

PD1_EksploracjaDanych

March 30, 2021

1 Eksploracja zbioru danych

1.1 Wczytanie danych

Aby ułatwić analizę zamienimy nazwy dni i miesięcy na wartości liczbowe.

```
[2]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from pandas_profiling import ProfileReport
sns.set()

import calendar

df_fire = pd.read_csv('forest_fires_dataset.csv')
df_fire_att = pd.read_csv('attributes_forest_fires.csv')

[3]: days_name = [each_string.lower() for each_string in calendar.day_abbr]
months_name = [each_string.lower() for each_string in calendar.month_abbr]

def dayToNum(day):
    return days_name.index(day) + 1
def monthToNum(month):
    return months_name.index(month)

df_fire.day = df_fire.day.apply(dayToNum)
df_fire.month = df_fire.month.apply(monthToNum)
```

Naszym zadaniem jest przeprowadzenie eksploracji danych zbioru dotyczącego pożarów w północno-wschodnim regionie Portugalii. Na początku rzucmy okiem na nasze dane

```
[4]: df_fire.head()
```

	X	Y	month	day	FFMC	DMC	DC	ISI	temp	RH	wind	rain	area
0	7	5	3	5	86.2	26.2	94.3	5.1	8.2	51.0	6.7	0.0	0.0
1	7	4	10	2	90.6	35.4	669.1	6.7	18.0	33.0	0.9	0.0	0.0
2	7	4	10	6	90.6	43.7	686.9	6.7	14.6	33.0	1.3	0.0	0.0
3	8	6	3	5	91.7	33.3	77.5	9.0	8.3	97.0	4.0	0.2	0.0
4	8	6	3	7	89.3	51.3	102.2	9.6	11.4	99.0	1.8	0.0	0.0

Opisy poszczególnych kolumn znajdują się w tabeli poniżej. Dokładne objaśnienie skrótów FFMC, DMC, DC, ISI znajduje się pod tym linkiem: <https://www.nwgc.gov/publications/pms437/cffdrs/fire-weather-index-system>

[5]: `df_fire_att`

```
[5]:      name      type      description
0      X  integer  x-axis spatial coordinate within the Montesinh...
1      Y  integer  y-axis spatial coordinate within the Montesinh...
2  month  string      month of the year: 'jan' to 'dec'
3    day  string      day of the week: 'mon' to 'sun'
4    FFMC  float      FFMC index from the FWI system: 18.7 to 96.20
5     DMC  float      DMC index from the FWI system: 1.1 to 291.3
6     DC   float      DC index from the FWI system: 7.9 to 860.6
7     ISI  float      ISI index from the FWI system: 0.0 to 56.10
8    temp  float      temperature in Celsius degrees: 2.2 to 33.30
9     RH   float      relative humidity in %: 15.0 to 100
10   wind  float      wind speed in km/h: 0.40 to 9.40
11   rain  float      outside rain in mm/m2 : 0.0 to 6.4
12   area  float      the burned area of the forest (in ha): 0.00 to...
```

[6]: `df_fire.describe()`

```
[6]:           X           Y      month      day      FFMC      DMC \
count  517.000000  517.000000  517.000000  517.000000  517.000000  517.000000
mean     4.669246    4.299807    7.475822    4.259188   90.644681  110.872340
std     2.313778    1.229900    2.275990    2.072929   5.520111   64.046482
min      1.000000    2.000000    1.000000    1.000000   18.700000    1.100000
25%      3.000000    4.000000    7.000000    2.000000   90.200000   68.600000
50%      4.000000    4.000000    8.000000    5.000000   91.600000  108.300000
75%      7.000000    5.000000    9.000000    6.000000   92.900000  142.400000
max      9.000000    9.000000   12.000000    7.000000   96.200000  291.300000

           DC           ISI      temp      RH      wind      rain \
count  517.000000  517.000000  517.000000  517.000000  517.000000  517.000000
mean   547.940039    9.021663   18.889168   44.288201    4.017602    0.021663
std   248.066192    4.559477    5.806625   16.317469    1.791653    0.295959
min     7.900000    0.000000    2.200000   15.000000    0.400000    0.000000
25%   437.700000    6.500000   15.500000   33.000000    2.700000    0.000000
50%   664.200000    8.400000   19.300000   42.000000    4.000000    0.000000
75%   713.900000   10.800000   22.800000   53.000000    4.900000    0.000000
max   860.600000   56.100000   33.300000  100.000000    9.400000    6.400000

           area
count  517.000000
mean    12.847292
std    63.655818
min      0.000000
25%      0.000000
```

```

50%      0.520000
75%      6.570000
max     1090.840000

```

Z powyższej tabeli możemy zauważyć, że kolumny rain i area są głównie wypełnione zerami.

```

[33]: print('Liczba rekordów, gdzie w kolumnie area nie ma zera: ', df_fire[df_fire.
        →area != 0 ].shape[0])
       print('Liczba rekordów, gdzie w kolumnie rain nie ma zera: ', df_fire[df_fire.
        →rain != 0 ].shape[0])

```

```
Liczba rekordów, gdzie w kolumnie area nie ma zera: 270
```

```
Liczba rekordów, gdzie w kolumnie rain nie ma zera: 8
```

Powysze obliczenia potwierdzają to. Mniej więcej połowa rekordów zawiera zero w kolumnie area i tylko 8 z 517 zawiera coś innego niż zero w kolumnie rain.

```
[6]: df_fire.info()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 517 entries, 0 to 516
Data columns (total 13 columns):
 #   Column  Non-Null Count  Dtype
---  -
 0    X      517 non-null    int64
 1    Y      517 non-null    int64
 2  month  517 non-null    int64
 3   day   517 non-null    int64
 4  FFMC   517 non-null    float64
 5  DMC    517 non-null    float64
 6   DC    517 non-null    float64
 7   ISI    517 non-null    float64
 8  temp   517 non-null    float64
 9   RH    517 non-null    float64
10  wind   517 non-null    float64
11  rain   517 non-null    float64
12  area   517 non-null    float64
dtypes: float64(9), int64(4)
memory usage: 52.6 KB

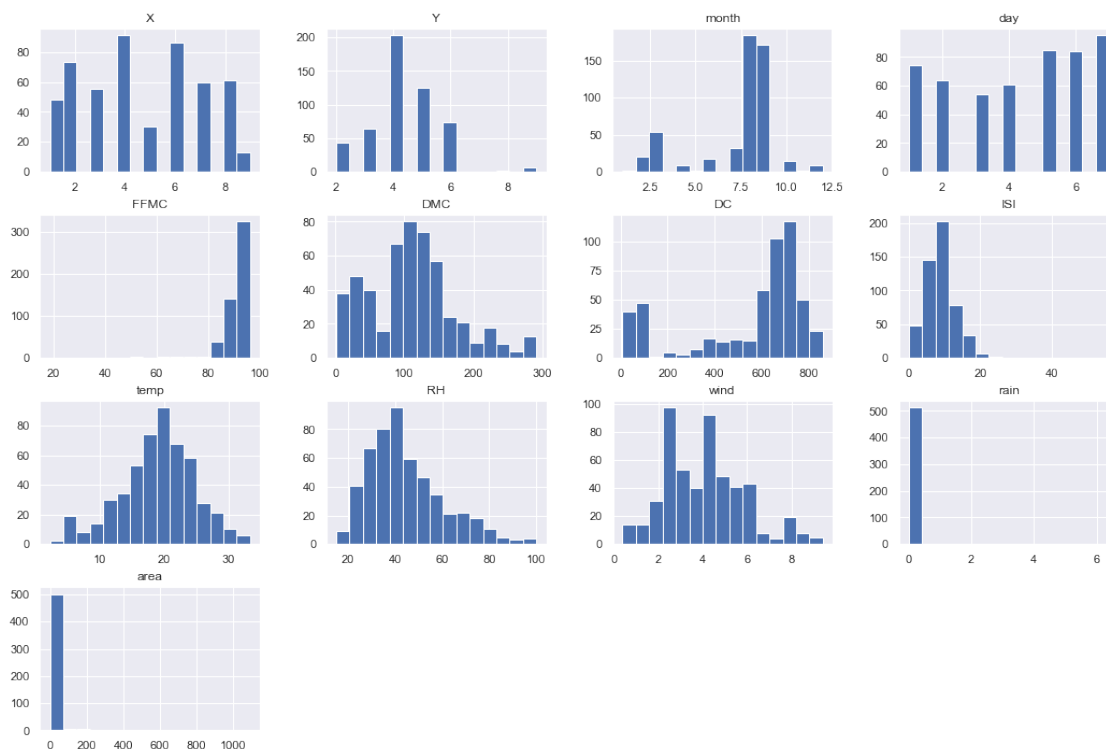
```

Dzięki funkcji info() wiemy, że w naszej ramce danych nie ma wartości null

```

[7]: df_fire.hist(bins = 15, figsize=(18, 12))
      plt.show()

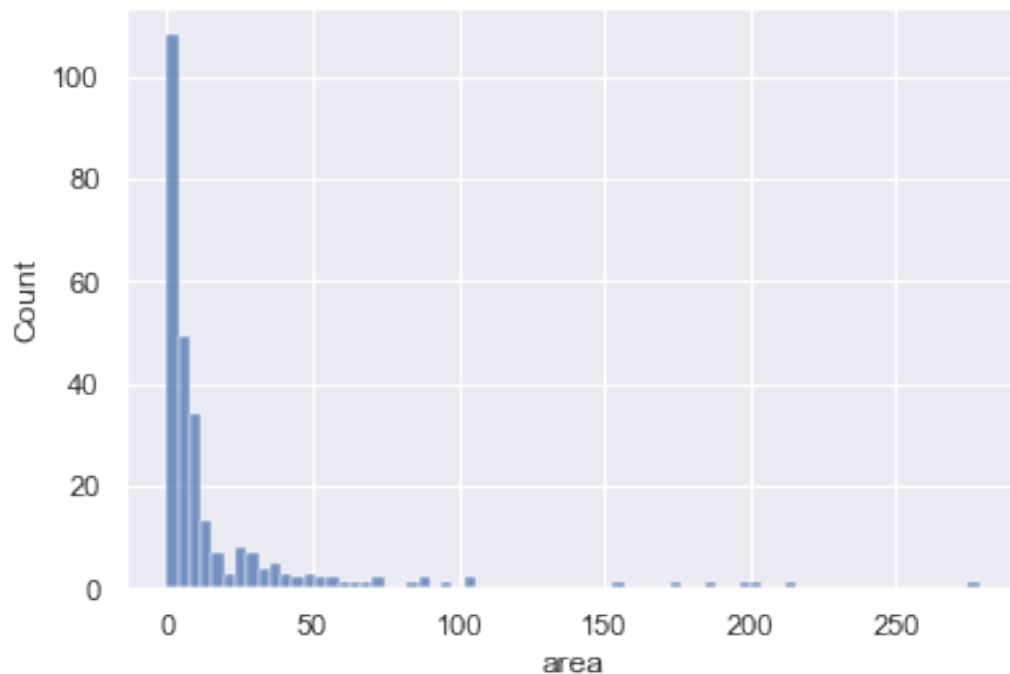
```



Powysze histogramy obrazuj rozklady poszczególnych kolumn. Wida w nich wyranie to co rzucio nam si w oczy przy uyciu funkcji describe(), tzn kolumny area i rain s wypenione prawie cakowicie zerami. Dodatkowo moemy zauway, e rozkad temperatury jest mocno zbliiony do rozkadu normalnego. Oprócz tego moemy zauway, e cecha dotyczca uczszczania do lasu przez ludzi ma rozkad skony prawostronny.

```
[8]: sns.histplot(df_fire.area[(df_fire.area < 300) & (df_fire.area != 0)])
```

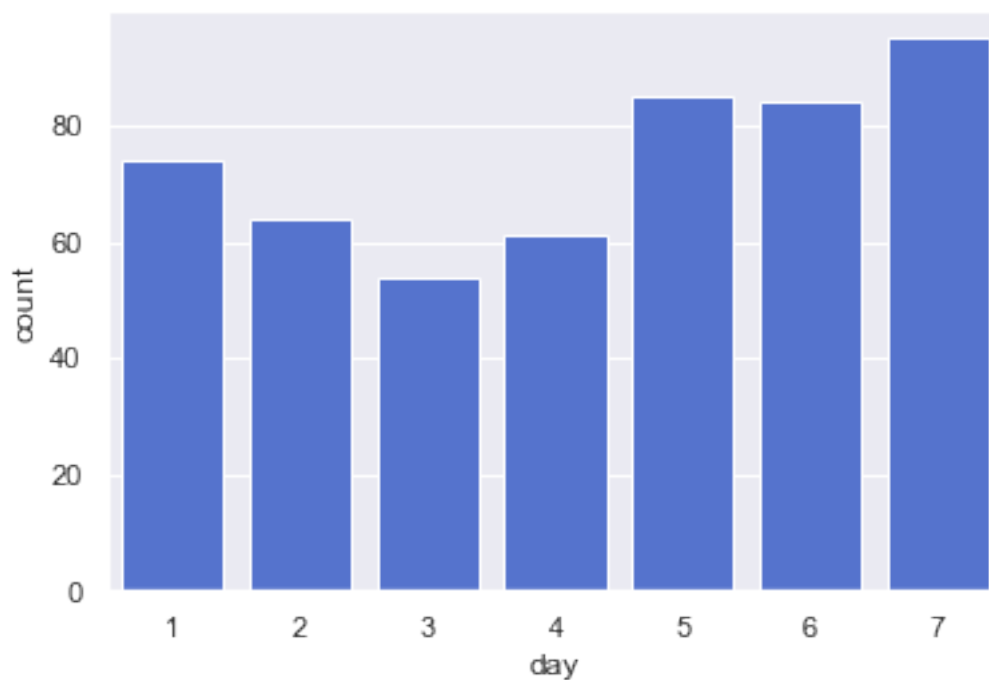
```
[8]: <AxesSubplot:xlabel='area', ylabel='Count'>
```



Po odrzuceniu skrajnie duzych wartosci i wartosci zerowych widzimy, e rozkad obszaru objtego poarem ma ju rozkad podobny do $1/x$.

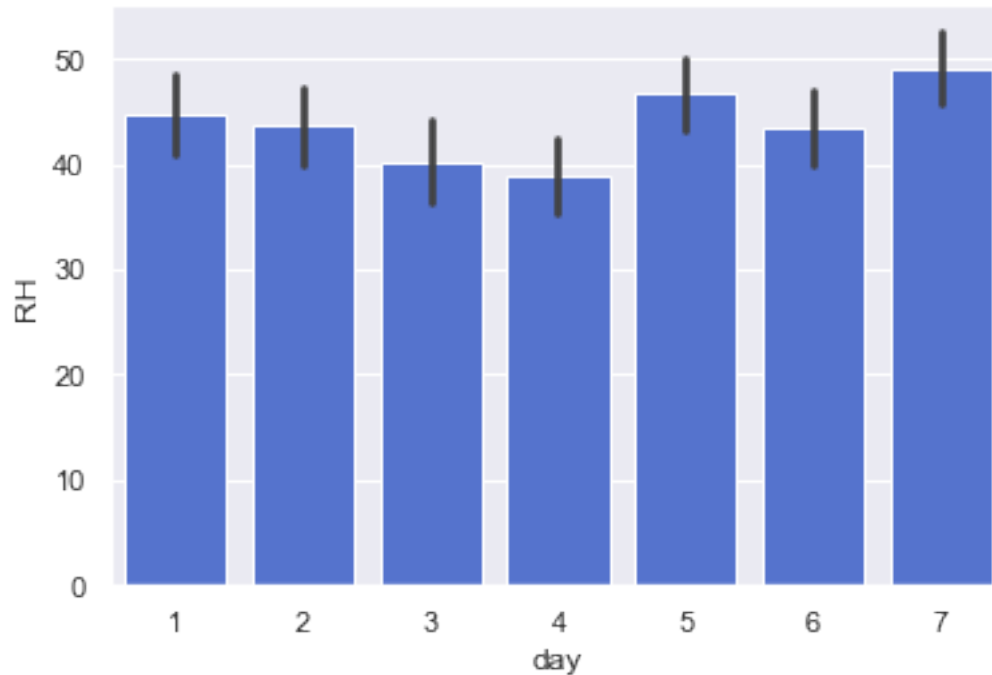
```
[9]: sns.countplot(data = df_fire, x = 'day', color="RoyalBlue")
```

```
[9]: <AxesSubplot:xlabel='day', ylabel='count'>
```



```
[10]: sns.barplot(data = df_fire, x = 'day', y = 'RH', color="RoyalBlue" )
```

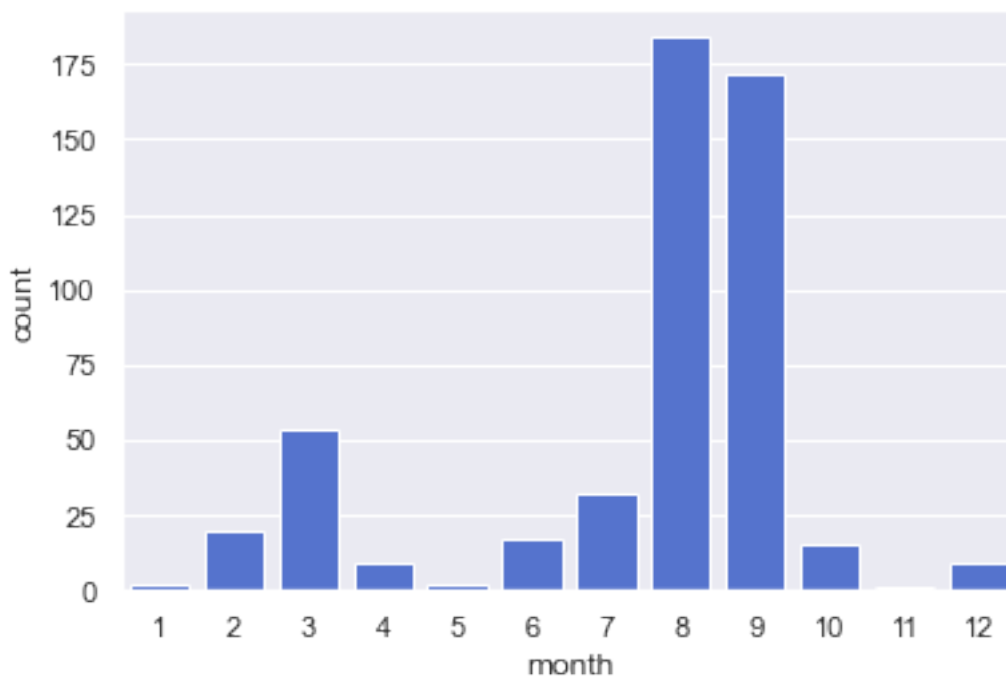
```
[10]: <AxesSubplot:xlabel='day', ylabel='RH'>
```



Z dwóch powyższych wykresów widać to, co byśmy widzieli wcześniej. Najwięcej poarów miało miejsce w piątek, sobotę i niedzielę. Jest to do logicznego zjawiska, gdy spojrzymy na drugi wykres. Widać z niego, że najwięcej ludzi przychodzi do lasu w weekendy. To w weekendy ludzie mogą być zatem przyczyną wielu poarów.

```
[11]: sns.countplot(data = df_fire, x = 'month', color="RoyalBlue")
```

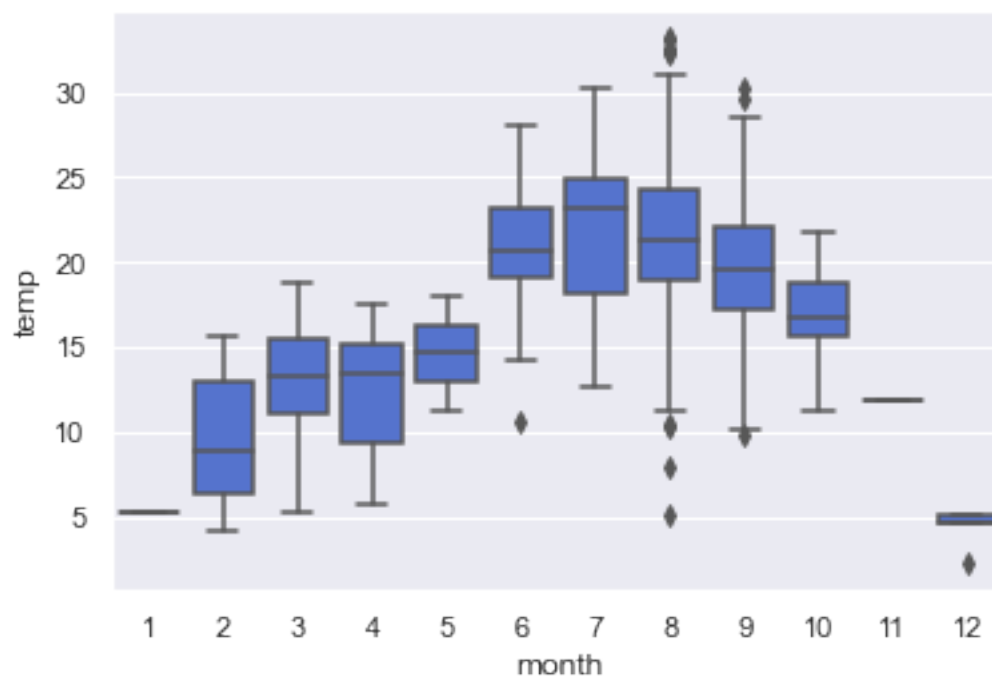
```
[11]: <AxesSubplot:xlabel='month', ylabel='count'>
```



Ten wykres pokazuje natomiast ilo parów w różnych miesiącach. Szczególnie dużo parów jest w sierpniu i wrześniu. Zaskakująco dużo parów jest również w marcu.

[25]: `sns.boxplot(x="month", y="temp", data=df_fire, color="RoyalBlue")`

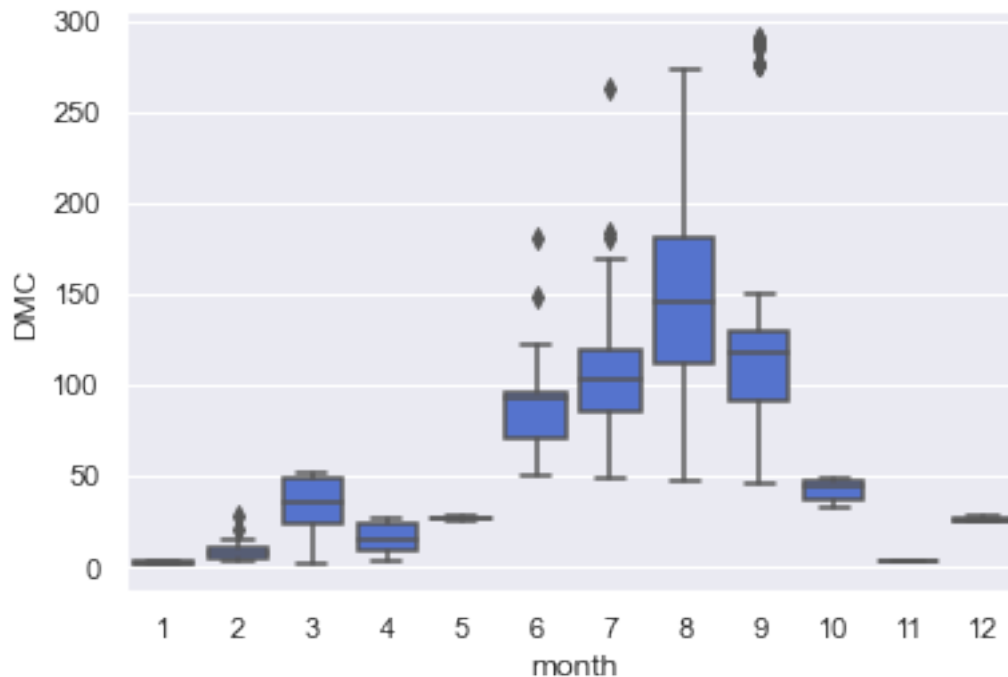
[25]: `<AxesSubplot:xlabel='month', ylabel='temp'>`



Widzimy, e temperatura jest zapewne istotnym czynnikiem w poarach, widzimy jednak, e nie zaley to tylko od temperatury. W czerwcu i libcu temperatury s porównywalne z temperaurami w sierpniu i wrzeniu, a jednak ilo poarów w tych miesiacach jest znaczca.

```
[32]: sns.boxplot(x="month", y="DMC", data=df_fire, color="RoyalBlue")
```

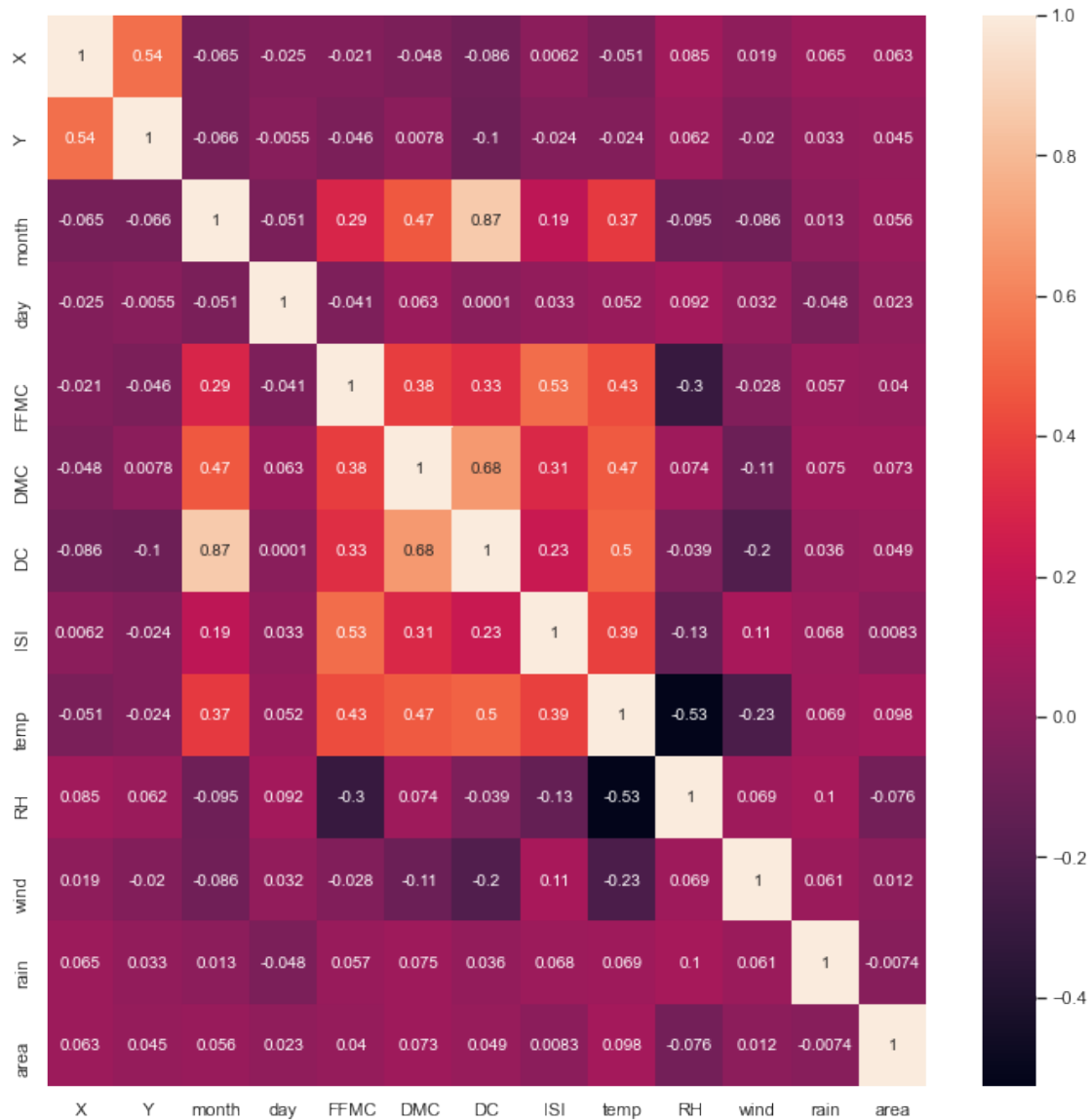
```
[32]: <AxesSubplot:xlabel='month', ylabel='DMC'>
```



Po obejrzeniu powyszego wykresu moemy wnioskowa, e to odpowiednia temperatura w poczeniu z nisk wilgotnoscia sprzyja poarom. Warto ta jest cakiem wysoka w marcu co zgadza si ze wzrotem poarów w tym miesicu.

```
[12]: correlation = df_fire.corr()  
fig, ax = plt.subplots(figsize=(12,12))  
sns.heatmap(correlation, annot = True)
```

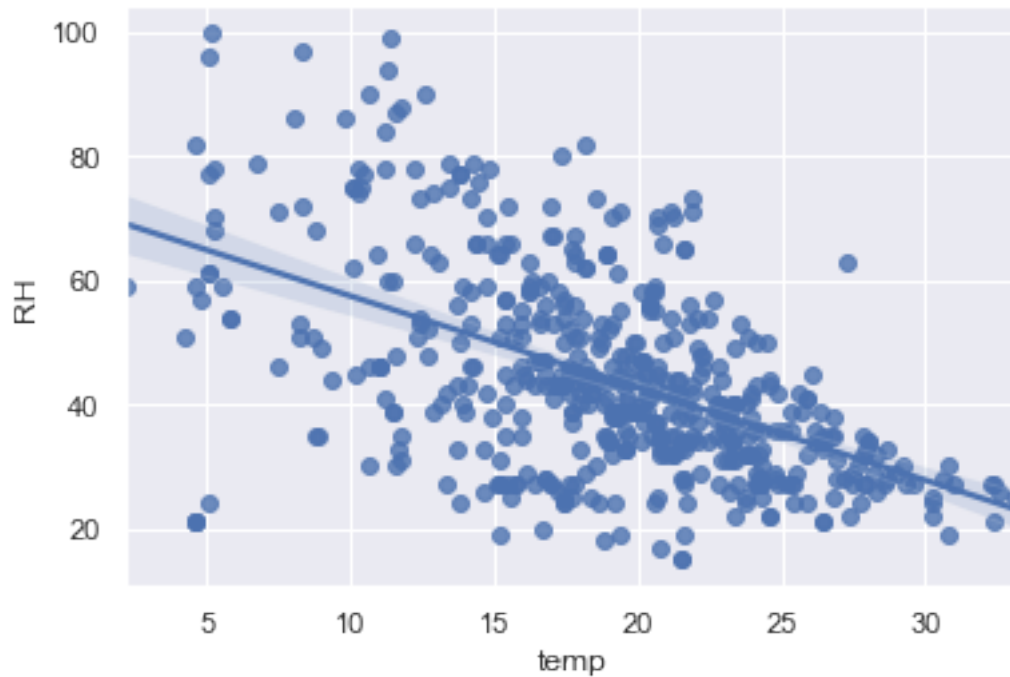
```
[12]: <AxesSubplot:>
```

Powysza macierz korelacji pokazuje nam które wartości mogłyby ze sobą skorelowane. Widać z niej na przykład, że temperatura jest skorelowana ujemnie z iloci ludzi w lasach. Oznacza to, że im większa temperatura tym mniej ludzi chodzi do lasu. Jest to dość zaskakujące, gdy w Polsce tendencja jest raczej odwrotna. Nie jest zaskoczeniem, że cechy oznaczające wilgotność ze sobą skorelowane. Na ogół jednak kolumny nie są jednak ze sobą skorelowane.

```
[13]: sns.regplot(data = df_fire, x = 'temp', y = 'RH')
```

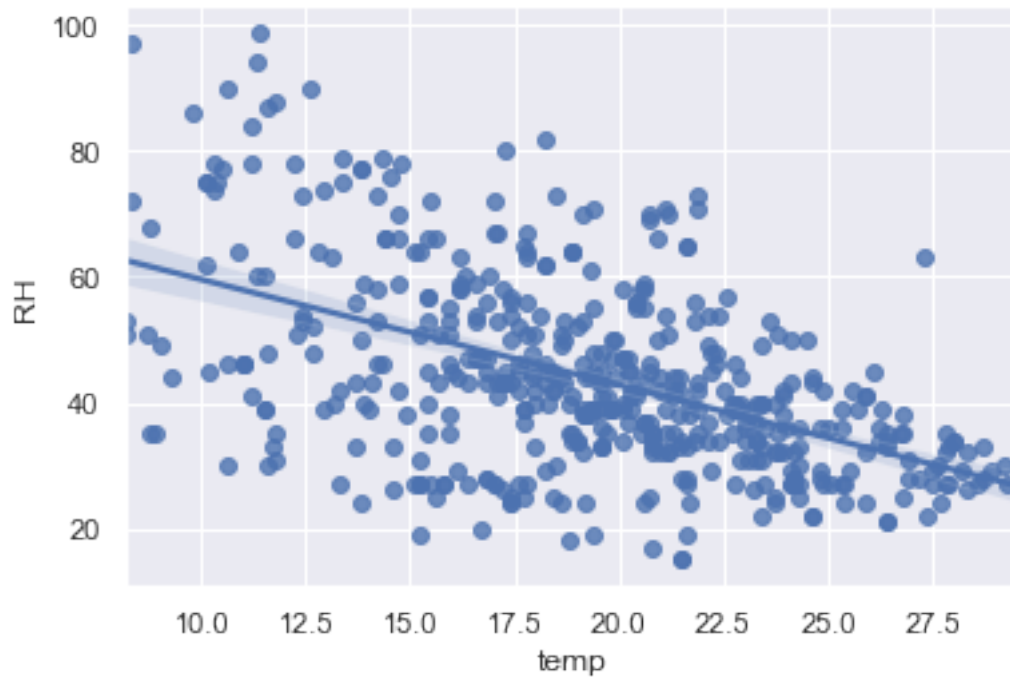
```
[13]: <AxesSubplot:xlabel='temp', ylabel='RH'>
```



Widzimy z wykresu e wraz ze wzrostem temperatury maleje liczba ludzi w lesie. Spróbujmy jednak pozby si wartoci odstajcych aby zobaczy czy to zmieni tendencje.

```
[14]: sns.regplot(data = df_fire[(df_fire.temp > 8) & (df_fire.temp <= 30)], x = 'temp', y = 'RH')
```

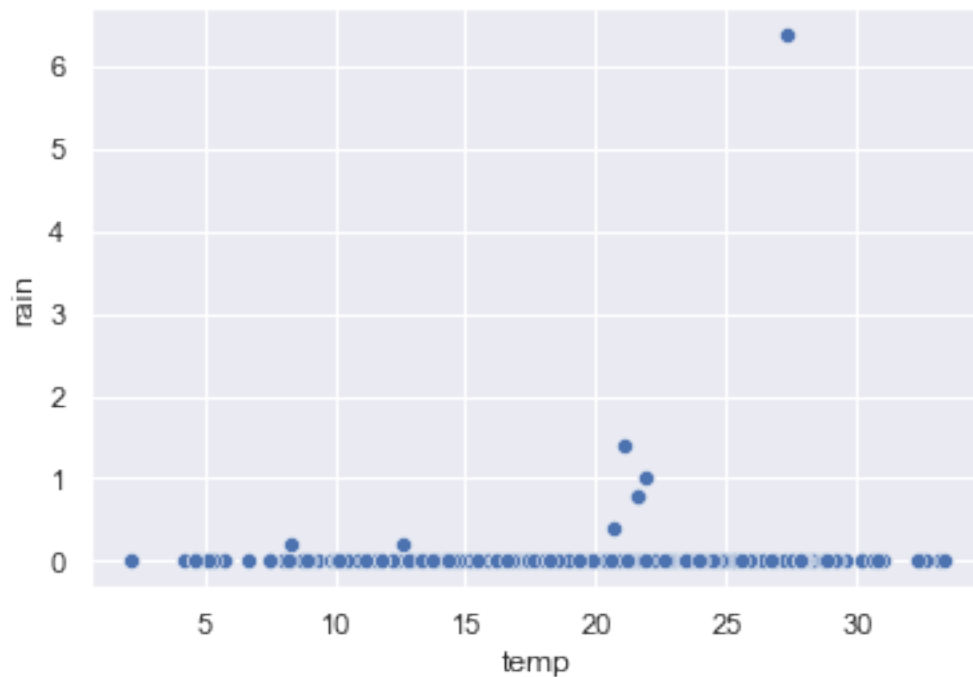
```
[14]: <AxesSubplot:xlabel='temp', ylabel='RH'>
```



Wypaszcza to nieco krzyw dopasowanie, jednak nieznaczenie. By moe, gdy jest cieplej to więcej pada i std mniejsze zainteresowanie spacerami do lasu.

```
[15]: sns.scatterplot(data = df_fire, x = 'temp', y = 'rain')
```

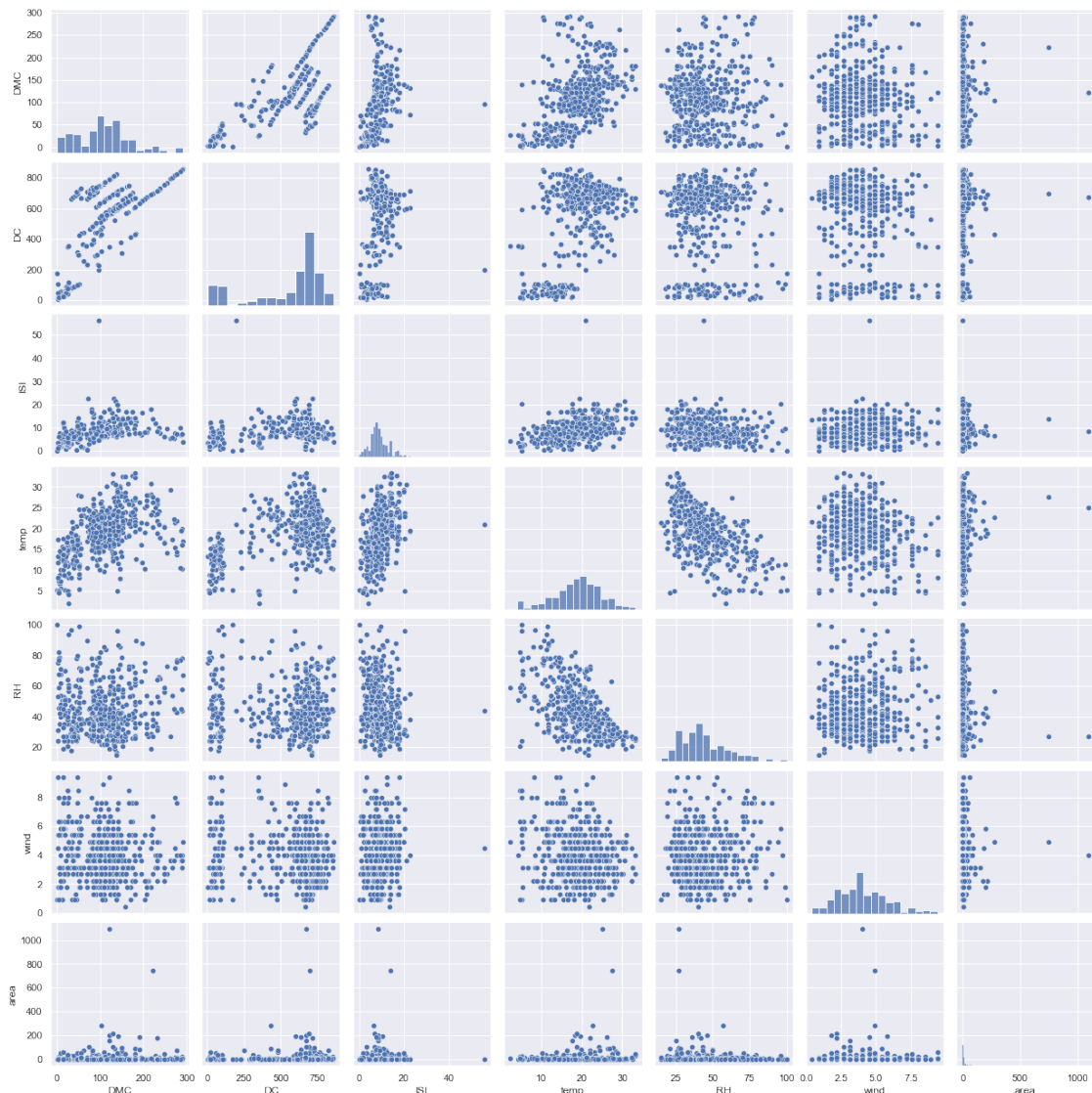
```
[15]: <AxesSubplot:xlabel='temp', ylabel='rain'>
```



Nie możemy tego stwierdzić, gdy w większości rekordów w kolumnie rain jest zerem.
Sprawdźmy jak korelują ze sobą inne zmienne.

```
[26]: sns.pairplot(df_fire[['DMC', 'DC', 'ISI', 'temp', 'RH', 'wind', 'area']])
```

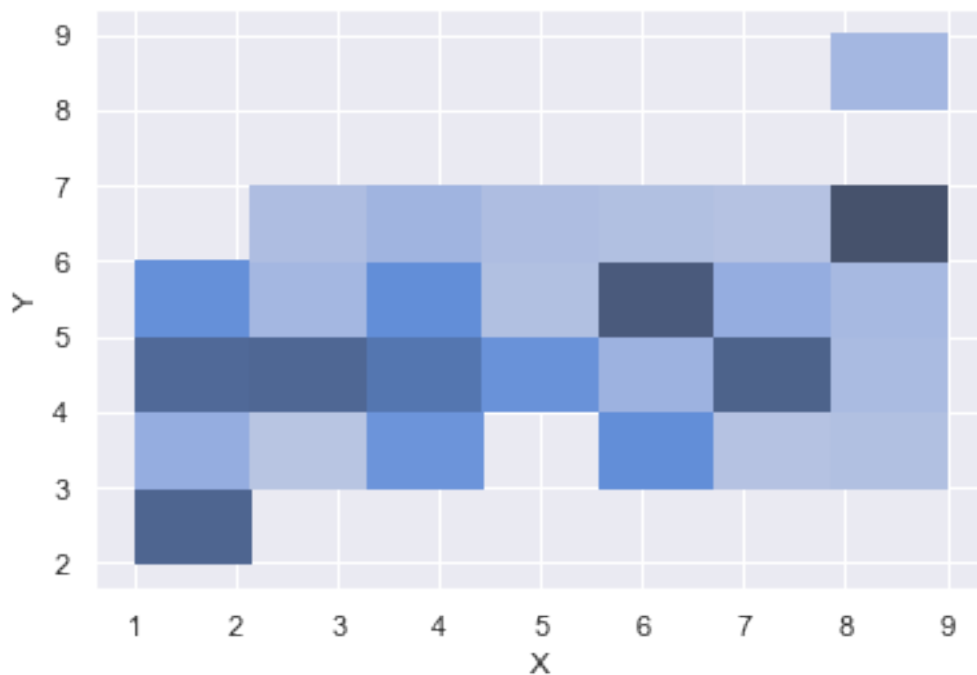
```
[26]: <seaborn.axisgrid.PairGrid at 0x13a1148d0>
```



Oprócz analizowanej wcześniej zależności temperatury i liczby ludzi w lesie, ciekawe zależności widzi się między DC a DMC. Podobnie wygląda również zależność między DMC a ISI. Relacja między temperaturą a ISI przypomina trochę relację liniową. Na wykresach wyżej widzi się również zależności pomiędzy DMC a temperaturą, potwierdza to nasze przypuszczenia dotyczące ich związku z parametrami w danych miesięcznych. Powyższe wykresy punktowe potwierdzają niejako to, co mówią macie korelacji.

```
[30]: sns.histplot(data = df_fire, x = 'X', y = 'Y', bins = 7)
```

```
[30]: <AxesSubplot:xlabel='X', ylabel='Y'>
```



Powysza mapa prezentuje nam, w których obszarach lasu najczęściej dochodzi do poarów

1.2 Narzdzia do automatycznej eksploracji danych

```
[18]: profile = ProfileReport(df_fire, title='Pandas Profiling Report')  
profile
```

```
Summarize dataset:  0%|          | 0/26 [00:00<?, ?it/s]
```

```
Generate report structure:  0%|          | 0/1 [00:00<?, ?it/s]
```

```
Render HTML:  0%|          | 0/1 [00:00<?, ?it/s]
```

```
<IPython.core.display.HTML object>
```

```
[18]:
```

Profile report ułatwia zdecydowanie eksploracji danych. Ma on jednak trochę ograniczeń. Za jego pomocą nie można na przykład bliżej przyjrzeć się poszczególnym rozkładom, nie byłoby w stanie zobaczyć, że rozkład zmiennej area ma postać $1/x$. Drugim problemem jest ograniczona wizualizacja zmiennych kategorycznych, nie można na przykład zobaczyć boxplotów, które lepiej odzwierciedlają rozkłady z podziałem na kategorie. Myślę, że pandas_profiling to świetne narzędzie, do rzucenia okiem na dane. Zapewni nam ono podstawowe statystyki i wykresy. Powinniśmy jednak zwrócić uwagę na to, co jest ciekawe lub niepasujące i samodzielnie to pogłębić.

[]: