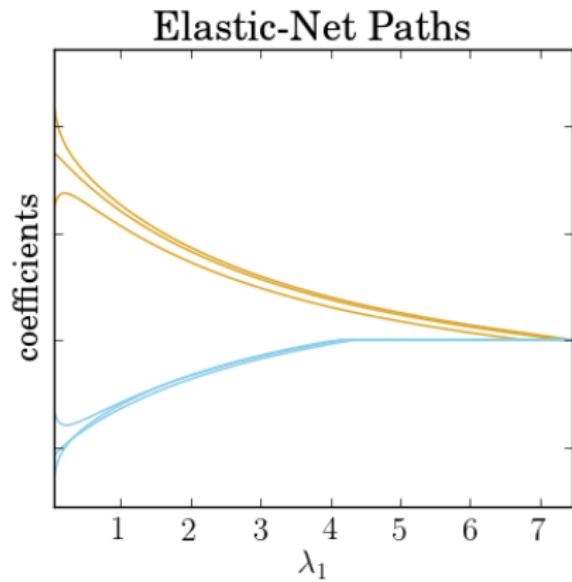
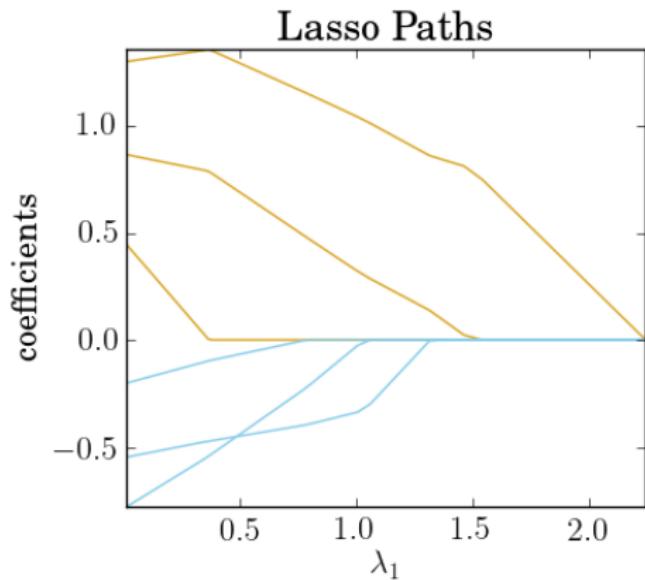
# Elastic Net

```
# generate data
require(MASS)
require(glmnet)
n = 100
z1 = rnorm(n)
z2 = rnorm(n)
e = mvtnorm(n, mu=rep(0,6), Sigma=diag(6))
x1 = z1 + e[,1:3]/5
x2 = z2 + e[,4:6]/5
x = cbind(x1,x2)
y = 3*z1 - 1.5*z2 + 2*rnorm(n)

# lasso
lassofit = glmnet(x,y,alpha=1)

# elastic net
enetfit = glmnet(x,y,alpha=0.3) #30% L1 and 70%L2
```

# Elastic Net



# Bayesian Interpretation of Regularization

The lasso and ridge regression estimates correspond to the *maximum a posteriori (MAP)* estimates with Laplace and Gaussian priors<sup>12</sup>.

- Ridge regression is the posterior mode<sup>13</sup> for  $\beta$  under a normal prior with mean zero.
- The lasso is the posterior mode for  $\beta$  under a double-exponential (Laplace) prior with mean zero.
- The standard deviation of the prior distribution corresponds to regularization strength ( $\lambda$ ).

---

<sup>12</sup>See Appendix I for a simple derivation

<sup>13</sup>as well as the posterior mean

# Bayesian Interpretation of Regularization

The **maximum à posteriori (MAP)** estimate is defined as

$$\begin{aligned}\hat{\theta}_{MAP} &= \arg \max_{\theta} p(\theta | \mathcal{D}) \\ &= \arg \max_{\theta} p(\mathcal{D} | \theta) p(\theta)\end{aligned}$$

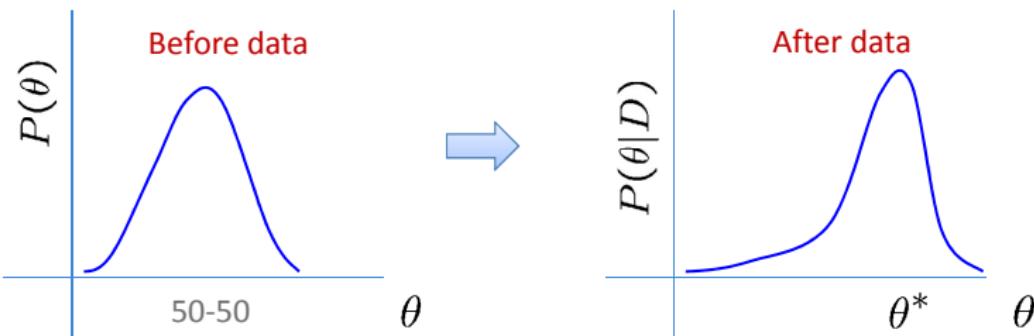
, i.e.,  $\hat{\theta}_{MAP}$  is the **mode** of the posterior distribution of  $\theta$ .

- Given a uniform prior  $p(\theta) = \text{constant}$ ,  $\hat{\theta}_{MAP} = \hat{\theta}_{MLE}$ <sup>14</sup>.

---

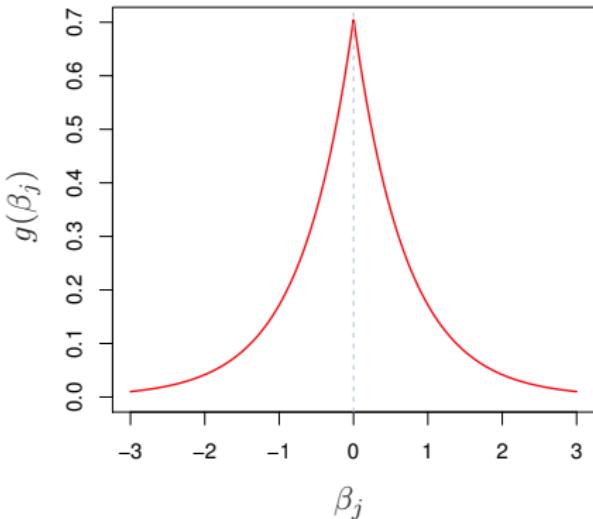
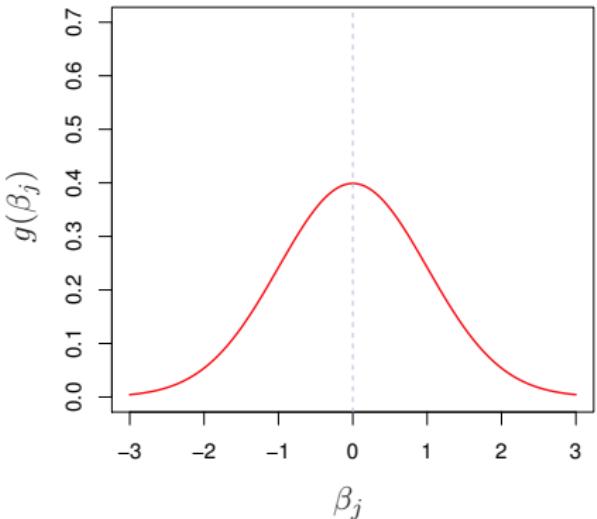
<sup>14</sup>In this case, the Bayesian estimate equals the frequentist estimate because the prior is *uninformative* (hence uniform priors are sometimes called **uninformative priors**). As a result,  $\mathcal{D}$  becomes the only source of our knowledge for  $\theta$ , as in the frequentist approach. If  $p(\theta) = \text{constant}$  over the entire real line, then it is also an **improper prior**, since it does not integrate to one.

# Bayesian Interpretation of Regularization



The Bayesian Way

# Bayesian Interpretation of Regularization



Left: Gaussian prior for Ridge; Right: Laplace prior for the Lasso

# Bayesian Interpretation of Regularization

- From a Bayesian perspective, regularization amounts to the use of *informative priors*, where we introduce our knowledge or belief about the target function in the form of priors, and use them to “regulate” the behavior of the hypothesis we choose<sup>15</sup>.

---

<sup>15</sup> Recall that in choosing a model, our choice should be based both on our belief about the true target function *and* on the amount of data we have (matching model complexity to data sources). Similarly, the choice of regularizers is based both on our belief about the target function (such as its smoothness), and the necessity of tempering down the variability in our estimates when  $p$  is large relative to  $N$ .

# Smoothing Splines

A smoothing spline is a function  $g(x)$  that solves the following problem:

$$\min_h \left\{ \sum_{i=1}^N (y_i - h(x_i))^2 + \lambda \int [h''(t)]^2 dt \right\} \quad (14)$$

, where  $\lambda \geq 0$  is a tuning parameter.

- The regularizer  $\Omega(h) = \int [h''(t)]^2 dt$  is a **roughness** penalty. Hence  $\lambda$  regulates the **smoothness** of the spline.
- $\lambda = 0$  :  $g(x)$  *interpolates* every  $(x_i, y_i)$ .
- $\lambda \rightarrow \infty$  :  $g(x)$  becomes the linear least squares line.

# Smoothing Splines

The solution to (14) is a natural cubic spline with knots at every  $x_i^{16, 17, 18}$ .

Recall that the least squares solution for natural cubic splines is  $\hat{\beta}^{LS} = (\Phi' \Phi)^{-1} \Phi' Y$ , where  $\Phi(x)$  is the set of basis functions for natural cubic splines. Solving (14) gives the *shrinkage* solution:

$$\hat{\beta} = (\Phi' \Phi + \lambda \Psi)^{-1} \Phi' Y$$

, where  $\Psi$  is an  $N \times N$  matrix with  $\Psi_{ij} = \int \Phi_i''(t) \Phi_j''(t) dt$ .

---

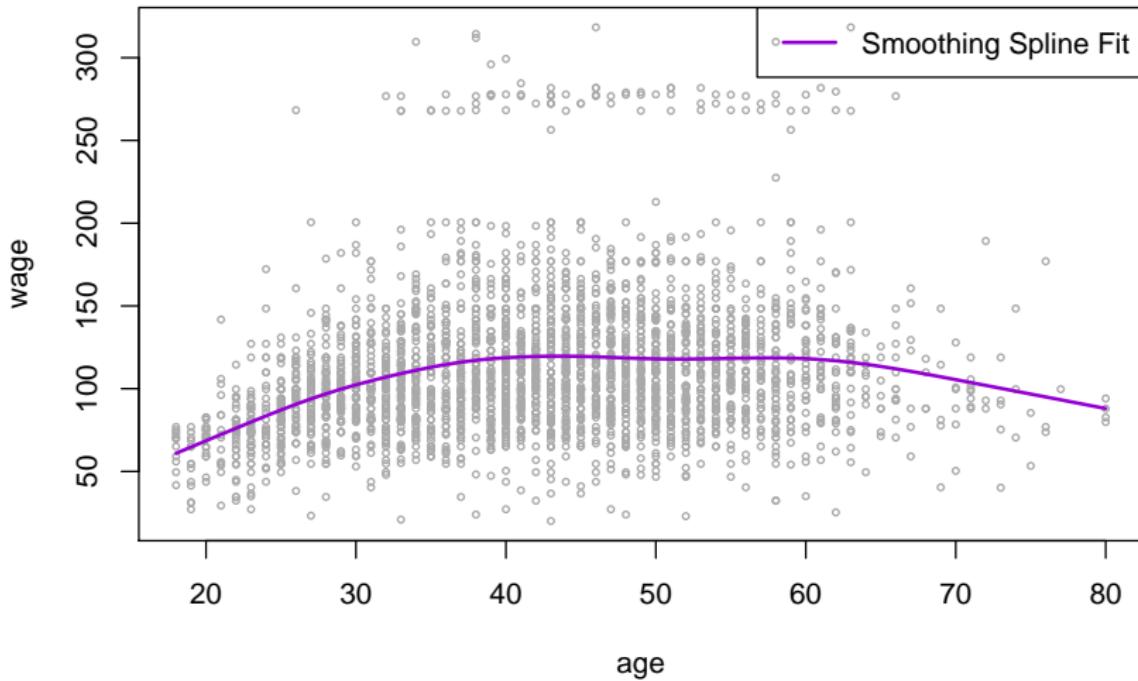
<sup>16</sup>Hence, we no longer have to choose knot placement!

<sup>17</sup>See Appendix II for proof.

<sup>18</sup>More generally, penalizing the squared  $k^{th}$  derivative leads to a natural spline of degree  $2k - 1$ .

# Wage Profile

```
fit = smooth.spline(Wage$age,Wage$wage)
```



## Appendix I: Gaussian Prior and L2 Regularization

Consider a normal linear model with Gaussian prior:

$$\begin{aligned}\beta_j &\sim \mathcal{N}(0, \tau^2), \quad j = 0, \dots, p \\ y &\sim \mathcal{N}(\beta' x, \sigma^2)\end{aligned}$$

Then we have:

$$\hat{\theta}_{MAP} = \arg \max_{\beta} \left\{ \left( \prod_{i=1}^N \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(y_i - \beta' x_i)^2}{2\sigma^2}} \right) \left( \prod_{j=1}^p \frac{1}{\tau \sqrt{2\pi}} e^{-\frac{\beta_j^2}{2\tau^2}} \right) \right\} \quad (15)$$

## Appendix I: Gaussian Prior and L2 Regularization

(15)  $\Rightarrow$

$$\begin{aligned}\hat{\theta}_{MAP} &= \arg \max_{\beta} \left\{ \log \left( \prod_{i=1}^N \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(y_i - \beta' x_i)^2}{2\sigma^2}} \right) + \log \left( \prod_{j=1}^p \frac{1}{\tau \sqrt{2\pi}} e^{-\frac{\beta_j^2}{2\tau^2}} \right) \right\} \\ &= \arg \max_{\beta} \left\{ - \sum_{i=1}^N \frac{(y_i - \beta' x_i)^2}{2\sigma^2} - \sum_{j=0}^p \frac{\beta_j^2}{2\tau^2} \right\} \\ &= \arg \min_{\beta} \left\{ \sum_{i=1}^N (y_i - \beta' x_i)^2 + \lambda \sum_{j=0}^p \beta_j^2 \right\}\end{aligned}$$

, where  $\lambda = \frac{\sigma^2}{\tau^2}$  <sup>19</sup>.

---

<sup>19</sup> Hence smaller values of  $\tau$  correspond to stronger regularization

## Appendix I: Laplace Prior and L1 Regularization

Consider a normal linear model with Laplace prior:

$$\begin{aligned}\beta_j &\sim \text{Laplace}(\mu, b), \quad j = 0, \dots, p \\ y &\sim \mathcal{N}(\beta' x, \sigma^2)\end{aligned}$$

Then we have:

$$\begin{aligned}\hat{\theta}_{MAP} &= \arg \max_{\beta} \left\{ \left( \prod_{i=1}^N \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(y_i - \beta' x_i)^2}{2\sigma^2}} \right) \left( \prod_{j=1}^p \frac{1}{2b} e^{-\frac{|\beta_j|}{2b}} \right) \right\} \\ &= \arg \min_{\beta} \left\{ \sum_{i=1}^N (y_i - \beta' x_i)^2 + \lambda \sum_{j=0}^p |\beta_j| \right\}\end{aligned}$$

, where  $\lambda = \frac{\sigma^2}{b}$ <sup>20</sup>.

---

<sup>20</sup>Smaller values of  $b$  correspond to stronger regularization

## Appendix II: Smoothing Splines

### Theorem

Let  $h$  be any differentiable function on  $[a, b]$  for which  $h(x_i) = z_i$  for  $i = 1, \dots, n$ . Suppose  $n > 2$ , and that  $g$  is the natural cubic spline interpolant to the values  $z_1, \dots, z_n$  at points  $x_1, \dots, x_n$  with  $a < x_1 < \dots < x_n < b$ . Then  $\int (h'')^2 \geq \int (g'')^2$  with equality if and only if  $h = g$ .

## Appendix II: Smoothing Splines

### Proof.

Let  $\ell = h - g$ . So  $\ell(x_i) = 0$  for  $i = 1, \dots, n$ .

$$\int (h'')^2 = \int (g'' + \ell'')^2 = \int (g'')^2 + 2 \int g'' \ell'' + \int (\ell'')^2 \quad (16)$$

Since

$$\begin{aligned} \int_a^b g''(x) \ell''(x) dx &= g''(x) \ell'(x) \Big|_a^b - \int_a^b \ell'(x) g^{(3)}(x) dx \\ &= - \sum_{i=1}^{n-1} g^{(3)}(x_j^+) \int_{x_j}^{x_{j+1}} \ell'(x) dx \\ &= - \sum_{i=1}^{n-1} g^{(3)}(x_j^+) (\ell(x_{j+1}) - \ell(x_j)) = 0 \end{aligned}$$

$$, \text{ (16)} \Rightarrow \int (h'')^2 \geq \int (g'')^2.$$



# Acknowledgement I

Part of this lecture is adapted from the following sources:

- Abu-Mostafa, Y. S., M. Magdon-Ismail, and H. Lin. 2012. *Learning from Data*. AMLBook.
- Blei, D. M. *Interacting with data*. Lecture at Princeton University, retrieved on 2017.01.01. [[link](#)]
- Dave, R. *Advanced Scientific Computing*, Lecture at Harvard Institute for Applied Computational Science, retrieved on 2019.01.01. [[link](#)]
- Hastie, T., R. Tibshirani, and J. Friedman. 2008. *The Elements of Statistical Learning* (2<sup>nd</sup> ed.). Springer.
- James, G., D. Witten, T. Hastie, and R. Tibshirani. 2013. *An Introduction to Statistical Learning: with Applications in R*. Springer.

## Acknowledgement II

- Liang, F. *Statistical Learning: Nonparametric and Regularization Approaches*, Lecture at Duke University, retrieved on 2019.01.01. [[link](#)]
- Rosenberg, D. S. *Machine Learning and Computational Statistics*, Lecture at NYU Center for Data Science, retrieved on 2019.01.01. [[link](#)]
- Tibshirani, R. *Data Mining*, Lecture at Carnegie Mellon University, retrieved on 2017.01.01. [[link](#)]