



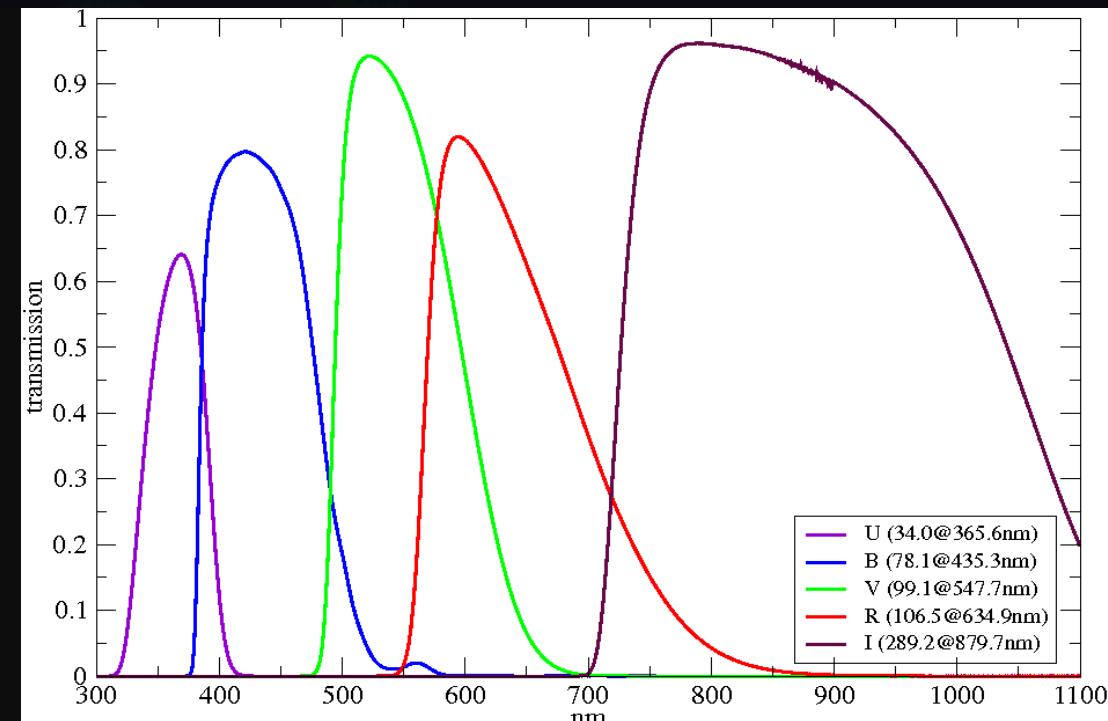
# *Galaxies' Redshift Estimates*

---

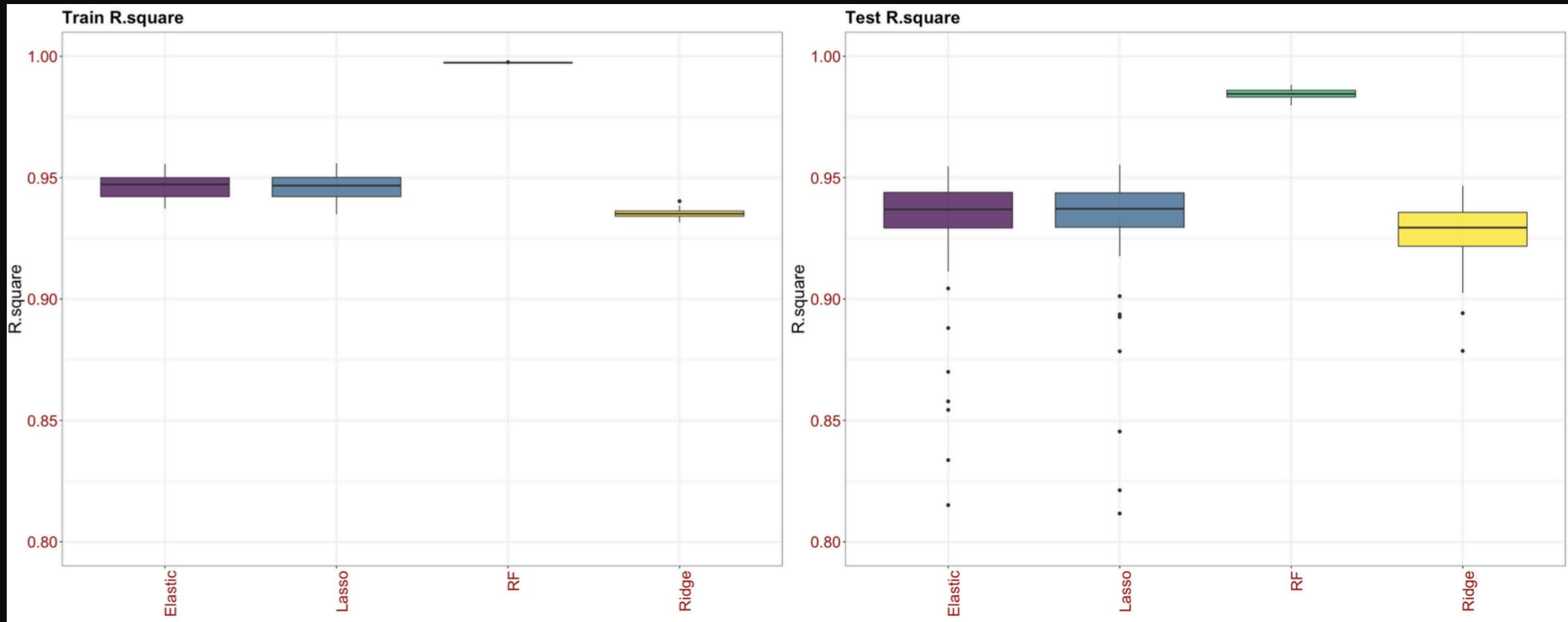
My Vien

## INTRODUCTION

- The target variable is photometric estimates of redshift, a measure of distance from us to an object.
- The predictors variables are mainly the measurements and errors of brightness of color bands of galaxies.
- $n = 3438$ , number of galaxies samples (after removing missing values)
- $p = 60$  (number of predictors)
- We will be able to predict the velocity of a galaxy based on the brightness of color bands of a galaxy



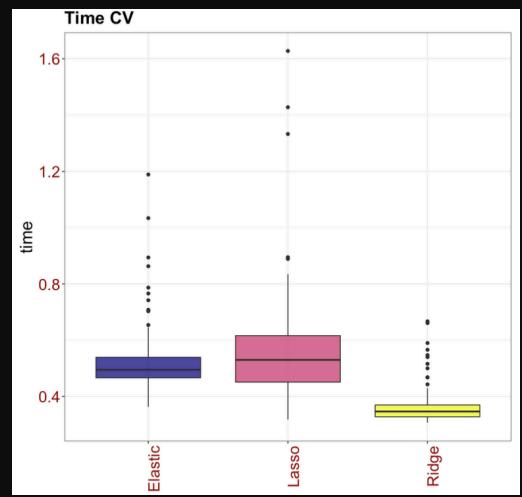
# R2



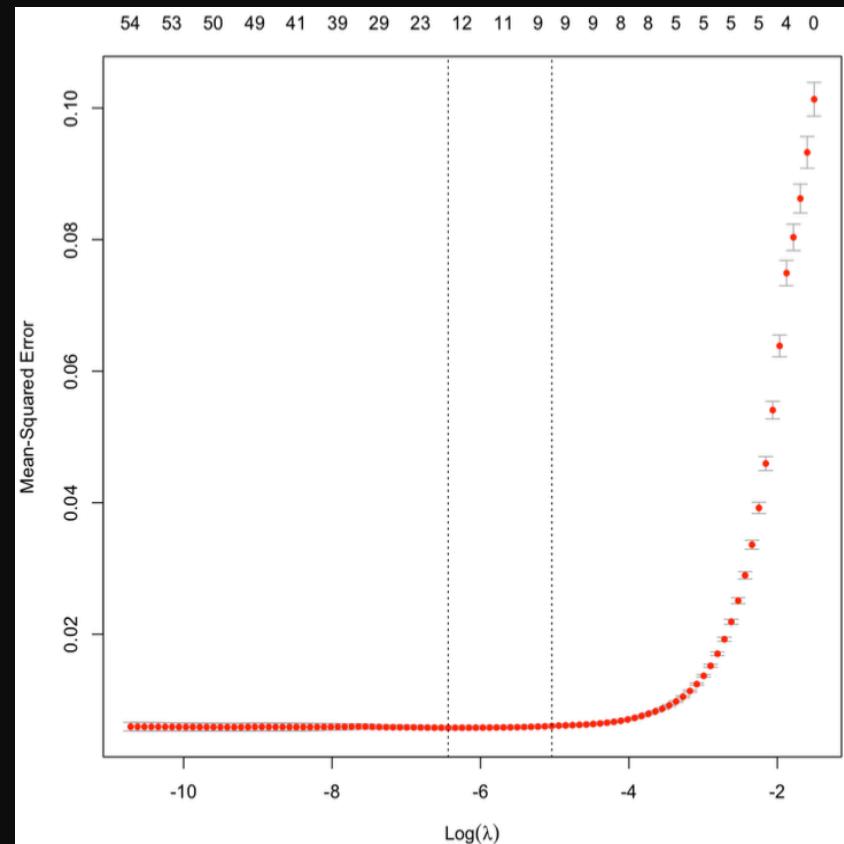
- Random Forest has the best train R-square and test R-square, while Ridge is the lowest on both train and test set.
- Test R-square values are lower and more spread on all models

## 10-FOLD CV CURVE

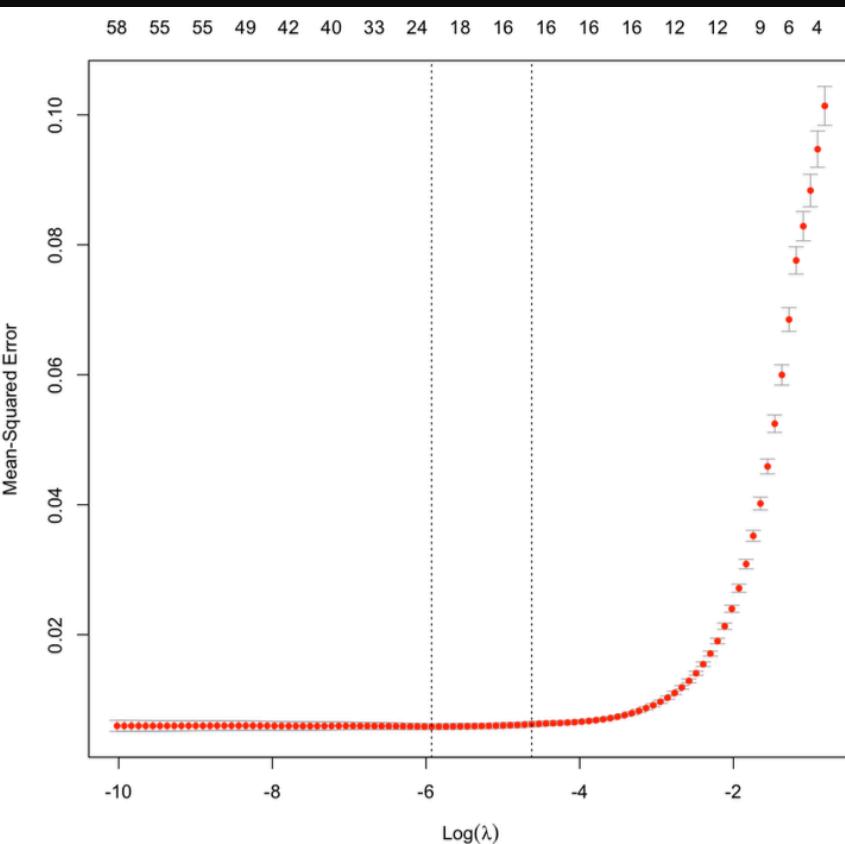
- Optimized Lasso has around 12 non-zero coefficients
- Optimized ElasticNet has around 22-23 non-zero coefficients
- As expected, Optimized Ridge has maintained all predictors coefficients



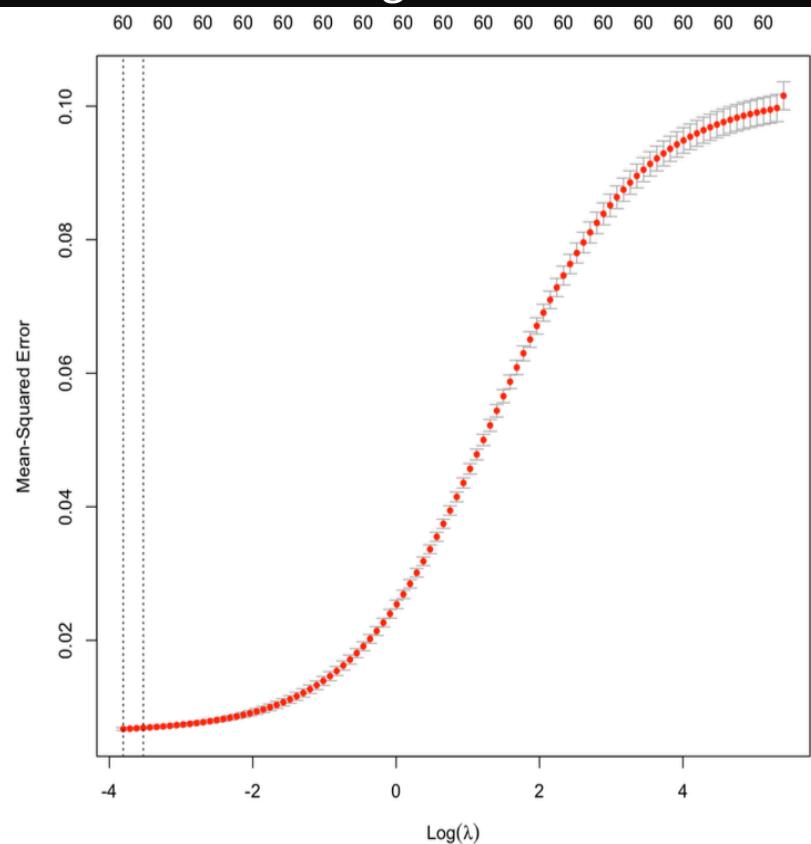
LASSO



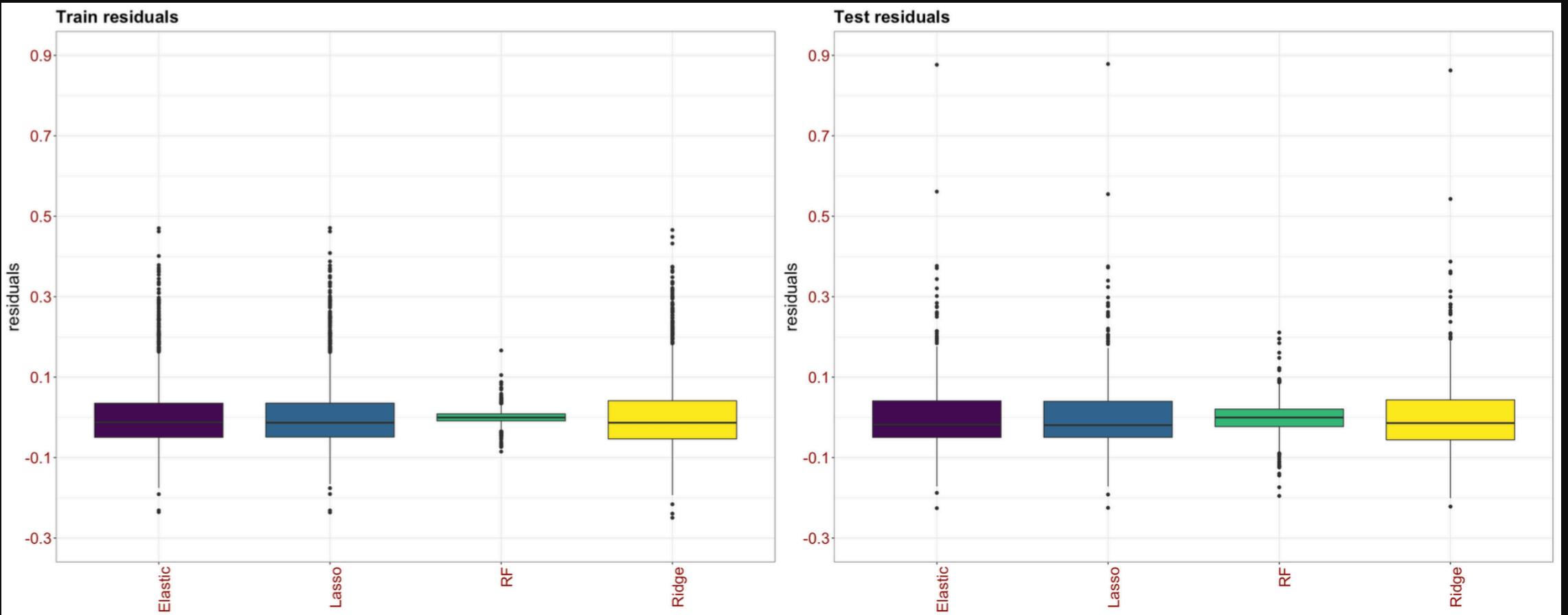
Elastic Net



Ridge



# RESIDUALS



- Train size is 2750 and test size is 688.
- Random Forest has smallest spread while all constrained models have larger spread.
- Extreme residual outliers are detected on test set of all 3 constrained models.

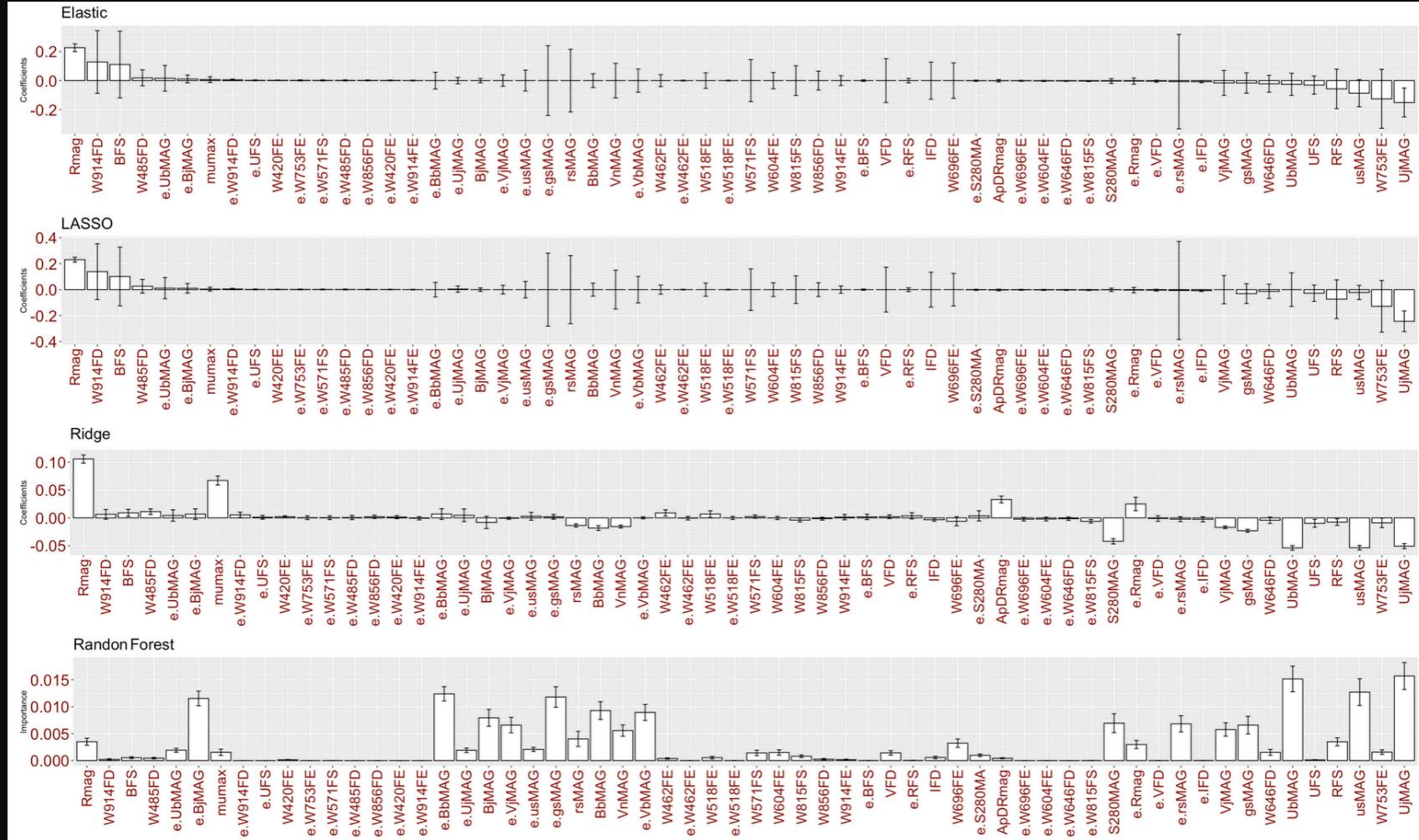
## Performance

- There is a trade off between time and model performance comparing between 3 constrained linear models and RF.
- Lasso takes longest time to run but does not perform much better than Elastic Net.
- Ridge is fastest, but its score is less comparable.

Models	90% test R2	Time (s)
Lasso	(0.92183, 0.947541)	0.623
ElasticNet	(0.91962, 0.948475)	0.595
Ridge	(0.91070, 0.93855)	<b>0.409</b>
RF	<b>(0.98144, 0.98680)</b>	39.352

# Coefficients (100 Bootstrap Samples)

- All linear models have similar feature importance: UjMAG (Johnson U luminosity) and Rmag (red band magnitude) are the two largest coefficients in 3 linear models.
- With Random Forest, it is UjMAG (Johnson U luminosity) and UbMAG (Bessell U luminosity)
- Lasso and Elastic Net seem to have similar coefficients magnitudes.
- All linear models have more spread in features importance than Random Forest.
- Some features in Elastic Net and Lasso have very wide standard deviation.
- Based on Elastic Net and Lasso model, when holding other predictors constant, redshift is predicted to increase by 0.25 if red band magnitude goes up by 1; it is the opposite with ultra-violet luminosity.



## Conclusion:

---

- There is not a big difference between train and test R-squared, train and test residuals which proves that all our predictive models have good performance.
- Even though RF takes longest time to run, it is still the best model considering its 90% test R-squared.
- The features importance between RF and linear models are different, there must be some non-linear relationship between the predictors and response variable so that RF can achieve significantly better score than the linear models.
- There is a positive linear relationship between redshift and red band magnitude; negative linear relationship between redshift and ultra-violet luminosity.