## POLITECNICO DI TORINO

Master of science ICT FOR SMART SOCIETIES

Master's Degree thesis

# Big Data analysis of Floating Car Data to identify traffic congestion in urban areas



Supervisors

Prof. Danilo Giordano Prof. Luca Vassio Prof. Marco Diana Candidate Seyedamirmohammad Sakaki

October 2022

## Summary

The inescapable phenomenon of traffic congestion has been steadily worsening over the last several decades, particularly in big urban regions that have been experiencing population growth all over the globe. In its most basic form, congestion occurs when there is a larger demand for space than there is capacity on the roadway to accommodate it. However, these activities, which force people to interact with one another, are necessary for the effective operation of economic systems and cannot be avoided. In today's societies, these demands are frequently generated at the same time, while a large number of people intend to commute to school or their workplace as well as other movements to provide goods and services. Although there is a potential for congestion with every method of transportation, the emphasis of this thesis will be on the congestion that occurs with automobiles on public highways and streets.

The identification of regions that are prone to traffic congestion has been and will continue to be an essential step in the process of developing urban transportation systems. The technologies used to gather data on traffic have seen significant development over the last several years, and access to real-time traffic information is increasingly becoming the norm around the globe. Floating Car Data (FCD) is one of the most well-liked ways, and it has been more widespread in usage over the last several years. FCD is often timestamped geo-localization and speed data that is directly acquired by moving cars utilizing "in-vehicle" devices via mobile phones or GPS. These devices are connected to the internet. Because of the growing quantity of data that are acquired by FCD, it is becoming more desired to find and extract useful traffic-related information from the accumulated historical dataset. Examples of this kind of information include patterns of congestion. Despite this, it is not an easy task since the dataset is rather large, and the characteristics of the traffic, such as its complexity and dynamics, make it quite difficult.

Using Data-Driven approaches and an already-collected large FCD dataset, this thesis aims to identify traffic congestion areas in an urban area. This will be accomplished by extracting the desired meaningful information for the purpose of the research from a large FCD dataset. The collected FCD includes information spanning the entirety of 2019 and includes millions of records from 290,185 distinct vehicles, of which 92% are declared to be personal vehicles and the remaining 8% fall under the category of fleet vehicles. In addition, the vast majority of the data are collected at intervals of one minute for each vehicle. The FCD were

collected over the course of the 2019 calendar year. From these data, the segments of the trips of distinct cars are extracted. Thanks to a grid map consisting of cells that are each uniquely identified, the extracted segments are assigned to their corresponding cells such that each cell could provide key required information for traffic congestion calculation at the next step. In this project, the relative and absolute traffic congestion of morning and evening peak hours compared to off-peak hours is discussed and analyzed based on the average speed (in kilometers per hour) of vehicles on their respective segments of the route. Finally, the traffic congestion severity is extracted and polished using Kernel Density Estimation to demonstrate in a convenient way, how much different zones of the urban area are suffering from traffic congestion 1, so the future transportation system developments could be planned accordingly.



Figure 1: Highlight of the most congested areas, Bandwidth = 0.001

## Acknowledgements

My studies at Politecnico di Torino with the purpose of acquiring a Master of Science (MSc) degree in ICT for smart societies have come to a close with the completion of this master thesis report. During the time that I was working on my thesis, I went through an experience that was challenging, but more than that, it was intellectually engaging. Without the substantial and unwavering aid and guidance provided by a number of different persons, it would not have been feasible to do this project. To begin, I would like to offer my most sincere gratitude to Professor Giordano at Politecnico di Torino, who serves as my principal advisor for providing comments on my work and aid me in finding solutions to the various problems I encountered throughout the development process. I also would like to express my sincere appreciation to professors Vassio and Diana for the direction they provided me in determining the objectives of my study and how to assess strategy.

Seyedamirmohammad Sakaki

## Contents

List of Tables 8				
$\mathbf{Li}$	st of	Figures	9	
A	crony	ms 1	2	
1	Intr	oduction 1	5	
	1.1	Traffic congestion as a growing problem	5	
	1.2	Traffic data collection methods	6	
	1.3	IoT sensors, big data and data analysis development	7	
	1.4	Thesis organization	.7	
<b>2</b>	Bac	ground and concept definition 1	9	
	2.1	The Floating Car Data (FCD)	9	
		2.1.1 GPS-based FCD	20	
		2.1.2 Cellular network based FCD	20	
	2.2	Topographical Concepts	22	
		2.2.1 Geographical coordinate system and WGS84	22	
		2.2.2 Haversine Distance	22	
		2.2.3 Geometric Dilution Of Precision	23	
		2.2.4 DOP value classification	24	
	2.3	ICT background	25	
		2.3.1 Python	25	
	2.4	Related works	27	
	2.5	Dataset	29	
3	Imp	lementation and Results	31	
	3.1	Input FCD analysis	31	
		3.1.1 Data cleaning	31	
		3.1.2 Area of interest $\ldots \ldots \ldots \ldots \ldots \ldots \ldots \ldots \ldots \ldots 3$	32	
		3.1.3 Data behaviour	33	

	3.2 Segment Extraction			36	
		3.2.1	Segment definition	36	
		3.2.2	Segment extraction methodology	37	
		3.2.3	Segment table	41	
		3.2.4	Segment distribution	43	
	3.3 Grid map		44		
		3.3.1	Grid size selection	44	
3.4 Grid table			able	46	
		3.4.1	Grid plots	47	
3.5 Traffic congestion index			congestion index	50	
		3.5.1	Absolute traffic congestion	52	
		3.5.2	Relative traffic congestion	54	
		3.5.3	Cells with negative congestion	55	
	3.6	6 Identification of congested areas			
		3.6.1	Kernel Density Estimation: a method for smoothing	61	
4	Conclusions 7			71	
Bibliography 7					
$\mathbf{A}$	A Source code 7				

## List of Tables

DOP value classification	25
FCD dataset description	29
Segment table description, Other properties derived directly from	
the latitude and longitude coordinates have been omitted for the	
purpose of brevity.	42
Grid table. Obtained by outer join of 3.1 and 3.3	47
	DOP value classification

## List of Figures

1	Highlight of the most congested areas, Bandwidth = $0.001 \dots$	3
2.1	Communication from GPS, source [14]	20
2.2	Communication from cellular phones, source $[14]$	21
3.1	The area of interest, including Turin and nearby cities	33
3.2	Cumulative Distribution Function of GPS accuracy DOP	34
3.3	Distribution of GPS accuracy DOP	34
3.4	Daily Records Distribution	35
3.5	Hourly Records Distribution	36
3.6	Cumulative distribution of Segment duration's	38
3.7	Distribution of Segment duration's	38
3.8	Cumulative distribution of Segment Speed (KM/H) of different groups	39
3.9	Cumulative distribution of Segment Speed (KM/H)	40
3.10	Distribution of Segment Speed (KM/H)	41
3.11	Daily Segments Distribution	43
3.12	Hourly Segments Distribution	44
3.13	Segment density VS cell size	45
3.14	Required processing time VS cell size	46
3.15	Number of segments per cell	48
3.16	Mean of the speeds of the segments per cell	49
3.17	Mean of GPS accuracy DOP per cell	50
3.18	Peak and off-peak hours selection	51
3.19	Cell-wise cumulative distribution of absolute traffic congestion	53
3.20	Cell-wise distribution of absolute traffic congestion	53
3.21	Cell-wise cumulative distribution of relative traffic congestion	55
3.22	Cell-wise distribution of relative traffic congestion	55
3.23	Cumulative distribution of cell-wise traffic congestion	56
3.24	Distribution of cell-wise traffic congestion	57
3.25	comparison of cell-wise cumulative distribution of relative traffic	
	congestion in filtered and not filtered scenarios	58
3.26	Filtered cell-wise cumulative distribution of relative traffic congestion	59

Filtered cell-wise distribution of relative traffic congestion	59
Relative traffic congestion	60
KDE of Relative traffic congestion, Bandwidth = $0.0005$	62
KDE of Relative traffic congestion, Bandwidth = $0.001 \dots \dots \dots$	63
KDE of Relative traffic congestion, Bandwidth = $0.0015$	64
KDE of Relative traffic congestion, Bandwidth = $0.002 \ldots \ldots$	65
KDE of Relative traffic congestion, Bandwidth = $0.0025$	66
KDE of Relative traffic congestion, Bandwidth = $0.003 \ldots \ldots$	67
KDE of Relative traffic congestion, Bandwidth = $0.0035$	68
Highlight of the most congested areas, Bandwidth = $0.001 \dots \dots$	69
	Filtered cell-wise distribution of relative traffic congestionRelative traffic congestionKDE of Relative traffic congestion, Bandwidth = 0.0005KDE of Relative traffic congestion, Bandwidth = 0.001KDE of Relative traffic congestion, Bandwidth = 0.0015KDE of Relative traffic congestion, Bandwidth = 0.0025KDE of Relative traffic congestion, Bandwidth = 0.0025KDE of Relative traffic congestion, Bandwidth = 0.003KDE of Relative traffic congestion, Bandwidth = 0.001

## Acronyms

### AI

artificial intelligence

### FCD

Floating Car Data

### IoT

Internet of Things

### $\mathbf{GPS}$

Global Positioning system

### $\mathbf{GSM}$

Global System for Mobile communications

### GPRS

General Packet Radio Service

### UMTS

Universal Mobile Telecommunications System

### ITS

Intelligent Transportation Systems

### WGS

World Geodetic System

### GDOP

Geometric Dilution Of Precision

### HDOP

Horizontal Dilution Of Precision

### GNSS

Global Navigation Satellite System

### $\mathbf{OTM}$

Open Transport Map

### $\mathbf{RTMS}$

Remote Traffic Microwave Sensor Data

### LOS

Level Of Service

### KDE

Kernel Density Estimation

# Chapter 1 Introduction

### 1.1 Traffic congestion as a growing problem

Congestion is a condition in traffic flow in which an increased volume of traffic on a highway, or "throughput" facility (such as a tunnel), leads to greater queue lengths and therefore increased travel time, lower average speed and longer trip times. Congestion occurs when traffic demand is high enough that interactions among vehicles reduce the velocity of the traffic stream. Since the 1950s, there has been a significant increase in traffic congestion on metropolitan road networks, [1] especially during rush hour. Extreme traffic congestion occurs when demand exceeds the capacity of a road (or of the junctions along the road), that leads to a condition under which the vehicles are completely halted for a period of time [2] [3]. A traffic jam can last anywhere from a few minutes to a several hours.

As the traffic congestion tends to be more severe year by year, there are some major reasons causing this increase. The most obvious one is population expansion. More people in an affluent country implies more vehicles; However, overall vehicle kilometers have increased far faster than the population. For instance, Between 1980 and 2000, the overall population of the United States increased by 24%, while total car kilometers driven increased by 80% due to more intensive usage of each vehicle [4]. Other main reasons as [5] mention, are the following:

- Economic growth
- Increase in the number of vehicles
- Insufficient road capacity
- Ineffective city planning and management
- Accidents

• Natural disasters

### **1.2** Traffic data collection methods

In recent decades, different approaches to traffic data collection have been employed to collect and analyze relevant data to extract congestion patterns in order to manage and decrease traffic. Data analytic has proven to be crucial in determining various causes of congestion based on data from sensors and cameras placed on roads, GNSS data collected from cars, cell phone signal strength in an area, which also reflects amount of cars and traffic around it, etc. In general there are two principal methods of traffic data collection including data collection by surveys and data collection by using the records of some specific devices [6]. The latter one is known as 'in-situ' method and itself is consist of two sub methods of intrusive and non-intrusive [7]. The intrusive approaches essentially consist of a data recorder and a sensor placed on or in the road:

- Pneumatic road tubes [8]
- Piezoelectric sensors [9]
- Inductive loop [10]

On the other hand, Remote observations are the foundation of non-intrusive approaches. Even while manual counting is the most often used approach, new technologies that appear to be highly promising have lately emerged:

- Manual counts
- Passive and active infra-red
- Passive magnetic
- Microwave radar
- Ultrasonic and passive acoustic
- Video image detection

Profile surveys are used by experts to assist evaluate the intensity of traffic on a route and the direction of traffic flow. These surveys are frequently set up at junctions to monitor one or more lanes of traffic going in different directions for the most effective data collecting 6.

- **Profile Surveys**: A profiling survey is a type of manual counting in which the data is collected by a human. Most surveyors will design a paper form (in some cases substituded by mobile apps) to assist differentiate between different sorts of vehicle classifications, the direction or lane in which the vehicle is driving, and other pre-selected time intervals.
- **Directional Surveys**: consist of several human surveyors covering various directions along a single length of road or crossroads counting pre-specified criteria

### 1.3 IoT sensors, big data and data analysis development

The Internet of Things, also known as IoT, is a system that consists of interconnected computer devices, mechanical and digital equipment, items, or people that are given the capacity to communicate data across a network with almost no involvement from humans [11]. This implies the installation of various types of sensors on a variety of objects that, thanks to the concept of data analysis, are involved in various life scenarios, such as those involving individuals and businesses: household appliances that make up the so-called Smart Home; manufacturing machines in Industry 4.0; vehicles and smart roads, all of which are integrated to form Smart Cities. The steady increase in the installation of Internet of Things (IoT) objects leads to the consequent generation of Big Data. Big Data is a collection of data that is so large, varied, and rapidly transmitted that it is extremely challenging, if not nearly impossible, to analyze using conventional techniques. It has been said that given the right paradigms, it is feasible to filter, arrange, analyze, and get value from them. Therefore, Internet of Things and Big Data are deeply intertwined in this setting.

The dissertation does a wonderful job of addressing the issues raised by the present case study. After a considerable amount of data has been gathered from a large number of cars, it is sent to a data storage facility over the internet. From there, it is eventually analyzed by means of data analysis tools.

### 1.4 Thesis organization

The thesis is structured as follows:

Backgrounds and definitions of topics used in this thesis are presented in Chapter 2, which includes mathematical background, relevant topological explanations for

the thesis purpose, and technological background to understand how the thesis was constructed.

The implementation of the thesis is discussed in Chapter 3, beginning with pre-processing on initial input data to extract required information for the following step, segment extraction. The projection of a grid of cells on the map is then explained, followed by the methodology for assigning each segment to a grid cell. Furthermore, how to use the cells to calculate and display traffic congestion on a map, and eventually extract the pattern of the congestion, so that it can be used for future transportation development projects.

### Chapter 2

# Background and concept definition

### 2.1 The Floating Car Data (FCD)

Floating car data, also known as FCD, is a term used to describe a technology that gathers information on the status of the traffic from a series of individual vehicles that float in the traffic. Each vehicle, which can be thought of as a moving sensor that operates in a distributed network, is outfitted with positioning (GPS) and communication (GSM, GPRS, UMTS, etc.) systems. These systems allow each vehicle to send information about its location, speed, and direction to a centralized control unit, which then uses this information to make decisions [12] or simply just stores in a databases.

Since FCD systems transcend the limits of fixed traffic monitoring methods (Poor flexibility, high installation and maintenance costs, etc.), they are being employed in a number of crucial application, therefore, FCD is a replacement or supplement source of high-quality data to existing technologies. They will contribute to the transportation system's safety, efficiency, and dependability. They are becoming increasingly important in the creation of future *Intelligent Transportation Systems* (ITS). In this Thesis, a dataset of collected FCD is being exploit as the initial entry data.

In its most basic forms, there are essentially two primary categories of FCDs, namely GPS and cellular-based systems [7]:

### 2.1.1 GPS-based FCD

Even though the Global Positioning System (GPS) is becoming more popular and more inexpensive, only a small percentage of vehicles, generally those used for fleet management services, are now equipped with this system (e.g. taxi drivers). Due to the usage of dual frequency receivers that make use of several constellations in order to determine position, the level of accuracy achieved is quite high and stable [13] and is often less than 30 meters.



Figure 2.1: Communication from GPS, source [14]

In most cases, traffic statistics collected from individual automobiles or trucks are preferable for use in highways and rural areas. In the event of urban traffic, taxi fleets are especially valuable owing to the large number of vehicles in the fleet as well as the on-board communication technologies that are already in place. In the present day, GPS probe data are widely used as a source of real-time information by a variety of service providers. Despite this widespread use, there are only a limited number of vehicles equipped with the necessary equipment, and the cost of such equipment is relatively high in comparison to those associated with floating cellular data.

### 2.1.2 Cellular network based FCD

Given that the vast majority of cars on the road nowadays are equipped with at least one or more mobile phones, it may be desirable to use mobile phones to carryout anonymous traffic surveys. The network receives periodic transmissions of the mobile phone's location, which are often accomplished by triangulation but may also be accomplished through other methods (e.g. handover). Mobile phones are required to be powered on, although this does not always imply that they are being used. Due to the shorter distance that separates each antenna in this system, it is especially well suited to the delivery of information that is generally precise in metropolitan areas, which is where traffic statistics are required the most.



Figure 2.2: Communication from cellular phones, source [14]

There is no need for any specialized hardware or software to be installed in moving vehicles, in contrast to fixed traffic detectors and systems that are based on GPS. Additionally, there is no requirement for any unique infrastructure to be erected along the route. As a result, it is more cost-effective than traditional detectors while also providing more capabilities in terms of coverage. Instead of obtaining data at discrete points, continuous data collection is done on the traffic. It takes far less time to set up, it is less difficult to install, and it significantly requires less maintenance. In spite of the fact that the location accuracy is normally rather poor (about 300 meters), this deficiency is largely compensated for by the huge number of devices. It is important to take into account that the UMTS technology should be used in order to collect more precise data. At the moment, FCD is participating in many applications all over the globe that deal with the management and information on real-time traffic.

### 2.2 Topographical Concepts

### 2.2.1 Geographical coordinate system and WGS84

The geographic coordinate system, also known as the GCS, is a spherical or ellipsoidal coordinate system that is used to measure and communicate locations directly on the surface of the Earth as latitude and longitude [15]. It is the simplest, the oldest, and the one that is utilized the most commonly out of all of the many spatial reference systems that are now in use, and it serves as the basis for the majority of the others [16]. The World Geodetic System, also abbreviated as WGS, is a standard that is used in the fields of cartography, geodesy, and satellite navigation, including GPS. The National Geospatial-Intelligence Agency of the United States has been responsible for the establishment and maintenance of WGS since 1984. The most recent modification, which took place in January 2021, is WGS 84, which is also known as the WGS 1984 ensemble [17] [18]. Decimal degrees, often known as DD, are a method of expressing geographic coordinates such as latitude and longitude as decimal fractions of a degree. GPS devices and various geographic information systems (GIS), online mapping programs like OpenStreetMap, and other mapping applications typically employ DD.

### 2.2.2 Haversine Distance

The Haversine formula is used to calculate the distance in terms of great circles that separates two locations on a sphere, given the longitudes and latitudes of those places. It is a specific example of a more general formula in spherical trigonometry called the rule of Haversines, which relates the sides and angles of spherical triangles. This law is important in navigation, as it relates the sides and angles of spherical triangles.

Haversine function  $hav(\theta)$  is:

$$hav(\theta) = \sin^2\left(\frac{\theta}{2}\right) = \frac{1 - \cos(\theta)}{2} \tag{2.1}$$

where:

 $hav(\theta)$  is the central angle between any two points on a sphere.

By performing *arcsine* on 2.1 the equation could be solved to obtain distance d:

$$d = 2 \cdot r \cdot \arcsin \sqrt{hav} \left(\varphi_2 - \varphi_1\right) + \left(1 - hav \left(\varphi_1 - \varphi_2\right) - hav \left(\varphi_1 + \varphi_2\right)\right) \cdot hav \left(\lambda_2 - \lambda_1\right)$$
(2.2)

where:

 $\lambda_1, \lambda_2$  are the longitude of point 1 and longitude of point 2  $\varphi_1, \varphi_2$  are the latitude of point 1 and latitude of point 2.

### 2.2.3 Geometric Dilution Of Precision

The term "dilution of precision," also known as "geometric dilution of precision" (GDOP), is used in satellite navigation and geomatics engineering to describe error propagation as a mathematical effect of navigation satellite geometry on positional measurement precision. The purpose of the GDOP is to specify how inaccuracies in the measurement will have an impact on the estimate of the final state. One possible explanation for this is:

$$GDOP = \frac{\Delta(OutputLocation)}{\Delta(MeasuredData)}$$
(2.3)

It is claimed that the geometry is weak and the DOP value is high when visible navigation satellites are close together in the sky; on the other hand, it is said that the geometry is strong and the DOP value is low when the satellites are far away. Imagine two circles, or annuli, that overlap one another and have different centers. When they overlap each other at right angles, the maximum extent of the overlap is much less than when they overlap each other in a nearly parallel fashion. Because of the greater angular separation between the satellites that are used to determine a unit's location, a lower DOP value indicates a higher level of positional accuracy than a higher number. Other factors, like nearby obstacles like mountains or buildings, may also cause the effective DOP to go up [**35**].

DOP can be expressed as a number of separate measurements:

- HDOP Horizontal Dilution Of Precision
- VDOP Vertical Dilution Of Precision
- PDOP Position (3D) Dilution Of Precision
- TDOP Time Dilution Of Precision
- GDOP Geometric Dilution Of Precision

DOP calculation for satellite i, consider:

$$R_i = \sqrt{(x_i - x)^2 + (y_i - y)^2 + (z_i - z)^2}$$
(2.4)

where:

x, y, z remark about the location of the receiver, and  $x_i, y_i, z_i(\theta)$  represent the location of satellite i. Matrix A is formulated as follow:

$$A = \begin{bmatrix} \frac{(x_1 - x)}{R_1} & \frac{(y_1 - y)}{R_1} & \frac{(z_1 - z)}{R_1} & -1\\ \frac{(x_2 - x)}{R_2} & \frac{(y_2 - y)}{R_2} & \frac{(z_2 - z)}{R_2} & -1\\ \frac{(x_3 - x)}{R_3} & \frac{(y_3 - y)}{R_3} & \frac{(z_3 - z)}{R_3} & -1\\ \frac{(x_4 - x)}{R_4} & \frac{(y_4 - y)}{R_4} & \frac{(z_4 - z)}{R_4} & -1 \end{bmatrix}$$
(2.5)

The first three elements in each row of A are the components of a unit vector from the receiver to the specified satellite. The last element of each row is the partial derivative of pseudorange with respect to the clock bias of the receiver. Formulate Q as the covariance matrix with the fewest squares that comes from the normal matrix:

The elements of Q are designated as:

$$Q = \begin{bmatrix} \sigma_x^2 & \sigma_{xy} & \sigma_{xz} & \sigma_{xt} \\ \sigma_{xy} & \sigma_y^2 & \sigma_{yz} & \sigma_{yt} \\ \sigma_{xz} & \sigma_{yz} & \sigma_z^2 & \sigma_{zt} \\ \sigma_{xt} & \sigma_{yt} & \sigma_{zt} & \sigma_t^2 \end{bmatrix}$$
(2.6)

Eventually, GDOP and HDOP are formulated as [36]:

$$GDOP = \sqrt{\sigma_x^2 + \sigma_y^2 + \sigma_z^2 + \sigma_t^2} \tag{2.7}$$

$$HDOP = \sqrt{\sigma_n^2 + \sigma_e^2} \tag{2.8}$$

depends on the coordinate system used To match the local horizon plane and vertical in a north, east, or up coordinate system.

### 2.2.4 DOP value classification

[37] interprets the obtained DOP value and classifies it as follows:

DOP value	Rating	Description
	Ideal	Finest confidence level
<1		to be utilized for applications
		that need the highest accuracy at all times.
	Excellent	At this confidence level
1-2		location measurements are deemed precise enough
		for all applications save the most sensitive ones.
		Represents a minimal acceptable threshold for making
25	Good	correct selections. Positional measurements might
2-0		be utilized to provide the user with credible in-route
		navigation recommendations.
		Although positional measurements might be utilized
5 - 10	Moderate	for computations, the quality of the fix could still
		be improved. A more open view of the sky is advised.
		Represents a low degree of confidence. Positional
10-20	Fair	measurements should be ignored or only utilized to
		provide an approximate estimation of the present position.
>20	Poor	At this level, readings using a 6-meter accurate equipment
ZU	1 001	are erroneous by up to 300 meters and should be ignored.

 Table 2.1: DOP value classification

### 2.3 ICT background

### 2.3.1 Python

Python is a general-purpose programming language, which means that it may be used to develop a wide range of different applications and is not specifically tailored to solve any particular issues. Because of its flexibility and the fact that it is easy to learn, it has quickly become one of the most widely used programming languages today. Developed by a group of programmers led by British scientist, Guido van Rossum, Python is actually an interpreted language. The design concept behind it places an emphasis on the readability of the code by making extensive use of indentation [19].

The language has applications in a variety of domains, including Web Development, Data Science and Machine Learning, and Simulation; all of these domains, as well as others, are being investigated as part of the same research effort that this thesis is a part of. As a result of this, Python has been preferred over other languages like R. There are several libraries that can be imported, which is another selling feature for the product. In direct relation to this body of work, the most significant ones that have been applied are as follows:

- **Pandas**: A software library designed for data manipulation and analysis. In particular, it provides data structures and functions for the manipulation of numerical tables and time series [20] [21].
- Geooandas: a free and open-source tool that simplifies the process of dealing with geographical data in Python. The datatypes that pandas uses are extended by GeoPandas so that spatial operations may be performed on geometric types. Shapely is responsible for carrying out geometric operations [22] [23].
- Numpy: Adding support for big, multi-dimensional arrays and matrices, as well as a vast set of high-level mathematical functions to work on these arrays. The "ndarray" data structure, which stands for "n-dimensional array," lies at the heart of NumPy's fundamental capabilities. These arrays have a consistent data type across the board [24].
- Matplotlib: a complete library for the creation of interactive, animated, and static visualizations. [25] [26]
- Shapely:Python package with the Shapely name that allows for the manipulation and analysis of planar geometric objects under the BSD license. Although Shapely is not concerned with specific data formats or coordinate systems, it can be easily integrated with other packages that are [27].
- Folium: Folium is a robust Python module that provides assistance in the development of a variety of different Leaflet maps. By default, Folium builds a map in a separate HTML file. Because Folium's findings may be interacted with, this library is an excellent resource for the construction of dashboards. In Folium, you also have the ability to construct inline Jupyter maps. Folium takes advantage of the strengths of the Python environment for working with data and the strengths of the Leaflet.js library for making maps. You can alter your data with Python using Folium, and then you can view the results using a Leaflet map [28] [29].
- **Pyspark**: PySpark is a Python-based interface for the Apache Spark data processing framework. It not only lets you build Spark applications with Python APIs, but it also gives you access to the PySpark shell, which lets you analyze your data in a distributed setting in an interactive way. In other words, Apache Spark is an analytical processing engine designed for use in

applications that need strong distributed data processing on a large scale. Machine learning applications also make use of Apache Spark [30] [31].

### 2.4 Related works

Studying the patterns of traffic and the conditions of traffic on roads and streets has become an extremely important part of the administration and planning of urban areas due to the ever-increasing amount of automobile traffic in urban areas. This is of utmost significance in large cities, which often have frustratingly slow moving traffic. Magnetic loops, video observation, laser and infrared vehicle detection are some of the more classic methods that have been employed in recent decades to gather information about traffic. However, in more recent times, a new set of other methods and data sources have also been utilized. In recent years, another data source utilized to investigate traffic situations is FCD. In contrast to traditional traffic data, which is typically collected at a fixed location by a stationary device or observer, Floating Car Data in traffic engineering and management is typically timestamped geo-localization and speed data directly collected by moving vehicles that are equipped with GNSS receivers using their on board unit in real-time.

Numerous publications have been published that investigate the efficiency of FCD in a variety of contexts, including but not limited to: costs, penetration rate into automobiles, dependability, and comparison of the acquired findings to more conventional methodologies. As was explained in [42], some traffic analysis might be performed using FCD, which would provide findings that are comparable to those of traffic sensors but would incur a far lower cost. In this study, the data that the FCD gathered over the course of 17 continuous hours and included 375,000 records that belonged to 10,000 cars were mapped to the nearest street using Open Transport Map (OTM) so that analysis could be carried out. [43] also described an experience of producing a comprehensive traffic analysis by utilizing coherent numbers of FCD and highlighting how this solution provides higher coverage than traffic sensors while also having lower costs. Using FCD and Remote Traffic Microwave Sensor Data (RTMS) data, [47] conducted an investigation of the peculiarities of traffic flow on Beijing's ring road expressways. While the average speeds were calculated using both sets of data, the volume calculations were carried out using RTMS data only. The flow-speed relationship in the basic figure was derived by combining the data that were previously collected. After doing a regression analysis on the average speeds obtained from RTMS and FCD, it was discovered that the values obtained from RTMS were, on average, 6% greater than those obtained from FCD. Another comparison was made by [48] between the average speeds of road segments acquired from FCD data and Bluetooth data. A statistical analysis was carried out for each

of the four speed categories, and it was discovered that the average Bluetooth speeds were not substantially different from the FCD speeds for any of the speed categories.

One of the disadvantages of using FCD data is the low penetration rate of on board units into cars, which is mentioned in [42] which could lead to the results having less reliability than they otherwise would due to the fact that the data were collected from a small percentage of cars. In spite of the fact that [44], on the basis of their simulation, suggests that the 1.5% penetration rate of vehicles with on board units is required to have an accurate estimation of travel time, [45] suggests that the penetration rate between 3%-5% in order to achieve the confidence level of over 90% would be sufficient. in addition, [46] came to the conclusion that the appropriate penetration rate should be anywhere from 5 to 10%, depending on the current traffic condition. Nonetheless, the scope of the investigation for [46] is limited to just one highway in Singapore. [46] also investigates the variations in the intervals of FCD collection, which revealed that the best results would acquire if the collection interval is between 15 and 20 seconds, despite the fact that this demands a great deal more processing power. Basically, extending the interval for every 15 seconds results in a considerable reduction in the needed amount of processing power.

[49] by studying just one parameter, which is average travel speed from FCD, the researchers came to the conclusion that it is still feasible to uncover key patterns along urban roadways, provided that the data is continuous and broad. This research, analyzes the FCD data obtained from automobiles in Ankara over a period of two months with an interval of one minute. The data was taken from the vehicles every minute. In this investigation, the average speed of the road segments is converted to a Level of Service (LOS), and the status of the traffic is determined by a comparison to the speed at which free-flow may occur on the same road segment. Moreover, by comparing the traffic states of two, three, and four successive road segments, one may determine where the bottlenecks are in the road system.

Alongside the implementation of FCD in [50] and [51], the idea of Big Data is also included. [50] blends big data with small data in the creation of elements of transport system models. This is done with the intention of increasing the ability of transport analysts and planners to assess, predict, and prepare for mobility phenomena. [50] concludes that the conventional techniques for constructing transport system models are improved with the use of big data for mobility. Similarly,[51] presents an original demand estimate framework that is comprised of two steps and argues that the use of Big Data and FCD results in the limitation of not having a strong initial demand matrix being overcome.

### 2.5 Dataset

The exploited dataset in the thesis is collected from an insurance company. It is a very large dataset, containing more than 700 million rows, that was collected from January of 2019 to January of 2020. The collection was made before the crisis of the COVID-19 pandemic, therefore the dataset represents the behavior of the vehicles in an ordinary time period.

The records inside the dataset are spread over an area, as vast as Piedmont province. However, since the goal of the thesis is to retrieve congestion pattern in a metropolitan area, thus only the records belonging to the city of Turin and some smaller nearby towns are conserved, and the other records are eliminated before proceeding to further steps.

The collected FCD dataset contains following columns and features:

Attribute	Type	Description
dowieoId	Tort	the unique ID associated to each device
devicera	Iext	1 to 20 characters
datoTimo	DateTime String	Time and date of submitted record
		format: YYYY-MMDD HH24:MI:SS
latitudo	Number $(9,6)$	expressed in degrees
		from -90 to 90
longitude	Number $(9,6)$	expressed in degrees
Iongitude		from -90 to 90
speedKmh	Integer	Instantaneous speed in Kilometer per Hour
		ranging from 0 to 250
heading	Integer	The direction that the vehicle is headed
		0 to 360 in 4 degree increments
	8 Integer	GPS signal quality
accuracyDop 2.8		expressed in tenths of HDOP $2.1$
		between 0 (excellent) and 150 (very bad)
EngineStatus	Integer	1 engine on
		0 engine off
Type	Integer	1 car
турс	III00gel	2 fleet

 Table 2.2:
 FCD dataset description

## Chapter 3

## **Implementation and Results**

### 3.1 Input FCD analysis

### 3.1.1 Data cleaning

The process of detecting and correcting (or removing) corrupt or inaccurate records from a record set, table, or database is known as data cleansing or data cleaning. It also means finding parts of the data that are missing, wrong, inaccurate, or not important, and then replacing, changing, or getting rid of the dirty or coarse data [**32**]. When numerous data sources are combined, there is a greater likelihood that some of the data may be duplicated or incorrectly categorized. Even though the results and methods seem to be right on the surface, their reliability is compromised if the underlying data is inaccurate. Because the procedures differ from dataset to dataset, there is no one method to prescribe the precise stages in the process of data cleaning. This is because there is no one definitive technique to do so. However, it is essential to develop a pattern for each data cleaning procedure in order to ensure that you are always doing the task in the appropriate manner [**33**].

As the main data cleaning steps taken on the input FCD data, below criterias, have been considered:

- **Duplication removal**: During the process of collecting the data, there will most likely be instances of duplicate observations. There is a possibility of creating duplicate data whenever data sets are joined from different locations, whenever data is scraped, whenever data is received from customers or numerous departments, and when scraping data. In this procedure, de-duplication is one of the most important aspects that needs to be addressed.
- Structural errors and missing data removal: When doing measurements or uploading data, it is possible to uncover irregular naming conventions,

typos, missing values, or incorrect capitalization. These inconsistencies may result in incorrect categorization of groups or classes.

The work flow of the thesis calls for a couple of data cleaning processes to be conducted at various phases of the development of the project. Each of these steps has its own unique set of requirements that must be met. The first stage in cleaning the data is based on the criteria that were specified before. Nevertheless, this would not be the only component of the thesis that involves cleaning up the data, and the other one is discussed in 3.2.

The original FCD input data spans 12 months of 2019 and consists of 738,969,717 records, of which 5.1% are duplicated rows; following the duplication elimination process, there are a total of 701,536,205 records. In addition, after the 'Not a Number' (NaN) values have been removed, the output comprises 655,798,896 entries. These steps are the very first initial steps of data cleaning.

### 3.1.2 Area of interest

Choosing the scope of the problem that will be investigated is one of the most important decisions to be made at the first phase of the process. Since the purpose of the thesis is to describe the traffic congestion pattern in an urban area, and the large FCD dataset contains records that are as large as the Piedmont province in Italy, and since the goal of the thesis is to describe the traffic congestion pattern, only the metropolitan city of this region and some small but nearby cities are investigated, and other zones are excluded from the investigation. Figure 3.1.3 shows the interest area for further analysis with the following coordinates that start from minimum latitude and longitude (44.96282106687191,7.502048016422193) to the maximum latitude and longitude (45.19265016665321,7.791812422724604). Eventually, the number of records that are inside this area is: 375,647,074.



Figure 3.1: The area of interest, including Turin and nearby cities

### 3.1.3 Data behaviour

### Vehicle types

From the total number of 375,647,074 records in the urban area, 310,264,034 records, meaning 92.99% of total, are belong to 'Personal' vehicle type and the rest with 23,367,845 records consist the 'Fleet' vehicle type. Since the 'Personal' type is dominant in the dataset, so the behavior of the vehicles are considered as a whole and not specifically by each type of vehicle.

### GPS accuracy DOP

As mentioned in 2.2.3, the accuracy of a GPS signal is affected by an error propagation factor, which is known as dilution of precision. Thus, the quantity of DOP can specify the signal quality of GPS. In figures 3.2 and 3.3 ,both cumulative and non cumulative distribution of DOP values of FCD dataset are demonstrated. according to [**37**], since a DOP value less or equal than 2 is considered as an "excellent" signal in terms of precision, further, based on the behavior on the cumulative distribution function; the value of  $HDOP \leq 2$  which contains 94.09%

of the whole dataset, is considered as a threshold to filter the records, and the records with a HDOP greater than 2 are dropped. This is also a part of the data cleaning part.



Figure 3.2: Cumulative Distribution Function of GPS accuracy DOP Histogram of Accuracy DOP



Figure 3.3: Distribution of GPS accuracy DOP

### **Records** daily distribution

Figure 3.4 shows how the FCD records are spread over the year 2019 in a daily time frame. However, for the sake of more readability, the plot is labeled in monthly intervals. As it can be seen, there is a repeating pattern in the behavior of the records. Basically, the pattern shows the weekly behavior. There is a peak on Fridays of each week, and a dip is found on Sundays, which are the weekend. However, in summer, specifically in the month of August, due to the summer holidays, a dramatic decrease in terms of submitted records is happening.



Figure 3.4: Daily Records Distribution

#### **Records hourly distribution**

Beside the daily and weekly time frames, it is also interesting to investigate the hourly distribution of the records as well. It could give information about peak hours that are most probably the time intervals in which the most congestion occurs. Figure 3.5 illustrates the diffusion of the FCD records in the hourly time frame. As it can be clearly seen, there are two peaks that exist. One in the morning and the other one, which has a higher peak, in the evening. Furthermore, there has been a drop in the early morning hours, when most people are sleeping and outside activity is at its lowest. Such behavior is quite reasonable. In the morning, when people usually go to school or work, a peak happens, and the other one, when people are returning to their houses after having finished the work activity.



Figure 3.5: Hourly Records Distribution

### **3.2** Segment Extraction

After conducting basic data cleaning and experimentation on the input FCD data, the first significant step is to extract the path segments of the taken route for each individual vehicle. This is done after examining the FCD data. This is the fundamental prerequisite for the subsequent procedures that will be undertaken to perform analysis.

### 3.2.1 Segment definition

The movement of an item may be segmented down into trajectories that span the distance between any two significant points. This division differs depending on the application. The movement of a vehicle used by a delivery service, for instance, may be segmented down not just into its daily travels, but also into its moves between individual clients [38].
In this thesis, the considered segmentation is defined as the smallest portion of the overall route that is traveled by each individual vehicle. These portions are identified as the traveled path between the time intervals during which the data has been collected.

#### 3.2.2 Segment extraction methodology

The first thing that has to be done in order to extract the path segment is to sort the records that belong to each unique vehicle (which has its own ID) according to the time stamp. Next, the path that was taken between every consecutive pair of recordings is what is known as the path segment for each and every unique vehicle. Then for each path segment, the following parameters are calculated:

#### Segment distance

Geopy is an extremely useful Python module that is used to determine the distance between segments. This package takes as its input the coordinates of two points, each of which has a latitude and a longitude, and its output is the Haversine distance 2.2.2, which it calculates using meters as the unit of measurement. Combining the segment distance with the calculated segment duration, enables the calculation of the segment speed at a later stage, which is one of the most important characteristics that must be determined.

#### Segment duration

Is the amount of time that elapses between the beginning of one segment and the beginning of the next. The calculation of the segment length provides valuable information on the time intervals at which each FCD report is sent.

As Figure 3.6 shows the cumulative distribution of duration of the segments, and also Figure 3.7 confirms as well, a majority of the path segments which means 56.93% of the whole, have been recorded with the intervals of 60 seconds. Considering a tolerance of 10 seconds, as a part of data cleaning, the segments with duration more than 70 seconds which form 8.5% of the dataset are identified as outlier and are dropped.



Figure 3.6: Cumulative distribution of Segment duration's



Figure 3.7: Distribution of Segment duration's

#### Segment speed

It is quite easy to calculate the segment's average speed; all that is required is a simple division of the segment's distance, which is measured in meters, by the segment's duration, which is measured in seconds. After that, a factor of 3.6 is applied to the segment average speed that was determined before in meters per second in order to convert it to the desired unit of kilometers per hour. To investigate the behavior of the path segments better and following by the applied filter in 3.6, the segment duration of records are grouped based on the different intervals as following:

- Segments with duration between 0-9 seconds
- Segments with duration between 10-19 seconds
- Segments with duration between 20-29 seconds
- Segments with duration between 30-39 seconds
- Segments with duration between 40-49 seconds
- Segments with duration between 50-59 seconds

Following this, the segment speed is computed for each of these groups in order to investigate the behavior of segment speed in each interval. The findings of this investigation are shown in the figure 3.8. It is clear that despite the fact that there are certain insignificant differences that are reasonable, no group is responsible for any kind of bias in the overall behavior, and groups could be considered similar to one another. To determine whether or not more data cleaning is required, this phase is carried out. However, because of the similarities, no special data cleaning action was performed at this phase. Instead, for the next steps, the segment duration of all of the records will be merged into a single group.



Figure 3.8: Cumulative distribution of Segment Speed (KM/H) of different groups

The cumulative distribution of the average speed of the segments is shown in figure 3.9, which is then followed by the distribution plot in figure 3.10. Because the highest speed that is permitted by law in Italy is 130 kilometers per hour, this limit has been selected as the upper bond in order to retain only the segments that have a value that is lower than or equal to this value. Because of this, 0.52% of the segments are discarded at this step of the data cleansing process. In addition, 2.72% of the segments have been found to have values that are equivalent to 0 on the absolute scale. Since a value of absolute zero indicates that there has been no movement at all for those particular segments, it is also necessary to eliminate them as part of the process of cleaning the data.

In a nutshell, there were two stages of data cleaning that took place in relation to the average speed of the segments: first, the segments whose average speed was equal to or less than 130 kilometers per hour were maintained as the upper bond; second, the segments whose average speed was greater than absolute zero were maintained as the lower bond.



**Figure 3.9:** Cumulative distribution of Segment Speed (KM/H)



Figure 3.10: Distribution of Segment Speed (KM/H)

### 3.2.3 Segment table

For the sake of ease and consistency throughout the work, the segment table is produced using the original FCD data as the input. This is based on the segment definition that has been supplied, as well as the technique that describes how the segments are extracted in 3.2.1 and 3.2.2. Due to the fact that it provides the computed average speed as well as the starting and ending location of each segment, the segment table is an extremely important component in the analysis of the thesis. This is because it leads the way for the discovery of traffic congestion. The most important columns of the segment table are shown in table 3.1.

Attribute	Type	Description
deviceId	Text e.g: 2507297	the unique ID associated to each device 1 to 20 characters
Туре	Integer	Vehicle type, 1: car , 2: fleet
dateTime	DateTime String e.g: 2019-04-02 10:15:28	Time and date of segment starting point format: YYYY-MMDD HH24:MI:SS
startLatitude	Number (9,6) e.g: 45.08672	latitude of segment starting point expressed in degrees from -90 to 90
startLongitude	Number (9,6) e.g: 7.60882	longitude of segment starting point expressed in degrees from -90 to 90
startEngineStatus	Integer	Engine status at segment starting point 1: engine on , 0: engine off
startAccuracyDop	Integer e.g: 10	GPS signal quality expressed in tenths of HDOP between 0 and 20
endAccuracyDop	Integer e.g: 10	GPS signal quality expressed in tenths of HDOP between 0 and 20
endEngineStatus	Integer	Engine status at segment ending point 1: engine on , 0: engine off
startLatitude	Number (9,6) e.g: 45.08683	latitude of segment ending point expressed in degrees from -90 to 90
startLongitude	Number (9,6) e.g: 7.60893	longitude of segment ending point expressed in degrees from -90 to 90
segmentDistance	Float, e.g: 139.25	travelled distance in the segment in unit of meters
segmentDuration	Float, e.g: 60.0	duration of the the segment in unit of seconds
segmentSpeedKmH	Float, e.g: 26.38	average speed of the segment in unit kilometers per hour

**Table 3.1:** Segment table description, Other properties derived directly from the latitude and longitude coordinates have been omitted for the purpose of brevity.

#### 3.2.4 Segment distribution

Another aspect to investigate regarding the behavior of the segments, is to examine the distribution of the segments over different time spans, for instance, monthly, weekly, daily, and hourly.

#### Segments monthly, weekly and daily distribution

Figure 3.11 illustrates how the segments are distributed during the course of 2019 on a daily basis. The figure, however, is labeled in monthly intervals for the convenience of reading. Observably, the behavior of the records follows a recurring pattern. The pattern essentially depicts the weekly activity. Each week has a weekly high on Friday and a weekly trough on Sunday, the weekend. However, throughout the summer, notably the month of August, owing to the summer vacations, the number of submitted records decreases dramatically.



Figure 3.11: Daily Segments Distribution

In addition to the daily and weekly time periods, it is also worthwhile to analyze the hourly distribution of the segments. It provides information on peak hours, which are most likely the most congested periods of time. Figure 3.12 depicts the distribution of segments across an hourly time period. Clearly visible are the existence of two peaks. The one with the greater peak occurs in the evening, while the one with the lower peak occurs in the morning. In the early morning hours, when most people are asleep and outdoor activity is at its lowest, there has also been a decline. Such conduct is very fair. A peak occurs in the morning, when people typically go to school or work, and another occurs in the evening, when people are coming home from work.



Figure 3.12: Hourly Segments Distribution

# 3.3 Grid map

Based on the goal of the thesis, the considered area 3.1 is broken down to a grid on the map. A network of horizontal and vertical lines that are equally spaced apart and used to designate places on a map is known as a grid. The row and column labels of a reference grid will often indicate locations that are mentioned in an index for the map [**39**].

#### 3.3.1 Grid size selection

When determining the appropriate cell size inside the grid, some criteria that should be considered include the resolution, the amount of time needed for computation, and the number of segments contained within each cell. There needs to be some sort of compromise between the various aspects that have been brought up, and they are discussed one by one:

#### Number of segments per cell

The target cell should have a significant number of subcellular segments packed inside of it. Because the influence of each segment's inherent bias would be reduced as the density increased, the cell would become both more informative and more trustworthy as the density increased. Figure 3.13 shows the average number of segments per cell for different cell sizes



Figure 3.13: Segment density VS cell size

#### Computation time

To obtain a greater resolution, the size of each cell must be reduced, which results in an increase in the number of cells that must be processed. As a direct consequence of this increase, the amount of calculation time that is necessary also rises. However, in the case of this thesis, the procedure may be completed in a reasonably short period of time because to the extremely powerful computing resources that are given by SmartData@PoliTO [40]. This is true even in difficult scenarios that need a significant level of compute power. The amount of time needed for calculation is broken down below in figure 3.14 according to the various cell sizes.



Figure 3.14: Required processing time VS cell size

#### Resolution

Resolution of cells is one of the most important considerations to make while selecting an appropriate cell size. If the resolution was higher, then there would be more specific information on every little piece of the area. On the other hand, there would be fewer segments for that area, which would result in the information being less accurate. Alternatively, if the resolution is low, the reliability will rise, but there will be less details as a consequence.

As a whole, and on the basis of the findings that were obtained, the size of the grid's individual cells has been decided to be 100 meters by 100 meters.

# 3.4 Grid table

Following the completion of two key phases, which are represented in 3.1 and 3.3, the segment table and the outcome of dividing the map into a grid of cells are combined into a single entity with the common key which is the identifier assigned to each cell. The consequence of this merge is having the possibility of visualizing on the map how the segments are distributed across the metropolitan area of Turin and some of the smaller cities that are located in the surrounding area in a variety of different aspects. These aspects for each cell of the grid including, the density of the segments, the speed of the segments, and the GPS accuracy DOP of the segments. For each component, one might just compute the minimum, maximum, and the mean; nevertheless, the end result of calculating the mean value is the

Attribute	Type	Description
polygon_id	String	20_63 means cell located in 20th
	e.g: $20_{63}$ means	row and 63th column on the grid
geometry	Polygon	The geometry of the
		corresponding cell
cell_id	String	$20_{63}$ means cell located in $20$ th
	e.g: $20_{63}$ means	row and 63th column on the grid
num_segments	Interger	Indicates the density
	e.g: 1594	of the segments inside the cell
minSpeed	Float	Minimum speed of the
	e.g: $1.54 \; ({\rm KM/H})$	segments inside the cell
avgSpeed	Float	Mean of speed of the
	e.g: 20.39 (KM/H)	segments inside the cell
avgSpeed	Float	Maximum speed of the
	e.g: 120.14 (KM/H)	segments inside the cell
minAccuracyDop	Integer e.g: 1	Best GPS signal quality
		of the segments inside the cell
		expressed in tenths of HDOP
avgAccuracyDop	Float e.g: 9.32	Mean of GPS signal quality
		of the segments inside the cell
		expressed in tenths of HDOP
maxAccuracyDop	Integer e.g: 20	Worst GPS signal quality
		of the segments inside the cell
		expressed in tenths of HDOP

point of interest because it relates to the purpose of the thesis.

Table 3.2: Grid table. Obtained by outer join of 3.1 and 3.3

## 3.4.1 Grid plots

This current section has, up until this point, provided an explanation of the mechanism behind the extraction of grid tables. As a direct consequence of this, the plots of the grid map in three different aspects—the density of segments in each cell, the mean of the speeds of the segments in each cell, and the GPS accuracy DOP—are presented respectively in 3.15, 3.16, and 3.17.

#### Segment density per cell



Figure 3.15: Number of segments per cell

It is clear by looking at 3.15 that the density of segments within is substantially higher along the main streets and highways than it is in any other section of the city. This demonstrates the predicted behavior of automobiles and demonstrates that main streets and highways are being utilized to a greater extent by automobiles.



#### Mean of the speeds of the segments per cell

Figure 3.16: Mean of the speeds of the segments per cell

it can be seen in 3.16, highways contain the cells within which, the average speed of the segments is at the highest. In the urban area, along the main streets, the average speed is relatively higher, with respect to the downtown and side streets. As expected, the calculated average speed of the segments is quite low in downtown, where it is expected to be the most crowded area.



#### Mean of GPS accuracy DOP per cell

Figure 3.17: Mean of GPS accuracy DOP per cell

As Figure 3.17 demonstrates, the accuracy of received GPS signal in urban areas and also in the hills located on west side of Turin, are relatively lower than the other pats. This could be due to the existence of the buildings which cause multipath errors. The receipt of signals that have come not only straight from satellites but also reflected or diffracted from the objects in the immediate vicinity is what causes GPS multipath that eventually resulting a shift in the location that was estimated [41].

# 3.5 Traffic congestion index

In this part of the article, the methodology behind calculating the traffic congestion index is broken down and explained. This is made possible as a direct result of the creation of the grid table 3.4, which makes this particular option available. Two independent measures of traffic congestion are used in the calculation of the traffic congestion index: absolute traffic congestion and relative traffic congestion. In spite of this, the focus will mostly be on the latter since it is the objective of the thesis. The reason for this is because the region being researched has several various kinds of roadways. As an example, on a roadway that has an average speed of 100 KM/H during off-peak hours and 80 KM/H during peak hours, the difference of 20 KM/H is not considered to be severe congestion. On the other hand, the difference of 20 KM/H is significant in a neighborhood in the downtown area where the average speed is 40 KM/H.

One key step to be taken before the traffic congestion index calculation, is to identify the peak and off-peak of hours over a complete day cycle of 24 hours. Based on the hourly distribution of the segments 3.12, two peaks the chart are considered as peak hours, morning and evening peak. morning peak hours start from 9AM to 11AM, while the evening peak period is from 6PM to 8PM. The off-peak interval is chosen for the interval of 1AM to 6AM. The reason for choosing the off-peak period longer than the peak hours is due to the low number of records submitted during this period that could lead the analysis to be less reliable. The selected intervals are depicted in figure 3.18.



Figure 3.18: Peak and off-peak hours selection

#### Dropping uncommon cells

One of the first steps toward calculation and comparison of the traffic congestion is to keep only the cells that contain segments in both peaks and off-peak intervals. This is obvious, since in the case of a lack of data in either interval, it is not possible to make a comparison and calculate the congestion.

#### Traffic index calculation

In this thesis, two different ways of traffic congestion index calculations consist of absolute traffic congestion and relative traffic congestion are experimented that are respectively explained in 3.5.1 and 3.5.2, however, as mentioned earlier the focus is mainly on the relative congestion index

## 3.5.1 Absolute traffic congestion

Absolute traffic congestion index is simply calculated as the difference of the speed between peak and off-peak hours for the same cell in the grid:

$$V_{OP} - V_P \tag{3.1}$$

Where:

 $V_{OP}$ , is the Speed (KM/H) in the Off-peak period  $V_P$ , is the Speed (KM/H) in the Peak period

#### Distribution

The cumulative distribution of cell-wise traffic congestion is shown in Figure 3.19, which compares the morning and evening peaks with the off-peak period individually. As can be seen, 34.74% of the cells in the morning peak and 29.78% of the cells in the evening peak are exhibiting negative congestion. This indicates that the calculated average speed of cells during the off-peak period is lower than the one during the peak period, which is behavior that is somewhat undesirable. It is to be anticipated that it will be more congested during the peak hours than during the off-peak hours in the regular scenario. In addition, the distribution of cell-wise absolute traffic congestion is displayed in Figure 3.20 below. It is abundantly evident that there is a behavior that is more favorable for the evening peak. The primary reason is because there are more data obtained for the evening peak in comparison to the data collected for the morning peak, which simply means that the results are more reliable. Therefore, the focus will be mostly the comparison between evening peak and the off-peak hours in the very early morning

 $3.5 - Traffic \ congestion \ index$ 



Figure 3.19: Cell-wise cumulative distribution of absolute traffic congestion



Figure 3.20: Cell-wise distribution of absolute traffic congestion

#### 3.5.2 Relative traffic congestion

The relative congestion index uses a simple formula 3.3 to measure the relative variations in speed during peak hours against off-peak hours:

$$\left(\frac{V_{OP} - V_P}{V_{OP}}\right) * 100\tag{3.2}$$

Where:

 $V_{OP}$ , is the Speed (KM/H) in the Off-peak period  $V_P$ , is the Speed (KM/H) in the Peak period

To standardize the scale and keep it between -100% and 100%, the greater velocity (either in the peak or off-peak interval) is always used as the base to compute the relative changes in speed. As a result, in circumstances when negative congestion is observed, the formula will be slightly different:

$$-\left(\frac{V_P - V_{OP}}{V_P}\right) * 100\tag{3.3}$$

#### Distribution

Similarly to 3.5.1, the distribution of the traffic congestion index indicates similar behavior. Negative congestion is present in 34.74% of the cells in the morning peak and 29.78% percent of the cells in the evening peak. The more reliable behavior of the evening peak interval is obvious in the relative traffic congestion index as well. Figures 3.21 and 3.22 show the distribution of relative traffic congestion in cumulative and normal ways respectively.

3.5 – Traffic congestion index



Figure 3.21: Cell-wise cumulative distribution of relative traffic congestion



Figure 3.22: Cell-wise distribution of relative traffic congestion

## 3.5.3 Cells with negative congestion

As indicated in 3.5.1 and 3.5.2, there exist some cells in the grid that demonstrate negative congestion, implying that the average speed in the off-peak interval is

lower than in the peak hours. To understand in which cells this scenario is more likely to happen, the behavior of the number of segments per cell is investigated as the main feature of each cell based on the purpose of the thesis. As shown in 3.23, cells with fewer segments inside are much more likely to exhibit negatively congested behavior. This is mainly due to the unreliability of the results with a low amount of input. In addition, as it can be seen also in figure 3.24 the distribution of the negatively congested cells, tend to fade out (green bar) as the number of the segments per cell increase and the blue and orange bars will be aligned eventually. Therefore, it would be reasonable to filter out the cells with fewer segments inside for further steps in the analysis.



Figure 3.23: Cumulative distribution of cell-wise traffic congestion

3.5 – Traffic congestion index



Figure 3.24: Distribution of cell-wise traffic congestion

#### Filtering cells

Based on the cells behavior that is discussed in 3.5.3, it is already observed that the cells with a lower number of segments inside are more prone to demonstrating unreliable results like negative congestion. As a result, the cells with less than 100 segments and also segments with less than 200 segments are filtered separately, and then they are plotted against the not filtered cumulative distribution and the results are shown in the figure 3.25



Figure 3.25: comparison of cell-wise cumulative distribution of relative traffic congestion in filtered and not filtered scenarios

As can be observed, the overall behavior of the relative congestion index is significantly better after doing the filtration because there are less cells with negative congestion. Furthermore, in the region with positive congestion, the non-filtered line (blue) is flatter and hence less informative. When cells with less than 100 segments per cell are filtered, the improvement in behavior is more noticeable. When the filtering criteria is increased to 200, the improvement is not significant, thus the threshold of 100 segments per cell is deemed to filter.

The next filtering criteria is to remove the cells that are showing negative congestion. Because achieving the goal of the thesis requires concentrating on locations that are experiencing the highest levels of traffic congestion, the point of interest is not a single cell but rather the behavior of a group of cells that are adjacent to one another. Figures 3.26 and 3.27 show the distribution of cell-wise relative speed difference after applying the filter

3.5 – Traffic congestion index



Figure 3.26: Filtered cell-wise cumulative distribution of relative traffic congestion



Figure 3.27: Filtered cell-wise distribution of relative traffic congestion

As it can be seen, the relative congestion index is peaked at the range between 15% and 20% that means the severity of the congestion is mostly occurs around this range.

# **3.6** Identification of congested areas

The criteria for filtering the cells in the grid are broken down and explained in 3.5.3. After the filters have been applied, everything is ready to show which cells and eventually which areas are the most congested based on the relative congestion index 3.5.2. It should not come as a shock that the downtown area is the most likely



Figure 3.28: Relative traffic congestion

zone to be highly congested given the evidence presented here; it is the downtown area, after all. However, this is not the only area with a significant amount of traffic congestion; it is also possible to see that the traffic congestion occurs on the highways, despite the fact that the average speed may not be particularly slow; however, the speed difference in terms of percentage demonstrates that the highways may also have a remarkable amount of traffic congestion. This was able to be seen because of the relative congestion calculations, which makes it possible to observe this phenomenon. On the other hand, such a thing was not able to be observed when calculating the absolute traffic congestion index.

#### 3.6.1 Kernel Density Estimation: a method for smoothing

Based on the fact that the structure of the project is a grid of cells, the overall image is still a grid of individual cells, each of which depicts the congestion in corresponding coordinates individually. This is the case even after the congestion index of the cells has been identified. However, According to the objective of the thesis, which is to identify the zones, and because each zone is most likely composed of a couple of cells, a method is required to demonstrate the behavior of a group of cells in an area rather than the behavior of individual cells. To smooth out the density, one of the approaches that might be used is to implement Kernel Density Estimation (KDE), which is a method that estimates density. This would be one of the available methods. In this instance, the relative congestion index, which is stated as a percentage, is taken into consideration as density, and several KDE models with varying amounts of bandwidth are trained. The amount of available bandwidth is the primary factor that determines how smooth the outcomes will be. The findings are shown in the charts that follow.



Figure 3.29: KDE of Relative traffic congestion, Bandwidth = 0.0005



Figure 3.30: KDE of Relative traffic congestion, Bandwidth = 0.001



Figure 3.31: KDE of Relative traffic congestion, Bandwidth = 0.0015



Figure 3.32: KDE of Relative traffic congestion, Bandwidth = 0.002



Figure 3.33: KDE of Relative traffic congestion, Bandwidth = 0.0025



Figure 3.34: KDE of Relative traffic congestion, Bandwidth = 0.003



Figure 3.35: KDE of Relative traffic congestion, Bandwidth = 0.0035

As can be seen in 3.29 3.30 3.31 3.32 3.33 3.34 3.35, as bandwidth grows, surrounding cells begin to merge more and more, and this continues until the whole of the city is combined into one single cell. A bandwidth of 0.001 is deemed to be an appropriate bandwidth for the purpose of highlighting the locations that are suffering the most from traffic congestion.

#### Plot of the congested areas

as the final step of the identification of congested areas, in order to highlight the congested areas, a filter is applied to the bandwidth of 0.001 to demonstrate the congested areas clearly. The result is shown in 3.36.



Figure 3.36: Highlight of the most congested areas, Bandwidth = 0.001

# Chapter 4 Conclusions

The subject that was suggested for this study was motivated mostly by the two key questions that are outlined below.

- How is it possible to extract traffic congestion pattern from FCD using a data-driven approach?
- How much is it scalable to expand the area under investigation? and what could be the best resolution to achieve?

We realized that it is possible to extract the traffic congestion pattern by making a comparison of the average speed of vehicles in the same area but at different time intervals throughout the day. We did this by integrating a large quantity of FCD that had already been collected into the analysis of this thesis. We came to this conclusion after realizing that it is possible to do so. In addition, we were aware that, despite the fact that Politecnico di Torino has a very robust hardware infrastructure, the experiment's reach could not be expanded beyond the confines of a metropolitan urban region. In addition, because to the limitations imposed by the computing capability of the hardware, the resolution may be reduced the greater the region that is the subject of the study.

A further consideration is the application of the grid over the top of the map. As we have seen, it is necessary to take into account the behavior of a group of cells that are adjacent to one another in order to draw a conclusion about the degree to which there is an issue with traffic congestion in the region. This is the case even though some of the individual cells of the grid have displayed undesirable behaviors. This could be done by doing a smoothing step on all of the grid cells as the last step. This would give a better picture of the areas with a lot of traffic. The existing findings might be used in a later piece of research on this subject in order to apply certain machine learning techniques in order to create predictions regarding the congestion index in the future.
## Bibliography

- [1] Caves, R. W. (2004). Encyclopedia of the City. Routledge. p. 141.
- [2] Treiber, Martin; Kesting, Arne (2012-10-11). Traffic Flow Dynamics: Data, Models and Simulation. Springer Science Business Media. ISBN 978-3-642-32459-8.
- [3] May, Adolf Darlington (1990). Traffic Flow Fundamentals. Prentice Hall. ISBN 9780139260728.
- [4] https://www.brookings.edu/research/traffic-why-its-getting-worse-what-government-can-do/.
- [5] Agyapong F and Ojo T K 2018 Managing traffic congestion in the accra central market, Ghana Journal of Urban Management. 7 pp 85-96
- [6] https://translineinc.com/traffic-data-collection/#:~:text=Other% 20methods%20of%20collecting%20traffic,the%20direction%20of% 20traffic%20flow.
- [7] Leduc, Guillaume. (2008). Road Traffic Data: Collection Methods and Applications.
- [8] G.S. Larue, C. Wullems, A new method for evaluating driver behavior and interventions for passive railway level crossings with pneumatic tubes [J]. J. Transportation Saf. Secur. 11(2), 150–166 (2019)
- [9] S. Rajab, M.O. Al Kalaa, H. Refai, Classification and speed estimation of vehicles via tire detection using single-element piezoelectric sensor[J]. J. Adv. Transportation 50(7), 1366–1385 (2016)
- [10] M. Grote, I. Williams, J. Preston, et al., A practical model for predicting road traffic carbon dioxide emissions using Inductive loop detector data[J]. Transportation Res. Part D Transp. Environ. 63, 809–825 (2018)
- [11] https://www.techtarget.com/iotagenda/definition/ Internet-of-Things-IoT
- [12] Llorca, D.F., Sotelo, M.A., Sánchez, S. et al. Traffic Data Collection for Floating Car Data Enhancement in V2I Networks. EURASIP J. Adv. Signal Process. 2010, 719294 (2010)

- [13] https://www.septentrio.com/en/learn-more/about-GNSS/why-multi-frequencyand-multi-constellation-matters
- [14] Travel Time Data Collection Handbook, FHWA report, chapter 5, ITS Probe Vehicle Techniques, 1998.
- [15] Chang, Kang-tsung (2016). Introduction to Geographic Information Systems (9th ed.). McGraw-Hill. p. 24. ISBN 978-1-259-92964-9.
- [16] https://en.wikipedia.org/wiki/Geographic\_coordinate\_system
- [17] https://www.ga.gov.au/scientific-topics/positioning-navigation/ wgs84
- [18] https://www.gpsworld.com/data-collection-of-wgs-84-information-or-is-it/
- [19] Kuhlman, Dave. "APythonBook:BeginningPython, AdvancedPython, andPythonExercises"
- [20] Wes McKinney et al. Data structures for statistical computing in python. In Proceedings of the 9th Python in Science Conference, volume 445, pages 51–56. Austin, TX, 2010.
- [21] https://pandas.pydata.org/
- [22] Kelsey Jordahl, Joris Van den Bossche, Martin Fleischmann, Jacob Wasserman, James McBride, Jeffrey Gerard, Jeff Tratner, Matthew Perry, Adrian Garcia Badaracco, Carson Farmer, Geir Arne Hjelle, Alan D. Snow, Micah Cochran, Sean Gillies, Lucas Culbertson, Matt Bartos, Nick Eubank, maxalbert, Aleksey Bilogur, Sergio Rey, Christopher Ren, Dani Arribas-Bel, Leah Wasser, Levi John Wolf, Martin Journois, Joshua Wilson, Adam Greenhall, Chris Holdgraf, Filipe, and François Leblanc. geopandas/geopandas: v0.8.1, July 2020.
- [23] https://geopandas.org/
- [24] https://numpy.org/
- [25] J. D. Hunter. Matplotlib: A 2d graphics environment. Computing in Science & Engineering, 9(3):90–95, 2007.
- [26] https://matplotlib.org/
- [27] https://pypi.org/project/Shapely/
- [28] https://www.dominodatalab.com/data-science-dictionary/folium
- [29] https://python-visualization.github.io/folium/
- [30] https://spark.apache.org/
- [31] https://sparkbyexamples.com/
- [32] Wu, Shaomin (2013) A Review on Coarse Warranty Data and Analysis. Reliability Engineering and System Safety, 114. pp. 1-11. ISSN 0951-8320.
- [33] https://www.tableau.com/learn/articles/what-is-data-cleaning
- [34] Dudek, Gregory; Jenkin, Michael (2000). Computational Principles of Mobile Robotics
- [35] https://en.wikipedia.org/wiki/Dilution\_of\_precision\_(navigation)
- [36] https://web.archive.org/web/20141122153439/http://www.gmat.unsw. edu.au/snap/gps/gps\_survey/chap1/149.htm

- [37] Isik, O.K.; Hong, J.; Petrunin, I.; Tsourdos, A. Integrity Analysis for GPS-Based Navigation of UAVs in Urban Environment. Robotics 2020, 9, 66. https://doi.org/10.3390/robotics9030066
- [38] Biljecki, F., Ledoux, H., Van Oosterom, P. (2013): Transportation mode-based segmentation and classification of movement trajectories. International Journal of Geographical Information Science, 27(2), pp. 385-407
- [39] https://desktop.arcgis.com/en/arcmap/latest/map/page-layouts/ what-are-grids-and-graticules-.htm
- [40] https://smartdata.polito.it/computing-facilities/
- [41] T. Kos, I. Markezic and J. Pokrajcic, "Effects of multipath reception on GPS positioning performance," Proceedings ELMAR-2010, 2010, pp. 399-402.
- [42] Ajmar, A., Arco, E., Boccardo, P., Perez, F. (2019). Floating car data (fcd) for mobility applications. The International Archives of Photogrammetry, Remote Sensing and Spatial Information Sciences, 42, 1517-1523.
- [43] Schäfer, R.P., Thiessenhusen, K.U., Wagner, P., 2002. A Traffic Information System by Means of Real-Time Floating- Car Data. In: 9th World Congress on Intelligent Transport Systems.
- [44] BS Kerner, C Demir, RG Herrtwich, SL Klenov, H Rehborn, A Haug, et al. Traffic state detection with floating car data in road networks. In Intelligent Transportation Systems, Proceedings. IEEE, pages 44–49. IEEE, 2005.
- [45] Xiaowen Dai, Martin Ferman, Robert P Roesser, et al. A simulation evaluation of a real-time traffic information system using probe vehicles. In Intelligent Transportation Systems, 2003. Proceedings. 2003 IEEE, volume 1, pages 475–480. IEEE, 2003
- [46] Sunderrajan, Abhinav, et al. "Traffic state estimation using floating car data." Procedia Computer Science 80 (2016): 2008-2018.
- [47] Zhao, N., Yu, L., Zhao H., Guo, J., and Wen, H., 2009. Analysis of Traffic Flow Characteristics on Ring Road Expressways in Beijing: Using Floating Car Data and Remote Traffic Microwawe Sensor Data. Transportation Research Record: Journal of the Transportation Research Board, pp. 178-185.
- [48] Haghani, A., Hamedi, M., Sadabadi, K.F., Young, S., and Tarnoff, P., 2010. Data Collection of Freeway Travel Time Ground Truth with Bluetooth Sensors. Journal of Transportation research record, pp. 60-68.
- [49] Oruc Altintasi, Hediye Tuydes-Yaman, Kagan Tuncay, Detection of urban traffic patterns from Floating Car Data (FCD), Transportation Research Procedia, Volume 22, 2017, Pages 382-391, ISSN 2352-1465, https://doi.org/10.1016/j.trpro.2017.03.057.
- [50] Croce, A.I.; Musolino, G.; Rindone, C.; Vitetta, A. Transport System Models and Big Data: Zoning and Graph Building with Traditional Surveys, FCD and GIS. ISPRS Int. J. Geo-Inf. 2019, 8, 187. https://doi.org/10.3390/ijgi8040187
- [51] Cantelmo, Guido, and Francesco Viti. "A big data demand estimation model for

urban congested networks." Transport and Telecommunication (2020).

## Appendix A Source code

Reading and filtering initial data

```
from pyspark.sql import functions as F
1
2 import matplotlib.pyplot as plt
3 from pyspark.sql.types import *
4 import seaborn as sns
5 import pandas as pd
6 import numpy as np
7 import geopy.distance
   import matplotlib
8
9 import datetime
10 import os
   from haversine import haversine, Unit
11
12
13
   import fastplot
14 |%matplotlib inline
15
   #schema of initial data
16
   schema_initial_data = 'idRequest_string, deviceId_string, 
17
      dateTime\_string, \_
   latitude_{\Box}double, \_longitude_{\Box}double, \_speedKmh_{\Box}Integer, \_
18
      heading \Box Integer, \backslash
   accuracyDrop \sqcup Integer, \sqcup EngineStatus \sqcup Integer, \sqcup Type \sqcup Integer'
19
   #reading csv
20
   path = os.path.abspath(os.getcwd())
21
   df = spark.read.csv('file:///%s../../shared/data/
22
      FCD_complete/*/* '%path, sep=", ", schema =
      schema initial data)
```

```
23
24
   #dropping duplicates
25
   df = df.distinct()
26
   df = df.na.drop("any")
27
28
29
   #dropping points with accuracy higher than 20
   df = df. filter (df [ 'accuracyDrop '] <= 20)
30
31
32
   #dateTime format
   df = df.withColumn('dateTime', F.to_timestamp('dateTime'))
33
34
35
              ——— Torino and countryside
   # =
   minLat = 44.96282106687191
36
37
   \min Lon = 7.502048016422193
38
   \max Lat = 45.19265016665321
39
   \max Lon = 7.791812422724604
   \# =
40
   @F.udf(returnType=BooleanType())
41
   def drop_outside_points(lat, lon):
42
43
       condition = (lat \geq minLat) & (lat \leq maxLat) & (lon \geq
            minLon) & (lon \ll maxLon)
       return condition
44
45
   df = df. filter (drop_outside_points ('latitude', 'longitude')
46
      )
47
48
   # Adding calculated segment parameters to initial data
49
   schema_inital_segment = 'index_Integer, idRequest_string,_
      deviceId \_ string , \_ dateTime\_ timestamp , \_ \
   latitude_double,_longitude_double,_speedKmh_Integer,_
50
      heading \Box Integer, \backslash
51
   accuracyDrop_Integer, EngineStatus_Integer, Type_Integer,
      segmentDistance double, \
52
   segmentDuration_double, _segmentSpeedKmH_double'
53
54
   @F.pandas_udf(schema_inital_segment, functionType=F.
      PandasUDFType.GROUPED MAP)
55
   def adding segment parameter(df):
       os.environ["ARROW PRE 0 15 IPC FORMAT"] = "1"
56
```

```
57
       df = df.sort_values(by="dateTime")
58
59
       df.reset_index(inplace=True)
60
        distance = [float(0)]
61
        duration = [float(0)]
62
       speed = [float(0)]
63
64
       for index, row in df.iterrows():
65
            if index == len(df) - 1:
66
                pass
67
            else:
                 time_difference = df.loc[index+1, 'dateTime'] -
68
                    df.loc[index, 'dateTime']
                 time_difference = time_difference.total_seconds
69
                    ()
                 duration.append(time_difference)
70
71
72
                org = (df.loc[index, 'latitude'], df.loc[index, '
                    longitude '])
73
                 des = (df.loc[index+1, 'latitude'], df.loc[index
                    +1, 'longitude'])
74
                 distance.append(float("{0:.2f}".format(geopy.
75
                    distance.distance(org,des).m)))
76
77
                \mathbf{try}:
78
                     speed.append(geopy.distance.distance(org,
                        des).m/time_difference)
79
                except:
80
                     speed.append(float("NaN"))
81
82
       df ['segmentDistance'] = distance
       df['segmentDuration'] = duration
83
        df ['segmentSpeedKmH'] = [v*3.6 \text{ if } not \text{ pd.isna}(v) \text{ else}
84
           float("NaN") for v in speed ]
85
       return df
86
   df_added_parameters = df.groupby("deviceId").apply(
87
      adding segment parameter)
88 df added parameters = df added parameters.drop('index')
```

89 df\_added\_parameters.write.csv('file:///%s/all\_merged/'%path , sep=",")

Creating segment table

```
from pyspark.sql import functions as F
1
2
  import matplotlib.pyplot as plt
  from pyspark.sql.types import *
3
  import seaborn as sns
4
5
  import pandas as pd
  import numpy as np
6
7
  import geopy.distance
  import matplotlib
8
9
  import datetime
10
  import os
11
12
  import fastplot
  %matplotlib inline
13
14
15
  #schema of initial data
16
   schema_inital_segment = 'idRequest_string, _deviceId_string,
     \Box dateTime\Box timestamp, \Box \setminus
   latitude_double,_longitude_double,_speedKmh_Integer,_
17
     heading \Box Integer, \backslash
18
   accuracyDrop_Integer, EngineStatus_Integer, Type_Integer,
     segmentDistance_double, \setminus
   segmentDuration_double, _segmentSpeedKmH_double'
19
20
  \#reading \ csv
21
   path = os.path.abspath(os.getcwd())
   df = spark.read.csv('file:///%s/all_merged/*'%path,sep=",",
22
      schema = schema_inital_segment)
23
24
  #dateTime format
   df = df.withColumn('dateTime', F.to_timestamp('dateTime'))
25
26
27
   schema_segment_inital ='deviceId_String,
28
   uuuuuuuuuuuuuuutype.Integer,
29
   dateTime string,
30
   uuuuuuuuuuuuuuuuustartLatitude.double,\
  startLongitude_double,
31
  32
33 | .....startAccuracyDrop_Integer, \
```

```
34 | .....endAccuracyDrop_Integer, \
35 | .....endEngineStatus_Integer, \
36
  endLatitude_double,
37
  endLongitude_double,
38
   segmentDistance_double,
   segmentDuration_double,
39
   \label{eq:segmentSpeedKmH} usual segmentSpeedKmH_double '
40
41
42
  @F.pandas_udf(schema_segment_inital, functionType=F.
     PandasUDFType.GROUPED MAP)
   def segment table maker(df):
43
       os.environ["ARROW_PRE_0_15_IPC_FORMAT"] = "1"
44
       df = df. sort values (by="dateTime")
45
       df.reset index(inplace=True)
46
       df segment temp=pd.DataFrame()
47
48
                                          , , ,
       '' Creating empty lists for columns
49
       deviceId = []; dateTime = [];
50
       startLat = []; startLon = []; startEngineStatus = [];
51
         startAccuracyDrop = [];
       endLat = []; endLon = []; endEngineStatus = [];
52
         endAccuracyDrop = [];
       distance = []; duration = []; speed = [];
53
       Type = df['Type'][0]
54
       for index, row in df.iterrows():
55
           if index = 0:
56
57
              pass
58
           else:
              deviceId.append(row['deviceId'])
59
60
              dateTime.append(str(df.loc[index-1,'dateTime'])
                 )
              startLat.append(df.loc[index-1,'latitude'])
61
              startLon.append(df.loc[index-1,'longitude'])
62
63
              startEngineStatus.append(df.loc[index-1,'
                 EngineStatus '])
              startAccuracyDrop.append(df.loc[index-1,'
64
                 accuracyDrop '])
              endLat.append(row['latitude'])
65
66
              endLon.append(row['longitude'])
              endEngineStatus.append(row['EngineStatus'])
67
              endAccuracyDrop.append(row['accuracyDrop'])
68
```

```
69
                distance.append(row['segmentDistance'])
70
                duration.append(row['segmentDuration'])
                speed.append(row['segmentSpeedKmH'])
71
72
73
74
        df_segment_temp['deviceId'] = deviceId
        df segment temp['type'] = Type
75
        df_segment_temp['dateTime'] = dateTime
76
        df_segment_temp['startLatitude'] = startLat
77
        df_segment_temp['startLongitude'] = startLon
78
79
        df_segment_temp['startLongitude'] = startLon
        df_segment_temp['startEngineStatus'] =
80
           startEngineStatus
        df segment temp['startAccuracyDrop'] =
81
           startAccuracyDrop
        df_segment_temp['endLatitude'] = endLat
82
        df_segment_temp['endLongitude'] = endLon
83
        df_segment_temp['endEngineStatus'] = endEngineStatus
84
        df_segment_temp['endAccuracyDrop'] = endAccuracyDrop
85
        df_segment_temp['segmentDistance'] = distance
86
87
        df_segment_temp['segmentDuration'] = duration
        df_segment_temp['segmentSpeedKmH'] = speed
88
89
90
        return df segment temp
91
    df segment temp = df.groupby("deviceId").apply(
92
       segment table maker)
93
    "'applying filter: dropping segment durations longer than
94
       70s
95
    as well as segments with speed higher than 130KmH or equal
       to absolute zero ','
96
    df_segment_temp = df_segment_temp.filter(df_segment_temp['
       segmentDuration ']<=70)
97
    df_segment_temp = df_segment_temp.filter(df_segment_temp['
       segmentSpeedKmH']<=130)
    df_segment_temp = df_segment_temp.filter(df_segment_temp['
98
       segmentSpeedKmH']!=0)
99
    "," adding columns related to grid table ","
100
    schema_segments_ultimate ='deviceId_String,
101
```

```
102
   type Integer,
   dateTime_string ,
103
   startLatitude_double,
104
   uuuuuuuuuuuuuuuuuuustartLongitude.double,\
105
   startEngineStatus_Integer,
106
   startAccuracyDrop_Integer,
107
   understeinen und end Accuracy Drop, Integer .
108
   \label{eq:endergineStatus} end EngineStatus Integer , \\ \dot{\ }
109
   undlatitude_double,
110
   unununununununun endLongitude_double,
111
   segmentDistance_double,
112
   segmentDuration_double,
113
   uuuuuuuuuuuuuuuuuuusegmentSpeedKmHudouble,\
114
   \label{eq:theta} xtmp_double , \label{eq:theta}
115
   ytmp.double,
116
   uuuuuuuuuuuuuuustart_x integer,\
117
   118
   119
   \label{eq:end_y_integer} \label{eq:end_y_integer} \label{eq:end_y_integer} \label{eq:end_y_integer}
120
   cell_id_start.string,
121
122
   uuuuuuuuuuuuuuuuuuuuuuucell_id_endustring '
123
   '''shift by 100m in degrees'''
124
   shiftInMeterLat = 0.0008983152841182118 #100m
125
   shiftInMeterLon = 0.001270644533487797
126
   # ==
      ====== Torino and countryside =====
127
128
   minLat = 44.96282106687191
129
   \min Lon = 7.502048016422193
   \max Lat = 45.19265016665321
130
   \max Lon = 7.791812422724604
131
132
   # ====
133
   @F.pandas udf(schema segments ultimate, functionType=F.
134
     PandasUDFType.GROUPED MAP)
   def segment_position_calculator(df_segments):
135
      os.environ["ARROW PRE 0 15 IPC FORMAT"] = "1"
136
137
      df_segments = df_segments.sort_values(by="dateTime")
138
      df segments.reset index(drop = True, inplace=True)
139
140
```

141	", calculating the cells that these segments belong to
142	df_segments['xtmp'] = (df_segments['startLongitude']- minLon)/shiftInMeterLon
143	df_segments['ytmp'] = (df_segments['startLatitude']- minLat)/shiftInMeterLat
144	$df\_segments['start\_x'] = (df\_segments['xtmp'].apply(lambda x : np.floor(x)))+1$
145	df_segments['start_y'] = (df_segments['ytmp'].apply( lambda y : np.floor(y)))+1
146	
147	df_segments['xtmp'] = (df_segments['endLongitude']- minLon)/shiftInMeterLon
148	df_segments['ytmp'] = (df_segments['endLatitude']- minLat)/shiftInMeterLat
149	df_segments['end_x'] = (df_segments['xtmp'].apply( lambda x : np_floor(x)))+1
150	df segments ['end v'] = ( $df$ segments ['ytmp'], apply(
100	lambda v : np. floor(v)) + 1
151	df segments ['cell id start'] = $[f'{int(df segments.}]$
	start_y.values[i])}_{int(df_segments.start_x.values[ i])}'
152	<b>for</b> i <b>in range</b> ( <b>len</b> (df segments))]
153	df segments ['cell id end'] = $[f' \{ int (df segments.end y) \}$
	values [i]) }_{int (df_segments.end_x.values [i]) }'
154	<b>for</b> i <b>in range</b> ( <b>len</b> (df_segments))]
155	return df_segments
156	
157	df_segment = df_segment_temp.groupby("deviceId").apply(
	$segment\_position\_calculator)$
158	
159	df_segment.write.csv('file:///%s/segment_table'%path, sep=",")

Grid creation on the map with cells with size  $100\mathrm{m}$  by  $100\mathrm{m}$ 

```
1 import pandas as pd
2 import geopandas as gpd
3 from shapely.geometry import Point, Polygon, shape
4 import math
5 
6 df_segments = pd.read_csv(r'segment_table.csv')
```

```
7
   \#\%\% boundary of interest area
8
9 \text{ minLat} = 44.96282106687191
10 \text{ minLon} = 7.502048016422193
11 | \max Lat = 45.19265016665321
   \max Lon = 7.791812422724604
12
13
14 #%% creating empty dataframes to initialize
15
   df_geojson = pd.DataFrame()
   df_geojson['polygon_id'] = df_geojson['polygon'] = ""
16
   df data = pd.DataFrame()
17
18
19 #%%
20
   shiftInMeterLat = 0.0008983152841182118 #100m
21
   shiftInMeterLon = 0.001270644533487797
22
   \# grid creation with 100m by 100m steps
23
24
   range_lat = (maxLat - minLat)/shiftInMeterLat
25
   range_lon = (maxLon - minLon)/shiftInMeterLon
   coords = \{\}
26
   \min \text{LonCopy} = \min \text{Lon}
27
28
   #creating cells with corresponding id
29
   for i in range(math.ceil(range lat)):
30
        \min Lon = \min Lon Copy
31
        new_latitude = minLat + shiftInMeterLat
32
        for j in range(math.ceil(range lon)):
33
            new_longitude = minLon + shiftInMeterLon
34
            \operatorname{coords} [f'\{i+1\}_{j+1}] = [[\min \operatorname{Lon}, \min \operatorname{Lat}]],
                new_longitude, minLat],
35
            [new_longitude, new_latitude], [minLon, new_latitude]]
36
            minLon = new_longitude
37
        minLat = new_{-} latitude
38
39 #%% storing in a dataframe
   df_geojson['polygon_id']=coords.keys()
40
   df_geojson['polygon'] = coords.values()
41
   df_geojson ['polygon'] = df_geojson ['polygon']. apply (Polygon
42
43
   df data ['id'] = coords.keys()
44
45 | df_data [ 'count' ] = 0
```

## 46 |df\_data.to\_csv('100m/data\_TO.csv', index=False)

Grid table and plot over the map

1	from pyspark.sql import functions as F
2	import matplotlib.pyplot as plt
3	<pre>from pyspark.sql.types import *</pre>
4	import seaborn as sns
5	import pandas as pd
6	import numpy as np
7	import geopy.distance
8	import matplotlib
9	import datetime
10	import pickle
11	import os
12	
13	#schema of segment table
14	schema_segments ='deviceId_String, $\langle$
15	
16	$ \qquad \qquad$
17	$\qquad \qquad $
18	$ = startLongitude_double, \$
19	$ \qquad \qquad$
20	$ \qquad \qquad$
21	uuuuuuuuuuuuuuuuuuuuuuuuuuuuuuuuuuuu
22	$\cdots \cdots $
23	uuuuuuuuuuuuuuuuendLatitudeudouble,\
24	uuuuuuuuuuuuuuuendLongitudeudouble,
25	$\qquad \qquad $
26	$ \qquad \qquad$
27	$\square \square $
28	$ \qquad \qquad$
29	$\cdots \cdots \cdots \cdots \cdots \cdots \cdots \cdots y tmp_{\cup} double , \setminus$
30	$\cdots \cdots \cdots \cdots \cdots \cdots \cdots \cdots s tart \_x \_integer , \land$
31	uuuuuuuuuuustart_yuinteger,\
32	uuuuuuuuuuuuuuuuend_xuinteger,\
33	$ \qquad \qquad$
34	$\cdots \cdots \cdots \cdots \cdots \cdots \cdots \cdots cell\_id\_start\_string, \setminus$
35	······································
36	#reading csv
37	path = os.path.abspath(os.getcwd())

```
df_segments = spark.read.csv('file:///%s/segment_table'%
38
      path, sep=", ", schema = schema_segments)
39
   df_segments = df_segments.withColumn('dateTime', F.
      to timestamp('dateTime'))
40
   \#converting dates to hour (for peak selection)
41
42 @F.udf()
   def hour_extraction(dateTime):
43
44
       hour = dateTime.hour
45
       return hour
46
47
   df_segments = df_segments.withColumn('dateTime',
      hour_extraction('dateTime'))
48
   ',''Peak 1',''
49
   @F.udf(returnType=BooleanType())
50
   def peak_1(hour):
51
52
       valid_period = [9, 10, 11]
53
       inside_peak = hour in valid_period
54
       return inside_peak
55
56
   df_peak_1 = df_segments.filter(peak_1('dateTime'))
57
   ','Peak 2','
58
   @F.udf(returnType=BooleanType())
59
   def peak_2(hour):
60
61
       valid_period = [18, 19, 20]
62
       inside_peak = hour in valid_period
63
       return inside_peak
64
65
   df_peak_2 = df_segments.filter(peak_2('dateTime'))
66
   ', 'Non-Peak', ',
67
68
   @F.udf(returnType=BooleanType())
69
   def off_peak(hour):
70
       valid_period = range(1,7)
71
       inside_peak = hour in valid_period
72
       return inside_peak
73
   df off peak = df segments. filter (off peak ('dateTime'))
74
75
```

```
76
   \# grid table
    schema grid calculations= StructType([
77
        StructField("cell_id",StringType(),True),
78
        StructField("num_segments", IntegerType(), True),
79
80
        StructField("minSpeed", DoubleType(), True),
81
82
        StructField("avgSpeed", DoubleType(), True),
        StructField("maxSpeed", DoubleType(), True),
83
84
85
        StructField("minAccuracyDrop", IntegerType(), True),
        StructField("avgAccuracyDrop", DoubleType(), True),
StructField("maxAccuracyDrop", IntegerType(), True),
86
87
88
      ])
89
90
    @F.pandas_udf(schema_grid_calculations, functionType=F.
       PandasUDFType.GROUPED_MAP)
    def cell data calculator(df segments):
91
        os.environ["ARROW PRE 0 15 IPC FORMAT"] = "1"
92
93
          df_segments = df_segments.sort_values(by="dateTime")
94
    #
95
        df_segments.reset_index(drop = True, inplace=True)
96
        df data = pd.DataFrame()
97
98
99
        #______ min, avg, max calculations for "SPEED
100
               _____#
101
        minSpeed = min(df_segments['segmentSpeedKmH'])
        avgSpeed = np.mean(df_segments['segmentSpeedKmH'])
102
        maxSpeed = max(df_segments['segmentSpeedKmH'])
103
        #_____
104
                                                   ___#
        #======= min, avg, max calculations for ACCURACY DROP
105
           _____#
        accuracyDrops = list (df_segments['startAccuracyDrop'])
106
           + list (df_segments ['endAccuracyDrop'])
        minAccuracyDrop = min(accuracyDrops)
107
        avgAccuracyDrop = np.mean(accuracyDrops)
108
        maxAccuracyDrop = max(accuracyDrops)
109
110
        #
                                                     #
            writing calculated parameters in
        #=
111
           df data _____#
```

```
112
        cell_id = df_segments['cell_id_start'][0]
        df_data['cell_id'] = [cell_id]
113
114
115
        num segments = len(df segments)
        df_data['num_segments'] = [num_segments]
116
        df_data['num_segments'] = num_segments
117
118
        df_data['minSpeed'] = round(minSpeed, 2)
119
        df_data['avgSpeed'] = round(avgSpeed, 2)
120
121
        df data ['maxSpeed'] = round (maxSpeed, 2)
122
123
        df_data ['minAccuracyDrop'] = minAccuracyDrop
124
        df_data['avgAccuracyDrop'] = round(avgAccuracyDrop, 2)
        df_data['maxAccuracyDrop'] = maxAccuracyDrop
125
126
127
        return df_data
128
    df_grid_peak_1 = df_peak_1.groupby("cell_id_start" or "
129
       cell_id_end").apply(cell_data_calculator)
    df_grid_peak_1.toPandas().to_csv('grids/grid_table_peak_1.
130
       csv', index = False)
131
    df_grid_peak_2 = df_peak_2.groupby("cell_id_start" or "
132
       cell_id_end").apply(cell_data_calculator)
    df_grid_peak_2.toPandas().to_csv('grids/grid_table_peak_2.
133
       csv', index = False)
134
135
    df_grid_off_peak = df_off_peak.groupby("cell_id_start" or "
       cell id end").apply(cell data calculator)
136
    df_grid_off_peak.toPandas().to_csv('grids/
       grid_table_off_peak.csv', index = False)
```

```
1 #Import Libraries
2 import geopandas as gpd
3 import pandas as pd
4 import numpy as np
5 import folium
6 from folium.features import GeoJsonTooltip
7
8 #Read the geoJSON file using geopandas
9 geojson = gpd.read_file(r'../../100m/coords_TO.geojson')
```

10	
11	#Read the grid tables
12	df_grid_peak_1=pd.read_csv(r'/grids/grid_table_peak_1.csv
	,)
13	df_grid_peak_2=pd.read_csv(r'/grids/grid_table_peak_2.csv ')
14	df_grid_off_peak=pd.read_csv(r'/grids/grid_table_off_peak .csv')
15	
16	df_final_peak_1 = geojson.merge(df_grid_peak_1, left_on=" polygon_id", right_on="cell_id", how="outer")
17	df_final_peak_1 = df_final_peak_1[~df_final_peak_1[' geometry'], isna()]
18	df final peak 1.dropna(inplace = True)
19	
20	df_final_peak_2 = geojson.merge(df_grid_peak_2, left_on="
21	df final peak $2 = df$ final peak $2 [\sim df$ final peak $2 ['$
	$\begin{bmatrix} - & -i \\ geometry' \end{bmatrix} \cdot isna() \end{bmatrix}$
22	df_final_peak_2.dropna(inplace = True)
23	
20	
23 24	df_final_off_peak = geojson.merge(df_grid_off_peak , left_on="polygon_id", right_on="cell_id", how="outer")
23 24 25	<pre>df_final_off_peak = geojson.merge(df_grid_off_peak , left_on="polygon_id", right_on="cell_id", how="outer") df_final_off_peak = df_final_off_peak [~df_final_off_peak[ 'geometry'].isna()]</pre>
23 24 25 26	<pre>df_final_off_peak = geojson.merge(df_grid_off_peak , left_on="polygon_id", right_on="cell_id", how="outer") df_final_off_peak = df_final_off_peak [~df_final_off_peak[ 'geometry'].isna()] df_final_off_peak .dropna(inplace = True)</pre>
<ul> <li>23</li> <li>24</li> <li>25</li> <li>26</li> <li>27</li> </ul>	<pre>df_final_off_peak = geojson.merge(df_grid_off_peak ,     left_on="polygon_id", right_on="cell_id", how="outer") df_final_off_peak = df_final_off_peak [~df_final_off_peak[     'geometry'].isna()] df_final_off_peak .dropna(inplace = True)</pre>
<ul> <li>23</li> <li>24</li> <li>25</li> <li>26</li> <li>27</li> <li>28</li> </ul>	<pre>df_final_off_peak = geojson.merge(df_grid_off_peak ,     left_on="polygon_id", right_on="cell_id", how="outer") df_final_off_peak = df_final_off_peak [~df_final_off_peak[</pre>
23 24 25 26 27 28 29	<pre>df_final_off_peak = geojson.merge(df_grid_off_peak , left_on="polygon_id", right_on="cell_id", how="outer") df_final_off_peak = df_final_off_peak [~df_final_off_peak[ 'geometry'].isna()] df_final_off_peak .dropna(inplace = True) #management before plot df_final_dict_peak_1 = df_final_peak_1.set_index('cell_id') ['avgSpeed']</pre>
<ol> <li>23</li> <li>24</li> <li>25</li> <li>26</li> <li>27</li> <li>28</li> <li>29</li> <li>30</li> </ol>	<pre>df_final_off_peak = geojson.merge(df_grid_off_peak , left_on="polygon_id", right_on="cell_id", how="outer") df_final_off_peak = df_final_off_peak [~df_final_off_peak[ 'geometry'].isna()] df_final_off_peak .dropna(inplace = True) #management before plot df_final_dict_peak_1 = df_final_peak_1.set_index('cell_id') ['avgSpeed'] df_final_dict_peak_2 = df_final_peak_2.set_index('cell_id') ['avgSpeed']</pre>
23 24 25 26 27 28 29 30 31	<pre>df_final_off_peak = geojson.merge(df_grid_off_peak , left_on="polygon_id", right_on="cell_id", how="outer") df_final_off_peak = df_final_off_peak [~df_final_off_peak[ 'geometry'].isna()] df_final_off_peak .dropna(inplace = True) #management before plot df_final_dict_peak_1 = df_final_peak_1.set_index('cell_id') ['avgSpeed'] df_final_dict_peak_2 = df_final_peak_2.set_index('cell_id') ['avgSpeed'] df_final_dict_off_peak = df_final_off_peak.set_index('</pre>
<ul> <li>23</li> <li>24</li> <li>25</li> <li>26</li> <li>27</li> <li>28</li> <li>29</li> <li>30</li> <li>31</li> <li>32</li> </ul>	<pre>df_final_off_peak = geojson.merge(df_grid_off_peak , left_on="polygon_id", right_on="cell_id", how="outer") df_final_off_peak = df_final_off_peak [~df_final_off_peak[ 'geometry'].isna()] df_final_off_peak .dropna(inplace = True) #management before plot df_final_dict_peak_1 = df_final_peak_1.set_index('cell_id') ['avgSpeed'] df_final_dict_peak_2 = df_final_peak_2.set_index('cell_id') ['avgSpeed'] df_final_dict_off_peak = df_final_off_peak.set_index(' cell_id')['avgSpeed']</pre>
<ol> <li>23</li> <li>24</li> <li>25</li> <li>26</li> <li>27</li> <li>28</li> <li>29</li> <li>30</li> <li>31</li> <li>32</li> <li>33</li> </ol>	<pre>df_final_off_peak = geojson.merge(df_grid_off_peak , left_on="polygon_id", right_on="cell_id", how="outer") df_final_off_peak = df_final_off_peak [~df_final_off_peak[ 'geometry'].isna()] df_final_off_peak .dropna(inplace = True) #management before plot df_final_dict_peak_1 = df_final_peak_1.set_index('cell_id') ['avgSpeed'] df_final_dict_peak_2 = df_final_peak_2.set_index('cell_id') ['avgSpeed'] df_final_dict_off_peak = df_final_off_peak.set_index(' cell_id')['avgSpeed'] df_final_dict_conf_peak = df_final_off_peak.set_index(' cell_id')['avgSpeed']</pre>
<ol> <li>23</li> <li>24</li> <li>25</li> <li>26</li> <li>27</li> <li>28</li> <li>29</li> <li>30</li> <li>31</li> <li>32</li> <li>33</li> </ol>	<pre>df_final_off_peak = geojson.merge(df_grid_off_peak , left_on="polygon_id", right_on="cell_id", how="outer") df_final_off_peak = df_final_off_peak [~df_final_off_peak[ 'geometry'].isna()] df_final_off_peak .dropna(inplace = True) #management before plot df_final_dict_peak_1 = df_final_peak_1.set_index('cell_id') ['avgSpeed'] df_final_dict_peak_2 = df_final_peak_2.set_index('cell_id') ['avgSpeed'] df_final_dict_off_peak = df_final_off_peak.set_index(' cell_id')['avgSpeed'] df_final_dict_count_peak_1 = df_final_peak_1.set_index(' cell_id')['num_segments']</pre>
<ol> <li>23</li> <li>24</li> <li>25</li> <li>26</li> <li>27</li> <li>28</li> <li>29</li> <li>30</li> <li>31</li> <li>32</li> <li>33</li> <li>34</li> </ol>	<pre>df_final_off_peak = geojson.merge(df_grid_off_peak , left_on="polygon_id", right_on="cell_id", how="outer") df_final_off_peak = df_final_off_peak [~df_final_off_peak[ 'geometry'].isna()] df_final_off_peak .dropna(inplace = True) #management before plot df_final_dict_peak_1 = df_final_peak_1.set_index('cell_id') ['avgSpeed'] df_final_dict_peak_2 = df_final_peak_2.set_index('cell_id') ['avgSpeed'] df_final_dict_off_peak = df_final_off_peak.set_index(' cell_id')['avgSpeed'] df_final_dict_count_peak_1 = df_final_peak_1.set_index(' cell_id')['num_segments'] df_final_dict_count_peak 2 = df_final_peak 2.set_index('</pre>
<ol> <li>23</li> <li>24</li> <li>25</li> <li>26</li> <li>27</li> <li>28</li> <li>29</li> <li>30</li> <li>31</li> <li>32</li> <li>33</li> <li>34</li> </ol>	<pre>df_final_off_peak = geojson.merge(df_grid_off_peak , left_on="polygon_id", right_on="cell_id", how="outer") df_final_off_peak = df_final_off_peak [~df_final_off_peak[ 'geometry'].isna()] df_final_off_peak .dropna(inplace = True) #management before plot df_final_dict_peak_1 = df_final_peak_1.set_index('cell_id') ['avgSpeed'] df_final_dict_peak_2 = df_final_peak_2.set_index('cell_id') ['avgSpeed'] df_final_dict_off_peak = df_final_off_peak.set_index(' cell_id')['avgSpeed'] df_final_dict_count_peak_1 = df_final_peak_1.set_index(' cell_id')['num_segments']</pre>
<ol> <li>23</li> <li>24</li> <li>25</li> <li>26</li> <li>27</li> <li>28</li> <li>29</li> <li>30</li> <li>31</li> <li>32</li> <li>33</li> <li>34</li> <li>35</li> </ol>	<pre>df_final_off_peak = geojson.merge(df_grid_off_peak , left_on="polygon_id", right_on="cell_id", how="outer") df_final_off_peak = df_final_off_peak [~df_final_off_peak[ 'geometry'].isna()] df_final_off_peak .dropna(inplace = True) #management before plot df_final_dict_peak_1 = df_final_peak_1.set_index('cell_id') ['avgSpeed'] df_final_dict_peak_2 = df_final_peak_2.set_index('cell_id') ['avgSpeed'] df_final_dict_off_peak = df_final_off_peak.set_index(' cell_id')['avgSpeed'] df_final_dict_count_peak_1 = df_final_peak_1.set_index(' cell_id')['num_segments'] df_final_dict_count_peak_2 = df_final_peak_2.set_index(' cell_id')['num_segments'] df_final_dict_count_peak_3 = df_final_peak_2.set_index(' cell_id')['num_segments']</pre>

36	
37	#adding calculated parameters for plot
38	$m = folium . Map(location = [45.0703, 7.6869], zoom_start=12, tiles=None_overlay=False)$
39	import branca
40	import branca, colormap as cm
41	<pre>colorScale_speed = cm.LinearColormap(['darkred', 'red', ' orange', 'yellow', 'lime', 'green'], vmin=0., vmax=130., index=(df_final_peak_2['avgSpeed'].quantile ((0,0.2,0.4,0.6,0.8,1))).tolist())</pre>
42	
43	$colorScale\_speed.caption = Average_speed_[KmH] $
44	colorScale_speed.add_to(m)
40	fr off peak - folium FeatureCroup(name-'Speed[KmH] off
40	$\operatorname{real}_{\operatorname{peak}}$ =
47	$f_{\sigma}$ peak 1 = folium FeatureGroup(name='Speed[KmH]peakhours
11	9-10-11, overlay=False) add to(m)
48	fg peak 2 = folium. FeatureGroup(name='Speed [KmH]_peak_hours
	18-19-20, overlay=False). add to (m)
49	
50	$\# off\_peak$
51	folium . features . GeoJson (
52	data=df_final_off_peak,
53	$name = 'test \Box for \Box segment \Box density \Box ',$
54	smooth_factor=2,
55	style_function=lambda feature: {
56	'fillColor': colorScale_speed(
	df_final_dict_off_peak[feature['
	properties 'j[ 'polygon_id 'j]),
57	color : black ,
58	fillOne eitre 2,
09 60	lintOpacity : .0
61	f, tooltin-folium features GeolsonTooltin(
62	fields = ['cell id']
63	'num_segments'.
64	'minSpeed',
65	'avgSpeed',
66	
00	'maxSpeed',

68	'avgAccuracyDrop',
69	'maxAccuracyDrop '
70	],
71	$a \operatorname{liases} = ["Cell_{\sqcup} \operatorname{id}:"],$
72	"Number $_{\cup}$ of $_{\cup}$ Segments : ",
73	$Minimum_{\cup}Speed_{\cup}(Km/H)$ ,
74	'Average Speed (Km/H)',
75	$Maximum_{\cup}Speed_{\cup}(Km/H)$ ,
76	$Best_{\Box}GPS_{\Box}accuracy'$ ,
77	$Average \Box GPS \Box accuracy'$ ,
78	'Worst_GPS_accuracy'
79	],
80	localize=True,
81	sticky = False,
82	labels=True,
83	style = """
84	background-color: #F0EFEF;
85	border: 2px solid black;
86	border-radius: 3px;
87	box-shadow: 3px;
88	<i>" " "</i> ,
89	$\max\_width=800,),$
90	highlight_function=lambda x: { '
	weight ':3, 'fillOpacity': 2
	$\left\{ \right. ,$
91	$).add_to(fg_off_peak)$
92	$\# p eak\_1$
93	folium . features . GeoJson (
94	$data=df_final_peak_1$ ,
95	$name = 'test \Box for \Box segment \Box density \Box ',$
96	$smooth\_factor=2,$
97	style_function=lambda feature: {
98	'fillColor ': colorScale_speed(
	$df_final_dict_peak_1[feature[']$
	<pre>properties '][ 'polygon_id ']]) ,</pre>
99	'color' : 'black',
100	'weight ' : .2,
101	'fillOpacity': .6
102	},
103	tooltip=folium.features.GeoJsonTooltip(
	r · · · · · · · · · · · · · · · · · · ·

105	'num segments'
106	'minSpeed'
107	'avgSpood'
107	'maxSpeed'
100	'min A coursey Drop '
1109	'aug A coursey Drop '
111	'mayAccuracyDrop',
111 110	
112	], aliases — ["Cell.id:"
114	"Number, of Segments:"
115	'Minimum, Speed, (Km/H) '
116	'Average, Speed, (Km/H) '
117	'Maximum, Speed, (Km/H) '
118	'Best_GPS_accuracy'
119	'Average GPS accuracy',
120	'Worst GPS accuracy'
121	],
122	localize=True,
123	sticky = False,
124	labels = True,
125	style = """
126	background-color: #F0EFEF;
127	border: 2px solid black;
128	border-radius: 3px;
129	box-shadow: 3px;
130	<i>"""</i> ,
131	$\max_{\text{width}=800,)}$ ,
132	highlight_function=lambda x: { '
	weight ':3, 'fillOpacity ': 2
100	$\},$
133	).add_to(fg_peak_1)
134 195	
100	#peak_2
130 $137$	data-df final peak 2
137	name-'tost for sogment density '
130	smooth factor $-2$
140	style function=lambda feature {
141	'fillColor' colorScale speed(
	df final dict peak 2 [feature]'
	properties ']['polygon_id']])

Source code

142	'color' : 'black',
143	'weight': .2,
144	'fillOpacity': .6
145	$\},$
146	tooltip=folium.features.GeoJsonTooltip(
147	$fields = ['cell_id',$
148	'num_segments',
149	'minSpeed ',
150	'avgSpeed',
151	'maxSpeed',
152	'minAccuracyDrop ',
153	`avgAccuracyDrop`,
154	'maxAccuracyDrop '
155	],
156	$a li a s e s = ["Cell_{\sqcup} id:"],$
157	$"Number_{\cup} of_{\cup} Segments:"$ ,
158	$\operatorname{Minimum}_{\Box}\operatorname{Speed}_{\Box}(\operatorname{Km}/\operatorname{H})$ ,
159	$\operatorname{Average}_{\sqcup}\operatorname{Speed}_{\sqcup}(\operatorname{Km/H})$ ,
160	$\operatorname{Maximum}_{\Box}\operatorname{Speed}_{\Box}(\operatorname{Km}/\operatorname{H})$ ',
161	$Best \Box GPS \Box accuracy',$
162	$Average \square GPS \square accuracy'$ ,
163	$'Worst_{\Box}GPS_{\Box}accuracy'$
164	],
165	localize=True,
166	sticky = False,
167	labels=True,
168	style = """
169	background-color: #F0EFEF;
170	border: 2px solid black;
171	border-radius: 3px;
172	box-shadow: 3px;
173	<i>"""</i> ,
174	$\max_{width=800,)}$ ,
175	highlight_function=lambda x: { '
	weight ':3, 'fillOpacity ': 2
	},
176	).add_to(fg_peak_2)
177	tolium. TileLayer ('openstreetmap', overlay=True, name="
	OpenStreetMap").add_to(m)
178	tolium . LayerControl ( collapsed=False ) . add_to (m)
179	m. save ( "peaks_100m_speed . html" )

Traffic congestion index

```
import geopandas as gpd
1
2 import pandas as pd
3 import numpy as np
  import matplotlib.pyplot as plt
4
   import matplotlib
5
6
7
   "'reading the grid tables for the peak and off peak hours
   off_peak_table = pd.read_csv(r'grids/grid_table_off_peak.
8
      csv')
9
   peak_table_1 = pd.read_csv(r'grids/grid_table_peak_1.csv')
10
   peak_table_2 = pd.read_csv(r'grids/grid_table_peak_2.csv')
   off_peak_cells = list (off_peak_table['cell_id'])
11
   peak 1 cells = list(peak_table_1['cell_id'])
12
   peak 2 cells = list (peak table 2 ['cell id'])
13
14
15
   '''finding the intersection of grid cells '''
16
   intersection = list (set (peak_1_cells) & set (peak_2_cells) &
       set(off_peak_cells))
17
   ''' dropping uncommon cells '''
18
   df_peak_1 = peak_table_1.loc[peak_table_1['cell_id'].isin(
19
      intersection)]
20
   df peak 1. reset index (inplace = True)
21
22
   df_{peak_2} = peak_{table_2.loc} [peak_{table_2} 'cell_id'].isin(
      intersection)]
23
   df_peak_2.reset_index(inplace = True)
24
25
   df off peak = off peak table.loc[off peak table['cell id'].
      isin(intersection)]
   df_off_peak.reset_index(inplace = True)
26
27
28
   #merging grid tables with useful columns
29
   df_merged = pd.merge(pd.merge(df_off_peak[['cell_id', '
      num segments', 'avgSpeed']],
                         df_peak_1[['cell_id', 'num_segments', '
30
                            avgSpeed ']], on='cell_id'),
                         df_peak_2[[ 'cell_id ', 'num_segments', '
31
                            avgSpeed ']], on='cell id')
```

```
32
   df_merged.columns = ['cell_id', 'num_segments_offPeak', '
33
      avgSpeed_offPeak',
34
                          'num segments peak1', 'avgSpeed peak1'
                          'num_segments_peak2', 'avgSpeed_peak2'
35
36
37
38
   min num segments = list(zip(df merged.num segments offPeak),
      df merged.num segments peak1))
39
   \min\_num\_segments = [\min(i[0], i[1]) \text{ for } i \text{ in }
      min num segments]
40
   df merged ['min num segments off p1'] = min num segments
41
42
   min_num_segments = list (zip (df_merged.num_segments_offPeak,
      df merged.num segments peak2))
43
   \min\_num\_segments = [\min(i[0], i[1]) \text{ for } i \text{ in }
      min_num_segments]
   df_merged [ 'min_num_segments_off_p2 '] = min_num_segments
44
45
46
   ''' Absolute traffic congestion
   df absCongestionKMH = df merged.copy()
47
   df_absCongestionKMH['speedDiffKMH_off_p1'] =
48
      df_absCongestionKMH['avgSpeed_offPeak'] -
      df_absCongestionKMH['avgSpeed_peak1']
49
   df_absCongestionKMH['speedDiffKMH_off_p2'] =
      df_absCongestionKMH['avgSpeed_offPeak'] -
      df absCongestionKMH['avgSpeed peak2']
50
51
   #CDF plot
   speedDiffs = { 'speedDiffKMH off p1' : df absCongestionKMH[ '
52
      speedDiffKMH off p1'],
        'speedDiffKMH_off_p2' : df_absCongestionKMH['
53
          speedDiffKMH_off_p2']}
   # matplotlib.rcParams.update(matplotlib.rcParamsDefault)
54
   plt.rcParams ["figure.figsize"] = (10,5)
55
   for key, value in speedDiffs.items():
56
57
       x = list(value)
58
       x.sort()
       y = (1. * np.arange(len(x)) / (len(x) - 1))*100
59
```

```
60
61
62
       plt.plot(x, y, zorder=2)
63
          plt. hlines(y=(len(df absCongestionKMH.loc[value <0]))
   #
       / len(df_absCongestionKMH))*100
                      xmin = min(x), xmax = max(x), colors =
64
   #
      colors [0], ls = ': '
65
   title = 'Cellwise \BoxSpeed \Box Difference \Box Cumulative \Box Distribution '
66
67
   plt.title(title)
   # title = 'Absolute Congestion'
68
69
   plt.xlabel('Absolute_Cellwise_speed_difference_(KM/H)')
   plt.ylabel('Percent<sub>1</sub>%')
70
   # plt.ylim(0,1)
71
72
   plt.xlim(-125,125)
73
   plt.grid(True)
   plt.yticks(np.arange(0,110,10))
74
75
   plt.xticks(np.arange(-130,131,20))
   plt.legend (['Morning_Peak_VS_off-Peak', 'Evening_Peak_VS_
76
      off-Peak'], loc ='lower_right')
77
78
   plt.axvline (x=0,ymin = 0, ymax = 1, color='black', linewidth
       = 1, zorder=1)
   plt.axhline(y=len(df_absCongestionKMH.loc[
79
      df_absCongestionKMH['speedDiffKMH_off_p1'] <0 ]) / len(
      df absCongestionKMH) *100
               xmin = 0, xmax = 1, color = '#1f77b4', ls=':')
80
81
   plt.axhline(y=(len(df_absCongestionKMH.loc[
      df_absCongestionKMH['speedDiffKMH_off_p2'] <0 ]) / len(
      df_absCongestionKMH))*100
82
               xmin = 0, xmax = 1, color = '#ff7f0e', ls=':')
   title='abs congestion p1 2'
83
   # plt.savefig(f'{ title }_CDF.jpg', bbox_inches='tight')
84
85
   plt.show()
86
   plt.rcParams ["figure.figsize"] = (10,5)
87
   |# matplotlib.rcParams.update(matplotlib.rcParamsDefault)
88
89 |x_p1 = df_absCongestionKMH['speedDiffKMH_off_p1']
90 | x p2 = df absCongestionKMH ['speedDiffKMH off p2']
   bins = np. arange(-125, 126, 5)
91
92 | plt.hist([x_p1,x_p2],bins=bins,edgecolor="lightblue")
```

```
plt.xticks (np.arange(-75,76,5))
93
94
    plt.xlim(-70,70)
    plt.yticks(np.arange(0,10001,1000))
95
96
    \# plt.xlim(-50,10)
    title='Cellwise_Speed_Difference_Distribution'
97
    plt.title(title)
98
    plt.xlabel('Cellwise_Speed_Difference(KM/H)')
99
100
101
    plt.ylabel('Frequency')
102
    plt.legend(['Morning_Peak_VS_off-Peak', 'Evening_Peak_VS_
       off-Peak'], loc ='best')
103
    plt.axvline (x=0, ymin = 0, ymax = 1, color='black', linewidth
        = 1.5)
    plt.grid()
104
    title = 'abs_congestion_p1_2'
105
106
    # plt.savefig(f'{ title }_Histogram.jpg', bbox_inches='tight')
    plt.show()
107
108
    ''' Relative traffic congestion
109
    df relCongestion = df merged.copy()
110
111
112
    #relative speed difference calculation
    def relative speed difference (row, avgSpeed):
113
        if row['avgSpeed_offPeak'] >= row[avgSpeed]:
114
115
            speedDiff = (row['avgSpeed_offPeak'] - row[avgSpeed
                ]) /row ['avgSpeed offPeak']
116
        else:
117
            speedDiff = -(row[avgSpeed] - row['avgSpeed_offPeak]
                '])/row[avgSpeed]
118
        return speedDiff*100
119
    df relCongestion ['speedDiff% p1'] = df relCongestion.apply(
120
       lambda row: relative speed difference(row,'
       avgSpeed_peak1'), axis=1)
    df_relCongestion['speedDiff%_p1'] = [float("{0:.2 f}".format
121
       (i)) for i in df_relCongestion['speedDiff%_p1']]
122
    df_relCongestion['speedDiff%_p2'] = df_relCongestion.apply(
123
       lambda row: relative speed difference(row,'
       avgSpeed_peak2'), axis=1)
```

```
df_relCongestion ['speedDiff%_p2'] = [float(" \{0:.2f\}".format
124
       (i)) for i in df relCongestion ['speedDiff% p2']]
125
126
    #CDF plot
    speedDiffs = { 'speedDiff% p1' : df relCongestion [ 'speedDiff
127
       %_p1'],
         'speedDiff% p2' : df relCongestion ['speedDiff% p2']}
128
    plt.rcParams ["figure.figsize"] = (10,5)
129
130
    for key, value in speedDiffs.items():
131
        x = list(value)
132
        x.sort()
        y = (1. * np.arange(len(x)) / (len(x) - 1))*100
133
134
    #
           matplotlib.rcParams.update(matplotlib.rcParamsDefault
135
136
        plt.plot(x, y)
137
    title = 'Cellwise \Box Speed \Box Difference \Box Cumulative \Box Distribution '
138
139
    plt.title(title)
    title = 'Relative_Congestion'
140
    plt.xlabel('Relative \Box Cellwise \Box speed \Box difference \Box(\%)')
141
142
    plt.ylabel('Percent_%')
    \# plt.ylim(0,1)
143
    # plt.xlim(−1000,160)
144
    plt.grid(True)
145
    plt.yticks(np.arange(0,110,10))
146
    plt.xticks(np.arange(-100,101,10))
147
148
    \# plt.hlines(y=negative_values_percentage, xmin = min(x),
       xmax = max(x), colors = 'red', ls = ':')
149
    plt.legend(['Morning_Peak_VS_off-Peak', 'Evening_Peak_VS_
       off-Peak'], loc ='lower_right')
150
151
    plt.axvline (x=0, ymin = 0, ymax = 1, color='black', linewidth
        = 1)
    plt.axhline(y=len(df_relCongestion.loc[df_relCongestion['
152
       speedDiff%_p1'] <0 ]) / len(df_relCongestion)*100</pre>
                xmin = 0, xmax = 1, color = '#1f77b4', ls=':')
153
    plt.axhline(y=(len(df_relCongestion.loc[df_relCongestion['
154
       speedDiff% p2'] <0 ]) / len(df absCongestionKMH))*100
                xmin = 0, xmax = 1, color = '#ff7f0e', ls=':')
155
   title='rel congestion p1 2'
156
```

```
plt.savefig(f'{ title }_CDF.jpg', bbox_inches='tight')
157
158
    plt.show()
159
160
    plt.rcParams ["figure.figsize"] = (10,5)
161
    # matplotlib.rcParams.update(matplotlib.rcParamsDefault)
162
163
    x p1 = df relCongestion ['speedDiff% p1']
    x_p2 = df_relCongestion['speedDiff%_p2']
164
165
    bins = np. arange(-100, 100, 10)
166
    plt.hist([x_p1,x_p2],bins=bins,edgecolor="lightblue")
167
    plt.xticks(np.arange(-100,101,10))
168
    # plt.yticks(np.arange(0,3000,500))
169
    # plt.xlim(-50,10)
170
    title='Cellwise_Speed_Difference_Distribution'
171
172
    plt.title(title)
    plt.xlabel('Relative \Box Cellwise \Box speed \Box difference \Box(\%)')
173
174
    title = 'Relative \Box Congestion'
175
    plt.ylabel('Frequency')
176
177
    plt.grid()
178
    plt.legend (['Morning_Peak_VS_off-Peak', 'Evening_Peak_VS_
       off-Peak'], loc ='best')
    plt.axvline(x=0,ymin = 0, ymax = 1, color='black', linewidth
179
        = 1.5)
    title='rel congestion p1 2'
180
    plt.savefig(f'{title}_Histogram.jpg',bbox_inches='tight')
181
182
    plt.show()
183
                                                           , , ,
184
    ''' Cells behavior based on density investigation
    \#negative congestion extraction
185
    df rel negative conge = df relCongestion [df relCongestion ['
186
       speedDiff\% p2'] < 0]
187
    df_rel_positive_conge = df_relCongestion [ df_relCongestion [ '
188
       speedDiff\%_p2' >= 0
189
190
    #CDF of negative congestion based on number of segments per
        cell [relative]
    matplotlib.rcParams.update(matplotlib.rcParamsDefault)
191
```

```
dfs = [df_relCongestion, df_rel_positive_conge,
192
       df rel negative conge]
193
    for df in dfs:
194
        \mathbf{x} = \mathbf{list} (df ['min num segments off p2'])
195
        x.sort()
    \# y = df\_negative\_conge.speedDiffKMH
196
        y = (1. * np.arange(len(x)) / (len(x) - 1))*100
197
198
        plt.rcParams ["figure.figsize"] = (10,5)
199
200
        plt.plot(x,y)
201
    plt.title('Cumulative_Distribution_of_Segments_per_cell')
202
    \# title = 'Absolute Congestion'
    plt.xlabel('Number_of_segments_per_cell')
203
    plt.ylabel('Percent<sub>1</sub>%')
204
205
    \# plt.ylim(0,1)
206
    plt.xlim(-10,2000)
207
    plt.grid(True)
208
    plt.yticks(np.arange(0,110,10))
209
    plt.xticks(np.arange(0,2001,100))
    plt.legend(['Evening_peak_VS_Off-peak', 'Evening_peak_VS_Off
210
       -peak [Positively congested]', 'Evening peak VS Off-peak
       [Negatively congested]'], loc='lower right')
211
    \# plt. vlines (x=70. ymin = 0. ymax = max(y), colors='red', ls
       = ': ')
    title='congestion_num_segments'
212
    plt.savefig(f'{title}_CDF.jpg', bbox_inches='tight')
213
214
    plt.show()
215
216
    bins = np.arange(0, \max(x)+1, 100)
217
    plt.hist([list(df['min_num_segments_off_p2']) for df in dfs
       ], bins=bins, edgecolor="lightblue")
218
    title = 'Distribution \cup of \cup Segments \cup per \cup cell '
219
    plt.title(title)
220
    plt.xlabel('Number_of_segments_per_cell')
221
    plt.ylabel('Frequency')
222
    plt.xticks(np.arange(0,2001,100))
    plt.xlim(0,2000)
223
224
    plt.grid()
225
    plt.legend (['Evening_peak_VS_Off-peak', 'Evening_peak_VS_Off
       -\text{peak}_{\Box} [Positively_congested]', 'Evening_peak_VS_Off-peak_
       [Negatively congested]'])
```

```
226
    title='congestion_num_segments'
    plt.savefig(f'{title}_Histogram.jpg',bbox_inches='tight')
227
228
    plt.show()
229
230
    #behavior comparison with and without filtering
    df_relCongestion_above200 = df_relCongestion.loc
231
       df relCongestion ['min num segments off p2']>200]
232
    #CDF plot
    speedDiffs = { 'total ' : df_relCongestion [ 'speedDiff%_p2'],
233
234
         'above 100' : df relCongestion above100 ['speedDiff% p2'
            ],
         'above 200' : df relCongestion above200 ['speedDiff% p2'
235
                     }
    plt.rcParams ["figure.figsize"] = (10,5)
236
    for key, value in speedDiffs.items():
237
238
        x = list(value)
239
        x.sort()
240
        y = (1. * np.arange(len(x)) / (len(x) - 1))*100
241
242
243
    title = 'Cellwise \Box Speed \Box Difference \Box Cumulative \Box Distribution '
244
    plt.title(title)
245
    title = 'Relative_Congestion'
    plt.xlabel('Relative Cellwise speed difference (\%)')
246
    plt.ylabel('Percent_%')
247
248
    plt.grid(True)
249
    plt.yticks(np.arange(0,110,10))
250
    plt.xticks(np.arange(-100,101,10))
251
    plt.legend (['Evening_Peak_VS_Off-peak_[Not_filtered]',
252
                  'Evening_Peak_VS_Off-peak_[cells_containing_
                    >100 Segments]', 'Evening Peak VS Off-peak
                    cells_{\Box}containing_{\Box} > 200_{\Box}Segments]'])
253
    plt.axvline(x=0, ymin = 0, ymax = 1, color='black', linewidth
254
        = 1)
255
    title='rel_congestion_p2_above100'
256
    plt.savefig(f'{ title }_CDF.jpg', bbox_inches='tight')
257
    plt.show()
```

Applying filters on relative congestion index

1 **import** geopandas as gpd

```
import pandas as pd
2
3 import numpy as np
4
   import matplotlib.pyplot as plt
5
   import matplotlib
6
   "'reading the grid tables for the peak and off peak hours
7
   off_peak_table = pd.read_csv(r'grid_table_off_peak.
8
      csv')
9
   peak table 1 = pd.read csv(r'grids/grid table peak 1.csv')
   peak_table_2 = pd.read_csv(r'grids/grid_table_peak_2.csv')
10
   off_peak_cells = list (off_peak_table [ 'cell_id '])
11
   peak 1 cells = list (peak table 1 ['cell id'])
12
   peak_2_cells = list(peak_table_2['cell_id'])
13
14
   ''' dropping uncommon cells '''
15
   df_peak_1 = peak_table_1.loc[peak_table_1['cell_id'].isin(
16
      intersection)]
17
   df_peak_1.reset_index(inplace = True)
18
   df_peak_2 = peak_table_2.loc[peak_table_2['cell_id'].isin(
19
      intersection)]
20
   df peak 2.reset index (inplace = True)
21
   df_off_peak = off_peak_table.loc[off_peak_table['cell_id'].
22
      isin(intersection)]
23
   df_off_peak.reset_index(inplace = True)
24
   "", "merging related columns of off-peak and peak grid tables
25
       , ,
   df_merged = pd.merge(df_off_peak[['cell_id', 'num_segments',
26
      'avgSpeed']],
27
                                    df_peak_2[['cell_id','
                                       num_segments', 'avgSpeed'
                                       ]], on='cell_id')
28
   #minimum number of segment in cellwise comparison
29
  \min\_num\_segments = list(zip(df\_merged.num\_segments\_x))
30
      df merged.num segments y))
   \min\_num\_segments = [\min(i[0], i[1]) \text{ for } i \text{ in }
31
      min_num_segments]
```

```
df_merged ['min_num_segments'] = min_num_segments
32
33
34
   def relative_speed_difference(row):
35
       if row ['avgSpeed x'] >= row ['avgSpeed y']:
            speedDiff = (row['avgSpeed_x'] - row['avgSpeed_y'])
36
               /row['avgSpeed_x']
37
       else:
            speedDiff = -(row['avgSpeed_y'] - row['avgSpeed_x'
38
               ) /row [ 'avgSpeed_y']
39
       return speedDiff*100
40
41
   df_relCongestion = df_merged.copy()
   df_relCongestion['speedDiff%'] = df_relCongestion.apply(
42
      lambda row: relative_speed_difference(row), axis=1)
   df_relCongestion ['speedDiff%'] = [float(" \{0:.2 f\}".format(i)
43
      ) for i in df_relCongestion['speedDiff%']]
44
   ''' Applying filters '''
45
46
   df_relCongestion = df_relCongestion.loc[df_relCongestion['
      min num segments ']>100]
47
48
   # dropping negative congested cells
   df relCongestion = df relCongestion [df relCongestion ['
49
      speedDiff\%' >= 0
50
   #CDF plot
51
52
   matplotlib.rcParams.update(matplotlib.rcParamsDefault)
53
   plt.rcParams ["figure.figsize"] = (10,5)
   x = list (df_relCongestion ['speedDiff%'])
54
   x.sort()
55
   y = (1. * np.arange(len(x)) / (len(x) - 1))*100
56
   plt.plot(x, y)
57
58
   title = 'Cellwise \Box Speed \Box Difference \Box Cumulative \Box Distribution \Box
      [Positive_side]'
59
   plt.title(title)
   title = 'Relative \Box Congestion'
60
   plt.xlabel('Relative Cellwise speed difference (\%)')
61
   plt.ylabel('Percent_%')
62
63
   plt.grid(True)
64
   plt.yticks(np.arange(0,110,10))
   plt.xticks(np.arange(0, 101, 5))
65
```

```
66
   plt.savefig(f'{ title }_CDF.jpg', bbox_inches='tight')
67
68
   plt.show()
69
70
   plt.rcParams ["figure.figsize"] = (10,5)
71
   bins = np.arange(0, 100, 5)
72
73
   plt.hist(x,bins=bins,edgecolor="lightblue")
   plt.xticks (np.arange (0, 100, 5))
74
   title='Cellwise_Speed_Difference_Distribution_[Positive_
75
      side]'
   plt.title(title)
76
   plt.xlabel('Relative_Cellwise_Speed_Difference(%)')
77
78
79
   plt.ylabel('Frequency')
   plt.grid()
80
   title = 'rel_pos_congestion'
81
   plt.savefig(f'{title}_Histogram.jpg',bbox_inches='tight')
82
83
   plt.show()
84
85
   "' creating grid table of positively congested areas"
86
   positive_cells = list (df_relCongestion['cell_id'])
87
   off_peak_grid = df_off_peak.loc[df_off_peak['cell_id'].isin
88
      (positive_cells)]
   off_peak_grid = pd.merge(off_peak_grid[['cell_id', '
89
      num_segments', 'minSpeed', 'avgSpeed', 'maxSpeed',
90
                                             'minAccuracyDrop', '
                                               avgAccuracyDrop',
                                                 'maxAccuracyDrop
                                                ']],
                                    df_relCongestion [['cell_id',
91
                                       'speedDiff%', 'avgSpeed_y'
                                       ]], on='cell id')
92
   off_peak_grid.to_csv(r'4_2-grid-congestion-peaks/
      off_peak_above_100.csv', index= False)
93
   peak_1_grid = df_peak_1.loc[df_peak_1['cell_id'].isin(
94
      positive cells)]
   peak_1_grid = pd.merge(peak_1_grid[['cell_id', '
95
      num_segments', 'minSpeed', 'avgSpeed', 'maxSpeed',
```

```
96 'minAccuracyDrop', '
avgAccuracyDrop', '
avgAccuracyDrop', '
maxAccuracyDrop '
]],
97 df_relCongestion [['cell_id',
'speedDiff%']], on='
cell_id')
98 peak_1_grid.to_csv(r'4_2-grid-congestion-peaks/
peak_2_above_100.csv', index= False)
```

Congestion plot on the map

```
#Import Libraries
1
2 import geopandas as gpd
  import pandas as pd
3
   import numpy as np
4
   import folium
5
   from folium.features import GeoJsonTooltip
6
7
   #Read the geoJSON file using geopandas
8
   geojson = gpd.read_file(r'../../100m/coords_TO.geojson')
9
10
11
   df_grid_peak_2=pd.read_csv(r'peak_2_above_100.csv')
   df grid off peak=pd.read csv(r'off peak above 100.csv')
12
13
14
   df_final_off_peak = geojson.merge(df_grid_off_peak ,
      left_on="polygon_id", right_on="cell_id", how="outer")
   df_final_off_peak = df_final_off_peak [~df_final_off_peak]
15
      'geometry'].isna()]
   df final off peak .dropna(inplace = True)
16
17
   df_final_dict_rel_congestion = df_final_off_peak.set_index(
18
      'cell id') ['speedDiff%']
19
20
  m2 = folium.Map(location = [45.0703, 7.6869], zoom_start = 12,
      tiles=None, overlay=False)
21
   colorScale_relative_congestion = cm.LinearColormap(['
22
      aquamarine', 'palegreen', 'greenyellow', 'yellow',
      orangered', 'darkred'], vmin=0., vmax=100.,
```

23	index=(
	dt_tinal_ott_peak['
	speedD111% ].
	$\begin{array}{c} \text{quantile} \\ ((0, 0, 2, 0, 4, 0, 6, 0, 8, 1)) \end{array}$
	((0, 0.2, 0.4, 0.0, 0.8, 1))
24	
$\frac{24}{25}$	colorScale relative congestion caption - 'Belative
20	congestion (%)
26	$colorScale$ relative congestion add $to(m^2)$
$\frac{20}{27}$	
$\frac{21}{28}$	fg_relative_congestion = folium.FeatureGroup(name='Evening
-0	$Peak_VS_Off_peak_Relative_Congestion_(\%)', overlay=False$
	). add to $(m2)$
29	
30	#relative congestion
31	folium . features . GeoJson (
32	data=df_final_off_peak,
33	name='test $\_$ for $\_$ segment $\_$ density $\_$ ',
34	$\mathrm{smooth\_factor} = 2,$
35	style_function=lambda feature: {
36	'fillColor ':
	colorScale_relative_congestion(
	df_final_dict_rel_congestion[
	feature ['properties']['cell_id'
37	'color' : 'black',
38	'weight ' : .2,
39	'fillOpacity': .6
40	
41	tooltip=folium.features.GeoJsonTooltip(
42	$fields = [cell_id', ]$
43	speed Diff%,
44 45	avgSpeed,
40 46	avgSpeed_y ,
$\frac{40}{47}$	num_segments ,
48	aliases - ["Cell.id:"
40 40	, Relative Speed Difference
10	
50	$Off-peak Speed (KM/H) \cdot $

```
51
                                        'Evening_peak_Speed_(KM/H)
                                           : ',
                                        "Number_of_Segments:",
52
53
                                       ],
                              localize=True,
54
55
                              sticky=False,
56
                              labels=True,
                              style = """
57
58
                                  background-color: #F0EFEF;
59
                                  border: 2px solid black;
                                  border-radius: 3px;
60
61
                                  box-shadow: 3px;
                              " " "
62
63
                              \max_{width=800,)}
                                  highlight_function=lambda x: { '
64
                                     weight ':3, 'fillOpacity ': 2
                                     },
                              ).add_to(fg_relative_congestion)
65
66
   folium.TileLayer('openstreetmap', overlay=True, name="
67
      OpenStreetMap").add_to(m2)
68
   folium.LayerControl(collapsed=False).add_to(m2)
69
  m2.save("peaks_100m_relative_congestion_filt100.html")
70
```

KDE implementation

1 *#Import Libraries* 2 **import** folium 3 from folium.features import GeoJsonTooltip **import** pandas as pd 4 import numpy as np 56 from sklearn.neighbors import KernelDensity from sklearn.model\_selection import GridSearchCV 7 **import** matplotlib.pyplot as plt 8 **import** branca.colormap as cm 9 10 import geojsoncontour 11 #Read the geoJSON file using geopandas 12geojson = gpd.read\_file(r'centroids\_wgs.geojson') 13df\_grid=pd.read\_csv(r'../4\_2-grid-congestion-peaks/ 14off peak above 100.csv')

108
```
15
   df_merge = geojson.merge(df_grid, left_on="polygon_id",
16
      right_on="cell_id", how="outer")
   df_merge = df_merge[~df_merge['geometry'].isna()]
17
   df merge.dropna(inplace = True)
18
19
20
   #convert float percentages to the closest integer
21
   df_merge['speedDiff%'] = [round(i) for i in df_merge['
      speedDiff%']]
22
23
   #relative congestion as density
24
   df_multiplies = pd.DataFrame()
25
   for index, row in df merge.iterrows():
       multiplier = row['speedDiff%']
26
27
       df_multiplies = df_multiplies.append([row]*multiplier,
          ignore_index=True)
28
29
   #Grid
30
   #%% min and max for coordinates
31 \mid \text{minLat} = 44.96282106687191
32 \text{ minLon} = 7.502048016422193
33 \mid \max Lat = 45.19265016665321
34 \mid \text{maxLon} = 7.791812422724604
35
36 |tmp= list (df_merge['cell_id'].str.split('_'))
37
   x_num = len(np.unique(np.array(tmp)[:,0].astype(int)))
38
  y_num = len(np.unique(np.array(tmp)[:,1].astype(int)))
39
40
   #building a grid of point for sampling the resulting KDEs
   xx = np.linspace(minLon, maxLon, x_num)
41
   yy = np.linspace(minLat, maxLat, y_num)
42
43
44
   xg, yg = np.meshgrid(xx, yy)
45
   grid\_coords = np.c\_[xg.ravel(), yg.ravel()]
46
47
   # KDE
   colors = ['navy', 'blue', 'royalblue', 'cyan', 'yellow', 'orange
48
      ', 'orangered', 'red', 'darkred']
49
   kde_cols = ['centroid_lon', 'centroid_lat']
50
   bandwidths = [.0005, .001, .00150, .002, .0025, .003, .0035]
51 \mid
```

```
52
53
   #train model
54
   maps = []
55
   for i, bandwidth in enumerate(bandwidths):
       m = folium.Map(location = [45.0703, 7.6869], zoom_start
56
           =12, overlay=False, tiles='openstreetmap')
       model = KernelDensity( metric="euclidean", kernel="
57
           gaussian ", bandwidth=bandwidth). fit (df_multiplies [
           kde_cols])
58
       kde samples = np.exp(model.score samples(grid coords)).
           reshape(*xg.shape))
59
        levels = np.linspace(kde_samples.min(), kde_samples.max
           (), 10)
        colorScale = cm.LinearColormap(colors, vmin=round(
60
           kde_samples.min(), 2),
                                          vmax=round(kde_samples.
61
                                             \max(), 2), caption = '
                                             Score Sample').
                                             to\_step(10)
        colorScale.caption = f'Score_{\Box}Sample_{\Box}of_{\Box}bandwidth_{\Box}=_{\Box}
62
           bandwidth } '
63
        colorScale.add_to(m)
64
        contourf = plt.contourf(xg, yg, kde_samples, levels =
           levels, alpha=0.6, colors = colorScale.colors)
65
66
       # Convert matplotlib contourf to geojson
67
       geojson = geojsoncontour.contourf_to_geojson(
68
            contourf=contourf,
69
            min angle \deg = 3.0,
70
            ndigits = 5,
            stroke_width = .5,
71
72
            fill opacity = 0.5)
73
74
       folium_geo = folium.GeoJson(
75
            geojson,
76
            style_function=lambda x: {
                 'color':
                               x['properties']['stroke'],
77
                               x['properties']['stroke-width'],
78
                 'weight':
79
                 'fillColor': x['properties']['fill'],
80
                 'opacity':
                               0.8.
            \}).add_to(m)
81
```

```
82 maps.append(m)
```

```
83 | for i,m in enumerate (maps):
```

```
84
```

m. save(f"bandwidth\_{bandwidths[i]}.html")

Filtering the KDE result with bandwidth = 0.001

```
1
  #Import Libraries
2 import folium
3 from folium.features import GeoJsonTooltip
4 import pandas as pd
5
   import numpy as np
6 from sklearn.neighbors import KernelDensity
7
   from sklearn.model_selection import GridSearchCV
8
   import matplotlib.pyplot as plt
   import branca.colormap as cm
9
10 import geojsoncontour
11
12
   #Read the geoJSON file using geopandas
13
   geojson = gpd.read file(r'drafts/centroids wgs.geojson')
14
   df_grid=pd.read_csv(r'../4_2-grid-congestion-peaks/
      off_peak_above_100.csv')
15
16 df_merge = geojson.merge(df_grid, left_on="polygon_id",
      right_on="cell_id", how="outer")
   df_merge = df_merge[~df_merge['geometry'].isna()]
17
   df merge.dropna(inplace = True)
18
19
20 \mid \# convert \ float \ percentages \ to \ the \ closest \ integer
   df_merge['speedDiff%'] = [round(i) for i in df_merge['
21
      speedDiff%']]
22
23
   #relative congestion as density
   df multiplies = pd.DataFrame()
24
   for index,row in df_merge.iterrows():
25
       multiplier = row['speedDiff%']
26
       df multiplies = df_multiplies.append([row]*multiplier,
27
          ignore_index=True)
28
29
  \#Grid
30 \mid \#\%\% min and max for coordinates
31 \mid \text{minLat} = 44.96282106687191
32 \min Lon = 7.502048016422193
```

```
\max Lat = 45.19265016665321
33
34
  \max Lon = 7.791812422724604
35
  tmp= list (df_merge['cell_id'].str.split('_'))
36
37
   x_num = len(np.unique(np.array(tmp)[:,0].astype(int)))
38
   y_num = len(np.unique(np.array(tmp)[:,1].astype(int)))
39
   #building a grid of point for sampling the resulting KDEs
40
   xx = np.linspace(minLon, maxLon, x_num)
41
42
   yy = np. linspace(minLat, maxLat, y num)
43
44
   xg, yg = np.meshgrid(xx, yy)
   grid\_coords = np.c\_[xg.ravel(), yg.ravel()]
45
46
47
   colors = ['black', 'red']
48
49
   #KDE
50
   kde_cols = ['centroid_lon', 'centroid_lat']
   bandwidth = .001 \# approx 100m
51
52
  m = folium . Map(location = [45.0703, 7.6869], zoom_start = 12,
53
      overlay=False , tiles='openstreetmap ')
54
55
   threshold = 74
   kde_samples [kde_samples < threshold] = 0
56
57
   levels = np.array(0)
58
   levels = np.append(levels, np.linspace(threshold,
      kde_samples.max(), 2))
   levels = [round(1,2) for 1 in levels]
59
60
   colorScale = cm.LinearColormap(colors, vmin=round(
      kde_samples.min(),2), vmax=round(kde_samples.max(),2)).
      to step(index=levels)
61
62
   colorScale.caption = f'Score_Sample_of_bandwidth_=_{1}
      bandwidth } '
63
   colorScale.add to (m)
   contourf = plt.contourf(xg, yg, kde_samples, levels = levels
64
      , alpha=0.9, colors = colorScale.colors)
65
66
   # Convert matplotlib contourf to geojson
67
   geojson = geojsoncontour.contourf_to_geojson(
```

```
68
          contourf=contourf,
          \min_{\text{angle}} \deg = 3.0,
69
70
          n digits = 5,
71
          stroke_width = .5,
          fill_opacity=1)
72
73
    folium_geo = folium.GeoJson(
74
75
          geojson,
76
          style_function=lambda x: {
                ``color ': x[ 'properties '][ 'stroke '],
'weight ': x[ 'properties '][ 'stroke-width '],
'fillColor ': x[ 'properties '][ 'fill '],
77
78
79
                'opacity': 1,
80
81
          \}).add_to(m)
82
83 m. save(f "congested areas.html")
```