

# PES University, Bangalore

Established under Karnataka Act No. 16 of 2013

UE20CS312 - Data Analytics - Worksheet 1b - Correlation Analysis

Ishita Bharadwaj, Dept. of CSE - PES1UG20CS648

Collaborated with Hita - PES1UG20CS645

## Correlation

Correlation is a measure of the strength and direction of relationship that exists between two random variables and is measured using correlation coefficient. Correlation can assist data scientists to choose the variables for model building that is used for solving an analytics problem.

### Solution 1

Find the total number of accidents in each state for the year 2016 and display your results. Make sure to display all rows while printing the dataframe. Print only the necessary columns. (Hint: use the grep command to help filter out column names).

```
library(ggpubr)
library(dplyr)
df <- read.csv('road_accidents_india_2016.csv', row.names=1)

acc_cols<-grep("Total.Accidents$",colnames(df),ignore.case = TRUE, value = TRUE)
print(acc_cols)

## [1] "Fine.Clear...Total.Accidents" "Mist..Foggy...Total.Accidents"
## [3] "Cloudy...Total.Accidents"    "Rainy...Total.Accidents"
## [5] "Snowfall...Total.Accidents"  "Hail.Sleet...Total.Accidents"
## [7] "Dust.Storm...Total.Accidents" "Others...Total.Accidents"

totalAccidents<-data.frame(state.ut=df$State..UT, total_acc=rowSums(df[,c(acc_cols)], na.rm = TRUE))
print(totalAccidents)

##      state.ut total_acc
## 0 Andhra Pradesh 24888
## 1 Arunachal Pradesh 249
## 2 Assam 7435
## 3 Bihar 8222
## 4 Chhattisgarh 13580
## 5 Goa 4304
## 6 Gujarat 21859
## 7 Haryana 11234
## 8 Himachal Pradesh 3168
## 9 Jammu & Kashmir 5501
## 10 Jharkhand 4932
## 11 Karnataka 44403
## 12 Kerala 39420
## 13 Madhya Pradesh 53972
## 14 Maharashtra 39878
## 15 Manipur 538
## 16 Meghalaya 620
## 17 Mizoram 83
## 18 Nagaland 75
## 19 Orissa 10532
## 20 Punjab 6952
## 21 Rajasthan 23066
## 22 Sikkim 210
## 23 Tamil Nadu 71431
## 24 Telangana 22811
## 25 Tripura 557
## 26 Uttarakhand 1591
## 27 Uttar Pradesh 35612
## 28 West Bengal 13580
## 29 A & N Islands 238
## 30 Chandigarh 428
## 31 D & N Haveli 70
## 32 Daman & Diu 71
## 33 Delhi 7375
## 34 Lakshadweep 1
## 35 Puducherry 1766
```

### Solution 2

Find the (fatality rate =  $\frac{\text{total number of deaths}}{\text{total number of accidents}}$ ) in each state. Find out if there is a significant linear correlation at a significance of  $\alpha = 0.05$  between the fatality rate of a state and the mist/foggy rate (fraction of total accidents that happen in mist/foggy conditions).

Correlation between two continuous RVs: Pearson's correlation coefficient. Pearson's correlation coefficient between two RVs  $x$  and  $y$  is given by:

$$\rho = \frac{\text{Covariance}(x,y)}{\sigma_x \sigma_y}$$

where  $\sigma$  is the standard deviation of a variable.

Plot the fatality rate against the mist/foggy rate. (Hint: use the ggscatter library to plot a scatterplot with the confidence interval of the correlation coefficient).

Plot the fatality rate and mist/foggy rate (see [this](#) and [this](#) for R plot customization).

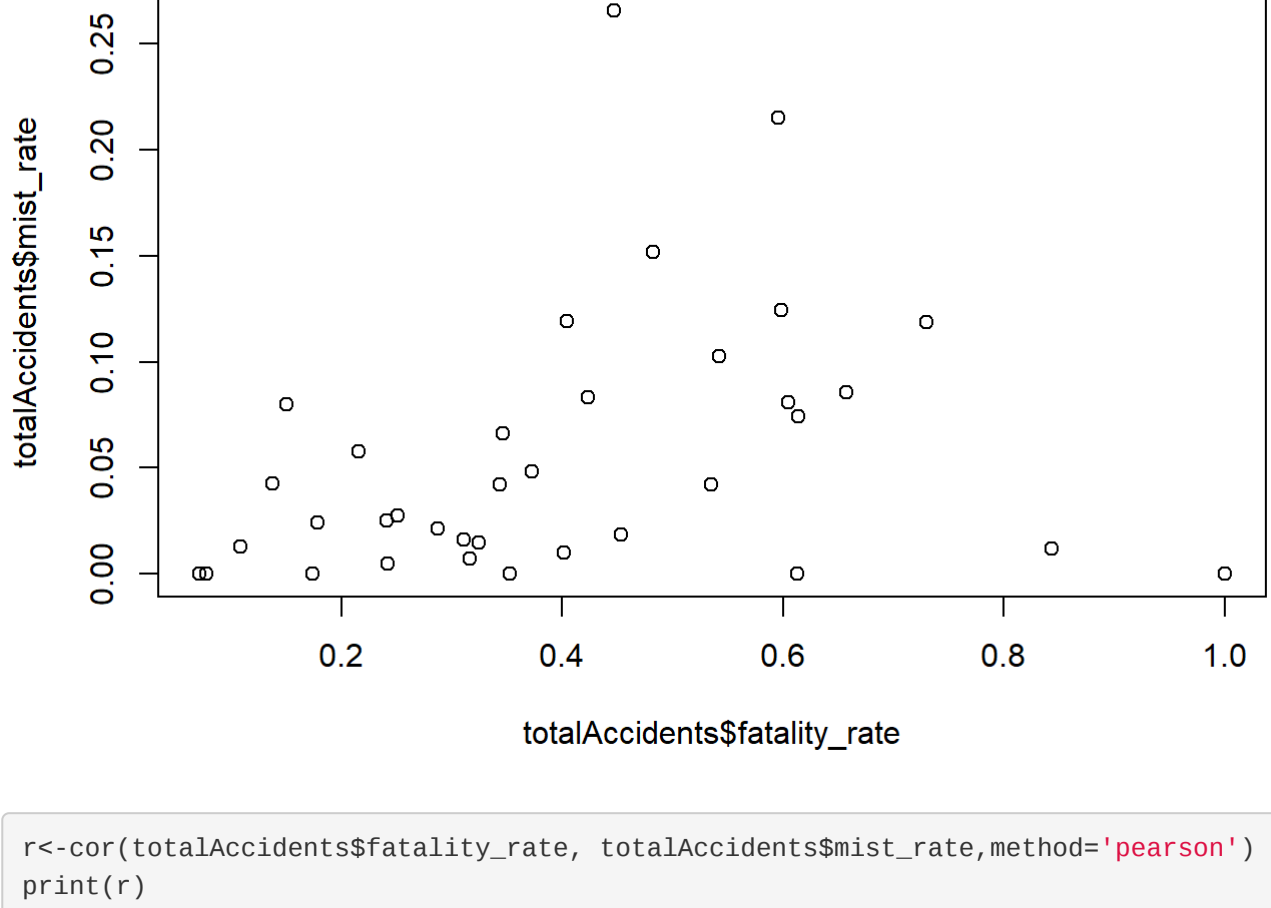
```
death_cols<-grep("Persons.Killed$",colnames(df),ignore.case = TRUE, value = TRUE)
print(death_cols)

## [1] "Fine.Clear...Persons.Killed" "Mist..Foggy...Persons.Killed"
## [3] "Cloudy...Persons.Killed"    "Rainy...Persons.Killed"
## [5] "Snowfall...Persons.Killed"  "Hail.Sleet...Persons.Killed"
## [7] "Dust.Storm...Persons.Killed" "Others...Persons.Killed"

totalAccidents$total_deaths<-NA
totalAccidents$fatality_rate<-NA
totalAccidents$mist_rate<-NA
totalAccidents$total_deaths<-c(rowSums(df[, c(death_cols)], na.rm=TRUE))
totalAccidents$fatality_rate<-c(totalAccidents$total_deaths/totalAccidents$total_acc)
totalAccidents$mist_rate<-df$Mist..Foggy...Total.Accidents/totalAccidents$total_acc
print(head(totalAccidents))

##      state.ut total_acc total_deaths fatality_rate mist_rate
## 0 Andhra Pradesh 24888 8541 0.34317743 0.04222919
## 1 Arunachal Pradesh 249 149 0.59839357 0.12449799
## 2 Assam 7435 2572 0.34593141 0.06603905
## 3 Bihar 8222 4901 0.59608368 0.21515446
## 4 Chhattisgarh 13580 3908 0.28777614 0.02120766
## 5 Goa 4304 336 0.07806691 0.00909000

plot(x=totalAccidents$fatality_rate, y=totalAccidents$mist_rate)
```



```
r<-cor(totalAccidents$fatality_rate, totalAccidents$mist_rate,method='pearson')
print(r)

## [1] 0.2935159

corr_test<-cor.test(totalAccidents$fatality_rate, totalAccidents$mist_rate,method='pearson')
print(corr_test)

##
## Pearson's product-moment correlation
##
## data: totalAccidents$fatality_rate and totalAccidents$mist_rate
## t = 1.7903, df = 34, p-value = 0.08231
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
## -0.03875722 0.56734253
## sample estimates:
## cor
## 0.2935159
```

Correlation coefficient is 0.29 which is not so high. This implies that the strength of linear relationship between fatality\_rate and mist\_rate is weak. This association is a weak association and not significant as p(0.08) is greater than alpha(0.05).

### Solution 3

Rank the states based on total accidents and total fatalities (give a rank of 1 to the state that has the highest value of a property). You are free to use any tie-breaking method for assigning ranks.

Find the Spearman-Rank correlation coefficient between the two rank columns and determine if there is any statistical significance at a significance level of  $\alpha = 0.05$ . Also test the hypothesis that the correlation coefficient is at least 0.2.

The t statistic is given by

$$t = \frac{r_s - \rho_s}{\sqrt{\frac{1 - r_s^2}{n - 2}}}$$

Where  $r_s$  is the calculated Spearman-Rank correlation coefficient and  $\rho_s$  is the value of the population correlation coefficient being tested against.

```
totalAccidents$acc_ranks<-NA
totalAccidents$acc_ranks<-rank(desc(totalAccidents$total_acc),ties.method = 'random')
totalAccidents$death_ranks<-rank(desc(totalAccidents$total_deaths),ties.method = 'random')
print(data.frame(totalAccidents))

##      state.ut total_acc total_deaths fatality_rate mist_rate acc_ranks
## 0 Andhra Pradesh 24888 8541 0.34317743 0.042229187 7
## 1 Arunachal Pradesh 249 149 0.59839357 0.124497992 29
## 2 Assam 7435 2572 0.34593141 0.066039005 16
## 3 Bihar 8222 4901 0.59608368 0.215154464 15
## 4 Chhattisgarh 13580 3908 0.28777614 0.021207658 12
## 5 Goa 4304 336 0.07806691 0.009090000 21
## 6 Gujarat 21859 8136 0.37229367 0.048446864 10
## 7 Haryana 11234 5024 0.44721382 0.265533203 23
## 8 Himachal Pradesh 3168 1271 0.40119949 0.009785354 12
## 9 Jammu & Kashmir 5501 958 0.17415015 0.009090000 19
## 10 Jharkhand 4932 3027 0.61374696 0.074412003 20
## 11 Karnataka 44403 11133 0.25072630 0.027520663 3
## 12 Kerala 39420 4287 0.10875190 0.012683917 5
## 13 Madhya Pradesh 53972 9646 0.17872230 0.024216260 2
## 14 Maharashtra 39878 12935 0.32436431 0.014820202 4
## 15 Manipur 538 81 0.15095762 0.079925651 27
## 16 Meghalaya 620 150 0.24193548 0.004838710 25
## 17 Mizoram 83 70 0.84337349 0.012048193 32
## 18 Nagaland 75 46 0.61333333 0.009090000 33
## 19 Orissa 10532 4463 0.42375617 0.003175005 14
## 20 Punjab 6952 5077 0.73029344 0.118670906 18
## 21 Rajasthan 23066 10465 0.45369808 0.018642157 8
## 22 Sikkim 210 85 0.40476190 0.119047619 31
## 23 Tamil Nadu 71431 17218 0.24104380 0.025353138 1
## 24 Telangana 22811 7219 0.31647012 0.007233352 9
## 25 Tripura 557 173 0.31059246 0.016157989 26
## 26 Uttarakhand 1591 962 0.60465116 0.081081081 24
## 27 Uttar Pradesh 35612 19320 0.54251376 0.102886667 6
## 28 West Bengal 13580 6544 0.48188513 0.151767305 11
## 29 A & N Islands 238 17 0.07142857 0.009090000 30
## 30 Chandigarh 428 151 0.35280374 0.009090000 28
## 31 D & N Haveli 70 46 0.65714286 0.085714286 35
## 32 Daman & Diu 71 38 0.53521127 0.042253521 14
## 33 Delhi 7375 1591 0.21572881 0.057627119 37
## 34 Lakshadweep 1 1 1.00000000 0.009090000 26
## 35 Puducherry 1766 244 0.13816535 0.042468856 33

##      death_ranks
## 0 7
## 1 28
## 2 18
## 3 13
## 4 16
## 5 23
## 6 8
## 7 12
## 8 20
## 9 22
## 10 17
## 11 4
## 12 15
## 13 6
## 14 3
## 15 30
## 16 27
## 17 31
## 18 33
## 19 14
## 20 11
## 21 5
## 22 29
## 23 2
## 24 9
## 25 25
## 26 21
## 27 1
## 28 10
## 29 35
## 30 26
## 31 32
## 32 34
## 33 19
## 34 36
## 35 24

rs<-cor(totalAccidents$acc_ranks, totalAccidents$death_ranks, method = 'spearman')
print(rs)

## [1] 0.957529

corr_test<-cor.test(totalAccidents$acc_ranks, totalAccidents$death_ranks, method = 'spearman')
print(corr_test)

##
## Spearman's rank correlation rho
##
## data: totalAccidents$acc_ranks and totalAccidents$death_ranks
## S = 330, p-value < 2.2e-16
## alternative hypothesis: true rho is not equal to 0
## sample estimates:
## rho
## 0.957529
```

Yes, the correlation coefficient is atleast 0.2

```
degrees<-nrow(totalAccidents)-2
t_stat<-(rs-0.2)/sqrt((1-rs*rs)/(nrow(totalAccidents)-2))
2*pt(q=t_stat, df=degrees, lower.tail=FALSE)

## [1] 7.921248e-17
```

### Solution 4

Convert the column Hail.Sleet...Total.Accidents to a binary column as follows. If a hail/sleet accident has occurred in a state, give that state a value of 1. Otherwise, give it a value of 0. Once converted, find out if there is a significant correlation between the hail\_accident\_occur binary column created and the number of rainy total accidents for every state.

Calculate the point bi-serial correlation coefficient between the two columns. (Hint: it is equivalent to calculating the Pearson correlation between a continuous and a dichotomous variable. You could also use the ltm package's biserial.cor function).

```
df$hail_accident_occurr<-factor(ifelse(df$Hail.Sleet...Total.Accidents==0,0,1))
cat("Hail.Sleet...Total.Accidents:\n",df$Hail.Sleet...Total.Accidents)

## Hail.Sleet...Total.Accidents:
## 92 12 15 0 0 0 64 70 0 0 160 47 0 426 0 10 0 0 0 47 0 28 9 0 0 32 199 6 0 0 0 0 171 0 1

#hail_accident_occurr
print(df$hail_accident_occurr)

## [1] 1 1 1 0 0 0 1 1 0 0 1 1 0 0 1 1 0 0 0 1 0 1 1 0 0 0 1 1 1 0 0 0 0 1 0 1
## Levels: 0 1

library(ltm)
bi<-biserial.cor(df$Rainy...Total.Accidents,df$hail_accident_occurr)
print(bi)

## [1] -0.1429725

Solution 5

Similar to in Problem 4, create a binary column to represent whether a dust storm accident has occurred in a state (1 = occurred, 0 = not occurred). Convert the two columns into a contingency table.

Calculate the phi coefficient of the two tables. (Hint: use the psych package).
```

```
df$dust_storm_occurr<-NA
df$dust_storm_occurr<-factor(ifelse(df$Dust.Storm...Total.Accidents==0,0,1))
cat("Dust.Storm...Total.Accidents:\n",df$Dust.Storm...Total.Accidents)

## Dust.Storm...Total.Accidents:
## 42 0 36 582 132 0 302 243 0 0 247 580 0 429 474 3 0 0 0 249 69 8 0 0 4 0 9 1521 196 0 0 0 0 122 0 3

#dust_storm_occurr
print(df$dust_storm_occurr)

## [1] 1 0 1 1 1 0 1 1 0 0 1 1 0 0 1 1 1 0 0 0 1 1 1 0 0 0 1 1 1 0 0 0 0 1 0 1
## Levels: 0 1

#Contingency Table For Hail accident occurrence
table1 = table(df$hail_accident_occurr)
print.table(table1)

##
## 0 1
## 19 17

#Contingency Table for Dust Storm Occurrence
table2=table(df$dust_storm_occurr)
print.table(table2)

##
## 0 1
## 16 20

#Contingency Table
library(psych)
conTable = table(df$hail_accident_occurr,df$dust_storm_occurr)
print(conTable)

##
## 0 1
## 0 14 5
## 1 2 15

##Phi Correlation Coefficient
phi(conTable)

## [1] 0.62
```