



Final Project

Communities and Crime dataset

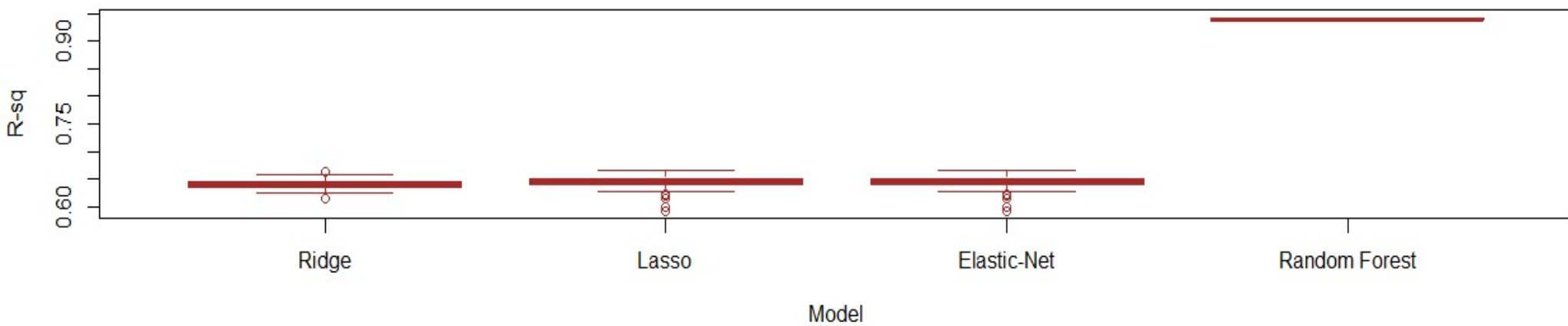
STA 9890 - Statistical Learning for Data Mining
Spring 2020 - May 20, 2020

Ariel Chajet

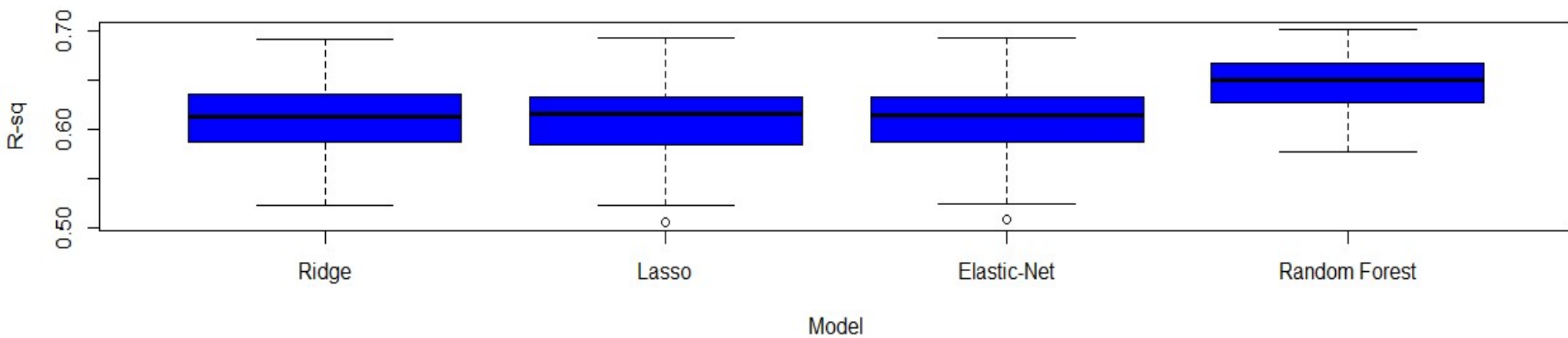
Data Description

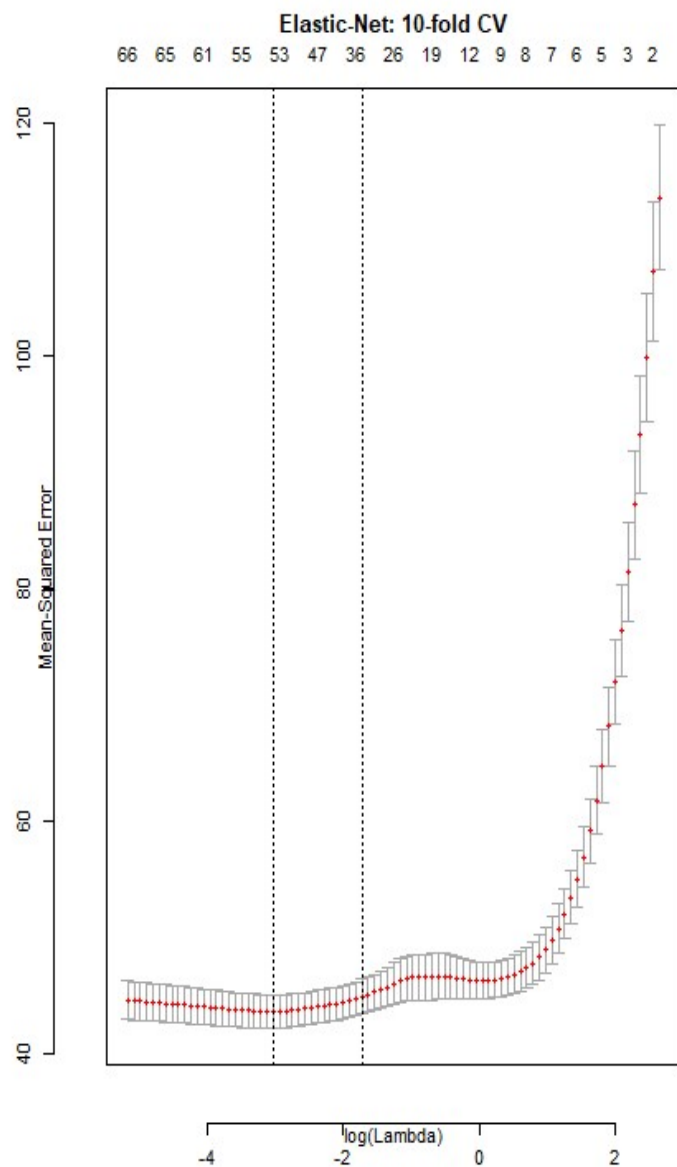
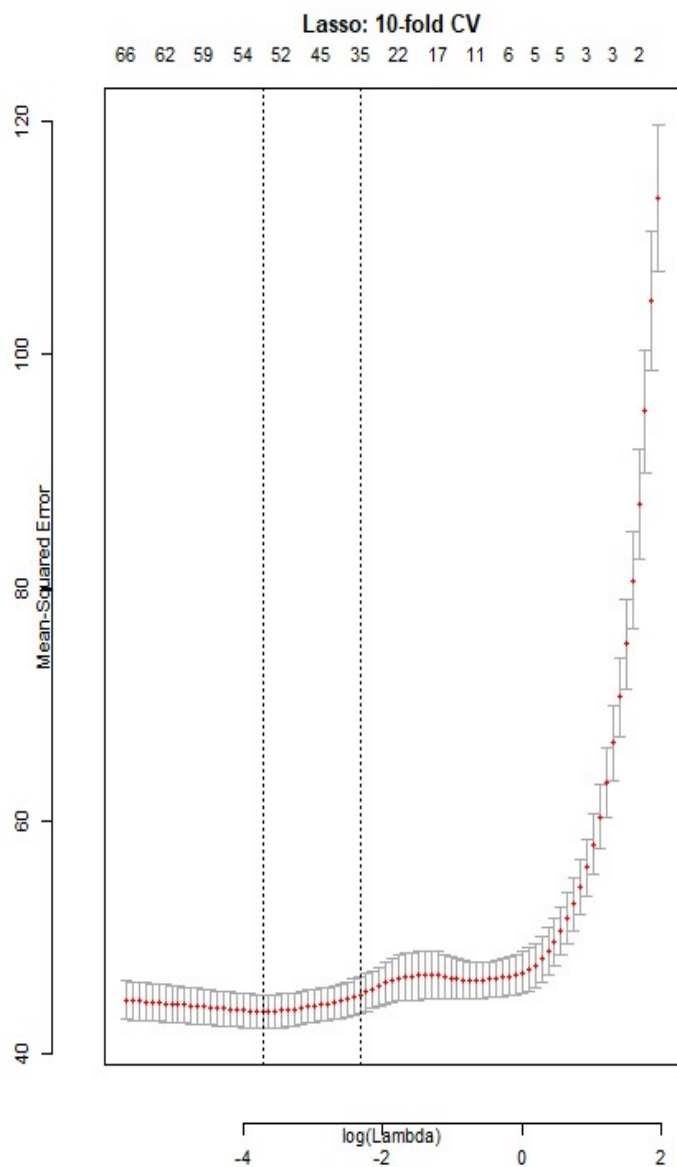
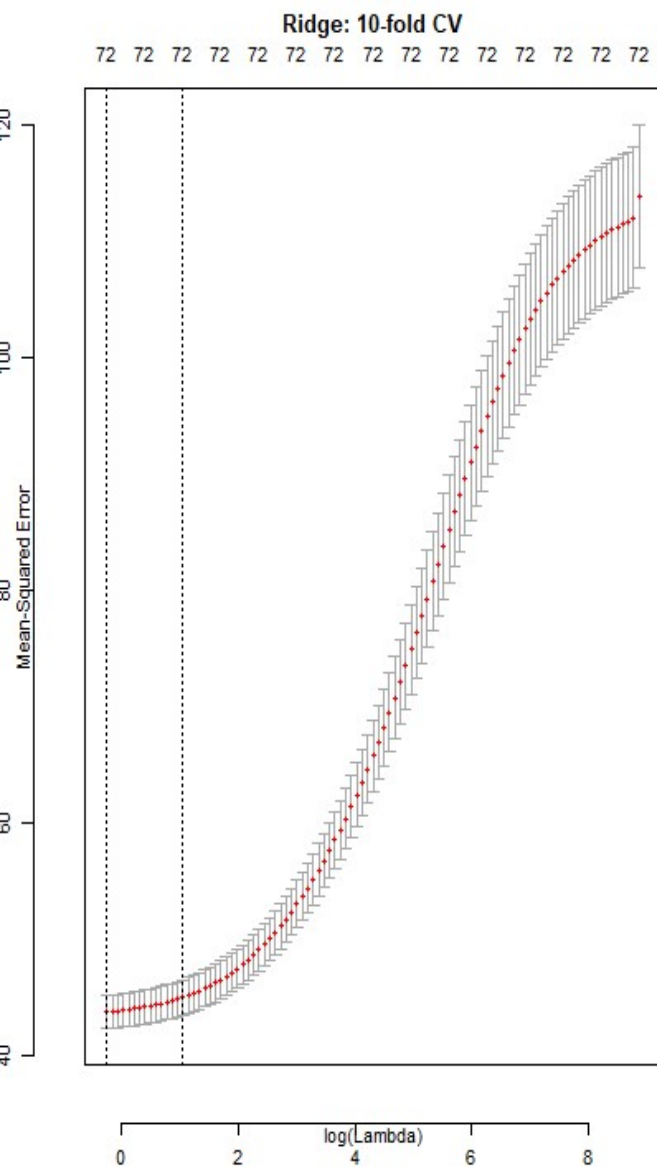
- ▶ The data compiles socio-economic data from the 1990 Census, law enforcement data from the 1990 Law Enforcement Management and Admin Stats survey, and crime data from the 1995 FBI UCR. (Source: <https://archive.ics.uci.edu/ml/datasets/Communities+and+Crime+Unnormalized>)
- ▶ Each record in the data represents an entire town/city and all related census and crime attributes pertaining to that city.
- ▶ The selected response variable represents the rate of violent crimes per 100K population, which is positively skewed. The square root of the rate of violent crimes was ultimately used to help center the distribution of the response.
- ▶ The raw data contains 2,215 observations with 147 features (5 non-predictive features, 125 predictive, and 18 potential response variables). All predictor and response variables are quantitative. Ultimately, the working data set was reduced to 72 predictors for modeling.
- ▶ The data was downloaded from the URL and saved locally as a CSV file. Values of “?” were converted to null values in excel before importing data into R. The column names were also added as headers to the CSV file.

Training R-sq For Each Model

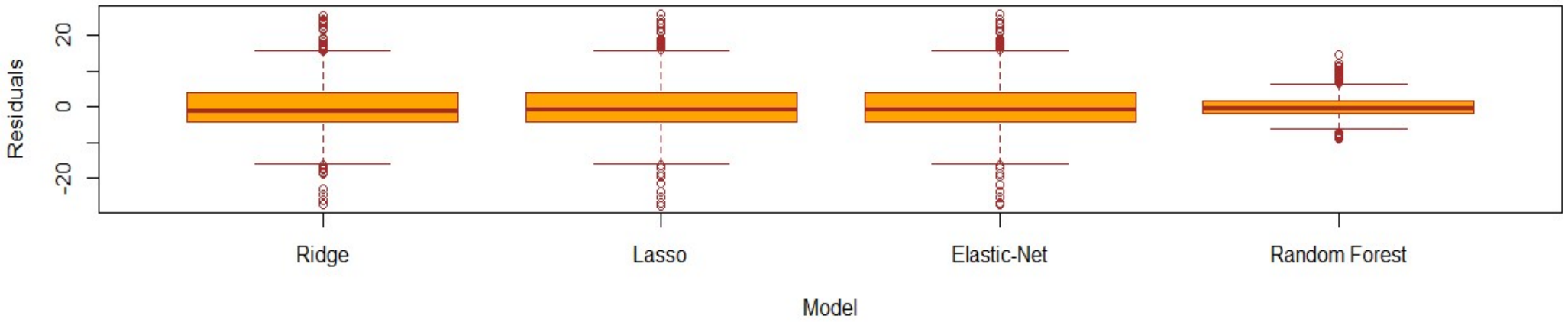


Test R-sq For Each Model

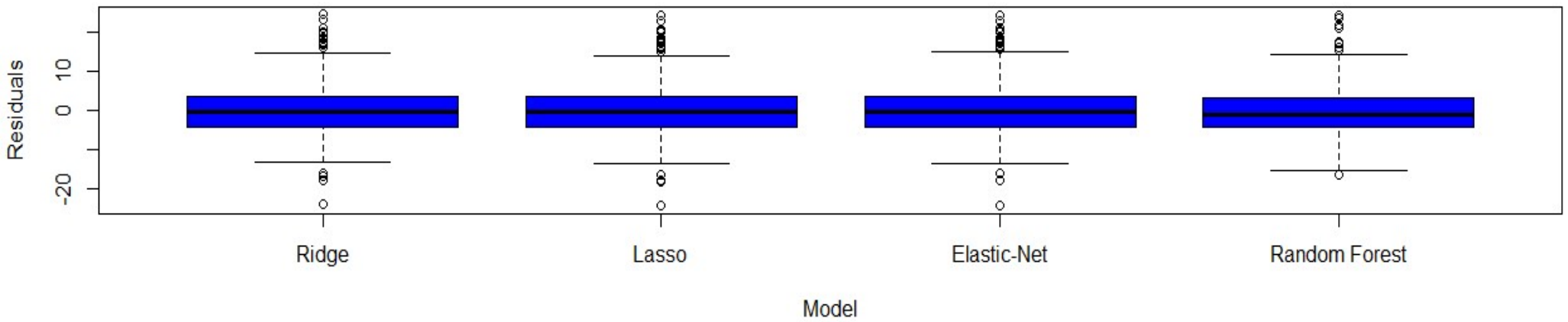




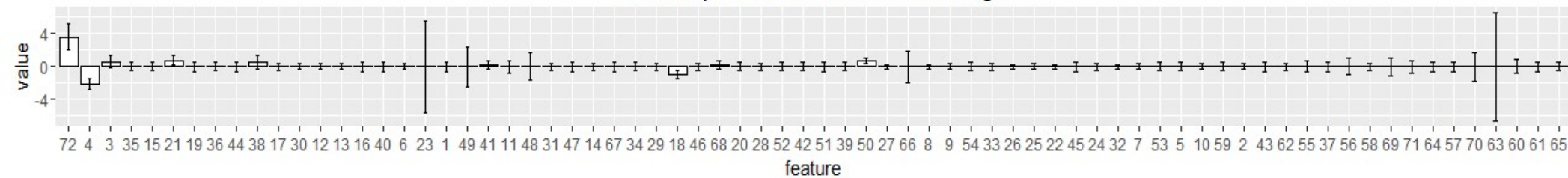
Training Residuals For Each Model



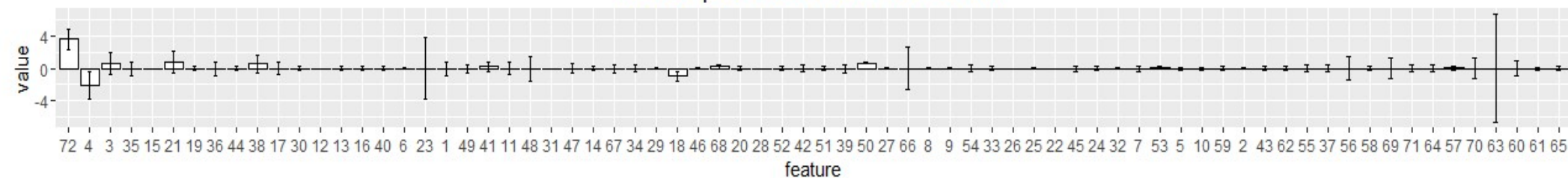
Test Residuals For Each Model



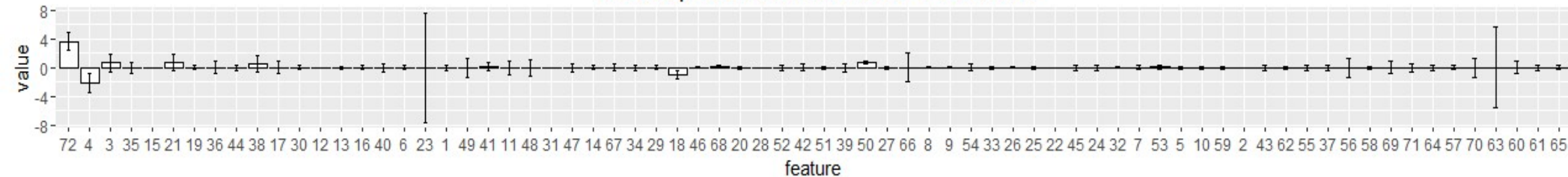
Bootstrap Coefficient Estimates - Ridge



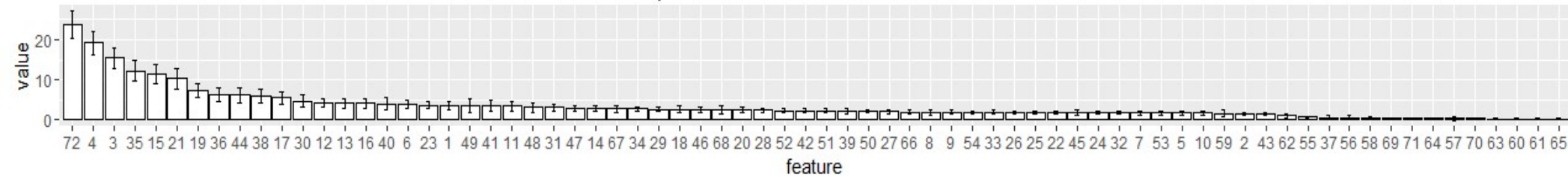
Bootstrap Coefficient Estimates - Lasso



Bootstrap Coefficient Estimates - Elastic Net



Bootstrap Coefficient Estimates - Random Forest



Predictor Index Legend

1	population	19	PctUnemployed	37	MedNumBR	55	PolicReqPerOffic
2	householdsize	20	PctEmploy	38	HousVacant	56	PolicPerPop
3	racepctblack	21	TotalPctDiv	39	PctHousOccup	57	RacialMatchCommPol
4	racePctWhite	22	PersPerFam	40	PctHousOwnOcc	58	PctPolicWhite
5	racePctAsian	23	NumImmig	41	PctVacantBoarded	59	PctPolicBlack
6	racePctHisp	24	PctRecentImmig	42	PctVacMore6Mos	60	PctPolicHisp
7	agePct12t21	25	PctReclImmig5	43	MedYrHousBuilt	61	PctPolicAsian
8	agePct12t29	26	PctReclImmig8	44	PctHousNoPhone	62	PctPolicMinor
9	agePct16t24	27	PctReclImmig10	45	PctWOFullPlumb	63	OfficAssgnDrugUnits
10	agePct65up	28	PctSpeakEnglOnly	46	MedRent	64	NumKindsDrugsSeiz
11	numbUrban	29	PctNotSpeakEnglWell	47	MedRentPctHousInc	65	PolicAveOTWorked
12	medIncome	30	PctLargHouseFam	48	NumInShelters	66	LandArea
13	medFamInc	31	PctLargHouseOccup	49	NumStreet	67	PopDens
14	perCapInc	32	PersPerOccupHous	50	PctForeignBorn	68	PctUsePubTrans
15	NumUnderPov	33	PersPerOwnOccHous	51	PctBornSameState	69	PolicCars
16	PctLess9thGrade	34	PersPerRentOccHous	52	PctSameHouse85	70	PolicOperBudg
17	PctNotHSGrad	35	PctPersDenseHous	53	PctSameCity85	71	PolicBudgPerPop
18	PctBSorMore	36	PctHousLess3BR	54	PctSameState85	72	nonViolPerPop

Conclusions

- ▶ The regularized regression methods experienced similar results across training and testing performance.
 - ▶ Lasso: reduced to 35 variables
 - ▶ Elastic-Net: reduced to 36 variables
- ▶ Random Forest showed signs of overfitting in training but only marginally better than the regularized regression methods in the model testing phase.
- ▶ Random Forest took approx. 15 times longer to train.
- ▶ All models were generally consistent in selecting their variables of importance

Model	Approx time to train
Ridge	3 min
Lasso	3 min
Elastic-Net	3 min
Random Forest	45 min

Important Variables for all models:		Other important variables in Reg. Reg. Methods:	
72	nonViolPerPop	21	TotalPctDiv
4	racePctWhite	38	HousVacant
3	racepctblack	23	NumImmig
35	PctPersDenseHous	18	PctBSorMore
15	NumUnderPov	50	PctForeignBorn