

MODEL SELECTION

-INTRODUCTION TO DATA SCIENCE-

ISL 6.1

Lecturer: Darren Homrighausen, PhD

Preamble:

- Armed with our risk estimators, we need ways to sift through the possible models
- Available techniques vary with the **absolute** and **relative** sizes of n and p
- Like most statistical techniques, model selection comes down to optimization...

Brief optimization and convexity detour

OPTIMIZATION

An optimization problem (program) can be generally formulated as

$$\begin{array}{ll}\text{minimize} & F(\beta) \\ \text{subject to} & f_j(\beta) \leq 0 \text{ for } j = 1, \dots, m\end{array}$$

Here

$\beta = (\beta_1, \dots, \beta_p)^\top$ are the **parameters**

$F : \mathbb{R}^p \rightarrow \mathbb{R}$ is the **objective function**

$f_j : \mathbb{R}^p \rightarrow \mathbb{R}$ are **constraint functions**

The **optimal solution** β^* is such that $F(\beta^*) \leq F(\beta)$ for any β^*, β that satisfies the constraints

CONVEXITY

The main dichotomy of optimization programs is **convex** vs. **nonconvex**

A **convex** program is one in which the objective and constraint functions are all convex

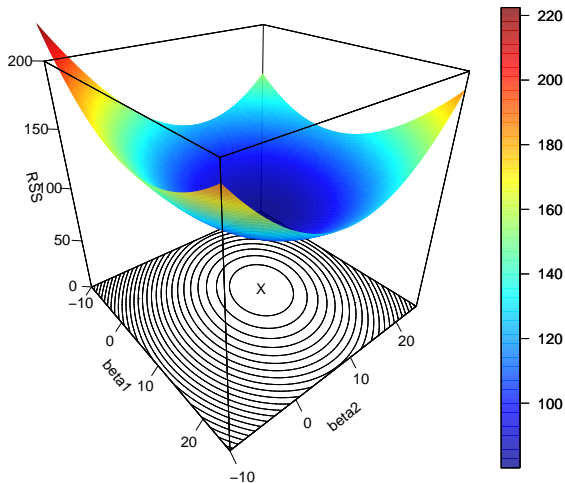
$$f(t\beta + (1 - t)\beta') \leq tf(\beta) + (1 - t)f(\beta') \quad \text{for any } t \in [0, 1]$$

This can be thought of (for smooth enough f)

$$f(\beta') \geq f(\beta) + (\nabla f|_{\beta})^{\top}(\beta' - \beta)$$

Intuition: This means that the function values at a point β' are **above** the supporting hyperplane given by the tangent space at **any** point β

CONVEXITY EXAMPLE



With $RSS = \|Y - \mathbb{X}\beta\|_2^2$ for $p=2$

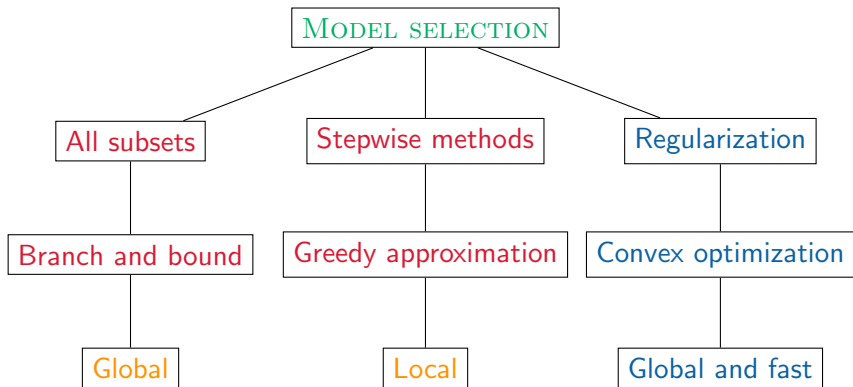
CONVEXITY

Methods for convex optimization programs are (roughly) always **global** and **fast**

For general nonconvex problems, we have to give up one of these:

- Local optimization methods that are fast, but need not find global solution
(So called **greedy** approximations)
- Global optimization methods that find global solutions, but are not always fast (indeed, are often slow)
(Usually exhaustive search type approaches)

Model selection



Some comments:

Non convex programs

Can be seen as a convex relaxation of the nonconvex program
giving all subsets

ALL SUBSETS REGRESSION

First, identify all considered features and transformations and put them in the feature matrix $\mathbf{X} \in \mathbb{R}^{n \times p}$

BEST SUBSET SELECTION ALGORITHM: For $k = 1, \dots, p$



1. Find \hat{R} for the $\binom{p}{k}$ models of size k
2. Save the model that minimizes \hat{R}

(This can be found in Algorithm 6.1 in ISL)

Now, report the model that minimizes a version of GIC over these p models

(Such as AIC, BIC, Mallows C_p , ...)

In general, this is a nonconvex problem, though some shortcuts can be taken

( A general idea known as “Branch and Bound”, see [branchBound.pdf](#) )

ALL SUBSETS REGRESSION IN R

We can use the function `regsubsets` in the package `leaps`

The syntax and associated objects look like:

```
allsubsets.out = regsubsets(Y~.,data=X,nvmax=min(c(n,p)))  
#or  
allsubsets.out = regsubsets(x=X,y=Y,nvmax=min(c(n,p)))
```

- The `nvmax = min(c(n,p))` controls the max size of models considered. The default is 8 and that is usually far too small.
- Now, we can pick among the `min(c(n,p))` models that minimize \hat{R} for a given model size using BIC or Cp

ALL SUBSETS REGRESSION: A BIG PROBLEM (LITERALLY)

If there are p features then there are 2^p possible models

If $p = 40$ (which is considered a modest problem these days), then the number of possible models is

$$2^{40} \approx 1,099,512,000,000 \Rightarrow \text{More than 1 trillion!}$$

If $p = 265$, then the number of possible models is more than the number of atoms in the universe¹

We must sift through the models in a computationally feasible way

¹It is estimated there are 10^{80} atoms in the universe.

Greedy approximations

FORWARD SELECTION

In the likely event that 2^p is too large to be searched over exhaustively, a common **greedy** approximation is the following

1. Find $\text{GIC}(\emptyset)$: The GIC of the intercept only model
2. Search over all p single feature models, computing GIC for each one. Say including X_j minimizes GIC with a value $\text{GIC}(X_j)$. If $\text{GIC}(X_j) < \text{GIC}(\emptyset)$, add X_j to the model and continue. Otherwise terminate
3. Now search over all $p - 1$ models that contain X_j and find the $X_{j'}$ that minimizes GIC. If $\text{GIC}(X_j, X_{j'}) < \text{GIC}(X_j)$, add $X_{j'}$ to the model and continue. Otherwise terminate
4. ...

(See Algorithm 6.2 in ISL)

FORWARD SELECTION

```
regsubsets(x=X,y=Y,nvmax=min(c(n,p)),method='forward')
```

PROS:

- This approach can be used effectively in either the **Big Data** or **High Dimensional** regimes
- It tends to produce sensible answers that are not too different from all-subsets

CONS:

- Can get trapped in a poor local minimum

GENERAL STEPWISE SELECTION

This algorithm can be adapted to..

- start with the full model and stepwise remove features. This is known as **backward selection**

```
regsubsets(x=X,y=Y,nvmax=min(c(n,p)),method='backward')
```

(useful if the full model isn't too large and a superset of the important features is desired)

- consider both adding and removing features at each step. This is known as **stepwise selection**

```
regsubsets(x=X,y=Y,nvmax=min(c(n,p)),method='stepAIC')
```


EXAMPLE: FORWARD SELECTION IN R

```
nvmax      = min(c(n,p))
out        = regsubsets(x=X,y=Y,nvmax=nvmax,method='forward')
out.sum     = summary(out)
out.model  = out.sum$which[which.min(out.sum$bic),]
out.S      = out.model[-1]  #get rid of the intercept entry
out.lm.for = lm(Y~X[,out.S])#regsubsets only scores models,
                                #not fit them

betaHat.for = coef(out.lm.for)
betaHat.for
```

(Note: the `coef` function in `leaps` (that is, using `coef(out,id=which.min(out.sum$bic))`) is unreliable in some cases. Hence, I always refit the model using `lm` and get the `coef` from the `lm` object)

We refer to `out.S` as the **active set**, it is notated $\mathcal{S} \subseteq \{1, 2, \dots, p\}$

Note that \mathcal{S} is a function of Y

IMPORTANT COMMENTS

After using any of these model selection approaches, we produce estimates $\hat{\beta}$ and predictions $\hat{Y} = \hat{\beta}^\top X_{\text{select}}$ where X_{select} includes only the selected features

This can be interpreted as these features are most important for predicting Y from the features included in $X \in \mathbb{R}^p$

(The usual caveats apply: linearity (correlation), there are surely some important coefficients left out/unimportant ones included)

If we run `out = lm(Y ~ Xselect)`, then `summary(out)` will produce the usual significance tests:

→ these are not valid after model selection

IMPORTANT COMMENTS

- If we want to be sure to include all the important features, then we can use AIC or C_p + backward selection
- If we want to be sure to only include important features, then we can use BIC + forward selection
- If we want to do predictions, use AIC or C_p , but it isn't clear what selection method is the best

IMPORTANT: As stated in “Building multiple regression models interactively” (1981) *Biometrics*

“The data analyst knows more than the computer...

failure to use that knowledge produces inadequate data analysis”



ADDITIONAL COMMENTS



Some famous comments about stepwise variable selection:

(From Frank Harrel, author of “Regression Modeling Strategies”)

- The F and chi-squared tests quoted next to each variable on the printout do not have the claimed distribution.
- The method yields confidence intervals for effects and predicted values that are falsely narrow
- It yields p-values that do not have the proper meaning, and the proper correction for them is a difficult problem.
- It gives biased regression coefficients that need shrinkage
- It has severe problems in the presence of collinearity.
- It is based on methods (e.g., F tests for nested models) that were intended to be used to test prespecified hypotheses.
- Increasing the sample size does not help very much
- It allows us to not think about the problem.

(As an aside, all subsets regression solves none of these problems)

BIG DATA

Selection methods can be used on very large data sets

In R, some pseudo-code for the first step:

```
AIC_star = Inf
j_star = 0
for(j in 1:p){
  grabVec = rep('NULL',p)
  grabVec[j] = NA
  X = read.csv('bigDataSet.csv', colClasses=grabVec)
  AIC_x = AIC(X) #Get the AIC value for this vector
  if(AIC_x < AIC_star){
    AIC_star = AIC_x
    j_star = j
  }
}
```

(Note: the ideas from branch and bound can be used here as well)

SOME MORE ADVANCED/RECENT DEVELOPMENTS

There has been a lot of recent research about these topics:

- ALL-SUBSETS:

- ▶ There are more modern approaches to all-subsets than “branch and bound” based on **mixed-integer programming**. This opens all-subsets up to noticeably larger problems

(See <https://projecteuclid.org/euclid.aos/1458245736>)

- ▶ Often, the perspective is that if we could all-subsets it would lead to the best results. This isn't necessarily the case and is only true for “easy” problems

(See <http://www.stat.cmu.edu/~ryantibs/papers/bestsubset.pdf> for a very readable paper)

- POST MODEL SELECTION INFERENCE:

- ▶ As noted, the reported p-values after model-selection are not to be trusted. There are ways that p-values can be computed after model selection

(See <http://www.stat.cmu.edu/~ryantibs/papers/lassoinf.pdf>)

Postamble:

- Armed with our risk estimators, we need ways to sift through the possible models
(This can be done either via a **global and slow** approach or a **local and fast(er)** approach)
- Available techniques vary with the **absolute** and **relative** sizes of n and p
(If n and p are very small, all subsets is possible. Larger values of either indicate all subsets is not feasible)
- Like most statistical techniques, model selection comes down to optimization...