

[\(/\)](#)[Fork\(1\) \(/fork/pUw773\)](#) [download \(/plain/pUw773\)](#)[copy](#)

```
1. # Kaggle Walmart recruiting competition 2014-02-20 to 2014-05-05.
2.
3. WORKING_DIRECTORY = "~/walmart"
4. options(stringsAsFactors = FALSE)
5. setwd(WORKING_DIRECTORY)
6.
7. library(Hmisc) # Hmisc is first so its summarize function does not mask plyr's
8. library(plyr)
9. library(testthat)
10. library(lubridate)
11. library(stringr)
12.
13. trend_sales <- function(v_sales, v_id, v_dt, id_num, trend_fctr) {
14.
15.   # Apply a trend factor to the historical sales, moving them to the
16.   # beginning of the test period.
17.   #
18.   # Args:
19.   #   v_sales: Vector of all sales in the test set.
20.   #   v_id: Vector of all store or department ids in the test set.
21.   #   v_dt: Vector of all dates in the test set.
22.   #   id_num: The store or department for which sales are to be trended.
23.   #   trend_fctr: The historical (not prospective) trend factor.
24.   #
25.   # Returns:
26.   #   The revised sales vector with trend applied to the
27.   #   components corresponding to id_num.
28.
29.   ind <- which(v_id == id_num)
30.   wks_between <- as.integer(difftime (http://www.opengroup.org/onlinepubs/009695
3199/functions/difftime.html) (v_dt[ind], min(v_dt[ind]), units="weeks"))
31.   fctr <- trend_fctr^(1/52 * (52 - wks_between))
32.   v_sales[ind] <- round(v_sales[ind] * fctr, 2)
33.   return(v_sales)
34. }
35.
36. blend_weeks <- function(next_yr_dt, coef1 = NULL, coef2 = NULL) {
37.
38.   # Given a date from the test set, the week ending on the corresponding date
39.   # in the training set will usually straddle two training weeks. This function
40.   # calculates an appropriate weighted average of the train weeks for
41.   # predicting the test week.
42.   #
43.   # Args:
44.   #   next_yr_dt: An end of week date (must be a Friday) from the test set
45.   #   coef1, coef2: Specify the weights rather than calculating them. Not used.
46.   #
47.   # Returns:
48.   #   A data frame with the test set id and predicted sales for next_yr_dt.
49.   #
50.   # Note:
51.   #   Dataframes test and train are used globally and are referenced within the
52.   #   blend_weeks function, although not passed as arguments.
53.
54.   stopifnot(wday(next_yr_dt) == 6) # End of week must be a Friday.
55.   dt <- next_yr_dt - years(1)
56.   stopifnot(wday(dt) != 6)
```

```

57.   days_to_friday <- (13 - wday(dt)) %% 7
58.   next_friday <- dt + days(days_to_friday)
59.   prev_friday <- next_friday - days(7)
60.   stopifnot(wday(next_friday) == 6)
61.   stopifnot(wday(prev_friday) == 6)
62.
63.   df1 <- subset(train, dt == next_friday)
64.   df2 <- subset(train, dt == prev_friday)
65.   df_valid <- subset(test, dt == next_yr_dt)[, c("Store", "Dept")]
66.
67.   df_both <- merge(df1[, 1:4], df2[, 1:4], by = c("Store", "Dept"),
68.     all = TRUE)
69.   df_both <- merge(df_valid, df_both, by = c("Store", "Dept"), all.x = T)
70.   df_both[, c("sales.x", "sales.y")] <-
71.     Hmisc::impute(df_both[, c("sales.x", "sales.y")], 0)
72.
73.   if(is.null(coef1)) coef1 <- 1 - days_to_friday/7
74.   if(is.null(coef2)) coef2 <- days_to_friday/7
75.   blended_sales <- round(with(df_both, coef1 * sales.x +
76.     coef2 * sales.y), 0)
77.   Id <- with(df_both, paste(Store, Dept, next_yr_dt, sep = "_"))
78.   df_ans <- data.frame(Id = Id, sales = blended_sales)
79.   return(df_ans)
80. }
81.
82. # Read and validate the data -----
83. train <- readRDS("train.rds") # Training data covers 2010-02-05 to 2012-11-01
84. test <- readRDS("test.rds") # Test data covers 2012-11-02 to 2013-07-26
85. expect_equal(nrow(train), 421570)
86. expect_equal(nrow(test), 115064)
87. expect_equal(with(train, length(unique(paste(Store, Dept, Date)))), nrow(train))
88. expect_equal(with(test, length(unique(paste(Store, Dept, Date)))), nrow(test))
89.
90. # Create derived variables -----
91. train <- mutate(train, dt = ymd(Date), yr = year(dt), wk = week(dt))
92. train <- rename(http://www.opengroup.org/onlinepubs/009695399/functions/rename.
93.   html)(train, replace = c("Weekly_Sales" = "sales"))
94. test <- mutate(test, dt = ymd(Date), yr = year(dt), wk = week(dt),
95.   prior_yr = yr - 1)
96.
97. # Map weeks of test period to corresponding weeks in train period -----
98. # Week Mapping Adjustments:
99. # Thanksgiving 2012 is in week 47, Thanksgiving 2011 in week 48,
100. # thus 47 is replaced with 48 and 48 is replaced by 49.
101. #
102. # Easter 2013 is on March 31 (week 13).
103. # Model week after Easter (14) by week after Easter (15).
104. # For Easter week wound up just doing the same blending as for other weeks.
105. test$wk <- plyr::mapvalues(test$wk, from = c(47, 48, 14), to = c(48, 49, 15))
106.
107. # Make initial predictions -----
108. # Construct the initial test set predictions (just a merge with train, lagging
109. # the test set by one year).
110. ans <- merge(test, train, by.x = c("Store", "Dept", "prior_yr", "wk"),
111.   by.y = c("Store", "Dept", "yr", "wk"), all.x = TRUE)
112. ans$sales[is.na(ans$sales)] <- 0
113. ans <- ans[, c("Store", "Dept", "Date.x", "sales")]
114. ans$Id <- with(ans, paste(Store, Dept, Date.x, sep = "_"))
115.
116. # Week blending adjustments -----
117. # Remove records in the test set that will be replaced by records derived
118. # from blending.


```

```

118. UNBLENDED_DATES <- c("2012-11-23", "2012-11-30", "2013-04-05")
119. BLEND_DATES <- setdiff(as.character(ymd("2012-11-02") + weeks(0:38)),
120.                          UNBLENDED_DATES)
121. ans <- subset(ans, !(Date.x %in% BLEND_DATES))
122. sub <- ans[, c("Id", "sales")]
123.
124. # Calculate the blended weeks and add them back to sub using plyr::rbind.fill.
125. blended_weeks <- plyr::rbind.fill(lapply(ymd(BLEND_DATES), blend_weeks))
126. sub <- rbind(sub, blended_weeks)
127.
128. # Reconstruct date, store, and department from the submission -----
129. # (awkward - could be cleaned up)
130. dt <- ymd(str_extract(sub$Id, ".{10}$" ))
131. store <- str_extract(sub$Id, "[0-9]+")
132. dept <- substr(str_extract(sub$Id, "[0-9]+"), 2, 3)
133.
134. # Make the trend adjustments (geometric mean of quarters). -----
135. store_trend_data <- list(c(1, 1.01), c(2, 1.01), c(3, 1.07), c(4, 1.02),
136.                          c(5, 1.05), c(6, 1.01), c(7, 1.03), c(8, 1.00),
137.                          c(9, 1.01), c(10, 0.97), c(11, 1.00), c(12, 0.99),
138.                          c(13, 1.01), c(14, 0.85), c(15, 0.95), c(16, 0.99),
139.                          c(17, 1.04), c(18, 1.03), c(19, 0.96), c(20, 0.99),
140.                          c(21, 0.90), c(22, 0.97), c(23, 1), c(24, 0.99),
141.                          c(25, 1.00), c(26, 1.00), c(27, 0.94), c(28, 0.95),
142.                          c(29, 0.98), c(30, 1.01), c(31, 0.96), c(32, 0.99),
143.                          c(33, 1.04), c(34, 1.01), c(35, 1.00), c(36, 0.80),
144.                          c(37, 0.97), c(38, 1.10), c(39, 1.07), c(40, 0.99),
145.                          c(41, 1.04), c(42, 1.00), c(43, 0.97), c(44, 1.08),
146.                          c(45, 0.97))
147. for(v in store_trend_data) {
148.   sub$sales <- trend_sales(sub$sales, store, dt, v[1], v[2])
149. }
150.
151. dept_trend_data <- list(c(1, 0.96), c(2, 0.98), c(3, 1.01), c(4, 1),
152.                          c(5, 0.91), c(6, 0.79), c(7, 0.99), c(8, 0.99),
153.                          c(9, 1.03), c(10, 0.99), c(11, 0.98), c(12, 0.98),
154.                          c(13, 0.98), c(14, 1.02), c(16, 0.95), c(17, 0.97),
155.                          c(18, 0.87), c(19, 1.06), c(20, 0.98), c(21, 0.94),
156.                          c(22, 1.01), c(23, 1.02), c(24, 1), c(25, 0.96),
157.                          c(26, 0.96), c(27, 1.02), c(28, 0.89), c(29, 1.02),
158.                          c(30, 0.92), c(31, 0.9), c(32, 0.97), c(33, 0.99),
159.                          c(34, 1.02), c(35, 0.92), c(36, 0.79), c(37, 0.97),
160.                          c(38, 0.98), c(40, 1.01), c(41, 0.94), c(42, 1.01),
161.                          c(44, 1.02), c(45, 0.53), c(46, 0.99), c(48, 1.96),
162.                          c(49, 0.96), c(50, 0.97), c(52, 0.93), c(54, 0.54),
163.                          c(55, 0.83), c(56, 0.93), c(58, 1.13), c(59, 0.7),
164.                          c(60, 1.02), c(65, 1.09), c(67, 1.02), c(71, 0.98),
165.                          c(72, 0.96), c(74, 0.97), c(79, 0.98), c(80, 0.96),
166.                          c(81, 0.98), c(82, 1.02), c(83, 1.01), c(85, 0.9),
167.                          c(87, 1.14), c(90, 0.98), c(91, 0.98), c(92, 1.04),
168.                          c(93, 1.02), c(94, 0.96), c(95, 0.99), c(96, 1.04),
169.                          c(97, 0.97), c(98, 0.95), c(99, 1.19))
170. for(v in dept_trend_data) {
171.   sub$sales <- trend_sales(sub$sales, dept, dt, v[1], v[2])
172. }
173.
174. # Save the submission -----
175. sub <- sub[, c("Id", "sales")]
176. names(sub) <- c("Id", "Weekly_Sales")
177. sub <- arrange(sub, Id)
178. expect_equal(nrow(sub), 115064)
179. z <- gzfile("submission.csv.gz")

```

<http://ideone.com/pUw773>

language: **R**
created: 1 hour ago
visibility:  public (/faq#visibility-of-a-code)

 Share or Embed source code

```
<script src="http://ideone.com/e.js/pUw773"
type="text/javascript"></script>
```

```
180. write.csv(sub, z, row.names = FALSE)
```


Runtime error #stdin #stdout #stderr 0.3s 22832KB

 comments (0)

 stdin

copy

Standard input is empty

 stdout

Standard output is empty

stderr

Error in setwd(WORKING_DIRECTORY) : cannot change working directory

Execution halted

Looking for Job?

14, 0.02s, 0.0 is

Sphere Research Labs (<http://sphere-research.com>). Ideone is powered by Sphere Engine™ (<http://sphere-engine.com/>)
[home \(/\)](#) [terms of use \(/terms\)](#) [api \(/sphere-engine\)](#) [language](#) [faq \(/faq\)](#) [credits \(/credits\)](#) [feedback & bugs \(/ideone/tools/bug/form/1/link/puw773/compiler/117\)](#) [desktop](#) [mobile \(/switch/mobile/l3bvdzc3mw==\)](#)