# Regression analysis on oceanographic dataset
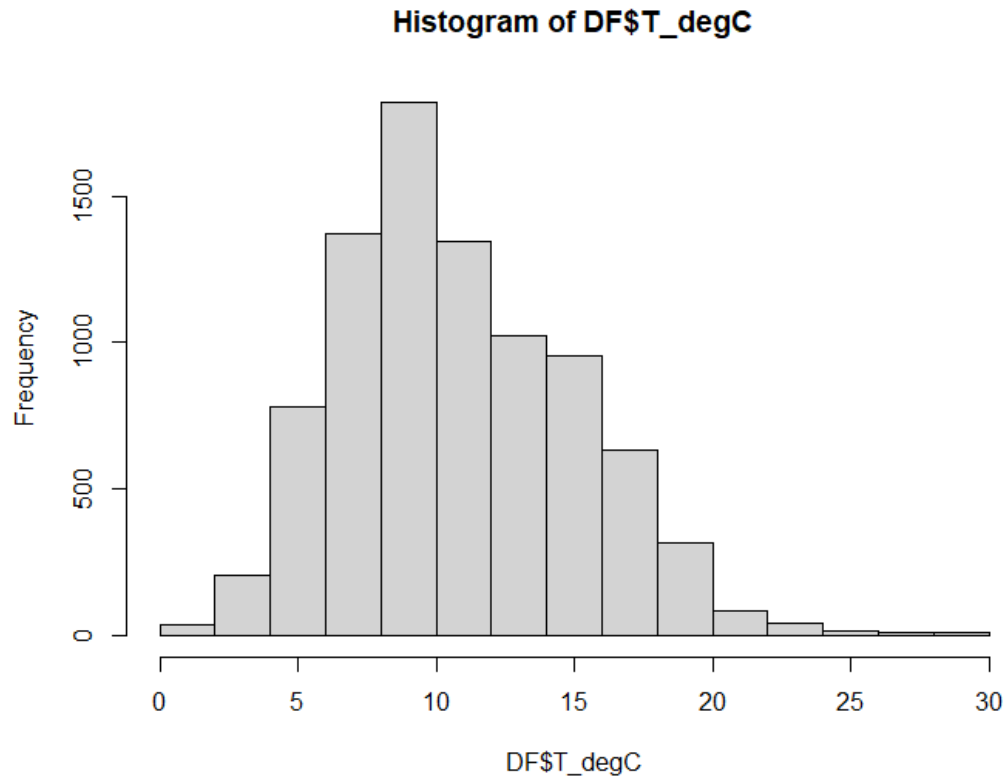
Group members: Diana Gorshechnikova and Tatiana Parashina

STA 9890

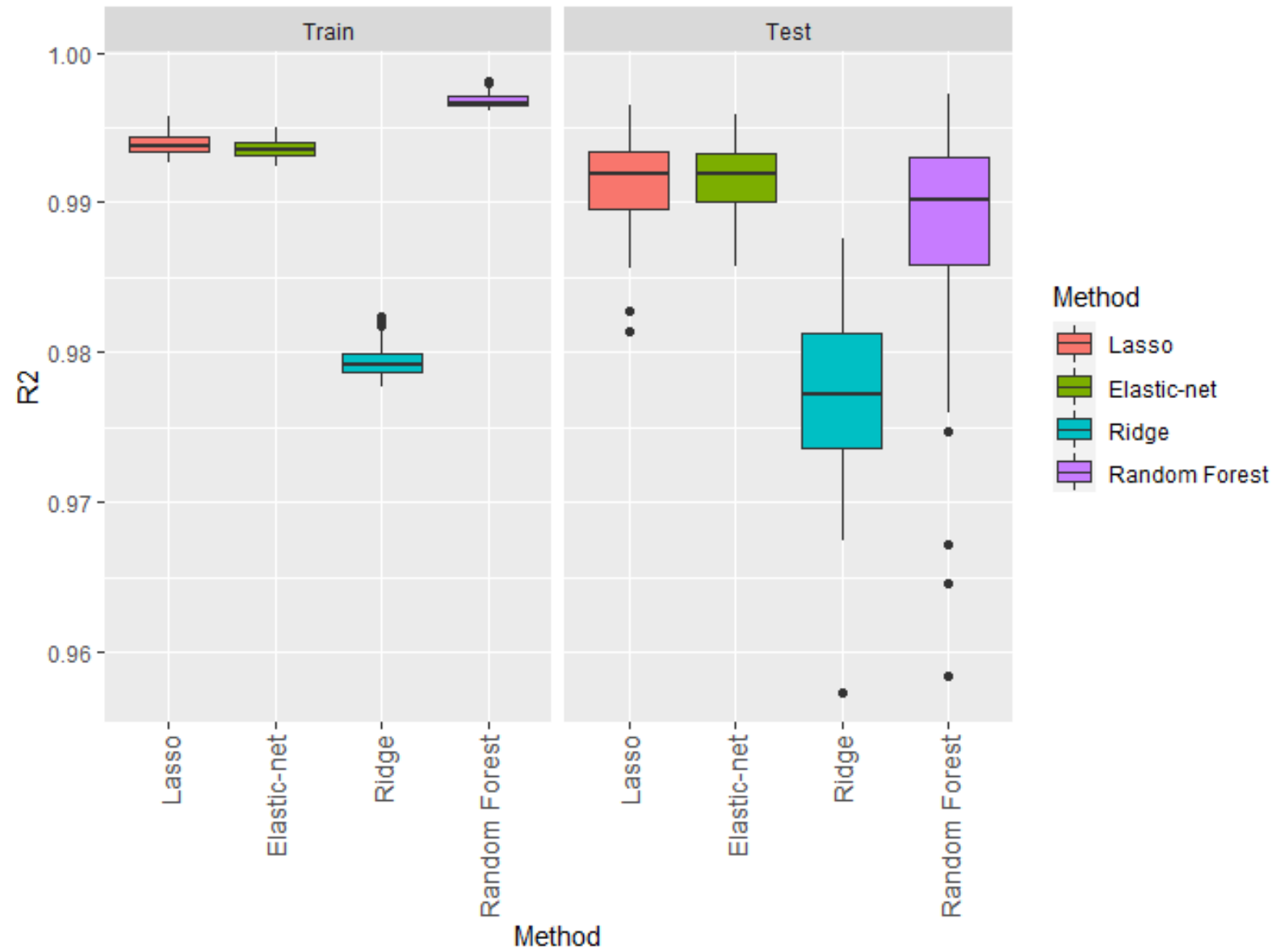# The California Cooperative Oceanic Fisheries Investigations (CalCOFI) dataset

- Original dataset has 74 columns and 864863 observations;

- Response variable is water temperature;

- Increasing ocean temperatures severely affect marine species and ecosystems;

- Rising temperatures can contribute to coral bleaching and the loss of breeding grounds for marine fishes and mammals;

- Machine learning can be useful to predict what contributes to water temperature increase and to mitigate the rising temperatures in a timely fashion.
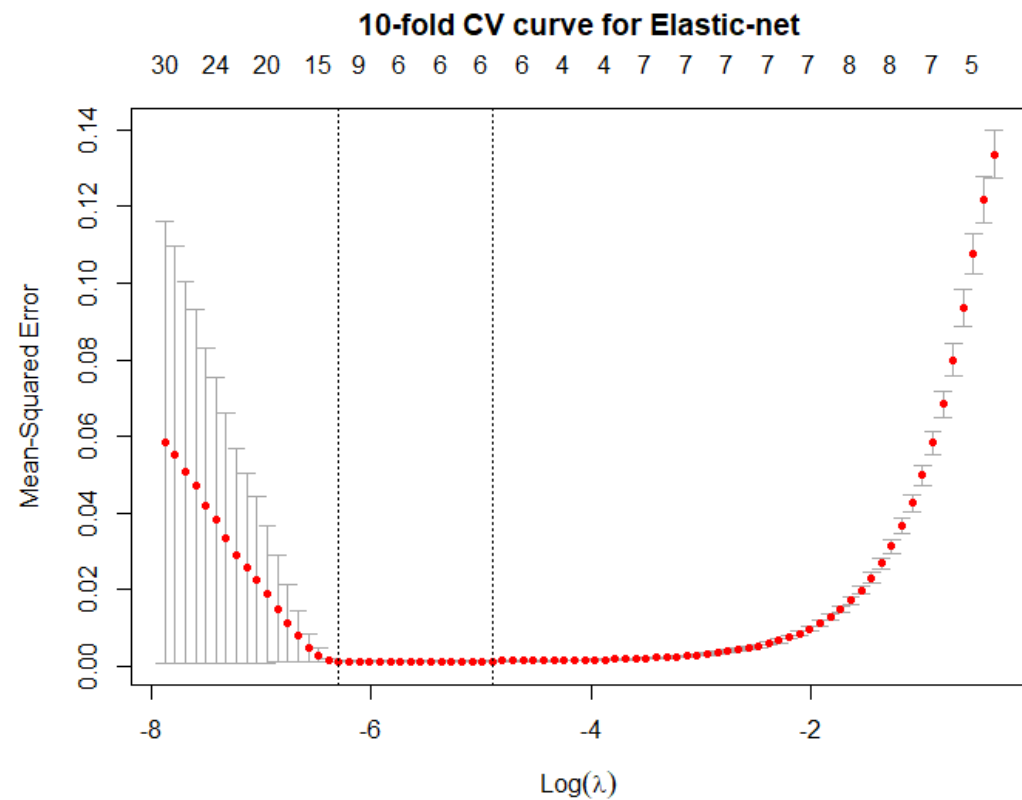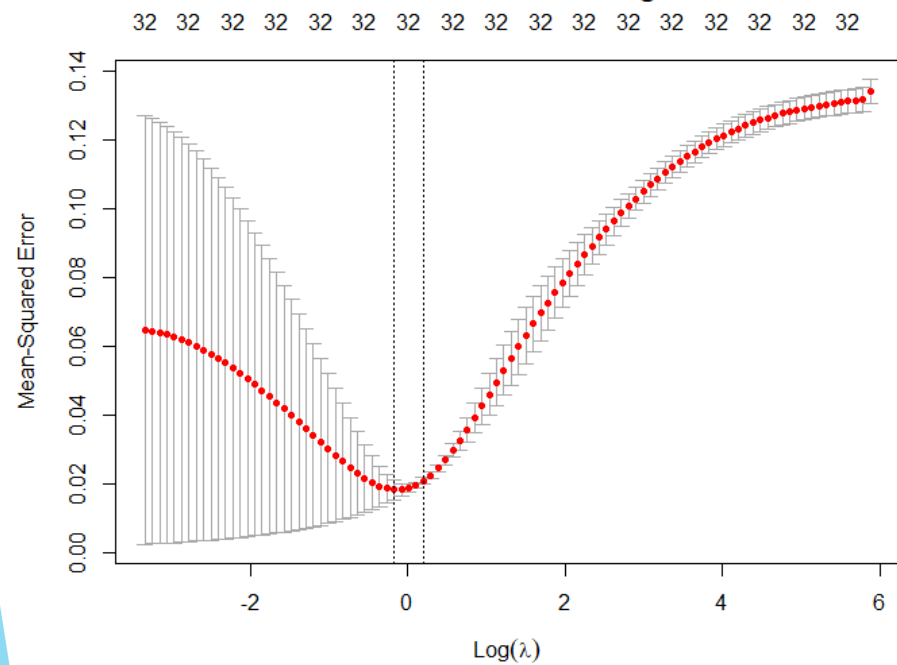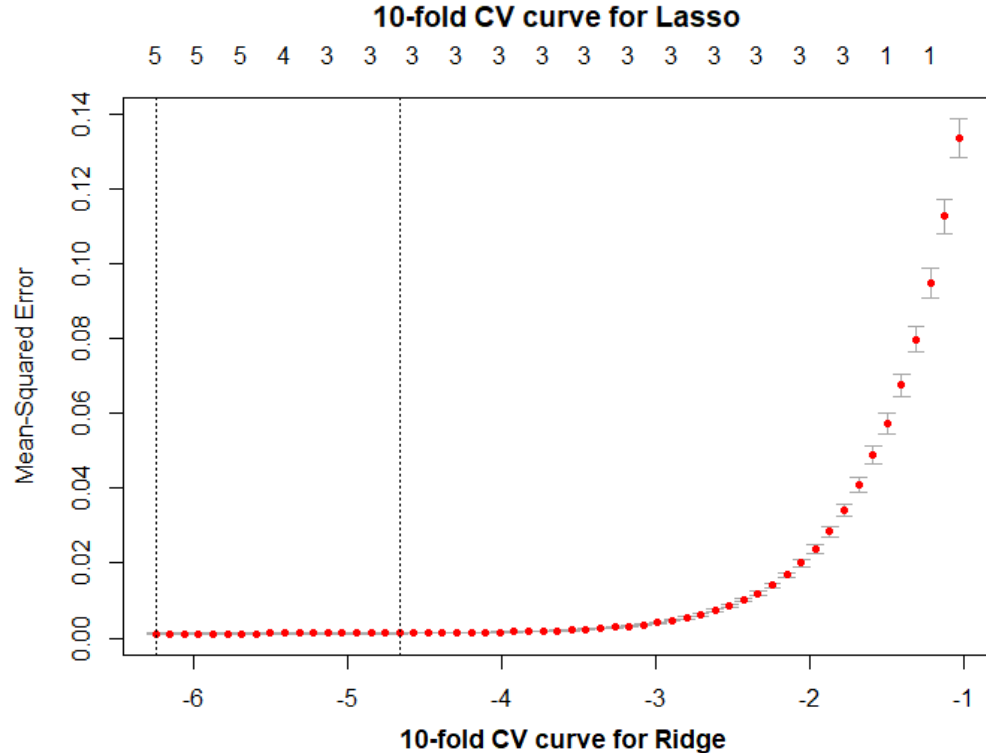
| Minimum t | Maximum t | Mean t |
|-----------|-----------|--------|
| 1.44 C | 31.14 C | 10.8 C |


Histogram of DF$T_degC

- The predictors are: salinity, oxygen, phosphate, silicate, nitrate and nitrite, chlorophyll, transmissometer, PAR, C14 primary productivity, phytoplankton biodiversity, zooplankton biomass, zooplankton biodiversity, etc.
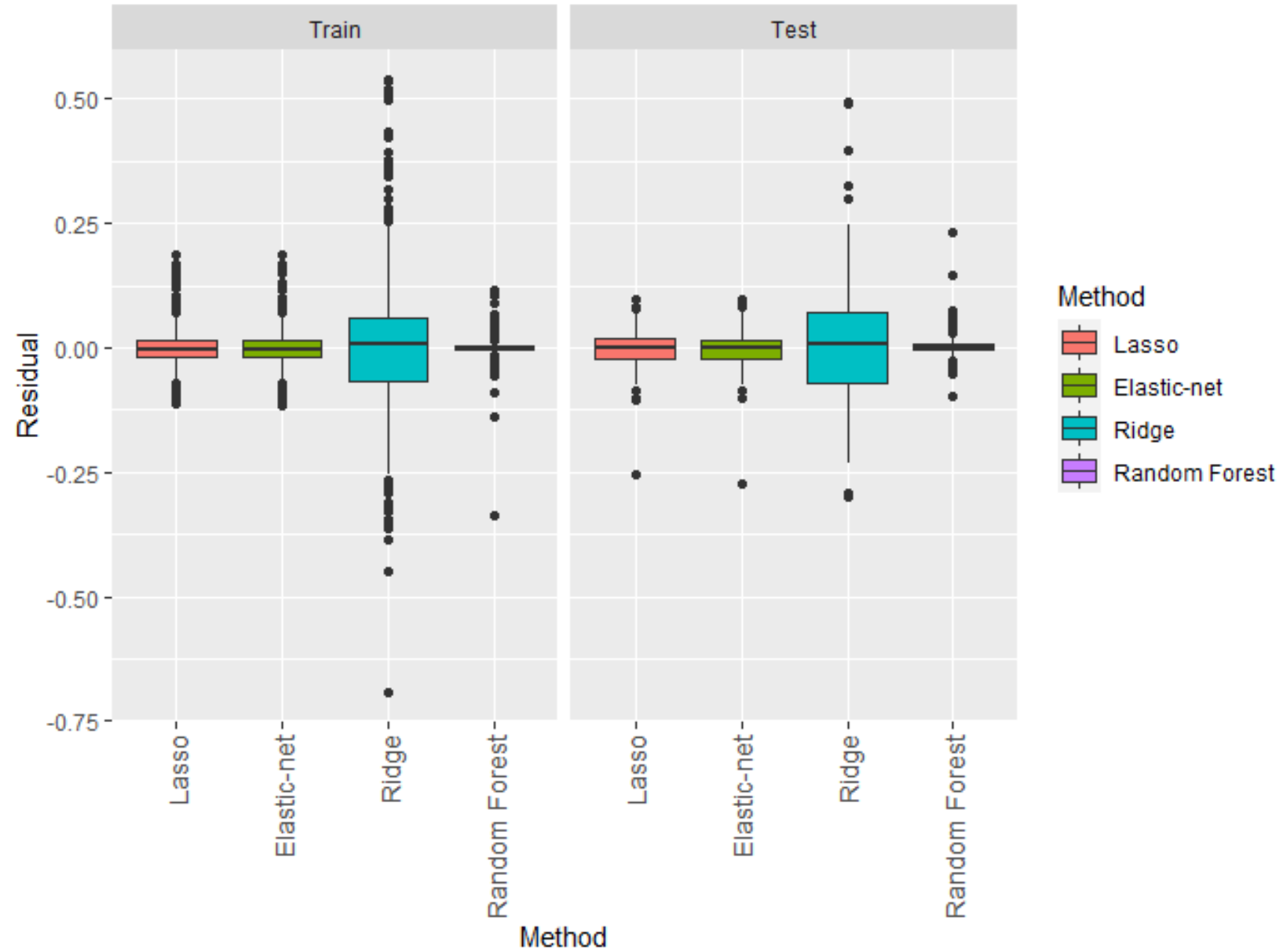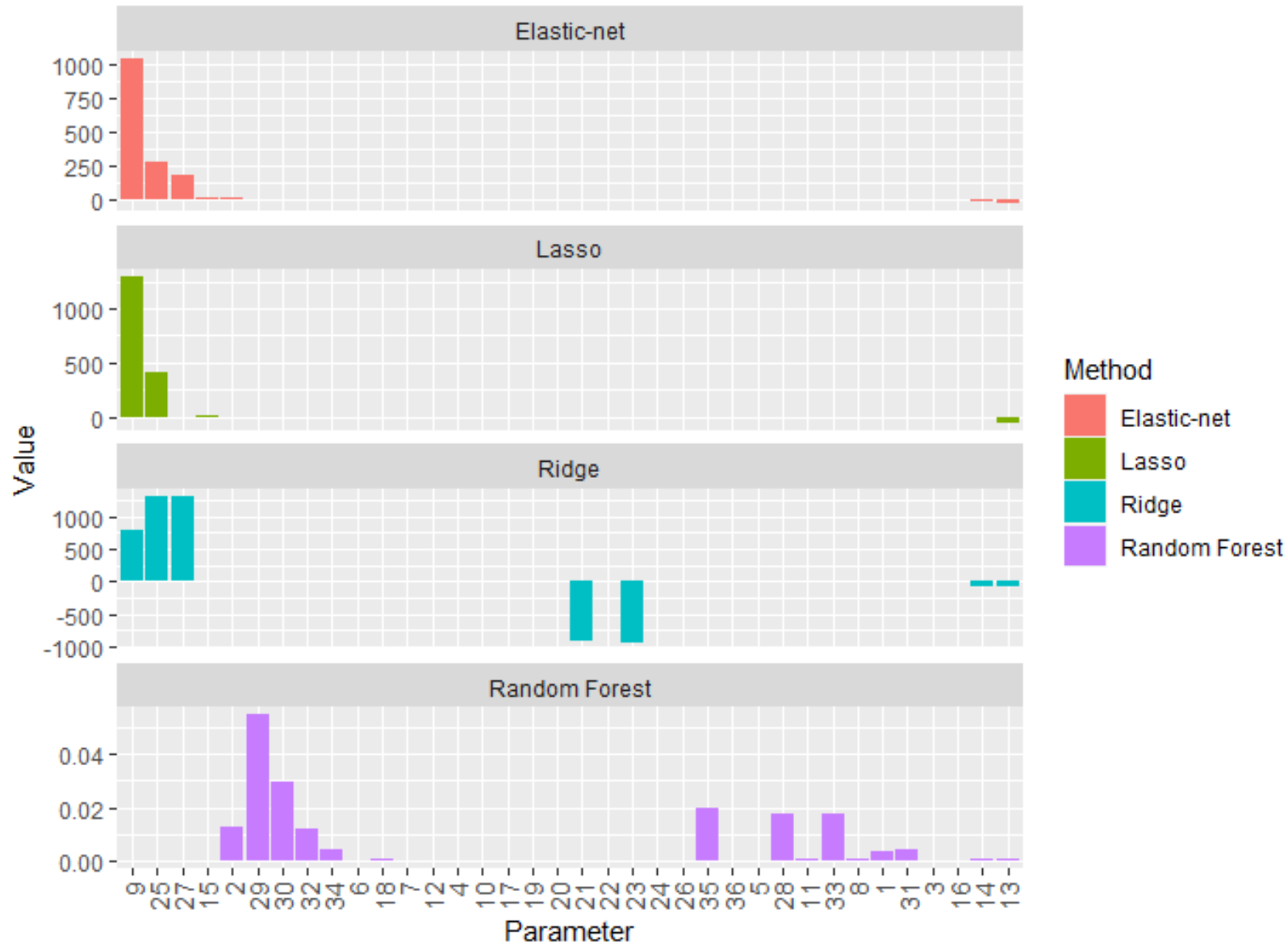- N = 1300
- P = 40

| | Time (sec) |
|---|---|
| Lasso | 0.870 |
| El-net | 1.050 |
| Ridge | 1.090 |

Train and Test Residuals

Importance of the parameters

# Test R2 and time

| | 90% Test R2 | Time (in sec) |
|---|---|---|
| Lasso | (0.85056 - 0.97669) | 1.36 |
| El-net | (0.8169 - 0.9473) | 1.30 |
| Ridge | (0.5393 - 0.7524) | 1.48 |
| RF | (0.99056 - 0.99206) | 3.94 |

# Concluding remarks

▶ In terms of time vs. efficiency Lasso gives the highest R2 and the fastest time which makes it the best model to predict water T;

▶ Ridge has the worst performance;

▶ Elastic net and lasso agree on most important features.