# KNN Tutorial

Data Mining for Business Intelligence

# OUTLINE

- KNN in R with IBk

- IBk package defaults and parameters.

- Exploring Train vs Test set results

- Cross-Validation Performance of various argument settings.

- Grid-Search and IBk best models

- Ablation Analysis

# IBk – Package Introduction

See IBk documentation and parameter settings at
http://weka.sourceforge.net/doc.dev/weka/classifiers/lazy/IBk.html
Parameters:

-I
Weight neighbours by the inverse of their distance  (use when k > 1)

-F
Weight neighbours by 1 - their distance (use when k > 1)

-K <number of neighbors>
Number of nearest neighbours (k) used in classification. (Default = 1)

-X
Select the number of nearest neighbours between 1  and the k value specified using hold-one-out evaluation   on the training data (use when k > 1)

-E
Minimize mean squared error rather than mean absolute  error when using -X option with numeric prediction.

# Data Import and Categorical Factorization

- Here we import the needed libraries

- Import the required dataset

- Create variables for use throughout the tutorial such as the metrics list and the seed value.

- Explore the structure and summarize the dataset.

```r
44  library(caret)
45  library(RWeka)
46  library(rminer)
47  library(matrixStats)
48  library(knitr)
49  library(tictoc)
50  library(tidyverse)
51  tic()
52
53  # import data
54  cloud_wd <- getwd()
55  setwd(cloud_wd)
56  insurance <- read.csv(file = "insurance.csv", stringsAsFactors = TRUE)
57
58  ### Set up cv parameters
59
60  df <- insurance
61  target <- 7
62  seedval <- 500
63  metrics_list <- c("MAE","RMSE","MAPE","RMSPE")
64  ```

65
66  ```{r}
67  str(insurance)
68  ```
```
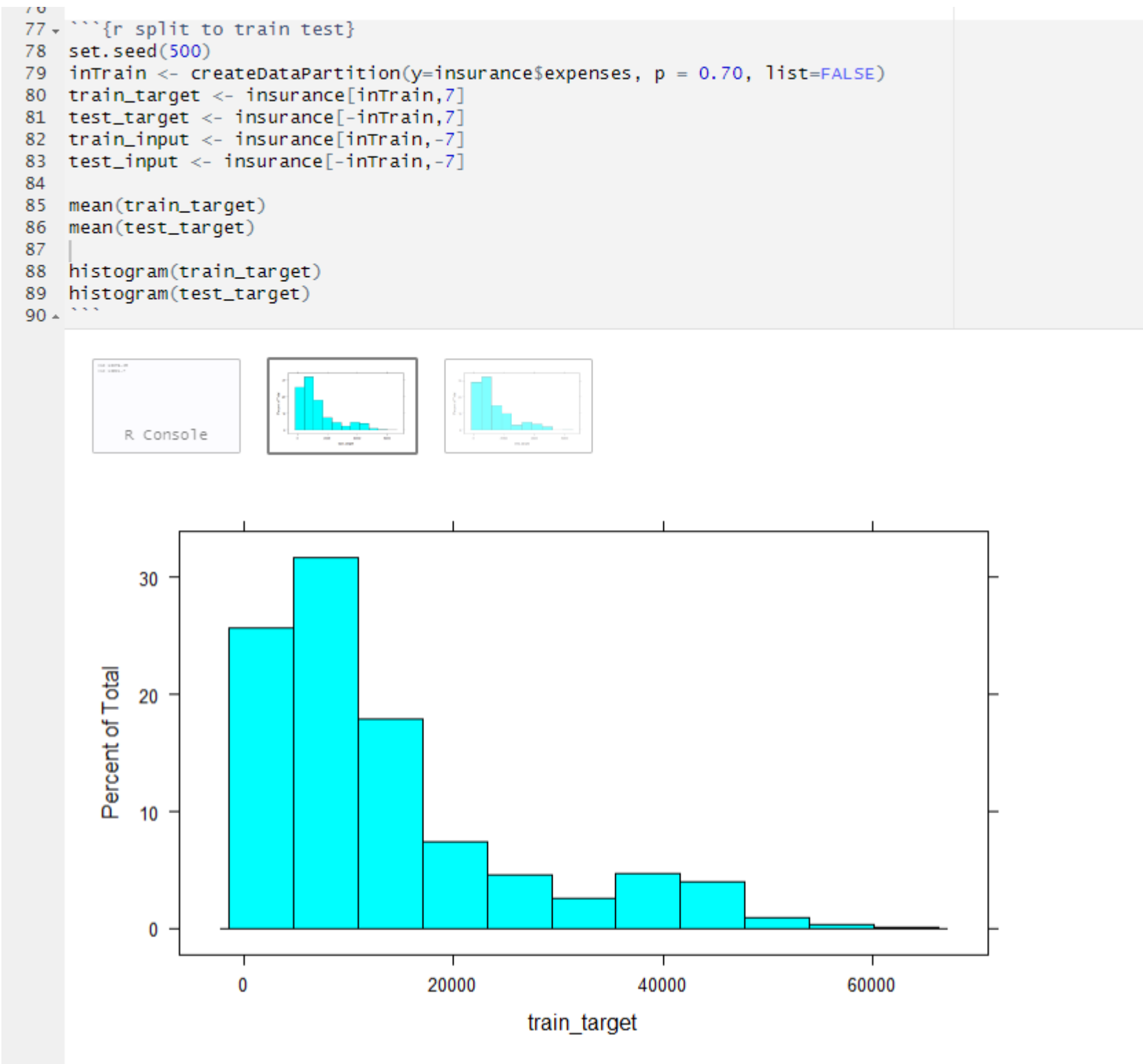
```
'data.frame':   1338 obs. of  7 variables:
 $ age     : int  19 18 28 33 32 31 46 37 37 60 ...
 $ sex     : Factor w/ 2 levels "female","male": 1 2 2 2 2 1 1 1 2 1 ...
 $ bmi     : num  27.9 33.8 33 22.7 28.9 25.7 33.4 27.7 29.8 25.8 ...
 $ children: int  0 1 3 0 0 0 1 3 2 0 ...
 $ smoker  : Factor w/ 2 levels "no","yes": 2 1 1 1 1 1 1 1 1 1 ...
 $ region  : Factor w/ 4 levels "northeast","northwest",..: 4 3 3 2 2 3 3 2 1 2 ...
 $ expenses: num  16885 1726 4449 21984 3867 ...
```

```
69  ```{r}
70  summary(insurance)
71  ```
```

```
      age             sex           bmi          children       smoker           region       expenses
 Min.   :18.00   female:662   Min.   :16.00   Min.   :0.000   no :1064   northeast:324   Min.   : 1122
 1st Qu.:27.00   male  :676   1st Qu.:26.30   1st Qu.:0.000   yes: 274   northwest:325   1st Qu.: 4740
 Median :39.00                Median :30.40   Median :1.000              southeast:364   Median : 9382
 Mean   :39.21                Mean   :30.67   Mean   :1.095              southwest:325   Mean   :13270
 3rd Qu.:51.00                3rd Qu.:34.70   3rd Qu.:2.000                              3rd Qu.:16640
 Max.   :64.00                Max.   :53.10   Max.   :5.000                              Max.   :63770
```

# Data Partition – Simple Holdout

- Split the data into test and train sets, extract the target variable into a new variable for both sets.

- Show the distribution of the target in both train and test to prove both sets are similar when comparing via the target.

```r
77 - ```{r split to train test}
78  set.seed(500)
79  inTrain <- createDataPartition(y=insurance$expenses, p = 0.70, list=FALSE)
80  train_target <- insurance[inTrain,7]
81  test_target <- insurance[-inTrain,7]
82  train_input <- insurance[inTrain,-7]
83  test_input <- insurance[-inTrain,-7]
84
85  mean(train_target)
86  mean(test_target)
87
88  histogram(train_target)
89  histogram(test_target)
90 - ```
```

# K = 1 on Train and Test

- Compare model performance from Train to Test using a K = 1 (1 nearest neighbor setting).
- Train set errors are 0 because each instance can be predicted using itself resulting in no errors at all.
- Test set errors are much higher. The nearest neighbor in Test is not the exact same observation.

```r
96  when using K = 1 on the training set we will then achieve an MAE = 0 (or if classification 100% accuracy)
97
98  ## MAE for K = 1
99
100 ```{r simple IBk examples}
101
102 # K = 1
103
104 knn_model_k_1 <- IBk(train_target ~ .,data = train_input,control = Weka_control(K=1))
105 knn_model_k_1
106 ```

    IB1 instance-based classifier
    using 1 nearest neighbour(s) for classification


107
108 ```{r k1 train}
109 insample_pred <- predict(knn_model_k_1, train_input)
110 mmetric(train_target, insample_pred, metrics_list)
111 ```

     MAE   RMSE   MAPE RMSPE
       0      0      0      0


112
113 ```{r k1 test}
114 test_pred <- predict(knn_model_k_1, test_input)
115 mmetric(test_target, test_pred, metrics_list)
116 ```

           MAE          RMSE          MAPE         RMSPE
     3803.667750   7217.775045     37.651870      9.425661
```

# K = 5 on Train and Test

- Higher errors on both
  Train (expected) and Test.

- It appears K = 5 using
  unweighted voting is not
  as good as K = 1.

```r
120
121  ## MAE for K = 5
122
123  ```{r}
124  knn_model_k_5 <- IBk(train_target ~ .,data = train_input,control = Weka_control(K=5))
125  knn_model_k_5
126  ```

     IB1 instance-based classifier
     using 5 nearest neighbour(s) for classification

127
128  ```{r k5 train}
129  insample_pred <- predict(knn_model_k_5, train_input)
130  mmetric(train_target, insample_pred, metrics_list)
131  ```

              MAE           RMSE          MAPE          RMSPE
     2827.042013 4698.609321     26.490765      5.015548

132
133  ```{r k5 test}
134  test_pred <- predict(knn_model_k_5, test_input)
135  mmetric(test_target, test_pred, metrics_list)
136  ```

             MAE           RMSE          MAPE          RMSPE
     3994.90551 6388.61020     41.18298      7.17385

137
```

# K = 1 and I = TRUE on Train and Test

- The errors shouldn't change from using only k=1 since we are using only 1 neighbor. Weighting is useful when k > 1.

```r
139  ## weighted distance
140
141  ## MAE for K = 1 and I = TRUE
142
143  ```{r Simple IBk examples I=TRUE}
144
145  # K = 1
146
147  knn_model_k1_itrue <- IBk(train_target ~ .,data = train_input,control = weka_control(K=1,I=TRUE))
148  knn_model_k1_itrue
149  ```

     IB1 instance-based classifier
     using 1 inverse-distance-weighted nearest neighbour(s) for classification


150
151  ```{r k1 train I=TRUE}
152  insample_pred <- predict(knn_model_k1_itrue, train_input)
153  mmetric(train_target, insample_pred, metrics_list)
154  ```

     MAE   RMSE   MAPE  RMSPE
       0     0     0     0


155
156  ```{r k1 test I=TRUE}
157  test_pred <- predict(knn_model_k1_itrue, test_input)
158  mmetric(test_target, test_pred, metrics_list)
159  ```

          MAE         RMSE        MAPE        RMSPE
     3803.667750 7217.775045   37.651870    9.425661
```

# K = 5 and I = TRUE on Train and Test

- Weighted voting has improved Train performance due to more weight on the exact same instance being the most heavily weighted neighbor.
- In test we've also improved our performance as compared to K = 5 without weighted voting (I=FALSE)

```
163
164 ## MAE for K = 5 and I = TRUE
165
166 ```{r}
167 knn_model_k5_itrue <- IBk(train_target ~ .,data = train_input,control = Weka_control(K=5,I=TRUE))
168 knn_model_k5_itrue
169 ```

    IB1 instance-based classifier
    using 5 inverse-distance-weighted nearest neighbour(s) for classification

170
171 ```{r k5 train I=TRUE}
172 insample_pred <- predict(knn_model_k5_itrue, train_input)
173 mmetric(train_target, insample_pred, metrics_list)
174 ```

            MAE          RMSE          MAPE         RMSPE
    239.6097554 475.9304220     2.6135382    0.7020244

175
176 ```{r k5 test I=TRUE}
177 test_pred <- predict(knn_model_k5_itrue, test_input)
178 mmetric(test_target, test_pred, metrics_list)
179 ```

            MAE          RMSE          MAPE         RMSPE
    3821.653457 6260.476947    39.918157     7.330998
```

# K = from 1 to 40 and I = TRUE and X=TRUE on Train and Test

- We now ask the model to eval k values from 1 to 40 using leave one out cross validation. The best K will be selected.
- The model has chosen a K = 3 based on its ability to reduce the MAE.

```
184  ## MAE for K = 40, I = TRUE, X = TRUE.
185
186  Try k from 1 to 40, using inverse of distance weights.
187
188  Demonstrating how to use the Leave one out cross-validation built into the IBk package to determine the best k.
189
190  ```{r X TRUE}
191  knn_model_itrue_xtrue <- IBk(train_target ~ .,data = train_input,control = Weka_control(K=40,I=TRUE,X=TRUE))
192  knn_model_itrue_xtrue
193  ```

     IB1 instance-based classifier
     using 3 inverse-distance-weighted nearest neighbour(s) for classification


194  The model has chosen k = 3 to minimize the MAE.
195
196  ```{r}
197  insample_pred <- predict(knn_model_itrue_xtrue, train_input)
198  mmetric(train_target, insample_pred, metrics_list)
199  ```

            MAE         RMSE         MAPE         RMSPE
     159.9757828  372.7041571    1.8370163    0.6670268


200
201  ```{r}
202  test_pred <- predict(knn_model_itrue_xtrue, test_input)
203  mmetric(test_target, test_pred, metrics_list)
204  ```

            MAE         RMSE         MAPE         RMSPE
     3813.619044  6464.275066    40.250375    8.281032
```

# IBk custom cross-validation function

- Create a custom cross validation function which allows us to explore k and I values.

- In addition, we change the end of this function to include columns for the k and i parameters as well as return a dataframe instead of a kable.

```r
216
217 ```{r Define cv_IBk}
218
219 cv_IBk <- function(df, target, nFolds, seedval, metrics_list, k, i)
220 {
221 # create folds using the assigned values
222
223 set.seed(seedval)
224 folds = createFolds(df[,target],nFolds)
225
226 # The lapply loop
227
228 cv_results <- lapply(folds, function(x)
229 {
230 # data preparation:
231
232   test_target <- df[x,target]
233   test_input <- df[x,-target]
234
235   train_target <- df[-x,target]
236   train_input <- df[-x,-target]
237   pred_model <- IBk(train_target ~ .,data = train_input,control = Weka_control(K=k,I=i))
238   pred <- predict(pred_model, test_input)
239   train_pred <- predict(pred_model, train_input)
240
241   return(mmetric(test_target,pred,metrics_list))
242 })
243
244 cv_results_m <- as.matrix(as.data.frame(cv_results))
245 cv_mean<- as.matrix(rowMeans(cv_results_m))
246 cv_sd <- as.matrix(rowSds(cv_results_m))
247 colnames(cv_mean) <- "Mean"
248 colnames(cv_sd) <- "Sd"
249 cv_df <- data.frame(t(cbind(cv_mean,cv_sd))) %>% round(2)
250 cv_df$param_K <- k
251 cv_df$param_I <- as.logical(i)
252 cv_df <- cv_df %>% rownames_to_column(var = "measure")
253 cv_df
254 }
255 ```
256
```

# IBk custom cross-validation function

- We run the new cv_function with various values of k and I = FALSE

```r
262  ```{r call cv_IBk_train and cv_IBk}
263  cv_IBk(df, target, 3, seedval, metrics_list, 1, FALSE)
264  ```
```

Description: df [2 × 7]

| measure<br><chr> | MAE<br><dbl> | RMSE<br><dbl> | MAPE<br><dbl> | RMSPE<br><dbl> | param_K<br><dbl> | param_I<br><lgl> |
|---|---|---|---|---|---|---|
| Mean | 3728.35 | 7169.88 | 39.04 | 11.91 | 1 | FALSE |
| Sd | 218.95 | 302.90 | 3.17 | 1.54 | 1 | FALSE |

2 rows

```r
265
266
267  ## Call cross-validation with k=2 and I=FALSE
268  ```{r}
269  cv_IBk(df, target, 3, seedval, metrics_list, 2, FALSE)
270  ```
```

Description: df [2 × 7]

| measure<br><chr> | MAE<br><dbl> | RMSE<br><dbl> | MAPE<br><dbl> | RMSPE<br><dbl> | param_K<br><dbl> | param_I<br><lgl> |
|---|---|---|---|---|---|---|
| Mean | 3497.67 | 6212.26 | 38.02 | 9.81 | 2 | FALSE |
| Sd | 68.79 | 88.36 | 3.54 | 1.55 | 2 | FALSE |

2 rows

```r
271
272  ## Call cross-validation with k=3 and I=FALSE
273  ```{r}
274  cv_IBk(df, target, 3, seedval, metrics_list, 3, FALSE)
275  ```
```

Description: df [2 × 7]

| measure<br><chr> | MAE<br><dbl> | RMSE<br><dbl> | MAPE<br><dbl> | RMSPE<br><dbl> | param_K<br><dbl> | param_I<br><lgl> |
|---|---|---|---|---|---|---|
| Mean | 3575.87 | 5986.29 | 38.02 | 8.01 | 3 | FALSE |
| Sd | 115.05 | 202.05 | 3.05 | 0.96 | 3 | FALSE |

2 rows

# Parameter Grid Setups

- Build a grid of k values from 1 to 40 and I options (TRUE,FALSE)
- Notice our dataframe is 80 rows long.

```r
313
314 # Part 3.
315
316 Use a custom grid search to find our best hyperparameter combination.
317
318 ## Build the grid
319
320 ```{r}
321 # Create multiple vectors
322 param_k <- c(seq(1,40))
323 param_i <- c(FALSE, TRUE)
324
325 # Generate a grid of all combinations
326 grid <- expand.grid(param_k, param_i, stringsAsFactors = FALSE)
327
328 colnames(grid) <- c("k","i")
329 # Print the grid
330
331 grid
332 ```
```

Description: df [80 × 2]

| k <int> | i <lgl> |
|---|---|
| 1 | FALSE |
| 2 | FALSE |
| 3 | FALSE |
| 4 | FALSE |
| 5 | FALSE |
| 6 | FALSE |
| 7 | FALSE |
| 8 | FALSE |
| 9 | FALSE |
| 10 | FALSE |

1-10 of 80 rows

# Parameter Grid Looping function

- Define a new named function called "run_grid_cv" which loops over each row in the dataframe and calls the cross validation function for the rows setups for I and K.

- This function will loop 80 times based on the previous grid setups.

```r
334 ▾ ## define a named function to loop through the grid
335
336    This function will take the grid as input and call the cross validation function for each row in the grid. Finally, each cross vali
       deviation will be returned as a dataframe.
337
338 ▾ ```{r}
339
340 ▾ run_grid_cv <- function(grid) {
341      results <- data.frame()|
342 ▾    for (i in 1:nrow(grid)) {
343        row <- grid[i,]  # Get the i-th row
344        cv_result <- cv_IBk(df = df,target = target,nFolds = 5,seedval =  seedval,metrics_list =  metrics_list,k = row$k,i = row$i)
345        results <- rbind(results, cv_result)
346 ▴    }
347      return(results)
348 ▴ }
349 ▴ ```
```

# Calling the grid search function

- We simply pass the grid as a parameter into the new named function.
- The function takes the grid and loops over each row calling the cross validation function each time.
- We store the resulting dataframe in a variable named cv_grid_results.

```r
351   ## run the grid search, return a dataframe with the results.
352
353   ```{r}
354   cv_grid_results <- run_grid_cv(grid)
355   ```
356
```

# Filter the result dataframe. Find the best 10 MAE.

- We now demonstrate how to filter the dataframe to find only the Mean results, and arrange the dataframe in ascending order to show the lowest MAE. We use the head function to show only the best 10 MAE results. k = 2 is the winner.

```r
358  ```{r}
359  # pull the top 3 based on MAE
360
361  cv_grid_results %>%
362    filter(measure == 'Mean') %>% |
363    arrange(MAE) %>%
364    head(10)
365  ```
```

Description: df [10 × 7]

| | measure <chr> | MAE <dbl> | RMSE <dbl> | MAPE <dbl> | RMSPE <dbl> | param_K <int> | param_I <lgl> |
|---|---|---|---|---|---|---|---|
| 1 | Mean | 3261.69 | 6012.29 | 32.22 | 8.18 | 2 | TRUE |
| 2 | Mean | 3325.64 | 5850.02 | 33.11 | 7.40 | 3 | TRUE |
| 3 | Mean | 3357.56 | 6021.81 | 33.29 | 8.00 | 2 | FALSE |
| 4 | Mean | 3366.06 | 5684.82 | 33.17 | 6.46 | 5 | TRUE |
| 5 | Mean | 3371.33 | 5772.18 | 33.48 | 6.92 | 4 | TRUE |
| 6 | Mean | 3388.12 | 5681.46 | 33.24 | 6.21 | 6 | TRUE |
| 7 | Mean | 3416.70 | 5697.26 | 33.10 | 6.00 | 7 | TRUE |
| 8 | Mean | 3456.47 | 5719.96 | 33.15 | 5.85 | 8 | TRUE |
| 9 | Mean | 3474.73 | 5743.77 | 33.10 | 5.72 | 9 | TRUE |
| 10 | Mean | 3483.40 | 5925.34 | 34.63 | 7.18 | 3 | FALSE |

1-10 of 10 rows

# Filter the result dataframe. Find the best 10 MAPE.

- We now demonstrate how to filter the dataframe to find only the Mean results, and arrange the dataframe in ascending order to show the lowest MAPE. We use the head function to show only the best 10 MAPE results. k = 2 is the winner.
- It's interesting to note that some of the top models are not the same based on MAE vs MAPE.

```r
367 ▾ ```{r}
368   # pull the top 10 based on MAPE
369
370   cv_grid_results %>%
371     filter(measure == 'Mean') %>%
372     arrange(MAPE) %>%
373     head(10)
374 ▴ ```
```

Description: df [10 × 7]

| | measure <chr> | MAE <dbl> | RMSE <dbl> | MAPE <dbl> | RMSPE <dbl> | param_K <int> | param_I <lgl> |
|---|---|---|---|---|---|---|---|
| 1 | Mean | 3261.69 | 6012.29 | 32.22 | 8.18 | 2 | TRUE |
| 2 | Mean | 3484.70 | 5757.88 | 32.80 | 5.57 | 10 | TRUE |
| 3 | Mean | 3498.20 | 5764.89 | 32.82 | 5.49 | 11 | TRUE |
| 4 | Mean | 3516.06 | 5795.59 | 33.04 | 5.45 | 12 | TRUE |
| 5 | Mean | 3416.70 | 5697.26 | 33.10 | 6.00 | 7 | TRUE |
| 6 | Mean | 3474.73 | 5743.77 | 33.10 | 5.72 | 9 | TRUE |
| 7 | Mean | 3325.64 | 5850.02 | 33.11 | 7.40 | 3 | TRUE |
| 8 | Mean | 3456.47 | 5719.96 | 33.15 | 5.85 | 8 | TRUE |
| 9 | Mean | 3366.06 | 5684.82 | 33.17 | 6.46 | 5 | TRUE |
| 10 | Mean | 3388.12 | 5681.46 | 33.24 | 6.21 | 6 | TRUE |

1-10 of 10 rows

# Ablation Analysis

- Do we see performance drop or improve when we remove a single column from the predictors?
- Here we remove each column one at a time and use cross validation to examine the results.
- We run a baseline first, then remove age, then sex, then bmi, then children then smoker and finally region.
- Based on MAE removing Region appears to lower the MAE, however it increases the MAPE.