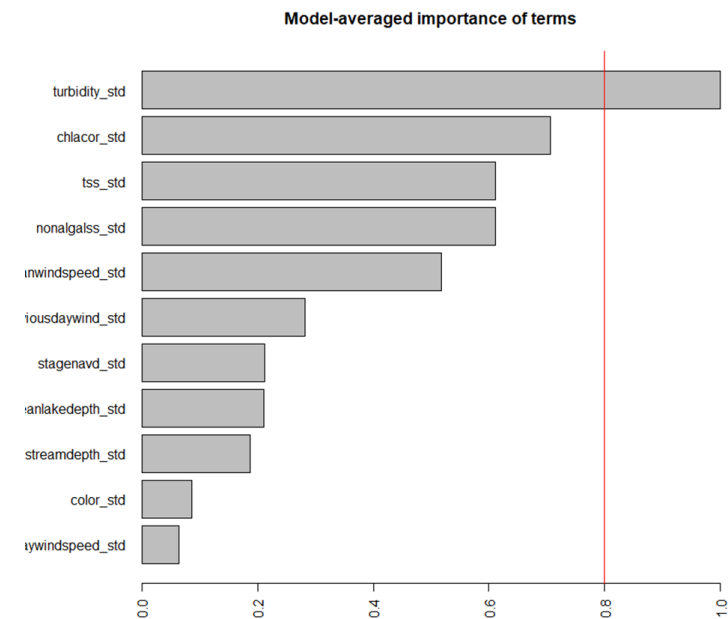
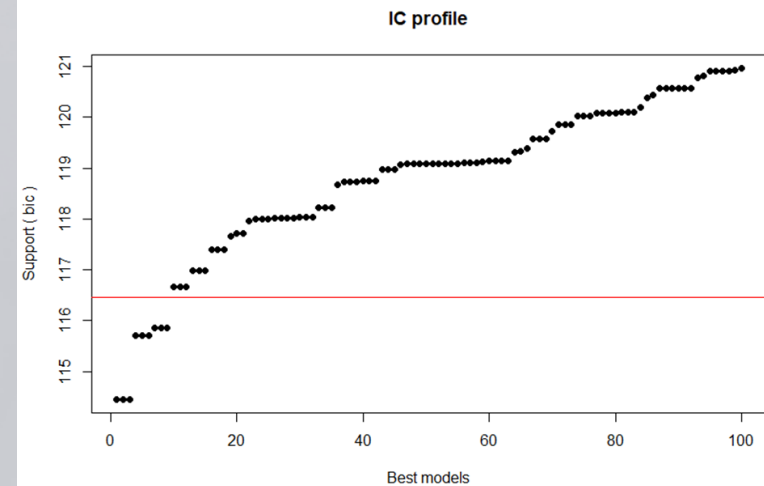


# Advanced R: Statistical Machine Learning

Dan Schmutz, MS  
Chief Environmental Scientist

Zoom Workshop for SJRWMD  
September 24, 2020



# Regularized Regression: Ridge, Lasso, and Elastic Net

# Regularization introduces some bias to reduce variance

- In statistical machine learning we frequently encounter datasets with multicollinearity: intercorrelated variables.
- Multicollinearity makes regression coefficients unstable: high variance between datasets.
- Normal regression seeks to minimize squared residuals
- Regularized regression adds either an L1 or L2 penalty to the regression coefficients (or a mixture of both).

# Ridge (L2 regularization)

- As lambda increases, the coefficients are forced to be smaller, introducing bias, but decreasing variance

$$\text{minimize} \left( SSE + \lambda \sum_{j=1}^p \beta_j^2 \right)$$

# Lasso (L1 Regularization)

- As lambda increases, the coefficients are forced to go to zero, effectively making the lasso a form of variable selection

$$\text{minimize } \left( SSE + \lambda \sum_{j=1}^p |\beta_j| \right)$$

# Elastic Net Regularization

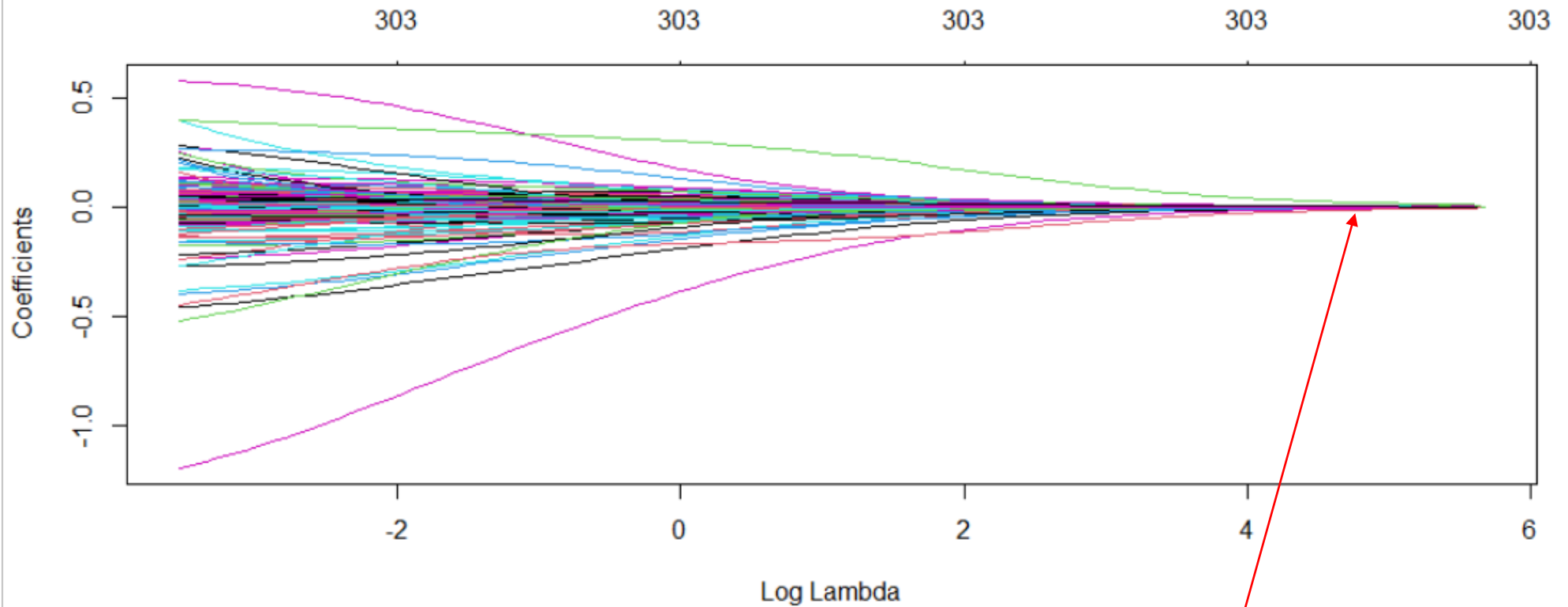
- Two hyper parameters, lambda you know, and alpha which controls the mixing ratio of L2 and L1

$$\text{minimize} \left( SSE + \lambda_1 \sum_{j=1}^p \beta_j^2 + \lambda_2 \sum_{j=1}^p |\beta_j| \right)$$

# Ridge on ames data

```
15 # Create training feature matrices
16 # we use model.matrix(...)[, -1] to discard the intercept
17 X <- model.matrix(Sale_Price ~ ., train_3)[, -1]
18 Xnotused<-model.matrix(Sale_Price ~.,train_3)
19
20 # transform y with log transformation
21 Y <- log(train_3$Sale_Price)
22
23 # Apply ridge regression to ames data
24 ridge <- glmnet(
25   x = X,
26   y = Y,
27   alpha = 0
28 )
29
30 plot(ridge, xvar = "lambda")
31
```

# Ridge on ames data



Higher lambda values crush the coefficients down to near but not reaching zero.



# Ridge on ames data

- glmnet chose range of 100 lambdas to investigate
- Compare coefficients at the largest and smallest lambdas

```
> ridge$lambda[1:100] %>% head()
[1] 291.4724 265.5788 241.9855 220.4882 200.9006 183.0532
> # small lambda results in large coefficients
> coef(ridge)[c("Latitude", "Overall_QualVery_Excellent"), 100]
      Latitude Overall_QualVery_Excellent
      0.3950246              0.1326520
> # large lambda results in small coefficients
> coef(ridge)[c("Latitude", "Overall_QualVery_Excellent"), 1]
      Latitude Overall_QualVery_Excellent
      6.203488e-36              9.744773e-37
```

# What is the optimal lamda for inference on out-of-sample data?

- Glmnet offers built-in crossvalidation (10-fold)

```
54 # Apply CV ridge regression to ames data
55 ridge <- cv.glmnet(
56   x = X,
57   y = Y,
58   alpha = 0
59 )
```

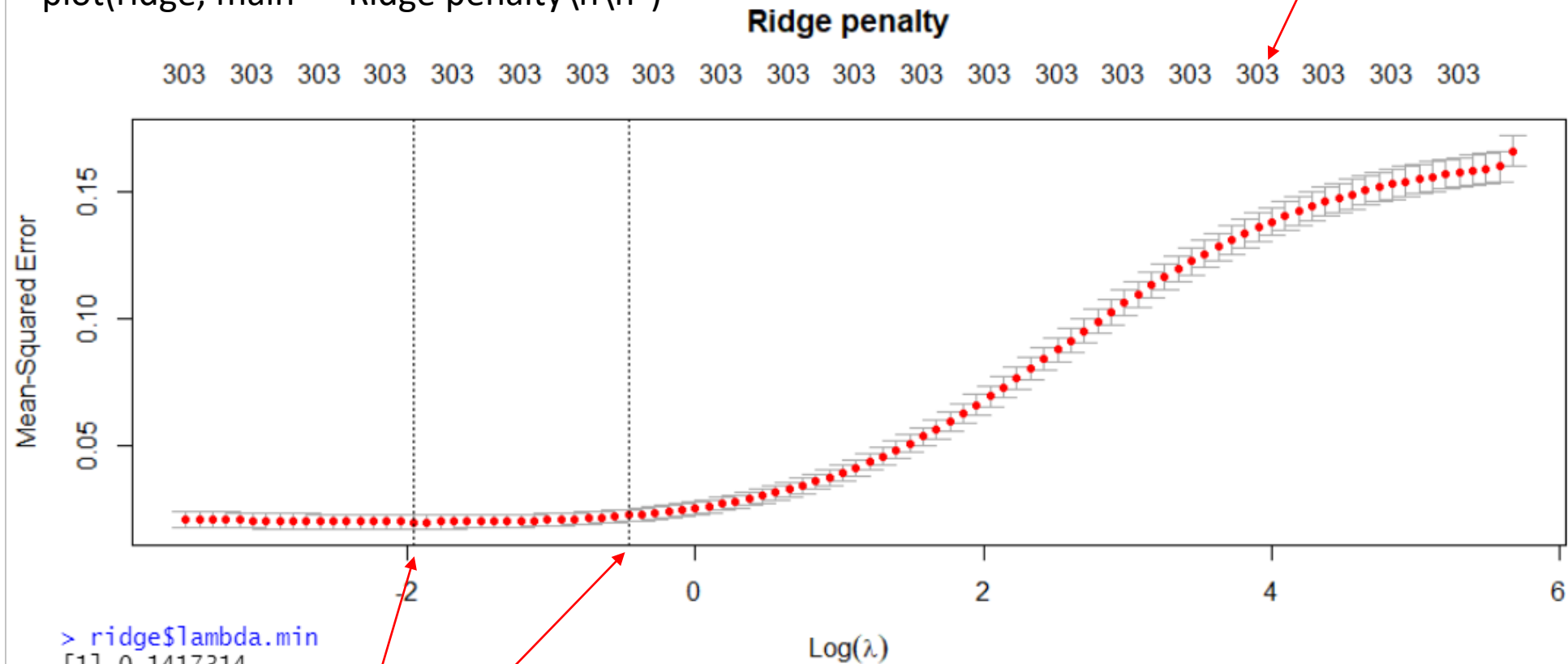
- Str(ridge)

```
## Ridge Regression Results
##- attr(*, "class")= chr [1:2] "elnet" "glmnet"
$ lambda.min: num 0.142
$ lambda.1se: num 0.628
- attr(*, "class")= chr "cv.glmnet"
```

# crossvalidated results for ridge penalty

All 303 coefficients  
stay in model.

```
plot(ridge, main = "Ridge penalty\n\n")
```



```
> ridge$lambda.min  
[1] 0.1417314  
> log(ridge$lambda.min)  
[1] -1.953821  
> ridge$lambda.1se  
[1] 0.6279583  
> log(ridge$lambda.1se)  
[1] -0.4652815  
>
```

# cv-selected ridge performance on train and test

- Looks much better than the glmulti-selected model (which only included numeric variables)

```
> # compute RMSE of transformed predicted  
> RMSE(exp(predridge), exp(Y))  
[1] 28870.87  
> # can we run it on test?  
> Xtest <- model.matrix(Sale_Price ~ ., test_3)[, -1]  
> predridgetest <- predict(ridge, Xtest)  
> Ytest <- log(test_3$Sale_Price)  
> RMSE(exp(predridgetest), exp(Ytest))  
[1] 43656.97
```

s

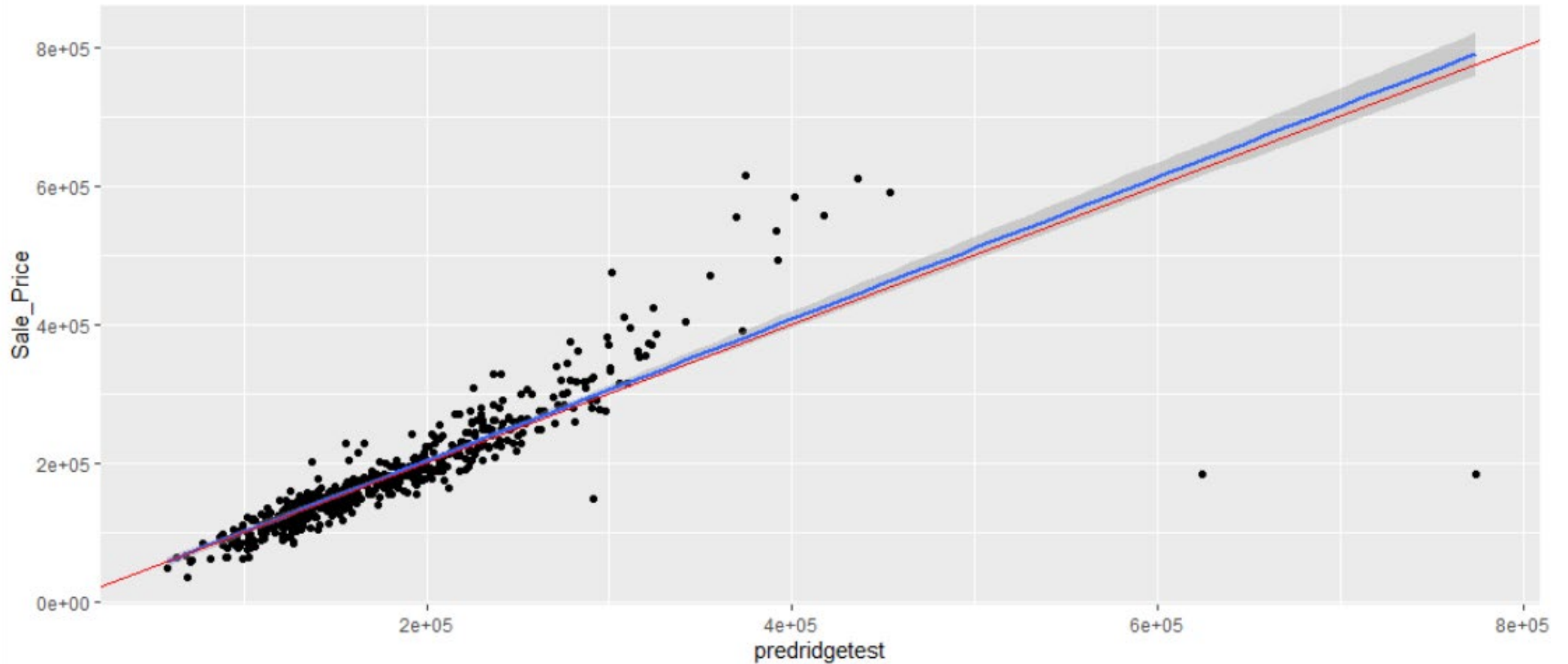
Value(s) of the penalty parameter `lambda` at which predictions are required. Default is the value `s="lambda.1se"` stored on the CV `object`. Alternatively `s="lambda.min"` can be used. If `s` is numeric, it is taken as the value(s) of `lambda` to be used. (For historical reasons we use the symbol 's' rather than 'lambda' to reference this parameter)

# Is lamda.min or lambda.1se better for inference on out-of-sample data?

- This is a question about parsimony. The lambda.min option refers to value of  $\lambda$  at the lowest CV error. The error at this value of  $\lambda$  is the average of the errors over the k folds and hence this estimate of the error is uncertain. The lambda.1se represents the value of  $\lambda$  in the search that was simpler than the best model (lambda.min), but which has error within 1 standard error of the best model. In other words, using the value of lambda.1se as the selected value for  $\lambda$  results in a model that is slightly simpler than the best model but which cannot be distinguished from the best model in terms of error given the uncertainty in the k-fold CV estimate of the error of the best model.
- The choice is yours:
- The best model that may be too complex or slightly overfitted: lambda.min
- The simplest model that has comparable error to the best model given the uncertainty: lambda.1se

<https://stats.stackexchange.com/questions/70249/feature-selection-model-with-glmnet-on-methylation-data-pn>

# Obs vs. exp plot for ames test\_3

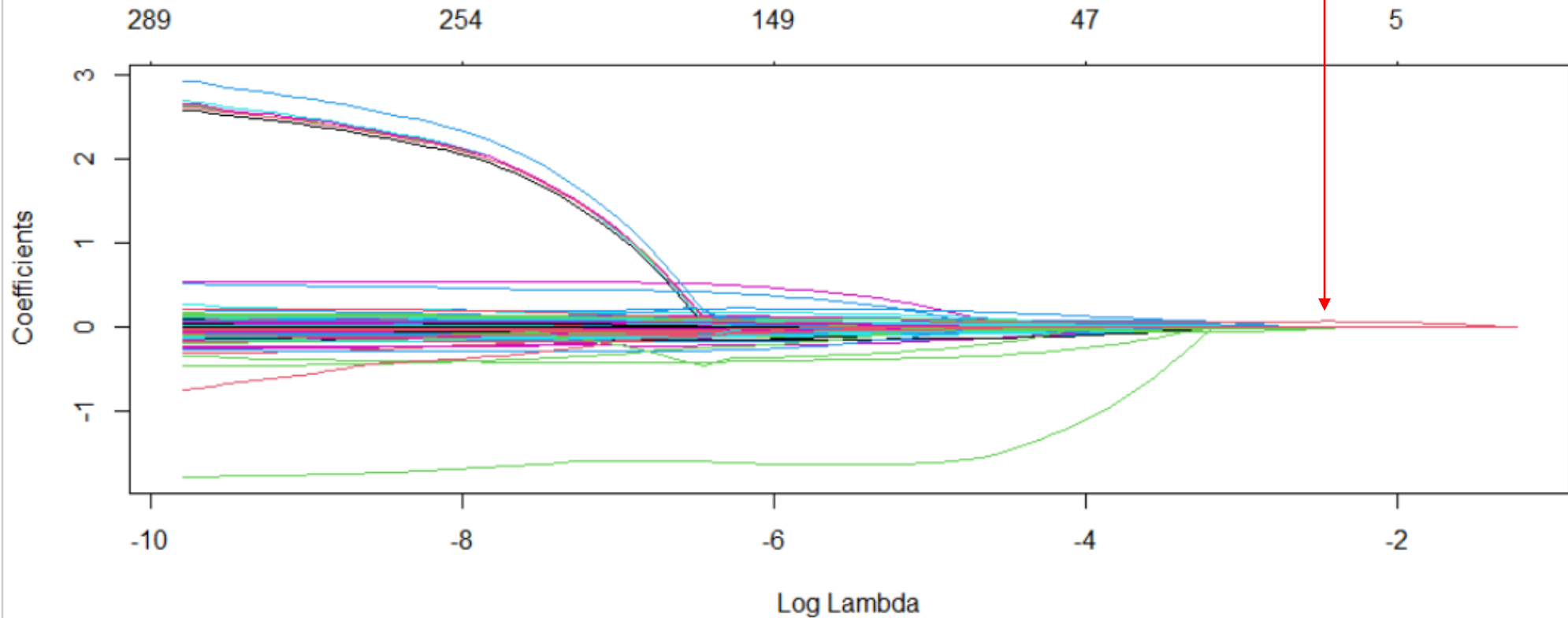


```
> cor(test_3aug$predridgetest, test_3aug$Sale_Price)^2  
[1]  
1 0.7264753  
> res_predridgetest <- test_3aug$Sale_Price - test_3aug$predridgetest  
> (mean((res_predridgetest)^2))^0.5  
[1] 43656.97
```

# Lasso on ames data

```
101  
102 #lasso on ames data  
103  
104 lasso <- glmnet(  
105   x = X,  
106   y = Y,  
107   alpha = 1  
108 )  
109  
110 plot(lasso, xvar = "lambda")  
111
```

Higher lambda values squeeze coefficients down to zero, creating a variable selection situation.



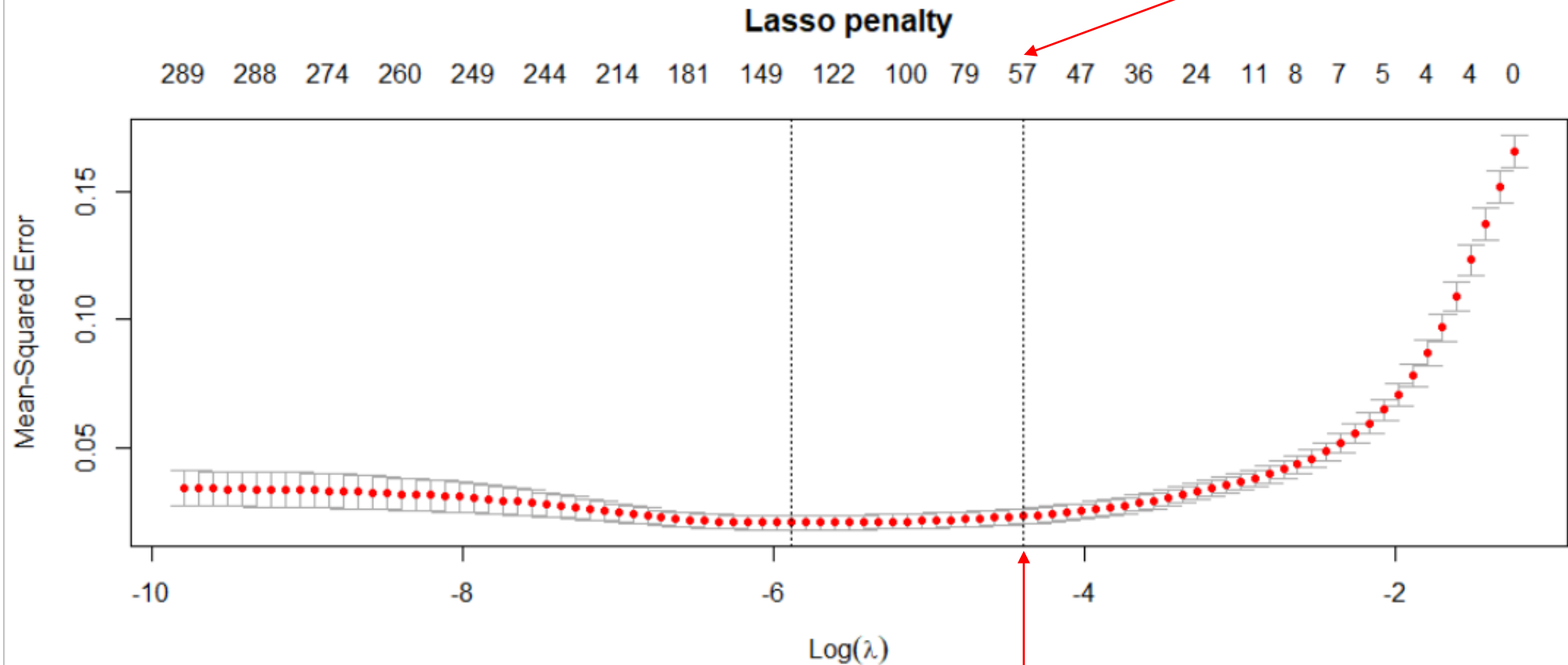
# Lasso on ames data

```
102 #lasso on ames data
103
104 # Apply CV lasso regression to Ames data
105 lasso <- cv.glmnet(
106   x = X,
107   y = Y,
108   alpha = 1
109 )
110 plot(lasso, main = "Lasso penalty\n\n")
111
112
```



# Lasso on ames data

Note coefficients being forced to 0 with L1 regularization, 57 left in model at 1SE rule



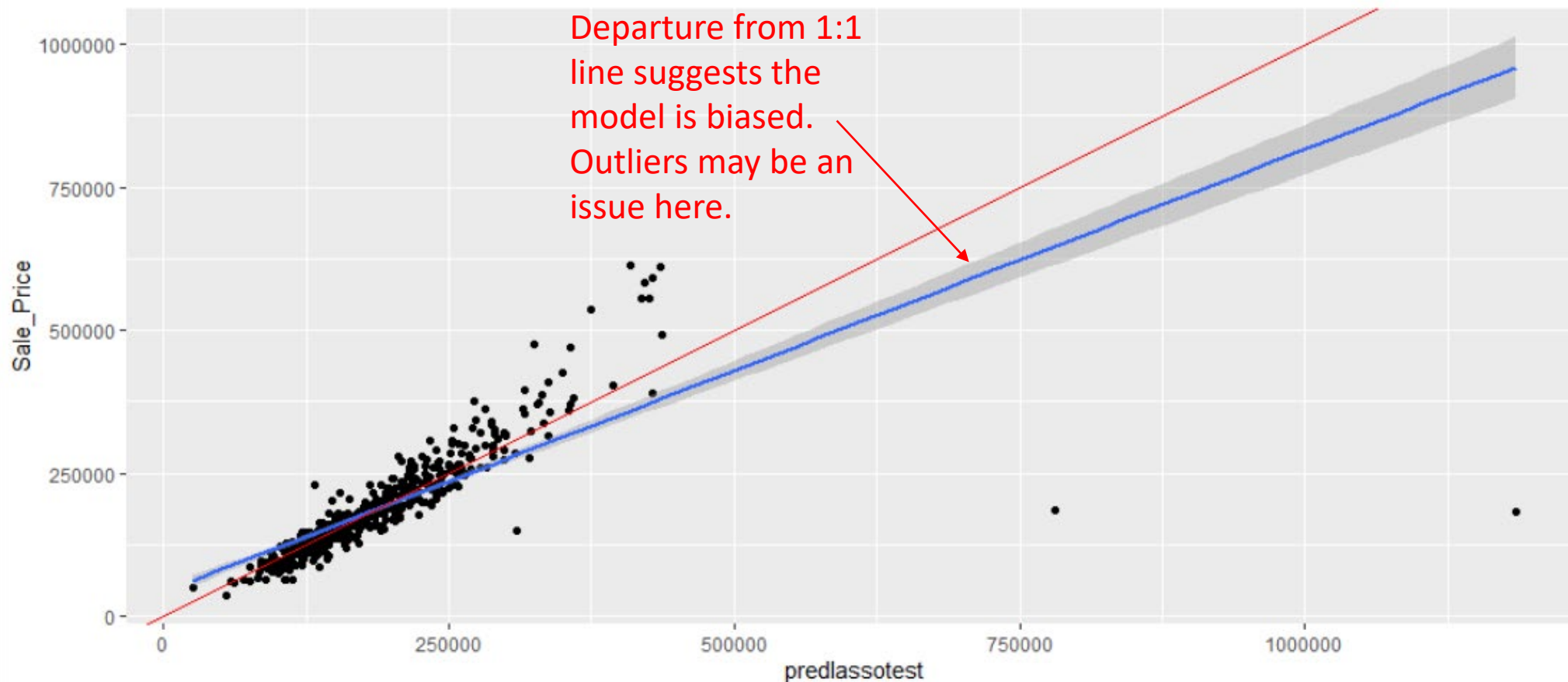
```
> lasso$lambda.1se  
[1] 0.01352895  
> #[1] 0.6279583  
> log(lasso$lambda.1se)  
[1] -4.302923  
> #[1] -0.4652815
```

# cv-selected lasso performance on train and test

- Not as good as ridge for this problem

```
> # evaluation on train_3 lasso
> predlasso <- predict(lasso, X)
> # compute RMSE of transformed predicted
> RMSE(exp(predlasso), exp(Y))
[1] 33823.32
> predlassotest <- predict(lasso, Xtest)
> RMSE(exp(predlassotest), exp(Ytest))
[1] 56674.53
.
```

# Obs vs. exp plot for lasso on test\_3



```
> cor(test_3aug$predlassotest, test_3aug$Sale_Price)^2  
[1] 0.5856687  
> res_predlassotest <- test_3aug$Sale_Price - test_3aug$predlassotest  
> (mean((res_predlassotest)^2))^0.5  
[1] 56674.53
```

# elastic net on ames data requires tuning of both lambda and alpha hyperparameters

- Grid search implemented using caret package

```
# elastic net model
# for reproducibility
set.seed(42)

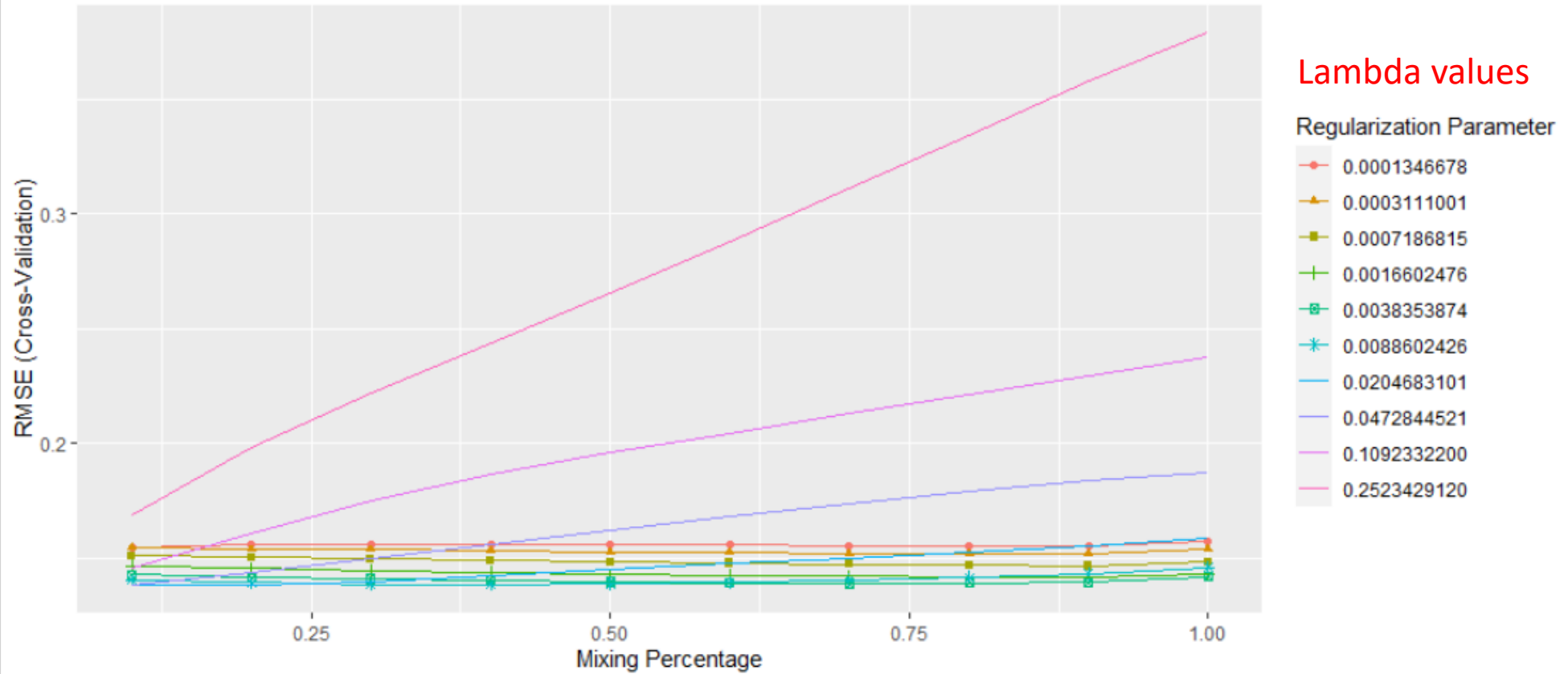
# grid search across
cv_glmnet <- train(
  x = X,
  y = Y,
  method = "glmnet",
  preProc = c("zv", "center", "scale"),
  trControl = trainControl(method = "cv", number = 10),
  tuneLength = 10
)

# model with lowest RMSE
cv_glmnet$bestTune
##   alpha    lambda
## 8    0.1 0.04728445

# results for model with lowest RMSE
cv_glmnet$results %>%
  filter(alpha == cv_glmnet$bestTune$alpha, lambda == cv_glmnet$bestTune$lambda)
##   alpha    lambda    RMSE  Rsquared    MAE    RMSESD  RsquaredSD
## 1    0.1 0.04728445 0.1382785 0.8852799 0.08427956 0.0273119 0.04179193
## MAESD
## 1 0.005330001
```

# elastic net tuning results

- `ggplot(cv_glmnet)`



Alpha values

# Elastic net performance on train and test

- Better than ridge on training set but worse on test!

```
> RMSE(exp(predelastic), exp(Y))  
[1] 25329.31  
> predelastictest <- predict(cv_glmnet, Xtest)  
> RMSE(exp(predelastictest), exp(Ytest))  
[1] 50741.77  
> test_3aug$predelastictest<-exp(predelastictest)  
> ggplot(test_3aug,aes(x=predelastictest,y=Sale_Price))+  
+   geom_point()+stat_smooth(method=lm)+geom_abline(slope=1, intercept=0  
+   `geom_smooth()` using formula 'y ~ x'  
> cor(test_3aug$predelastictest,test_3aug$Sale_Price)^2  
[1] 0.6622576  
> res_predelastictest<-test_3aug$Sale_Price-test_3aug$predelastictest  
> (mean((res_predelastictest)^2))^0.5  
[1] 50741.77  
<
```

# Variable importance for elastic net (top 20 variables)

- `vip(cv_glmnet, num_features = 20, geom = "point")`

