



School of Business

Multivariate Data Analytics

Model Comparison and Selection

Prof. Feng Mai
School of Business

For academic use only.



Key concepts

- Objective of least squared method
- Also called **sum of squared estimate of errors (SSE)**
- **Mean squared error**
- n is the sample size, p is the number of parameters (number of x variables + 1)

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

loss of OLS
 obj
 MLE

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

$$MSE = \frac{RSS}{n-p}$$

estimator σ^2

- **Total sum of squares**
- Total amount of variance in y
- **Coefficient of Determination, R-squared**
- The proportion of variance explained by the model, from 0 to 1.
- Always increases when more variables are added in the model. ↙

Other model comparison criteria

- Adjusted R²

$$R_{adj}^2 = 1 - (1 - R^2) \frac{n - 1}{n - P - 1}$$

- Higher the better

- AIC (Akaike information criterion)

$$AIC = 2P + n \log(RSS/n)$$

RSS (min)

- Lower the better

- BIC (Bayesian information criteria)

$$BIC = P \log(n) + n \log(RSS/n)$$

- Lower the better

most "conservative"

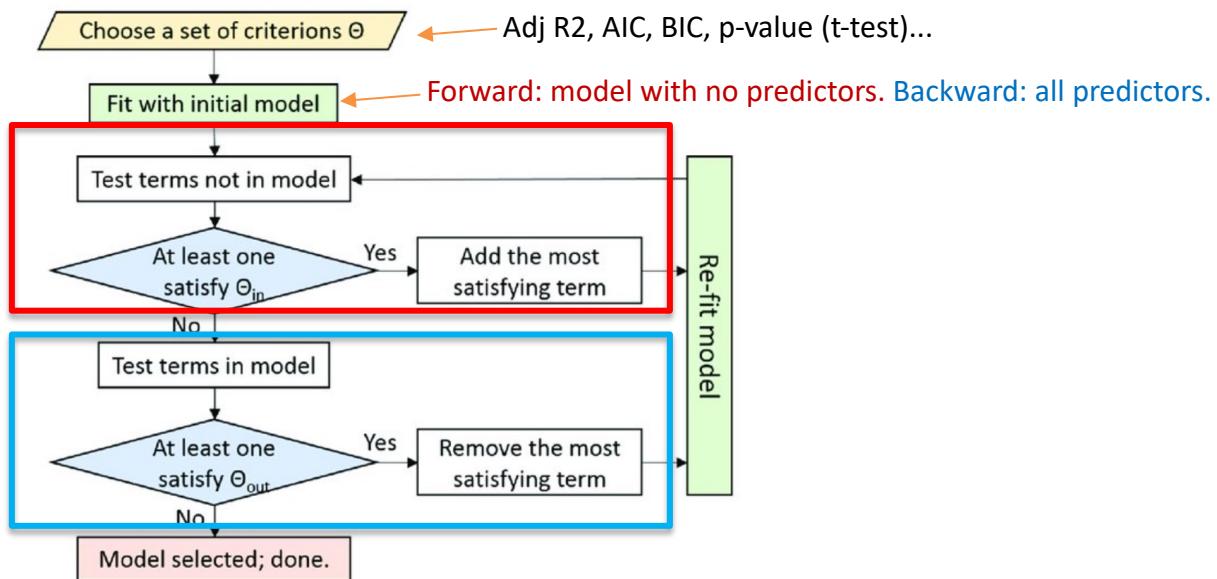
of P (X variables)

$$BIC < AIC < \text{Adj } R^2$$

Model selection

- **Best subset regression:** fit all the possible models from all of the possible combinations of the predictors, then compare adj R² / AIC /BIC
 - 2^k possible models for k predictors
 - Over 1 million models with 20 predictors ↪
 - **Stepwise regression:** forward selection, backward elimination, stepwise

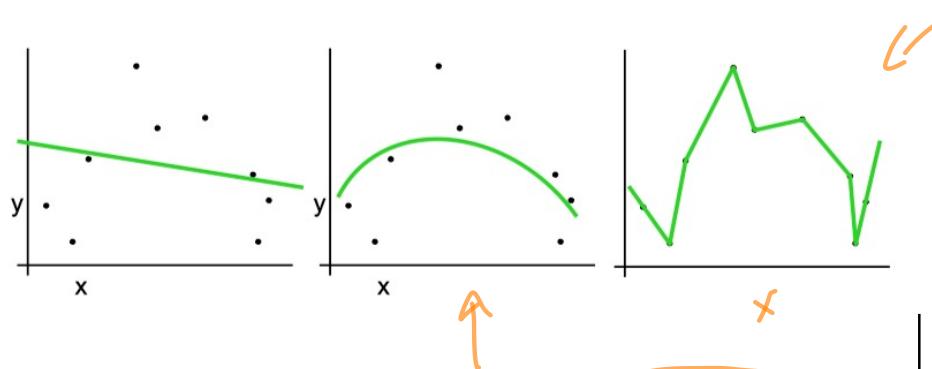
$$\# \text{ of males} = \underbrace{k + k-1 + k-2 \dots}_{\text{number of males}}$$



Source: Yang et al 2017. [10.3390/nano7100319](https://doi.org/10.3390/nano7100319)

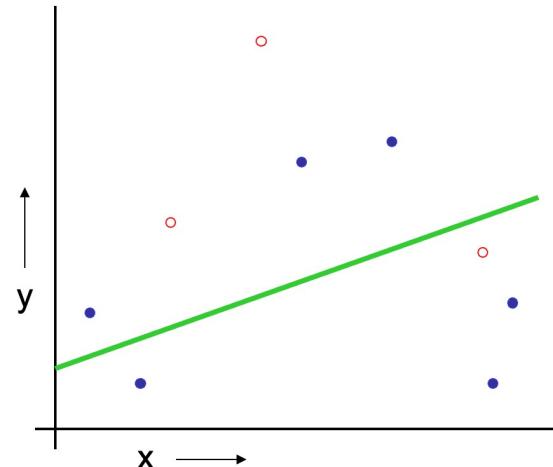
Model Comparison Revisited

- A good model should be generalizable (useful for prediction)
- **Adj R², AIC, BIC**: penalizes complicated models after fitting
- **Test Set Approach**: Find the best model by explicitly testing models' ability to generalize

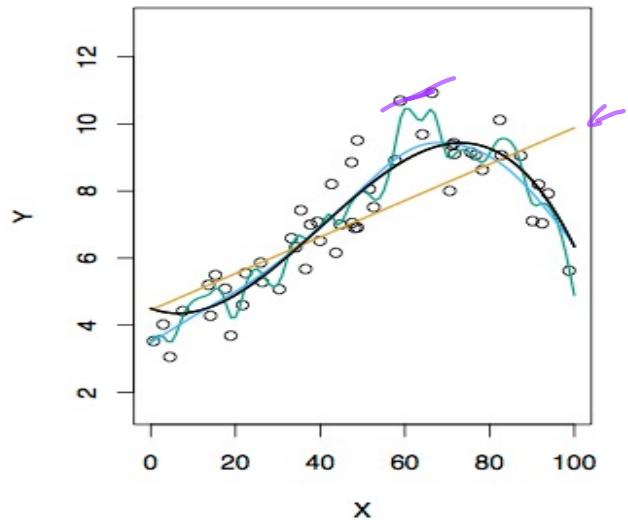


1. Randomly choose 30% of the data to be in a **test set**
2. The remainder is a **training set**
3. Fit regression models on the training set
4. For each model, estimate future performance (MSE) with the test set

$$\frac{SSE}{n}$$
5. Use the model with the best test set performance



Source: Adapted from Andrew W. Moore's Tutorial



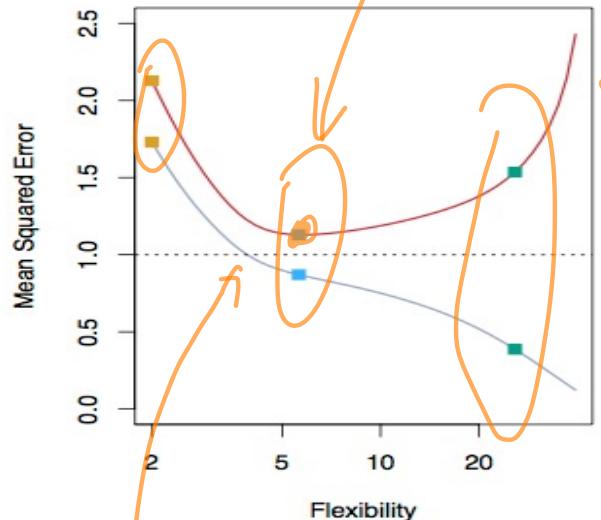
LEFT

Black: Truth

Orange: Linear Estimate

Blue: smoothing spline

Green: smoothing spline (more flexible)

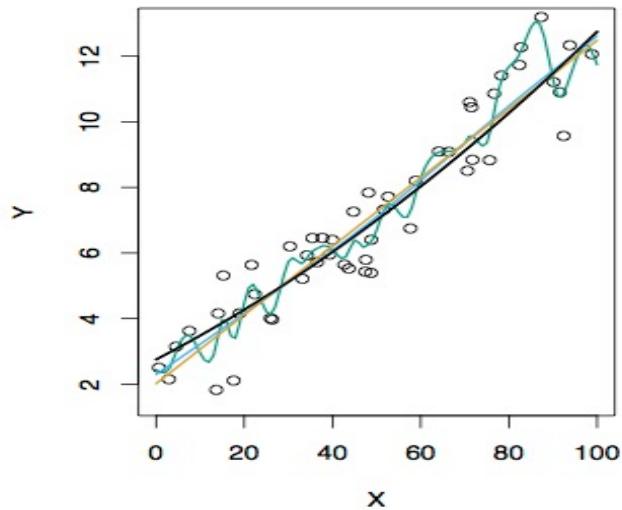


RIGHT

RED: Test MSE

Grey: Training MSE

Dashed: Minimum possible test MSE (irreducible error)



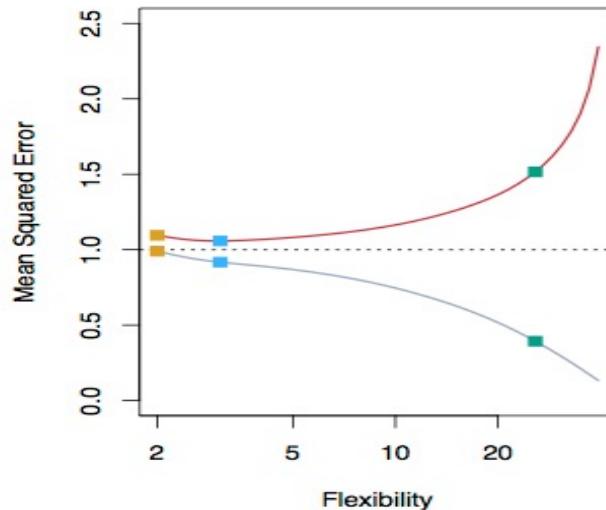
LEFT

Black: Truth

Orange: Linear Estimate

Blue: smoothing spline

Green: smoothing spline (more flexible)



RIGHT

RED: Test MES

Grey: Training MSE

Dashed: Minimum possible test MSE (irreducible error)

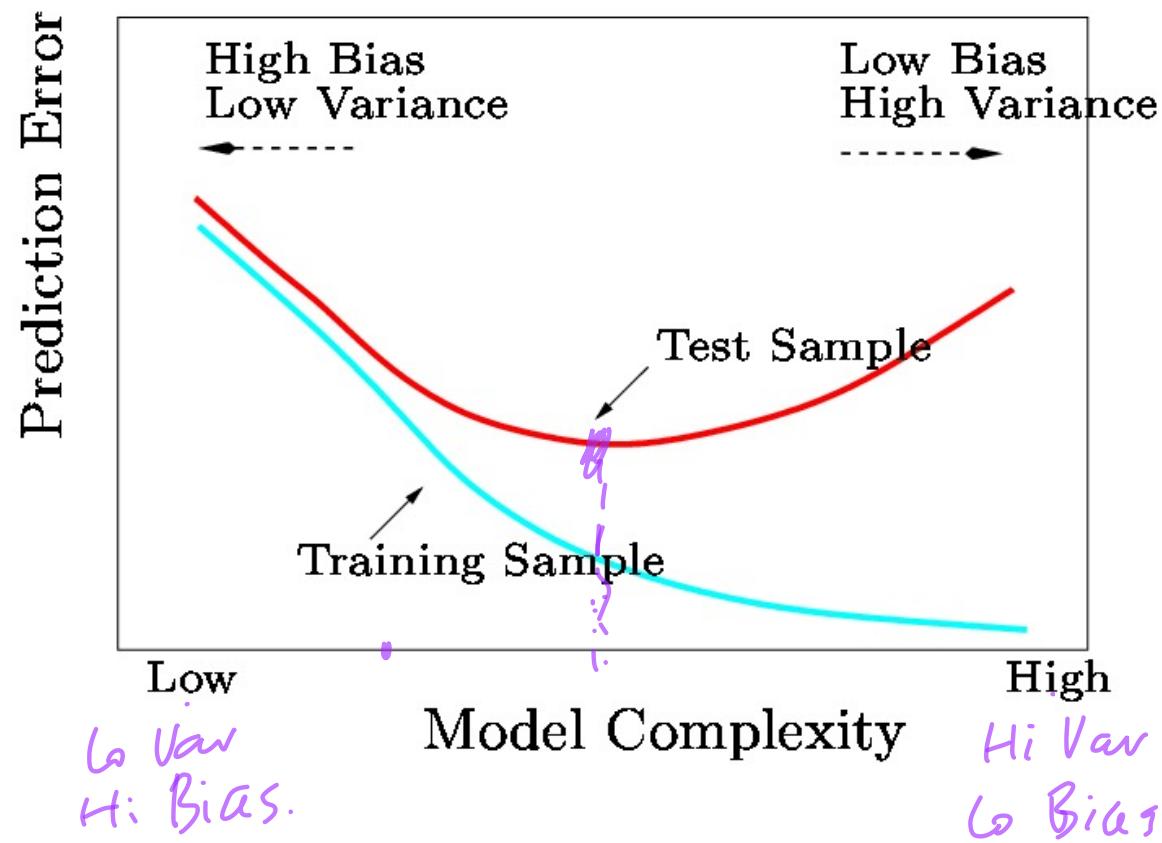
Bias/ Variance Tradeoff

- The previous graphs of test versus training MSE's illustrates a very important tradeoff that governs the choice of statistical learning methods.
- There are always two competing forces that govern the choice of learning method i.e. bias and variance.
- Bias refers to the error that is introduced by modeling a real life problem (that is usually extremely complicated) by a much simpler problem.
 - For example, linear regression assumes that there is a linear relationship between Y and X. It is unlikely that, in real life, the relationship is exactly linear so some bias will be present.
 - The more flexible/complex a method is the less bias it will generally have.
- Variance refers to how much your estimate for f would change by if you had a different training data set.
 - Generally, the more flexible a method is the more variance it has.
- The Trade-Off

$$\text{Expected Test MSE} = E(Y - f(x_0))^2 = \underbrace{\text{Bias}^2}_{\text{Irreducible Error}} + \underbrace{\text{Var}}_{\text{Irreducible Error}} + \underbrace{\sigma^2}_{\text{Irreducible Error}}$$

- As a method gets more complex the bias will decrease and the variance will increase but expected test MSE may go up or down!

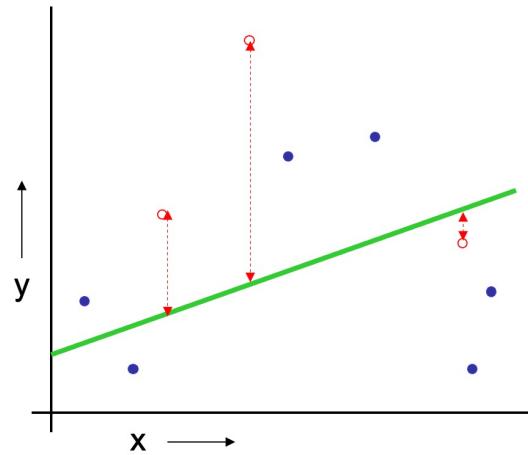
- In general training errors will always decline.
- However, test errors will decline at first (as reductions in bias dominate) but will then start to increase again (as increases in variance dominate).



We must always keep this picture in mind when choosing a learning method. More flexible/complicated is not always better!

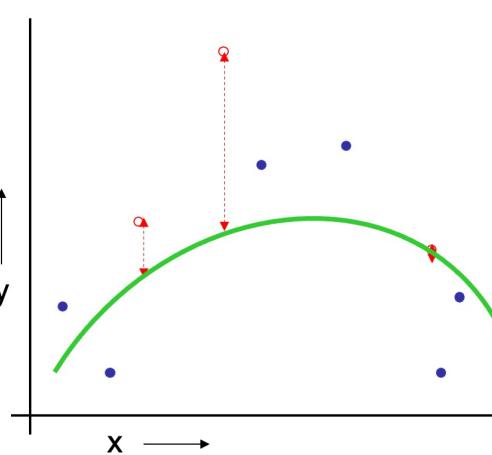
Test Set Performance

All models are fitted using **train set (blue points)** but evaluated on the **test set (red points)**



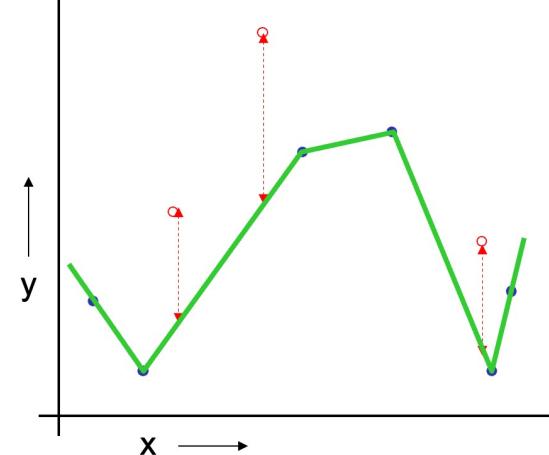
(Linear regression example)

Mean Squared Error = 2.4



(Quadratic regression example)

Mean Squared Error = 0.9



(Join the dots example)

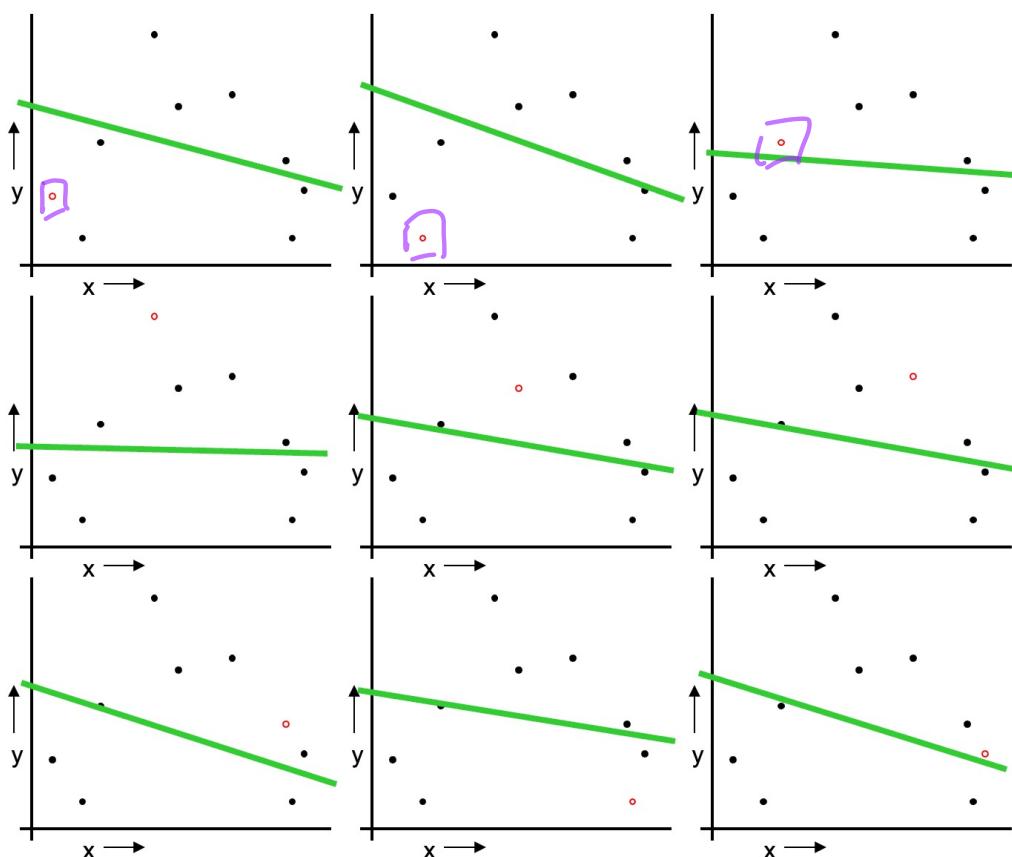
Mean Squared Error = 2.2

Best Model!

But...

- We wasted data.
- Performance estimation on small test set has high variance

LOOCV (Leave-one-out Cross Validation)



For $k=1$ to n

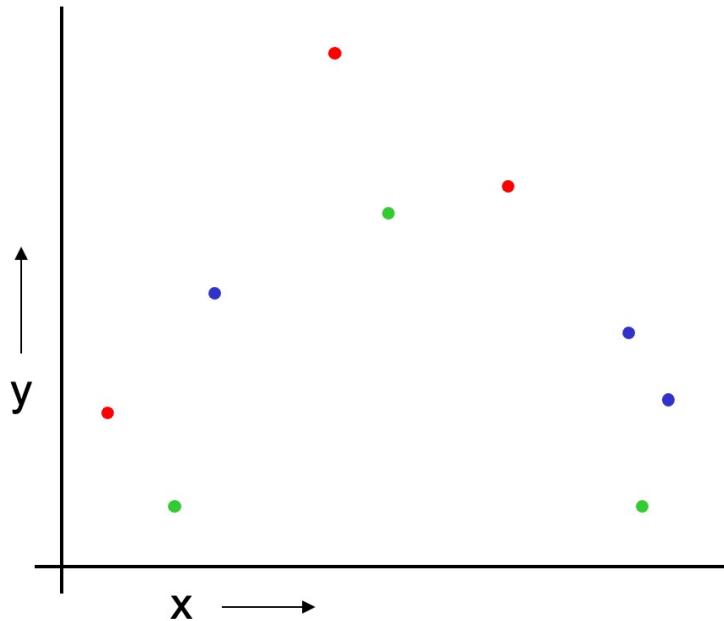
1. Let (x_k, y_k) be the k^{th} record
2. Temporarily remove (x_k, y_k) from the dataset (use this single point as the test set)
3. Train on the remaining $n-1$ datapoints
4. Note your error (x_k, y_k)

When you've done all points, report the average error for all points.

Pros: Wasted less data and more stable

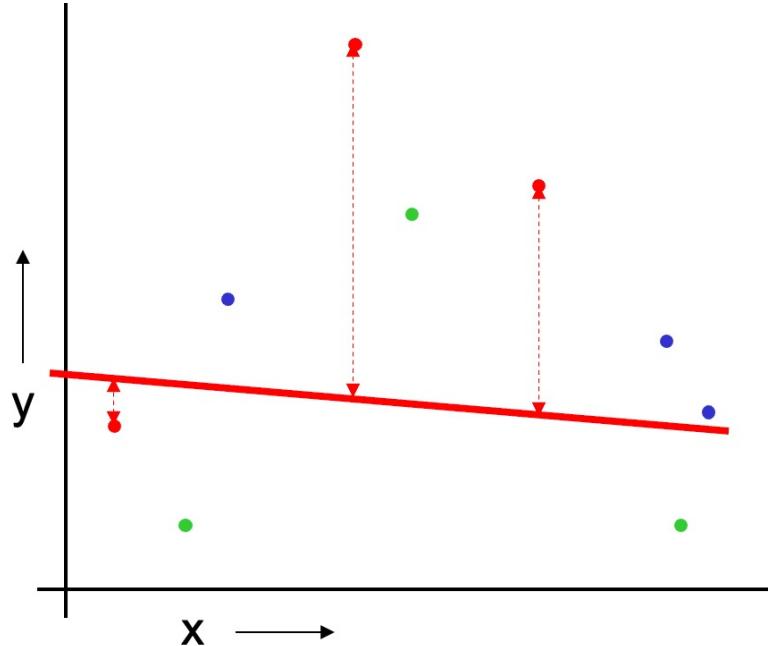
Cons: Time consuming

k-Fold Cross Validation



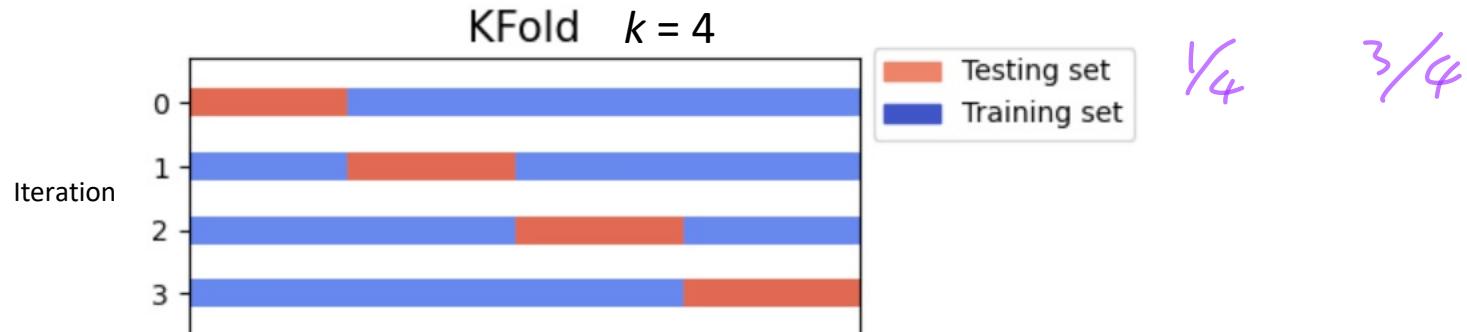
- Randomly break the dataset into k partitions (in our example we'll have k=3 partitions colored Red, Green and Blue)
 - For the red partition: Train on all the points not in the red partition (2/3 of the data). Find the test-set sum of errors on the red points (1/3 of the data).
 - For the green partition: Train on all the points not in the green partition (2/3 of the data). Find the test-set sum of errors on the green points (1/3 of the data).
 - For the blue partition: Train on all the points not in the blue partition (2/3 of the data). Find the test-set sum of errors on the blue points (1/3 of the data).
- Then report the average error

k-Fold Cross Validation

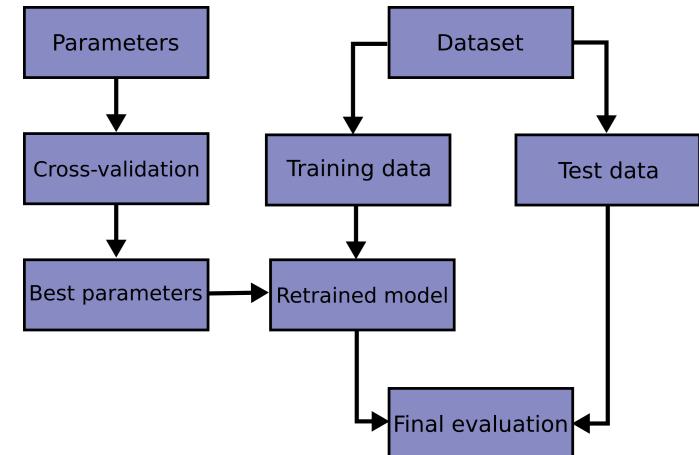


- Randomly break the dataset into k partitions (in our example we'll have $k=3$ partitions colored Red, Green and Blue)
 - For the red partition: Train on all the points not in the red partition (2/3 of the data). Find the test-set sum of errors on the red points (1/3 of the data).
 - For the green partition: Train on all the points not in the green partition (2/3 of the data). Find the test-set sum of errors on the green points (1/3 of the data).
 - For the blue partition: Train on all the points not in the blue partition (2/3 of the data). Find the test-set sum of errors on the blue points (1/3 of the data).
- Then report the average error

k-Fold Cross Validation



- Widely used to compare models and tune model parameters
- Less computationally expensive than LOOCV
 - Still more computationally expensive than AIC/BIC
- More stable than test set approach (average performance of the k partitions)
- Typically choose $k = 10$ or 5



Source: https://scikit-learn.org/stable/modules/cross_validation.html

Review

- In a linear regression (OLS), the goal is to **minimize** the residual sum of squares:

$$\text{RSS} = \sum_{i=1}^n \left(y_i - \underbrace{\beta_0 + \sum_{j=1}^p \beta_j x_{ij}}_{\hat{y}_i} \right)^2. \quad \text{obj.}$$

- If a dataset has many predictors, we can fit a full model with all predictors.
 - Not all variables are useful
 - More complicated models are less generalizable (overfit)
- We can use best subset regression or stepwise regression to find a simpler and more generalizable model with a subset of predictors.
 - best subset regression: slow
 - stepwise regression: inconsistent results with different model selection algorithms

Ridge Regression

- Modifying the objective

Min

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

OLS objective

Tuning parameter

Model complexity penalty

Tuning parameter λ is positive. When $\lambda = 0$, we get the OLS objective

- Since we are minimizing the objective, model with high β s = bad! ↩
- This has the effect of “shrinking” large values of β s towards zero. ↩
- Improve the fit, because shrinking the coefficients can significantly reduce their variance

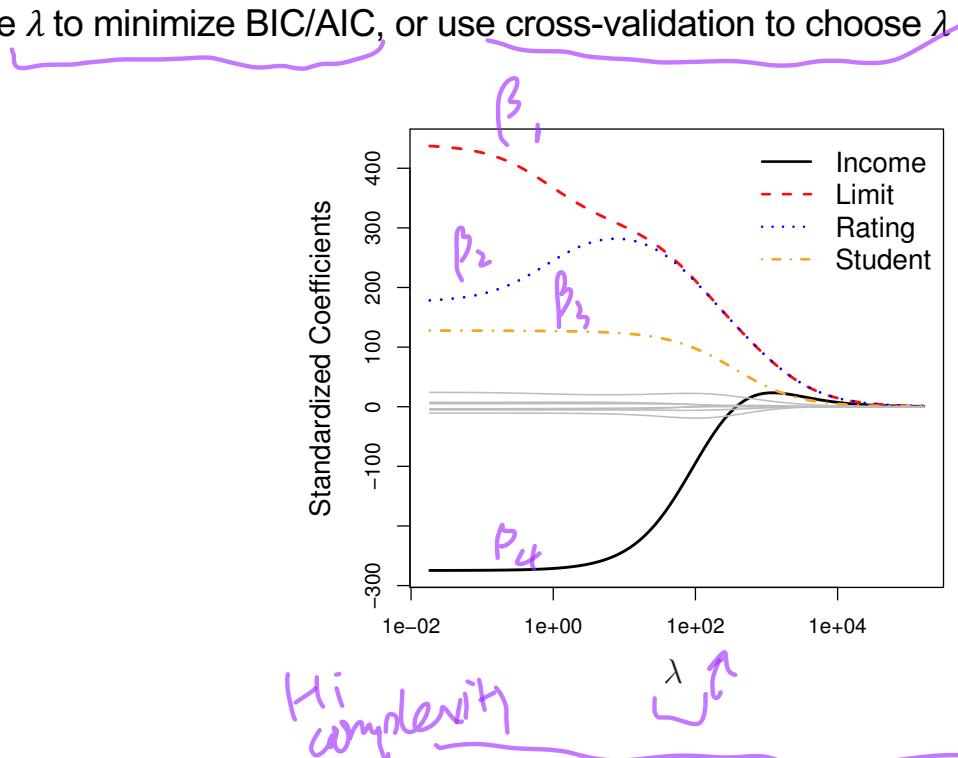
- !!
- Also handles multicollinearity
 - Also called L2 regularization

X are highly correlated
 β s have high variance
 High Bias
 Low Var

Complexity
 High

Tuning Parameter λ for Ridge

- When $\lambda = 0$, we get the OLS objective
- As λ increases, the standardized coefficients shrinks towards zero.
- Small $\lambda \rightarrow$ More complicated models (high variance, low bias).
- Large $\lambda \rightarrow$ Simpler models (low variance, high bias).
- Choose the λ to minimize BIC/AIC, or use cross-validation to choose λ



Source: James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An introduction to statistical learning* (Vol. 112, p. 18). New York: Springer.

LASSO (least absolute shrinkage and selection operator)

- Modifying the objective

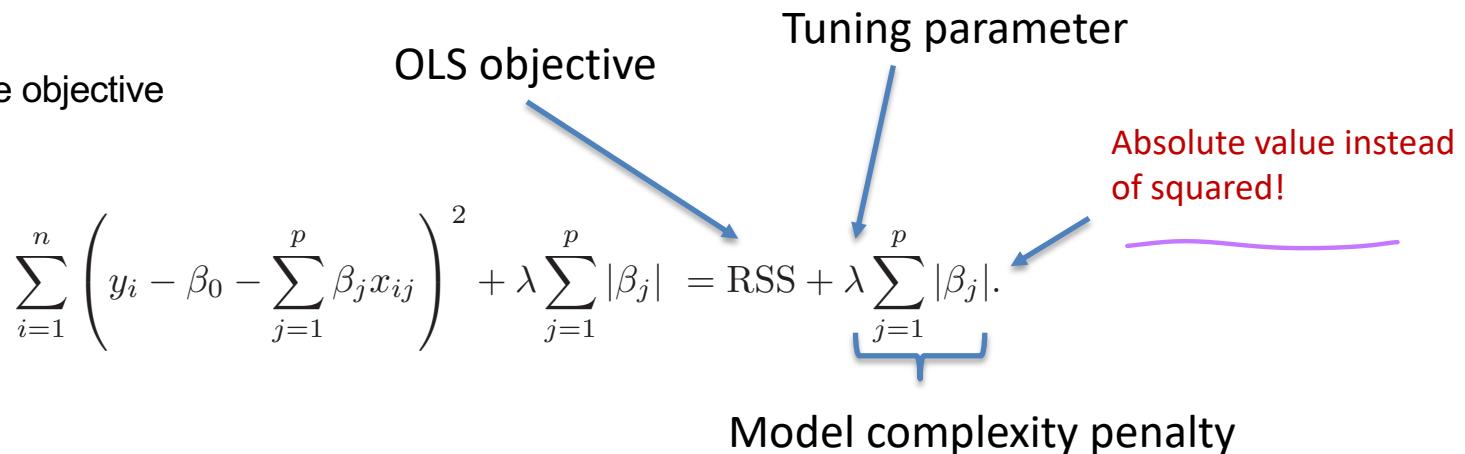
OLS objective

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

Tuning parameter

Absolute value instead of squared!

Model complexity penalty

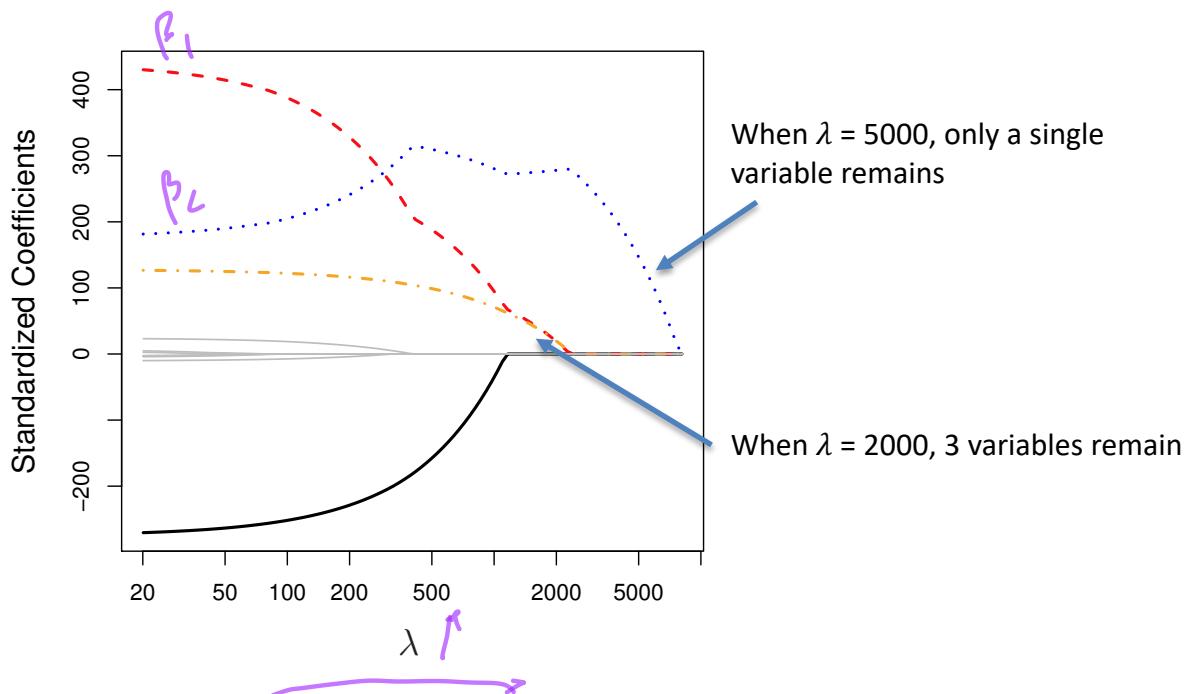


Similar idea to Ridge Regression, with a big difference in results

- Using this penalty, **some coefficients end up being set to exactly zero**
- LASSO does **variable selection!** Variables with β coefficients = 0 are eliminated.
- Also called **L1 regularization**

Tuning Parameter λ for LASSO

- When $\lambda = 0$, we get the OLS objective
- As λ increases, some coefficients become zero.
- Small $\lambda \rightarrow$ Fewer variables eliminated, more complicated models (high variance, low bias).
- Large $\lambda \rightarrow$ More variables eliminated, simpler models (low variance, high bias).
- Choose the λ to minimize BIC/AIC, or use cross-validation to choose λ



Source: James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An introduction to statistical learning* (Vol. 112, p. 18). New York: Springer.

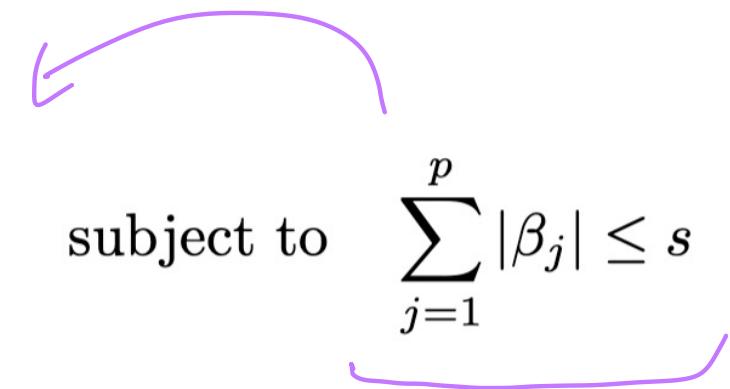
The Variable Selection Property of the Lasso

Why is it that the lasso, unlike ridge regression, results in coefficient estimates that are exactly equal to zero?

One can show that the lasso and ridge regression coefficient estimates solve the problems

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

OLS

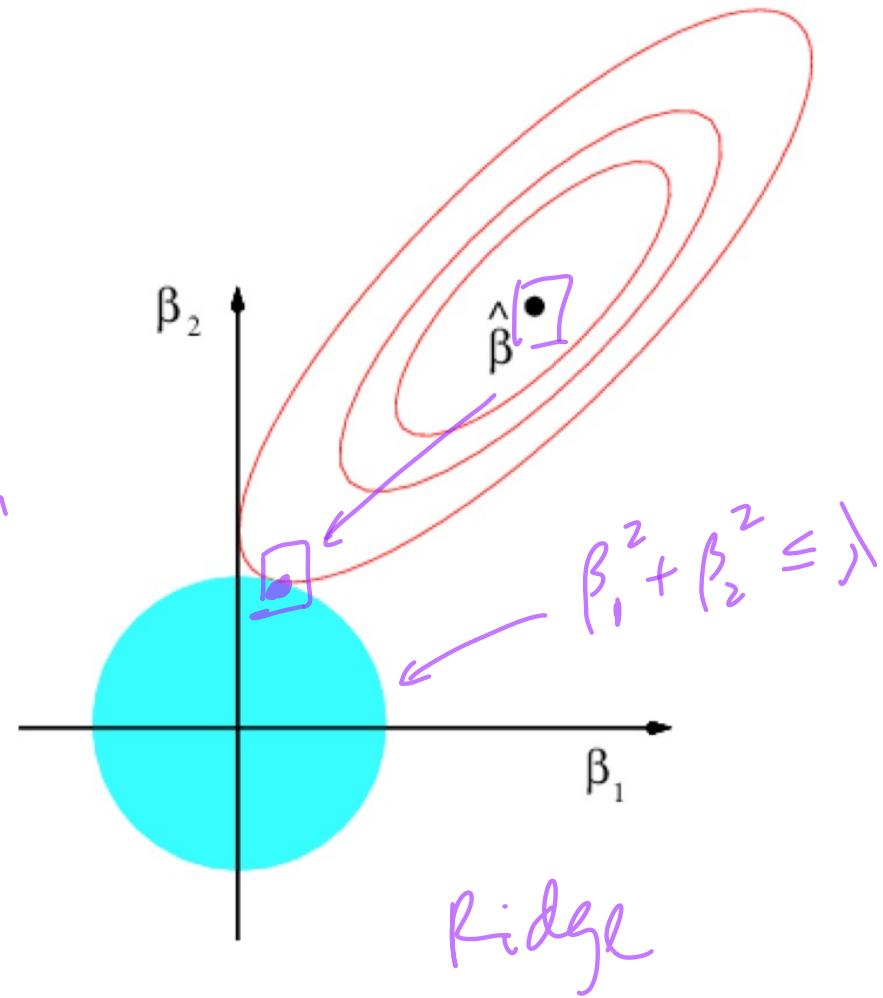
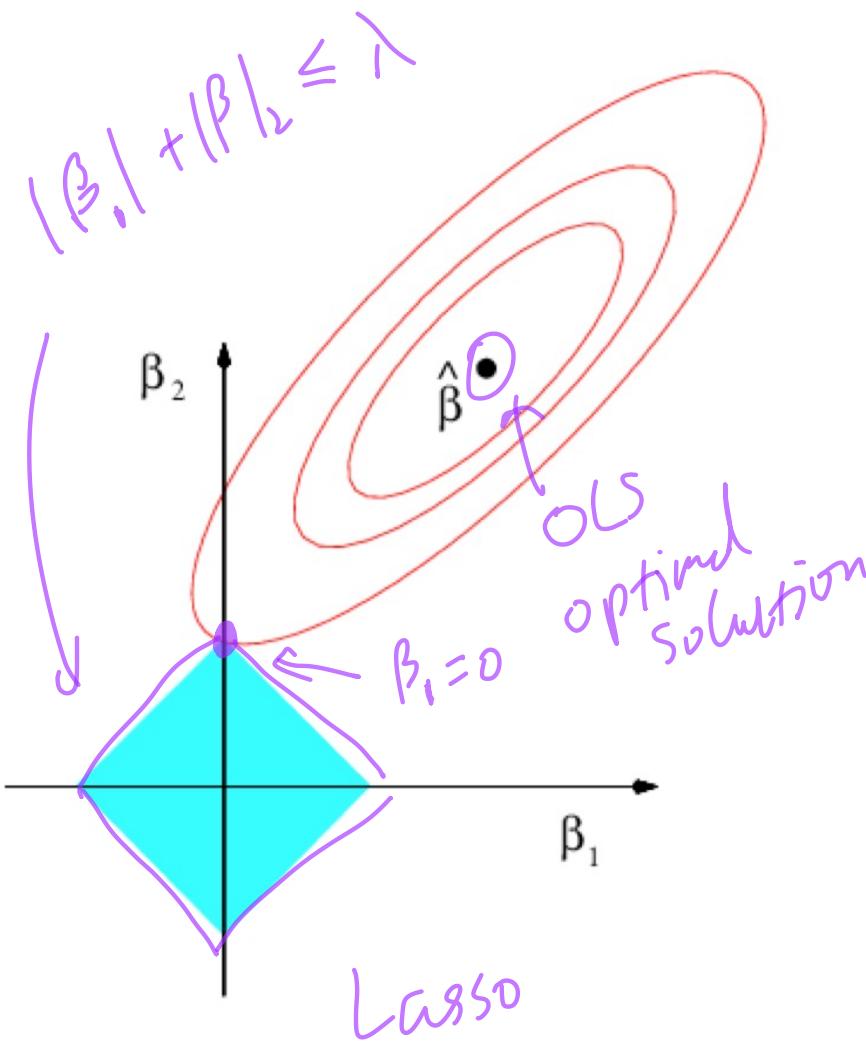


and

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

respectively.

The Lasso Picture





Thank you!

Prof. Feng Mai
School of Business

For academic use only.