

Tipologia i cicle de vida de les dades

Pràctica 2. Neteja i anàlisi de dades

Solució

Aitor Ferrus Blasco [aferrus]
Alonso López i Vicente [alopezvic]

05/01/2021

Contents

1. Descripció del dataset.	2
2. Integració i selecció de les dades d'interès a analitzar.	2
3. Neteja de les dades	4
3.1. Les dades contenen zeros o elements buits? Com gestionaries aquests casos?	4
3.2. Identificació i tractament de valors extrems.	4
4. Anàlisi de les dades.	4
4.1. Selecció dels grups de dades.	4
4.2. Comprovació de la normalitat i homogeneïtat de la variància.	5
4.3. Aplicació de proves estadístiques.	5
5. Representació dels resultats.	5
6. Resolució del problema.	5
7. Codi.	5
8. Contribucions	6

1. Descripció del dataset.

Perquè és important i quina pregunta/problema pretén respondre?

Resposta

El dataset que hem escollit és *Rain in Australia* (<https://www.kaggle.com/jsphyg/weather-dataset-rattle-package>).

Conté 10 anys de dades d'observacions diàries del clima en diferents llocs d'Austràlia. Conté una variable objectiu (RainTomorrow) per predir el temps del dia següent. Si la variable és Yes indica que ha plogut el dia següent 1mm o més. Amb aquesta variable podem entrenar models per tal de predir si plourà el dia següent.

Les variables que inclou el dataset són les següents:

Variable	Descripció
Date	La data de l'observació
Location	El nom de la localització de l'estació meteorològica.
MinTemp	La temperatura mínima en graus Celsius
MaxTemp	La temperatura màxima en graus Celsius
Rainfall	La quantitat de pluja registrada durant el dia en mm
Evaporation	La denominada Class A pan evaporation (mm) durant 24 hores a les 9am
Sunshine	El nombre d'hores de sol durant el dia.
WindGustDir	La direcció de la ratxa de vent més forta en les 24 hores fins la mitjanit
WindGustSpeed	La velocitat (km/h) de la ratxa de vent més forta en les 24 hores fins a mitjanit
WindDir9am	Direcció del vent a les 9am
WindDir3pm	Direcció del vent a les 3pm
WindSpeed9am	Mitjana de la Velocitat del vent (km/hr) 10 minuts abans de les 9am
WindSpeed3pm	Mitjana de la Velocitat del vent (km/hr) 10 minuts abans de les 3pm
Humidity9am	Humitat (percentatge) a les 9am
Humidity3pm	Humitat (percentatge) a les 3pm
Pressure9am	Pressió atmosfèrica (hpa) reduïda al nivell mitjà del mar a les 9am
Pressure3pm	Pressió atmosfèrica (hpa) reduïda al nivell mitjà del mar a les 3pm
Cloud9am	Fracció del cel enfosquida pels núvols a les 9am. Es mesura en "oktas", els quals són una unitat de vuitens. Registre quants hi ha
Cloud3pm	Fracció del cel enfosquida pels núvols a les 3pm. Es mesura en "oktas", els quals són una unitat de vuitens. Registre quants hi ha
Temp9am	Temperatura (graus Celsius) a les 9am
Temp3pm	Temperatura (graus Celsius) a les 3pm
RainToday	Booleà: 1 si la precipitació (mm) en les 24 hores anteriors a les 9am és superior a 1mm, sinó 0
RainTomorrow	La quantitat de pluja al dia següent en mm. Utilitzada per crear la variable resposta RainTomorrow. Un tipus de mesura del "risc".

2. Integració i selecció de les dades d'interès a analitzar.

Resposta

Hem seleccionat les dades de Melbourne, ja que tenen poques NA. Creiem que l'anàlisi que es pot realitzar en aquesta localització és pot adaptar ràpidament a qualsevol de les altres estacions que inclou el dataset.

```
library(readr)
weatherAUS <- read_csv("weatherAUS.csv",
  col_types = cols(Date = col_date(format = "%Y-%m-%d"),
    Evaporation = col_double(), Sunshine = col_double()))
```

```
weatherMelb <- weatherAUS[weatherAUS$Location == "Melbourne",]
summary(weatherMelb)
```

```
##      Date      Location      MinTemp      MaxTemp
## Min.   :2008-07-01 Length:3193 Min.    : 1.40 Min.    : 9.70
## 1st Qu.:2010-09-07 Class :character 1st Qu.: 8.70 1st Qu.:16.10
## Median :2013-01-13 Mode  :character Median :11.40 Median :19.50
## Mean   :2013-01-02 Mean   :11.78 Mean   :20.77
## 3rd Qu.:2015-04-19 3rd Qu.:14.60 3rd Qu.:24.20
## Max.   :2017-06-25 Max.   :28.60 Max.   :46.40
##                                     NA's   :480 NA's   :481
##      Rainfall      Evaporation      Sunshine      WindGustDir
## Min.    : 0.00 Min.    : 0.00 Min.    : 0.000 Length:3193
## 1st Qu.: 0.00 1st Qu.: 2.20 1st Qu.: 3.100 Class :character
## Median : 0.00 Median : 4.00 Median : 6.500 Mode  :character
## Mean    : 1.87 Mean    : 4.65 Mean    : 6.385
## 3rd Qu.: 1.20 3rd Qu.: 6.40 3rd Qu.: 9.600
## Max.    :82.20 Max.    :23.80 Max.    :13.900
## NA's    :758 NA's    :3 NA's    :1
## WindGustSpeed      WindDir9am      WindDir3pm      WindSpeed9am
## Min.    : 11.00 Length:3193 Length:3193 Min.    : 0.00
## 1st Qu.: 33.00 Class :character Class :character 1st Qu.:11.00
## Median : 43.00 Mode  :character Mode  :character Median :17.00
## Mean    : 45.61 Mean    :19.13
## 3rd Qu.: 56.00 3rd Qu.:26.00
## Max.    :122.00 Max.    :67.00
## NA's    :14 NA's    :2
## WindSpeed3pm      Humidity9am      Humidity3pm      Pressure9am
## Min.    : 0.0 Min.    :14.00 Min.    : 6.00 Min.    : 988.9
## 1st Qu.:15.0 1st Qu.: 58.00 1st Qu.: 41.00 1st Qu.:1012.6
## Median :20.0 Median : 68.00 Median : 51.00 Median :1017.9
## Mean    :22.1 Mean    : 67.55 Mean    : 51.18 Mean    :1017.6
## 3rd Qu.:28.0 3rd Qu.: 78.00 3rd Qu.: 61.00 3rd Qu.:1023.0
## Max.    :76.0 Max.    :100.00 Max.    :100.00 Max.    :1039.0
##                                     NA's    :482 NA's    :487 NA's    :480
## Pressure3pm      Cloud9am      Cloud3pm      Temp9am
## Min.    : 988.3 Min.    :0.000 Min.    :0.000 Min.    : 2.90
## 1st Qu.:1010.7 1st Qu.:3.000 1st Qu.:4.000 1st Qu.:11.28
## Median :1016.1 Median :7.000 Median :6.000 Median :14.10
## Mean    :1015.8 Mean    :5.314 Mean    :5.336 Mean    :14.60
## 3rd Qu.:1021.1 3rd Qu.:7.000 3rd Qu.:7.000 3rd Qu.:17.40
## Max.    :1035.8 Max.    :8.000 Max.    :8.000 Max.    :35.50
## NA's    :483 NA's    :1034 NA's    :1106 NA's    :481
## Temp3pm      RainToday      RainTomorrow
## Min.    : 7.20 Length:3193 Length:3193
## 1st Qu.:14.90 Class :character Class :character
## Median :18.20 Mode  :character Mode  :character
## Mean    :19.26
## 3rd Qu.:22.50
## Max.    :45.40
## NA's    :484
```

3. Neteja de les dades

3.1. Les dades contenen zeros o elements buits? Com gestionaries aquests casos?

Resposta

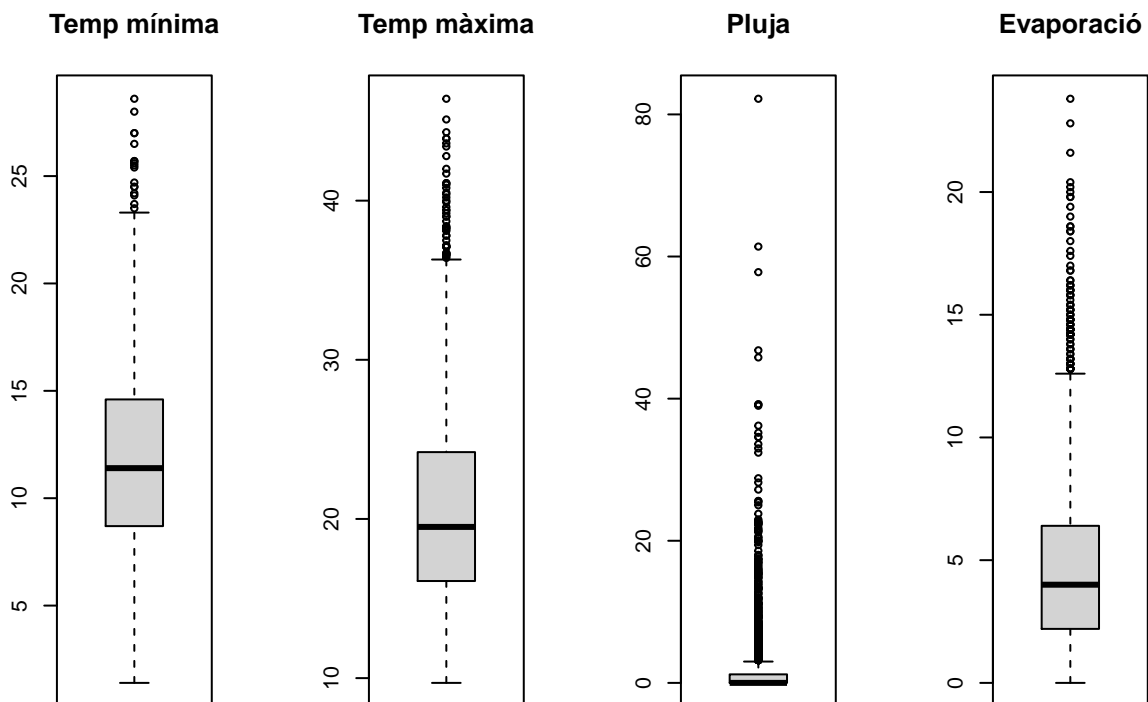
Pendent. Hem parlat d'aplicar kNN.

3.2. Identificació i tractament de valors extrems.

Resposta

Podem identificar valors extrems utilitzant boxplots. Els fem tots?

```
par(mfrow=c(1,4))
boxplot(weatherMelb$MinTemp, na.rm=TRUE, main="Temp mínima")
boxplot(weatherMelb$MaxTemp, na.rm=TRUE, main="Temp màxima")
boxplot(weatherMelb$Rainfall, na.rm=TRUE, main="Pluja")
boxplot(weatherMelb$Evaporation, na.rm=TRUE, main="Evaporació")
```



4. Anàlisi de les dades.

4.1. Selecció dels grups de dades.

Selecció dels grups de dades que es volen analitzar/comparar (planificació dels anàlisis a aplicar).

Resposta

Pendent. Cal establir quines dades utilitzarem per predir la pluja. Cerquem si una (o més) variables ens serveixen per contruir un model que predigui la pluja al dia següent. L'objectiu és determinar si hi ha una relació entre les variables.

4.2. Comprovació de la normalitat i homogeneïtat de la variància.

Resposta

Utilitzem el test de Shapiro-Wilk per comprovar la normalitat. Si el pvalor és inferior a 0.05, el nivell de significació, podem rebutjar la hipòtesi nul·la i concloure que les dades no tenen una distribució normal. En cas contrari, si el pvalor és major que 0.05 podem concloure que les dades segueixen una distribució normal.

```
# Suposem que hem escollit la variable  
shapiro.test(weatherMelb$WindSpeed9am)
```

```
##  
## Shapiro-Wilk normality test  
##  
## data: weatherMelb$WindSpeed9am  
## W = 0.93538, p-value < 2.2e-16
```

Per comprovar la homoscedasticitat, és a dir, la igualtat de variàncies, podem utilitzar el test de Levene si les dades segueixen una distribució normal, o el de Fligner-Killen si les dades no segueixen una distribució normal.

```
# Utilitzem Fligner-Killen perquè les dades no són normals.  
fligner.test(Rainfall ~ WindSpeed9am, data = weatherMelb)
```

```
##  
## Fligner-Killeen test of homogeneity of variances  
##  
## data: Rainfall by WindSpeed9am  
## Fligner-Killeen:med chi-squared = 237.69, df = 36, p-value < 2.2e-16
```

4.3. Aplicació de proves estadístiques.

Per comparar els grups de dades. En funció de les dades i de l'objectiu de l'estudi, aplicar proves de contrast d'hipòtesis, correlacions, regressions, etc. Aplicar almenys tres mètodes d'anàlisi diferents.

Resposta

5. Representació dels resultats.

A partir de taules i gràfiques.

Resposta

6. Resolució del problema.

A partir dels resultats obtinguts, quines són les conclusions? Els resultats permeten respondre al problema?

Resposta

7. Codi.

Cal adjuntar el codi, preferiblement en R, amb el que s'ha realitzat la neteja, anàlisi i representació de les dades. Si ho preferiu, també podeu treballar en Python.

Resposta

8. Contribucions

Contribucions	Firma
Investigació prèvia	Aitor Ferrus Blasco, Alonso López i Vicente
Redacció de les respostes	Aitor Ferrus Blasco, Alonso López i Vicente
Desenvolupament codi	Aitor Ferrus Blasco, Alonso López i Vicente