# ML01 – Introduction to Machine Learning
## Model Selection

Thierry Denœux

tdenoeux@utc.fr
https://www.hds.utc.fr/~tdenoeux

Université de technologie de Compiègne

Spring 2021

# Overview

# Need for model selection

- Consider, for instance, a regression problem with a response variable $Y$ and 3 predictors $X_1, X_2, X_3$.
- We can consider many (an infinity of) models, such as

$$Y = \beta_0 + \beta_1 X_1 + \epsilon$$

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \epsilon$$

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_1^2 + \beta_5 X_2^2 +$$
$$\beta_6 X_1 X_2 + \beta_7 X_3^2 + \beta_8 X_1 X_3 + \beta_9 X_2 X_3 + \epsilon$$

$$\vdots$$

Which model to choose?

# Bias-variance trade-off

- We have seen that a more complex model will not always have a smaller error when applied to test data.
- This is due to the bias-variance trade-off: when the number of parameters increases, the bias of the model decreases, but the variance increases.
- Furthermore, a simpler model often has a distinct advantages in terms of its interpretability.
- In this chapter, we discuss some tools to select models that will be
  - Complex enough to fit the data, but
  - Not too complex to avoid overfitting and to be interpretable.
- We focus mainly on linear regression, but the tools can be adapted to classification.

## Three classes of methods

Subset Selection. We identify a subset of the $p$ predictors that we believe to be related to the response. We then fit a model using the reduced set of variables.

Regularization. We fit a model involving all $p$ predictors, but the estimated coefficients are shrunken towards zero to obtain a smoother prediction function. This regularization (also known as shrinkage) has the effect of reducing variance and can also perform variable selection.

Dimension Reduction. We project the $p$ predictors into a $q$-dimensional subspace, where $q < p$. This is achieved by computing $q$ different linear combinations of the variables. Then these $q$ new variables are used as predictors to fit a linear model. One such technique will be studied later in this course.

# Overview

# Overview

# Best subset selection

1. Let $\mathcal{M}_0$ denote the null model, which contains no predictors. This model simply predicts the sample mean for each observation.

2. For $k = 1, 2, \ldots, p$:
   1. Fit all $\binom{p}{k} = \frac{p!}{(p-k)!k!}$ models that contain exactly $k$ predictors.
   2. Pick the best among these $\binom{p}{k}$ models, and call it $\mathcal{M}_k$. Here "best" is defined as having the smallest RSS, or equivalently the largest $R^2$.

3. Select a single best model from among $\mathcal{M}_0, \ldots, \mathcal{M}_p$. (How? to be seen later).
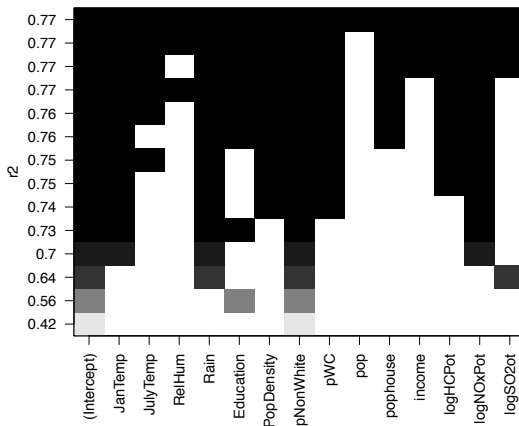
# Example: air pollution and mortality

- Data are from McDonald and Schwing (1973), "Instabilities of Regression Estimates Relating Air Pollution to Mortality", *Technometrics*, 15, 463-481.
- This data set of 15 predictors and a measure of mortality on 60 US metropolitan areas in 1959-1961.

# Variables

- Response: Total Age Adjusted Mortality Rate
- Predictors:
  1. Mean annual precipitation in inches
  2. Mean January temperature in degrees Fahrenheit
  3. Mean July temperature in degrees Fahrenheit
  4. Percent of 1960 SMSA population that is 65 years of age or over
  5. Population per household, 1960 SMSA
  6. Median school years completed for those over 25 in 1960 SMSA
  7. Percent of housing units that are found with facilities
  8. Population per square mile in urbanized area in 1960
  9. % of 1960 urbanized area population that is non-white
  10. % employment in white-collar occupations in 1960 urbanized area
  11. % of families with income under 3,000 in 1960 urbanized area
  12. Relative population potential of hydrocarbons, HC
  13. Relative pollution potential of oxides of nitrogen, NOx
  14. Relative pollution potential of sulfur dioxide, SO2
  15. Percent relative humidity, annual average at 1 p.m.

# Best subset selection in R

```
library('leaps')
reg.fit<-regsubsets(Mortality~.-logNOx,data=pollution,method='exhaustive',nvmax=15)
plot(reg.fit,scale="r2")
```

# Extension to other models

- Although we have presented best subset selection here for least squares regression, the same ideas apply to other types of models, such as logistic regression.
- The deviance, $-2\ell(\widehat{\theta})$, plays the role of RSS for a broader class of models.

# Stepwise selection

- For computational reasons, best subset selection cannot be applied with very large p.

- Best subset selection may also suffer from statistical problems when $p$ is large: larger the search space, the higher the chance of finding models that look good on the training data, even though they might not have any predictive power on future data.

- Thus an enormous search space can lead to overfitting and high variance of the coefficient estimates.

- For both of these reasons, stepwise methods, which explore a far more restricted set of models, are attractive alternatives to best subset selection.

# Forward stepwise selection

- Forward stepwise selection begins with a model containing no predictors, and then adds predictors to the model, one-at-a-time, until all of the predictors are in the model.

- In particular, at each step the variable that gives the greatest additional improvement to the fit is added to the model.
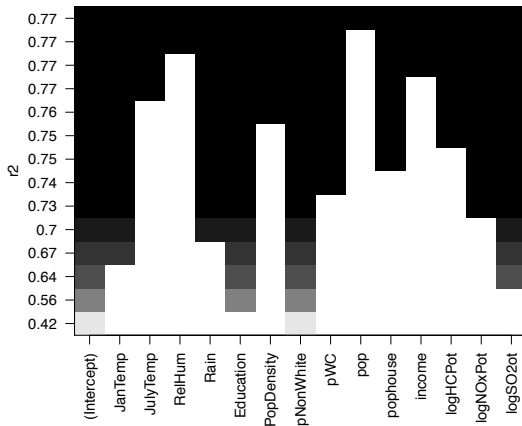
# Forward stepwise selection in detail

1. Let $\mathcal{M}_0$ denote the null model, which contains no predictors.
2. For $k = 0, \ldots, p - 1$:
   1. Consider all $p - k$ models that augment the predictors in $\mathcal{M}_k$ with one additional predictor.
   2. Choose the best among these $p - k$ models, and call it $\mathcal{M}_{k+1}$. Here "best" is defined as having highest $R^2$.
3. Select a single best model from among $\mathcal{M}_0, \ldots, \mathcal{M}_p$.

# More on forward stepwise selection

- Computational advantage over best subset selection is clear.
- It is not guaranteed to find the best possible model out of all $2^p$ models containing subsets of the $p$ predictors.

# Forward stepwise selection in R

```
reg.fit<-regsubsets(Mortality~.-logNOx,data=pollution,method='forward',nvmax=15)
plot(reg.fit,scale="r2")
```

# Backward stepwise selection

- Like forward stepwise selection, backward stepwise selection provides an efficient alternative to best subset selection.
- However, unlike forward stepwise selection, it begins with the full least squares model containing all $p$ predictors, and then iteratively removes the least useful predictor, one-at-a-time.

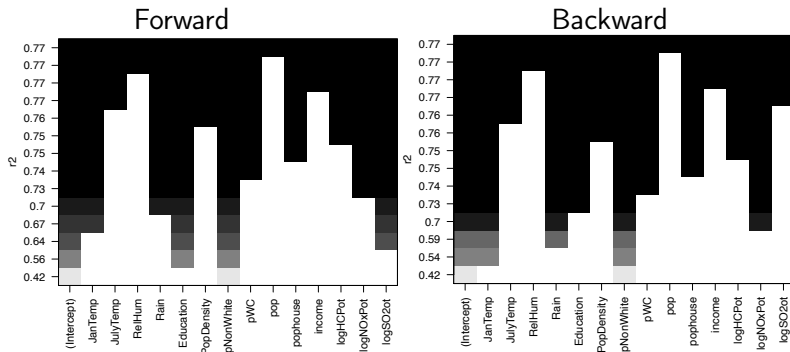# Backward stepwise selection in detail

1. Let $\mathcal{M}_p$ denote the full model, which contains all $p$ predictors.
2. For $k = p, p - 1, \ldots, 1$:
   1. Consider all $k$ models that contain all but one of the predictors in $\mathcal{M}_k$, for a total of $k - 1$ predictors.
   2. Choose the best among these $k$ models, and call it $\mathcal{M}_{k-1}$. Here "best" is defined as having highest $R^2$.
3. Select a single best model from among $\mathcal{M}_0, \ldots, \mathcal{M}_p$.

# More on backward stepwise selection

- Like forward stepwise selection, the backward selection approach searches through only $1 + p(p+1)/2$ models, and so can be applied in settings where $p$ is too large to apply best subset selection

- Like forward stepwise selection, backward stepwise selection is not guaranteed to yield the best model containing a subset of the $p$ predictors.

- 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). In contrast, forward stepwise can be used even when $n < p$, and so is the only viable subset method when $p$ is very large.

# Backward stepwise selection in R

```
reg.fit<-regsubsets(Mortality~.-logNOx,data=pollution,method='backward',nvmax=15)
plot(reg.fit,scale="r2")
```

# Overview

# Choosing the optimal model

- The model containing all of the predictors will always have the smallest RSS and the largest $R^2$, since these quantities are related to the training error.

- We wish to choose a model with low test error, not a model with low training error. Recall that training error is usually a poor estimate of test error.

- Therefore, RSS and $R^2$ are not suitable for selecting the best model among a collection of models with different numbers of predictors.
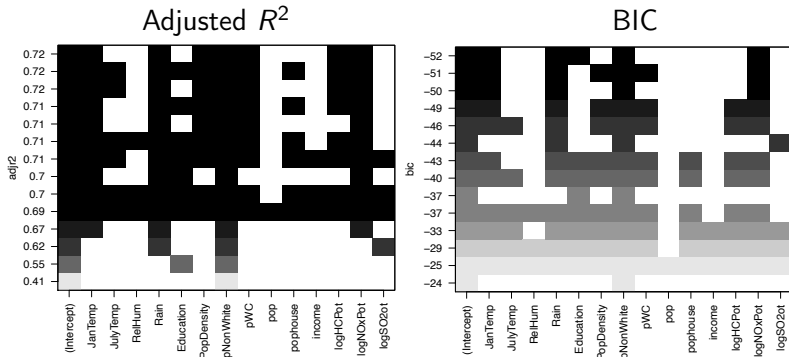
# Estimating test error: two approaches

- We can
    1. Indirectly estimate test error by making an adjustment to the training error to account for the bias due to overfitting, or
    2. Directly estimate the test error, using either a hold-out approach or a cross-validation approach.
- We illustrate both approaches next.

# Training error adjustment techniques

- These techniques adjust the training error for the model size, and can be used to select among a set of models with different numbers of variables.
- Three criteria:
  1. Adjusted $R^2$
  2. Akaike information criterion (AIC)
  3. Bayesian information criterion (BIC)
- The next figure displays BIC, and adjusted $R^2$ for the best model of each size produced by best subset selection on the pollution data set.

# Example

```
reg.fit<-regsubsets(Mortality~.-logNOx,data=pollution,method='exhaustive',nvmax=15)
plot(reg.fit,scale="adjr2") plot(reg.fit,scale="bic")
```



Adjusted $R^2$                    BIC

# Adjusted $R$-squared

- Idea: introduce the "population $R^2$" as

$$R^2_{pop} = 1 - \frac{\sigma^2}{\text{Var}(Y)}$$

- The usual $R^2$ is

$$R^2 = 1 - \frac{RSS/n}{TSS/n}$$

  It is based on biased estimates of the residual and total variances.

- The adjusted $R^2$ is based on unbiased estimates:

$$\boxed{\overline{R}^2 = 1 - \frac{RSS/(n-p-1)}{TSS/(n-1)}}$$

  where $p$ is the number of predictors used.

- This criterion is specific to regression.

# AIC

- The AIC criterion is defined for a large class of models fit by maximum likelihood:

$$AIC = -2\ell(\widehat{\theta}) + 2r$$

  where $\ell(\widehat{\theta})$ is the maximized value of the log-likelihood function for the estimated model, and $r$ is the number of parameters.

- The best model has the smallest AIC value.

- For linear regression with $p$ variables and a constant term, $r = p + 1$.

# BIC

- Definition:

$$BIC = -2\ell(\widehat{\theta}) + r\log(n)$$

  where $r$ is the number of parameters.
- Like AIC, BIC will tend to take on a small value for a model with a low test error, and so generally we select the model that has the lowest BIC value.
- Notice that BIC replaces the $2r$ used by AIC with a $r\log(n)$ term, where $n$ is the number of observations.
- 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 AIC.

# Direct estimation of the prediction error

- Each of the subset selection procedures returns a sequence of models $\mathcal{M}_k$ indexed by model size $k = 0, 1, 2, \ldots, p$. Our job here is to select $\widehat{k}$. Once selected, we will return model $\mathcal{M}_{\widehat{k}}$.

- We compute an estimate of the prediction error for each model $\mathcal{M}_k$ under consideration, and then select the $k$ for which the resulting estimated prediction error is smallest.

- This procedure has an advantage relative to AIC, BIC, and adjusted $R^2$, in that it provides a direct estimate of the prediction error.

- It can also be used in a wider range of model selection tasks, even in cases where it is hard to define the model degrees of freedom.

# Direct estimation of the test error

Two methods:

1. Validation-set (hold-out) approach
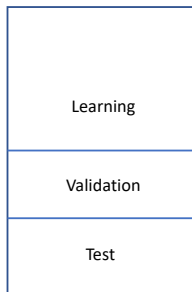2. Cross-validation

# Validation-set approach



- Here we randomly divide the available set of samples into two parts:
  1. a training set and
  2. a validation set

- The model is fit on the training set, and the fitted model is used to predict the responses for the observations in the validation set.

- The resulting validation error provides an estimate of the prediction error. This is typically assessed using MSE in the case of regression and misclassification rate in the case of classification.

# Hold-out approach (continued)

- After the best model has been selected, it is usually fit on the whole data (training+validation).
- The validation error for the best model is biased (optimistic).

| Learning |
|----------|
| Validation |
| Test |

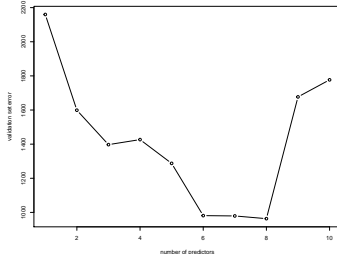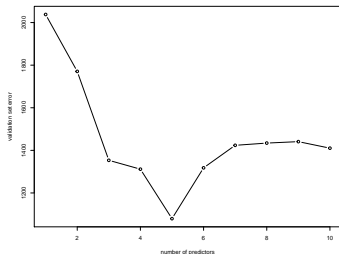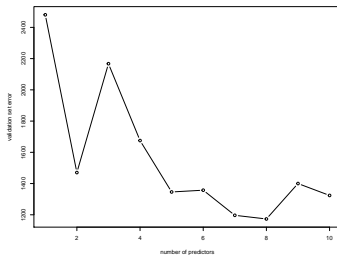- The error of the best model has to be estimated using an independent test set.

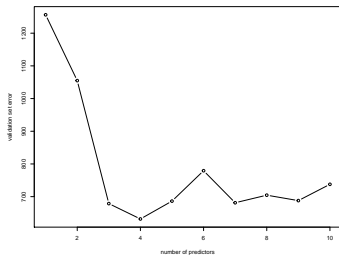# Example

```
n<-nrow(pollution)
napp=45
ntst=n-napp
train<-sample(1:n,napp)
pollution.train<-pollution[train,]
pollution.test<-pollution[-train,]

Formula<-c(Mortality ~ pNonWhite,
Mortality ~ Education + pNonWhite,
Mortality ~ Rain + pNonWhite + logSO2ot,
Mortality ~ JanTemp+ Rain +pNonWhite +logNOxPot,
...
)

for(i in 1:10){
reg<-lm(Formula[[i]],data=pollution.train)
pred<-predict(reg,newdata=pollution.test)
err[i]<-mean((pollution.test$Mortality-pred)^2)
}
```

# Results with 4 different splits (Air pollution data)

# Limitations of the hold-out approach

- The validation estimate of the test error can be highly variable, depending on precisely which observations are included in the training set and which observations are included in the validation set.

- In the hold-out approach, only a subset of the observations – those that are included in the training set rather than in the validation set – are used to fit the model.

- This suggests that the validation-set error may tend to overestimate the test error for the model fit on the entire data set.

# K-fold cross-validation

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

| 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|
| Validation | Train | Train | Train | Train |

# $K$-fold cross-validation in detail

- Let the $K$ subsets be $C_1, C_2, \ldots, C_K$, where $C_k$ denotes the indices of the observations in subset $k$. There are $n_k$ observations in subset $k$: if $n$ is a multiple of $K$, then $n_k = n/K$.

- Compute

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

where

$$\mathsf{MSE}_k = \frac{1}{n_k} \sum_{i \in C_k} \left( y_i - \widehat{y}_i^{(-k)} \right)^2$$

and $\widehat{y}_i^{(-k)}$ is the fit for observation $i$, obtained from the data with subset $k$ removed.

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

# Special case

- With least-squares linear or polynomial regression, a shortcut makes the cost of LOOCV the same as that of a single model fit!
- The following formula holds:

$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^{n} \left( \frac{y_i - \widehat{y}_i}{1 - h_i} \right)^2,$$

where $\widehat{y}_i$ is the $i$th fitted value from the original least squares fit, and $h_i$ is the leverage (diagonal term of the "hat" matrix). (Reminder: $\mathbf{H} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T$, and $\widehat{\mathbf{y}} = \mathbf{H}\mathbf{y}$).

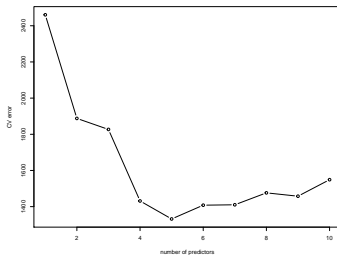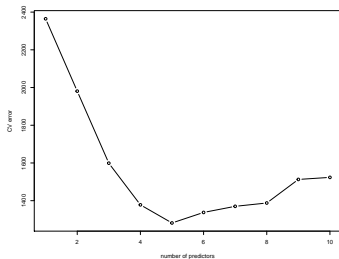- This is like the ordinary MSE, except the $i$-th residual is divided by $1 - h_i$.

# Choice of $K$

- Since each training set is only $(K-1)/K$ as big as the original training set, the estimates of prediction error will typically be biased upward.

- This bias is minimized when $K = n$ (LOOCV), but this estimate has high variance, because the estimates from each fold are highly correlated.

- $K = 5$ or 10 provides a good compromise for this bias-variance tradeoff.

# Example of 10-fold cross-validation

```
K<-10
folds=sample(1:K,n,replace=TRUE)
CV<-rep(0,10)

for(i in (1:10)){
for(k in (1:K)){
reg<-lm(Formula[[i]],data=pollution[folds!=k,])
pred<-predict(reg,newdata=pollution[folds==k,])
CV[i]<-CV[i]+ sum((pollution$Mortality[folds==k]-pred)^2)
}
CV[i]<-CV[i]/n
}
```

# Result (4 trials)

# Final remarks on cross-validation

- The CV error rates can be averaged over $r$ repetitions of $K$-fold cross-validation with different random partitions, to reduce the variance of the CV error estimates.

- After the best model has been selected, we usually re-estimate the model parameters using the whole training set.

- To obtain an unbiased estimate of the best model's error, we need an independent test set.

# Overview

# Shrinkage methods

- By retaining only a subset of the predictors, subset selection produces a model that is interpretable and has possibly lower prediction error than the full model.
- However, because it is a discrete process – variables are either retained or discarded – it often exhibits high variance, and so does not always reduce the prediction error of the full model.
- Shrinkage methods are more continuous, and do not suffer as much from high variability.
- Two main methods:
  1. Ridge regression
  2. Lasso

# Overview

# Ridge regression

- Ridge regression shrinks the regression coefficients by imposing a penalty on their size. The ridge coefficients minimize a penalized residual sum of squares:

$$\widehat{\beta}^{\text{ridge}} = \arg\min_{\beta} \left\{ \sum_{i=1}^{n} \left( y_i - \beta_0 - \sum_{j=1}^{p} x_{ij}\beta_j \right)^2 + \lambda \sum_{j=1}^{p} \beta_j^2 \right\}$$

- Here $\lambda \geq 0$ is a complexity parameter that controls the amount of shrinkage: the larger the value of $\lambda$, the greater the amount of shrinkage. The coefficients are shrunk toward zero (and each other), i.e., to the simplest model (with only the constant term).

- Selecting a good value for $\lambda$ is critical; cross-validation can be used for this.

# Equivalent form

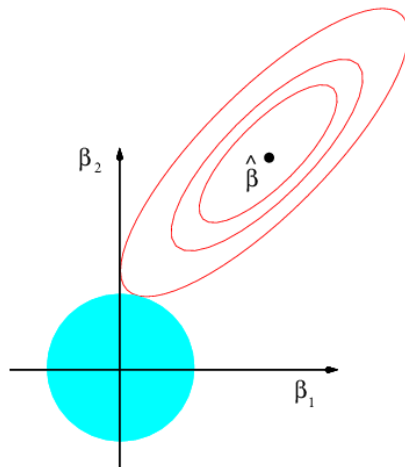- An equivalent way to write the ridge problem is

$$\widehat{\beta}^{\text{ridge}} = \arg\min_{\beta} \left\{ \sum_{i=1}^{n} \left( y_i - \beta_0 - \sum_{j=1}^{p} x_{ij}\beta_j \right)^2 \right\}$$
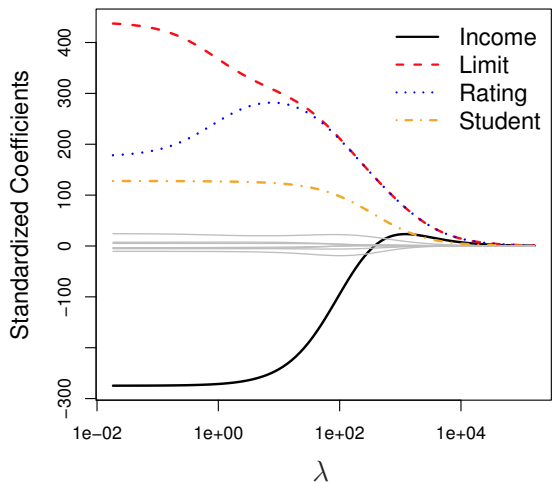
$$\text{subject to } \sum_{j=1}^{p} \beta_j^2 \leq t,$$

which makes explicit the size constraint on the parameters. (See next slide).

- There is a one-to-one correspondence between parameters $t$ and $\lambda$ in the previous formulation.

# Ridge regression as a constrained optimization problem

# The effect of ridge regression

# Derivation of the ridge regression estimates

- We can show that $\widehat{\beta}^{\text{ridge}}$ can be found by separating the minimization problem into two parts, after centering the inputs (replacing $x_{ij}$ by $x_{ij} - \overline{x}_j$):
  1. We estimate $\beta_0$ by $\overline{y} = \frac{1}{n} \sum_{i=1}^{n} y_i$
  2. The remaining coefficients get estimated by a ridge regression without intercept, using the centered $x_{ij}$ and the centered $y_i$.

- We assume that both the inputs and the output have been centered, so that the input matrix $\mathbf{X}$ has $p$ (rather than $p+1$) columns, and $\mathbf{y}$ is the $n$-vector of centered outputs.

- The criterion can be written in matrix form

$$\text{RSS}_\lambda(\beta) = (\mathbf{y} - \mathbf{X}\beta)^T(\mathbf{y} - \mathbf{X}\beta) + \lambda\beta^T\beta.$$

# Derivation of the ridge regression estimates (continued)

- The criterion can be rewritten as

$$\text{RSS}_\lambda(\beta) = \mathbf{y}^T\mathbf{y} - 2\beta^T\mathbf{X}^T\mathbf{y} + \beta^T(\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I}_p)\beta$$

- Differentiating with respect to $\beta$ we obtain

$$\frac{\partial \text{RSS}_\lambda(\beta)}{\partial \beta} = -2\mathbf{X}^T\mathbf{y} + 2(\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I}_p)\beta$$
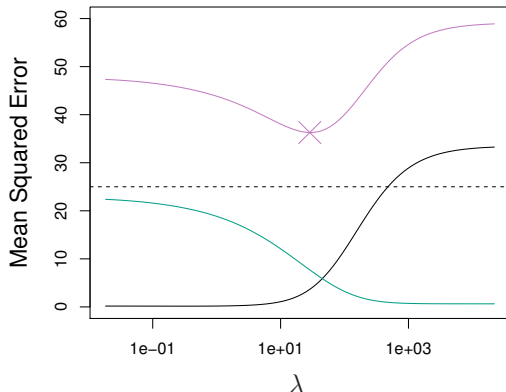
- The solution of the equation $\frac{\partial \text{RSS}_\lambda(\beta)}{\partial \beta} = 0$ is

$$\boxed{\widehat{\beta}^{\text{ridge}} = (\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I}_p)^{-1}\mathbf{X}^T\mathbf{y}}$$

# Ridge regression: scaling of predictors

- The standard least squares coefficient estimates are scale equivariant: multiplying $X_j$ by a constant $c$ simply leads to a scaling of the least squares coefficient estimates by a factor of $1/c$. In other words, regardless of how the $j$th predictor is scaled, $X_j\widehat{\beta}_j$ will remain the same.

- In contrast, the ridge regression coefficient estimates can change substantially when multiplying a given predictor by a constant, due to the sum of squared coefficients term in the penalty part of the ridge regression objective function.

- Therefore, it is best to apply ridge regression after standardizing the predictors (dividing each centered variable by its standard deviation).

# Why does ridge regression improve over least squares?



Simulated data with $n = 50$ observations, $p = 45$ predictors, all having nonzero coefficients. Squared bias (black), variance (green), and test MSE (purple) for the ridge regression predictions, as a function of $\lambda$. The horizontal dashed lines indicate the minimum possible MSE.
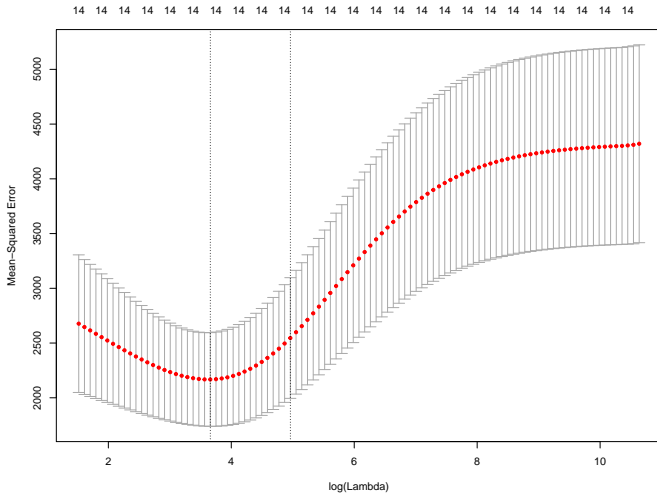
# Ridge regression in R

```
library(glmnet)

x<-model.matrix(Mortality~.-logNOx,pollution)
y<-pollution$Mortality[-21] # obs 21 has 2 missing values
n<-nrow(x)
napp=45
ntst=n-45
train<-sample(1:n,napp)
xapp<-x[train,]
yapp<-y[train]
xtst<-x[-train,]
ytst<-y[-train]

cv.out<-cv.glmnet(xapp,yapp,alpha=0)
plot(cv.out)

fit<-glmnet(xapp,yapp,lambda=cv.out$lambda.min,alpha=0)
ridge.pred<-predict(fit,s=cv.out$lambda.min,newx=xtst)
print(mean((ytst-ridge.pred)^2))
2421.136
```

# CV error as a function of $\lambda$

# Coefficients

```
fit$beta
s0
(Intercept) .
JanTemp -2.641635e-01
JulyTemp 7.231499e-01
RelHum -1.443636e-01
Rain 9.618201e-01
Education -1.154417e+01
PopDensity 2.066547e-03
pNonWhite 1.478269e+00
pWC -1.105875e+00
pop 2.629839e-06
pophouse 3.057905e+01
income -1.008305e-03
logHCPot 2.311552e+00
logNOxPot 6.616369e+00
logSO2ot 3.966114e+00
```

# Overview

# The lasso

- Ridge regression has one obvious disadvantage: unlike subset selection, it includes all $p$ predictors in the final model
- The lasso is a relatively recent alternative to ridge regression that overcomes this disadvantage. The lasso coefficients, $\widehat{\beta}^{\text{lasso}}$ minimize the quantity
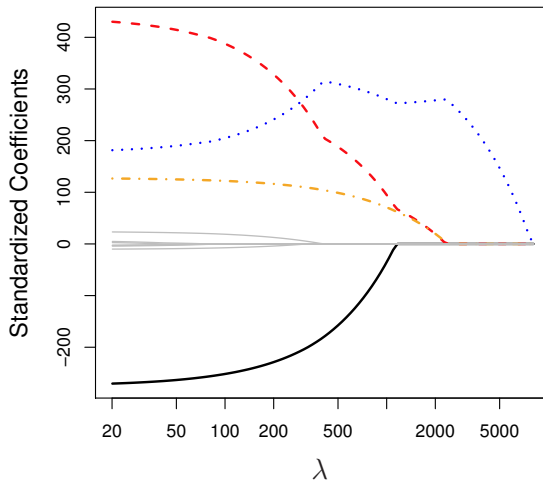
$$\widehat{\beta}^{\text{lasso}} = \arg\min_{\beta} \left\{ \sum_{i=1}^{n} \left( y_i - \beta_0 - \sum_{j=1}^{p} x_{ij}\beta_j \right)^2 + \lambda \sum_{j=1}^{p} |\beta_j| \right\},$$

i.e., the $L_2$ norm is replaced by the $L_1$ norm in the penalty term.

(Reminder: the $L_p$ norm is defined as $\|\beta\|_p = \left( \sum_j |\beta_j|^p \right)^{1/p}$).

# The lasso (continued)

- As with ridge regression, the lasso shrinks the coefficient estimates towards zero.

- However, in the case of the lasso, the $L_1$ penalty has the effect of forcing some of the coefficient estimates to be exactly equal to zero when the tuning parameter $\lambda$ is sufficiently large.

- Hence, much like best subset selection, the lasso performs variable selection.

- We say that the lasso yields sparse models – that is, models that involve only a subset of the variables.

- As in ridge regression, selecting a good value of $\lambda$ for the lasso is critical; cross-validation is again the method of choice.
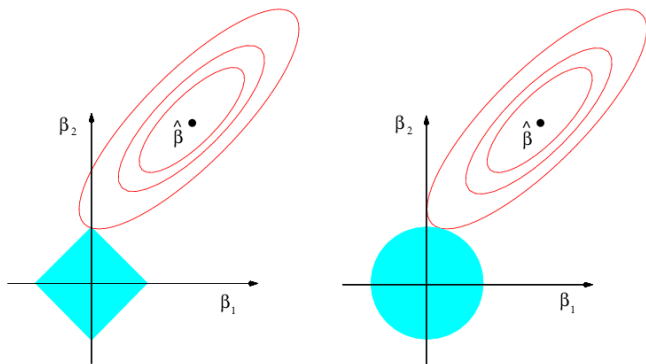
# Example

# Equivalent form

- As in the case of ridge problem, the previous unconstrained optimization problem is equivalent to the following constrained one:

$$\widehat{\beta}^{\text{lasso}} = \arg\min_{\beta} \left\{ \sum_{i=1}^{n} \left( y_i - \beta_0 - \sum_{j=1}^{p} x_{ij}\beta_j \right)^2 \right\}$$

$$\text{subject to } \sum_{j=1}^{p} |\beta_j| \leq t,$$

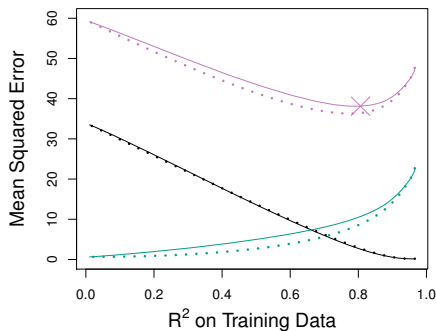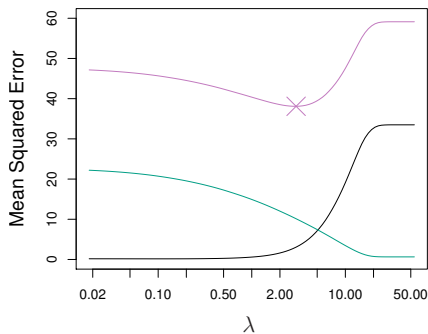- This problem can be solved using a quadratic programming algorithm.

# Why does the lasso eliminate variables?



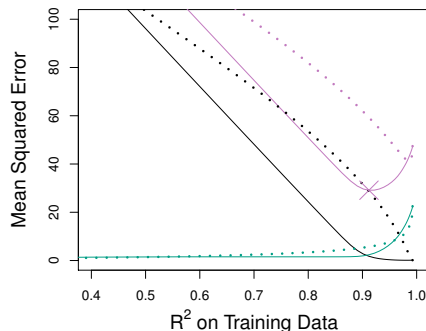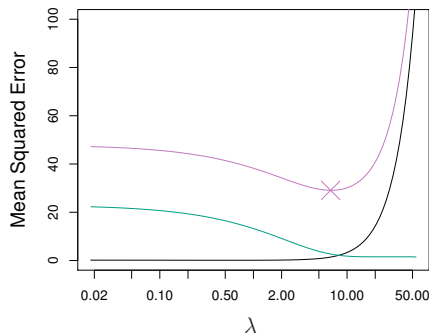When $p = 2$, the feasibility region is a diamond, which has corners; if the solution occurs at a corner, then it has one parameter $\beta_j$ equal to zero. When p > 2, the feasibility region has many corners, flat edges and faces there are many more opportunities for the estimated parameters to be zero.

# Comparing the lasso and ridge regression



Left: Plots of squared bias (black), variance (green), and test MSE (purple) for the lasso on simulated data set of Slide 54. Right: Comparison of squared bias, variance and test MSE between lasso (solid) and ridge (dashed). Both are plotted against their $R^2$ on the training data, as a common form of indexing.

# Comparing the lasso and ridge regression (continued)



Left: Plots of squared bias (black), variance (green), and test MSE (purple) for the lasso. The simulated data are similar to those in the previous slide, except that now only two predictors are related to the response. Right: Comparison of squared bias, variance and test MSE between lasso (solid) and ridge (dashed). Both are plotted against their $R^2$ on the training data, as a common form of indexing.
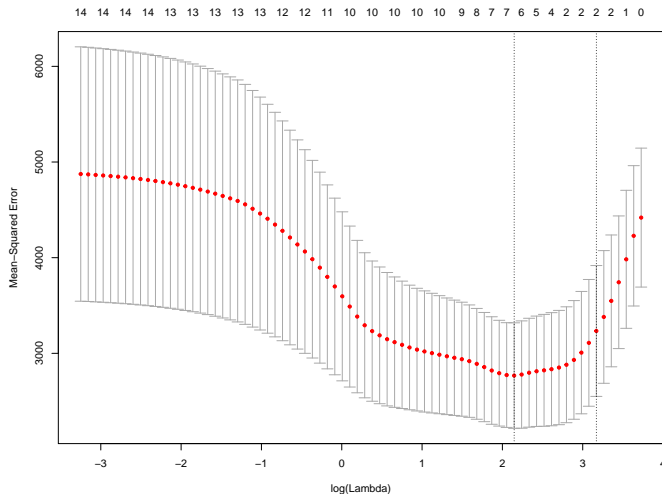
# The lasso in R

```
cv.out<-cv.glmnet(xapp,yapp,alpha=1)
plot(cv.out)

fit.lasso<-glmnet(xapp,yapp,lambda=cv.out$lambda.min,alpha=1)

lasso.pred<-predict(fit.lasso,s=cv.out$lambda.min,newx=xtst)
print(mean((ytst-lasso.pred)^2))
1946.667
```

# CV error as a function of $\lambda$ (lasso)

# Coefficients

```
> print(fit.lasso$beta)
s0
(Intercept) .
JanTemp -1.157095e+00
JulyTemp .
RelHum .
Rain 1.404239e+00
Education -1.796084e+01
PopDensity .
pNonWhite 2.880287e+00
pWC -9.421496e-01
pop 2.141275e-06
pophouse .
income -4.655832e-04
logHCPot .
logNOxPot 1.392387e+01
logSO2ot 3.461564e-01
```