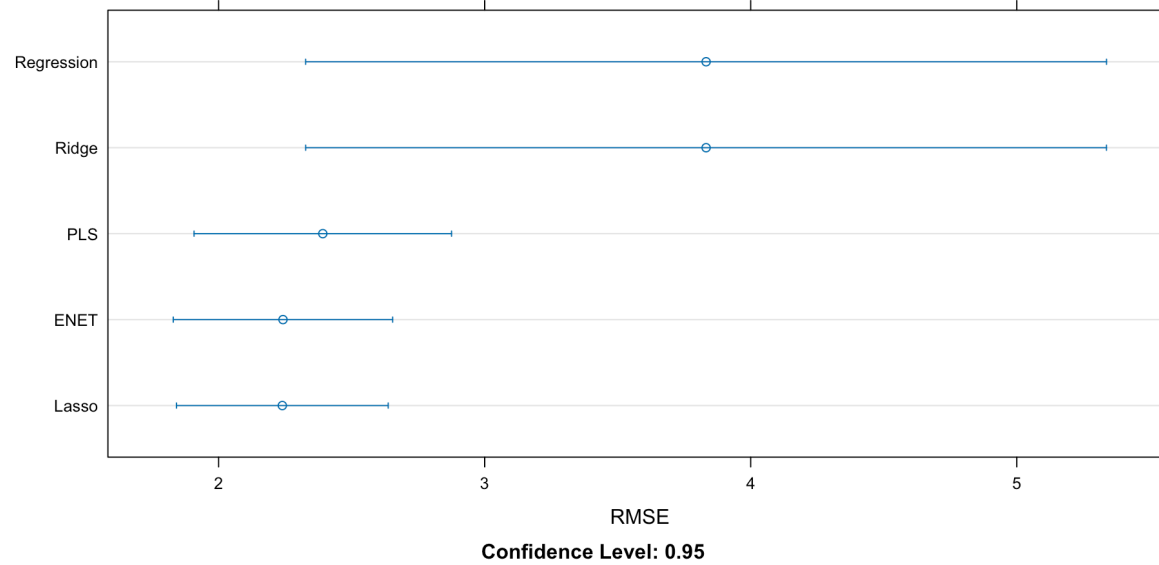


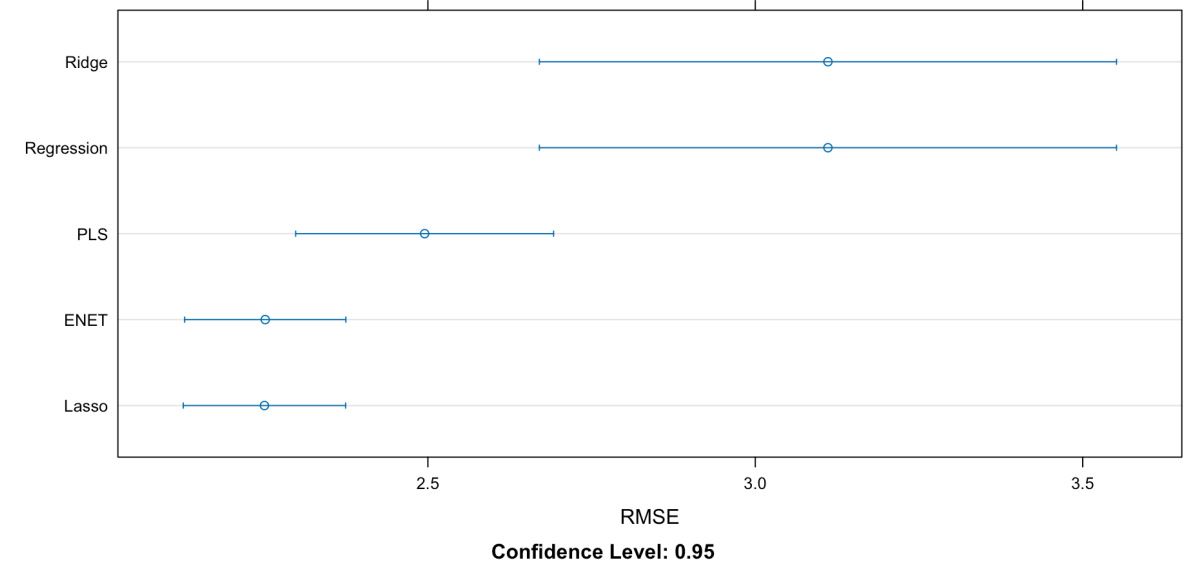
CV

```
1 dotplot(train_metrics_cv, metric = 'RMSE')
```



Repeated CV

```
1 dotplot(train_metrics, metric = 'RMSE')
```



ADS 503 - Applied Predictive Modeling (M4)

Summer 2024 - Week 4

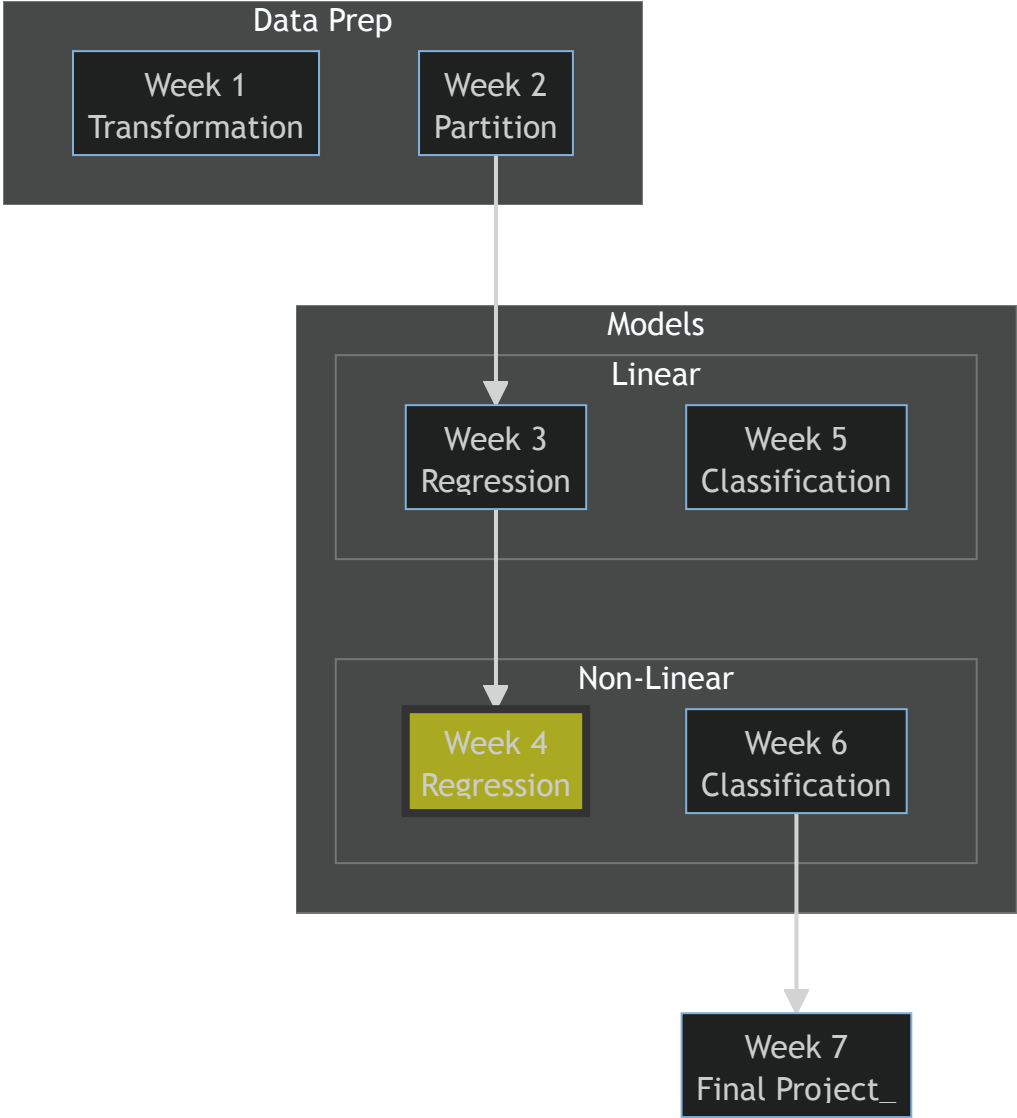
Dave Hurst

Start Recording!

Agenda

- Course Map
- Assignment 3 Review
 - Hyperparameter Tuning
- Assignment 4 Tips
- QA

Course Map



Assignment 3 Review

3.1.c with `method = "cv"`

```
1 seed <- 503
2 data(tecator)
3
4 # Extract the predictors (absorbance) and response (fat content)
5 absorbance <- as.data.frame(absorp)
6 fat_content <- endpoints[,2] # fat is the second column
7
8 # Split the data into training and validation sets
9 set.seed(seed) # for reproducibility
10 train_index <- createDataPartition(fat_content, p = 0.8, list = FALSE)
11 train_data <- absorbance[train_index, ]
12 train_fat <- fat_content[train_index]
13 test_data <- absorbance[-train_index, ]
14 test_fat <- fat_content[-train_index]
15
16 #CV
17 train_cv <- trainControl(method = "cv")
18
19 tic('31c')
20 elapsed <- numeric()
21 # Linear Regression
22 set.seed(seed); tic('lm') # for reproducibility/timing
23 lm_model_cv <- train(train_data, train_fat,
24                      method = "lm",
25                      preProcess = c("center", "scale"),
26                      trControl = train_cv
```

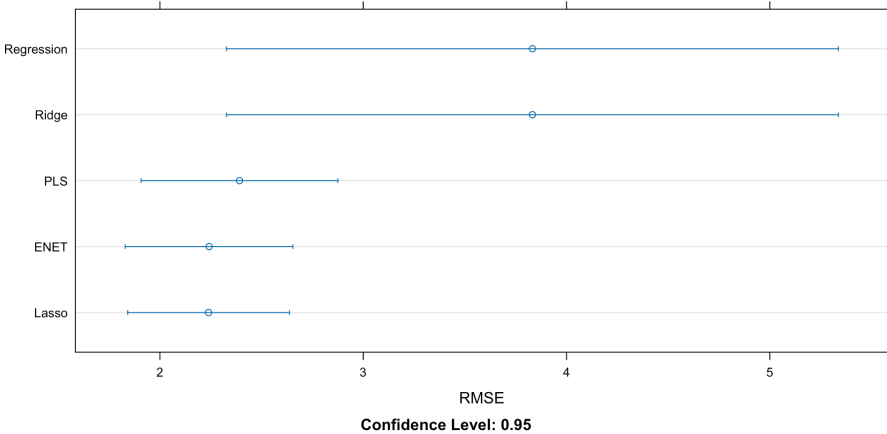
Models/CV	Times
LinReg	2.110
PLS	2.260
Ridge	35.777
Lasso	6.734
ENET	18.011

3.1.c with `method = "cv"`

```
1 train_metrics_cv <- resamples(list(
2   Regression = lm_model_cv,
3   PLS = pls_model_cv,
4   Ridge = ridge_model_cv,
5   Lasso = lasso_model_cv,
6   ENET = enet_model_cv))
7 summary(train_metrics_cv)$statistics$RMSE
```

	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	NA's
Regression	2.154094	2.416813	3.067239	3.832277	4.194495	8.925659	0
PLS	1.077971	2.070936	2.547476	2.391443	2.874441	3.182084	0
Ridge	2.154062	2.416863	3.067253	3.832244	4.194643	8.924921	0
Lasso	1.217658	1.971209	2.236171	2.239400	2.560048	3.040155	0
ENET	1.116217	1.971931	2.241003	2.241807	2.638687	3.000962	0

```
1 dotplot(train_metrics_cv, metric = 'RMSE')
```



```
1 diff(train_metrics_cv, metric = "RMSE") |> summary()
```

Call:
summary.diff.resamples(object = diff(train_metrics_cv, metric = "RMSE"))

p-value adjustment: bonferroni
Upper diagonal: estimates of the difference
Lower diagonal: p-value for H0: difference = 0

RMSE	Regression	PLS	Ridge	Lasso	ENET
Regression		1.441e+00	3.316e-05	1.593e+00	1.590e+00
PLS	0.6874		-1.441e+00	1.520e-01	1.496e-01
Ridge	1.0000	0.6872		1.593e+00	1.590e+00
Lasso	0.4395	1.0000	0.4394		-2.407e-03
ENET	0.4471	1.0000	0.4470	1.0000	

3.1.c with `method = "repeatedcv"`

```
1 #Repeated CV
2 train_repeated_cv <- trainControl(method = "repeatedcv", repeats = 5)
3
4 tic('31c')
5 elapsed <- numeric()
6 # Linear Regression
7 set.seed(seed); tic('lm') # for reproducibility/timing
8 lm_model <- train(train_data, train_fat,
9                   method = "lm",
10                  preProcess = c("center", "scale"),
11                  trControl = train_repeated_cv
12                  )
13 tics_lm <- toc(quiet = TRUE); elapsed <- c(elapsed, tics_lm$tic - tics_lm$toc)
14
15
16 # Partial Least Squares
17 set.seed(seed); tic('pls') # for reproducibility
18 pls_model <- train(train_data, train_fat,
19                   method = "pls",
20                   tuneLength = 40,
21                   preProcess = c("center", "scale"),
22                   trControl = train_repeated_cv
23                   )
24 tics_pls <- toc(quiet = TRUE); elapsed <- c(elapsed, tics_pls$tic - tics_pls$toc)
25
26 # Ridge Regression
```

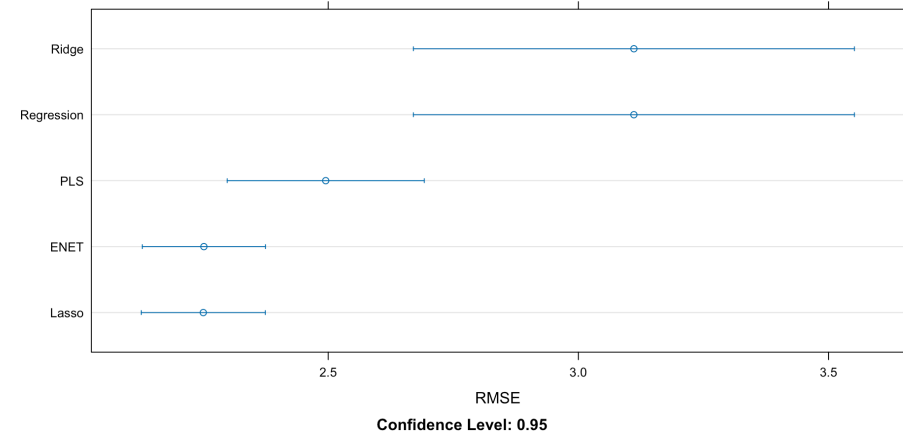
Models/RCV	Times
LinReg	8.212
PLS	8.598
Ridge	182.964
Lasso	30.580
ENET	77.334

3.1.c with `method = "repeatedcv"`

```
1 train_metrics <- resamples(list(
2   Regression = lm_model,
3   PLS = pls_model,
4   Ridge = ridge_model,
5   Lasso = lasso_model,
6   ENET = enet_model))
7 summary(train_metrics)$statistics$RMSE
```

	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	NA's
Regression	1.274747	2.063995	2.590646	3.110938	3.568557	8.925659	0
PLS	1.077971	2.053263	2.496077	2.495050	2.903803	4.966348	0
Ridge	1.274757	2.063859	2.590538	3.110968	3.568723	8.924921	0
Lasso	1.217658	1.960208	2.240318	2.250353	2.535707	3.229732	0
ENET	1.240802	1.974018	2.232149	2.251569	2.522635	3.257020	0

```
1 dotplot(train_metrics, metric = 'RMSE')
```



```
1 diff(train_metrics, metric = "RMSE") |> summary()
```

Call:
summary.diff.resamples(object = diff(train_metrics, metric = "RMSE"))

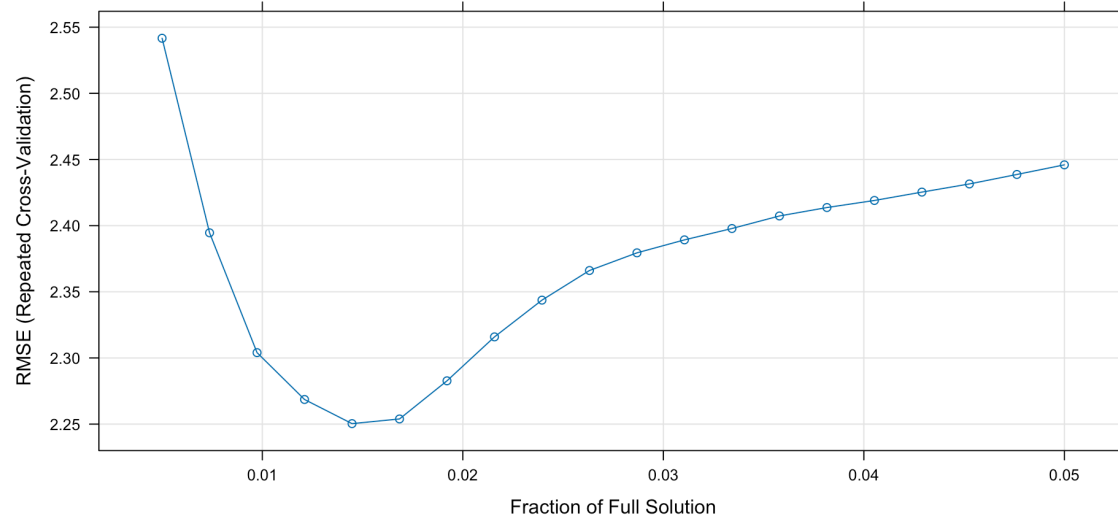
p-value adjustment: bonferroni
Upper diagonal: estimates of the difference
Lower diagonal: p-value for H0: difference = 0

RMSE	Regression	PLS	Ridge	Lasso	ENET
Regression		6.159e-01	-3.037e-05	8.606e-01	8.594e-01
PLS	0.092884		-6.159e-01	2.447e-01	2.435e-01
Ridge	1.000000	0.092848		8.606e-01	8.594e-01
Lasso	0.003122	0.023851	0.003120		-1.216e-03
ENET	0.003152	0.030100	0.003151	1.000000	

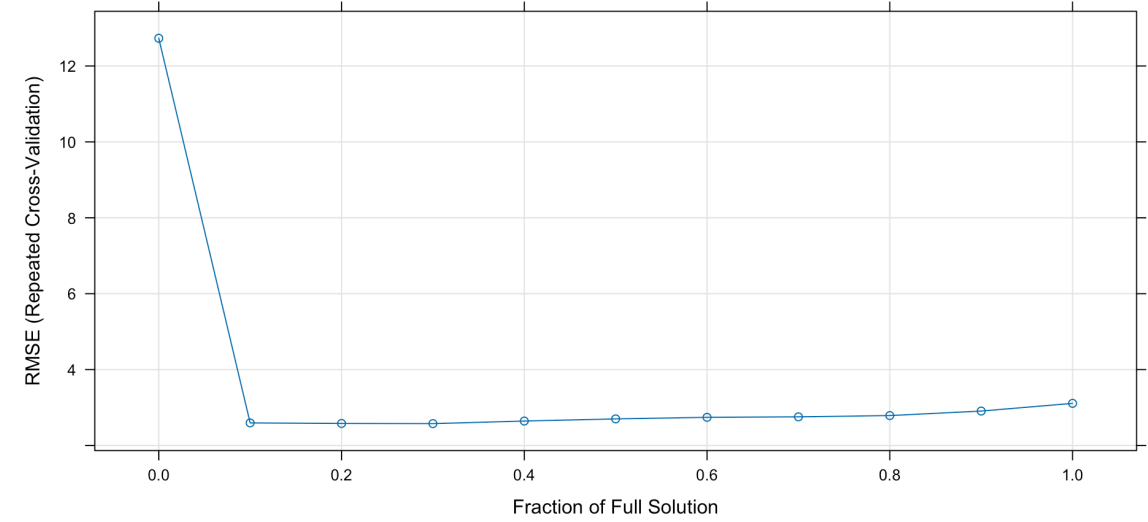
Hyperparameter Optimizaton

```
1 lambda_grid <- data.frame(.fraction = seq(0.005, 0.05, length = 20))
2 lambda_grid2 <- data.frame(.fraction = seq(0, 1, length = 11))
3
4 # Lasso Regression
5 set.seed(seed)
6 lasso_model2 <- train(train_data, train_fat,
7                       method = "lasso",
8                       tuneGrid = lambda_grid2,
9                       preProcess = c("center", "scale"),
10                      trControl = train_repeated_cv
11 )
```

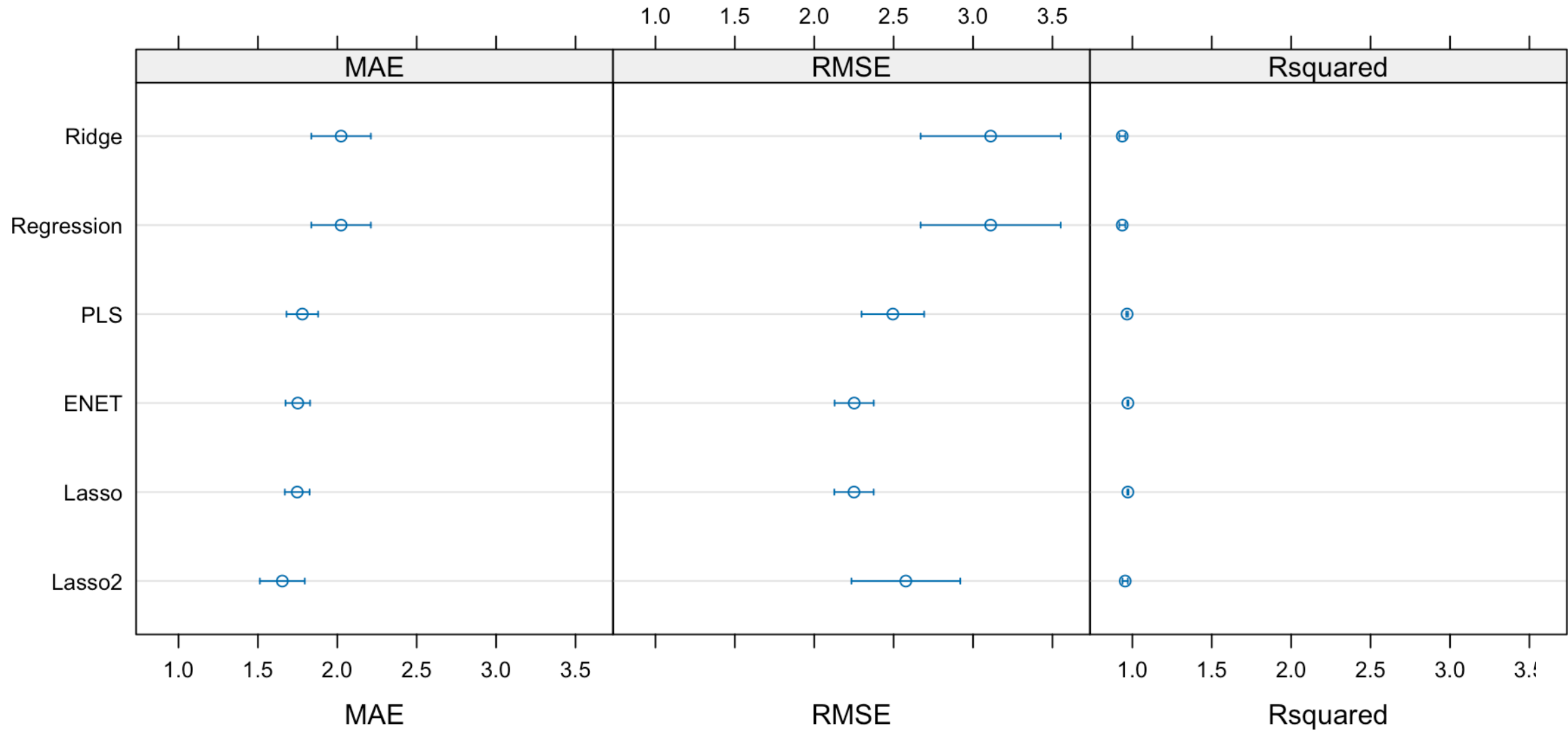
```
1 plot(lasso_model1)
```



```
1 plot(lasso_model2)
```



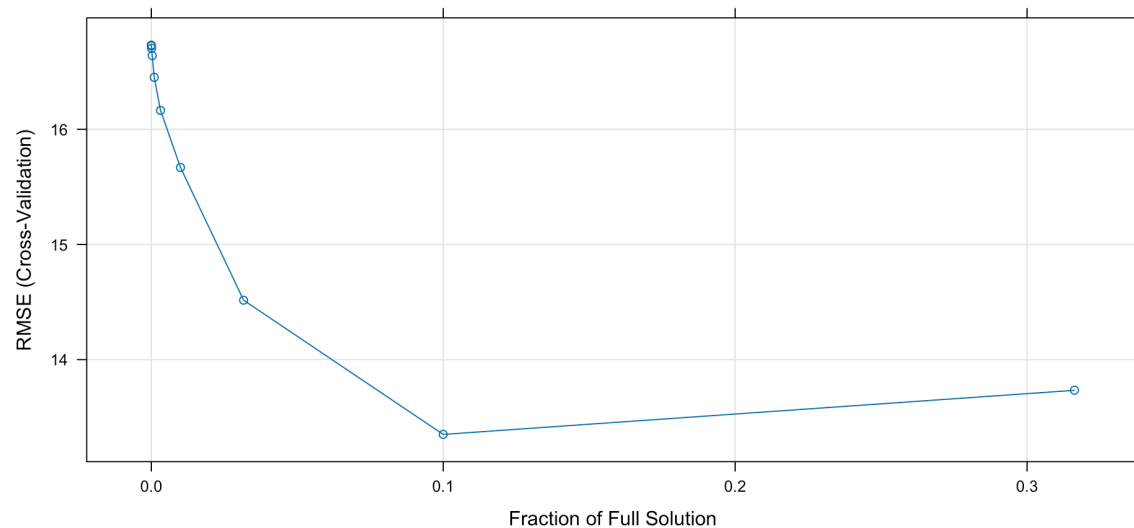
```
1 dotplot(train_metrics2)
```



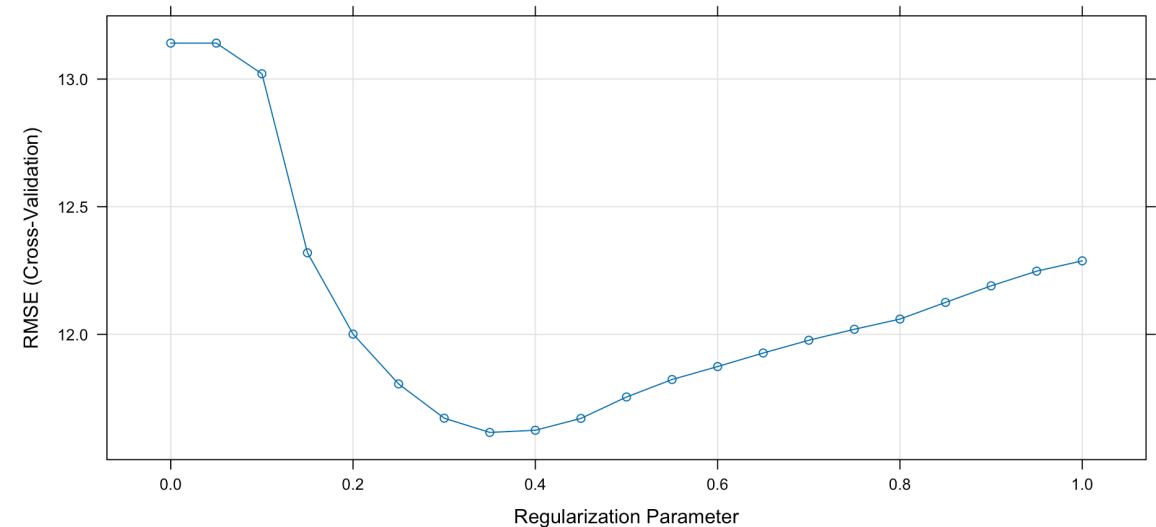
Confidence Level: 0.95

3.2.e - Lasso Issue

```
1 # Lasso Regression
2 set.seed(seed)
3 lambda_grid <- data.frame(.fraction = 10^seq(-5, -.5,
4 fp_lasso_model <- train(train_fingerprints_trans, tra
5                         method = "lasso",
6                         tuneGrid = lambda_grid,
7                         trControl = trainControl(method =
8 )
9 plot(fp_lasso_model)
```



```
1 # alternate Lasso
2 lasso_grid <- expand.grid(alpha = 1, lambda = seq(0,
3 fp_lasso_model2 <- train(train_fingerprints_trans, tra
4                         method = "glmnet",
5                         tuneGrid = lasso_grid,
6                         trControl = trainControl(method =
7 )
8 plot(fp_lasso_model2)
```



Assignment 4 Tips

- Use the A4-M1.qmd template (whether you use posit or not)
- Plot your model (hyperparameters) (page limit extended to 21)
- posit.cloud solutions should solve quickly (< 90 s for entire notebook)

Q&A

Is there a difference between using `preProcess` within `train()` versus manually transforming the data?

```
1 # Apply pre-processing with `caret::train()` (repeated
2 # Split the data into training and validation sets
3 set.seed(seed) # for reproducibility
4 train_index <- createDataPartition(fat_content, p = 0.
5 train_data <- absorbance[train_index, ]
6 train_fat <- fat_content[train_index]
7 test_data <- absorbance[-train_index, ]
8 test_fat <- fat_content[-train_index]
9
10 #CV
11 train_cv <- trainControl(method = "cv")
12
13 # Partial Least Squares
14 set.seed(seed);
15 pls_model_auto <- train(train_data, train_fat,
16                         method = "pls",
17                         tuneLength = 40,
18                         preProcess = c("center", "scale"),
19                         trControl = train_cv
20 )
```

```
1 # Apply pre-processing manually
2 absorb_prep <- preProcess(absorbance, method = c("cent
3 absorb_trans <- predict(absorb_prep, newdata = absorba
4 train_data_trans <- absorb_trans[train_index, ]
5 test_data_trans <- absorb_trans[-train_index, ]
6
7 # Partial Least Squares
8 set.seed(seed);
9 pls_model_manual <- train(train_data_trans, train_fat,
10                           method = "pls",
11                           tuneLength = 40,
12                           preProcess = c("center", "sc
13                           trControl = train_cv
14 )
```

```
1 # Compare the outputs.
2 predictions <- tibble(
3   auto = predict(pls_model_auto, newdata = test_data
4   manual = predict(pls_model_manual, newdata = test_
5 ) |>
6   mutate(abs_error = abs(manual - auto))
7 sum(predictions$abs_error)
```

[1] 6.375721e-08

No difference. In summary, `caret` knows to apply the same transformations automatically in the `predict()` step, ensuring that the data is processed consistently throughout the model lifecycle.