# Regularized regression

Akitaka Matsuo

# Regularized regression

- We saw linear regression in the previous lecture.
- Linear regression is BLUE for the train set, but might be overly-sensitive to the train data.
- We can adjust the problem by using penalized regression.
- Methods
  - Ridge regression
  - LASSO
  - (Elastic net)

# Regularized regression, objective function

- Linear regression:

$$\text{argmin}_\beta \sum_i (Y_i - (\beta_0 + \sum_j \beta_j X_{ij}))^2$$

- Regularized regression:

$$\text{argmin}_\beta \sum_i (Y_i - (\beta_0 + \sum_j \beta_j X_{ij}))^2 + \lambda g(\beta_{-0})$$
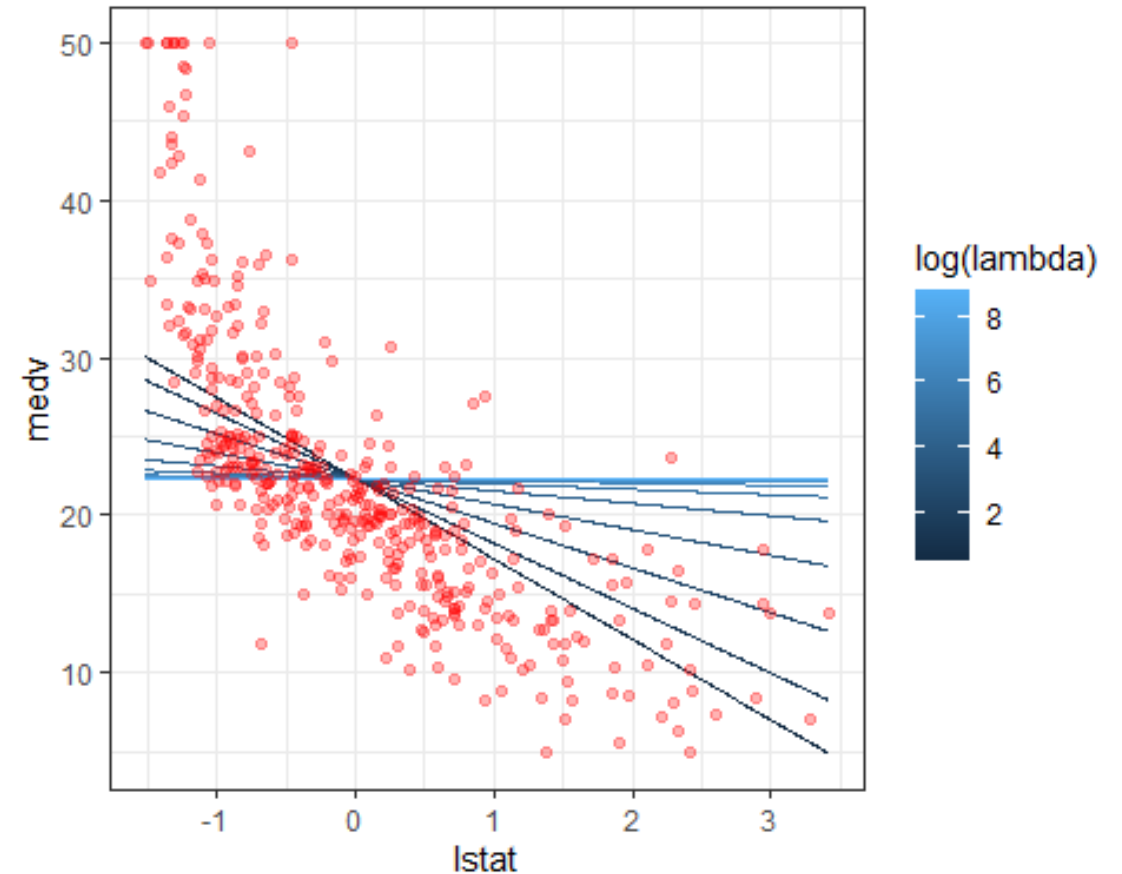
- The shape of $g(\beta)$ is different across methods

# Ridge regression

$$\text{argmin}_\beta \sum_i (Y_i - (\beta_0 + \sum_j \beta_j X_{ij}))^2 + \lambda \sum_j \beta^2$$

- $\lambda \sum_j \beta^2$ is the penalty term (i.e. shrinkage penalty)
  - L2-penalty
  - sum of the squared $\beta$s multiplied by $\lambda$
  - $\lambda = 0$: OLS
  - $\lambda = \infty$: completely shrunken $\beta$
- $\lambda$ is an only tuning parameter in ridge regression

# Ridge regression, diferent lambda

- This is an illustration of fitted line with different $\lambda$ value

- When $\lambda$ gets bigger, the line gets flatter

- The best $\lambda$ value:
  - Enough shrinkage without too much bias (see example later)

# LASSO Regression

The objective function is similar but slightly different
- – LASSO

$$\text{argmin}_\beta \sum_i (Y_i - (\beta_0 + \sum_j \beta_j X_{ij}))^2 + \lambda \sum_j |\beta|$$

- – Ridge regression

$$\text{argmin}_\beta \sum_i (Y_i - (\beta_0 + \sum_j \beta_j X_{ij}))^2 + \lambda \sum_j \beta^2$$

- – LASSO penalty
  - L1-penalty
  - sum of the absolute value of $\beta$s multiplied by $\lambda$
  - $\lambda = 0$: OLS
  - $\lambda = \infty$: completely shrunken $\beta$
- – $\lambda$ is an only tuning parameter in LASSO regression

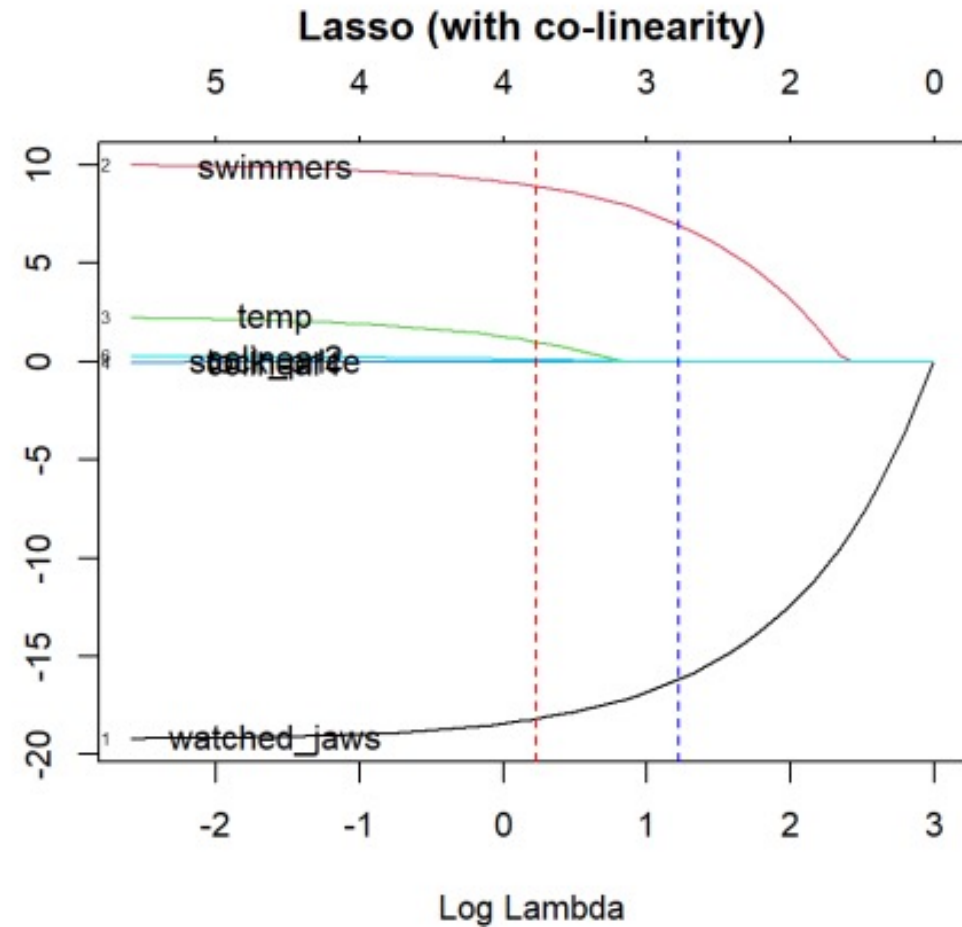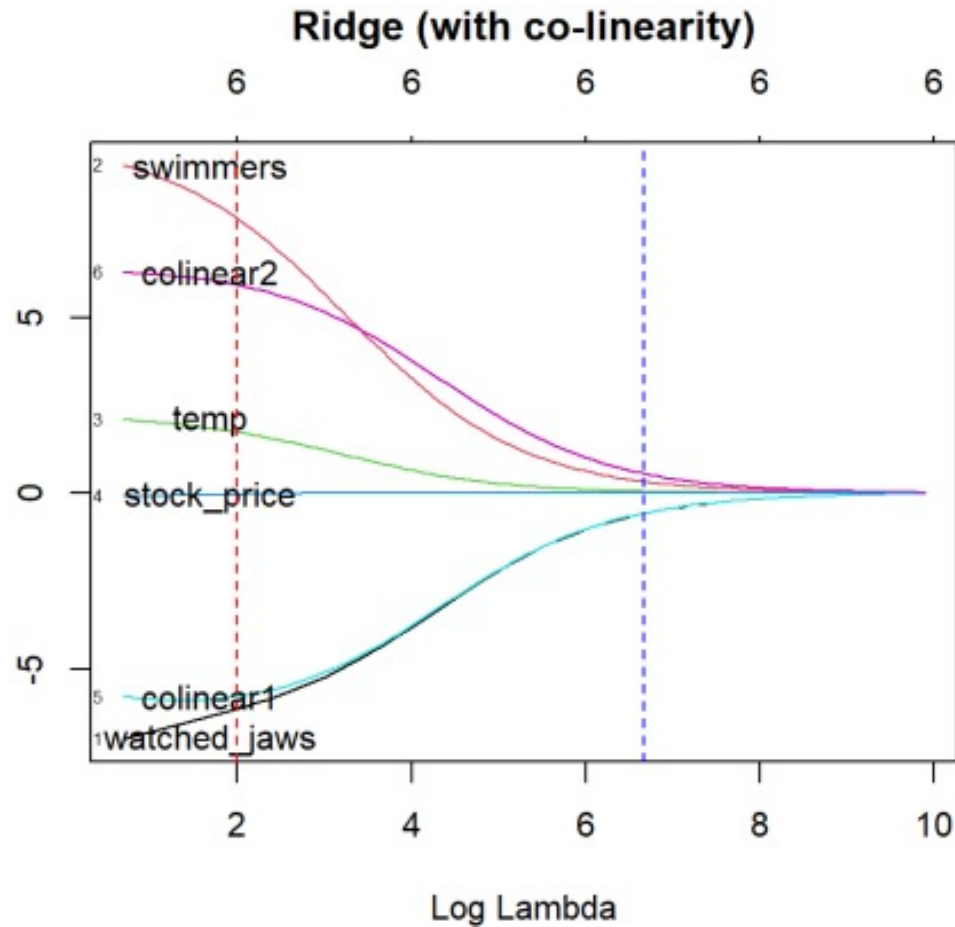# Similarity/Difference between Ridge and LASSO

## Similarity

– Both penalize

– Can be estimated when more variables than observations

## Differences

– The way of shrinking:

- Ridge: Make all beta smaller but rarely gets to 0
- LASSO: Quickly shrink $\beta$ for meaningless variables to 0

So, Ridge is powerful when a lot of weak/meaningful predictors, while LASSO is useful when a lot of junk variables. That's why LASSO is used for variable selection.

# Regularization paths for LASSO and Ridge

# Elastic net

– Elastic net is the combination of Lasso and ridge regression with both L1 and L2 norm

– One formulation is:

$$\text{argmin}_\beta \sum_i (Y_i - (\beta_0 + \sum_j \beta_j X_{ij}))^2 + \lambda(\alpha \sum_j |\beta| + (1-\alpha) \sum_j \beta^2)$$

– Two tuning parameters:

- $\alpha$: weight of L1 and L2
  – $\alpha = 1$: LASSO
  – $\alpha = 0$: Ridge regression

– If tuned well, could perform the best

# Summary

– Regularized regression: Methods to reduce model variance

– Two methods:

- Ridge regression
    – Shrink everything smaller, basically keep all variables

- LASSO regression
    – Variable selection