

Chose Your Own Project - Machine Learning Submission

HarvardX Data Science Capstone - PH125.9x

Simon Gibson

2023-12-03

Contents

Introduction	1
Method	1

Introduction

For the 9th Course in the HarvardX Data Science course we have been asked to create two recommendation systems. The first was a Movie Recommendation System using the MovieLens dataset. The second is a “Choose your Own Project.” For this a we are targetting a Workforce Recommendation System - mixing weather forecasts with Police 911 call information to see if it is possible to predict Police staffing requirements based on weather based trends.

We are using the Seattle Police Department 911 Incident Response data set found here : <https://www.kaggle.com/datasets/sohier/seattle-police-department-911-incident-response>

For Weather data we will use National Oceanic and Atmospheric Administration (NOAA) data. Michael Minns’ tutorial is inciteful for weather analysis. It can be found here: <https://michaelminn.net/tutorials/r-weather/index.html> This weather data does not appear to be available via an api call or similar and is quite a manual download process. Due to download constraints we will be using a locally sourced dataset covering the years 2001 to 2002.

In order to test the results of the recommendation system we are using the root-mean-square error (RMSE) to measure the difference between the values predicted by the model and the observed values.

Method

The first step is to clear any set variables so we do not introduce anything unexpected into the data we are working with.

Then we install the packages required to manipulate the data.

```
#####
# This code is divided into the following sections #
# 1. Install required packages                      #
# 2. edx code for creating data sets                 #
# 3. Data set exploration                           #
#####

#####
# 1. Install required packages and download data    #
#####

# Note: this process takes a couple of minutes
```

```

if(!require(tidyverse)) install.packages("tidyverse", repos = "https://cran.us.r-project.org")
if(!require(caret)) install.packages("caret", repos = "https://cran.us.r-project.org")
if(!require(dplyr)) install.packages("dplyr", repos = "https://cran.us.r-project.org")
if(!require(kableExtra)) install.packages("kableExtra", repos = "https://cran.us.r-project.org")
if(!require(lubridate)) install.packages("lubridate", repos = "https://cran.us.r-project.org")
if(!require(scales)) install.packages("scales", repos = "https://cran.us.r-project.org")
if(!require(stringr)) install.packages("stringr", repos = "https://cran.us.r-project.org")
if(!require(readr)) install.packages("readr", repos = "https://cran.us.r-project.org")

library(tidyverse)
library(caret)
library(dplyr)
library(kableExtra)
library(lubridate)
library(scales)
library(stringr)
library(readr)

```

Following that, the data is downloaded and then divided into 2 sets. The first set is used to train the algorithm and the second set is used to validate the algorithm. By dividing the data the problem of over-training and thus producing skewed results can be avoided.

The creation of the 2 sets involves the following steps. Initially required packages are installed if not installed and then loaded. Next the data is downloaded if the zip files are not found. Column names are set and the data is converted into forms more easily processed. Then the data is joined. Finally the joined data is split into 2 sets - the edx set used to train the algorithm and the final_holdout_test set that will be used to validate the algorithm and calculate the final RMSE score.

```

#Seattle Police Department 911 Incident Response
#https://www.kaggle.com/datasets/sohier/seattle-police-department-911-incident-response/download?dataset=

#National Oceanic and Atmospheric Administration (NOAA) data
#https://www.ncei.noaa.gov/orders/cdo/3533326.csv

options(timeout = 120)

dl <- "archive.zip"
if(!file.exists(dl))
  download.file("https://www.kaggle.com/datasets/sohier/seattle-police-department-911-incident-response", dl)

dl <- "3533326.csv"
if(!file.exists(dl))
  download.file("https://www.ncei.noaa.gov/orders/cdo/3533326.csv", dl)

#Load Seattle 0911 Call data
Seattle_911 <- read_csv("Seattle_Police_Department_911_Incident_Response.csv")
#Load weather data
Weather <- read_csv("3533326.csv")

##Data Investigation

head(Seattle_911)

## # A tibble: 6 x 19
##   `CAD CDW ID` `CAD Event Number` General Offense Numbe~1 `Event Clearance Code`
##   <chr>                <dbl>                <dbl> <chr>

```

```
## 1 15736          10000246357          2010246357 242
## 2 15737          10000246471          2010246471 065
## 3 15738          10000246255          2010246255 250
## 4 15739          10000246473          2010246473 460
## 5 15740          10000246330          2010246330 250
## 6 15741          10000246477          2010246477 281
## # i abbreviated name: 1: `General Offense Number`
## # i 15 more variables: `Event Clearance Description` <chr>,
## #   `Event Clearance SubGroup` <chr>, `Event Clearance Group` <chr>,
## #   `Event Clearance Date` <chr>, `Hundred Block Location` <chr>,
## #   `District/Sector` <chr>, `Zone/Beat` <chr>, `Census Tract` <chr>,
## #   Longitude <dbl>, Latitude <dbl>, `Incident Location` <chr>,
## #   `Initial Type Description` <chr>, `Initial Type Subgroup` <chr>, ...
```

```
head(Weather)
```

```
## # A tibble: 6 x 30
##   STATION NAME DATE      PRCP SNOW TAVG TMAX TMIN TSUN WT01 WT02 WT03
##   <chr> <chr> <date> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1 USC004~ BREM~ 2000-01-01 0.23 0 NA 44 38 NA NA NA NA
## 2 USC004~ BREM~ 2000-01-02 0 0 NA 44 31 NA NA NA NA
## 3 USC004~ BREM~ 2000-01-03 0.1 0 NA 45 32 NA NA NA NA
## 4 USC004~ BREM~ 2000-01-04 1.38 0 NA 47 35 NA NA NA NA
## 5 USC004~ BREM~ 2000-01-05 0.02 0 NA 51 30 NA NA NA NA
## 6 USC004~ BREM~ 2000-01-06 0.01 0 NA 44 34 NA NA NA NA
## # i 18 more variables: WT04 <lgl>, WT05 <dbl>, WT06 <lgl>, WT07 <lgl>,
## #   WT08 <dbl>, WT09 <lgl>, WT11 <lgl>, WT13 <dbl>, WT14 <dbl>, WT15 <lgl>,
## #   WT16 <dbl>, WT17 <lgl>, WT18 <lgl>, WT19 <lgl>, WT21 <dbl>, WT22 <lgl>,
## #   WV01 <dbl>, WV03 <lgl>
```

```
summary(Seattle_911)
```

```
##   CAD CDW ID      CAD Event Number   General Offense Number
## Length:1433853   Min.   :9.000e+09   Min.   :2.011e+04
## Class :character 1st Qu.:1.200e+10   1st Qu.:2.010e+09
## Mode  :character Median :1.400e+10   Median :2.012e+09
##                  Mean  :1.366e+10   Mean  :1.641e+09
##                  3rd Qu.:1.600e+10   3rd Qu.:2.015e+09
##                  Max.   :1.700e+10   Max.   :2.012e+10
##
## Event Clearance Code Event Clearance Description Event Clearance SubGroup
## Length:1433853      Length:1433853      Length:1433853
## Class :character    Class :character      Class :character
## Mode  :character    Mode  :character      Mode  :character
##
##
##
## Event Clearance Group Event Clearance Date Hundred Block Location
## Length:1433853      Length:1433853      Length:1433853
## Class :character    Class :character      Class :character
## Mode  :character    Mode  :character      Mode  :character
##
##
##
```

```

##
## District/Sector      Zone/Beat      Census Tract      Longitude
## Length:1433853      Length:1433853      Length:1433853      Min.      :-122.4
## Class :character     Class :character     Class :character     1st Qu.:-122.3
## Mode  :character     Mode  :character     Mode  :character     Median :-122.3
##                                     Mean  :-122.3
##                                     3rd Qu.:-122.3
##                                     Max.  :-122.2
##                                     NA's   :1
## Latitude      Incident Location  Initial Type Description
## Min.      :47.45      Length:1433853      Length:1433853
## 1st Qu. :47.59      Class :character     Class :character
## Median :47.61      Mode  :character     Mode  :character
## Mean   :47.62
## 3rd Qu. :47.66
## Max.   :47.78
## NA's    :1
## Initial Type Subgroup Initial Type Group At Scene Time
## Length:1433853      Length:1433853      Length:1433853
## Class :character     Class :character     Class :character
## Mode  :character     Mode  :character     Mode  :character
##
##
##
##

```

summary(Weather)

```

## STATION      NAME      DATE      PRCP
## Length:17773      Length:17773      Min.      :2000-01-01      Min.      :0.0000
## Class :character     Class :character     1st Qu. :2000-10-03      1st Qu. :0.0000
## Mode  :character     Mode  :character     Median   :2001-07-07      Median   :0.0000
##                                     Mean     :2001-07-04      Mean     :0.1278
##                                     3rd Qu. :2002-04-07      3rd Qu. :0.1200
##                                     Max.     :2002-12-31      Max.     :4.3000
##                                     NA's     :110
## SNOW      TAVG      TMAX      TMIN
## Min.      : 0.000      Min.      : 0.00      Min.      : 0.00      Min.      :-16.00
## 1st Qu. : 0.000      1st Qu. :44.00      1st Qu. :50.00      1st Qu. : 36.00
## Median : 0.000      Median :51.00      Median :58.00      Median : 43.00
## Mean   : 0.042      Mean   :52.27      Mean   :59.14      Mean   : 43.21
## 3rd Qu. : 0.000      3rd Qu. :60.00      3rd Qu. :68.00      3rd Qu. : 50.00
## Max.   :24.000      Max.   :82.00      Max.   :99.00      Max.   : 77.00
## NA's   :7233      NA's   :11397      NA's   :2511      NA's   :2537
## TSUN      WT01      WT02      WT03
## Min.      : 0.00      Min.      :1      Min.      :1      Min.      :1
## 1st Qu. : 0.00      1st Qu. :1      1st Qu. :1      1st Qu. :1
## Median : 0.00      Median :1      Median :1      Median :1
## Mean   : 31.76      Mean   :1      Mean   :1      Mean   :1
## 3rd Qu. : 0.00      3rd Qu. :1      3rd Qu. :1      3rd Qu. :1
## Max.   :931.00      Max.   :1      Max.   :1      Max.   :1
## NA's   :14935      NA's   :16900      NA's   :17707      NA's   :17755
## WT04      WT05      WT06      WT07      WT08
## Mode:logical      Min.      :1      Mode:logical      Mode:logical      Min.      :1
## TRUE:7      1st Qu. :1      TRUE:1      NA's:17773      1st Qu. :1

```

##	NA's:17766	Median :1	NA's:17772	Median :1
##		Mean :1		Mean :1
##		3rd Qu.:1		3rd Qu.:1
##		Max. :1		Max. :1
##		NA's :17761		NA's :17752
##	WT09	WT11	WT13	WT14
##	Mode:logical	Mode:logical	Min. :1	Min. :1
##	NA's:17773	NA's:17773	1st Qu.:1	1st Qu.:1
##			Median :1	Median :1
##			Mean :1	Mean :1
##			3rd Qu.:1	3rd Qu.:1
##			Max. :1	Max. :1
##			NA's :17286	NA's :17688
##	WT16	WT17	WT18	WT19
##	Min. :1	Mode:logical	Mode:logical	Mode:logical
##	1st Qu.:1	TRUE:1	TRUE:25	TRUE:2
##	Median :1	NA's:17772	NA's:17748	NA's:17771
##	Mean :1			Median :1
##	3rd Qu.:1			Mean :1
##	Max. :1			3rd Qu.:1
##	NA's :17209			Max. :1
##				NA's :17725
##	WT22	WV01	WV03	
##	Mode:logical	Min. :1	Mode:logical	
##	TRUE:7	1st Qu.:1	TRUE:2	
##	NA's:17766	Median :1	NA's:17771	
##		Mean :1		
##		3rd Qu.:1		
##		Max. :1		
##		NA's :17767		