# Machine Learning for Official Statistics and SDGs

# Regression



UNITED NATIONS
SIAP
Statistical Institute for
Asia and the Pacific

[ LINEAR REGRESSION ]

Multivariate Linear Regression is one of the most popular tool

[ LINEAR REGRESSION ]

Multivariate Linear Regression is one of the most popular tool
  ► Can be used with many variables

## [ LINEAR REGRESSION ]

Multivariate Linear Regression is one of the most popular tool

►  Can be used with many variables
►  Can be used with both continuous or discrete variables
    (categories)

[ LINEAR REGRESSION ]

Multivariate Linear Regression is one of the most popular tool

- ▶ Can be used with many variables
- ▶ Can be used with both continuous or discrete variables (categories)
- ▶ Can be a very efficient tool

## [ LINEAR REGRESSION ]

Multivariate Linear Regression is one of the most popular tool

- ▶ Can be used with many variables
- ▶ Can be used with both continuous or discrete variables (categories)
- ▶ Can be a very efficient tool
- ▶ Has to be well defined, need to verify some hypothesis

# [ LINEAR REGRESSION ]

Multivariate Linear Regression is one of the most popular tool

- ► Can be used with many variables
- ► Can be used with both continuous or discrete variables (categories)
- ► Can be a very efficient tool
- ► Has to be well defined, need to verify some hypothesis
  Expressed as:

$$y = \beta_0 + x'\beta + \varepsilon \qquad E(\varepsilon|x) = 0$$

# [ LINEAR REGRESSION ]

Multivariate Linear Regression is one of the most popular tool

► Can be used with many variables

► Can be used with both continuous or discrete variables (categories)

► Can be a very efficient tool

► Has to be well defined, need to verify some hypothesis
   Expressed as:

$$y = \beta_0 + x'\beta + \varepsilon \qquad E(\varepsilon|x) = 0$$

with possibly many regressors $x_j$

# [ LINEAR REGRESSION: CENTERING VARIABLES]

$$y = \beta_0 + \beta_1 x_1 + \ldots + \beta_k x_k + \varepsilon \qquad E(\varepsilon|x) = 0$$

[ LINEAR REGRESSION: CENTERING VARIABLES]

$$y = \beta_0 + \beta_1 x_1 + \ldots + \beta_k x_k + \varepsilon \qquad E(\varepsilon|x) = 0$$

▶ $\beta_j$ is the "*ceteris paribus*" marginal effect of $x_j$ on $y$.

[ LINEAR REGRESSION: CENTERING VARIABLES]

$$y = \beta_0 + \beta_1 x_1 + \ldots + \beta_k x_k + \varepsilon \qquad E(\varepsilon|x) = 0$$

▶ $\beta_j$ is the "*ceteris paribus*" marginal effect of $x_j$ on $y$.

↪ *when $x_j$ increases by one unit, then $y$ increase by $\beta_j$ units.*

[ LINEAR REGRESSION: CENTERING VARIABLES]

$$y = \beta_0 + \beta_1 x_1 + \ldots + \beta_k x_k + \varepsilon \qquad E(\varepsilon|x) = 0$$

▶ $\beta_j$ is the "*ceteris paribus*" marginal effect of $x_j$ on $y$.
↪ *when $x_j$ increases by one unit, then $y$ increase by $\beta_j$ units.*
▶ $\beta_0$ is the mean of $y$ **if all $x_j$ are equal to zero**

$$\beta_0 = E(y) - \beta_1 E(x_1) - \ldots - \beta_k E(x_k)$$

# [ LINEAR REGRESSION: CENTERING VARIABLES]

$$y = \beta_0 + \beta_1 x_1 + \ldots + \beta_k x_k + \varepsilon \qquad E(\varepsilon|x) = 0$$

▶ $\beta_j$ is the "*ceteris paribus*" marginal effect of $x_j$ on $y$.
↪ *when $x_j$ increases by one unit, then $y$ increase by $\beta_j$ units.*
▶ $\beta_0$ is the mean of $y$ **if all $x_j$ are equal to zero**

$$\beta_0 = E(y) - \beta_1 E(x_1) - \ldots - \beta_k E(x_k)$$

▶ Centering the variables has no effect on the coefficients

$$y = \alpha_0 + \beta_1(x_1 - E(x_1)) + \ldots + \beta_k(x_k - E(x_k)) + \varepsilon \qquad E(\varepsilon|x) = 0$$

# [ LINEAR REGRESSION: CENTERING VARIABLES]

$$y = \beta_0 + \beta_1 x_1 + \ldots + \beta_k x_k + \varepsilon \qquad E(\varepsilon|x) = 0$$

▶ $\beta_j$ is the "*ceteris paribus*" marginal effect of $x_j$ on $y$.

↪ *when $x_j$ increases by one unit, then $y$ increase by $\beta_j$ units.*

▶ $\beta_0$ is the mean of $y$ **if all $x_j$ are equal to zero**

$$\beta_0 = E(y) - \beta_1 E(x_1) - \ldots - \beta_k E(x_k)$$

▶ Centering the variables has no effect on the coefficients

$$y = \alpha_0 + \beta_1(x_1 - E(x_1)) + \ldots + \beta_k(x_k - E(x_k)) + \varepsilon \qquad E(\varepsilon|x) = 0$$

**Except:** $\alpha_0$ is the mean of $y$ **if all $x_j$ are equal their mean**

## [ LINEAR REGRESSION: SCALING ]

We can also scale each variable by its own standard deviation to obtain

$$y = \alpha_0 + \gamma_1 \tilde{x_1} + \ldots + \gamma_k \tilde{x_k} + \varepsilon \qquad E(\varepsilon|x) = 0$$

# [ LINEAR REGRESSION: SCALING ]

We can also scale each variable by its own standard deviation to obtain

$$y = \alpha_0 + \gamma_1 \tilde{x_1} + \ldots + \gamma_k \tilde{x_k} + \varepsilon \qquad E(\varepsilon|x) = 0$$

where

$$\tilde{x} = \frac{x - E(x)}{\sigma_x}$$

# [ LINEAR REGRESSION: SCALING ]

We can also scale each variable by its own standard deviation to obtain

$$y = \alpha_0 + \gamma_1 \tilde{x_1} + \ldots + \gamma_k \tilde{x_k} + \varepsilon \qquad E(\varepsilon | x) = 0$$

where

$$\tilde{x} = \frac{x - E(x)}{\sigma_x}$$

► Now $\gamma_j$ is the *ceteris paribus* marginal effect of $\tilde{x}_j$ on $y$.

## [ LINEAR REGRESSION: SCALING ]

We can also scale each variable by its own standard deviation to obtain

$$y = \alpha_0 + \gamma_1 \tilde{x_1} + \ldots + \gamma_k \tilde{x_k} + \varepsilon \qquad E(\varepsilon|x) = 0$$

where

$$\tilde{x} = \frac{x - E(x)}{\sigma_x}$$

▶ Now $\gamma_j$ is the *ceteris paribus* marginal effect of $\tilde{x}_j$ on $y$.

↪ *when $\tilde{x}_j$ increases by one standard deviation, $y$ increases by $\gamma_j$ units*

## [ LINEAR REGRESSION: SCALING ]

We can also scale each variable by its own standard deviation to obtain

$$y = \alpha_0 + \gamma_1 \tilde{x_1} + \ldots + \gamma_k \tilde{x_k} + \varepsilon \qquad E(\varepsilon|x) = 0$$

where

$$\tilde{x} = \frac{x - E(x)}{\sigma_x}$$

► Now $\gamma_j$ is the *ceteris paribus* marginal effect of $\tilde{x}_j$ on $y$.

$\hookrightarrow$ *when $\tilde{x}_j$ increases by one standard deviation, $y$ increases by $\gamma_j$ units*

► The goal is to have variables and coefficients that are comparable

[ LINEAR REGRESSION: EXAMPLE ]

[ LINEAR REGRESSION: EXAMPLE ]

► Example on a regression model with few variables

|  | Est. | S.E. | t val. | p |
|---|---|---|---|---|
| (Intercept) | 4.390 | 0.136 | 32.204 | 0.000 |
| Trans | -0.002 | 0.001 | -1.603 | 0.110 |
| HighTrans | 0.013 | 0.003 | 3.945 | 0.000 |
| Checks | 0.008 | 0.005 | 1.534 | 0.126 |
| Years | 0.096 | 0.008 | 11.579 | 0.000 |

Initial regression with original values

## [ LINEAR REGRESSION: EXAMPLE ]

▶ Example on a regression model with few variables

|  | Est. | S.E. | t val. | p |
|---|---|---|---|---|
| (Intercept) | 5.927 | 0.039 | 150.243 | 0.000 |
| Trans | -0.002 | 0.001 | -1.603 | 0.110 |
| HighTrans | 0.013 | 0.003 | 3.945 | 0.000 |
| Checks | 0.008 | 0.005 | 1.534 | 0.126 |
| Years | 0.096 | 0.008 | 11.579 | 0.000 |

Regression with centred variables

[ LINEAR REGRESSION: EXAMPLE ]

▶ Example on a regression model with few variables

|  | Est. | S.E. | t val. | p |
|---|---|---|---|---|
| (Intercept) | 5.927 | 0.039 | 150.243 | 0.000 |
| Trans | -0.243 | 0.152 | -1.603 | 0.110 |
| HighTrans | 0.586 | 0.149 | 3.945 | 0.000 |
| Checks | 0.074 | 0.048 | 1.534 | 0.126 |
| Years | 0.462 | 0.040 | 11.579 | 0.000 |

Regression with scaled variables

# [ LINEAR REGRESSION: EXAMPLE ]

▶ Example on a regression model with few variables
The goal is to have *comparable* effects (same range)



Visual regression with scaled variables

[ PROBLEMS IN LINEAR REGRESSION ]

**Collinearity** of regressors is a big issue in regression

## [ PROBLEMS IN LINEAR REGRESSION ]

**Collinearity** of regressors is a big issue in regression

▶ Redundant predictors add more complexity than information

## [ PROBLEMS IN LINEAR REGRESSION ]

**Collinearity** of regressors is a big issue in regression

► Redundant predictors add more complexity than information

► Highly correlated predictors result in unstable estimation

## [ PROBLEMS IN LINEAR REGRESSION ]

**Collinearity** of regressors is a big issue in regression

▶ Redundant predictors add more complexity than information

▶ Highly correlated predictors result in unstable estimation

▶ Highly correlated predictors worsen predictability

## [ PROBLEMS IN LINEAR REGRESSION ]

**Collinearity** of regressors is a big issue in regression

- ▶ Redundant predictors add more complexity than information
- ▶ Highly correlated predictors result in unstable estimation
- ▶ Highly correlated predictors worsen predictability
- ↪ Correlation plot

# [ PROBLEMS IN LINEAR REGRESSION ]

**Collinearity** of regressors is a big issue in regression

► Redundant predictors add more complexity than information

► Highly correlated predictors result in unstable estimation

► Highly correlated predictors worsen predictability

$\hookrightarrow$ Correlation plot

| Introduction | Scaling variables | **Problems** | Classic methods | Penalization Methods | wrap-up |
| :--- | :--- | :--- | :--- | :--- | :--- |
| o | ooo | o●oo | oooooooooo | ooooooooooooo | oo |

[ PROBLEMS IN LINEAR REGRESSION ]

Which variable should be removed?

[ PROBLEMS IN LINEAR REGRESSION ]

Which variable should be removed?

**Variance Inflation Factor**

[ PROBLEMS IN LINEAR REGRESSION ]

Which variable should be removed?

**Variance Inflation Factor**

► Measure of multi-collinearity between variables

## [ PROBLEMS IN LINEAR REGRESSION ]

Which variable should be removed?

**Variance Inflation Factor**

▶ Measure of multi-collinearity between variables

▶ Measure how much the variance of the coefficient of $x_j$ is inflated due to the presence of other regressors.

[ PROBLEMS IN LINEAR REGRESSION ]

Which variable should be removed?

**Variance Inflation Factor**

▶ Measure of multi-collinearity between variables

▶ Measure how much the variance of the coefficient of $x_j$ is inflated due to the presence of other regressors.

▶ VIF for $x_j$ is calculated by running a regression of $x_j$ on all other regressors, computing the $R_j^2$ and use the formula:

# [ PROBLEMS IN LINEAR REGRESSION ]

Which variable should be removed?

### **Variance Inflation Factor**

▶ Measure of multi-collinearity between variables

▶ Measure how much the variance of the coefficient of $x_j$ is inflated due to the presence of other regressors.

▶ VIF for $x_j$ is calculated by running a regression of $x_j$ on all other regressors, computing the $R_j^2$ and use the formula:

$$VIF_j = \frac{1}{1 - R_j^2}$$

[ PROBLEMS IN LINEAR REGRESSION ]

Which variable should be removed?

**Variance Inflation Factor**

▶ Measure of multi-collinearity between variables

▶ Measure how much the variance of the coefficient of $x_j$ is inflated due to the presence of other regressors.

▶ VIF for $x_j$ is calculated by running a regression of $x_j$ on all other regressors, computing the $R_j^2$ and use the formula:

$$VIF_j = \frac{1}{1 - R_j^2}$$

$VIF_j = 1$ indicates no collinearity; a $VIF \geq 10$ is considered as large and problematic

[ SOLUTIONS TO MULTI-COLLINEARITY ]

2 solutions to multi-colinearity:

[ SOLUTIONS TO MULTI-COLLINEARITY ]

2 solutions to multi-colinearity:
► Create new variables from the ones that are collinear

[ SOLUTIONS TO MULTI-COLLINEARITY ]

2 solutions to multi-colinearity:
- ▶ Create new variables from the ones that are collinear
- ↪ using *e.g.* Principal Components Analysis

## [ SOLUTIONS TO MULTI-COLLINEARITY ]

2 solutions to multi-colinearity:
- ▶ Create new variables from the ones that are collinear
- ↪ using *e.g.* Principal Components Analysis
- ▶ Remove some variables

# [ USING *VIF* TO REMOVE COLLINEARITY]

[ USING *VIF* TO REMOVE COLLINEARITY]

Computing VIFs for all $x_j$s

# [ USING *VIF* TO REMOVE COLLINEARITY]

Computing VIFs for all $x_j$s

|  | Est. | S.E. | t val. | p | VIF |
| :--- | :--- | :--- | :--- | :--- | :--- |
| (Intercept) | 5.93 | 0.04 | 150.24 | 0.00 | NA |
| Trans | -0.24 | 0.15 | -1.60 | 0.11 | 14.71 |
| HighTrans | 0.59 | 0.15 | 3.94 | 0.00 | 14.15 |
| Checks | 0.07 | 0.05 | 1.53 | 0.13 | 1.47 |
| Years | 0.46 | 0.04 | 11.58 | 0.00 | 1.02 |

# [ USING *VIF* TO REMOVE COLLINEARITY]

Computing VIFs for all $x_j$s

|  | Est. | S.E. | t val. | p | VIF |
|---|---|---|---|---|---|
| (Intercept) | 5.93 | 0.04 | 150.24 | 0.00 | NA |
| Trans | -0.24 | 0.15 | -1.60 | 0.11 | 14.71 |
| HighTrans | 0.59 | 0.15 | 3.94 | 0.00 | 14.15 |
| Checks | 0.07 | 0.05 | 1.53 | 0.13 | 1.47 |
| Years | 0.46 | 0.04 | 11.58 | 0.00 | 1.02 |

► *Trans* has the highest VIF

| Introduction | Scaling variables | **Problems** | Classic methods | Penalization Methods | wrap-up |
| :--- | :--- | :--- | :--- | :--- | :--- |
| o | ooo | oooo | oooooooooo | oooooooooooo | oo |

[ USING *VIF* TO REMOVE COLLINEARITY]

Conclusion:

# [ USING *VIF* TO REMOVE COLLINEARITY]

Conclusion:

Omitting one variable (*Trans:*)

# [ USING *VIF* TO REMOVE COLLINEARITY]

Conclusion:

Omitting one variable (*Trans:*)

|  | Est. | S.E. | t val. | p | VIF |
|---|---|---|---|---|---|
| (Intercept) | 5.93 | 0.04 | 149.79 | 0.00 | NA |
| HighTrans | 0.36 | 0.05 | 7.69 | 0.00 | 1.40 |
| Checks | 0.06 | 0.05 | 1.24 | 0.22 | 1.41 |
| Years | 0.46 | 0.04 | 11.63 | 0.00 | 1.02 |

[ USING *VIF* TO REMOVE COLLINEARITY]

Conclusion:

Omitting one variable (*Trans:*)

|  | Est. | S.E. | t val. | p | VIF |
|---|---|---|---|---|---|
| (Intercept) | 5.93 | 0.04 | 149.79 | 0.00 | NA |
| HighTrans | 0.36 | 0.05 | 7.69 | 0.00 | 1.40 |
| Checks | 0.06 | 0.05 | 1.24 | 0.22 | 1.41 |
| Years | 0.46 | 0.04 | 11.63 | 0.00 | 1.02 |

► does not change the fit of the model

# [ USING *VIF* TO REMOVE COLLINEARITY]

Conclusion:

Omitting one variable (*Trans:*)

|  | Est. | S.E. | t val. | p | VIF |
| :--- | :--- | :--- | :--- | :--- | :--- |
| (Intercept) | 5.93 | 0.04 | 149.79 | 0.00 | NA |
| HighTrans | 0.36 | 0.05 | 7.69 | 0.00 | 1.40 |
| Checks | 0.06 | 0.05 | 1.24 | 0.22 | 1.41 |
| Years | 0.46 | 0.04 | 11.63 | 0.00 | 1.02 |

► does not change the fit of the model
► does not change the coefficients of the uncorrelated regressors

# [ USING *VIF* TO REMOVE COLLINEARITY]

Conclusion:

Omitting one variable (*Trans:*)

|  | Est. | S.E. | t val. | p | VIF |
|---|---|---|---|---|---|
| (Intercept) | 5.93 | 0.04 | 149.79 | 0.00 | NA |
| HighTrans | 0.36 | 0.05 | 7.69 | 0.00 | 1.40 |
| Checks | 0.06 | 0.05 | 1.24 | 0.22 | 1.41 |
| Years | 0.46 | 0.04 | 11.63 | 0.00 | 1.02 |

► does not change the fit of the model
► does not change the coefficients of the uncorrelated regressors
► reduces all the *VIF*s

## [REAL LIFE EXAMPLE]

In real life, one may have many variables

# [REAL LIFE EXAMPLE]

In real life, one may have many variables



Correlation between all numerical variables

# [REAL LIFE EXAMPLE]

In real life, one may have many variables



Correlation between all numerical variables

↪ *Automatic selection of regressors*

[AUTOMATIC SELECTION OF REGRESSORS]

Classic (but still alive) methods based on the variations of RSS

[AUTOMATIC SELECTION OF REGRESSORS]

Classic (but still alive) methods based on the variations of RSS

► Automatic Forward selection

[AUTOMATIC SELECTION OF REGRESSORS]

Classic (but still alive) methods based on the variations of RSS

▶ Automatic Forward selection

▶ Automatic Backward selection

[AUTOMATIC SELECTION OF REGRESSORS]

Classic (but still alive) methods based on the variations of RSS

- ► Automatic Forward selection
- ► Automatic Backward selection
- ► Stepwise selection

## [AUTOMATIC SELECTION OF REGRESSORS]

Classic (but still alive) methods based on the variations of RSS

► Automatic Forward selection
► Automatic Backward selection
► Stepwise selection

**Remark:**

## [AUTOMATIC SELECTION OF REGRESSORS]

Classic (but still alive) methods based on the variations of RSS

- ► Automatic Forward selection
- ► Automatic Backward selection
- ► Stepwise selection

**Remark:**

- ► The optimal number of regressors is unknown!

## [AUTOMATIC SELECTION OF REGRESSORS]

Classic (but still alive) methods based on the variations of RSS

- ▶ Automatic Forward selection
- ▶ Automatic Backward selection
- ▶ Stepwise selection

**Remark:**

- ▶ The optimal number of regressors is unknown!
- ↪ *The number of possible combinations with k regressors is* $2^k$

## [AUTOMATIC SELECTION OF REGRESSORS]

Classic (but still alive) methods based on the variations of RSS

- ► Automatic Forward selection
- ► Automatic Backward selection
- ► Stepwise selection

**Remark:**

- ► The optimal number of regressors is unknown!
- ↪ *The number of possible combinations with k regressors is* $2^k$
- ► Compute the optimal nb of regressors before testing which regressors to include with Cross Validation

[APPLICATION ON AN EXAMPLE]

*To reduce the computational burden we restrict our choice to 8 variables in the final regression.*

# [Application on an Example]

*To reduce the computational burden we restrict our choice to 8 variables in the final regression.*

**Forward selection variables**



Visual representation of variables used (Forward)

# [APPLICATION ON AN EXAMPLE]

*To reduce the computational burden we restrict our choice to 8 variables in the final regression.*



**Backward selection variables**

Visual representation of variables used (Backward)

[CRITERIONS FOR NUMBER OF REGRESSORS]

[CRITERIONS FOR NUMBER OF REGRESSORS]

▶ Forward & Backward selection provide different solutions:

# [CRITERIONS FOR NUMBER OF REGRESSORS]

▶ Forward & Backward selection provide different solutions:

# [CRITERIONS FOR NUMBER OF REGRESSORS]

► Forward & Backward selection provide different solutions:



► *Great Need for Criteria*

[AUTOMATIC SELECTION OF REGRESSORS]

From the linear model in matrix form:

$$y = X\beta + \varepsilon,$$

[AUTOMATIC SELECTION OF REGRESSORS]

From the linear model in matrix form:

$$y = X\beta + \varepsilon,$$

▶ Consider a partition of $X$ in $X_p$ ($p + 1$ regressors) and $X_r$ ($k - p - 1$ regressors)

## [AUTOMATIC SELECTION OF REGRESSORS]

From the linear model in matrix form:

$$y = X\beta + \varepsilon,$$

▶ Consider a partition of $X$ in $X_p$ ($p + 1$ regressors) and $X_r$ ($k - p - 1$ regressors)

The goal is to find the "*best*" $p$-model

$$y = X_p\beta_p + \varepsilon$$

## [AUTOMATIC SELECTION OF REGRESSORS]

From the linear model in matrix form:

$$y = X\beta + \varepsilon,$$

▶ Consider a partition of $X$ in $X_p$ ($p + 1$ regressors) and $X_r$ ($k - p - 1$ regressors)

The goal is to find the "*best*" $p$-model

$$y = X_p\beta_p + \varepsilon$$

▶ "Best" means best in prediction

## [AUTOMATIC SELECTION OF REGRESSORS]

From the linear model in matrix form:

$$y = X\beta + \varepsilon,$$

▶ Consider a partition of $X$ in $X_p$ ($p + 1$ regressors) and $X_r$ ($k - p - 1$ regressors)

The goal is to find the "*best*" $p$-model

$$y = X_p\beta_p + \varepsilon$$

▶ "Best" means best in prediction

↪ *Mean Squared Error of Prediction* or MSEP:

$$MSEP = n^{-1}E\|y_{new} - X_p\widehat{\beta_p}\|^2$$

## [AUTOMATIC SELECTION OF REGRESSORS]

The MSEP can be decomposed to help define criterions:

$$
\begin{aligned}
MSEP &= n^{-1}E\|y_{new} - X_p\widehat{\beta}_p\|^2 \\
&= n^{-1}\left\{ E\|y_{new} - X\beta\|^2 + E\|X\beta - X_p\widehat{\beta}_p\|^2 \right\} \\
&= (1 + (p+1)/n)\sigma^2 + (1/n)\beta'X'M_pX\beta
\end{aligned}
$$

## [AUTOMATIC SELECTION OF REGRESSORS]

The MSEP can be decomposed to help define criterions:

$$
\begin{aligned}
MSEP &= n^{-1} E\|y_{new} - X_p \widehat{\beta}_p\|^2 \\
&= n^{-1} \left\{ E\|y_{new} - X\beta\|^2 + E\|X\beta - X_p \widehat{\beta}_p\|^2 \right\} \\
&= (1 + (p+1)/n)\sigma^2 + (1/n)\beta' X' M_p X\beta
\end{aligned}
$$

▶ Mallow's $\mathbf{Cp} = \frac{RSS_p}{n} + \frac{2(p+1)}{n} \frac{RSS_k}{n-k-1}$

## [AUTOMATIC SELECTION OF REGRESSORS]

The MSEP can be decomposed to help define criterions:

$$
\begin{aligned}
MSEP &= n^{-1}E\|y_{new} - X_p\widehat{\beta}_p\|^2 \\
&= n^{-1}\left\{E\|y_{new} - X\beta\|^2 + E\|X\beta - X_p\widehat{\beta}_p\|^2\right\} \\
&= (1 + (p+1)/n)\sigma^2 + (1/n)\beta'X'M_pX\beta
\end{aligned}
$$

- Mallow's $\mathbf{Cp} = \frac{RSS_p}{n} + \frac{2(p+1)}{n}\frac{RSS_k}{n-k-1}$
- Akaike Information Criterion (**AIC**) $\propto C_p$ for linear regression

## [AUTOMATIC SELECTION OF REGRESSORS]

The MSEP can be decomposed to help define criterions:

$$
\begin{aligned}
MSEP &= n^{-1}E\|y_{new} - X_p\widehat{\beta}_p\|^2 \\
&= n^{-1}\left\{E\|y_{new} - X\beta\|^2 + E\|X\beta - X_p\widehat{\beta}_p\|^2\right\} \\
&= (1 + (p+1)/n)\sigma^2 + (1/n)\beta'X'M_pX\beta
\end{aligned}
$$

- Mallow's $\mathbf{Cp} = \frac{RSS_p}{n} + \frac{2(p+1)}{n}\frac{RSS_k}{n-k-1}$
- Akaike Information Criterion (**AIC**) $\propto C_p$ for linear regression
- Bayesian Information Criterion (**BIC**):

$$
BIC \propto \frac{RSS_p}{n} + \frac{(p+1)\log n}{n}\frac{RSS_k}{n-k-1}
$$

# [APPLICATION ON THE EXAMPLE]

# [APPLICATION ON THE EXAMPLE]

Selection using Mallow's $C_p$ ($\rightarrow$ 8 variables)

# [APPLICATION ON THE EXAMPLE]

Selection using Mallow's $C_p$ ($\to$ 8 variables)

# [APPLICATION ON THE EXAMPLE]

Selection using BIC ($\rightarrow$ 6 variables)

# [Application on the Example]
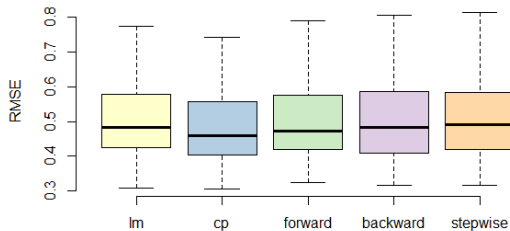
### Selection using BIC ($\rightarrow$ 6 variables)

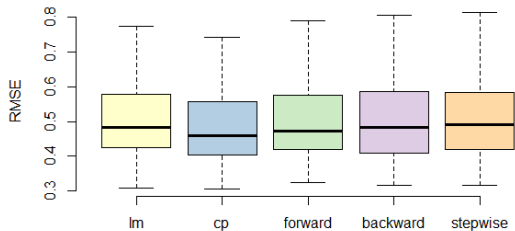## [CROSS VALIDATION ]

Now we can use cross Validation on all models

# [CROSS VALIDATION ]

Now we can use cross Validation on all models

# [CROSS VALIDATION ]

Now we can use cross Validation on all models



► The model selected with $C_p$ is doing the best job
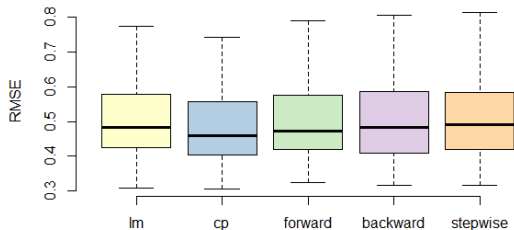
## [CROSS VALIDATION ]

Now we can use cross Validation on all models



▶ The model selected with $C_p$ is doing the best job

↪ **Technical issue**:

# [CROSS VALIDATION ]

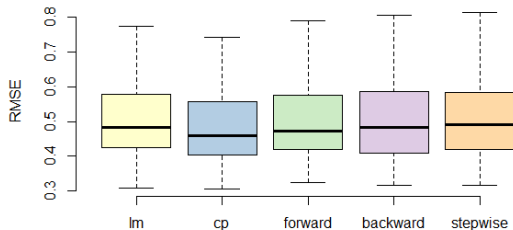Now we can use cross Validation on all models



► The model selected with $C_p$ is doing the best job

↪ **Technical issue**:

   When using CV with stepwise selection different p-models
   are selected for each K-fold

# [CROSS VALIDATION ]

Now we can use cross Validation on all models



- ▶ The model selected with $C_p$ is doing the best job
- ↪ **Technical issue**:

  When using CV with stepwise selection different p-models are selected for each K-fold
- ▶ Modern methods are available...

[QUIZ TIME]

| Introduction | Scaling variables | Problems | Classic methods | Penalization Methods | wrap-up |
| :--- | :--- | :--- | :--- | :--- | :--- |
| o | ooo | oooo | oooooooo●o | oooooooooooo | oo |

[QUIZ TIME]

[TAKEAWAYS]

- In linear regression, scaling allows to compare coefficients and to measure variable importance.

[TAKEAWAYS]

▶ In linear regression, scaling allows to compare coefficients and to measure variable importance.
▶ Multi-collinearity should be investigated beforehand.

# [TAKEAWAYS]

- ▶ In linear regression, scaling allows to compare coefficients and to measure variable importance.
- ▶ Multi-collinearity should be investigated beforehand.
- ▶ Mallows $C_p$ is a simple method to select regressors

# [TAKEAWAYS]

- ▶ In linear regression, scaling allows to compare coefficients and to measure variable importance.
- ▶ Multi-collinearity should be investigated beforehand.
- ▶ Mallows $C_p$ is a simple method to select regressors
- ▶ Stepwise methods with CV have three drawbacks:

## [TAKEAWAYS]

- ▶ In linear regression, scaling allows to compare coefficients and to measure variable importance.
- ▶ Multi-collinearity should be investigated beforehand.
- ▶ Mallows $C_p$ is a simple method to select regressors
- ▶ Stepwise methods with CV have three drawbacks:
    - ▶ they do not necessarily select the "best" model

## [TAKEAWAYS]

▶ In linear regression, scaling allows to compare coefficients and to measure variable importance.

▶ Multi-collinearity should be investigated beforehand.

▶ Mallows $C_p$ is a simple method to select regressors

▶ Stepwise methods with CV have three drawbacks:
  ▶ they do not necessarily select the "best" model
  ▶ the choice of variables can be sensitive to the number of repetitions

# [TAKEAWAYS]

- ▶ In linear regression, scaling allows to compare coefficients and to measure variable importance.
- ▶ Multi-collinearity should be investigated beforehand.
- ▶ Mallows $C_p$ is a simple method to select regressors
- ▶ Stepwise methods with CV have three drawbacks:
    - ▶ they do not necessarily select the "best" model
    - ▶ the choice of variables can be sensitive to the number of repetitions
    - ▶ Variable selection/elimination is done variable by variable (no interactions)

# [TAKEAWAYS]

- ▶ In linear regression, scaling allows to compare coefficients and to measure variable importance.
- ▶ Multi-collinearity should be investigated beforehand.
- ▶ Mallows $C_p$ is a simple method to select regressors
- ▶ Stepwise methods with CV have three drawbacks:
  - ▶ they do not necessarily select the "best" model
  - ▶ the choice of variables can be sensitive to the number of repetitions
  - ▶ Variable selection/elimination is done variable by variable (no interactions)
- ▶ Other (modern) methods exist

[PENALIZATION METHODS]

[PENALIZATION METHODS]

► Problems when many variables are available:

## [PENALIZATION METHODS]

- ▶ Problems when many variables are available:
  - ▶ Multi-collinearity

## [PENALIZATION METHODS]

- ▶ Problems when many variables are available:
    - ▶ Multi-collinearity
    - ▶ Model complexity

## [PENALIZATION METHODS]

▶ Problems when many variables are available:
  ▶ Multi-collinearity
  ▶ Model complexity
  ▶ High variance of the estimator

## [PENALIZATION METHODS]

- ▶ Problems when many variables are available:
    - ▶ Multi-collinearity
    - ▶ Model complexity
    - ▶ High variance of the estimator
- ↪ Methods *"penalizing"* model complexity (over fitting)

## [PENALIZATION METHODS]

- ► Problems when many variables are available:
    - ► Multi-collinearity
    - ► Model complexity
    - ► High variance of the estimator
- ↪ Methods *"penalizing"* model complexity (over fitting)
- ► Have intensionally a **higher bias** and a **lower variance**

# [PENALIZATION METHODS]

- ▶ Problems when many variables are available:
    - ▶ Multi-collinearity
    - ▶ Model complexity
    - ▶ High variance of the estimator
- ↪ Methods *"penalizing"* model complexity (over fitting)
- ▶ Have intensionally a **higher bias** and a **lower variance**
- ▶ Solutions of a *penalized least-squares problem*

$$\min_{\beta} \frac{1}{2n} \sum_{i=1}^{n} \left(y_i - \beta_0 - x_i'\beta\right)^2 + \lambda \cdot J(\beta_1, \cdots, \beta_k)$$

# [PENALIZATION METHODS]

- ▶ Problems when many variables are available:
    - ▶ Multi-collinearity
    - ▶ Model complexity
    - ▶ High variance of the estimator
- ↪ Methods *"penalizing"* model complexity (over fitting)
- ▶ Have intensionally a **higher bias** and a **lower variance**
- ▶ Solutions of a *penalized least-squares problem*

$$\min_{\beta} \frac{1}{2n} \sum_{i=1}^{n} \left( y_i - \beta_0 - x_i'\beta \right)^2 + \lambda \cdot J(\beta_1, \cdots, \beta_k)$$

$\lambda$ is a *hyper parameter*, $J(\cdot)$ is the penalization function

## [PENALIZATION METHODS]

- ▶ Problems when many variables are available:
    - ▶ Multi-collinearity
    - ▶ Model complexity
    - ▶ High variance of the estimator
- ↪ Methods *"penalizing"* model complexity (over fitting)
- ▶ Have intensionally a **higher bias** and a **lower variance**
- ▶ Solutions of a *penalized least-squares problem*

$$\min_{\beta} \frac{1}{2n} \sum_{i=1}^{n} \left( y_i - \beta_0 - x_i'\beta \right)^2 + \lambda \cdot J(\beta_1, \cdots, \beta_k)$$

$\lambda$ is a *hyper parameter*, $J(\cdot)$ is the penalization function

NB: *If $\lambda = 0$ the regression is just the classic OLS estimator*

## [PENALIZATION METHODS: RIDGE REGRESSION]

Ridge regression is the solution of:

$$\min_{\beta} \frac{1}{2n} \sum_{i=1}^{n} \left( y_i - \beta_0 - x_i'\beta \right)^2 + \lambda \frac{\|\beta\|^2}{2}$$

[PENALIZATION METHODS: RIDGE REGRESSION]

Ridge regression is the solution of:

$$\min_{\beta} \frac{1}{2n} \sum_{i=1}^{n} \left( y_i - \beta_0 - x_i'\beta \right)^2 + \lambda \frac{\|\beta\|^2}{2}$$

▶ Ridge regression shrinks parameters towards zero and thus avoids too large parameters

## [PENALIZATION METHODS: RIDGE REGRESSION]

Ridge regression is the solution of:

$$\min_{\beta} \frac{1}{2n} \sum_{i=1}^{n} \left( y_i - \beta_0 - x_i'\beta \right)^2 + \lambda \frac{\|\beta\|^2}{2}$$

► Ridge regression shrinks parameters towards zero and thus avoids too large parameters

► The solution is biased intensionally to reduce variance

# [PENALIZATION METHODS: RIDGE REGRESSION]

Ridge regression is the solution of:

$$\min_{\beta} \frac{1}{2n} \sum_{i=1}^{n} \left(y_i - \beta_0 - x_i'\beta\right)^2 + \lambda \frac{\|\beta\|^2}{2}$$

- ▶ Ridge regression shrinks parameters towards zero and thus avoids too large parameters
- ▶ The solution is biased intensionally to reduce variance
- ▶ **Penalization** methods prevent these problems by *"penalizing"* model complexity

## [PENALIZATION METHODS: RIDGE REGRESSION]

Ridge regression is the solution of:

$$\min_{\beta} \frac{1}{2n} \sum_{i=1}^{n} \left( y_i - \beta_0 - x_i'\beta \right)^2 + \lambda \frac{\|\beta\|^2}{2}$$

- ▶ Ridge regression shrinks parameters towards zero and thus avoids too large parameters
- ▶ The solution is biased intensionally to reduce variance
- ▶ **Penalization** methods prevent these problems by *"penalizing"* model complexity
- ▶ It is important to **center and scale** each of the $x$ to ensure the comparability of $\beta$s

[PENALIZATION METHODS: RIDGE REGRESSION]

Ridge regression is the solution of:

$$\min_{\beta} \frac{1}{2n} \sum_{i=1}^{n} \left( y_i - \beta_0 - x_i'\beta \right)^2 + \lambda \frac{\|\beta\|^2}{2}$$

## [PENALIZATION METHODS: RIDGE REGRESSION]

Ridge regression is the solution of:

$$\min_{\beta} \frac{1}{2n} \sum_{i=1}^{n} \left( y_i - \beta_0 - x_i'\beta \right)^2 + \lambda \frac{\|\beta\|^2}{2}$$

▶ The solution (ridge estimator) is $\widehat{\beta_R} = \left( \frac{X'X}{n} + \lambda I \right)^{-1} \frac{X'y}{n}$

# [PENALIZATION METHODS: RIDGE REGRESSION]

Ridge regression is the solution of:

$$\min_{\beta} \frac{1}{2n} \sum_{i=1}^{n} \left(y_i - \beta_0 - x_i'\beta\right)^2 + \lambda \frac{\|\beta\|^2}{2}$$

▶ The solution (ridge estimator) is $\widehat{\beta_R} = \left(\frac{X'X}{n} + \lambda I\right)^{-1} \frac{X'y}{n}$

▶ To compare with OLS estimator $\widehat{\beta_{OLS}} = \left(\frac{X'X}{n}\right)^{-1} \frac{X'y}{n}$

## [PENALIZATION METHODS: RIDGE REGRESSION]

Ridge regression is the solution of:

$$\min_{\beta} \frac{1}{2n} \sum_{i=1}^{n} \left(y_i - \beta_0 - x_i'\beta\right)^2 + \lambda \frac{\|\beta\|^2}{2}$$

▶ The solution (ridge estimator) is $\widehat{\beta_R} = \left(\frac{X'X}{n} + \lambda I\right)^{-1} \frac{X'y}{n}$

▶ To compare with OLS estimator $\widehat{\beta_{OLS}} = \left(\frac{X'X}{n}\right)^{-1} \frac{X'y}{n}$

↪ When there is collinearity, $X'X$ cannot be inverted, while if $\lambda > 0$ the matrix $(X'X + \lambda I)$ is invertible

# [RIDGE REGRESSION IN PRACTICE]

How to choose $\lambda$?

[RIDGE REGRESSION IN PRACTICE]

How to choose $\lambda$?

▶ We compute the (CV-averaged) *RMSE* for many values of $\lambda$

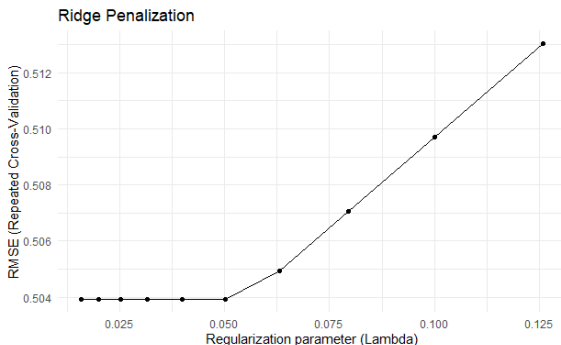# [RIDGE REGRESSION IN PRACTICE]

How to choose $\lambda$?

► We compute the (CV-averaged) *RMSE* for many values of $\lambda$

## [RIDGE REGRESSION IN PRACTICE]

How to choose $\lambda$?

▶ We compute the (CV-averaged) *RMSE* for many values of $\lambda$



*Remark: $\lambda$ is typically small, and we usually select it on the log scale.*

| Introduction | Scaling variables | Problems | Classic methods | Penalization Methods | wrap-up |
|:---|:---|:---|:---|:---|:---|
| o | ooo | oooo | oooooooooo | oooo●oooooo | oo |

[RIDGE REGRESSION IN PRACTICE]

Ridge regression with optimal $\lambda^*$

[RIDGE REGRESSION IN PRACTICE]

Ridge regression with optimal $\lambda^*$

▶ Variance inflation factor (VIF) for this model

# [RIDGE REGRESSION IN PRACTICE]

Ridge regression with optimal $\lambda^*$

► Variance inflation factor (VIF) for this model



VIF (Ridge Penalization)

## [PENALIZATION METHODS: LASSO]

LASSO or *Least Absolute Shrinkage and Selection Operator*, is
another common penalization method, is the solution of:

$$\min_{\beta} \frac{1}{2n} \sum_{i=1}^{n} \left(y_i - \beta_0 - x_i'\beta\right)^2 + \lambda|\beta| \qquad |\beta| = \sum_{j=1}^{k} |\beta_j|$$

## [PENALIZATION METHODS: LASSO]

LASSO or *Least Absolute Shrinkage and Selection Operator*, is another common penalization method, is the solution of:

$$\min_{\beta} \frac{1}{2n} \sum_{i=1}^{n} \left( y_i - \beta_0 - x_i'\beta \right)^2 + \lambda|\beta| \qquad |\beta| = \sum_{j=1}^{k} |\beta_j|$$

$\hookrightarrow$ If $\lambda = 0$, we obtain OLS. If $\lambda = \infty$, all parameters are zero.

## [PENALIZATION METHODS: LASSO]

LASSO or *Least Absolute Shrinkage and Selection Operator*, is another common penalization method, is the solution of:

$$\min_{\beta} \frac{1}{2n} \sum_{i=1}^{n} \left( y_i - \beta_0 - x_i'\beta \right)^2 + \lambda|\beta| \qquad |\beta| = \sum_{j=1}^{k} |\beta_j|$$
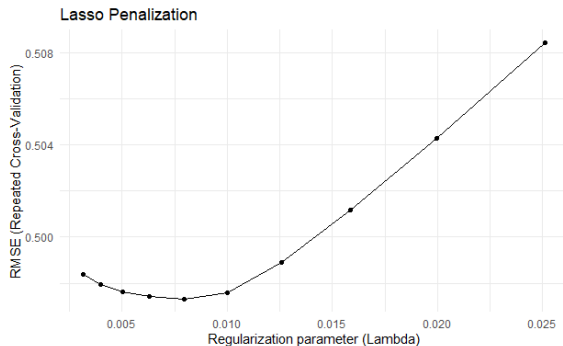
↪ If $\lambda = 0$, we obtain OLS. If $\lambda = \infty$, all parameters are zero.

▶ Lasso does automatic variable selection: if $\lambda$ is large enough, the solution put some parameters to zero

## [PENALIZATION METHODS: LASSO]

LASSO or *Least Absolute Shrinkage and Selection Operator*, is another common penalization method, is the solution of:

$$\min_{\beta} \frac{1}{2n} \sum_{i=1}^{n} \left(y_i - \beta_0 - x_i'\beta\right)^2 + \lambda|\beta| \qquad |\beta| = \sum_{j=1}^{k} |\beta_j|$$

$\hookrightarrow$ If $\lambda = 0$, we obtain OLS. If $\lambda = \infty$, all parameters are zero.

▶ Lasso does automatic variable selection: if $\lambda$ is large enough, the solution put some parameters to zero

▶ It is important to **center and scale** each of the $x$ to ensure the comparability of $\beta$s

[LASSO IN PRACTICE]

We compute the (CV-averaged) *RMSE* for many values of $\lambda$

# [LASSO IN PRACTICE]

We compute the (CV-averaged) *RMSE* for many values of $\lambda$

# [LASSO IN PRACTICE]

We compute the (CV-averaged) *RMSE* for many values of $\lambda$



*Remark: $\lambda$ is typically small, and we usually select it on the log scale.*

# [LASSO IN PRACTICE]

Regression with optimal $\lambda^*$ for LASSO
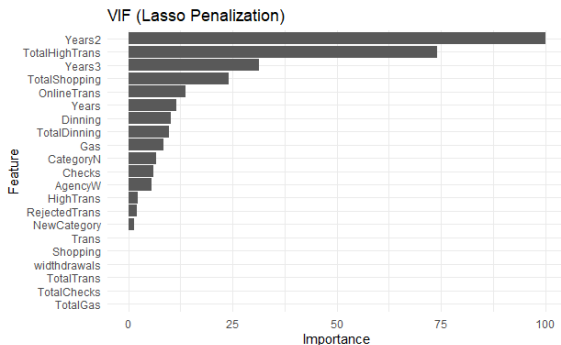
## [LASSO IN PRACTICE]

Regression with optimal $\lambda^*$ for LASSO

▶ Variance inflation factor (VIF) for this model

# [LASSO IN PRACTICE]

Regression with optimal $\lambda^*$ for LASSO

▶ Variance inflation factor (VIF) for this model

## [PENALIZATION METHODS: ELASTIC NET]

*Elastic Net* combines both Lasso and Ridge regression:

$$\min_{\beta} \frac{1}{2n} \sum_{i=1}^{n} \left(y_i - \beta_0 - x_i'\beta\right)^2 + \lambda \left((1-\alpha)\frac{\|\beta\|^2}{2} + \alpha|\beta|\right)$$
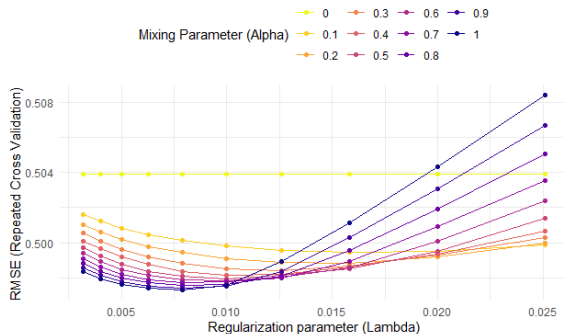
[PENALIZATION METHODS: ELASTIC NET]

*Elastic Net* combines both Lasso and Ridge regression:

$$\min_{\beta} \frac{1}{2n} \sum_{i=1}^{n} \left(y_i - \beta_0 - x_i'\beta\right)^2 + \lambda \left((1-\alpha)\frac{\|\beta\|^2}{2} + \alpha|\beta|\right)$$

► If $\alpha = 1$ we have the Lasso estimator, if $\alpha = 0$, the Ridge regression.

[PENALIZATION METHODS: ELASTIC NET]

*Elastic Net* combines both Lasso and Ridge regression:

$$\min_{\beta} \frac{1}{2n} \sum_{i=1}^{n} \left( y_i - \beta_0 - x_i'\beta \right)^2 + \lambda \left( (1 - \alpha)\frac{\|\beta\|^2}{2} + \alpha|\beta| \right)$$

► If $\alpha = 1$ we have the Lasso estimator, if $\alpha = 0$, the Ridge regression.

► Through $\alpha$ we balance variable selection (Lasso) and coefficient reduction (Ridge)

[ELASTIC NET IN PRACTICE]

Compute the (CV-averaged) *RMSE* on a grid of $(\lambda , \alpha)$

# [ELASTIC NET IN PRACTICE]

Compute the (CV-averaged) *RMSE* on a grid of $(\lambda, \alpha)$

# [ELASTIC NET IN PRACTICE]

Compute the (CV-averaged) *RMSE* on a grid of $(\lambda, \alpha)$



$\hookrightarrow$ Optimal value: $\alpha^* = 1$ & $\lambda = 0.008 \hookrightarrow$ Lasso estimator.

# [ELASTIC NET IN PRACTICE]

Compute the (CV-averaged) *RMSE* on a grid of $(\lambda, \alpha)$



↪ Optimal value: $\alpha^* = 1$ & $\lambda = 0.008$ ↪ Lasso estimator.

▶ Elastic net *encompass* both Ridge and Lasso estimators.

## [ELASTIC NET IN PRACTICE]

Elastic net with optimal $\lambda^*$ and $\alpha^*$ is LASSO since $\alpha^* = 1$!

## [ELASTIC NET IN PRACTICE]

Elastic net with optimal $\lambda^*$ and $\alpha^*$ is LASSO since $\alpha^* = 1$!

▶ Variance inflation factor (VIF) for Elastic net

# [ELASTIC NET IN PRACTICE]

Elastic net with optimal $\lambda^*$ and $\alpha^*$ is LASSO since $\alpha^* = 1$!

▶ Variance inflation factor (VIF) for Elastic net

## [BEST MODEL?]

In Machine Learning, focus on prediction (RMSE)
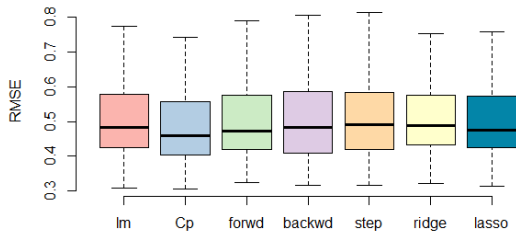
## [BEST MODEL?]

In Machine Learning, focus on prediction (RMSE)

► Cross-Validation performance (RMSE):

# [BEST MODEL?]

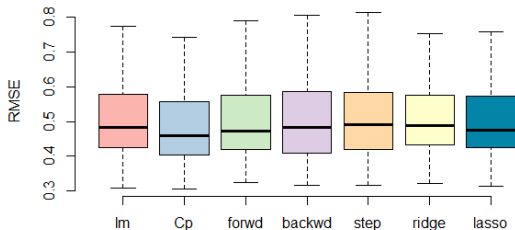In Machine Learning, focus on prediction (RMSE)

- ▶ Cross-Validation performance (RMSE):

## [BEST MODEL?]

In Machine Learning, focus on prediction (RMSE)

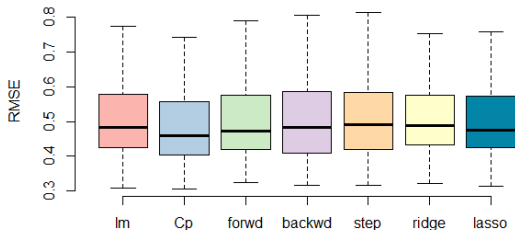► Cross-Validation performance (RMSE):



► LASSO is probably the best (lower RMSE)

## [BEST MODEL?]

In Machine Learning, focus on prediction (RMSE)

► Cross-Validation performance (RMSE):



► LASSO is probably the best (lower RMSE)
► Model selected by Mallow's $C_p$ (**8 regressors**) is almost as good!

# [QUIZ TIME]

# [QUIZ TIME]

## [TAKEAWAYS]

▶ Penalized least-squares methods can be used with multicollinearity or with a large number of regressors.

## [TAKEAWAYS]

- ▶ Penalized least-squares methods can be used with multicollinearity or with a large number of regressors.
- ▶ All solutions of a *penalized least-squares problem*

$$\min_{\beta} \frac{1}{2n} \sum_{i=1}^{n} \left( y_i - \beta_0 - x_i'\beta \right)^2 + \lambda \cdot J_\alpha(\beta_1, \cdots, \beta_k)$$

## [TAKEAWAYS]

▶ Penalized least-squares methods can be used with multicollinearity or with a large number of regressors.

▶ All solutions of a *penalized least-squares problem*

$$\min_{\beta} \frac{1}{2n} \sum_{i=1}^{n} \left(y_i - \beta_0 - x_i'\beta\right)^2 + \lambda \cdot J_\alpha(\beta_1, \cdots, \beta_k)$$

$\lambda$ is a *hyperparameter*, $J_\alpha(\cdot)$ is the penalization function

## [TAKEAWAYS]

▶ Penalized least-squares methods can be used with multicollinearity or with a large number of regressors.

▶ All solutions of a *penalized least-squares problem*

$$\min_{\beta} \frac{1}{2n} \sum_{i=1}^{n} \left( y_i - \beta_0 - x_i'\beta \right)^2 + \lambda \cdot J_\alpha(\beta_1, \cdots, \beta_k)$$

$\lambda$ is a *hyperparameter*, $J_\alpha(\cdot)$ is the penalization function

▶ It is important to **center and scale** each of the $x$ to ensure the comparability of $\beta$s

## [TAKEAWAYS]

▶ Penalized least-squares methods can be used with multicollinearity or with a large number of regressors.

▶ All solutions of a *penalized least-squares problem*

$$\min_{\beta} \frac{1}{2n} \sum_{i=1}^{n} \left( y_i - \beta_0 - x_i'\beta \right)^2 + \lambda \cdot J_\alpha(\beta_1, \cdots, \beta_k)$$

$\lambda$ is a *hyperparameter*, $J_\alpha(\cdot)$ is the penalization function

▶ It is important to **center and scale** each of the $x$ to ensure the comparability of $\beta$s

▶ Selection of *hyper parameters* is based on CV and RMSE

## [TAKEAWAYS]

▶ Penalized least-squares methods can be used with multicollinearity or with a large number of regressors.

▶ All solutions of a *penalized least-squares problem*

$$\min_{\beta} \frac{1}{2n} \sum_{i=1}^{n} \left( y_i - \beta_0 - x_i'\beta \right)^2 + \lambda \cdot J_\alpha(\beta_1, \cdots, \beta_k)$$

$\lambda$ is a *hyperparameter*, $J_\alpha(\cdot)$ is the penalization function

▶ It is important to **center and scale** each of the $x$ to ensure the comparability of $\beta$s

▶ Selection of *hyper parameters* is based on CV and RMSE

▶ The elastic net "encompass" both Ridge and Lasso estimators.

## [TAKEAWAYS]

▶ Penalized least-squares methods can be used with multicollinearity or with a large number of regressors.

▶ All solutions of a *penalized least-squares problem*

$$\min_{\beta} \frac{1}{2n} \sum_{i=1}^{n} \left( y_i - \beta_0 - x_i'\beta \right)^2 + \lambda \cdot J_\alpha(\beta_1, \cdots, \beta_k)$$

$\lambda$ is a *hyperparameter*, $J_\alpha(\cdot)$ is the penalization function

▶ It is important to **center and scale** each of the $x$ to ensure the comparability of $\beta$s

▶ Selection of *hyper parameters* is based on CV and RMSE

▶ The elastic net "encompass" both Ridge and Lasso estimators.

▶ Penalization methods are very popular in practice