

Assignment 3

Practical Introduction to Data Science (SS1-SEM2)

Candidate Exam Number: B136325

Department of Informatics
University of Edinburgh
Scotland
9th June 2019

Contents

Question 1	1
Introduction	1
Business Understanding	1
Pids.wellbeing.weather	2
GitHub	2
RStudio	2
Structure	2
File and exported function names	2
Exported function names and questions	3
Function question sequence number	3
Function name sequence number	3
Summary	3
Data understanding	3
Statistical value chain	3
Raw data	4
Technically correct: weather .txt data	4
Technically correct: weather .dsv data	4
Technically correct: weather data frames	5
Technically correct: weather data frames complete	5
Technically correct: weather single data frame	6
Consistently correct	6
Exploratory data analysis (EDA)	7
Box and whisker	7
Remove outliers	7
Weather variables by latitude	7
Weather variables by latitude and longitude	8
Cluster tendency	8
VAT	8
Data preparation	9

K choice	9
Algorithm choice	9
Linkage choice	9
Analysis	10
Hierarchical	10
K-means	10
Evaluation	10
Hierarchical	10
Deployment	11
Question 2	12
Business Understanding	12
Data Understanding	12
Append latitude categories	12
Exploratory data analysis (EDA)	12
Number of weather stations per latitude category	12
Weather stations per latitude category	12
Mean weather features per latitude category	13
Latitude category pairwise comparisons	13
Latitude and longitude per latitude category	13
Data Preparation	13
Training and test data split	13
Algorithm choice	14
Analysis	14
K Nearest Neighbours	14
Question 2	15
Business Understanding	15
Data Understanding	15
Statistical value chain	15
Technically correct: well-being	15
Consistently correct: weather and well-being	15
Exploratory data analysis (EDA)	15
Data Preparation	16
Analysis	16
Evaluation	16
Appendices	18
Raw data output	19

Question 1

Introduction

The demonstration will be structured with regard to the CRISP-DM [1] data science project methodology, which describes such projects in terms of the following six sequential phases: Business Understanding, Data Understanding, Data Preparation, Analysis, Evaluation and Deployment. In addition, the CRISP-DM Data Understanding phase will contain the data cleaning components of the Statistical Value Chain, namely, Technically Correct and Consistently Correct data. Nevertheless, the first of the CRISP-DM phases, Business Understanding, will be presented immediately below.

Business understanding

The goal of this question will be to demonstrate that weather data can be clustered into two or more groups with each containing 'similar' weather conditions. The data contains results from thirty six monitoring stations across the UK (as made available online by the Met Office). While the process of retrieving and then cleaning the data will be described within the Data Understanding stage further below, it is worth noting that all of the process described within this answer have been implemented as an R package called **pids.wellbeing.weather**. The package's name begins with an abbreviation of the Practical Introduction to Data Science course, namely **pids**, followed by the names of the two primary datasets within the package: **wellbeing.weather**. Further information about the package will be provided immediately below.

pids.wellbeing.weather

GitHub

The **pids-wellbeing-weather** package has been made available on GitHub [2], using this student's exam number, **b136325**, as the account name. In addition, the package makes use of Packrat [3] dependency management. This means that the package contains all of its dependencies. This approach increases the package's size. Importantly, however, it ensures portability, and the package can be cloned (and used immediately) from GitHub with the following command.

```
git clone https://github.com/b136325/pids.wellbeing.weather.git
```

RStudio

Alternatively, the package can be installed into RStudio using the commands below.

```
install.packages("devtools")  
library(devtools)  
install_github("b136325/pids.wellbeing.weather")
```

Structure

The **pids-wellbeing-weather** package has been structured in accordance with R best practice [4]. For example, raw data can be found within the **./data-raw/** directory. In addition, both *technically* and *consistently correct* data can be found within the **./data/** directory. The code, which was linted using Linter [5], and which was semantically versioned [6] using Git [7], can be found within the **./R/** directory.

Files and exported function names

Broadly, each code file exports only one function, and the name of the exported function accords with the associated file name. For example, the first code file relating to question 1, which downloads raw data, and whose address can be found directly below, exports a function called *question_1_001_svc_raw-data*. In almost all cases, the names of exported functions match their associated files. This approach was adopted to ensure that exported functions could be found quickly within the code.

```
./R/question-1-001-svc-raw-data.R
```

Exported function names and questions

In addition, the name of each export function begins with the related question number. For example, the exported function ***question_1_001_svc_raw-data***, as described above, relates to question 1. In contrast, the exported function ***question_2_007_analysis_knn***, which provides K Nearest Neighbours analysis, relates to question 2. This approach has been adopted for all exported functions, and it provides a structure for the code in relation to the questions.

Function sequence number

The name of each exported function (and the name of each associated code file) also contains a three digit number. This is the **function sequence number**. It can be found immediately after the question number, and it represents the function's sequence of use in relation to a specific question. For example, the exported function *question_1_001_svc_raw-data*, as, described above, is the first function relating to question 1. In contrast, *question_1_036_eva_charts_hier_latitude_longitude_min_temp* is the 36th exported function relation to the same question.

Function names and question phases

Lastly, each exported function name contains a **question phase** description. This can be found immediately after the **function sequence number**, and it describes the phase (or part) of a question to which the function is related. For example, the exported function *question_1_023_analysis_hierarchical*, which returns a hierarchical cluster, is a part of the analysis phase for question 1. In contrast, *question_2_003_edata_charts_weather_features* is an exported function, which returns a chart of weather features, and which is related to the Exploratory Data Analysis (EDA) phase of question 2.

Summary

Data understanding

Statistical value chain

1. Package directory structure - raw to data
2. Within data (svc) sub-directories
3. sub directory naming
4. Command naming
5. single exported command per file

Raw data: weather and well-being

This process downloads the required files into the *pids-wellbeing-weather* and stores them within the *data-raw/*weather and the *data-raw/*wellbeing directories, respectively. The web address from which the data are downloaded are constructed dynamically using constants defined within ‘constants-weather-station-names.R’ and ‘constants-wellbeing’. The function associated with this process has been exported from the *pids.wellbeing.weather* package. When run it returns a list of files that have been downloaded successfully, along with any that have failed, as illustrated in Figure ???. The function can be run using the command below:

```
pids.wellbeing.weather::question_1_001_svc_raw_data()
```

Within the *pids.wellbeing.weather* package the function described above can be found at the following address:

```
./R/question-1-001-svc-raw-data.R
```

Technically correct: weather .txt data

In order to ensure a strong data related ‘separation of concerns’, such that, no amendments should be made to the data stored within the ‘data-raw’ directory, this process moves the raw weather data from ‘data-raw’ into ‘data’. In addition, it adds a closing new line character to each file, facilitating subsequent processing. More specifically, the weather data is moved into a directory representing the first stage of the ‘technically correct’ Statistical Value Chain process, which can be found at ‘data/weather/stage-010-technically-correct-text’. The function associated with this process has been exported from the package and it can be run using the command below. The function returns a list containing the destination paths of the successfully moved files.

```
pids.wellbeing.weather::question_1_002_svc_tech_weather_txt()
```

Within the *pids.wellbeing.weather* package the function described above can be found at the following address:

```
./R/question-1-002-svc-tech-weather-txt.R
```

Technically correct: weather .dsv data

The next stage of technically correct data processing transforms the .txt files into a .dsv white space delimited file format. This processes involves the removal of non column related header items from the .txt files. It also involves ensuring that the all of the files have a standard number of columns, which enables them to be

reliably converted into a delimited file format. In addition, invalid data items, such as "—" are converted into a common data standard, such as "NA". The function associated with this process has been exported from the package and it can be run using the command below. The function returns a list containing the destination paths of the successfully moved files, along with any errors.

```
pids.wellbeing.weather::question_1_003_svc_tech_weather_dsv()
```

Within the *pids.wellbeing.weather* package the function described above can be found at the following address:

```
./R/question-1-003-svc-tech-weather-dsv.R
```

Technically correct: weather data frames

The stage of technically correct data processing transforms the .dsv files into data frames, saved in .Rds formats within the data/weather/stage-012-technically-correct-dataframe. The function associated with this process has been exported from the package and it can be run using the command below. The function returns a list containing the destination paths of the successfully saved files, along with any errors.

```
pids.wellbeing.weather::question_1_004_svc_tech_weather_df()
```

Within the *pids.wellbeing.weather* package the function described above can be found at the following address:

```
./R/question-1-004-svc-tech-weather-df.R
```

Technically correct: weather data frames completet

The stage of technically correct data processing transforms the .Rds files from the data/weather/stage-012-technically-correct-dataframe directory. It ensures that the data frames have common lower case column names, along with appropriate data types. The transformed data frames are stored within the 'stage-013-technically-correct-complete' directory. The function associated with this process has been exported from the package and it can be run using the command below. The function returns a list containing the destination paths of the successfully saved files, along with any errors.

```
pids.wellbeing.weather::question_1_005_svc_tech_weather_complete()
```

Within the *pids.wellbeing.weather* package the function described above can be found at the following address:

```
./R/question-1-005-svc-tech-weather-complete.R
```


Technically correct: weather single data frame

The stage of technically correct data processing transforms the .Rds files from the ‘stage-013-technically-correct-complete’ directory into a single data frame, which is stored within ‘stage-014-technically-complete-single-dataframe’. The function associated with this process has been exported from the package and it can be run using the command below. The function returns the destination path of the successfully saved file, along with any errors.

```
pids.wellbeing.weather::question_1_006_svc_tech_weather_single_df()
```

Within the *pids.wellbeing.weather* package the function described above can be found at the following address:

```
./R/question-1-006-svc-tech-weather-single-df.R
```

Consistently correct

Summary Type	Sunshine (hours)	Rain (mm)	Temp max (c)	Temp min (c)
Mean	119.119	75.032	12.816	5.998
Min	89.526	45.803	9.496	2.734
Max	155.310	151.117	14.968	8.386
SD	16.196	26.958	1.343	1.194

Table 1: Mean, min, max and SD for the weather features (3sf)

```
pids.wellbeing.weather::question_1_007_svc_cons_summary()  
pids.wellbeing.weather::question_1_008_svc_cons_grouped_data()  
pids.wellbeing.weather::question_1_009_svc_cons_grouped_data_summary()
```

Within the *pids.wellbeing.weather* package the functions described above can be found at the following address:

```
./R/question-1-007-svc-cons-summary.R  
./R/question-1-008-svc-cons-grouped-data.R  
./R/question-1-009-svc-cons-grouped-data-summary.R
```

Exploratory data analysis (EDA)

Box and whisker

```
pids.wellbeing.weather::question_1_010_eda_charts_box_whisker_hours_sun()  
pids.wellbeing.weather::question_1_010_eda_charts_box_whisker_hours_rain()  
pids.wellbeing.weather::question_1_010_eda_charts_box_whisker_max_temp()  
pids.wellbeing.weather::question_1_010_eda_charts_box_whisker_min_temp()
```

Within the *pids.wellbeing.weather* package the functions described above can be found at the following address:

```
./R/question_1_010_eda_charts_box_whisker_hours_sun.R  
./R/question_1_010_eda_charts_box_whisker_hours_rain.R  
./R/question_1_010_eda_charts_box_whisker_max_temp.R  
./R/question_1_010_eda_charts_box_whisker_min_temp.R
```

Remove outliers

Summary Type	Sunshine (hours)	Rain (mm)	Temp max (c)	Temp min (c)
Mean	118.43	69.602	12.816	5.998
Min	91.320	47.399	9.496	2.734
Max	149.349	127.368	14.968	8.386
SD	14.032	18.998	1.343	1.194

Table 2: Mean, min, max and SD for the weather features (3sf)

```
pids.wellbeing.weather::question_1_011_eda_remove_outliers()  
pids.wellbeing.weather::question_1_012_eda_remove_outliers_summary()
```

Within the *pids.wellbeing.weather* package the functions described above can be found at the following address:

```
./R/question-1-011-eda-remove-outliers.R  
./R/question-1-012-eda-remove-outliers-summary.R
```

Weather variables by latitude

```
pids.wellbeing.weather::question_1_013_eda_charts_latitude()
```

Within the *pids.wellbeing.weather* package the function described above can be found at the following address:

```
./R/question-1-013-eda-charts-latitude.R
```

Weather variables by latitude and longitude

```
pids.wellbeing.weather::question_1_014_eda_charts_longitude_latitude_hours_sun()  
pids.wellbeing.weather::question_1_015_eda_charts_longitude_latitude_rain()  
pids.wellbeing.weather::question_1_016_eda_charts_longitude_latitude_max_temp()  
pids.wellbeing.weather::question_1_017_eda_charts_longitude_latitude_min_temp()
```

Within the *pids.wellbeing.weather* package the functions described above can be found at the following address:

```
./R/question-1-014-eda-charts-longitude-latitude-hours-sun.R  
./R/question-1-015-eda-charts-longitude-latitude-rain.R  
./R/question-1-016-eda-charts-longitude-latitude-max-temp.R  
./R/question-1-017-eda-charts-longitude-latitude-min-temp.R
```

Cluster tendency

The clustering tendency of the data has been calculated using the Hopkins statistic (H). It assesses the probability that the data contains non random structures. The statistic has been calculated using the *factoextra* dependency. Using the data with outliers removed, the result of H was **0.352**.

- (1) high score—uniform distribution—no cluster
- (2) low score—not uniform distribution—(may be not) cluster.

The function used to generate this statistic can be run using the command below:

```
pids.wellbeing.weather::question_1_018_prep_cluster_tendency()  
pids.wellbeing.weather::question_1_018_prep_cluster_tendency(show_chart = TRUE)
```

Within the *pids.wellbeing.weather* package the function described above can be found at the following address:

```
./R/question-1-018-prep-cluster-tendency.R
```

VAT

```
pids.wellbeing.weather::question_1_019_prep_charts_vat()
```

Within the *pids.wellbeing.weather* package the function described above can be found at the following address:

```
./R/question-1-019-prep-charts-vat.R
```

Data preparation

K choice

```
pids.wellbeing.weather::question_1_020_prep_k_choice()
```

Within the *pids.wellbeing.weather* package the function described above can be found at the following address:

```
./R/question-1-020-prep-k-choice.R
```

Algorithm choice

Model type	Speed	Independence	Unequal cluster sizes	Unusual density
K-Means			No	No
K-Medoids			Yes	No
Hierarchical				
Distribution based			Yes	Yes

Table 3: Summary of clustering algorithms by feature)

```
pids.wellbeing.weather::question_1_021_prep_algorithm_choice()
```

Within the *pids.wellbeing.weather* package the function described above can be found at the following address:

```
./R/question-1-021-prep-algorithm-choice.R
```

Linkage choice

```
pids.wellbeing.weather::question_1_022_linkage_choice()
```

Within the *pids.wellbeing.weather* package the function described above can be found at the following address:

```
./R/question-1-022-linkage-choice.R
```

Analysis

Hierarchical

```
pids.wellbeing.weather::question_1_023_analysis_hierarchical()  
pids.wellbeing.weather::question_1_024_analysis_charts_dendrogram()  
pids.wellbeing.weather::question_1_025_analysis_attach_pruned_cluster_values()
```

Within the *pids.wellbeing.weather* package the functions described above can be found at the following address:

```
./R/question-1-023-analysis-hierarchical.R  
./R/question-1-024-analysis-charts-dendrogram.R  
./R/question-1-025-analysis-attach-pruned-cluster-values.R
```

k-means

```
pids.wellbeing.weather::question_1_026_analysis_kmeans()  
pids.wellbeing.weather::question_1_027_analysis_charts_sum_squares()
```

Within the *pids.wellbeing.weather* package the functions described above can be found at the following address:

```
./R/question-1-026-analysis-kmeans.R  
./R/question-1-027-analysis-charts-sum-squares.R
```

Evaluation

Hierarchical

```
pids.wellbeing.weather::question_1_032_eval_charts_hierarchical_latitude()  
pids.wellbeing.weather::question_1_033_eval_charts_hier_latitude_longitude_hours_sun  
pids.wellbeing.weather::question_1_034_eval_charts_hier_latitude_longitude_rain_mm  
pids.wellbeing.weather::question_1_035_eval_charts_hier_latitude_longitude_max_temp  
pids.wellbeing.weather::question_1_036_eval_charts_hier_latitude_longitude_min_temp
```

Within the *pids.wellbeing.weather* package the functions described above can be found at the following address:

```
./R/question-1-032-eval-charts-hier-latitude.R  
./R/question-1-033-eval-charts-hier-latitude-longitude-hours-sun.R  
./R/question-1-034-eval-charts-hier-latitude-longitude-rain-mm.R  
./R/question-1-035-eval-charts-hier-latitude-longitude-max-temp.R  
./R/question-1-036-eval-charts-hier-latitude-longitude-min-temp.R
```

Deployment

Question 2

Business Understanding

Data Understanding

Append latitude categories

```
pids.wellbeing.weather::question_2_001_bu_append_latitude_categories()
```

Within the *pids.wellbeing.weather* package the functions described above can be found at the following address:

```
./R/question-2-001-bu-append-latitude-categories()
```

Exploratory data analysis (EDA)

Number of weather stations per latitude category

```
pids.wellbeing.weather::question_2_002_eda_latitude_category_summary()$nrow
```

The function described above can be found within the *pids.wellbeing.weather* package at the following address:

```
./R/question-2-002-eda-latitude-category-summary.R
```

Weather stations per latitude category

```
pids.wellbeing.weather::question_2_002_eda_latitude_category_summary()$weather_st
```

The function described above can be found within the *pids.wellbeing.weather* package at the following address:

```
./R/question-2-002-eda-latitude-category-summary.R
```

Mean weather features per latitude category

```
pids.wellbeing.weather::question_2_002_eda_latitude_category_summary()$means  
pids.wellbeing.weather::question_2_003_eda_charts_weather_features()
```

The function described above can be found within the *pids.wellbeing.weather* package at the following address:

```
./R/question-2-002-eda-latitude-category-summary.R  
./R/question-2-003-eda-charts-weather-features.R
```

Latitude category pairwise comparisons

```
pids.wellbeing.weather::question_2_004_eda_charts_latitude_category_pairwise()
```

The function described above can be found within the *pids.wellbeing.weather* package at the following address:

```
./R/question-2-004-eda-charts-latitude-category-pairwise.R
```

Latitude and longitude per latitude category

```
pids.wellbeing.weather::question_2_005_eda_charts_latitude_longitude_hours_sun()  
pids.wellbeing.weather::question_2_005_eda_charts_latitude_longitude_rain()  
pids.wellbeing.weather::question_2_005_eda_charts_latitude_longitude_temp_max  
pids.wellbeing.weather::question_2_005_eda_charts_latitude_longitude_temp_min
```

The functions described above can be found within the *pids.wellbeing.weather* package at the following address:

```
./R/question-2-005-eda-charts-latitude-longitude.R
```

Data Preparation

Training and test data split

```
pids.wellbeing.weather::question_2_006_prep_split()
```

The function described above can be found within the *pids.wellbeing.weather* package at the following address:

```
./R/question-2-006-prep-split.R
```


Model type	Speed	Normal	Independence	Conditional prob
K-NN				
Multiple Logistic Regression				
Naive Bayes				

Table 4: Summary of clustering algorithms by feature)

Algorithm choice

Analysis

K Nearest Neighbours

```
pids.wellbeing.weather::question_2_007_analysis_knn()
```

The function described above can be found within the *pids.wellbeing.weather* package at the following address:

```
./R/question-2-007-analysis-knn.R
```

Question 3

Business Understanding

Data Understanding

Statistical value chain

Technically correct: well-being

```
pids.wellbeing.weather::question_3_001_svc_tech_wellbeing()
```

```
./R/question-3-001-svc-tech-wellbeing.R
```

Consistently correct: weather and well-being

```
pids.wellbeing.weather::question_3_002_svc_cons_weather_add_boundaries()
```

```
pids.wellbeing.weather::question_3_003_svc_cons_weather_wellbeing_join()
```

```
pids.wellbeing.weather::question_3_004_svc_cons_weather_wellbeing_summary()
```

```
./R/question-3-002-svc-cons-weather-add-boundaries.R
```

```
./R/question-3-003-svc-cons-weather-wellbeing-join.R
```

```
./R/question-3-004-svc-cons-weather-wellbeing-summary.R
```

Exploratory data analysis (EDA)

```
pids.wellbeing.weather::question_3_005_eda_charts_wellbeing_by_region()
```

```
pids.wellbeing.weather::question_3_006_eda_charts_correlation()
```

```
./R/question-3-005-eda-charts-wellbeing-by-region.R
```

```
./R/question-3-006-eda-charts-correlation.R
```

Data Preparation

Analysis

```
pids.wellbeing.weather::question_3_007_analysis_regression()  
pids.wellbeing.weather::question_3_008_analysis_regression_test()
```

```
./R/question-3-007-analysis-regression.R  
./R/question-3-008-analysis-regression-test.R
```

Evaluation

```
pids.wellbeing.weather::question_3_009_eval_regression_summary()  
pids.wellbeing.weather::question_3_010_eval_charts_predicted_happinesst()
```

```
./R/question-3-009-eval-regression-summary.R  
./R/question-3-010-eval-charts-predicted-happiness.R
```

Bibliography

- [1] P. Chapman, J. Clinton, R. Kerber, T. Khabaza, T. Reinartz, C. Shearer, and R. Wirth, “Crisp-dm 1.0.” <https://www.the-modeling-agency.com/crisp-dm.pdf>, January 2000. [Online].
- [2] “pids.wellbeing.weather.” <https://github.com/b136325/pids.wellbeing.weather/>. [Online].
- [3] “Packrat.” <https://rstudio.github.io/packrat/>. [Online].
- [4] H. Wickham, “R pckages.” <http://r-pkgs.had.co.nz/intro.html>. [Online].
- [5] “Lintr.” <https://github.com/jimhester/lintr>. [Online].
- [6] “Semvar.” <https://semver.org/>. [Online].
- [7] “Git.” <https://git-scm.com/>. [Online].

Appendices

A. Raw data output