

Weekly Summary Template

Brady Miller

Table of contents

Tuesday, March 14	1
Thursday, March 16	10

Tuesday, March 14

! TIL

Include a *very brief* summary of what you learnt in this class here.

Today, I learnt the following concepts in class:

1. More with k-fold cross validation
2. Using caret package with cross validation
3. Using caret package with LASSO regression

```
# opening some necessary libraries
library(ISLR2)
library(dplyr)
```

Attaching package: 'dplyr'

The following objects are masked from 'package:stats':

filter, lag

The following objects are masked from 'package:base':

intersect, setdiff, setequal, union

```
library(tidyr)
library(purrr)
library(readr)
library(glmnet)
```

Loading required package: Matrix

Attaching package: 'Matrix'

The following objects are masked from 'package:tidyr':

expand, pack, unpack

Loaded glmnet 4.1-6

```
library(caret)
```

Loading required package: ggplot2

Loading required package: lattice

Attaching package: 'caret'

The following object is masked from 'package:purrr':

lift

```
library(car)
```

Loading required package: carData

Attaching package: 'car'

The following object is masked from 'package:purrr':

some

The following object is masked from 'package:dplyr':

recode

```
library(torch)
```

```
# opening Boston data set which will be used in examples
df <- Boston
attach(Boston)
```

k-fold cross validation

```
k <- 5
folds <- sample(1:k, nrow(df), replace = T)

df_folds <- list()

# define list of data frame where every list has train and test
for (i in 1:k){
  df_folds[[i]] <- list()
  df_folds[[i]]$train = df[which(folds == i), ]
}

# finding the mean squared error for each fold variation
kfold_mspe <- c()
for (i in 1:k) {
  model <- lm(medv ~ ., df_folds[[i]]$train)
  y_hat <- predict(model, df_folds[[i]]$test)
  kfold_mspe[i] <- mean((y_hat - df_folds[[i]]$test$medv)^2)
}

make_folds <- function(df, k){
  folds <- sample(1:k, nrow(df), replace = T)
  df_folds <- list()
  for (i in 1:k){
    df_folds[[i]] <- list()
  }
}
```

```

    df_folds[[i]]$train = df[which(folds != i), ]
    df_folds[[i]]$test = df[which(folds == i), ]
  }
  return(df_folds)
}

# cross validation mean squared prediction error function
cv_mspe <- function(formula, df_folds){
  kfold_mspe <- c()
  # going through each fold to generate predictions and find the error for each
  for (i in 1:length(df_folds)){
    model <- lm(formula, df_folds[[i]]$train)
    y_hat <- predict(model, df_folds[[i]]$test)
    kfold_mspe <- mean((y_hat - df_folds[[i]]$test$medv)^2)
  }
  return(mean(kfold_mspe))
}

```

- This function can then be used to find the prediction error generated for a given set of variables over a certain amount of folds. This can help you determine which set of variables would be best to create the least error

cross validation with caret package

By using the caret package, we can write code to cross validate a dataset in much fewer lines as shown below

```

# specifying you want to cross validate using 5 folds (the number)
ctrl <- trainControl(method = 'cv', number = 5)
ctrl

```

```

$method
[1] "cv"

```

```

$number
[1] 5

```

```

$repeats
[1] NA

```

```

$search

```

```

[1] "grid"

$p
[1] 0.75

$initialWindow
NULL

$horizon
[1] 1

$fixedWindow
[1] TRUE

$skip
[1] 0

$verboseIter
[1] FALSE

$returnData
[1] TRUE

$returnResamp
[1] "final"

$savePredictions
[1] FALSE

$classProbs
[1] FALSE

$summaryFunction
function (data, lev = NULL, model = NULL)
{
  if (is.character(data$obs))
    data$obs <- factor(data$obs, levels = lev)
  postResample(data[, "pred"], data[, "obs"])
}
<bytecode: 0x000001bc7b3fe480>
<environment: namespace:caret>

$selectionFunction

```

```
[1] "best"

$preProcOptions
$preProcOptions$thresh
[1] 0.95

$preProcOptions$ICAcomp
[1] 3

$preProcOptions$k
[1] 5

$preProcOptions$freqCut
[1] 19

$preProcOptions$uniqueCut
[1] 10

$preProcOptions$cutoff
[1] 0.9

$sampling
NULL

$index
NULL

$indexOut
NULL

$indexFinal
NULL

$timingSamps
[1] 0

$predictionBounds
[1] FALSE FALSE

$seeds
[1] NA
```

```
$adaptive
$adaptive$min
[1] 5
```

```
$adaptive$alpha
[1] 0.05
```

```
$adaptive$method
[1] "gls"
```

```
$adaptive$complete
[1] TRUE
```

```
$trim
[1] FALSE
```

```
$allowParallel
[1] TRUE
```

```
# creating model that uses linear regression model to predict med. house price
# trControl attribute uses the ctrl object to specify 5-fold cross validation
model <- train(medv ~ ., data = df, method = 'lm', trControl = ctrl)
summary(model)
```

```
Call:
lm(formula = .outcome ~ ., data = dat)
```

```
Residuals:
      Min       1Q   Median       3Q      Max
-15.1304  -2.7673  -0.5814   1.9414  26.2526
```

```
Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  41.617270   4.936039   8.431 3.79e-16 ***
crim         -0.121389   0.033000  -3.678 0.000261 ***
zn           0.046963   0.013879   3.384 0.000772 ***
indus        0.013468   0.062145   0.217 0.828520
chas         2.839993   0.870007   3.264 0.001173 **
nox        -18.758022   3.851355  -4.870 1.50e-06 ***
rm           3.658119   0.420246   8.705 < 2e-16 ***
```

```

age          0.003611    0.013329    0.271 0.786595
dis         -1.490754    0.201623   -7.394 6.17e-13 ***
rad          0.289405    0.066908    4.325 1.84e-05 ***
tax         -0.012682    0.003801   -3.337 0.000912 ***
ptratio     -0.937533    0.132206   -7.091 4.63e-12 ***
lstat       -0.552019    0.050659  -10.897 < 2e-16 ***

```

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 4.798 on 493 degrees of freedom

Multiple R-squared: 0.7343, Adjusted R-squared: 0.7278

F-statistic: 113.5 on 12 and 493 DF, p-value: < 2.2e-16

This model created by cross validation using the caret package indicates that the age and indus variables are not significant in predicting the median house price of a given home. The R^2 value indicates that there is a somewhat strong positive correlation between the covariates and the outcome variable.

```

# creates predictions for each piece of data in the Boston data set
predictions <- predict(model,df)
sample(predictions,10)

```

```

      368      394      330      5      82      125      317      49
12.826665 19.318410 24.287581 27.928791 27.071925 20.496661 17.414382  8.743247
      33      207
 9.859984 23.549027

```

Above we show 5 randomly selected predictions made on the data set. Each row number has the associated predicted median house value with it.

LASSO regression with caret package

This deals with bias-variance trade off. As it states there is a trade off between having more bias or more variance. By having a model that has more variables, bias decreases but variance increases. Having too much variance can lead to overfitting data and create good performance on train data but bad performance on test data. On the other hand, a model with only a few variables increases bias and reduces variance. Increasing bias can lead to creating a model that underfits data, resulting in poor performance on both the train and test data. This relates to LASSO as it is a form of variable selection, so you want to make sure you are selecting the right amount of variables so you have a balanced trade off between bias and variance


```

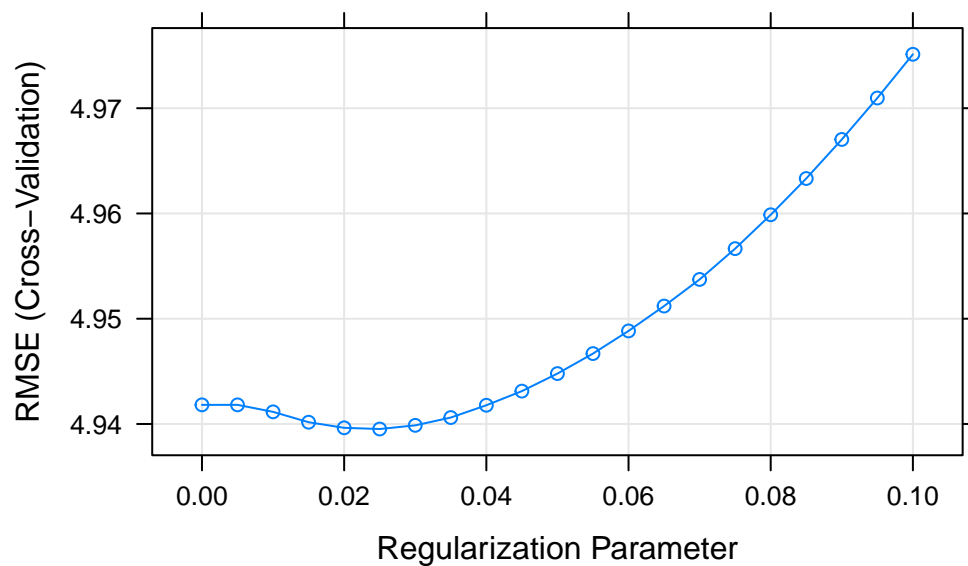
ctrl <- trainControl(method = 'cv', number = 5)

# Defining the tuning grid
grid <- expand.grid(alpha = 1, lambda = seq(0, 0.1, by = 0.005))

# Train the model using LASSO regression with cross validation
lasso_fit <- train(medv ~ ., data = df, method = 'glmnet', trControl = ctrl,
                  tuneGrid = grid, standardize = TRUE, family = 'gaussian')

plot(lasso_fit)

```



This plot shows the mean squared error value for each regularization parameter tried in the LASSO regression with the cross validation. From this, we can see that the regularization parameter that minimizes the mean squared error value is at about 0.03, so λ should be 0.03 when doing LASSO regression.

Thursday, March 16

! TIL

Include a *very brief* summary of what you learnt in this class here.

Today, I learnt the following concepts in class:

1. Why linear regression won't work with binary/categorical response variables for classification problems
2. Logistic regression using sigmoid function
3. Interpreting logistic regression/model summary with a real data set

Why linear regression won't work for binary/categorical classification

```
url <- 'https://archive.ics.uci.edu/ml/machine-learning-databases/breast-cancer-wisconsin/  
col_names <- c('id', 'diagnosis', paste0('feat', 1:30))  
df <- read_csv(  
  url, col_names, col_types = cols()  
) %>%  
  select(-id) %>%  
  mutate(outcome = ifelse(diagnosis == 'M', 1, 0)) %>%  
  select(-diagnosis)
```

```
reg_model <- lm(outcome ~ ., data = df)  
summary
```

standardGeneric for "summary" defined from package "base"

```
function (object, ...)  
standardGeneric("summary")  
<environment: 0x000001bc7d37ce98>  
Methods may be defined for arguments: object  
Use showMethods(summary) for currently available ones.
```

Model summary indicates that most features are not significant

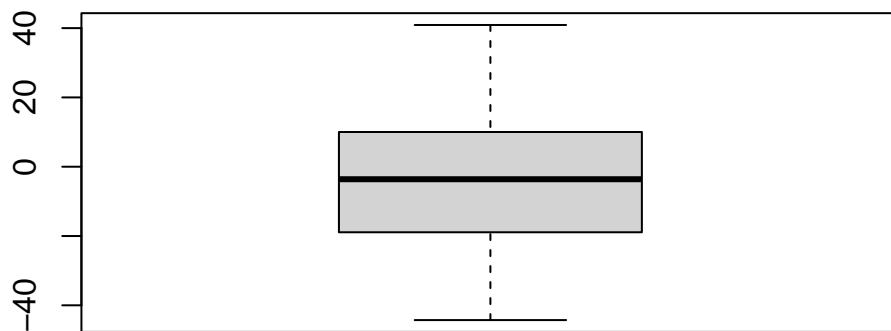
```
n <- 100  
new_patients <- data.frame(matrix(rnorm(30 * n), nrow = n))  
colnames(new_patients) <- paste0('feat', 1:30)  
new_predictions <- predict(reg_model, newdata = new_patients, type = 'response')
```

```
print(new_predictions %>% head())
```

	1	2	3	4	5	6
	8.996807	-11.319676	2.263436	-25.904301	9.840753	-33.990923

- The responses should be 0 or 1 (tumor is malignant or benign), but what we get from this is values that don't make sense → can't represent probability that someones tumor is malignant
- logistic regression would be better as this model can end up giving negative numbers which doesn't make sense (meaningless predictions)
- linear regression for binary models will give you bad values for extrapolation

```
boxplot(new_predictions)
```



The range of values going negative shows that there are bad/inaccurate responses

Logistic regression

*odds = $p/(1-p)$

*useful to interpret for binary responses

*if p is between (0,1) then odds is between (0, infinity)

*log-odds = $\log(\text{odds}) = \log(p/(1-p))$. log-odds takes values between (-infinity, infinity).
creates continuous scale which you can do linear regression on (transforming scale so you can do linear reg, and then transforming it back to get value)

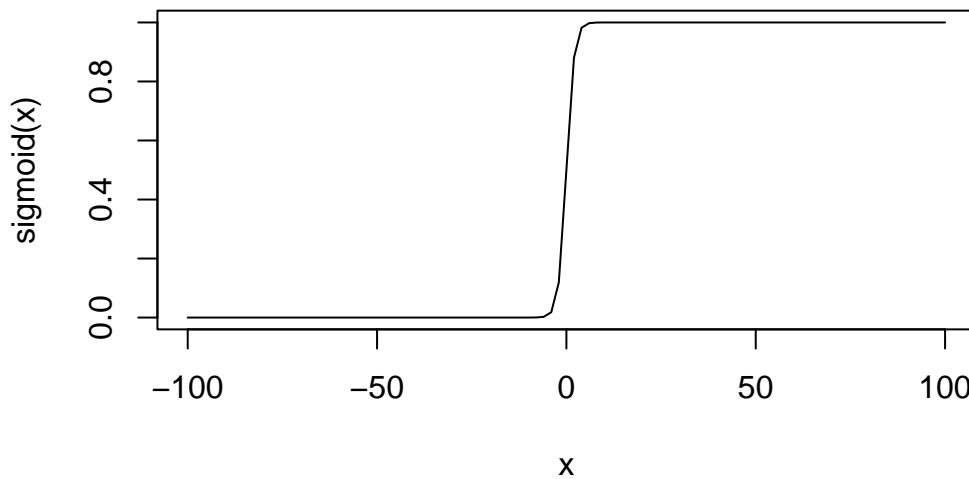
- will specify a regression model and then use log-odds to get back what the probabilities are

$$\frac{p}{1-p} = \exp(\beta_0 + \beta_1 x_1)$$
$$\text{logodds}(\beta(x)) = \beta_0 + \beta_1 x_1$$

$$p(x) = \frac{e^{(\beta_0 + \beta_1 x_1)}}{1 + e^{(\beta_0 + \beta_1 x_1)}} > \frac{1}{1 + e^{(\beta_0 + \beta_1 x_1)}}$$

SIGMOID FUNCTION - logistic regression

```
sigmoid <- \(x) 1/(1+exp(-x))  
curve(sigmoid, -100, 100, ylab = 'sigmoid(x)')
```



$$p(x) = \text{sigmoid}(\beta_0 + \beta_1 x_1) = \frac{1}{1 + \exp(-\beta_0 - \beta_1 x_1)}$$

Logistic regression for breast cancer data set

```
df <- df %>% mutate_at('outcome', factor)

model <- glm(outcome ~., data =df, family = binomial())
```

Warning: glm.fit: algorithm did not converge

Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

```
summary(model)
```

Call:

```
glm(formula = outcome ~ ., family = binomial(), data = df)
```

Deviance Residuals:

Min	1Q	Median	3Q	Max
-8.49	-8.49	-8.49	8.49	8.49

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-2.881e+06	2.816e+05	-10.233	< 2e-16 ***
feat1	2.427e+06	2.693e+05	9.014	< 2e-16 ***
feat2	1.958e+05	1.471e+04	13.313	< 2e-16 ***
feat3	1.473e+06	2.464e+04	59.791	< 2e-16 ***
feat4	-1.301e+05	3.907e+03	-33.301	< 2e-16 ***
feat5	-1.525e+08	8.361e+06	-18.234	< 2e-16 ***
feat6	-6.428e+06	3.213e+06	-2.001	0.04539 *
feat7	1.042e+06	1.408e+06	0.740	0.45959
feat8	-1.716e+07	5.382e+06	-3.188	0.00143 **
feat9	4.049e+07	7.772e+05	52.093	< 2e-16 ***
feat10	-4.233e+07	2.169e+06	-19.519	< 2e-16 ***
feat11	3.328e+07	1.169e+06	28.478	< 2e-16 ***
feat12	6.368e+06	2.005e+05	31.763	< 2e-16 ***
feat13	1.701e+06	4.720e+04	36.032	< 2e-16 ***
feat14	-6.393e+05	1.835e+04	-34.840	< 2e-16 ***

feat15	7.492e+08	1.224e+07	61.213	< 2e-16	***
feat16	-1.773e+08	5.732e+06	-30.931	< 2e-16	***
feat17	1.529e+08	5.340e+06	28.624	< 2e-16	***
feat18	-1.260e+09	4.012e+07	-31.398	< 2e-16	***
feat19	2.890e+08	4.126e+06	70.055	< 2e-16	***
feat20	1.512e+09	6.597e+07	22.921	< 2e-16	***
feat21	-6.130e+06	2.143e+05	-28.606	< 2e-16	***
feat22	-5.832e+05	2.437e+04	-23.935	< 2e-16	***
feat23	-3.538e+05	1.219e+04	-29.023	< 2e-16	***
feat24	8.950e+04	2.741e+03	32.658	< 2e-16	***
feat25	-2.161e+07	3.298e+06	-6.553	5.66e-11	***
feat26	8.986e+06	3.999e+05	22.470	< 2e-16	***
feat27	-3.028e+07	1.523e+06	-19.875	< 2e-16	***
feat28	1.431e+08	5.471e+06	26.162	< 2e-16	***
feat29	-2.474e+07	3.392e+05	-72.923	< 2e-16	***
feat30	-3.698e+07	5.340e+06	-6.926	4.32e-12	***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 751.44 on 568 degrees of freedom
 Residual deviance: 32006.76 on 538 degrees of freedom
 AIC: 32069

Number of Fisher Scoring iterations: 25

- Variables we initially thought were insignificant are now shown to be important
- Feature 7 is not significant at all → is actually significant, just very collinear with another variable
- Logistic regression suffers from multicollinearity
- There is no R^2 value for logistic regression as there is for linear regression → no line fitting x to y
- Computes deviance instead (similar but different)
- Also shows AIC