

Chose Your Own Project - Machine Learning Submission

HarvardX Data Science Capstone - PH125.9x

Simon Gibson

2024-06-02

Contents

Introduction	1
Method	1
Data Investigation	2
Weather Dataset	2
References	5

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 have chosen 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.

At a high level, the data is downloaded, analysed, the data merged and the merged 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 modified (for example spaces are removed to make the data easier to work with) 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.

Data Investigation

Weather Dataset

Looking at the first 5 rows of data we can see the following:

STATION	NAME	DATE	PRCP	SNOW	TAVG	TMAX	TMIN	TSUN	WT01	WT02	WT03	WT04	WT05	WT06	WT07	WT08	WT09	WT11	WT13	WT14	WT15	WT16	WT17	WT18	WT19	WT21	WT22	WV01	WV03
USC00450872	BREMERTON, WA US	2000-01-01	0.23	0	NA	44	38	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA
USC00450872	BREMERTON, WA US	2000-01-02	0.00	0	NA	44	31	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA
USC00450872	BREMERTON, WA US	2000-01-03	0.10	0	NA	45	32	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA
USC00450872	BREMERTON, WA US	2000-01-04	1.38	0	NA	47	35	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA
USC00450872	BREMERTON, WA US	2000-01-05	0.02	0	NA	51	30	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA
USC00450872	BREMERTON, WA US	2000-01-06	0.01	0	NA	44	34	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA

The weather data starts from 2000-01-01 and ends 2002-12-31. It has the following column headers: STATION, NAME, DATE, PRCP, SNOW, TAVG, TMAX, TMIN, TSUN, WT01, WT02, WT03, WT04, WT05, WT06, WT07, WT08, WT09, WT11, WT13, WT14, WT15, WT16, WT17, WT18, WT19, WT21, WT22, WV01, WV03. If we look at a single row, row 100 in this case, we can see the following: USC00450872, BREMERTON, WA US, 2000-04-09, 0, 0, NA, 66, 41, NA.

For our investigation the following columns may be of interest - STATION, NAME, DATE, PRCP, SNOW, TMAX TMIN. The PRCP (Precipitation) and SNOW (Snowfall) columns are in inches and the TMAX (Maximum Temperature) TMIN (Minimum Temperature) are in Fahrenheit. These will be converted to their metric equivalents using the measurements package. We are interested in determining if there is any relationship between crime reports and weather so we will determine the closest weather station using latitude and longitude bearings returned by the ggmap package.

We have data from 20 stations:

	NAME	STATION
X2000.01.01	BREMERTON, WA US	USC00450872
X2000.01.01.1	EVERETT, WA US	USC00452675
X2000.01.01.2	MONROE, WA US	USC00455525
X2000.01.01.3	TOLT SOUTH FORK RESERVOIR, WA US	USC00458508
X2000.01.01.4	RENTON MUNICIPAL AIRPORT, WA US	USW00094248
X2000.01.01.5	KENT, WA US	USC00454169
X2000.01.01.6	TACOMA NUMBER 1, WA US	USC00458278
X2000.01.01.7	LANDSBURG, WA US	USC00454486
X2000.01.01.8	CEDAR LAKE, WA US	USC00451233
X2000.01.01.9	SNOQUALMIE FALLS, WA US	USC00457773
X2000.01.01.10	WAUNA 3 W, WA US	USC00459021
X2000.01.01.11	PALMER 3 ESE, WA US	USC00456295
X2000.01.01.12	TACOMA NARROWS AIRPORT, WA US	USW00094274
X2000.01.01.13	EVERETT SNOHOMISH CO AIRPORT, WA US	USW00024222
X2000.01.01.14	SEATTLE TACOMA AIRPORT, WA US	USW00024233
X2000.01.01.15	SEATTLE SAND POINT WEATHER FORECAST OFFICE, WA US	USW00094290
X2000.01.01.16	SEATTLE BOEING FIELD, WA US	USW00024234
X2000.11.22.4	GIG HARBOR 3.4 NW, WA US	US1WAPR0075
X2001.08.11.3	OLALLA 1.4 WNW, WA US	US1WAKP0013
X2001.12.02.10	WOODINVILLE 0.9 ENE, WA US	US1WAKG0078

Of 17773 rainfall measurements, 7869 recorded rainfall, and 9794 recorded no rainfall. The maximum rainfall during this period was 109.22mm which fell on 2002-01-07 at BREMERTON, WA US. Heavy rainfall is defined by NIWA as rainfall of over 100mm in 24 hours¹ and this occurred 3 times during the period we have data for.

¹<https://niwa.co.nz/natural-hazards/extreme-weather-heavy-rainfall>

Of 17773 snowfall entries, 158 recorded snowfall, and 10382 recorded no snowfall, with 7233 not recording data. The maximum snowfall during this period was 609.6mm which fell on 2002-02-01 at TOLT SOUTH FORK RESERVOIR, WA US.

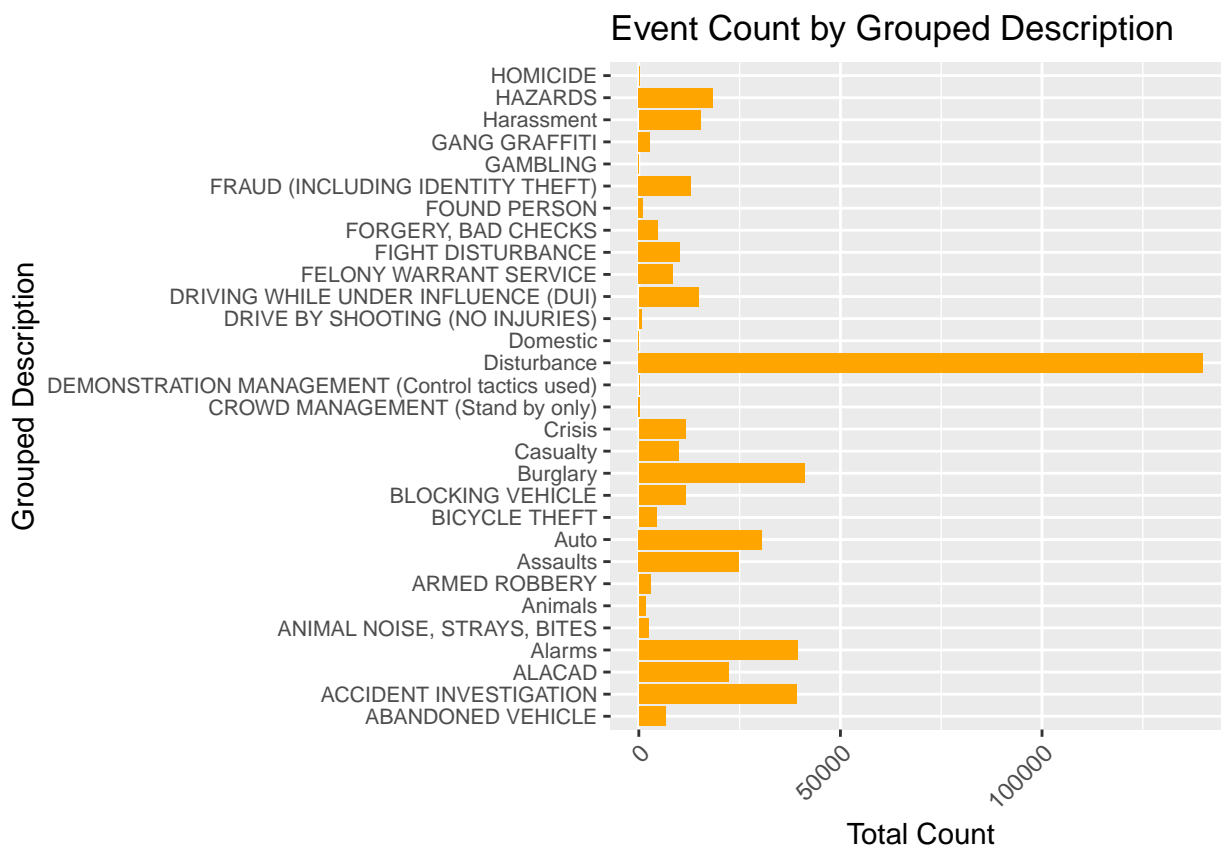
Over the period we have data for we have a maximum temperature of 37.22 at LANDSBURG, WA US and a minimum of -26.67 degrees Celsius at LANDSBURG, WA US. The mean maximum temperature was 15.08 while the mean minimum temperature was 6.23 degrees Celsius.

```
#We have data from `r n_distinct(Weather$STATION)` stations:
```

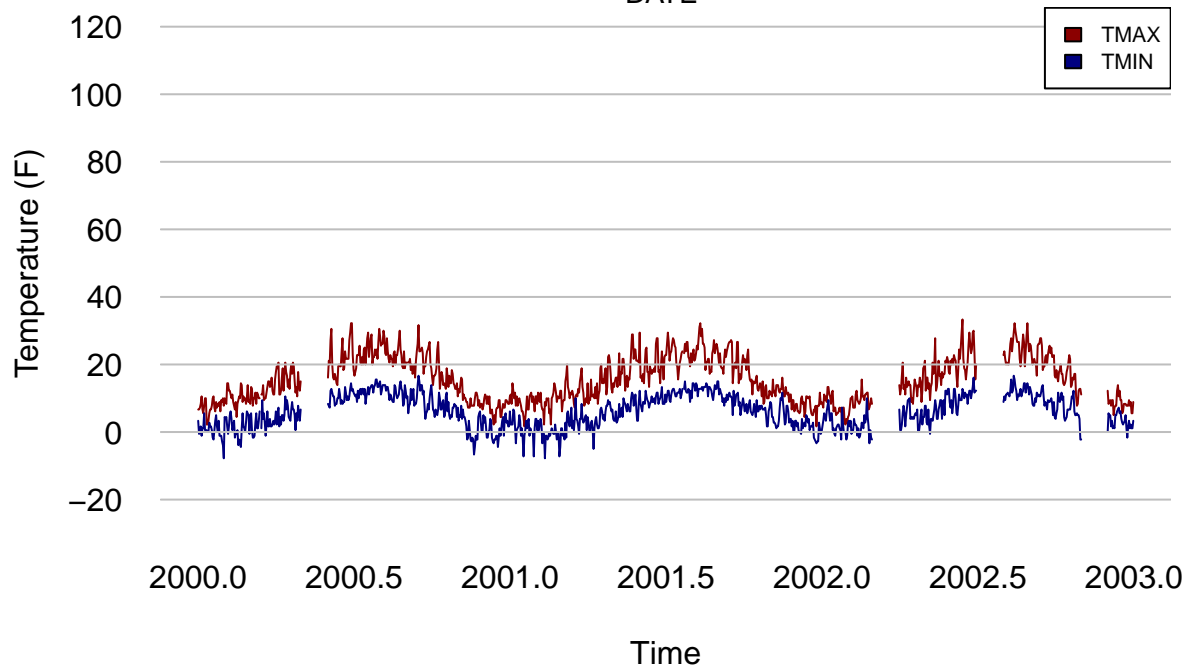
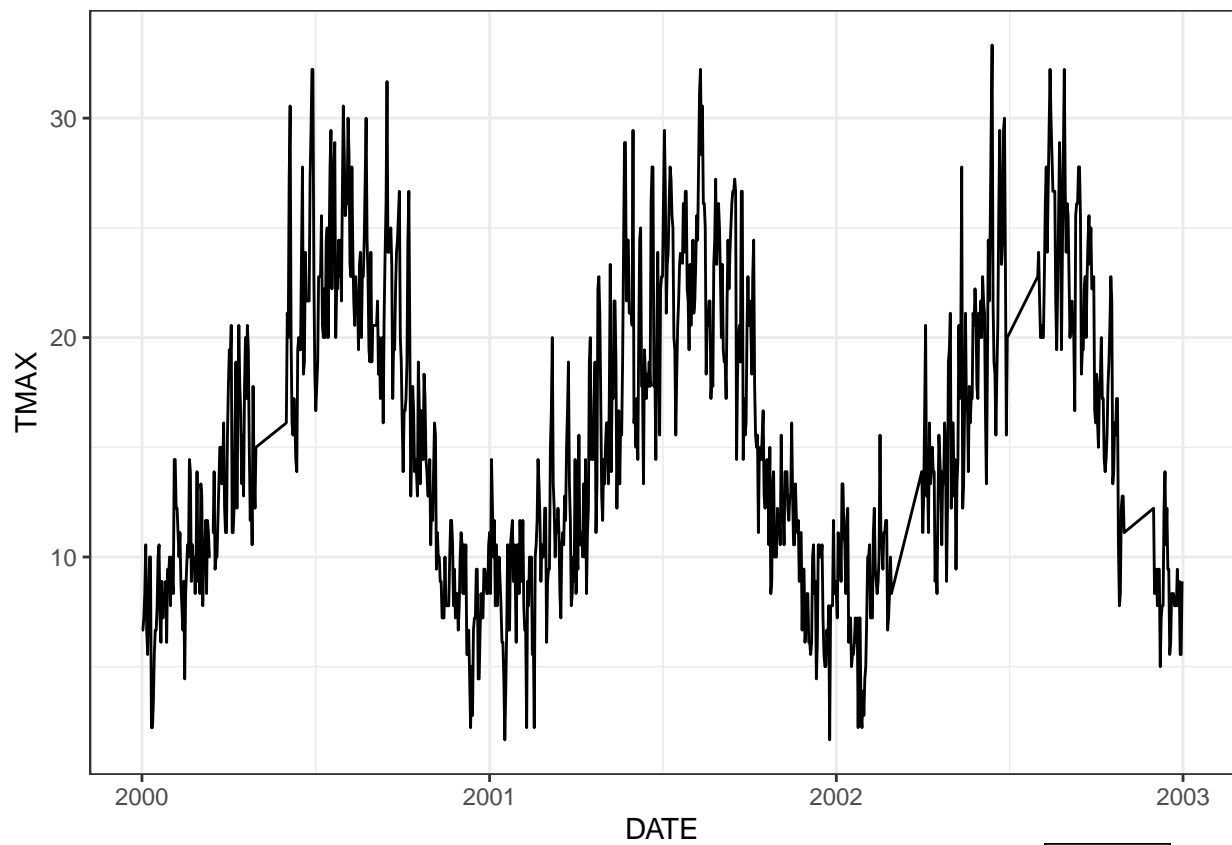
```
#r kable(Seattle_Stations, format = "markdown")`  
#Seattle_Weather %>% group_by(Seattle_Weather$STATION)
```

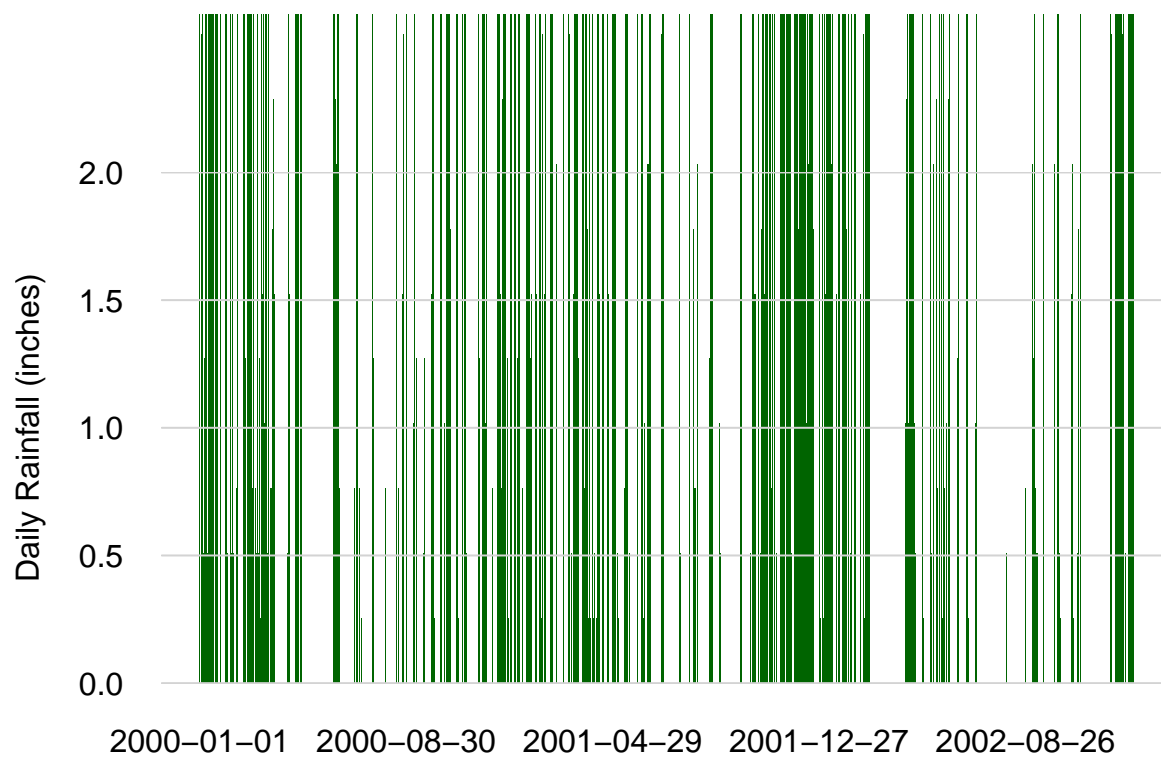
```
## Warning: Expected 2 pieces. Missing pieces filled with `NA` in 1 rows  
## [1288471].
```

The Seattle 911 dataset contains data from 2009-06-17 16:14:00 through to 2017-08-29 11:44:01. During this period, 1433853 CAD events were recorded.



To do - investigation of police data map weather station locations correlate weather station locations with police data





References

- 1.
- 2.
- 3.
4. <https://www.neonscience.org/resources/learning-hub/tutorials/da-viz-coop-precip-data-r>