

$b_0 \rightarrow 6$

Linear Model Selection & Regularization

In this chapter, we consider some approaches for extending the linear model framework.

In this chapter we will generalize the linear model in order to accommodate non-linear but still additive relationships.

→ Despite linear models simplicity, linear model has distinct advantages in terms of interpretability & often shows good predictive performance.

Hence in this chapter ~~as well~~ we discuss some ways in which the simple linear model can be improved by replacing ordinary least squares fitting with some alternative fitting procedures.

Why do we even need an alternative to least squares?

- Prediction Accuracy: especially when trying to control the variance
- Model Interpretability: by removing irrelevant features that is, by setting the corresponding coefficient estimate to zero - we can find a model that is more easily interpretable. In this chapter we'll learn the approaches of feature selection.

Feature Selection can be done in 3 ways:

- Subset Selection: identifies a subset of predictors that not only reduces least squares but also variance & bias.
- Shrinkage: fits a model with all " β " predictors but the estimated coefficients are shrunk to 0. This shrinkage (also known as regularization) reduces the variance & also does variable selection.
- Dimension Reduction: we project our " β " predictors on a " M " dimensional subspace, where $M \ll p$.

1. Subset Selection: (Complexity: 2^P)

i.) Best Subset Selection:

Ans. Let M_0 denote the null Model, which contains no predictors. This will simply predict the sample mean of the observations.

2. For $K=1, 2, 3, \dots, p$:

- (a) fit all $\binom{p}{K}$ models that contain exactly K predictors
- (b) pick the best amongst the $\binom{p}{K}$ models & call it M_K . Here best is defined as smallest RSS or largest R^2 .

3. Select a single best model from among, M_0, M_1, \dots, M_p using cross validated predicted error (e.g. AIC), BIC or adjusted R^2 .

ii) Opposing Reasons Disadvantages of Best Subset Selection:

- ① For computational reasons, best subset selection can not be applied with very large " p ". Say $p=10$, then total number of subsets ≈ 1000 , but when $p=40$, subsets = 1,000,000,000,000

②

& that is computationally very expensive.
It will also suffer from statistical problems when "p" is
the bigger the search date the bigger the chances of
a model that works good on training data but might
have any power on future data.

Thus an enormous search space can lead to overfitting.

ii) STEPWISE SELECTION: Complexity: P^2

For the reasons mentioned above stepwise methods work better.

① Forward Stepwise Selection:

→ This begins with a model containing no predictors, adds predictors to the model one at a time, until all predictors are in model. At each step the variable gives the best improvement to the fit is added to model.

→ STEPS:

1. Let M_0 denote a null Model, which contains no predictors

2. For $K = 0, 1, 2, \dots, p-1$:

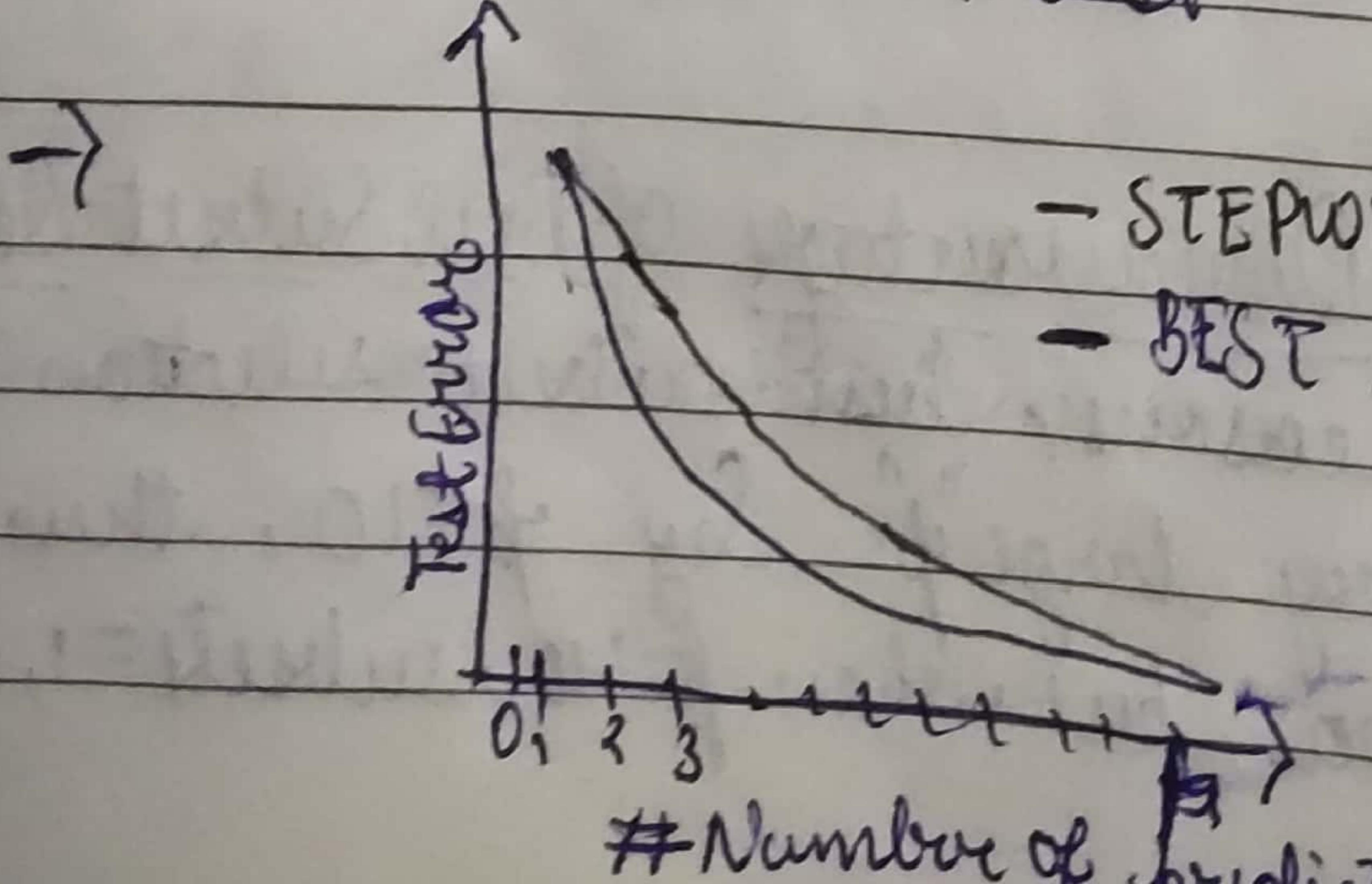
2.1 Consider all " $p-K$ " models that augments the predictors in M_K with one additional predictor.

2.2. Choose the best amongst all the models available.

3. Select a single best model from M_0 to M_p .

→ Computational advantage over best subset selection.

→ But this is not guaranteed that the best model selected will be the actual best.



- STEPWISE SELECTION

- BEST SUBSET SELECTION

→ Why does best subset selection provides lower test error? Reason being there might be some positive correlation between features which cannot be found in best subset selection.

→ Why does the graph start & end at the same point because they both start with 0 predictors & ends with all n predictors.

② Backward Stepwise Selection :

→ Backward Stepwise selection provides an efficient alternative to best subset selection like forward stepwise selection.

→ However, it starts with all "p" predictors & removes one at a time.

→ STEPS:

1. Let M_p denote the full model with all "p" predictors.

2. For $k = p, p-1, p-2, \dots, 1$:

- 2.1 Consider all K models with one removed predictor.

- 2.2 Choose the Best model with least RSS or highest R^2 .

3. Select a single best model from M_0, \dots, M_p

CHOOSING THE OPTIMAL MODEL :

→ The chosen model must have lowest RSS & highest R^2 .

→ We choose to with a model with ~~the~~ low test error not train error.

→ Therefore, RSS & R^2 are not good predictor for best model.

ESTIMATING THE TEST ERROR:

- We can indirectly estimate test error by making an adjustment to the training error to account for bias to the overfitting.
- We can directly estimate test error using a validation set approach or the cross-validation approach.

Adjusting Loss Function: (CP, AIC, BIC)

- Mallows' CP:

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

where, $d = \# \text{ of parameters}$

$\hat{\sigma}^2$ = estimate of variance of error term.

- AIC: (Akaike Information Criteria)

The AIC is defined by a large number of models fit by maximum likelihood.

$$AIC = -2 \log L + 2 \cdot d$$

where, L is the maximized value of the likelihood function

- BIC: (Bayesian Information Criteria)

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

- Adjusted R^2 :

$$1 - \frac{RSS / (n-d-1)}{TSS / (n-1)}$$

Validation & Cross Validation -

- We can use validation or cross validation for selecting the model with least test error.
- This has an advantage over AIC, BIC, Cp, adjusted R² that it provides direct estimate of test error.

② Shrinkage Methods:

The subset selection uses a subset of predictor to fit a linear model with least square, shrinkage methods on the other hand fits the model with all "p" predictors using a technique that constrains or regularizes the coefficient estimates, or equivalently shrinks the coefficient estimates towards zero.

i) Ridge Regression:

Recall, that least squares is given by:

$$RSS = \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2$$

On the other hand ridge regression adds a penalty to RSS, & β_R tries to minimize least square for the equation below:

$$\sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 = RSS + \lambda \sum_{j=1}^p \beta_j^2$$

where, $\lambda > 0$ is a tuning parameter, i.e. a hyperparameter.

- as with least squares, ridge regression seeks coefficient estimates that fit the data well by making RSS small.
- however, second term in ridge, i.e. $\lambda \sum \beta_j^2$, called a shrinkage penalty is small when $\beta_1, \beta_2, \beta_3, \dots, \beta_p$ are close to zero, so it has the effect of shrinking β_j to 0.
- so, this will automatically shrink less important parameters towards zero.
- λ here serves as tuning parameter, higher the lambda, lower the value of coefficients.
- The penalty part is also called l_2 -norm & is given by $\|\beta\|_2$

$$\|\beta\|_2 = \sqrt{\sum \beta_j^2}$$

Scaling the predictors in ridge regression is pretty important, why?

- Multiplying x_j with a constant "c" simply leads a scaling of least square coefficient estimates by a factor of $1/c$.
- So, because of the scaling in ridge regression coefficient estimates will lead to change in the output of the loss function because of the presence of the term $\sum \beta_j^2$ in the penalty.
- Scaling wasn't important in ordinary least squares but is important here because of the presence of penalty terms.

∴ it is best to apply ridge regression after standardization.

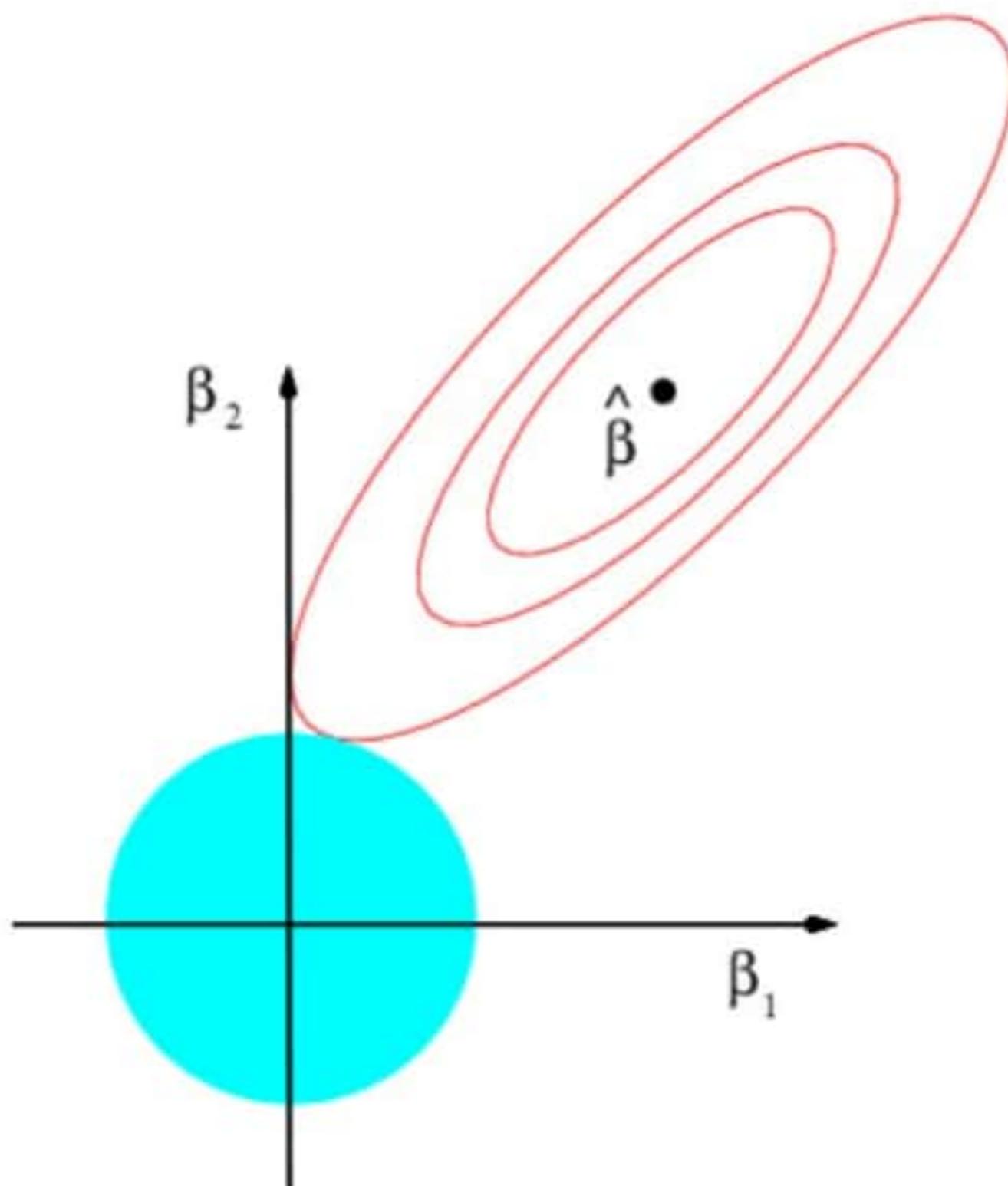
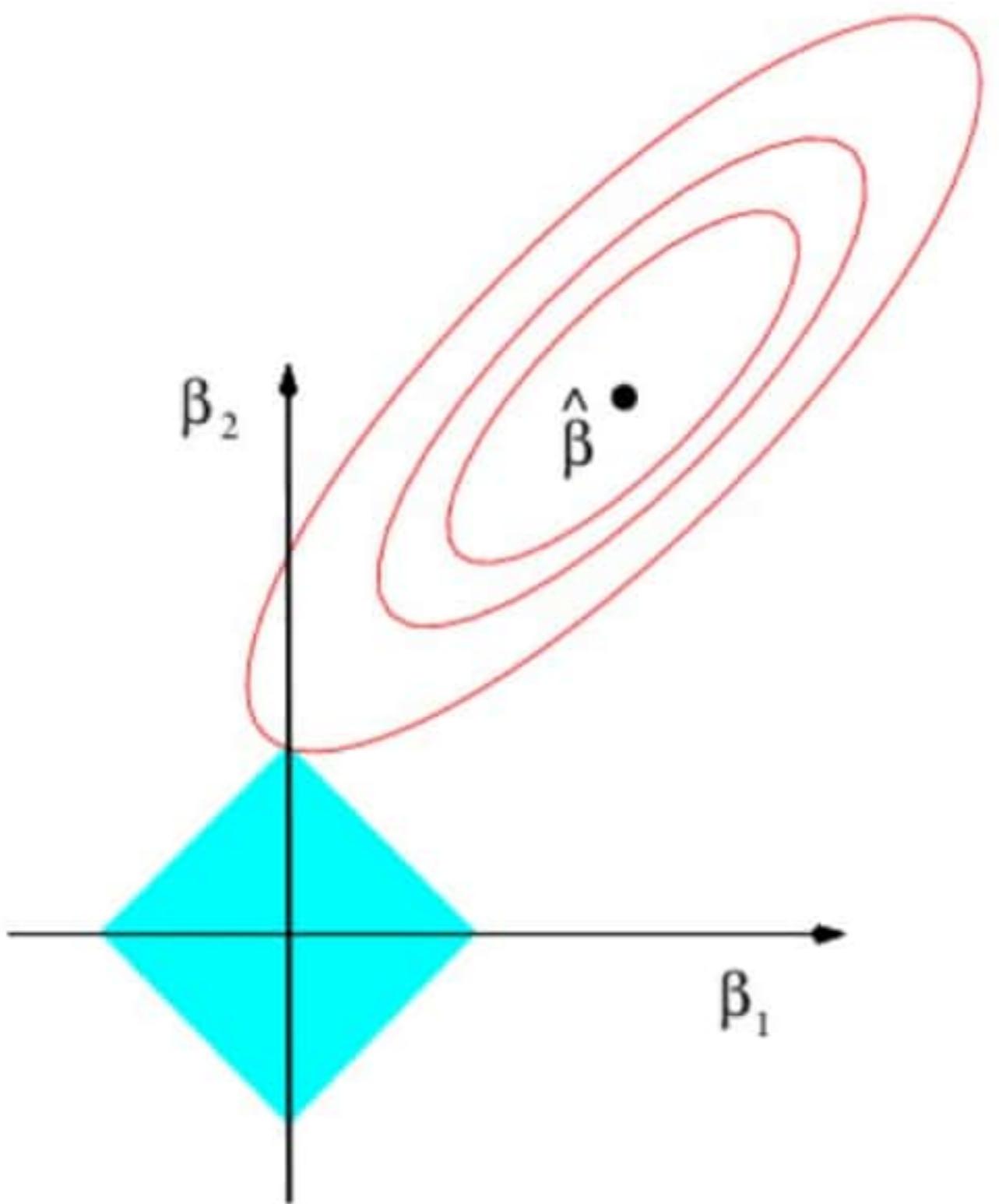
$$\tilde{x}_{ij} = \frac{x_{ij}}{\sqrt{\frac{1}{n} \sum (x_{ij} - \bar{x}_j)^2}} = \frac{x_{ij}}{\sqrt{s^2}}$$

ii) The LASSO:

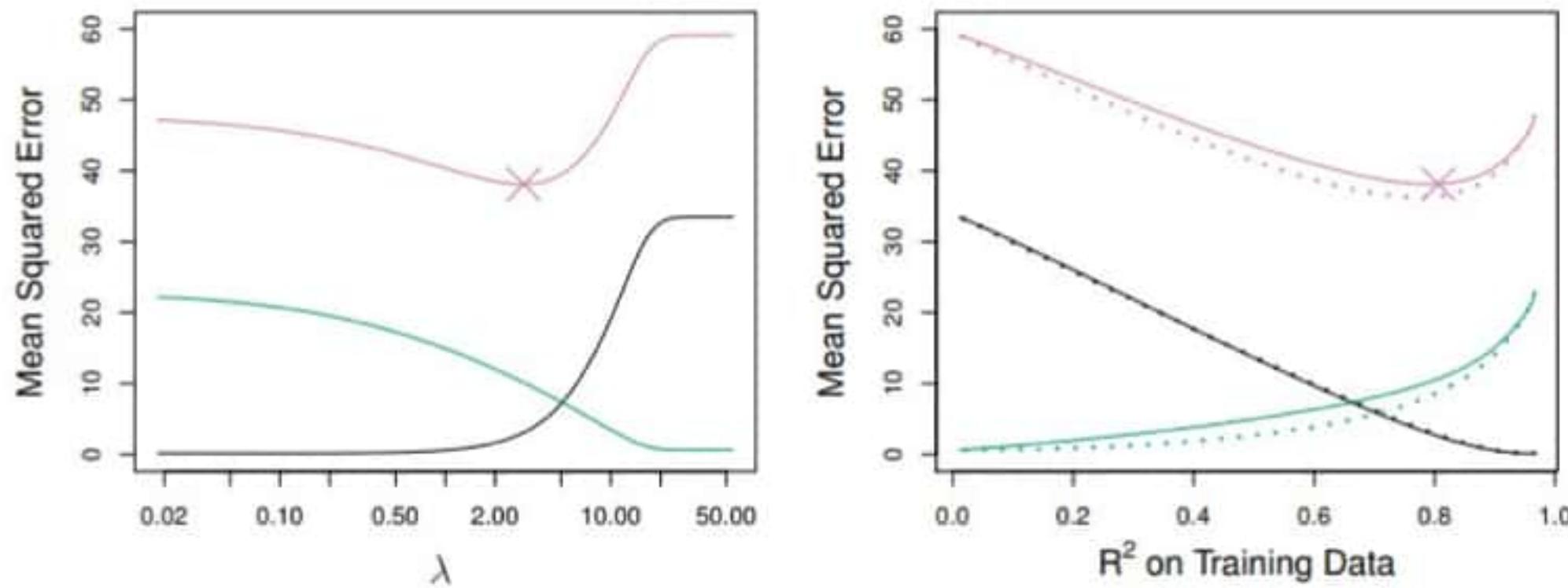
- Ridge Regression has one obvious disadvantage:
 Unlike subset selection, which selects "m" parameters out of "p" (~~(m < p)~~), ridge might reduce the coefficients close to zero but not exactly zero.
- The Lasso is a relatively recent alternative to ridge that overcomes this disadvantage. The lasso coefficients, $\hat{\beta}_j$, minimize the quantity

$$RSS + \lambda \sum_{j=1}^p |\beta_j|$$
- In statistical parlance, the lasso uses an L_1 penalty ($\|\beta_j\|_1$) instead of L_2 penalty
- However, both ridge & lasso shrink coefficient estimates towards zero, but lasso has an effect of forcing some of the coefficients towards zero.
- Hence like subset selection lasso does variable selection.
- We say lasso is a sparse model while ridge is a density model.

The Lasso Picture



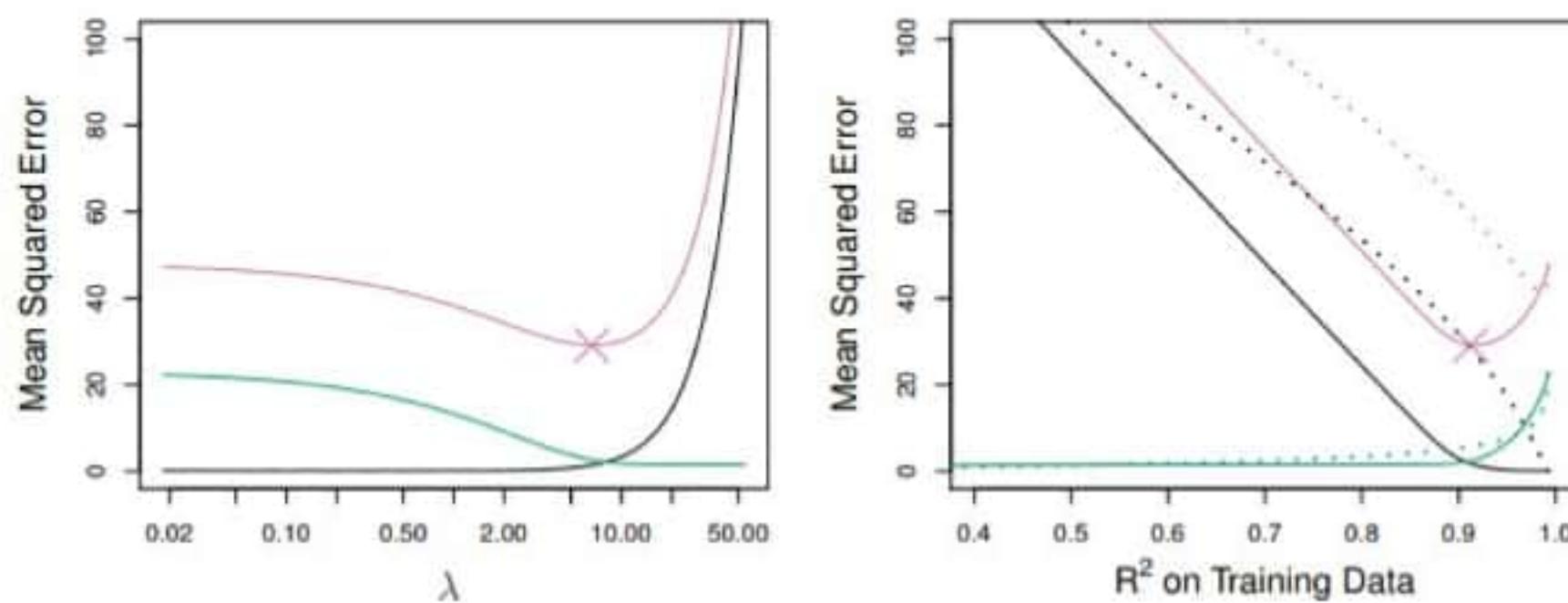
Comparing the Lasso and Ridge Regression



Left: Plots of squared bias (black), variance (green), and test MSE (purple) for the lasso on simulated data set of Slide 32.

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 crosses in both plots indicate the lasso model for which the MSE is smallest.

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 is similar to that in Slide 38, 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 crosses in both plots indicate the lasso model for which the MSE is smallest.