

1 Exploratory Data Analysis using R

Panagiotis Lamprakis^a

^aPhD Candidate, AI|LS Laboratory, National Technical University of Athens, panoslamprakis@mail.ntua.gr, 03003142

Tetuan is a city located in the north of Morocco which occupies an area of around 10375 km² and its population is about 550.374 inhabitants, according to the last Census of 2014, and is increasing rapidly, approximately 1.96% annually. Since it is located along the Mediterranean Sea, its weather is mild and rainy in the winter, hot and dry during the summer months. The power consumption data was collected from Supervisory Control and Data Acquisition System (SCADA) of Amendis which is a public service operator and in charge of the distribution of drinking water and electricity since 2002. The distribution network is powered by 3 source stations, namely: Quads, Smir and Boussafou. The three stations power 3 different areas of the city (zones).

1 Dataset description

We will use the "Tetuan City power consumption.csv" dataset, located in [course's website](#) (click [here](#) to directly download it). The dataset contains seven columns and has 52416 rows of data. The R code for this analysis can be found in [github](#) (repo is private, please ask for access).

Column Name	Description
DateTime	Time window of ten minutes: 01/01/2017 00:00 to 30/12/2017 23:50
Temperature	Weather Temperature in °C
Humidity	Weather Humidity in %
Wind Speed	Wind Speed in km/h
Zone 1 Power Consumption	in KiloWatts (KW)
Zone 2 Power Consumption	in KiloWatts (KW)
Zone 3 Power Consumption	in KiloWatts (KW)

2 Reading the file and transforming the data

In the following code snippet, we read the data in a data table structure. In order to achieve better performance, we read the first 100 rows, determine the type of the columns from those 100 rows and then we read the entire dataset. We make sure that the dataset is in data table representation.

```
1 library(data.table)
2 # Read dataset sample to determine the data table column classes
3 dataSample <- read.table(DATASET_PATH, nrows=100, sep=";", header = TRUE)
4
5 # Read the data by specifying each column type and make sure it is a data
  table
6 classes <- sapply(dataSample, class)
7 data <- read.table(DATASET_PATH, colClasses = classes, sep=";", header =
  TRUE)
8 setDT(data, keep.rownames=T)
```

For the EDA, we want to know which days of the year are bank holidays. For this reason, we will use data from (holidayapi.com 2022). We manually produce a csv file, which we read and create a new column *isBankHoliday* in the original data table:

```
1 officialBankHolidays <- read.csv(BANK_HOLIDAYS_PATH_OF, header = F)
```

We will create new features (*isBankHoliday*) and extract features from current data (*Date*, *DoW*, *Month*, *DoM*, *Hour*, *isWeekend*, *season*).

```
1 dateTimeSplitList <- strsplit(data$DateTime, " ")
2
3 # create new date and time columns.
4 data[, "DateStr"] = sapply(dateTimeSplitList, "[[", 1)
5 data[, "Time"] = sapply(dateTimeSplitList, "[[", 2)
6
7 # Create a column to see if the date was a bank holiday
8 data[, "isBankHoliday"] = rep(1, nrow(data))
9 data$isBankHoliday = ifelse(data$DateStr %in% officialBankHolidays[1]$V1,
  1, 0) # kept as numeric for the correlation map
10
11 # Convert Date column to date type, so to extract the month, the day of
  the week and the day of the month
12 data$Date <- as.Date(data$DateStr, format = "%m/%d/%Y")
13
14
15 data[, "DoW"] = wday(data$Date, week_start=1) # kept as numeric for the
  correlation map
16 data[, "Month"] = as.numeric(format(data$Date, "%m")) # kept as numeric
  for the correlation map
17 data[, "DoM"] = as.numeric(format(data$Date, "%d")) # kept as numeric for
  the correlation map
18 data[, "Hour"] <- as.numeric(format(as_datetime(data$DateTime, format = "%
  m/%d/%Y %H:%M"), "%H"))
```

```

19
20 # Produce info if the day was during the weekend or it was a workday
21 data[, "isWeekend"] = rep(1, nrow(data))
22 data$isWeekend = ifelse(data$DoW %in% c(6, 7), 1, 0) # 6: Sat, 7: Sun
23
24 # Determine the season of the year for the data sample
25 data[, "season"] = sapply(data$Month, getSeason) # convert to numeric for
    the correlation map
26 data[, "seasonF"] = sapply(data$season, encodeSeason)

```

Moreover, in order to produce some other plots we will expand our dataset from its current form. As of now, the dataset contains three columns with the consumption per zone on each column. The transformation we will put in place is to have one column with holding the power consumption information and another one showing to which city zone the consumption refers. The drawback of this approach is that we will duplicate the information we have and that we will end up increasing the memory used by approximately 200%. To achieve this result, we followed the following process:

```

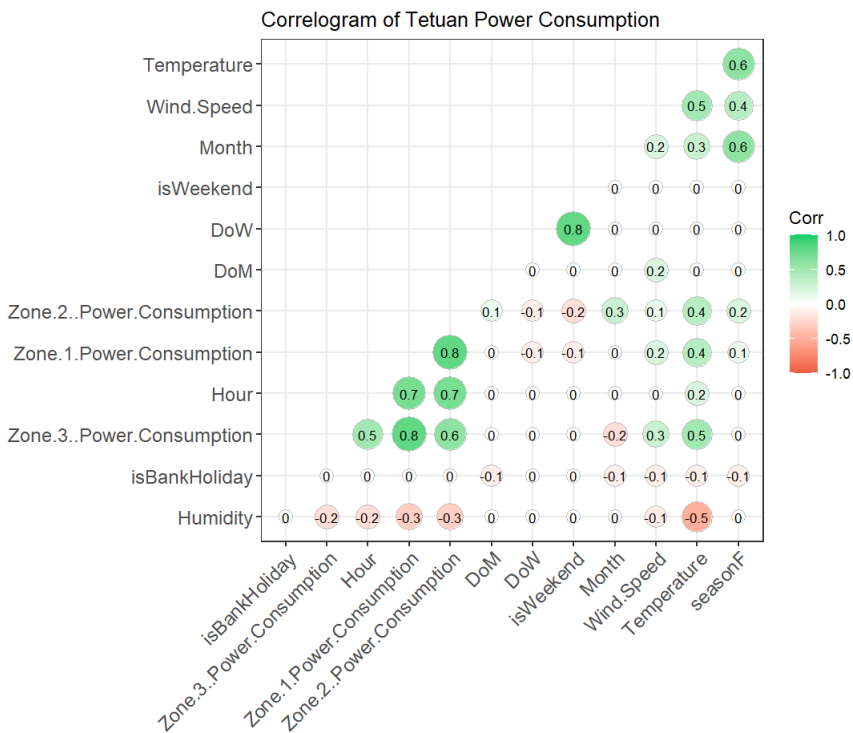
1 dataFiltered = data[, -c("rn", "DateTime", "DateStr", "Time", "season")]
2 dataFiltered$DateF <- as.numeric(factor(data$Date))
3 dataFiltered$TimeF <- as.numeric(factor(data$Time))
4
5 # Expand dataset and have the consumption in one column and another column
    with the zone
6 # This will increase the dataset size by 200% but will help us with the
    plotting
7 dataZone1 <- dataFiltered[, -c("Zone.2..Power.Consumption", "Zone.3..Power
    .Consumption")]
8 setnames(dataZone1, "Zone.1.Power.Consumption", "ZoneConsumption")
9 dataZone1$Zone <- rep(1, nrow(dataZone1))
10
11 dataZone2 <- dataFiltered[, -c("Zone.1.Power.Consumption", "Zone.3..Power.
    Consumption")]
12 setnames(dataZone2, "Zone.2..Power.Consumption", "ZoneConsumption")
13 dataZone2$Zone <- rep(2, nrow(dataZone2))
14
15 dataZone3 <- dataFiltered[, -c("Zone.1.Power.Consumption", "Zone.2..Power.
    Consumption")]
16 setnames(dataZone3, "Zone.3..Power.Consumption", "ZoneConsumption")
17 dataZone3$Zone <- rep(3, nrow(dataZone3))
18
19 dataZoneLevelWithDate <- rbind(dataZone1, dataZone2, dataZone3)
20 dataZoneLevelWithDate$DateF <- as.numeric(factor(dataZoneLevel$Date))
21 dataZoneLevelWithDate$TimeF <- as.numeric(factor(dataZoneLevel$Time))
22
23 dataZoneLevel <- dataZoneLevelWithDate[, -c("Date")]
24 dataZoneLevel[, "ZoneName"] <- factor(sapply(dataZoneLevel$Zone,
    decodeZone))

```

3 Data Analysis

3.1 Correlation Analysis

To begin with the analysis, we will first examine the correlation between the dataset features. From the produced correlation map, it is easy to observe that features as *DoM*, *isBankHoliday* are not related to the power consumption, since their correlation score is zero. These are features that we produced/extracted in the previous step. Features such as *DoW*, *isWeekend*, *season* have weak correlation with the power consumption (correlation coefficient's absolute value ≤ 0.2), while features such as *Wind.Speed*, *Month*, *Humidity* have moderate correlation with the power consumption. Finally, *Temperature* and *Hour* seem to be related to the power consumption (correlation score ≥ 0.4). It's also easy to observe that the power consumption in one zone is highly related to the power consumption in the other two zones (correlation score ≥ 0.6).

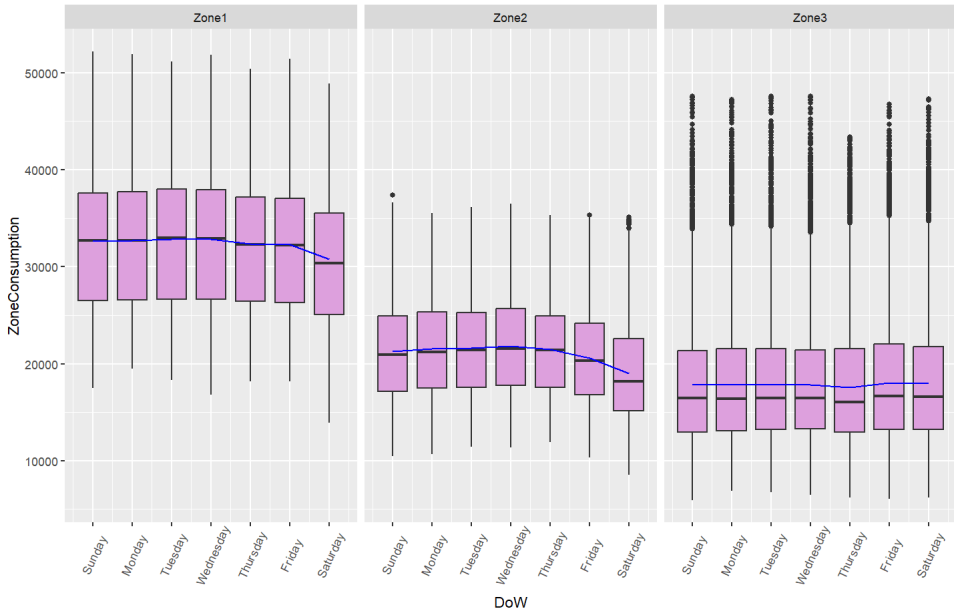


```

1 dataOnlyNumeric <- data[, -c("rn", "DateStr", "DateTime", "Date", "
2 Time", "season")]
3
4 corr <- round(cor(dataOnlyNumeric), 1)
5 ggcorrplot(corr, hc.order = TRUE,
6 type = "lower",
7 lab = TRUE,
8 lab_size = 3,
9 method="circle",
10 colors = c("tomato2", "white", "springgreen3"),
11 title="Correlogram of Tetuan Power Consumption",
12 ggtheme=theme_bw)

```

We establish our argument that *Dow* is not related to the power consumption, with the following plot. We see that the average power consumption per day is almost constant in each city zone. The code used to produce the plot follows it.

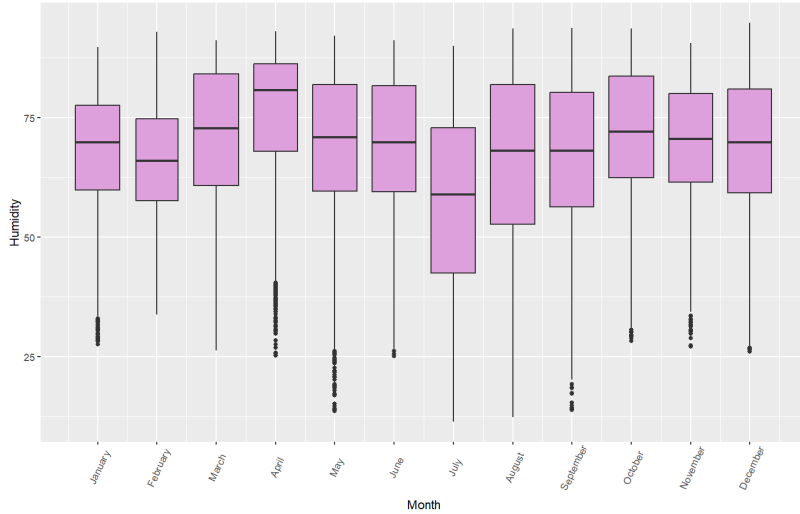


```

1 # Facet per zone per DoW consumption
2 facetBoxPlotBreaksDoW <- seq(1, 7)
3 dataZoneLevel[, "DoWName"] <- sapply(dataZoneLevel$DoW, decodeDoW)
4 gg <- ggplot(dataZoneLevel, aes (x = DoW , y = ZoneConsumption )) +
5   geom_boxplot(varwidth=T, fill="plum", aes(group=DoWName)) +
6   stat_summary(aes(group=ZoneName), fun=mean, geom="line", color="blue
7   ") +
8   facet_grid(~ZoneName) +
9   scale_x_continuous(labels = unique(dataZoneLevel$DoWName), breaks
10  = facetBoxPlotBreaksDoW) +
11   theme(axis.text.x = element_text (angle=65, vjust =0.6))
12 plot(gg)

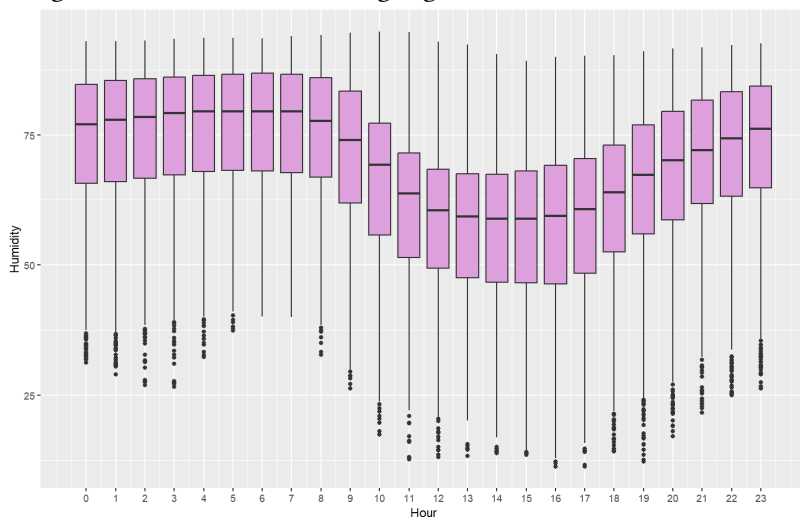
```

Another interesting insight we get from the correlation map, is that humidity is not related to month. From the plot, we easily see that the average humidity level is almost constant across months, with the exception of April and July.



```
1 gg <- ggplot(dataZoneLevel, aes(x = Month, y = Humidity)) +  
2   geom_boxplot(varwidth=T, fill="plum", aes(group=MonthName)) +  
3   scale_x_continuous(labels = unique(dataZoneLevel$MonthName),  
4     breaks = facetBoxPlotBreaksMonth) +  
5   theme(axis.text.x = element_text(angle=65, vjust = 0.6))  
6   plot(gg)
```

However, humidity is related to the hour of the day. Humidity levels decrease during sunlight and are increased during night hours.

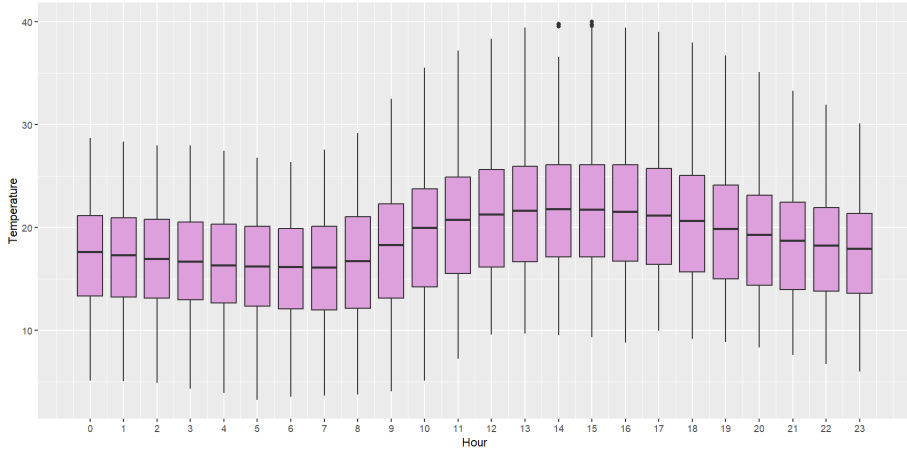


```

1 gg <- ggplot(dataZoneLevel, aes (x = Hour , y = Humidity )) +
2   geom_boxplot(varwidth=T, fill="plum", aes(group=Hour)) +
3   scale_x_continuous(labels = facetBoxPlotLabelsHour, breaks =
4     facetBoxPlotBreaksHour)
5   plot(gg)

```

The below box plot visualizes that Temperature is related to hour (correlation coefficient = 0.2). We observed that temperature increases during afternoon hours.

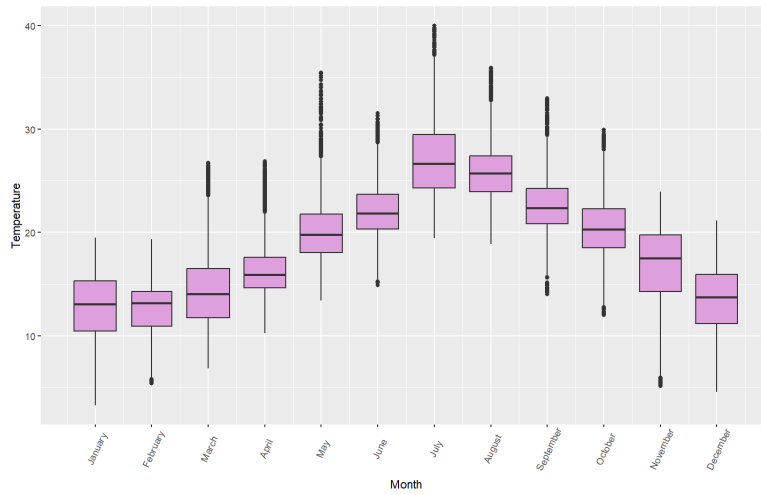


```

1 gg <- ggplot(dataZoneLevel, aes (x = Hour , y = Temperature )) +
2   geom_boxplot(varwidth=T, fill="plum", aes(group=Hour)) +
3   scale_x_continuous(labels = facetBoxPlotLabelsHour, breaks =
4     facetBoxPlotBreaksHour)
5   plot(gg)

```

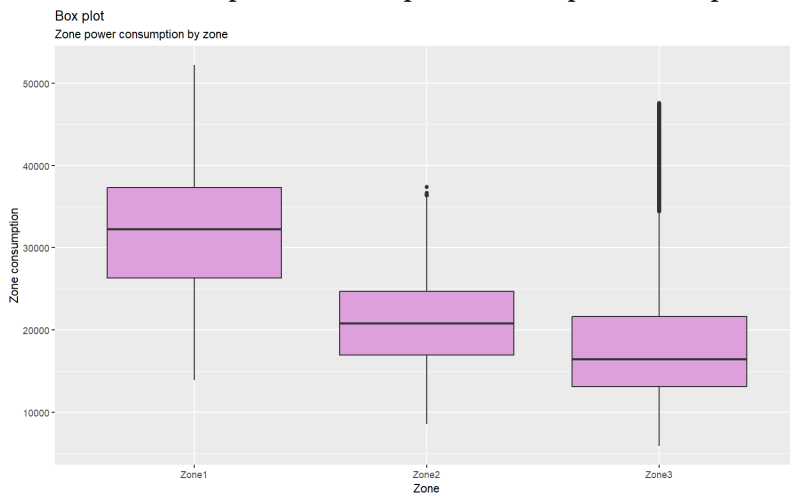
More interesting is the box plot for the Temperature compared to Month. May to October are the most hot months of the year for Tetuan, with July to be the hottest among all.



```
1 gg <- ggplot(dataZoneLevel, aes (x = Month , y = Temperature )) +
2   geom_boxplot(varwidth=T, fill="plum", aes(group=MonthName)) +
3   scale_x_continuous(labels = unique(dataZoneLevel$MonthName),
4     breaks = facetBoxPlotBreaksMonth) +
5   theme(axis.text.x = element_text (angle=65, vjust =0.6))
6   plot (gg)
```

3.2 Power Consumption Analysis

From now and on, our feature of interest will be the power consumption. In order to study the distribution of power consumption we will plot its Box plot.




```

1 g <- ggplot(dataZoneLevel, aes(Zone, ZoneConsumption)) +
2   geom_boxplot(varwidth=T, fill="plum") +
3   labs(title="Box plot",
4        subtitle="Zone power consumption by zone",
5        x="Zone",
6        y="Zone consumption")
7 plot(g)
8 ld <- layer_data(g) # retrieve box plot data

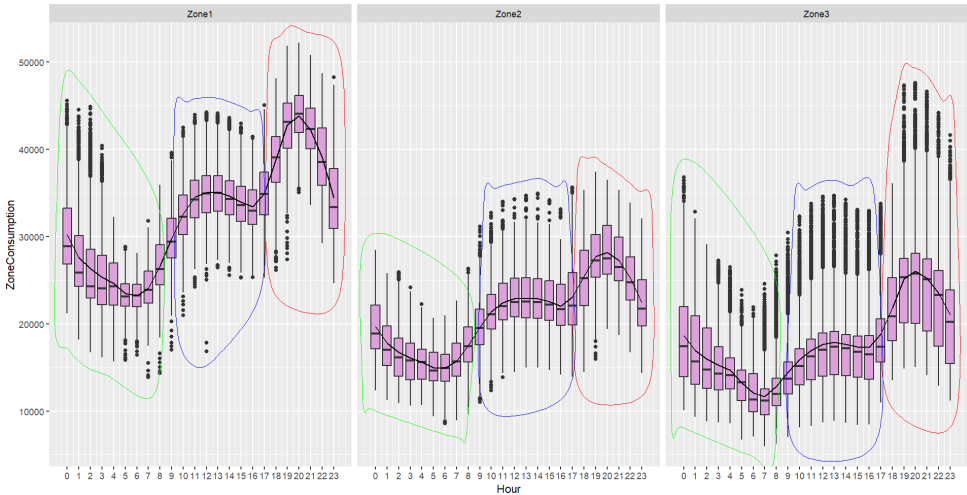
```

From the box plot itself, it is easy to see that the highest average power consumption is in Zone 1 while Zone 3 has the lowest consumption. This leads us to think that Zone 1 may be a rich neighborhood in the city, while Zone 3 is the least rich. Table 1 describes the box plot in precision.

Table 1: Box plot data

	Mean	25%	75%	Whisker(High)	Whisker(Low)
Zone 1	32265.92	26310.67	37309.02	52204.40	13895.696
Zone 2	20823.17	16980.77	24713.72	36201.48	8560.081
Zone 3	16415.12	13129.33	21624.10	34361.17	5935.174

We plot the power consumption per zone per hour. We split semantically the plot into three subsections. We assume that people sleep in time window 00:00 - 8:00 (*green circle*), they work in time window 9:00 - 17:00 (*blue circle*) and they rest at home in time window 18:00 - 23:00 (*red circle*).

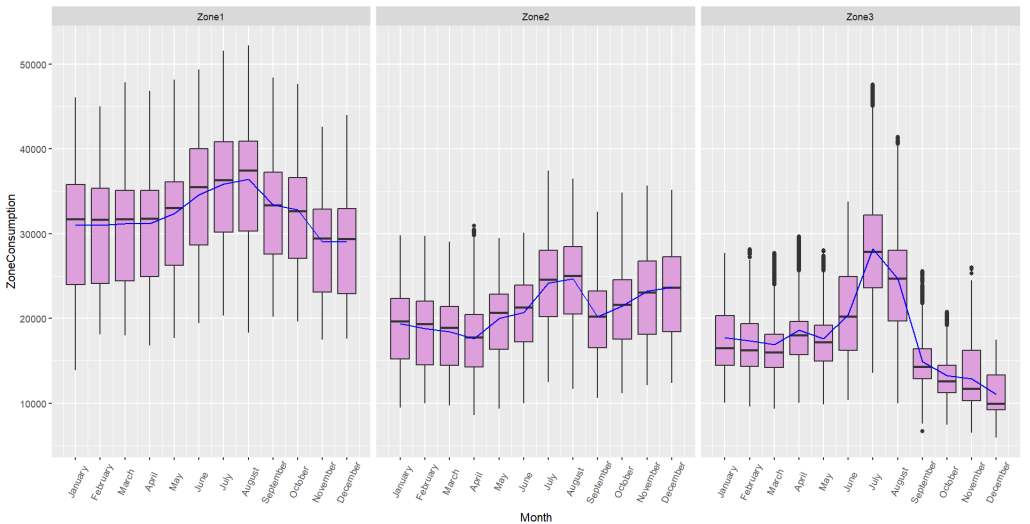


```

1 # Facet per zone per Hour consumption
2 facetBoxPlotLabelsHour <- seq(0, 23)
3 facetBoxPlotBreaksHour <- seq(0, 23)
4 gg <- ggplot(dataZoneLevel, aes (x = Hour , y = ZoneConsumption )) +
5   geom_boxplot(varwidth=T, fill="plum", aes(group=Hour)) +
6   stat_summary(aes(group=ZoneName), fun=mean, geom="line") +
7   facet_grid(~ZoneName) +
8   scale_x_continuous(labels = facetBoxPlotLabelsHour, breaks =
9     facetBoxPlotBreaksHour) +
10   geom_encircle(data = dataZoneLevel[Hour >= 18 & Hour <= 23 ],
11     aes(x = Hour , y = ZoneConsumption), colour="red") +
12   geom_encircle(data = dataZoneLevel[Hour > 9 & Hour < 17 ], aes(x
13     = Hour , y = ZoneConsumption), colour="blue") +
14   geom_encircle(data = dataZoneLevel[Hour >= 0 & Hour < 8 ], aes(x
15     = Hour , y = ZoneConsumption), colour="green")
16 plot(gg)

```

In the 3.1 section, we saw that in time window 18:00-10:00 (*green* and *red* circles), temperature is lower and humidity is higher. In the ZoneConsumption vs Hour plot, we see that in all zones, the power consumption is higher in the time period 18:00-23:00. This observation makes sense, since people leave their jobs (decentralize heating/cooling from corporate environments to home, shops, malls, etc) in the evening and they need to do household chores (cooking, laundry, etc). Also, electricity could be used for dehumidifier appliances at home. Power consumption drops in time window defined by *green* circle, since people sleep at that hours.



```

1 # Facet per zone per Month consumption
2 facetBoxPlotBreaksMonth <- seq(1, 12)
3 dataZoneLevel[, "MonthName"] <- sapply(dataZoneLevel$Month,
4   decodeMonth)

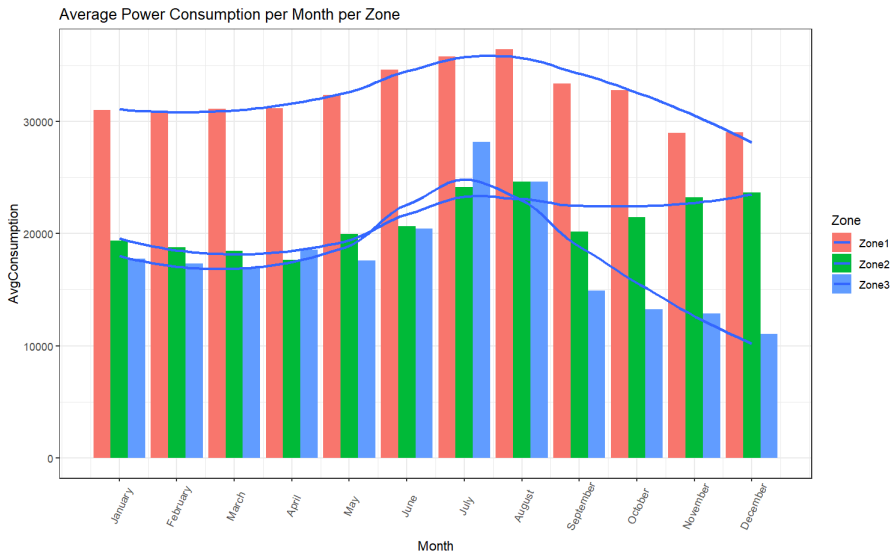
```

```

4 gg <- ggplot(dataZoneLevel, aes (x = Month , y = ZoneConsumption )) +
5   geom_boxplot (varwidth=T, fill="plum", aes (group=Month)) +
6   stat_summary (aes (group=ZoneName), fun=mean, geom="line", color="
  blue") +
7   facet_grid (~ZoneName) +
8   scale_x_continuous (labels = unique (dataZoneLevel$MonthName),
9     breaks = facetBoxPlotBreaksMonth) +
10  theme (axis.text.x = element_text (angle=65, vjust =0.6))
  plot (gg)

```

We also saw that May to October are the hottest months of the year for Tetuan, while humidity is not related to the month (humidity levels are almost constant across months). From the below plot, we see that the power consumption during summertime is higher for all three city zones (positive correlation ratio).



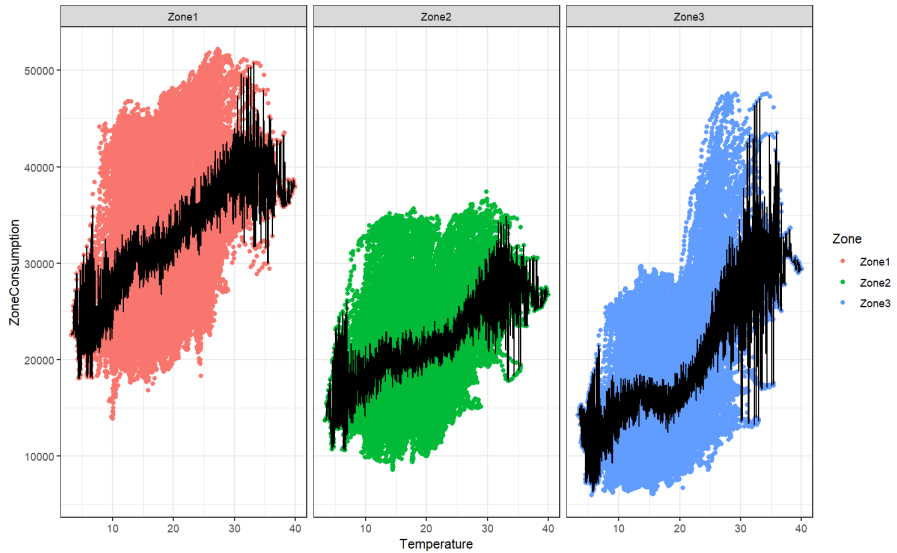
```

1 # Average consumption PER MONTH PER ZONE
2 avgConPerMonPerZoneData = dataZoneLevel[, c("ZoneName", "Month", "
  ZoneConsumption")]
3 avgConPerMonPerZoneData[, "MonthName"] = sapply(
4   avgConPerMonPerZoneData$Month, decodeMonth)
5
6 avgConPerMonPerZoneData = avgConPerMonPerZoneData[, .(mean(
7   ZoneConsumption)), by = .(Month, ZoneName)]
8 setnames (avgConPerMonPerZoneData, "V1", "AvgConsumption")
9
10 # It is aggregated per month per zone, so first 12 rows
11 # contain all discrete months in the proper order
12 months <- sapply (avgConPerMonPerZoneData$Month, decodeMonth)
13
14 gg <- ggplot (avgConPerMonPerZoneData, aes (x = Month, y =
15   AvgConsumption , fill = ZoneName )) +
16   geom_bar (position = "dodge", stat = "identity") +

```

```
14     geom_smooth(method="loess", se=F) +  
15     scale_x_continuous(breaks=seq(1, 12), labels=months[1:12]) +  
16     theme(axis.text.x = element_text(angle=65, vjust = 0.6)) +  
17     labs(title="Average Power Consumption per Month per Zone")  
18     plot(gg)
```

In principle, we see that the power consumption increases as temperature increases in all zones.



```
1  # Temperature Vs Consumption Facet per Zone  
2  gg <- ggplot(dataZoneLevel, aes (x = Temperature , y = ZoneConsumption  
3  )) +  
4  geom_point(aes (color=ZoneName)) +  
5  stat_summary(aes (group=ZoneName), fun=mean, geom="line") +  
6  facet_grid(~Zone)  
7  plot(gg)
```

4 Appendix

Auxiliary functions declared in scope of the analysis:

```

1 # Declare auxiliary functions
2 getSeason <- function(month) {
3   if (month %in% c(12, 1, 2)) {
4     "winter"
5   } else if (month %in% c(3, 4, 5)) {
6     "spring"
7   } else if (month %in% c(6, 7, 8)) {
8     "summer"
9   } else if (month %in% c(9, 10, 11)) {
10    "autumn"
11  } else {
12    "unknown"
13  }
14 }
15
16 encodeSeason <- function(season) {
17   if (season == "winter") {
18     1
19   } else if (season == "spring") {
20     2
21   } else if (season == "summer") {
22     3
23   } else if (season == "autumn") {
24     4
25   } else {
26     5
27   }
28 }
29
30 decodeMonth <- function(month) {
31   if (month == 1) {
32     "January"
33   } else if (month == 2) {
34     "February"
35   } else if (month == 3) {
36     "March"
37   } else if (month == 4) {
38     "April"
39   } else if (month == 5) {
40     "May"
41   } else if (month == 6) {
42     "June"
43   } else if (month == 7) {
44     "July"
45   } else if (month == 8) {
46     "August"
47   } else if (month == 9) {
48     "September"
49   } else if (month == 10) {
50     "October"
51   } else if (month == 11) {

```

```
52   "November"
53 } else {
54   "December"
55 }
56 }
57
58 decodeZone <- function(zoneId) {
59   if (zoneId == 1) {
60     "Zone1"
61   } else if (zoneId == 2) {
62     "Zone2"
63   } else {
64     "Zone3"
65   }
66 }
67
68 encodeZone <- function(zoneName) {
69   if (zoneName == "Zone1") {
70     1
71   } else if (zoneName == "Zone2") {
72     2
73   } else {
74     3
75   }
76 }
```

References

holidayapi.com. 2022. <https://holidayapi.com/countries/ma-01/2017> (27 December, 2022).