

General information:

Author:

Mikko Saukkoripi 013877851

Date:

14 Dec 2020

Background of the project:

This project is done for the course Distributed Data Infrastructure fall 2020. The course is part of the Masters of Data Science Programme, and the primary purpose of the project was to practice the use of the GraphLab, which was replaced by Turi Create.

General structure:

This notebook follows the structure of the given questions in the same order as they are. In other words, each question has been answered in its own chapter. At the end of each chapter, there is a text answer to the given research question.

Project questions:

Distributed Data Infrastructures, Fall 2020, Project 2

Assignment In this project, you are supposed to use GraphLab to analyze data sets. GraphLab is a high performance, open-source, big data framework which can be discussed in the course. We provide you with a data set and you should run GraphLab either on AWS or locally to analyze the data set.

The dataset was released by Telecom Italia for their Open Big Data Challenge in 2014. It contains telecommunication records, weather, air quality, electricity consumption for city of Milan and province of Trentino in Italy in November and December 2013. You can find detailed description of the dataset from the paper: <https://www.nature.com/articles/data201555>.

Requirements You need write a program that uses GraphLab to provide answers to the following questions.

First three questions are worth 2 points total.

- Find the most congested communication period of the day in Milan and Trentino.
- Find top 5 Italian provinces which are most called by residents of Milan and Trentino on average.
- List top 5 languages tweeted by distinct users in Milan. How popular is Finnish as a tweeting language in Milan?

The following questions are worth 3 points total.

- Compare call and internet activity between 24th, 25th and 26th December to 26th, 27th, 28th November for Milan. Plot the distribution.
- Find correlation between user communication activity and different weather conditions (e.g. rain, snow etc.) in Milan and Trentino.

The final questions are worth 2 points each.

- Plot the heatmap of user telecommunication activity for both Milan and Trentino. Do you observe any shift in communication pattern of users between day and night? (A typical day time is between 8AM to 8PM)
- Investigate and plot the correlation between air quality and weather (temperature, sunshine, precipitation, etc.).

Read needed libraries before starting to answer the question

```
In [242]. # Import
import pandas as pd
from turicreate import SFrame, SGraph, Vertex, Edge
import arrow
import glob2
import time
import datetime
import matplotlib.pyplot as plt
import seaborn as sns
import geopandas as gpd
import descartes
import datetime
```

1. Find the most congested communication period of the day in Milan and Trentino.

Answer to this question based on the telecommunication data for Milan and Trentino on 24th December 2013. Answer will be the most congested hour in the given day. We will use combination of SMS, call and internet traffic data to find most congested communication period.

I have used a single day in this exercise because I have answered the question on my laptop. The size of the single-day telecom data for one data is the range between 270Mb to 350Mb and using multiple days in the calculation is a heavy and slow process for a laptop. Simultaneously, using multiple days would only change the file reading, but otherwise, code would be same.

Text answer to the question is given at the end of the code.

Telecommunication data has following information:

- Square id: identification string of a given square of Milan/Trentino GRID;
- Time Interval: start interval time expressed in milliseconds. The end interval time can be obtained by adding 600,000 milliseconds (10min) to this value;
- Country code: the phone country code of the nation.
- SMS-in activity: activity proportional to the amount of received SMSs inside a given Square id and during a given time interval. The SMSs are sent from the nation identified by the Country code.
- SMS-out activity: activity proportional to the amount of sent SMSs inside a given Square id during a given time interval. The SMSs are received in the nation identified by the Country code;
- Call-in activity: activity proportional to the amount of received calls inside the Square id during a given time interval. The calls are issued from the nation identified by the Country code;
- Call-out activity: activity proportional to the amount of issued calls inside a given Square id during a given time interval. The calls are received in the nation identified by the Country code;
- Internet traffic activity: number of CDRs generated inside a given Square id during a given time interval. The Internet traffic is initiated from the nation identified by the Country code;

```
In [243]. # Define function to find peak hour for the day
def peak_n_hours(sframe, n):
    # Rename columns
    sframe = sframe.rename({'0': 'SquareID', '1': 'Time', '2': 'country_code', '3': 'SMSIn', '4': 'SMSOut', '5': 'CallIn', '6': 'CallOut', '7': 'InternetTraffic'})
    # Add hour column
    sframe['hour'] = sframe['Time'].apply(lambda t: dt.datetime.strptime(t, '%Y-%m-%d %H:%M:%S').strftime('%H'))
    # Calculate hourly activity
    sf_hourly = sframe.groupby(key_column_names = 'hour', operations= ('sum', 'SMSIn', 'SMSOut', 'CallIn', 'CallOut', 'InternetTraffic'))
    # Calculate sum of the activity columns
    sf_hourly['total'] = sf_hourly['SMSIn'] + sf_hourly['SMSOut'] + sf_hourly['CallIn'] + sf_hourly['CallOut'] + sf_hourly['InternetTraffic']
    # Drop other columns except hour and total
    sf = sf_hourly[['hour', 'total']]
    # Sort sf_hourly and get 5 largest values
    sf = sf.sort('total', ascending=False)
    # Return n first hours and their total activity
    return sf.head(n)
```

Read file and then use peak_n_hours function to find most congested hours.

File is initially read with Pandas, because Turi is not able to read single csv file with different length of rows. After file is read with pandas it is converted to Turi SFrame.

```
In [244]. # Set path to files
milano_dec24_path = 'data/milano_telecom_dec/sms-call-internet-mi-2013-12-24.txt'
trentino_dec24_path = 'data/trento_telecom_dec/sms-call-internet-tr-2013-12-24.txt'

# Read the files. SFrame is not able to handle rows with different lengths,
# so first read data with pandas and then convert to SFrame
milano = pd.read_csv(milano_dec24_path, sep='\t', header=None).fillna(0)
milano = SFrame(milano)
trentino = pd.read_csv(trentino_dec24_path, sep='\t', header=None).fillna(0)
trentino = SFrame(trentino)

# Run Function and print n=5 busiest communication hours
print("The most congested communication hours of the 24th Dec 2013 in Milan")
print(peak_n_hours(milano, 5))
print("The most congested communication hours of the 24th Dec 2013 in Trentino")
print(peak_n_hours(trentino, 5))

The most congested communication hours of the 24th Dec 2013 in Milan
-----+-----+
| hour | total |
|-----+-----+
| 16 | 5416257.583677531 |
| 11 | 5371330.078746542 |
| 15 | 5344031.76565861 |
| 10 | 537381.37456588 |
| 15 | 5284227.558210146 |
|-----+-----+
[5 rows x 2 columns]

The most congested communication hours of the 24th Dec 2013 in Trentino
-----+-----+
| hour | total |
|-----+-----+
| 16 | 1036600.378461629 |
| 10 | 1029579.4795295168 |
| 17 | 1011900.8027854251 |
| 11 | 1007262.7792211829 |
| 15 | 1006217.5522763213 |
|-----+-----+
[5 rows x 2 columns]
```

Answer to research question 1:

We can see from the results, that in the 24th December 2013 the most congested communication hour in both cities was from 4pm to 5pm UTC time, which is in Italian time 5pm to 6pm.

2. List top 5 Italian provinces which are most called by residents of Milan and Trentino on average.

I will answer this question based on the "Milan/Trentino to provinces" telecommunication dataset on 24th December 2013. Using only a single day is the same as in exercise 1. The reason is that I used only a laptop for computation, and the fact that using multiple days would not change the code. This code is also convertible to use multiple days dataset by only changing the file reading process.

Milan/trentino to provinces dataset has following columns:

- Square id: identification string of a given square of Milan/Trentino GRID;
- Time Interval: start interval time expressed in milliseconds. The end interval time can be obtained by adding 600,000 milliseconds (10 min) to this value;
- Square to Province Inter: Value representing the interaction between the Square id and the Province. It is proportional to the number of calls exchanged between callers, which are located in the Square id, and receivers located in the Province;
- Province to Square Inter: Value representing the interaction between the Square id and the Province. It is proportional to the number of calls exchanged between callers, which are located in the Province, and receivers located in the Square id.
- Province: the name of the Italian province.

Define function to find most called n number of provinces.

Used information is defined in the column "square to province".

```
In [245]. def most_called_provinces(sframe, n):
    # Rename columns
    sframe = sframe.rename({'0': 'SquareID', '1': 'Province', '2': 'Time', '3': 'SqrtToP'})
    # Drop columns where ProvinceToSqrt=0. These rows are incoming calls from province
    # and we want only count calls to provinces.
    sframe = sframe[sframe['SqrtToP'] > 0]
    # Count by provinces
    famous_prov = sframe.groupby(key_column_names = 'Province', operations= ('count',
    # Sort by count and return top provinces
    return famous_prov.sort('count', ascending=False).head(n)
```

Read files and use most_called_provinces to find most called provinces.

As in the question 1, we will first read file to Pandas, because Turi is not able to handle different length of the rows. After file is successfully read to Pandas DataFrame, we will convert it to Turi SFrame.

```
In [246]. # Measure how long it takes
start = time.time()

# Set path to files
milano_dec24_path = 'data/milano_provinces_telecom/ml-to-provinces-2013-12-24.txt'
trentino_dec24_path = 'data/trento_provinces_telecom/tr-to-provinces-2013-12-24.txt'

# Read the files. SFrame is not able to handle rows with different lengths,
# so first read data with pandas and then convert to SFrame
milano = pd.read_csv(milano_dec24_path, sep='\t', header=None).fillna(0)
milano = SFrame(milano)
trentino = pd.read_csv(trentino_dec24_path, sep='\t', header=None).fillna(0)
trentino = SFrame(trentino)

# Run Function and print top 5 provinces most called by residents
start = time.time()
print("The most called provinces from Milano on the 24th Dec 2013")
print(most_called_provinces(milano, 5))
print("The most called provinces from Trentino on the 24th Dec 2013")
print(most_called_provinces(trentino, 5))

# Print time how long it took
end = time.time()
print("Time to run was", round(end - start), "seconds")

The most called provinces from Milano on the 24th Dec 2013
-----+-----+
| Province | count |
|-----+-----+
| MILANO | 970483 |
| MONZA E DELLA BRIANZA | 217406 |
| PAVIA | 175463 |
| VARESE | 187539 |
| NAPOLI | 84398 |
|-----+-----+
[5 rows x 2 columns]

The most called provinces from Trentino on the 24th Dec 2013
-----+-----+
| Province | count |
|-----+-----+
| TRENTO | 527030 |
| BOLZANO/BOZEN | 107719 |
| BRESCIA | 80951 |
| VERONA | 58680 |
| MILANO | 54876 |
|-----+-----+
[5 rows x 2 columns]

Time to run was 2 seconds
```

Answer to research question 2:

We can see from the results that most of the calls happen inside the province. After this, there seem to be large neighboring provinces

3. List top 5 languages tweeted by distinct users in Milan. How popular is Finnish as a tweeting language in Milan?

To answer this question, I will use SocialPulse dataset. This is a Twitter dataset with all the needed information.

The SocialPulse dataset contains geolocalized tweets originated from Milan between November 1, 2013 and January 1st, 2014.

- user: anonymized Twitter username;
- language: language of the Tweet, where und means undefined;
- municipality: the municipality in which the tweet has been probably created.
- timestamp: Tweet timestamp;
- geometry latitude: approximate position of the tweet, in geoJSON format. Error <600m.
- geometry longitude: approximate position of the tweet, in geoJSON format. Error <600m.

Read the file, filter by municipality and then group and count by language

```
In [247]. # Set path to file
milano_path = 'data/social_pulse/social_pulse_milano.csv'

# Read the files
milano = SFrame.read_csv(milano_path, sep='\t', header=None)

# Rename columns
milano = milano.rename({'X1': 'user', 'X2': 'language', 'X3': 'municipality', 'X4': 'time'})

# Print unique values in municipality column
#print(list(milano['municipality'].unique()))

# Filter municipalities column to include only Milano
milano = milano[milano['municipality'] == 'Milano']

# Count language
languages_count = milano.groupby(key_column_names = 'language', operations= ('count',
languages_count = languages_count.sort('count', ascending=False)

# Print 15 most popular tweeting languages
languages_count.print_rows(num_rows=15)
```

Finished parsing file /Users/sauk/Desktop/distributed data infrastructure/GraphLab project/data/social_pulse/social_pulse_milano.csv

Parsing completed. Parsed 100 lines in 0.166448 sec.

-----+-----+
| language | count |
|-----+-----+
it	146546
en	42598
es	6971
tl	6502
pt	4611
fr	3650
de	2949
und	1800
tr	2005
no	1540
nl	1510
fi	1467
ro	1540
-----+-----+	
[45 rows x 2 columns]

Answer to research question 3:

From the table, we can see that the Italian language is the most common in the municipality of Milan. The second most popular language is English, that is vastly more popular than the rest of the languages. Finnish is in the list in place 15.

4. Compare call and internet activity between 24th, 25th and 26th December to 26th, 27th, 28th November for Milan. Plot the distribution.

To answer this question, I will calculate hourly call and internet activity for each day, and then make histogram and plot of hourly activities in given days in November and December.

Telecommunication data has following information:

- Square id: identification string of a given square of Milan/Trentino GRID;
- Time Interval: start interval time expressed in milliseconds. The end interval time can be obtained by adding 600,000 milliseconds (10min) to this value;
- Country code: the phone country code of the nation.
- SMS-in activity: activity proportional to the amount of received SMSs inside a given Square id and during a given time interval. The SMSs are sent from the nation identified by the Country code.
- SMS-out activity: activity proportional to the amount of sent SMSs inside a given Square id during a given time interval. The SMSs are received in the nation identified by the Country code;
- Call-in activity: activity proportional to the amount of received calls inside the Square id during a given time interval. The calls are issued from the nation identified by the Country code;
- Call-out activity: activity proportional to the amount of issued calls inside a given Square id during a given time interval. The calls are received in the nation identified by the Country code;
- Internet traffic activity: number of CDRs generated inside a given Square id during a given time interval. The Internet traffic is initiated from the nation identified by the Country code;

Define function to read multiple files in the same folder in given date interval

```
In [248]. def create_hourly_telecom_SFrame(paths, start_date, end_date):
    # Measure time how long it takes. Set start time.
    start = time.time()

    # Filter paths by start date and end date
    wanted_paths = []
    for path in paths:
        date = int(path[-6:-4])
        if (date <= end_date) and (date >= start_date):
            wanted_paths.append(path)

    # Read the files and append dataframes to result_df
    result_df = pd.DataFrame(columns = ['day', 'SquareID', 'Time', 'country_code', 'SMSIn', 'SMSOut', 'CallIn', 'CallOut', 'InternetTraffic'])
    for path in wanted_paths:
        day = path[-6:-4]
        df = pd.read_csv(path, sep='\t', header=None)
        df.columns = ['SquareID', 'Time', 'country_code', 'SMSIn', 'SMSOut', 'CallIn', 'CallOut', 'InternetTraffic']
        result_df = result_df.append(df)

    # Convert Pandas dataframe to SFrame
    sf = SFrame(data=result_df)

    # Drop columns SquareID and country_code
    sf = sf.remove_column('SquareID')
    sf = sf.remove_column('country_code')

    # Add hour column and add 1 hour to time (from UTC to Rome time).
    sf['hour'] = sf['Time'].apply(lambda timestamp: arrow.get(timestamp+3600000).format('%H'))

    # Group by day, hour, and sum of each activity in each hour and day
    sf_hourly = sf.groupby(['day', 'hour'], ('callin', 'agg.SUM('call_in', 'callout', 'agg.SUM('call_out', 'internet_traffic', 'agg.SUM('internet_traffic'))

    # Calculate sums of callin and callout
    sf_hourly['callin'] = sf_hourly['callin'] + sf_hourly['callout']

    # Drop columns SquareID and country_code
    sf_hourly = sf_hourly.remove_column('callin')
    sf_hourly = sf_hourly.remove_column('callout')

    # Sort by day and hour
    sf_hourly = sf_hourly.sort(['day', 'hour'], ascending=True)

    # Set hour and day type to int
    sf_hourly['day'] = sf_hourly['day'].astype(int)
    sf_hourly['hour'] = sf_hourly['hour'].astype(int)

    # Print time how long it took
    end = time.time()
    print("Time to run function was", round(end - start), "seconds")

    # Return results
    return sf_hourly
```

Read the files using create_hourly_telecom_SFrame with given days

```
In [249]. # Execute the glob function that returns a list of filepaths
data_paths_milano_nov = glob2.glob("data/milano_telecom_nov/**")
data_paths_milano_dec = glob2.glob("data/milano_telecom_dec/**")

# Use create_telecom_SFrame to read paths and join the files within given dates
milano_nov = create_hourly_telecom_SFrame(data_paths_milano_nov, 26, 28)
milano_dec = create_hourly_telecom_SFrame(data_paths_milano_dec, 24, 26)

# Print first 5 rows
milano_nov.print_rows(num_rows=5)
milano_dec.print_rows(num_rows=5)

Time to run function was 88 seconds
Time to run function was 76 seconds
-----+-----+
| day | hour | int_traffic | calls |
|-----+-----+
| 26 | 0 | 3081805.3779902484 | 40192.7323245973 |
| 26 | 1 | 2572701.778685294 | 17706.488070584164 |
| 26 | 2 | 2231849.13408716 | 11778.29143263028 |
| 26 | 3 | 2032399.175571133 | 10176.70211725669 |
| 26 | 4 | 1925251.8709400047 | 11383.31819843374 |
|-----+-----+
[72 rows x 4 columns]

-----+-----+
| day | hour | int_traffic | calls |
|-----+-----+
| 24 | 0 | 2478071.569021002 | 54213.38695708153 |
| 24 | 1 | 2100878.8674795628 | 26984.249997574632 |
| 24 | 2 | 1831355.198863678 | 17032.58246028943 |
| 24 | 3 | 1680874.742091802 | 11930.500057642956 |
| 24 | 4 | 1549273.9723427454 | 12418.661768335747 |
|-----+-----+
[72 rows x 4 columns]
```

Make call and internet activity histograms of both months to see the distribution

```
In [250]. fig, ax = plt.subplots(1, 2, tight_layout=True, figsize=(10, 5))

plt.subplot(1, 2, 1)
bins = list(range(0, 100000 + 2500, 2500))
plt.rcParams["figure.figsize"] = (6, 6)
plt.hist(milano_nov['calls'], bins=bins, alpha=0.5)
plt.hist(milano_dec['calls'], bins=bins, alpha=0.5)
plt.legend(['Milano Nov 26th to 28th', 'Milano Dec 24th to 26th'], loc = 'upper right')
plt.title('Milano calls activity histogram by month')

plt.subplot(1, 2, 2)
bins = list(range(0, 5000000 + 200000, 200000))
plt.rcParams["figure.figsize"] = (6, 6)
plt.hist(milano_nov['int_traffic'], bins=bins, alpha=0.5)
plt.hist(milano_dec['int_traffic'], bins=bins, alpha=0.5)
plt.legend(['Milano Nov 26th to 28th', 'Milano Dec 24th to 26th'], loc = 'upper left')
plt.title('Milano internet activity histogram by month')
plt.show()
```



Define function to make graphs of hourly internet and call activity

```
In [251]. # Make histograms
def make_graphs(SFrame, title):
    # Make days list and sort it
    days = list(SFrame['day'].unique())
    days = sorted(days)

    # Make Trentino histograms. One for each precipitation intensity class.
    fig, ax = plt.subplots(1, 2, tight_layout=True, figsize=(10, 5))

    # Graph1
    plt.subplot(1, 2, 1)
    for day in days:
        sf_day = SFrame[SFrame['day'] == day]
        plt.plot(sf_day['hour'], sf_day['calls'])
        plt.ylim(0, 1000000)
        plt.xlabel('Time (Italian time) \n Gray lines are at 8am and 8pm')
        plt.axvline(x=8, color='gray', linestyle='--', lw=1)
        plt.axvline(x=20, color='gray', linestyle='--', lw=1)
        plt.title('Hourly calls activity (' + day + ')')
        plt.legend([days[0], days[1], days[2]], loc = 'upper left')
        plt.show()

    # Graph2
    plt.subplot(1, 2, 2)
    for day in days:
        sf_day = SFrame[SFrame['day'] == day]
        plt.plot(sf_day['hour'], sf_day['int_traffic'])
        plt.ylim(0, 6000000)
        plt.xlabel('Time (Italian time) \n Gray lines are at 8am and 8pm')
        plt.axvline(x=8, color='gray', linestyle='--', lw=1)
        plt.axvline(x=20, color='gray', linestyle='--', lw=1)
        plt.title('Hourly internet activity (' + day + ')')
        plt.legend([days[0], days[1], days[2]], loc = 'upper left')
        plt.show()
```

Run the make_graphs function

```
In [252]. make_graphs(milano_nov, "Milano Nov 26th to 28th")
make_graphs(milano_dec, "Milano Dec 24th to 26th")

-----+-----+
| SquareID | total | timestamp |
|-----+-----+
| 10000 | 0.571041869401956 | 201311272300 |
| 10000 | 0.631737324861012 | 201311272310 |
| 10000 | 0.936505211409811 | 201311272320 |
| 10000 | 0.0416022755629953 | 201311272330 |
| 10000 | 0.14672898719976832 | 201311272330 |
|-----+-----+
[3049709 rows x 3 columns]

-----+-----+
| SquareID | total | timestamp |
|-----+-----+
| 10000 | 0.571041869401956 | 201311272300 |
| 10000 | 0.631737324861012 | 201311272310 |
| 10000 | 0.936505211409811 | 201311272320 |
| 10000 | 0.0416022755629953 | 201311272330 |
| 10000 | 0.14672898719976832 | 201311272330 |
|-----+-----+
[36823804 rows x 3 columns]
```

Read weather data and join it to telecom data

```
In [253]. # Set paths to weather files
milano_weather_path = 'data/MeteoMilano_01-11-13_01-14.csv'
trentino_weather_path = 'data/Precipitation-Trentino.csv'

# Read weather files
milano_weather = SFrame.read_csv(milano_weather_path, sep=',', header=None)
trentino_weather = SFrame.read_csv(trentino_weather_path, sep=',', header=None)

# Rename columns of the weather SFrames
milano_weather = milano_weather.rename({'X1': 'timestamp', 'X2': 'SquareID', 'X3': 'intensity', 'X4': 'coverage'})
trentino_weather = trentino_weather.rename({'X1': 'timestamp', 'X2': 'SquareID', 'X3': 'intensity', 'X4': 'coverage'})

# Drop coverage and type from Milano weather data, because we will only study
# how intensity affects the communication activity.
milano_weather = milano_weather.remove_column('coverage')
milano_weather = milano_weather.remove_column('type')

# Print first rows
```



```
milano_weather.print_rows(num_rows=5)
trentino_weather.print_rows(num_rows=5)
```

Finished parsing file /Users/saukk/Desktop/Distributed data infrastructure/Grap hLab project/data/MeteoMilano_01-11-13_01-01-14.csv
Parsing completed. Parsed 100 lines in 0.011216 secs.

```
-----
Inferred types from first 100 line(s) of file as
column_type_hints(int,int,int,int,int)
If parsing fails due to incorrect types, you can correct
the inferred type list above and pass it to read_csv in
the column_type_hints argument
-----
```

Finished parsing file /Users/saukk/Desktop/Distributed data infrastructure/Grap hLab project/data/MeteoMilano_01-11-13_01-01-14.csv
Parsing completed. Parsed 100 lines in 0.018 secs.

Finished parsing file /Users/saukk/Desktop/Distributed data infrastructure/Grap hLab project/data/precipitation-trentino.csv
Parsing completed. Parsed 100 lines in 0.314223 secs.

```
-----
Inferred types from first 100 line(s) of file as
column_type_hints(int,int,int,int,int)
If parsing fails due to incorrect types, you can correct
the inferred type list above and pass it to read_csv in
the column_type_hints argument
-----
```

Read 2489903 lines. Lines per second: 4.98608e+06
Finished parsing file /Users/saukk/Desktop/Distributed data infrastructure/Grap hLab project/data/precipitation-trentino.csv
Parsing completed. Parsed 3489417 lines in 0.583172 secs.

```
-----+-----+
| timestamp | SquareID | intensity |
+-----+-----+
| 201311010000 | 1 | 0 |
| 201311010000 | 2 | 0 |
| 201311010000 | 3 | 0 |
| 201311010000 | 1 | 0 |
| 201311010010 | 1 | 0 |
+-----+-----+
[33316 rows x 3 columns]

-----+-----+
| timestamp | SquareID | intensity |
+-----+-----+
| 201311010000 | 2383 | 1 |
| 201311010020 | 9452 | 1 |
| 201311010020 | 8278 | 1 |
| 201311010020 | 6871 | 2 |
| 201311010020 | 5689 | 3 |
+-----+-----+
[3489417 rows x 3 columns]
```

Join telecom and weather data

```
In [256]. # First set SFrames column timestamp and SquareID to int
milano_tele["timestamp"] = milano_tele["timestamp"].astype(int)
milano_tele["SquareID"] = milano_tele["SquareID"].astype(int)
trentino_tele["timestamp"] = trentino_tele["SquareID"].astype(int)
trentino_tele["SquareID"] = trentino_tele["SquareID"].astype(int)

trentino_tele["timestamp"] = trentino_tele["timestamp"].astype(int)
trentino_tele["SquareID"] = trentino_tele["SquareID"].astype(int)
trentino_weather["timestamp"] = trentino_weather["timestamp"].astype(int)
trentino_weather["SquareID"] = trentino_weather["SquareID"].astype(int)

# Join weather and telecom datasets
milano = milano_weather.join(milano_tele, on=["timestamp", "SquareID"], how='inner')
trentino = trentino_weather.join(trentino_tele, on=["timestamp", "SquareID"], how='inner')

# Print first rows
milano.print_rows(num_rows=5)
trentino.print_rows(num_rows=5)

-----+-----+
| timestamp | SquareID | intensity | total |
+-----+-----+
| 201311252300 | 1 | 0 | 0.05701250935247165 |
| 201311252300 | 1 | 0 | 0.5391165522386262 |
| 201311252310 | 1 | 0 | 0.003574620210998875 |
| 201311252310 | 1 | 0 | 0.544478482551247 |
| 201311252320 | 1 | 0 | 0.405479801987943 |
+-----+-----+
[10023 rows x 4 columns]

-----+-----+
| timestamp | SquareID | intensity | total |
+-----+-----+
| 201311280010 | 10043 | 1 | 0.0 |
| 201311281530 | 10095 | 1 | 0.41318456643054674 |
| 201311281530 | 10095 | 1 | 0.4934586718961125 |
| 201311281530 | 10095 | 1 | 0.0 |
| 201311281640 | 10100 | 1 | 1.4242910096578831 |
+-----+-----+
[3581663 rows x 4 columns]
```

Check the intensity groups distribution in weather data

We can see that in Milan, there are no datapoints in intensity class 3 and in Trentino in intensity classes between 14 and 18. These classes mean really heavy raining.

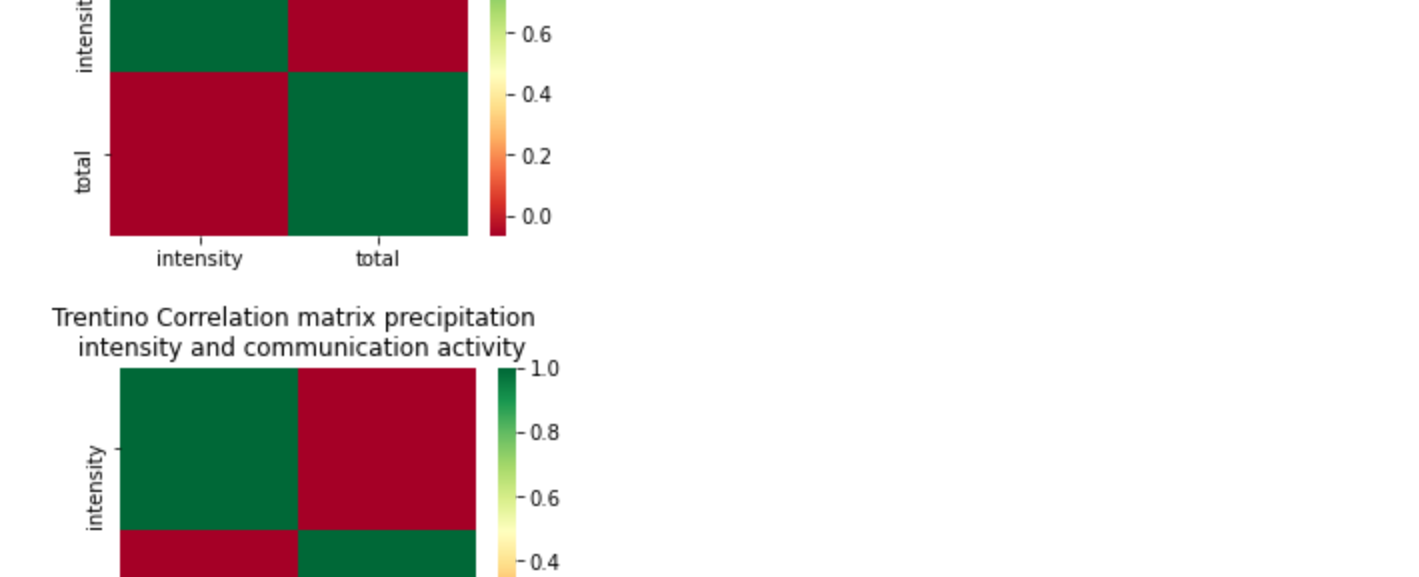
We can also see that intensity value 0 is the most common class.

```
In [257]. # Print unique intensity values in Milano and Trentino
# Used to make sure, that data description is correct.
print("Milano intensity classes: {}".format(list(milano["intensity"].unique())))
print("Trentino intensity classes: {}".format(list(trentino["intensity"].unique())))

# Make intensity histograms to see intensity distribution in Milan and Trentino.
# Milan and Trentino precipitation intensity has different scales.
# Milan scale is 0-4 and Trentino 1-18.
fig, ax = plt.subplots(2, 2, tight_layout=True, figsize=(10, 5))

plt.subplot(1, 2, 1)
plt.hist(milano["intensity"], bins=3, edgecolor='black')
plt.title("Milano precipitation intensity histogram. \n Nbr of datapoints: {}".format(len(milano["intensity"])))
plt.subplot(1, 2, 2)
plt.hist(trentino["intensity"], bins=13, edgecolor='black')
plt.title("Trentino precipitation intensity histogram. \n Nbr of datapoints: {}".format(len(trentino["intensity"])))

Milano intensity classes:
[0, 1, 2]
Trentino intensity classes:
[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13]
```



Try to find correlation values between precipitation intensity and total communication activity

Result: Correlation value is -0.063, which is really low. According to this, there is no correlation between precipitation intensity and communication activity.

```
In [258]. # Milano
milano_corr = pd.DataFrame(milano[["intensity", "total"]])
print("Milano intensity vs activity correlation: {}".format(milano_corr().corr()))

# Trentino
trentino_corr = pd.DataFrame(trentino[["intensity", "total"]])
print("Trentino intensity vs activity correlation: {}".format(trentino_corr().corr()))

Milano intensity vs activity correlation:
intensity total
intensity -0.000000 -0.063826
total -0.063826 1.000000

Trentino intensity vs activity correlation:
intensity total
intensity -0.000000 -0.012159
total -0.012159 1.000000
```

```
In [259]. # Milano
corrmat_milano = df.milano.corr()
corr_features_milano = corrmat_milano.index

# Trentino
corrmat_trentino = df.trentino.corr()
corr_features_trentino = corrmat_trentino.index

plt.figure(figsize=(4,3))
plt.title("Milano Correlation matrix precipitation \n intensity and communication activity")
g = sns.heatmap(df.milano[corr_features_milano].corr(), annot=False, cmap="RdYlGn")

plt.figure(figsize=(4,3))
plt.title("Trentino Correlation matrix precipitation \n intensity and communication activity")
g = sns.heatmap(df.trentino[corr_features_trentino].corr(), annot=False, cmap="RdYlGn")
```



Make histograms of total activity with different precipitation intensities

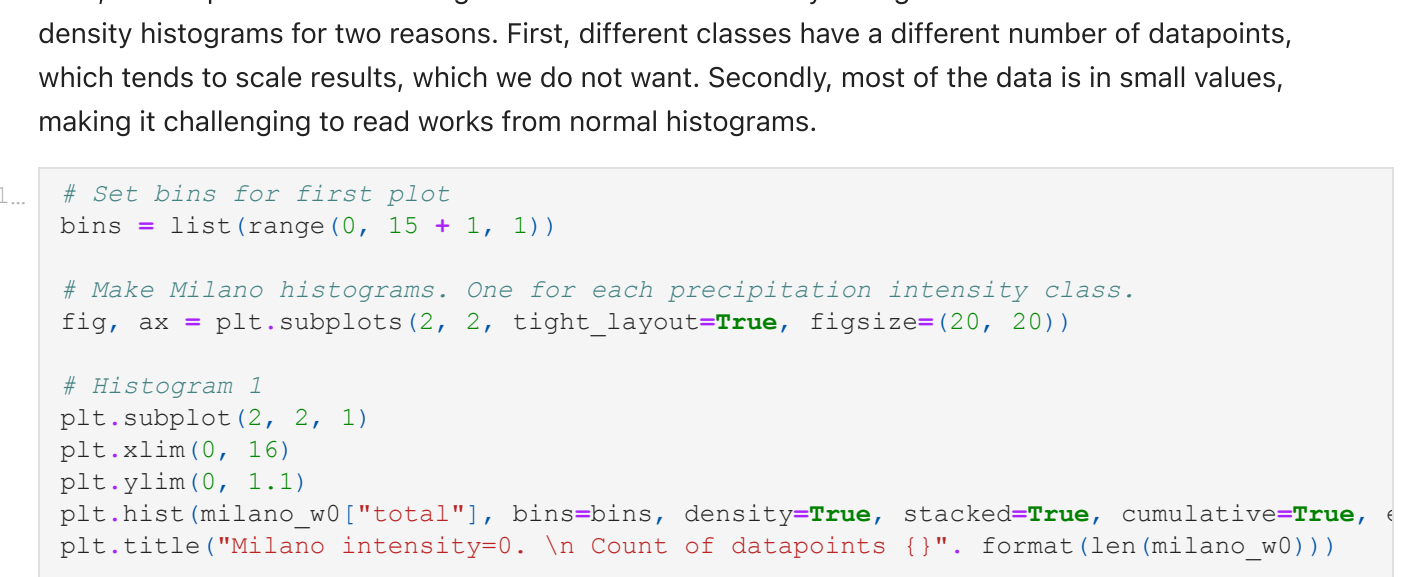
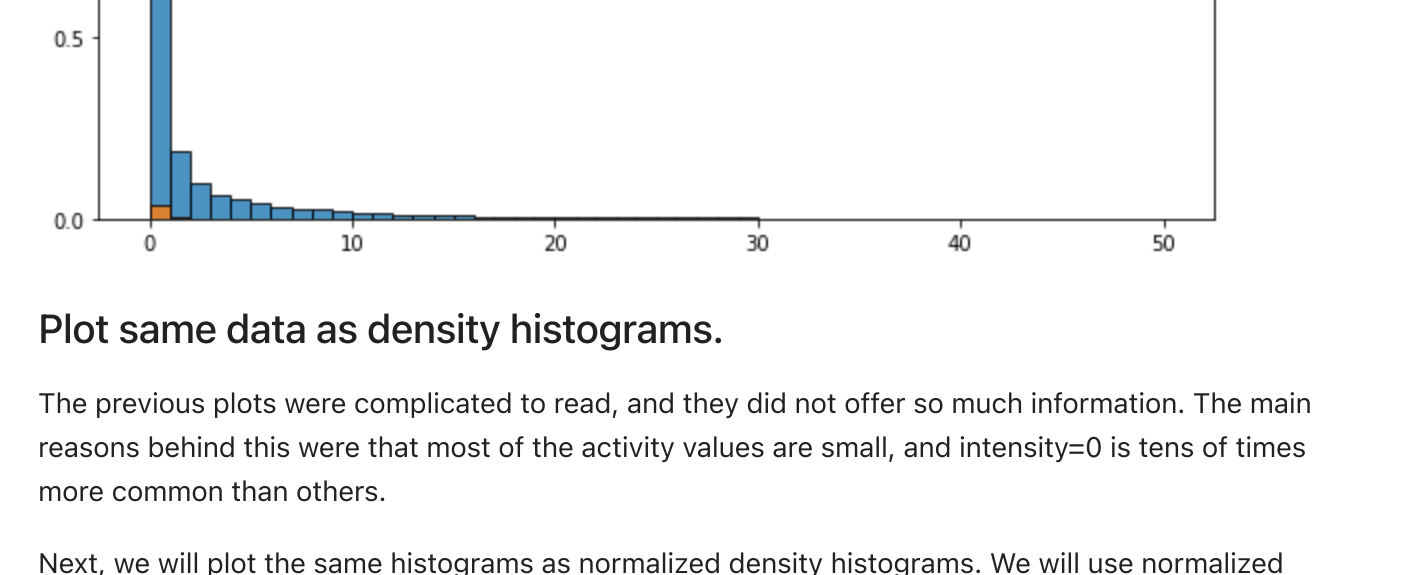
Check if there are differences in total communication activity during different levels of precipitation intensity.

```
In [260]. # Split Milano weather data, bases on intensity, to 4 groups (intensity classes 0 to 3)
milano_w0 = milano.filter(by(0, 'intensity'))
milano_w1 = milano.filter(by(1, 'intensity'))
milano_w2 = milano.filter(by(2, 'intensity'))
milano_w3 = milano.filter(by(3, 'intensity'))

# Split Milano weather data, bases on intensity, to 6 groups (intensity classes 1 to 18)
trentino_w1 = trentino.filter(trentino['intensity'] >= 1) & (trentino['intensity'] <= 3)
trentino_w2 = trentino.filter(trentino['intensity'] >= 4) & (trentino['intensity'] <= 6)
trentino_w3 = trentino.filter(trentino['intensity'] >= 7) & (trentino['intensity'] <= 9)
trentino_w4 = trentino.filter(trentino['intensity'] >= 10) & (trentino['intensity'] <= 12)
trentino_w5 = trentino.filter(trentino['intensity'] >= 13) & (trentino['intensity'] <= 15)
trentino_w6 = trentino.filter(trentino['intensity'] >= 16) & (trentino['intensity'] <= 18)

# Plot Milano telecom activity histogram
bins = list(range(0, 15 + 1, 1))
plt.rcParams["figure.figsize"] = (10, 10)
plt.hist(milano_w0["total"], bins=bins, alpha=0.8, edgecolor='black')
plt.hist(milano_w1["total"], bins=bins, alpha=0.8, edgecolor='black')
plt.hist(milano_w2["total"], bins=bins, alpha=0.8, edgecolor='black')
plt.hist(milano_w3["total"], bins=bins, alpha=0.8, edgecolor='black')
plt.title("Milano telecom activity distribution by precipitation intensity")
plt.legend(("intensity=0", "intensity=1", "intensity=2", "intensity=3"), loc="upper right")
plt.show()

# Plot Trentino telecom activity histogram
bins = list(range(0, 50 + 1, 1))
plt.rcParams["figure.figsize"] = (10, 10)
plt.hist(trentino_w1["total"], bins=bins, alpha=0.8, edgecolor='black')
plt.hist(trentino_w2["total"], bins=bins, alpha=0.8, edgecolor='black')
plt.hist(trentino_w3["total"], bins=bins, alpha=0.8, edgecolor='black')
plt.hist(trentino_w4["total"], bins=bins, alpha=0.8, edgecolor='black')
plt.hist(trentino_w5["total"], bins=bins, alpha=0.8, edgecolor='black')
plt.hist(trentino_w6["total"], bins=bins, alpha=0.8, edgecolor='black')
plt.title("Trentino telecom activity distribution by precipitation intensity")
plt.legend(("intensity=0", "intensity=1", "intensity=2", "intensity=3", "intensity=4", "intensity=5", "intensity=6"), loc="upper right")
plt.show()
```



Plot same data as density histograms.

The previous plots were complicated to read, and they did not offer so much information. The main reasons behind this were that most of the activity values are small, and intensity=0 is tens of times more common than others.

Next, we will plot the same histograms as normalized density histograms. We will use normalized density histograms for two reasons. First, different classes have a different number of datapoints, which tends to scale results, which we do not want. Secondly, most of the data is in small values, making it challenging to read works from normal histograms.

```
In [261]. # Set bins for first plot
bins = list(range(0, 15 + 1, 1))

# Make Milano histograms. One for each precipitation intensity class.
fig, ax = plt.subplots(2, 2, tight_layout=True, figsize=(20, 20))

# Histogram 1
plt.subplot(2, 2, 1)
plt.xlim(0, 16)
plt.ylim(0, 1.1)
plt.hist(milano_w0["total"], bins=bins, density=True, stacked=True, cumulative=True, label="intensity=0")
plt.title("Milano intensity=0. \n Count of datapoints {}".format(len(milano_w0)))

# Histogram 2
plt.subplot(2, 2, 2)
plt.xlim(0, 16)
plt.ylim(0, 1.1)
plt.hist(milano_w1["total"], bins=bins, density=True, stacked=True, cumulative=True, label="intensity=1")
plt.title("Milano intensity=1. \n Count of datapoints {}".format(len(milano_w1)))

# Histogram 3
plt.subplot(2, 2, 3)
plt.xlim(0, 16)
plt.ylim(0, 1.1)
plt.hist(milano_w2["total"], bins=bins, density=True, stacked=True, cumulative=True, label="intensity=2")
plt.title("Milano intensity=2. \n Count of datapoints {}".format(len(milano_w2)))

# Histogram 4
plt.subplot(2, 2, 4)
plt.xlim(0, 16)
plt.ylim(0, 1.1)
plt.hist(milano_w3["total"], bins=bins, density=True, stacked=True, cumulative=True, label="intensity=3")
plt.title("Milano intensity=3. \n Count of datapoints {}".format(len(milano_w3)))

plt.suptitle('Milano Telecom activity by precipitation intensity')
plt.show()

# Set bins for second plot
bins = list(range(0, 50 + 1, 1))

# Make Trentino histograms. One for each precipitation intensity class.
fig, ax = plt.subplots(3, 2, tight_layout=True, figsize=(20, 20))

# Histogram 1
plt.subplot(3, 2, 1)
plt.xlim(0, 51)
plt.ylim(0, 1.1)
plt.hist(trentino_w1["total"], bins=bins, density=True, stacked=True, cumulative=True, label="intensity=1")
plt.title("Trentino intensity=1. \n Count of datapoints {}".format(len(trentino_w1)))

# Histogram 2
plt.subplot(3, 2, 2)
plt.xlim(0, 51)
plt.ylim(0, 1.1)
plt.hist(trentino_w2["total"], bins=bins, density=True, stacked=True, cumulative=True, label="intensity=2")
plt.title("Trentino intensity=2. \n Count of datapoints {}".format(len(trentino_w2)))

# Histogram 3
plt.subplot(3, 2, 3)
plt.xlim(0, 51)
plt.ylim(0, 1.1)
plt.hist(trentino_w3["total"], bins=bins, density=True, stacked=True, cumulative=True, label="intensity=3")
plt.title("Trentino intensity=3. \n Count of datapoints {}".format(len(trentino_w3)))

# Histogram 4
plt.subplot(3, 2, 4)
plt.xlim(0, 51)
plt.ylim(0, 1.1)
plt.hist(trentino_w4["total"], bins=bins, density=True, stacked=True, cumulative=True, label="intensity=4")
plt.title("Trentino intensity=4. \n Count of datapoints {}".format(len(trentino_w4)))

# Histogram 5
plt.subplot(3, 2, 5)
plt.xlim(0, 51)
plt.ylim(0, 1.1)
plt.hist(trentino_w5["total"], bins=bins, density=True, stacked=True, cumulative=True, label="intensity=5")
plt.title("Trentino intensity=5. \n Count of datapoints {}".format(len(trentino_w5)))

# Histogram 6
plt.subplot(3, 2, 6)
plt.xlim(0, 51)
plt.ylim(0, 1.1)
plt.hist(trentino_w6["total"], bins=bins, density=True, stacked=True, cumulative=True, label="intensity=6")
plt.title("Trentino intensity=6. \n Count of datapoints {}".format(len(trentino_w6)))

plt.suptitle('Trentino Telecom activity by precipitation intensity')
plt.show()
```

/Users/saukk/opt/anaconda3/lib/python3.8/site-packages/matplotlib/axes/_axes.py:6664: RuntimeWarning: invalid value encountered in true_divide

tops = (tops / np.diff(bins)) / (tops[-1].sum())

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

Trentino Telecom activity by precipitation intensity

Milano Telecom activity by precipitation intensity

3. Station latitude
4. Station longitude
5. Pollution_type
6. Pollution_measure
7. Date format

Define function to read pollution data

Define function to read pollution data. Time is in format: 2013/12/31 04:00, but if using Turi read_csv to read multiple files, it drops the hours and minutes. This is because, some of the files have only date and no time, but we want to use both to increase accuracy.

Also change timestamp format to match format used in weather data: 2013/12/31 04:00 -> 201312310400.

```
In [265]: # Define function to read multiple telecom files and append them to one SFrame
def create_pollution_SFrame(paths):
    # Initialize results sf
    sf = SFrame({'id':[], 'timestamp':[], 'pollution':[]})
    sf['id'] = sf['id'].astype(int)
    sf['timestamp'] = sf['timestamp'].astype(str)
    sf['pollution'] = sf['pollution'].astype(float)

    # Read the files and append dataframes to result_sf
    for path in paths:
        sf_tmp = SFrame.read_csv(path, sep=',', header=None, column_type_hints=sf)
        sf_tmp = sf_tmp.rename({'X1': 'id', 'X2': 'timestamp', 'X3': 'pollution'})
        if len(sf_tmp['timestamp']) > 12:
            sf_tmp['id'] = sf_tmp['id'].astype(int)
            sf_tmp['timestamp'] = sf_tmp['timestamp'].astype(str)
            sf_tmp['pollution'] = sf_tmp['pollution'].astype(float)
        sf = sf.append(sf_tmp)

    # Change timestamp format
    sf['timestamp'] = sf['timestamp'].apply(lambda time: time.replace("/", ""))
    sf['timestamp'] = sf['timestamp'].apply(lambda time: time.replace("-", ""))
    sf['timestamp'] = sf['timestamp'].apply(lambda time: time.replace(":", ""))
    sf['timestamp'] = sf['timestamp'].astype(int)

    # Return results
    return sf
```

Read pollution data, weather data and pollution information data.

Pollution information data has coordinates of the stations and pollution types.

Grid data has information of weather data locations.

```
In [266]: # Define path to folder of Milano pollution data files
pollution_paths = glob2.glob("data/milano_pollution_data/*")

# Run milano_pollution function and get all Milano pollution data
pollution = create_pollution_SFrame(pollution_paths)

# Read weather files and rename columns
weather = SFrame.read_csv('data/MeteoMilano_01-11-13_01-01-14.csv', sep=',', header=None)
weather = weather.rename({'X1': 'timestamp', 'X2': 'quadID', 'X3': 'intensity', 'X4': 'cov'})

# Read pollution information file
pollution_info = SFrame.read_csv('data/milano_airQuality/milano_pollution_legend2.csv')
pollution_info = pollution_info.rename({'X1': 'id', 'X2': 'name', 'X3': 'lat', 'X4': 'lon', 'X5': 'pollution_type', 'X6': 'pollution_measure', 'X7': 'date_format'})

# Print first rows
pollution.print_rows(num_rows=5)
pollution_info.print_rows(num_rows=5)
weather.print_rows(num_rows=5)
```

Finished parsing file /Users/saukk/Desktop/Distributed data infrastructure/Grap hLab project/data/milano_pollution_data/mi_pollution_6956.csv
Parsing completed. Parsed 54 lines in 0.004946 secs.

Finished parsing file /Users/saukk/Desktop/Distributed data infrastructure/Grap hLab project/data/milano_pollution_data/mi_pollution_5722.csv
Parsing completed. Parsed 1353 lines in 0.005888 secs.

Finished parsing file /Users/saukk/Desktop/Distributed data infrastructure/Grap hLab project/data/milano_pollution_data/mi_pollution_6372.csv
Parsing completed. Parsed 1460 lines in 0.006122 secs.

Finished parsing file /Users/saukk/Desktop/Distributed data infrastructure/Grap hLab project/data/milano_pollution_data/mi_pollution_6366.csv
Parsing completed. Parsed 1312 lines in 0.006145 secs.

Finished parsing file /Users/saukk/Desktop/Distributed data infrastructure/Grap hLab project/data/milano_pollution_data/mi_pollution_5531.csv
Parsing completed. Parsed 1312 lines in 0.006096 secs.

Finished parsing file /Users/saukk/Desktop/Distributed data infrastructure/Grap hLab project/data/milano_pollution_data/mi_pollution_5725.csv
Parsing completed. Parsed 1312 lines in 0.012066 secs.

Finished parsing file /Users/saukk/Desktop/Distributed data infrastructure/Grap hLab project/data/milano_pollution_data/mi_pollution_10278.csv
Parsing completed. Parsed 1459 lines in 0.006753 secs.

Finished parsing file /Users/saukk/Desktop/Distributed data infrastructure/Grap hLab project/data/milano_pollution_data/mi_pollution_10279.csv
Parsing completed. Parsed 1461 lines in 0.008044 secs.

Finished parsing file /Users/saukk/Desktop/Distributed data infrastructure/Grap hLab project/data/milano_pollution_data/mi_pollution_5542.csv
Parsing completed. Parsed 1462 lines in 0.006125 secs.

Finished parsing file /Users/saukk/Desktop/Distributed data infrastructure/Grap hLab project/data/milano_pollution_data/mi_pollution_6328.csv
Parsing completed. Parsed 1460 lines in 0.006164 secs.

Finished parsing file /Users/saukk/Desktop/Distributed data infrastructure/Grap hLab project/data/milano_pollution_data/mi_pollution_6062.csv
Parsing completed. Parsed 1461 lines in 0.006669 secs.

Finished parsing file /Users/saukk/Desktop/Distributed data infrastructure/Grap hLab project/data/milano_pollution_data/mi_pollution_20020.csv
Parsing completed. Parsed 743 lines in 0.006328 secs.

Finished parsing file /Users/saukk/Desktop/Distributed data infrastructure/Grap hLab project/data/milano_pollution_data/mi_pollution_10280.csv
Parsing completed. Parsed 1352 lines in 0.006164 secs.

Finished parsing file /Users/saukk/Desktop/Distributed data infrastructure/Grap hLab project/data/milano_pollution_data/mi_pollution_10278.csv
Parsing completed. Parsed 1461 lines in 0.006261 secs.

Finished parsing file /Users/saukk/Desktop/Distributed data infrastructure/Grap hLab project/data/milano_pollution_data/mi_pollution_17127.csv
Parsing completed. Parsed 1484 lines in 0.006678 secs.

Finished parsing file /Users/saukk/Desktop/Distributed data infrastructure/Grap hLab project/data/milano_pollution_data/mi_pollution_20084.csv
Parsing completed. Parsed 692 lines in 0.005799 secs.

Finished parsing file /Users/saukk/Desktop/Distributed data infrastructure/Grap hLab project/data/milano_pollution_data/mi_pollution_6328.csv
Parsing completed. Parsed 1462 lines in 0.006378 secs.

Finished parsing file /Users/saukk/Desktop/Distributed data infrastructure/Grap hLab project/data/milano_pollution_data/mi_pollution_10283.csv
Parsing completed. Parsed 61 lines in 0.005863 secs.

Finished parsing file /Users/saukk/Desktop/Distributed data infrastructure/Grap hLab project/data/milano_pollution_data/mi_pollution_5855.csv
Parsing completed. Parsed 1462 lines in 0.00758 secs.

Finished parsing file /Users/saukk/Desktop/Distributed data infrastructure/Grap hLab project/data/milano_pollution_data/mi_pollution_17126.csv
Parsing completed. Parsed 1461 lines in 0.006261 secs.

Finished parsing file /Users/saukk/Desktop/Distributed data infrastructure/Grap hLab project/data/milano_pollution_data/mi_pollution_10273.csv
Parsing completed. Parsed 1484 lines in 0.006678 secs.

Finished parsing file /Users/saukk/Desktop/Distributed data infrastructure/Grap hLab project/data/milano_pollution_data/mi_pollution_5823.csv
Parsing completed. Parsed 1462 lines in 0.006112 secs.

Finished parsing file /Users/saukk/Desktop/Distributed data infrastructure/Grap hLab project/data/milano_pollution_data/mi_pollution_5834.csv
Parsing completed. Parsed 1462 lines in 0.006105 secs.

Finished parsing file /Users/saukk/Desktop/Distributed data infrastructure/Grap hLab project/data/milano_pollution_data/mi_pollution_17122.csv
Parsing completed. Parsed 59 lines in 0.005916 secs.

Finished parsing file /Users/saukk/Desktop/Distributed data infrastructure/Grap hLab project/data/milano_pollution_data/mi_pollution_6354.csv
Parsing completed. Parsed 1462 lines in 0.006828 secs.

Finished parsing file /Users/saukk/Desktop/Distributed data infrastructure/Grap hLab project/data/milano_pollution_data/mi_pollution_6340.csv
Parsing completed. Parsed 1061 lines in 0.005991 secs.

Finished parsing file /Users/saukk/Desktop/Distributed data infrastructure/Grap hLab project/data/milano_pollution_data/mi_pollution_5843.csv
Parsing completed. Parsed 1461 lines in 0.006239 secs.

Finished parsing file /Users/saukk/Desktop/Distributed data infrastructure/Grap hLab project/data/milano_pollution_data/mi_pollution_5804.csv
Parsing completed. Parsed 1461 lines in 0.006116 secs.

Finished parsing file /Users/saukk/Desktop/Distributed data infrastructure/Grap hLab project/data/milano_pollution_data/mi_pollution_5806.csv
Parsing completed. Parsed 1460 lines in 0.006011 secs.

Finished parsing file /Users/saukk/Desktop/Distributed data infrastructure/Grap hLab project/data/milano_pollution_data/mi_pollution_6344.csv
Parsing completed. Parsed 1461 lines in 0.007431 secs.

Finished parsing file /Users/saukk/Desktop/Distributed data infrastructure/Grap hLab project/data/MeteoMilano_01-11-13_01-01-14.csv
Parsing completed. Parsed 100 lines in 0.008735 secs.

Inferred types from first 100 line(s) of file as column_type_hints=[int,int,int,int,int]
If parsing fails due to incorrect types, you can correct the inferred type list above and pass it to read_csv in the column_type_hints argument

Finished parsing file /Users/saukk/Desktop/Distributed data infrastructure/Grap hLab project/data/MeteoMilano_01-11-13_01-01-14.csv
Parsing completed. Parsed 33316 lines in 0.012019 secs.

Finished parsing file /Users/saukk/Desktop/Distributed data infrastructure/Grap hLab project/data/milano_pollution_data/mi_pollution_5823.csv
Parsing completed. Parsed 36 lines in 0.005307 secs.

Inferred types from first 100 line(s) of file as column_type_hints=[int,str,float,float,str,str,str]
If parsing fails due to incorrect types, you can correct the inferred type list above and pass it to read_csv in the column_type_hints argument

Finished parsing file /Users/saukk/Desktop/Distributed data infrastructure/Grap hLab project/data/milano_pollution_data/mi_pollution_5823.csv
Parsing completed. Parsed 36 lines in 0.006151 secs.

```
| id | pollution | timestamp |
| 5722 | 10.0 | 201312310100 |
| 5722 | 9.0 | 201312310200 |
| 5722 | 9.0 | 201312310300 |
| 5722 | 10.0 | 201312310400 |
| 5722 | 10.0 | 201312310500 |
[40994 rows x 3 columns]
```

```
| id | name | lat | lon | pollution_type |
| 20020 | Milano -via Carlo Pascal | 45.478452 | 9.235016 | Ammonia |
| 17127 | Milano -via Mario Marce | 45.496067 | 9.193023 | Benzene |
| 17126 | Milano -via Carlo Pascal | 45.478452 | 9.235016 | Benzene |
| 6057 | Milano - via Senato | 45.47078 | 9.19718 | Benzene |
| 6062 | Milano - P.zza Zavattari | 45.476089 | 9.143509 | Benzene |
```

```
| pollution_measure | date_format |
| g/m | YYYY/MM/DD |
| g/m | YYYY/MM/DD HH24:MI |
| g/m | YYYY/MM/DD HH24:MI |
| g/m | YYYY/MM/DD HH24:MI |
| g/m | YYYY/MM/DD HH24:MI |
[36 rows x 7 columns]
```

```
| timestamp | quadID | intensity | coverage | prec_type |
| 201311010000 | 1 | 0 | 0 | 0 |
| 201311010000 | 2 | 0 | 0 | 0 |
| 201311010000 | 3 | 0 | 0 | 0 |
| 201311010000 | 4 | 0 | 0 | 0 |
| 201311010010 | 1 | 0 | 0 | 0 |
[33316 rows x 5 columns]
```

Join pollution and pollution_info to pollution, calculate weather quadrants average and join pollution and weather by timestamps.

```
In [267]: # Join pollution and pollution_info
pollution2 = pollution.join(pollution_info, on='[id]', how='inner')
```

```
# Drop columns from pollution2, that we do not need.
pollution2 = pollution2.remove_column('name')
pollution2 = pollution2.remove_column('pollution_measure')
pollution2 = pollution2.remove_column('date_format')
pollution2 = pollution2.remove_column('lat')
pollution2 = pollution2.remove_column('lon')
```

```
# Group weather by timestamp and take average of intensity and max of type.
# Weather data was in 4 quadrants, but because we do not have their
# coordinates, we will use avg intensity and max type to whole Milano.
weather2 = weather.groupby(['timestamp'], ['intensity':agg.AVG('intensity'),
'prec_type':agg.MAX('prec_type')])
```

```
# Join weather and pollution data
milano = pollution2.join(weather2, on='timestamp', how='inner')
```

```
# Add columns weekday, time and month. Monday and Sunday=7.
milano['timestamp'] = milano['timestamp'].astype(str)
milano['weekday'] = milano['timestamp'].apply(lambda time: (datetime.date(int(time[0:4]),int(time[5:6]),int(time[7:8]))).weekday()+1)
milano['month'] = milano['timestamp'].apply(lambda time: time[4:6])
milano['month'] = milano['month'].astype(int)
milano['time'] = milano['timestamp'].apply(lambda time: time[-4:])
milano['time'] = milano['time'].astype(int)
```

```
# Drop id and timestamp
milano = milano.remove_column('id')
milano = milano.remove_column('timestamp')
```

```
# Print first 5 rows of both
milano.print_rows(num_rows=5)
```

```
| pollution | pollution_type | intensity | prec_type | weekday | month | time | |
| 10.0 | Ozono | 0.0 | 0 | 0 | 12 | 100 |
| 9.0 | Ozono | 0.0 | 0 | 0 | 2 | 12 | 200 |
| 9.0 | Ozono | 0.0 | 0 | 0 | 2 | 12 | 300 |
| 10.0 | Ozono | 0.0 | 0 | 0 | 2 | 12 | 400 |
| 9.0 | Ozono | 0.0 | 0 | 0 | 2 | 12 | 500 |
[37827 rows x 7 columns]
```

Make correlation matrix for each pollution type and find if precipitation has correlation to pollution

```
In [268]: # Convert to dataframe and drop coordinates
df_milano = pd.DataFrame(milano)

# Ozono is in the data with names Ozono and Ozono. Rename Ozono to be Ozono.
df_milano = df_milano.replace({'pollution_type': 'Ozono'}, 'Ozono')
```

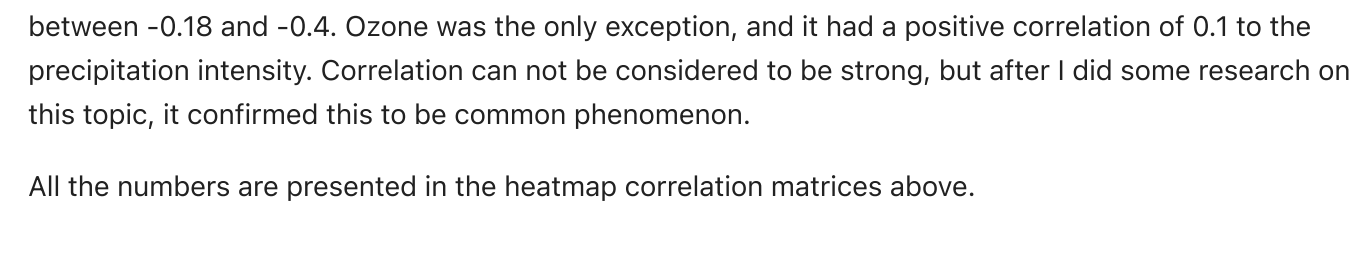
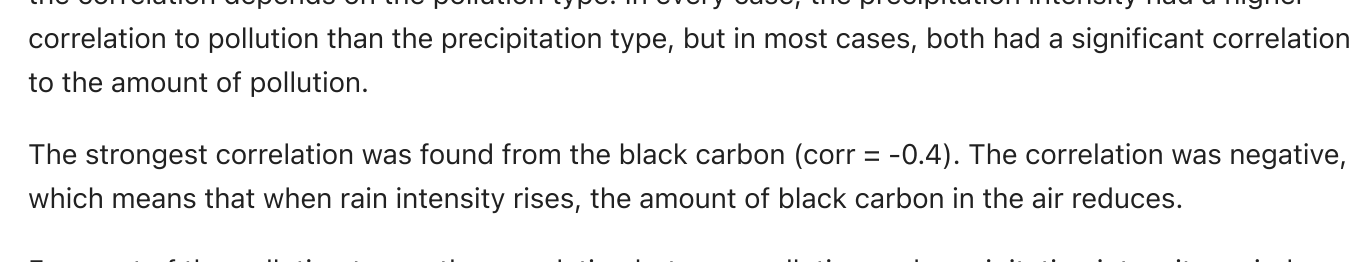
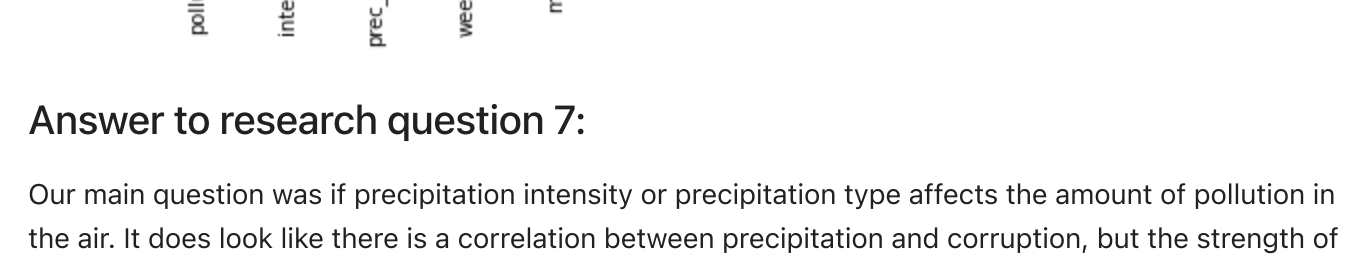
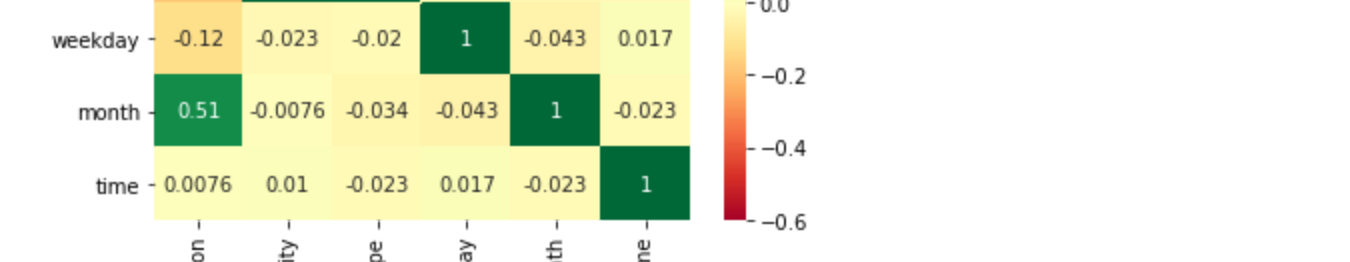
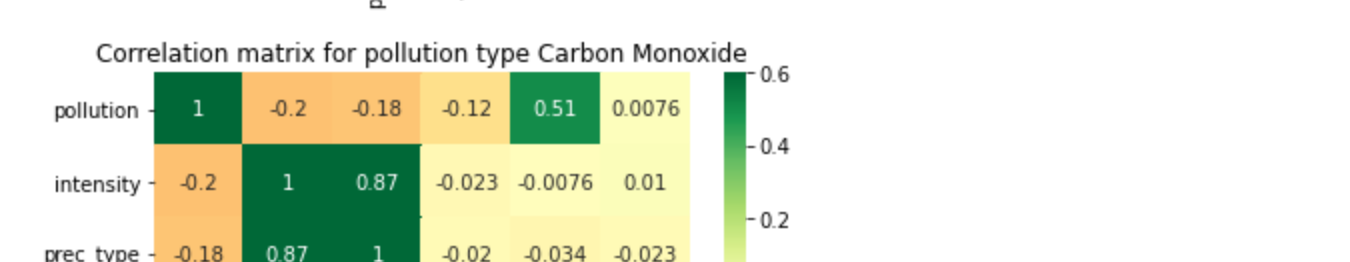
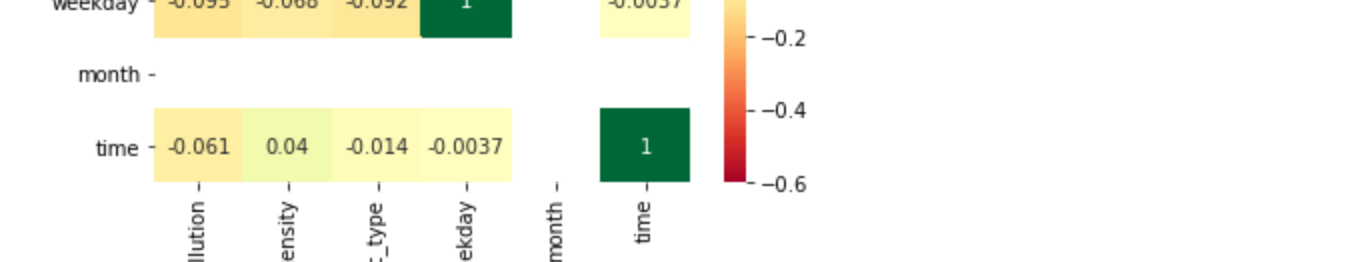
```
# Make list of pollution types
pol_types = list(df_milano['pollution_type'].unique())

# Loop through pollution types and make correlation matrix for each type
for pol_type in pol_types:
    # Make tmp dataframe with given pollution type
    df_tmp = df_milano[df_milano['pollution_type'] == pol_type]
    df_tmp = df_tmp.drop(columns=['pollution_type'])

    # Make correlation matrix
    corr_mat = df_tmp.corr()

    # Make correlation matrix
    corr_features = corr_mat.index

    plt.figure(figsize=(6,4))
    plt.title("Correlation matrix for pollution type {}".format(pol_type))
    sns.heatmap(df_tmp[corr_features].corr(), annot=True, cmap="RdYlGn", vmin=-0.6,
```



Answer to research question 7:

Our main question was if precipitation intensity or precipitation type affects the amount of pollution in the air. It does look like there is a correlation between precipitation and corruption, but the strength of the correlation depends on the pollution type. In every case, the precipitation intensity had a higher correlation to pollution than the precipitation type, but in most cases, both had a significant correlation to the amount of pollution.

The strongest correlation was found from the black carbon (corr = -0.4). The correlation was negative, which means that when rain intensity rises, the amount of black carbon in the air reduces.

For most of the pollution types, the correlation between pollution and precipitation intensity varied between -0.18 and -0.4. Ozono was the only exception, and it had a positive correlation of 0.1 to the precipitation intensity. Correlation can not be considered to be strong, but after I did some research on this topic, it confirmed this to be common phenomenon.

All the numbers are presented in the heatmap correlation matrices above.