

ReadMe

Coleen Smith

4/2/2019

Contents

Assignment	2
Overview of Analysis	2
Part 1:	2
Part 2:	2
Part 3:	3
Part 4:	3
Part 5:	3
Analysis	3
UCI HAR Data Set	3
1. Merge the training & the test sets to create one data set	5
A.1 Read test files into a list of tables	5
A.2 Combine test tables into a single table	7
A.3 Read train files into a list of tables	7
A.4 Combine train tables into a single table	8
B For each test_merge and train_merge table, add a new column that records the original file name.	9
C. Combine test_add and train_add tables into a single table	9
Tidy data	10
2. Extract only the measurements on the mean and standard deviation for each measurement.	11
3. Use descriptive activity names to name the activities in the data set	11
Check for missing values	11
Add new descriptive variable activityName	12
4. Appropriately label the data set with descriptive variable names.	12
How descriptive?	12
Standardize variable names from features.txt	12
5. From the data set in step 4, create a second, independent tidy data set with the average of each variable for each activity and each subject.	13
Not So Tidy - Do measured features represent more than one value?	14
Tidy Data - Measured features represent a single value	14
Summary: A tidy data set	14
Output file: final_tidyframe.csv	15
Resources	15
Afterword: Tidy Data, Tidy Pantry	15

Assignment

This assignment will create one R script called `run_analysis.R` that does the following.

1. Merges the training and the test sets to create one data set.
2. Extracts only the measurements on the mean and standard deviation for each measurement.
3. Uses descriptive activity names to name the activities in the data set.
4. Appropriately labels the data set with descriptive variable names.
5. From the data set in step 4, creates a second, independent tidy data set with the average of each variable for each activity.

The assignment must also meet the following criteria:

1. The submitted data set is tidy.
2. The Github repo contains the required scripts.
3. Github contains a code book that modifies and updates the available codebooks with the data to indicate which variables were used in the analysis.
4. The README that explains the analysis files is clear and understandable.
5. The work submitted for this project is the work of the student who submitted it.

Overview of Analysis

The assignment included data set from the Human Activity Recognition Using Smartphones Data Set was “built from the recordings of 30 subjects performing activities of daily living (ADL) while carrying a waist-mounted smartphone with embedded inertial sensors.” (<http://archive.ics.uci.edu/ml/datasets/Human+Activity+Recognition+Using+Smartphones>)

The files were downloaded and explored to review the type of data presented and the organization. I used Atom to determine line counts, and the occurrences of key words: mean, std, activity numbers, subject numbers, and features. As David Hood mentioned, based on the data dimensions, the merging the files could be seen a fitting “Lego bricks” together.

This section is just an summary. For specifics on the analysis and code, please see the appropriate, labeled sections. In general, text files were read into data frames with the same name. `X_test.txt` was read into `X_test` data frame.

Part 1:

I applied the Lego brick methodology, using `cbind()` to create the `test_merged` and `train_merged` data frames. The three smaller data files in the test folder (`subject_test`, `X_test`, `Y_test`) were combined before the three larger versions in the train folder. If something went wrong, it would be easier to trace the problem with the smaller test data set. Then, the two were combined to create the `test_train` data frame of 10299 rows. Although not required in this portion of the assignment, I added variable names to the data frames before `cbind()`. This was a useful way to visually check to make sure the columns were bound in the correct order. As recommended by Wickham, when combining observations from multiple sources, I added a fixed variable `original_file` to distinguish test data from train data.

Part 1 final output: tidyframe

Part 2:

Since `test_train` included column names, I could use `grep` to extract the columns with “mean” and “std” in their names. I did not include those with “Mean” since they were input to angle features rather than an actual mean.

Part 2 final output: `mean_std_tidyframe`

Part 3:

My background is with relational databases, so that influences the way I think about data structures. I created an `activity_label` data frame and used `inner_join()` to add the activity name to the `mean_std_tidyframe`. I checked to ensure that there were no missing values before the join. Since activity number was not required later, it was omitted from final output.

Part 3 final output: `activity_tidyframe`

Part 4:

By the end of Part 3, `activity_tidyframe` already included descriptive variable names, but there was room to improve. I removed the parentheses and replaced the dashes with underscores. I also removed the leading digits from the mean and std variable names. Not all the original variable names in `features.txt` were unique, so I retained the original feature number with the feature name to make sure that the variable names would be unique in Part 1. None of these were required in Part 2, so the digits could be removed from the `activity_tidyframe` variable names. Part 4 final output: `descriptive_tidyframe`.

Part 5:

To prepare for Part 5, I tried to make sure that the output from each part of the assignment was tidy. My core assumption was that each of the variables listed in `features.txt` was a discrete measurement. The question remained: was a variable like `tBodyAcc_mean_X` a single observed measurement or could it be segmented to create a more narrow and “tidy-er” data set? After playing with the various permutations, segmenting the variables was possible, but introduced NAs into the data. Introducing NAs to a data set that had recorded 561 measurements for each of six activities for every subject seemed inappropriate. Without better knowledge of the 561 variables measured, I did not feel qualified to create a data set that might highlight “missing” data.

Based on my core assumption, `activity_tidyframe` was tidy data. To average each measured variable for each combination of subject and activity, I used `select` (to exclude the source of the original file), `group_by` to organize the data and `summarize_all(mean)` to average each measured variable. Part 5 final output: `final_tidyframe`

Assignment File Submitted: `final_tidyframe.csv` To read the file: `my_final_tidyframe <- read.csv(“final_tidyframe.csv”, header = TRUE, sep = “,”)`

Analysis

UCI HAR Data Set

Data collected from the accelerometers from the Samsung Galaxy S smartphone was made available to download. From: Davide Anguita, Alessandro Ghio, Luca Oneto, Xavier Parra and Jorge L. Reyes-Ortiz. A Public Domain Dataset for Human Activity Recognition Using Smartphones. 21th European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning, ESANN 2013. Bruges, Belgium 24-26 April 2013.

Full description [Link to data](#)

After downloading and unzipping files, I renamed the directory from “UCI HAR Dataset” to “UCI_HAR_Dataset” to avoid possible issues with spaces in a pathname.

```
dir("UCI_HAR_Dataset")
```

```
## [1] "activity_labels.txt" "features_info.txt" "features.txt"
## [4] "README.txt"          "test"          "train"
```

Parent Directory Files

- activity_labels.txt - six activities, two columns (number & name), all uppercase, no headers
- features.txt - 561 variables, two columns (number & name), lowerCamelCase, no headers
 - 46 occurrences of “mean”
 - 10 occurrences of “Mean”, all as part of an ordered pair, in 7 rows (555-561)
 - 33 occurrences of “std”
- features_info.txt - information on the variables and calculations listed in features.txt
- README.txt - overview of project
- test (folder) - data from first 30 subjects
- train (folder)- data from remaining 70 subjects

Based on the information included in features_info.txt, the files in the test and train and train folders are formatted and organized the same way. Since the test folder files are smaller with data from just 30 subjects, they became the starting point to understand the data and prepare the files for merging.

Test Folder Files

```
dir("UCI_HAR_Dataset/test")
```

```
## [1] "Inertial Signals" "subject_test.txt" "X_test.txt"
## [4] "y_test.txt"
```

To get a sense of how the data were arranged, I also looked at each of the files using atom.

test (folder)

- inertial signals (folder)
- each .txt - a set of 35 element vectors
- subject_test.txt - a single column of numbers, no column headers
- x_test.txt - rows of 561 element vectors, no column headers
- y_test.txt - a single column of numbers, no column headers
- “1” 496 occurrences
- “2” 471 occurrences
- “3” 420 occurrences
- “4” 491 occurrences
- “5” 532 occurrences
- “6” 537 occurrences
- Total = 2947

Train Folder Files

```
dir("UCI_HAR_Dataset/train")
```

```
## [1] "Inertial Signals" "subject_train.txt" "X_train.txt"
## [4] "y_train.txt"
```

Train (folder) mirrors the file structure and organization seen in test files, only larger because files represent 70 subjects rather than 30.

1. Merge the training & the test sets to create one data set

Based on this assignment's requirements, I decided to follow Wickham approach to tidying data when there is one type of data in multiple tables (Ref: Wickham, "Tidy Data," Section 3.5).

- A. Read the files into a list of tables.
- B. For each table, add a new column that records the original file name.
- C. Combine all tables into a single table.

In this case, test and train data are one type of data observed for separate groups of people. Thinking ahead to tidy data,

1. Each variable forms a column.
 - subject (100)
 - activity (6)
 - features (561)
2. Each observation forms a row.
 - An observation is of a given subject
 - performing a specific activity
 - represented by 561 measured features
3. Each type of observational unit forms a table.
 - Interim: test and train in separate tables
 - Final: test and train data in single table

Since the subject, activity and measurements are recorded in separate tables within test and train, the first step would be to create a test data file with subject, activity and measurement observations. Then, do the same for the train data. Since they have the same variables, combine the interim test and train data files using `rbind()`.

A.1 Read test files into a list of tables

The number of data lines in the source file should equal the number of lines in the destination file. [Back in the dark ages when file transfer protocol included a magnetic tape and FedEx Overnight, I spent a day looking for 13 "lost" lines out of 25,000 - lesson burned.]

```
## library(R.utils)
countLines("UCI_HAR_Dataset/features.txt")
```

```
## [1] 561
## attr(,"lastLineHasNewline")
## [1] TRUE
```

```
countLines("UCI_HAR_Dataset/test/subject_test.txt")
```

```
## [1] 2947
## attr(,"lastLineHasNewline")
## [1] TRUE
```

```
countLines("UCI_HAR_Dataset/test/X_test.txt")
```

```
## [1] 2947
## attr("lastLineHasNewline")
## [1] TRUE
```

```
countLines("UCI_HAR_Dataset/test/Y_test.txt")
```

```
## [1] 2947
## attr("lastLineHasNewline")
## [1] TRUE
```

According to README.txt and features_info.txt, features.txt represents the names of the 561 measured variables. Since the 561 features variable list already contains the target mean & std calculations to be extracted in Part 2 of the assignment, the merge for Part 1 will not include the additional detail observations contained in the inertial signals folders for either test or train data.

Using read.table, create a data frame from each file: subject_test.txt, X_test.txt and Y_test.txt.

```
subject_test <- read.table("UCI_HAR_Dataset/test/subject_test.txt", header = FALSE)
str(subject_test)
```

```
## 'data.frame': 2947 obs. of 1 variable:
## $ V1: int 2 2 2 2 2 2 2 2 2 2 ...
```

```
X_test <- read.table("UCI_HAR_Dataset/test/X_test.txt", header = FALSE)
## This is a loooooong output, just need highlights
str(X_test, list.len = 5)
```

```
## 'data.frame': 2947 obs. of 561 variables:
## $ V1 : num 0.257 0.286 0.275 0.27 0.275 ...
## $ V2 : num -0.0233 -0.0132 -0.0261 -0.0326 -0.0278 ...
## $ V3 : num -0.0147 -0.1191 -0.1182 -0.1175 -0.1295 ...
## $ V4 : num -0.938 -0.975 -0.994 -0.995 -0.994 ...
## $ V5 : num -0.92 -0.967 -0.97 -0.973 -0.967 ...
## [list output truncated]
```

```
Y_test <- read.table("UCI_HAR_Dataset/test/Y_test.txt", header = FALSE)
str(Y_test)
```

```
## 'data.frame': 2947 obs. of 1 variable:
## $ V1: int 5 5 5 5 5 5 5 5 5 5 ...
```

Create a character vector from features.txt to use later in naming the Y_test data frame columns. Note, the features.txt file has some duplicate variable names, differing only by row/line number index. For example, “fBodyAcc-bandsEnergy()-1,8” appears at row 303, 317, and 331. There are three groups of 14 variables that start with “fBodyAcc-bandsEnergy()” and end with an ordered pair in rows 303-344. Use readLines to keep line number index with variable name.

```
features <- readLines("UCI_HAR_Dataset/features.txt")
str(features)
```

```
## chr [1:561] "1 tBodyAcc-mean()-X" "2 tBodyAcc-mean()-Y" ...
```

```
print (features[303:316])
```

```
## [1] "303 fBodyAcc-bandsEnergy()-1,8" "304 fBodyAcc-bandsEnergy()-9,16"
## [3] "305 fBodyAcc-bandsEnergy()-17,24" "306 fBodyAcc-bandsEnergy()-25,32"
## [5] "307 fBodyAcc-bandsEnergy()-33,40" "308 fBodyAcc-bandsEnergy()-41,48"
## [7] "309 fBodyAcc-bandsEnergy()-49,56" "310 fBodyAcc-bandsEnergy()-57,64"
## [9] "311 fBodyAcc-bandsEnergy()-1,16" "312 fBodyAcc-bandsEnergy()-17,32"
## [11] "313 fBodyAcc-bandsEnergy()-33,48" "314 fBodyAcc-bandsEnergy()-49,64"
```

```
## [13] "315 fBodyAcc-bandsEnergy()-1,24" "316 fBodyAcc-bandsEnergy()-25,48"
```

For readability, set variable names before combining tables. Use conventions described in README.txt and features_info.txt files. This will make it easier later to filter “mean” and “std” variables.

```
names(subject_test) <- "subject"
names(Y_test) <- "activityNumber"

## features.txt contains two columns feature number and feature name
names(X_test) <- features
```

With the exception of adding variable names, no data have been transformed or filtered. This is useful in case something goes wrong downstream and I need to reset to original values. Each file (subject_test.txt, X_test.txt and Y_test.txt) has 2947 rows, and represent different measured variables of a single observation.

A.2 Combine test tables into a single table

Each table represents a unique subset of variables, and are linked by the order of the rows rather than a specific, shared variable. Use cbind() rather than one of the dplyr join() commands to merge the three train data files. . Bind in order of subject (subject_test), activity (Y_test), then recorded measurements (X_test). Review the result.

```
test_merged <- cbind(subject_test, Y_test, X_test)
str(test_merged, list.len = 5)

## 'data.frame':    2947 obs. of  563 variables:
## $ subject                : int  2 2 2 2 2 2 2 2 2 2 ...
## $ activityNumber          : int  5 5 5 5 5 5 5 5 5 5 ...
## $ 1 tBodyAcc-mean()-X     : num  0.257 0.286 0.275 0.27 0.275 ...
## $ 2 tBodyAcc-mean()-Y     : num  -0.0233 -0.0132 -0.0261 -0.0326 -0.0278 ...
## $ 3 tBodyAcc-mean()-Z     : num  -0.0147 -0.1191 -0.1182 -0.1175 -0.1295 ...
## [list output truncated]
```

With test_merged’s 2947 observations, no rows lost. 563 column variables represent subject and activityNumber columns added to 561 feature measurements.

A.3 Read train files into a list of tables

As mentioned before, the test and train files have identical structure and organization. Although larger, they can be processed in the same way: count lines, name the columns, create a data frame for each. There are probably many ways to count the lines in a text file using R. The countLines() command is found in the R.utils package.

```
## library(R.utils)
countLines("UCI_HAR_Dataset/train/subject_train.txt")

## [1] 7352
## attr(,"lastLineHasNewline")
## [1] TRUE

countLines("UCI_HAR_Dataset/train/X_train.txt")

## [1] 7352
## attr(,"lastLineHasNewline")
## [1] TRUE
```

```
countLines("UCI_HAR_Dataset/train/Y_train.txt")
```

```
## [1] 7352
## attr(,"lastLineHasNewline")
## [1] TRUE
```

Each train data files has 7352 lines.

Using read.table, create a data frame from each file: subject_train.txt, X_train.txt and Y_train.txt.

```
subject_train <- read.table("UCI_HAR_Dataset/train/subject_train.txt", header = FALSE)
str(subject_train)
```

```
## 'data.frame':    7352 obs. of  1 variable:
## $ V1: int  1 1 1 1 1 1 1 1 1 1 ...
```

```
X_train <- read.table("UCI_HAR_Dataset/train/X_train.txt", header = FALSE)
## This is a loooooong output, just need highlights
str(X_train, list.len = 5)
```

```
## 'data.frame':    7352 obs. of  561 variables:
## $ V1 : num  0.289 0.278 0.28 0.279 0.277 ...
## $ V2 : num -0.0203 -0.0164 -0.0195 -0.0262 -0.0166 ...
## $ V3 : num -0.133 -0.124 -0.113 -0.123 -0.115 ...
## $ V4 : num -0.995 -0.998 -0.995 -0.996 -0.998 ...
## $ V5 : num -0.983 -0.975 -0.967 -0.983 -0.981 ...
## [list output truncated]
```

```
Y_train <- read.table("UCI_HAR_Dataset/train/Y_train.txt", header = FALSE)
str(Y_train)
```

```
## 'data.frame':    7352 obs. of  1 variable:
## $ V1: int  5 5 5 5 5 5 5 5 5 5 ...
```

For readability, set column names to match for subject_test, X_test, Y_test. Since X_test and X_train have identical structure, use the features vector created earlier to name columns.

```
names(subject_train) <- "subject"

names(Y_train) <- "activityNumber"

## features.txt contains two columns feature number and feature name
names(X_train) <- features
```

With the exception of adding column names, no data has been transformed or filtered. This is useful in case something goes wrong downstream and need to reset to original values. Each file (subject_train.txt, X_train.txt and Y_train.txt) has 7352 rows, and represent different measured variables of a single observation.

A.4 Combine train tables into a single table

Like the test data, each train table represents a unique subset of variables, and are linked by the order of the rows rather than a specific, shared variable. Use cbind() rather than one of the dplyr join() commands to merge the three train data files. Bind in order of subject (subject_train), activityNumber (Y_train), then recorded measurements (X_train). Review the result.

```
train_merged <- cbind(subject_train, Y_train, X_train)
str(train_merged, list.len = 5)
```



```
## 'data.frame': 7352 obs. of 563 variables:
## $ subject : int 1 1 1 1 1 1 1 1 1 1 ...
## $ activityNumber : int 5 5 5 5 5 5 5 5 5 5 ...
## $ 1 tBodyAcc-mean()-X : num 0.289 0.278 0.28 0.279 0.277 ...
## $ 2 tBodyAcc-mean()-Y : num -0.0203 -0.0164 -0.0195 -0.0262 -0.0166 ...
## $ 3 tBodyAcc-mean()-Z : num -0.133 -0.124 -0.113 -0.123 -0.115 ...
## [list output truncated]
```

With `train_merged`'s 7352 observations, no rows lost. 563 variables represent subject and activityNumber columns added to 561 feature measurements.

B For each `test_merge` and `train_merge` table, add a new column that records the original file name.

Install and load dplyr packages. Use `mutate()` to create a new variable to identify the original file (test or train).

```
test_add <- mutate(test_merged, original_file = "test")
```

```
## sanity check ... expecting 564
ncol(test_add)
```

```
## [1] 564
```

```
tail(test_add$original_file)
```

```
## [1] "test" "test" "test" "test" "test" "test"
```

```
train_add <- mutate(train_merged, original_file = "train")
```

```
## sanity check ... expecting 564
ncol(test_add)
```

```
## [1] 564
```

```
tail(train_add$original_file)
```

```
## [1] "train" "train" "train" "train" "train" "train"
```

Using `ncol` confirms that each new table now has 564 columns. Explicitly specifying new column by name in `tail` operation, demonstrates `original_name` was added to `test_merged` and `train_merged` and was filled.

C. Combine `test_add` and `train_add` tables into a single table

Since the `test_add` and `train_add` tables have the same column organization, use `rbind` to combine into a single table.

```
test_train <- rbind(test_add, train_add)
```

```
## 2947 from test + 7352 from train = 10299 in test_train
## matches the 10299 instances recorded in raw data
nrow(test_train)
```

```
## [1] 10299
```

Row binding `test_add` and `train_add` tables creates a single data frame with 10299 rows: 2947 from test plus 7352 from train. Total rows bound matches the 10299 instances recorded in raw data source's webpage (<http://archive.ics.uci.edu/ml/datasets/Human+Activity+Recognition+Using+Smartphones>).

Tidy data

The `test_train` data set is tidy, if not pretty.

Each variable forms a column.

The final `test_train` data frame includes 564 variables for subject, activity, origin file plus 561 measured features. While some of the features represent computed values such as mean, they are a single observed variable.

It might be tempting to split the columns by parsing some of the variable names. For example, the first 40 rows could be reshaped based on feature (`tBodyAcc`), calculation (`mean()`, `std()`, `min()`, `max()`, ...) and axis (`X`, `Y`, `Z`). With 561 measurements split into columns based on feature, calculation and axis, the data set would become very wide. Furthermore, not every variable includes an axis designation. In David Hood's blog, he writes, "[W]ithin the broad set theory background it becomes what is the best tidy form of data to answer a specific question." The next part of the assignment is to extract only the `mean()` and `std()` observations. The rows can be extracted using simple `grep` commands that do not reshaping or risk changing the data.

For this application the variables are tidy.

Each observation forms a row.

For this data set, an observation is of a given subject, performing a specific activity, described by 561 measured features, and labeled by origin (test or train).

Each type of observational unit forms a table.

The collection of variables, although from different samples, represent the same type of observations (subject, activity, measured features) and includes a variable that identifies the original source of the observation.

Tidy Data, plus

Based on Wickham's standards, the data frame `test_train` is tidy. Although not explicitly mentioned in the definition of tidy data, Wickham also discusses the order in which variables should appear in a data set:

"Fixed variables describe the experimental design and are known in advance. . . Measured variables are what we actually measure in the study. Fixed variables should come first, followed by measured variables, each ordered so that related variables are contiguous."

By this standard, subject, activity and original file would be fixed variables and should appear contiguously and first. Measured variables would be the items described by the feature measurements. This could be accomplished by moving the `original_file` column to the first column of the data frame.

```
## new dataframe with fixed variable first
tidyframe <- select(test_train, 564, 1:563)
```

Tidy data addresses the shape and organization of data that make it software readable rather than what makes it human readable: tidyframe is tidy data. Descriptive variable names and possible reshaping will be addressed later in the assignment (see sections 3 & 4).

2. Extract only the measurements on the mean and standard deviation for each measurement.

Based on features_info.txt, observed features representing mean or standard deviation include the string “mean” or “std” in their name: mean, meanFreq and std. Features that included “Mean” as part of an ordered pair were excluded since they are inputs to the angle measurements and not discrete entities.

The tidyframe data includes variable names, use grep to extract a list of columns containing the target strings.

```
## use grep to find target strings
mean_std_columns <- grep("mean()|std()", names(tidyframe), value = TRUE)
mean_std_tidyframe <- select(tidyframe, original_file, activityNumber, subject, mean_std_columns)
str(mean_std_tidyframe, list.len = 5)

## 'data.frame': 10299 obs. of 82 variables:
## $ original_file : chr "test" "test" "test" "test" ...
## $ activityNumber : int 5 5 5 5 5 5 5 5 5 5 ...
## $ subject : int 2 2 2 2 2 2 2 2 2 2 ...
## $ 1 tBodyAcc-mean()-X : num 0.257 0.286 0.275 0.27 0.275 ...
## $ 2 tBodyAcc-mean()-Y : num -0.0233 -0.0132 -0.0261 -0.0326 -0.0278 ...
## [list output truncated]
```

There were 79 total occurrences of the strings “mean” and “std” in the original features.txt file. Including the three fixed variables (original_file, activityNumber, subject) the final data frame correctly includes 82 variables.

3. Use descriptive activity names to name the activities in the data set

Six activity numbers and names are listed in activity_labels.txt. Read the activity_labels.txt file into a table to use with join to replace activity numbers with names in tidyframe.

```
## An activity id/number and name is listed in activity_labels.txt for each unique activity
activity_labels <- read.table("UCI_HAR_Dataset/activity_labels.txt", header = FALSE)

## name the columns
names(activity_labels) <- c("activityNumber", "activityName")
```

Check for missing values

Before doing any join based on a key, like activity, double check to make sure there will be a mapping for each. An outer join will substitute an NA for missed joins, but it would be better to know ahead of time about missing data.

```
uactivity <- unique(mean_std_tidyframe$activityNumber)
str(uactivity)

## int [1:6] 5 4 6 1 3 2
```

Add new descriptive variable activityName

From the above, there are 6 unique activities in tidyframe and mean_std_tidyframe, with values ranging from 1-6 and no NA. Add activity name to tidyframe via inner join, which is more descriptive than a number.

```
## apply names to mean_std_tidyframe using join
activity_tidyframe <- inner_join(mean_std_tidyframe, activity_labels, by = "activityNumber")

## move column to be contiguous with fixed variables
activity_tidyframe <- select (activity_tidyframe, original_file, activityName, subject, 4:82)

str(activity_tidyframe, list.len = 5)

## 'data.frame': 10299 obs. of 82 variables:
## $ original_file : chr "test" "test" "test" "test" ...
## $ activityName : Factor w/ 6 levels "LAYING","SITTING",...: 3 3 3 3 3 3 3 3 3 3 ...
## $ subject : int 2 2 2 2 2 2 2 2 2 2 ...
## $ 1 tBodyAcc-mean()-X : num 0.257 0.286 0.275 0.27 0.275 ...
## $ 2 tBodyAcc-mean()-Y : num -0.0233 -0.0132 -0.0261 -0.0326 -0.0278 ...
## [list output truncated]
```

The new data frame activity_tidyframe replaced a numeric activity number with a name, to better describe the activities in the data set.

4. Appropriately label the data set with descriptive variable names.

The column names for the subject, activity and measured variables were assigned earlier to the subject_, Y_, and X_ tables just before the three were merged to create the test_merged and train_merged data frames.

Naming the columns before the cbind() merges, allowed for better visual inspection of new data frames. Used str() to review the merge and make sure that columns were bound in the correct order. With subject and activityNumber both being numeric, it might be more challenging to scroll through the new data frames to make sure the columns had been bound in the correct order.

How descriptive?

There is always a trade off between how descriptive to make a variable name, and how practical it would be to use a very long name for data exploration and study. Whether the current labels are descriptive enough, depends on the end user. To one who extracts and studies this type of data, the names are very descriptive.

Assuming the end user of the final data set is knowledgeable about the field, the variable names listed in features.txt would be satisfactory. However, they would benefit from some standardization and fine tuning.

Standardize variable names from features.txt

The current version of activity_tidyframe has some quirks within the variable names:

- dashes and parentheses in the variable names
- leading digits leftover from earlier measured features variable name assignment

Remove dashes and parentheses

There are several approaches to creating variable names. The UCI HAR Data Set authors chose the lowerCamelCase style for variables and underscore_separated for text file names. Data frames and vectors used in run_analysis.R use underscore_separated.

Because special characters like commas, dashes and parentheses can create problems, use gsub to replace or remove them. While there are dashes and parentheses in the activity_tidyframe variable names, there are no commas.

```
## check for commas
grep(",", names(activity_tidyframe), value = TRUE)
```

```
## character(0)
```

```
## replace dashes w/ underscores and remove parentheses
names_dash_to_under <- gsub("-", "_", names(activity_tidyframe))
names_no_paren <- gsub("\\(\\)", "", names_dash_to_under)
```

Remove leading digits from measured variables

Earlier, I discovered there are three identical groups of 14 variables that start with “fBodyAcc-bandsEnergy()” and end with an ordered pair in lines 303 through 344 of features.txt. At the time, without changing the feature names, the only way to make these variable names unique was to include the row number. None of these variables are included in activity_tidyframe.

```
## Variables w/duplicate names followed the same pattern
grep("fBodyAcc-bandsEnergy()", names(activity_tidyframe), value = TRUE)
```

```
## character(0)
```

None of these variables are included with the mean() or std() calculated measurements extracted to create activity_tidyframe. The leading digits can be removed. Use gsub to remove leading digits and space leftover from Part 1. Update activity_tidyframe variable names.

```
names_no_digits <- gsub("^([0-9]{1,3}) ", "", names_no_paren)

## descriptive_tidyframe will be identical to activity_tidyframe except for variable names
descriptive_tidyframe <- activity_tidyframe
names(descriptive_tidyframe) <- names_no_digits
```

The measured variable names are now standardized, descriptive text for a end user familiar with the types of calculations used.

5. From the data set in step 4, create a second, independent tidy data set with the average of each variable for each activity and each subject.

The data set submitted will be wide form, with each column representing a fixed or measured variable.

Not So Tidy - Do measured features represent more than one value?

Looking at the data frame extracted in Part 4, the `mean()` and `std()` variables contain three segments: feature, calculation and reference. Furthermore, each variable starts with a “t” or “f.” For example with `tBodyAcc_mean_X`, `tBodyAcc` would be the feature variable, `mean` would be the calculation and `X` would be the reference. Summarized:

- features as measured variable (11): `fBodyAcc`, `fBodyAccJerk`, `fBodyBodyAccJerk`, `fBodyGyro`, `fBodyBodyGyro`, `fBodyBodyGyroJerk`, `tBodyAcc`, `tBodyAccJerk`, `tBodyGyro`, `tBodyGyroJerk`, `tGravityAcc`
- Calculation (3): `mean`, `meanFreq`, `std`
- Reference (4): X-axis, Y-axis, Z-Axis Mag

The data frame could be reshaped to allow feature variable grouping by subject and activity, as well as calculation and/or reference. Grouping by all four would yield summary data equivalent to the `final_tidyframe` above. For a look at one possible outcome, see `variable.xlsx` where I played with reshaping based on segmenting the variable names.

Not all possible combination of feature, calculation and reference are included in the original data. Introducing NAs creates the appearance the original researchers missed key calculations when perhaps that permutation of feature, calculation and frame was impossible, impractical or irrelevant.

That would not be tidy.

This is a good example of where I would spend time with the end user to truly understand the physics behind the data.

Tidy Data - Measured features represent a single value

The most direct approach for creating a second data set with the average of each activity and each subject would be to compute the column mean for each of the measured variables as listed in `features.txt`, grouped by subject and activity.

```
## final output, avg the means and std for each combination of subject and activity
## omit original_file (test or train)
final_tidyframe <- descriptive_tidyframe %>%
  select (2:82) %>%
  group_by(subject, activityName) %>%
  summarise_all(mean)
```

With thirty subjects and six activities, `final_tidyframe` should have 180 rows. The `original_file` column was not required, thus not included in `final_tidyframe`, resulting in 81 columns.

This approach is tidy and appropriate if the data set end user considers each of the 79 measured variables a single unit of data. Grouping this way by subject and then activity creates a data set that could be used as is, or to further summarize data by additional combinations of subject or activity.

Summary: A tidy data set

To meet the rubric, the data set may be narrow or wide, provided it is tidy. A tidy data set must meet all three of Wickham’s conditions. The `final_tidyframe` data set meets these conditions.

1. Each variable forms a column.
2. Each observation forms a row.
3. Each type of observational unit forms a table.

Wickham and Hood also discussed that “tidy-ness” is a function of how well the data set fits the task. The task was, from the the data set in step 4, to create a second independent tidy data set with the average of each variable for each activity and each subject. The data set `final_tidyframe` presents the average of each variable, in a tidy, rectangular format that could be further summarised by subject or activity.

Segmenting the feature variables without a solid knowledge of how the feature variable were computed and used, may not add to the value of the data set. There are likely alternate tidy arrangements of the data, but without a better understanding of the variable inputs, I chose not to reshape the data.

Output file: `final_tidyframe.csv`

```
## write to a csv for submission
write.table(final_tidyframe, "final_tidyframe.csv", row.names = FALSE, col.names = TRUE, sep = ",")

## to read file:
##my_final_tidyframe <- read.table("final_tidyframe.csv", header = TRUE, sep = ",")
```

Resources

In addition to lecture notes, several resources were very helpful in strategizing and working with the assignment.

Henry, Lionel, and Hadley Wickham. “Tidy Evaluation.” Tidy Evaluation, tidyeval.tidyverse.org/index.html.

Su, Xing, “Data Science Specialization Course Notes,” 16 Feb 2016, <https://sux13.github.io/DataScienceSpCourseNotes/>

Thoughtfulbloke. “Getting and Cleaning the Assignment.” Thoughtfulbloke Aka David Hood, 26 Jan. 2016, thoughtfulbloke.wordpress.com/2015/09/09/getting-and-cleaning-the-assignment/.

Wickham, Hadley, and Garrett Golemund. R For Data Science: Import, Tidy, Transform, Visualize and Model Data. O’Reilly, 2017.

Wickham, Hadley. Tidy Data. Journal of Statistical Software, 2014, vita.had.co.nz/papers/tidy-data.pdf.

Xie, Yihui, Allaire, J. J. and Golemund, Garrett, “RMarkdown: The Definitive Guide,” 19 Mar. 2019, <https://bookdown.org/yihui/rmarkdown/>

R Markdown Cheatsheet at <https://www.rstudio.com/wp-content/uploads/2016/03/rmarkdown-cheatsheet-2.0.pdf>

Afterword: Tidy Data, Tidy Pantry

As with creating descriptive variables, there is a trade off between using data as originally recorded and data as segmented. R offers powerful tools for reshaping data, but just because one can reshape the data, does not mean the reshaped version is meaningful. Reshaping data sets must be accompanied by a solid knowledge of the derivation of the data and its applications.

’Tis the season for spring cleaning and emerging critters. A good time to make sure all secure and tidy in the pantry. Consider an example, an observer’s top to bottom view of my pantry:

- shelf 1: square based, plastic bins, medium
- shelf 2: square based, plastic bins, large, medium, small

- shelf 3: glass jars and plastic bags, mixed inside square baskets
- shelf 4: canned goods, plastic bags
- floor: canned goods

It might be tempting to re-organize my pantry storage by storage type: bags, bins, baskets, cans, jars. It would be very tidy, but a total disaster for meal prep because it is not the container that is critical, but the contents and application.

A more valid recording of my pantry:

- shelf 1: sweeteners, chocolate, garnishes for desserts
- shelf 2: flours, rising agents for baking
- shelf 3: nuts, seeds, cereals and dried fruits for granola & trail mix
- shelf 4: beans, pasta, grains for side dishes
- floor: dog food

I definitely would not want the dog food next to the refried beans (cans), or the pumpkin seeds with the mung beans (glass jars)!