



HELSINGIN YLIOPISTO
HELSINGFORS UNIVERSITET
UNIVERSITY OF HELSINKI

Master's Thesis
Geography
Geoinformatics

PATTERNS OF AGGREGATED COMMUTING TIMES
IN THE HELSINKI CAPITAL REGION

Pinja-Liina Jalkanen

2020

Supervisor: Petteri Muukkonen
Co-supervisor: Olle Järv

UNIVERSITY OF HELSINKI
FACULTY OF SCIENCE
DEPARTMENT OF GEOSCIENCES AND GEOGRAPHY
GEOGRAPHY

PL 64 (Gustaf Hällströmin katu 2)
00014 Helsingin yliopisto



Tiedekunta – Fakultet – Faculty Faculty of Science		Osasto – Institution – Department Department of Geosciences and Geography	
Tekijä – Författare – Author Pinja-Liina Jalkanen			
Tutkielman otsikko – Avhandlingens titel – Title of thesis Patterns of Aggregated Commuting Times in the Helsinki Capital Region			
Koulutusohjelma ja opintosuunta – Utbildningsprogram och studierinriktning – Programme and study track Geography, Geoinformatics			
Tutkielman taso – Avhandlingens nivå – Level of the thesis Master's thesis, 40 credits	Aika – Datum – Date June 2020	Sivumäärä – Sidoantal – Number of pages 64 + 4 appendices	
Tiivistelmä – Referat – Abstract <p>Large-scale transport infrastructure projects change our daily mobility patterns, as they change the geographical accessibility of the places where we spend most of our time, such as our homes and work-places. Thus, there is a clear need for advance evaluation of the effects of those projects. Traditionally, however, the available methods have imposed severe limitations for both measuring accessibility and surveying mobility, and despite modern data collection methods enabled by the ever-present mobile phones, surveying mobility remains challenging due to data accessibility restrictions. Furthermore it would not enable any advance evaluation of mobility changes. However, using a modern accessibility dataset instead of a mobility one does offer a possible answer.</p> <p>In my study, I set out to investigate this possibility. I combined a modern, multimodal and longitudinal accessibility dataset, the Helsinki Region Travel Time Matrix (TTM), with a spatially compatible, census-based longitudinal commuting dataset to evaluate the aggregated journey times in the Helsinki Capital Region (HCR), the area covered by the TTM, and estimated the shares of different transport modes based on a previously published travel survey.</p> <p>Armed with this combined dataset, I assessed the changes in aggregated journey times between the three years that were included in the TTM dataset – 2013, 2015 and 2018 – by statistical district to estimate its usability for these kind of advance mobility evaluations. As a small subset of the commuting dataset was classified by industry, I also assessed regional differences between industries.</p> <p>My results demonstrate that for travel by public transport, the effects of new transport projects are plausibly identifiable in these aggregated patterns, with a number of areas served by several new, large-scale public transport infrastructure projects – the Ring Rail, the trunk bus lane 560 and the Western extension of the metro line – being outliers in the results. For travel by private car and for the industry-level changes, the results are more inconclusive, possibly due to absence of massive projects affecting the road network throughout the dataset timeframe, potential inaccuracies in the source data and limitations of the industry-classified part of the dataset.</p> <p>In conclusion, a modern accessibility dataset such as the TTM can be plausibly used to estimate the mobility effects of large-scale public transport infrastructure projects, although the final accuracy of the results is likely to be heavily dependent of the precision of the original datasets, which should be taken into account when such assessments are made. Further research is clearly needed to assess the effects of diurnal variations in travel times and the effects of more precise transport mode preference data.</p>			
Avainsanat – Nyckelord – Keywords Commuting, accessibility, multimodality, mobility, travel time, journey time, GIS, Helsinki Capital Region			
Säilytyspaikka – Förvaringställe – Where deposited University of Helsinki electronic theses library E-thesis/HELDA			
Muita tietoja – Övriga uppgifter – Additional information			



Tiedekunta – Fakultet – Faculty Matemaattis-luonnontieteellinen		Osasto – Institution – Department Geotieteiden ja maantieteen osasta	
Tekijä – Författare – Author Pinja-Liina Jalkanen			
Tutkielman otsikko – Avhandlings titel – Title of thesis Pääkaupunkiseudun matka-aikakertymien alueelliset muutokset			
Koulutusohjelma ja opintosuunta – Utbildningsprogram och studieriktning – Programme and study track Maantiede, Geoinformatiikka			
Tutkielman taso – Avhandlings nivå – Level of the thesis Pro gradu -tutkielma	Aika – Datum – Date Kesäkuu 2020	Sivumäärä – Sidoantal – Number of pages 64 + 4 liitettä	
Tiivistelmä – Referat – Abstract <p>Suuret joukkoliikenne- ja väylähankkeet vaikuttavat päivittäiseen liikkuvuuteemme muuttaessaan meille arjessamme tärkeimpien paikkojen, kuten kotiemme ja työpaikkojemme maantieteellistä saavutettavuutta. Siksi on olemassa selvä tarve arvioida etukäteen näiden hankkeiden vaikutuksia ihmisten liikkuvuuteen. Perinteisesti käytettävissä olevat menetelmät ovat kuitenkin asettaneet suuria rajoitteita saavutettavuuden ja liikkuvuuden mittaamiselle ja tutkimiselle, ja huolimatta kaikkialle levinneiden matkapuhelinten mahdollistamista moderneista datankeruumenetelmistä liikkuvuuden tutkiminen on edelleen haasteellista datan saatavuusongelmien vuoksi. Puhdas liikkuvuustutkimus ei myöskään mahdollista muutosten vaikutusten arviointia ennakkoon. Modernin saavutettavuusdatan käyttäminen liikkuvuusdatan asemasta on kuitenkin yksi mahdollinen ratkaisu.</p> <p>Tutkimuksessani selvitin kyseisen ratkaisun hyödynnettävyyttä. Yhdistin modernin, multimodaalisen ja monivuotisen saavutettavuusaineiston, Pääkaupunkiseudun matka-aikamatriisin, sen kanssa spatiaalisesti yhteensopivaan yhdyskuntarakennetilaston työmatka-aineistoon arvioidakseni pääkaupunkiseudun matka-aikakertymiä. Eri matkustustapojen osuuksia arvioin aiemmin julkaistun kyselytutkimuksen perusteella.</p> <p>Näin yhdistetyn aineiston avulla tutkin pääkaupunkiseudun työmatkojen matka-aikakertymien muutoksia tilastoalueittain matka-aikamatriisin kolmen eri aineistovuoden – 2013, 2015 ja 2018 – välillä, selvittääkseni tämänkaltaisten ennakoarvioiden käyttökelpoisuutta hankkeiden aiheuttamia liikkuvuusmuutoksia arvioitaessa. Koska pieni osuus työmatka-aineistosta oli toimialaluokiteltua, tutkin myös toimialojen alueellisia eroja.</p> <p>Tutkimukseni tulokset osoittavat, että joukkoliikenteellä tehtyjen matkojen osalta liikennehankkeiden vaikutukset ovat uskottavasti nähtävissä matka-aikakertymien alueellisista muutoksista, useiden suurien joukkoliikennehankkeiden luomien uusien yhteyksien – Kehäradan, runkolinja 560:n ja Länsimetron – palvelemien alueiden erottuessa muusta aineistosta selvästi. Yksityisautoilla tehtyjen matkojen ja toimialakohtaisten erojen osalta tulokset eivät kuitenkaan ole näin selkeitä, mahdollisesti suurten tiehankkeiden puutteen, lähdedatan epätarkkuuksien ja toimialaluokittelun työmatka-aineiston osan rajoitusten vuoksi.</p> <p>Johtopäätökseni on, että modernia, matka-aikamatriisin tapaista saavutettavuusaineistoa voidaan uskotavasti käyttää suurten joukkoliikennehankkeiden liikkuvuusvaikutusten ennakoarviointiin, vaikka tulosten lopullinen tarkkuus vaikuttaakin riippuvan voimakkaasti lähdeaineistojen tarkkuudesta, mikä tulisi ottaa huomioon tämänkaltaisia ennakoarvioita laadittaessa. Matka-aikojen vuorokaudenaikaisen vaihtelun sekä itselläni käytössä ollut aineistoa alueellisesti tarkempien matkustustapamielityksiä koskevien tietojen vaikutusten selvittämiseksi on kuitenkin selkeä jatkotutkimuksen tarve.</p>			
Avainsanat – Nyckelord – Keywords Työmatkat, saavutettavuus, multimodaalisuus, liikkuvuus, matka-aika, GIS, pääkaupunkiseutu			
Säilytyspaikka – Förvaringställe – Where deposited Helsingin yliopiston opinnäytetietokanta E-thesis/HELDA			
Muita tietoja – Övriga uppgifter – Additional information			

Table of Contents

List of figures.....	vi
List of tables.....	viii
List of abbreviations.....	ix
1 Introduction.....	1
2 Background.....	2
2.1 Accessibility.....	3
2.2 Mobility.....	5
2.3 Urbanisation, commuting and residential location-allocation.....	7
2.4 Human perspective and research rationale.....	8
3 Research area and research data.....	8
3.1 Research area.....	8
3.2 Commuting data.....	9
3.3 Travel time data.....	11
3.4 Travel survey data.....	13
3.5 Other ancillary datasets.....	14
4 Data processing and analysis methods.....	14
4.1 Tools.....	14
4.2 Workflow.....	15
4.3 Data preprocessing.....	15
4.3.1 Importing the datasets.....	17
4.3.2 Joining the different datasets to each other in the database.....	18
4.4 Data aggregation.....	19
4.4.1 Initial aggregation of the joined datasets.....	19
4.4.2 Final aggregation results.....	20
4.5 Visual map analysis.....	22
4.5.1 Initial considerations for visualisation.....	22
4.5.2 Visualising all journeys.....	22

4.5.3	<i>Visualising the journeys classified by industry</i>	24
4.6	<i>Statistical analysis</i>	25
4.6.1	<i>Statistical analyses of all journeys</i>	27
4.6.2	<i>Statistical analysis of journeys classified by industry</i>	29
5	<i>Results</i>	29
5.1	<i>Results of the analysis of all journeys</i>	29
5.1.1	<i>Results of the visual map analysis</i>	29
5.1.1.1	<i>Public transport maps</i>	29
5.1.1.2	<i>Private car maps</i>	30
5.1.2	<i>Results of the statistical analysis</i>	35
5.1.2.1	<i>Comparisons between individual journey data–TTM pairs</i>	35
5.1.2.2	<i>Assessing the effects of population changes</i>	39
5.1.2.3	<i>Comparison of comparisons</i>	40
5.2	<i>Results of the analysis of the journeys classified by industry</i>	43
5.2.1	<i>Results of the visual map analysis</i>	43
5.2.2	<i>Results of the statistical analysis</i>	43
6	<i>Discussion</i>	46
6.1	<i>Factors affecting the changes in aggregated times</i>	46
6.1.1	<i>Effects of population changes</i>	46
6.1.2	<i>Plausible effects of transport network changes</i>	47
6.1.3	<i>Regional changes affecting particular industries</i>	49
6.2	<i>Potential sources of errors and missing observations</i>	49
6.2.1	<i>Obvious error sources</i>	49
6.2.2	<i>Potential errors in original measurements</i>	52
6.2.3	<i>Processing errors</i>	53
7	<i>Conclusions</i>	53
8	<i>Acknowledgements</i>	55
9	<i>References</i>	56
	<i>Appendices</i>	65

List of figures

Figure 1. An example of a space-time path as described by Hägerstrand (1970).....	3
Figure 2. The interlinking of accessibility, mobility and land use.....	3
Figure 3. A description of the components of accessibility.....	4
Figure 4. Map of the research area, its municipalities and their travel mode factors. 9	
Figure 5. A 250 × 250 m cell of the SSUF grid cell example.....	10
Figure 6. A door-to-door approach to travel times.....	12
Figure 7. Euclidian route between two endpoints of a commute.....	12
Figure 8. Workflow of my research study.....	16
Figure 9. Bash shell import script for the SSUF MDB files.....	17
Figure 10. Bash shell script for the TTM files.....	17
Figure 11. A 250 × 250 m grid cell demonstrating the five-metre North-South discrepancy between the TTM-supplied and SSUF-supplied grids.....	18
Figure 12. An example SQL query used to create an analysis DB view.....	23
Figure 13. An example of the SQL queries executed by Python scripts.....	25
Figure 14. A map of the changes in aggregated travel times by PT, 2013–2015.....	31
Figure 15. A population-weighted map of the changes in aggregated travel times by PT, 2013–2015.....	31
Figure 16. A map of the changes in aggregated travel times by PT, 2015–2018.....	32
Figure 17. A population-weighted map of the changes in aggregated travel times by PT, 2015–2018.....	32
Figure 18. A map of the changes in aggregated travel times by car, 2013–2015.....	33
Figure 19. A population-weighted map of the changes in aggregated travel times by car, 2013–2015.....	33
Figure 20. A map of the changes in aggregated travel times by car, 2015–2018.....	34
Figure 21. A population-weighted map of the changes in aggregated travel times by car, 2015–2018.....	34
Figure 22. Kernel densities of distributions of the unweighted changes.....	36
Figure 23. Kernel densities of distributions of the population-weighted changes.....	37
Figure 24. Q–Q plots of the individual journey data–TTM pairs.....	37
Figure 25. Density plots of the <i>t</i> -tests of the individual journey data–TTM pairs.....	38

Figure 26. Plots of the residuals of the linear models.....	39
Figure 27. Q-Q plots of the 2013–2015 and 2015–2018 changes compared to each other.....	41
Figure 28. Density plots of the <i>t</i> -test distributions of the 2013–2015 and 2015–2018 changes compared to each other.....	41
Figure 29. An example of industry-classified maps (example industry: Human health and social work).....	42
Figure 30. An example of the Q-Q plots for IC data (example industry: Human health and social work).....	45
Figure 31. An example of the histogram plots for IC data (example industry: Human health and social work).....	45
Figure 32. The route of the trunk bus line No 560.....	48

List of tables

Table 1. All the subdataset classes of the complete SSUF dataset.....	10
Table 2. Industry classification of the SSUF dataset.....	11
Table 3. Travel modes and mean, median and standard travel times according to the Helsinki Region Travel Time Matrices.....	13
Table 4. List of the result tables created by the preparatory R script.....	21
Table 5. An example result of the aggregation of the IC part of the dataset.....	22
Table 6. Numerical changes visualised on the all-journeys maps.....	28
Table 7. Descriptive statistics and <i>t</i> -test results of the individual journey data-TTM pairs.....	36
Table 8. Results of the linear models.....	38
Table 9. Descriptive statistics and <i>t</i> -test results of the 2013–2015 and 2015–2018 changes compared to each other.....	40
Table 10. A list of industries where relative change between industries in at least one district exceeds ten percentage points.....	43
Table 11. Results of the statistical analysis of the IC part of the dataset.....	44
Table 12. Counts and relative shares of all journeys and IC journeys.....	51

List of abbreviations

API	Application Programming Interface.
DB	Database. In the context of this study, usually my research database.
GIS	Geographic information system.
GNSS	Global Navigation Satellite System, such as Galileo or GPS.
HCR	Helsinki Capital Region. An area consisting of the municipalities of Espoo, Helsinki, Kauniainen and Vantaa.
HMA	Helsinki Metropolitan Area. Another known name for the HCR. In this thesis, I use the abbreviation HCR wherever possible.
HSL	Helsinki Region Traffic. An inter-municipal PT authority of the HCR (Finnish: Helsingin Seudun Liikenne).
IC	Industry classification. Used as a reference to the classification by industry (or lack of it) of commuting data records.
ID	An identity number, a number or a group of numbers identifying a single database record.
I/O	Input/output (of data).
OLS	Ordinary least squares, a method for estimating the unknown parameters of linear regression models.
PT	Public Transport. Consists of any type of mass transportation system, e.g. rail (metro/train/tram), bus and ferry connections.
SQL	Structured Query Language. A computer language used for managing data in relational databases.
SSUF	(Monitoring System of) Spatial Structure and Urban Form. The equivalent abbreviation in Finnish is YKR.
TTM	(Helsinki Region) Travel Time Matrix. A matrix of travel times between various locations in the HCR.
WFS	Web Feature Service, a standard protocol for serving georeferenced vectorised map features over the Internet.

1 Introduction

Large-scale transport infrastructure projects change our mobility patterns. The introduction of new transport connections, such as railways, metro lines, bridges, bypass highways and even priority bike paths changes what routes we can use daily and, by our choices of routes and transport modes, affects the amount of time that each of us spends daily commuting between our homes and workplaces; in short, how reachable they ultimately are for us. Put in geographical terms, it means that the changes in the transport network change the *accessibility* of those locations – our ability to reach them (Hansen, 1959; Burns & Golob, 1976; Geurs & van Wee, 2004) – that are the most important ones for our daily lives. For many of us, that means our homes and workplaces. These changes in accessibility in turn affect our personal, realised *mobility* patterns; how do we actually end up moving between places (Kaufmann et al., 2004; Sheller & Urry, 2006; Salonen, 2014; Barbosa et al., 2018).

In the Helsinki Capital Region (HCR) two large-scale public transportation projects have been completed during the past decade and new connections introduced accordingly: the Ring Rail (Kiiskilä et al., 2017) and the Western extension of the metro line (HSL, 2019). These projects have brought a rail connection to places where previously taking a bus was the only public transit option. At the same time there has been a number of smaller scale changes and upgrades to HCR road networks, affecting those who commute by a private car.

In geographical terms, this means that the accessibility of various locations in the HCR is, to some extent, very likely to have changed over time; thus, the mobility patterns of the people are likely to have changed as well. Accessibility and land use are known to be closely linked (Hansen, 1959; Kenworthy & Laube, 1996; Priemus et al., 2001; Wegener & Fuerst, 2004; Bertolini et al., 2005), and the HCR – generally understood to be consisting of the municipal areas of the cities of Espoo, Helsinki, Kauniainen and Vantaa – is also known to be a relatively sparsely built urban area, if compared to various big cities around the world (City of Helsinki, 2018; United Nations, 2014; University of Oxford, 2014), which in turn makes people more depend-

ent of various motorised forms of transportation (Kenworthy & Laube, 1996; Priemus et al., 2001; Wegener & Fuerst, 2004; Bertolini et al., 2005).

As the population of the HCR has grown steadily over the past decade (Aluesarjat, 2020), I hypothesise that the accessibility changes have affected people's mobility patterns enough to cause observable changes in aggregated commuting times. Because a large part of the population commutes daily to work and back again, these small changes in personal mobility patterns of the individuals might result in significant changes in time spent travelling as aggregated to the level of the total population. As observing mobility patterns of the population is not a trivial task and requires utilising either data with a limited sample size, such as travel diaries, demographically biased datasets such as social media data (Ruths & Pfeffer, 2014; Heikinheimo et al., 2017) or data that is hard to obtain due to privacy and availability reasons such as mobile phone data (Ahas et al., 2008; Barbosa et al., 2018), I have combined existing travel survey information (Brandt et al., 2019) with an existing accessibility dataset (Tenkanen & Toivonen, 2020) and an existing national census-based dataset about commutes between homes and workspaces (SYKE, 2016b) to estimate the realised mobility patterns; further I have analysed the aggregated changes in estimated travel times to answer the following research questions:

1. Have the recent changes in the transport infrastructure had any significant effect to aggregated commuting times in any particular subregion of the HCR?
2. Are there any significant regional commuting time changes between industries?

2 Background

The basic concept of the movement of an individual in time and space – mobility – has perhaps best been illustrated by Hägerstrand (1970) with his time geography and space-time paths (figure 1). Planners of land use and transport, however, have a clear need to estimate those movements in advance, which has given rise to accessibility research. The policy decisions taken during the planning stages in turn change the established mobility patterns, having a human cost. Thus, accessibility,

mobility and land use are (figure 2) in a complex fashion, and understanding the terms thoroughly is a necessity for understanding the challenges related to policy changes.

2.1 Accessibility

Accessibility, in the geographical sense of the term, has historically been a somewhat contested concept, with a variety of alternative definitions (Gould, 1969, p. 64; Ingram, 1971; Geurs & van Wee, 2004). Hansen's

(1959) way of defining it – *'the potential of opportunities for interaction'*, i.e. the potential reachability – is one of the most common. Ingram (1971) defines it as a *'degree of interconnexion with all other points of the same surface'* – an integral – and defines that integral as the sum of the relative accessibilities of all the point pairs of that surface. Burns & Golob (1976) described¹ it as *'the ease with which any land-use activity can be reached from a location using a particular transportation system'*. However, in modern research the more precise definition of Geurs & van

Wee (2004), which builds on the aforementioned previous definitions, is usually ap-

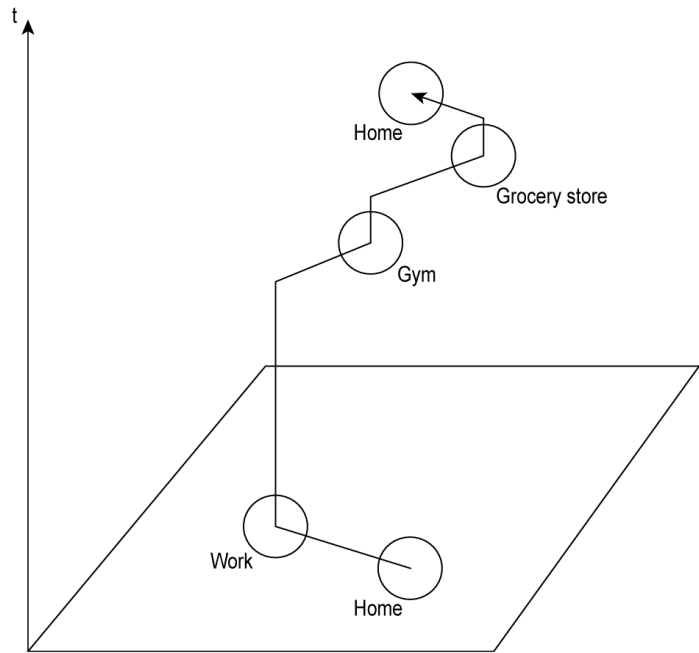


Figure 1. An example of a space-time path as described by Hägerstrand (1970), showing the daily movements of a single individual in space and time.

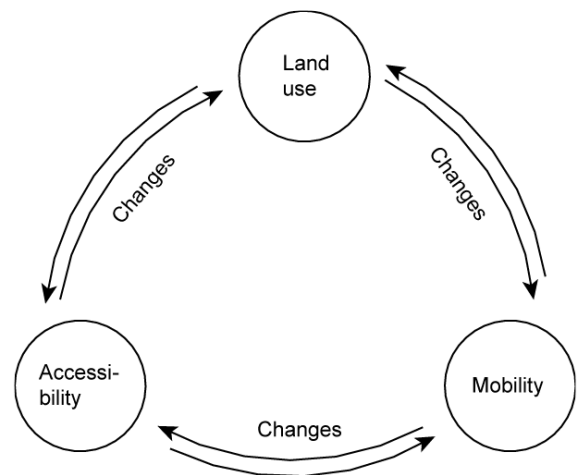


Figure 2. Accessibility, mobility and land use are interlinked in a complex fashion. Changes in one affect the other two, as policy decisions are made based on people's behaviour, and people's behavior is shaped by their consequences.

¹ A citation frequently – but erroneously – attributed to Dalvi & Martin (1976) instead, e.g. by Geurs & van Wee (2004).

plied: 'the extent to which land-use and transport systems enable (groups of) individuals to reach activities or destinations by means of a (combination of) transport mode(s)'. The latter definition is particularly useful, because it stresses the importance of land use and transport while defining geographical accessibility.

In addition, Geurs & van Wee (2004) further clarified their definition by identifying four different components of accessibility, and classified the various measures used to evaluate it. While transport and land use are widely known to be closely linked with accessibility (Hansen, 1959; Kenworthy & Laube, 1996; Priemus et al., 2001; Wegener & Fuerst, 2004; Bertolini et al., 2005) and can thus be considered to be components of it, Geurs & van Wee (2004) also identified temporality and the needs and abilities of an individual as additional components. As for the measures, they classified them to four perspectives: infrastructure-based, location-based, person-

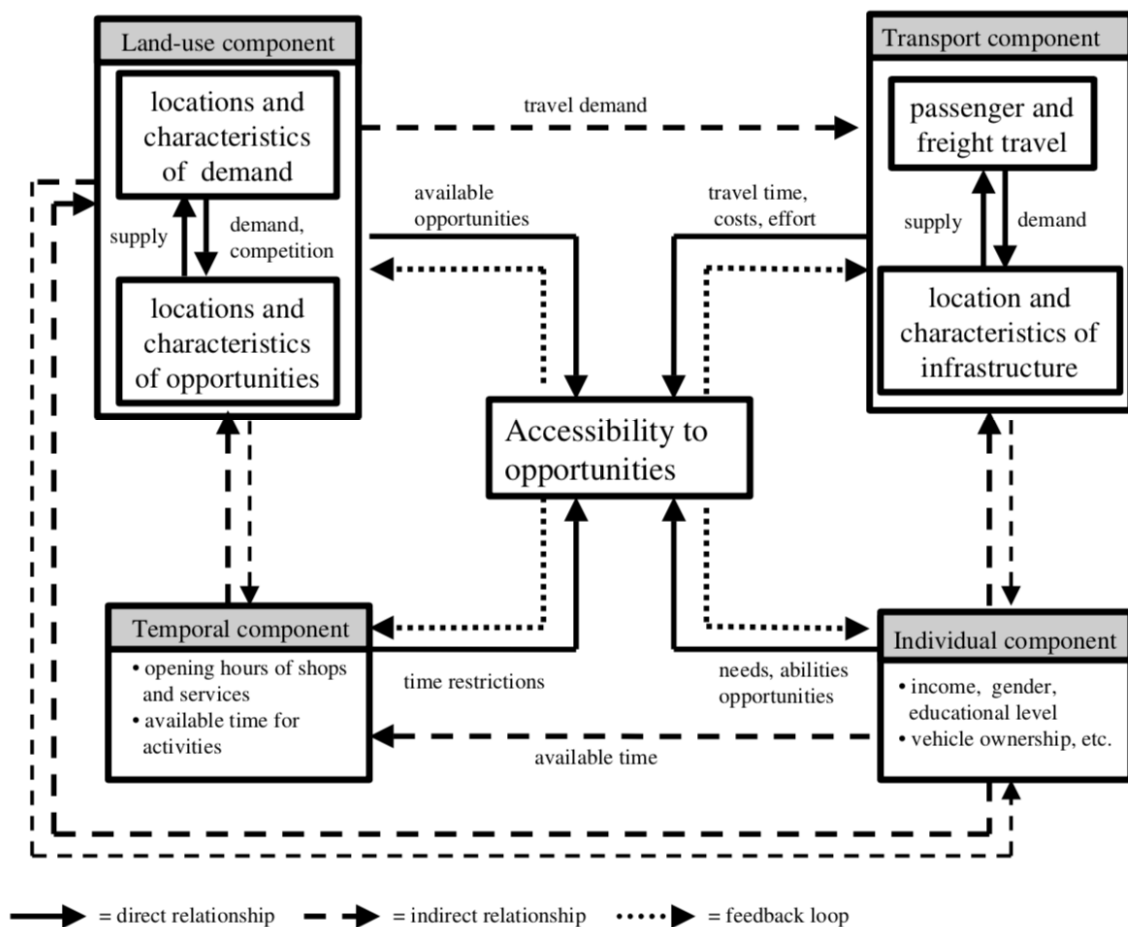


Figure 3. A description of the components of accessibility. Adopted from Geurs & van Wee (2004).

based and utility-based measures. In my research, I refer to accessibility and its components and measures using the definitions of Geurs & van Wee (2004).

As transport mode is a component of accessibility, it may consist of a single mode of transport, or a combination of them. When referring to single journeys utilising a combination of transport modes, I use the term *multimodal* transport, as used by e.g. Salonen & Toivonen (2013), Tenkanen et al. (2016), and Barbosa et al. (2018); an example would be walking to a bus stop, boarding a bus, and then walking to the intended destination from an another bus stop after disembarking the bus. While multimodality has traditionally not been given a careful consideration in literature (Salonen & Toivonen, 2013; Tenkanen et al., 2016), most vehicular transport is clearly multimodal by nature, unless the vehicle used for the travel can be parked right at the door of the source and destination points, so that no significant amount of walking is involved. In a densely built city environment, this is rarely possible.

Of the measures of accessibility, travel time is the most common one, as it can be used as a measure of both the transport and the temporal component, and it can be utilised from a location-based as well as from a person-based perspective. Geurs & Östh (2016) and Tenkanen & Toivonen (2020) note that the technological advances have greatly increased the precision and complexity of accessibility models that utilise location-based measures, giving raise to models that aim to enable very realistic comparisons of travel times between various travel modes and hours of day.

2.2 Mobility

Whereas accessibility refers to potentiality as an enabling factor, the term most commonly used to refer to the actual, realised movement patterns of people between various points in space – on the surface of the planet Earth – is mobility (Kaufmann et al., 2004; Sheller & Urry, 2006; Salonen, 2014; Barbosa et al., 2018). As with accessibility, non-geographical definitions of the term exist (Kaufmann et al., 2004); to distinguish between them, the term spatial mobility may be used. In my research, I use the term mobility to refer exclusively to spatial mobility. Ultimately, this definition of the term builds to Hägerstrand's (1970) time geography (see figure 1).

The most traditional method of surveying mobility patterns has been census (Palmer et al., 2013; Barbosa et al., 2018) and taxation (Barbosa et al., 2018) data, but due to their periodical nature they are only useful while surveying national migration patterns and not day-to-day mobility (Palmer et al., 2013; Barbosa et al., 2018). As a more creative traditional method, the movements of the banknotes in circulation have been tracked as well (Barbosa et al., 2018). Local migration and day-to-day mobility patterns have traditionally been tracked by travel surveys, in more recent times occasionally augmented by GNSS tracking, but the small scale of those surveys severely limit their usefulness in data validation and in creating a dynamic picture of mobility patterns (Palmer et al., 2013; Barbosa et al., 2018).

Arguably the most important new data source alleviating the aforementioned limitations is mobile phone data (e.g. Ahas et al., 2008; Candia et al., 2008; Palmer et al., 2013; Barbosa et al., 2018), which has even been called ‘*game-changing*’ (Barbosa et al., 2018). Analysing mobile phone data has helped to demonstrate that – contrary to previous studies – the movements of individuals ‘*show a high degree of temporal and spatial regularity*’ (González et al., 2008), and movements of phones have been described as being ‘*proxies for people*’ (Järv et al., 2017). The tracked phones themselves do not have to be actively participating in the tracking process, as their locations can be estimated by tracking their connections to cellular base stations, whose locations are known. (Ahas et al., 2008; Candia et al., 2008; Palmer et al., 2013; Barbosa et al., 2018). It is also possible to further improve the precision of such interpolation by utilising ancillary data sources. (Järv et al., 2017).

While mobile phone data may truly be a remarkable improvement over the traditional methods, utilising it still brings a share of its own problems. Privacy concerns over the research subjects is the most commonly mentioned problem (Ahas et al., 2008; Candia et al., 2008; Palmer et al., 2013; Barbosa et al., 2018), but there are also problems with accessing the data at all and sharing even anonymised versions of such datasets, due to trade secrets involved (Ahas et al., 2008; Barbosa et al., 2018). In addition, utilising mobile phone data is not unproblematic even quality-wise, due to location accuracy and temporal frequency problems especially on rural

areas (Järv et al., 2014; Barbosa et al., 2018). The geographical extent of any research area is usually also limited to the extent of the geographical area of the mobile operator's network (Hawelka et al., 2014), commonly to the area of a single nation. In some cases, social media data has been utilised instead (e.g. Hawelka et al., 2014; Tenkanen et al., 2017; Heikinheimo et al., 2017), but the results are affected by popularity of the chosen social media platform, proprietary data collection methods of the platforms and demographical biases of their user-bases, limiting the usability of such data (Ruths & Pfeffer, 2014; Heikinheimo et al., 2017; Tenkanen et al., 2016; Barbosa et al., 2018).

2.3 Urbanisation, commuting and residential location-allocation

Traditional economics models have assumed that choice of the residential location is based on cost minimisation (e.g. Mills, 1972; Muth, 1969). As Hamilton (1982) however found out and Small & Song (1992) confirmed, those models may significantly underestimate the actual length of the commutes, which has led to further criticism of the models by e.g. Giuliano (1991) and Giuliano & Small (1993). The former identified a number of factors more significant than commuting cost in the residential location-allocation decisions of individual households, concluding that jobs-housing balance is not a transportation issue and recommending transportation-centric solutions as solutions for transportation-specific problems instead of land-use policies aimed at changing jobs-housing balance (Giuliano, 1991); the latter confirmed that land-use policies aimed at changing the jobs-housing balance only have a minor effect on the actual commuting times.

While the Giuliano & Small's (1993) answer to their title question '*Is the Journey to Work Explained by Urban Structure?*' was an unambiguous no, Levinson (1998) counters this from the accessibility viewpoint, as he found out that '*17–38% of the variation in travel time to work of individuals can be explained by attributes of urban structure*'. He also notes that according to previous literature – specifically Gordon et al. (1991) and Levinson & Kumar (1994) – commuting times have remained stable or even shortened. He attributes this to suburbanisation of jobs and the consequent use of faster suburban road network and concludes – in contradiction to Gi-

uliano & Small's (1993) findings – that *'policies favoring a properly defined jobs/housing balance will, at the margins, reduce commuting duration'* (Levinson, 1998). Levinson & Wu (2005) confirmed his findings, and further concluded that *'commuting time clearly depends on metropolitan spatial structure'* (Levinson & Wu, 2005), further stressing the aforementioned (see p. 4) link between land use, transport and accessibility. While recent research has also emphasised the importance of equity questions (Neutens et al., 2010; Delafontaine et al., 2011; Neutens, 2015; Cui et al., 2019), the basic interdependency between accessibility, transport and land use still applies.

2.4 Human perspective and research rationale

Throughout recent decades, the average commuting distances have generally been increasing in several European countries (Lyons & Chatterjee, 2008; Helminen & Ristimäki, 2007; Sandow & Westin, 2010). The average length of a commute has increased steadily e.g. in Finland (Lintunen, 2000; Helminen et al., 2003; SYKE, 2016a), in Sweden (SOU, 2003; Hedberg, 2005) and in the United Kingdom (Lyons & Chatterjee, 2008). At the same time, longer commuting times are known to induce stress (Kluger, 1998; Evans et al., 2002).

Information about the total time our societies use for commuting is thus valuable, yet obtaining suitable mobility datasets remains a challenge, and the need remains to analyse the effects of land use and transportation policy changes before the mobility changes induced by them are measurable at all. In my research, I aim to tackle these challenges.

3 Research area and research data

3.1 Research area

I chose the Helsinki Capital Region (HCR) for my research area. The HCR, also known as Helsinki Metropolitan Area, consists of the areas of the municipalities of Espoo, Helsinki, Kauniainen and Vantaa. (The HCR should not be confused with the larger area known as Helsinki Region or Greater Helsinki Region.) This choice was partly induced by the spatial limitations of my ancillary datasets, but it is also a well-

known moderately large urban area, with a total population of 1,169,455 persons (Aluesarjat, 2020) at the end of the year 2018.

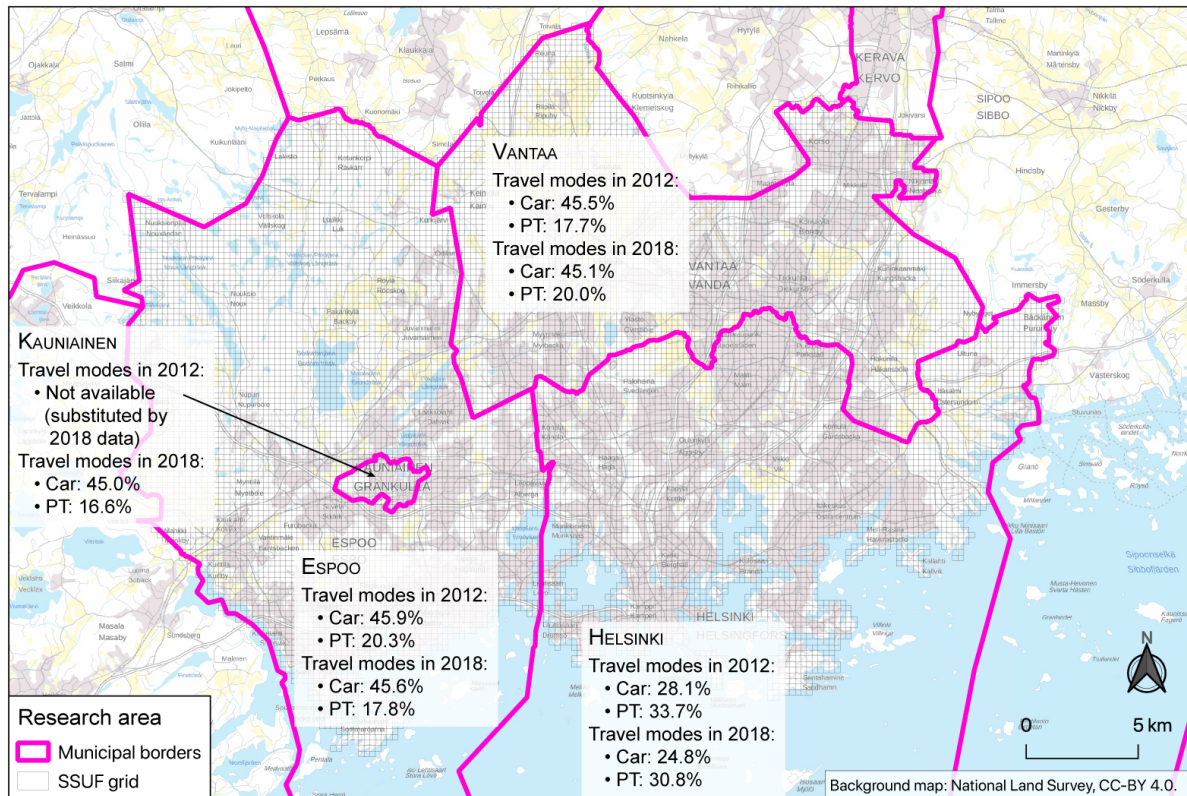


Figure 4. Map of the research area, depicting the municipalities of the area and the travel mode factors of motorised transport in each municipality.

3.2 Commuting data

I use a national census-based, grid-aligned spatial dataset, the Monitoring System of Spatial Structure and Urban Form (SSUF) dataset (SYKE, 2016b), as my main commuting dataset. This dataset includes various variables about the structure of population and land use, including information about individual commutes. A list of all the information included in the dataset is listed in table 1. In my research, I utilise this commuting information. The commuting part of the dataset includes information about the start and end points as well as the euclidian distance of each commute. The data has been anonymised to some extent by aggregating all spatial data to 250×250 m grid cells (figure 5). Henceforth, I refer to this dataset as the SSUF data. This dataset is not publicly available, but all Finnish higher education institutes have access to it for research purposes.

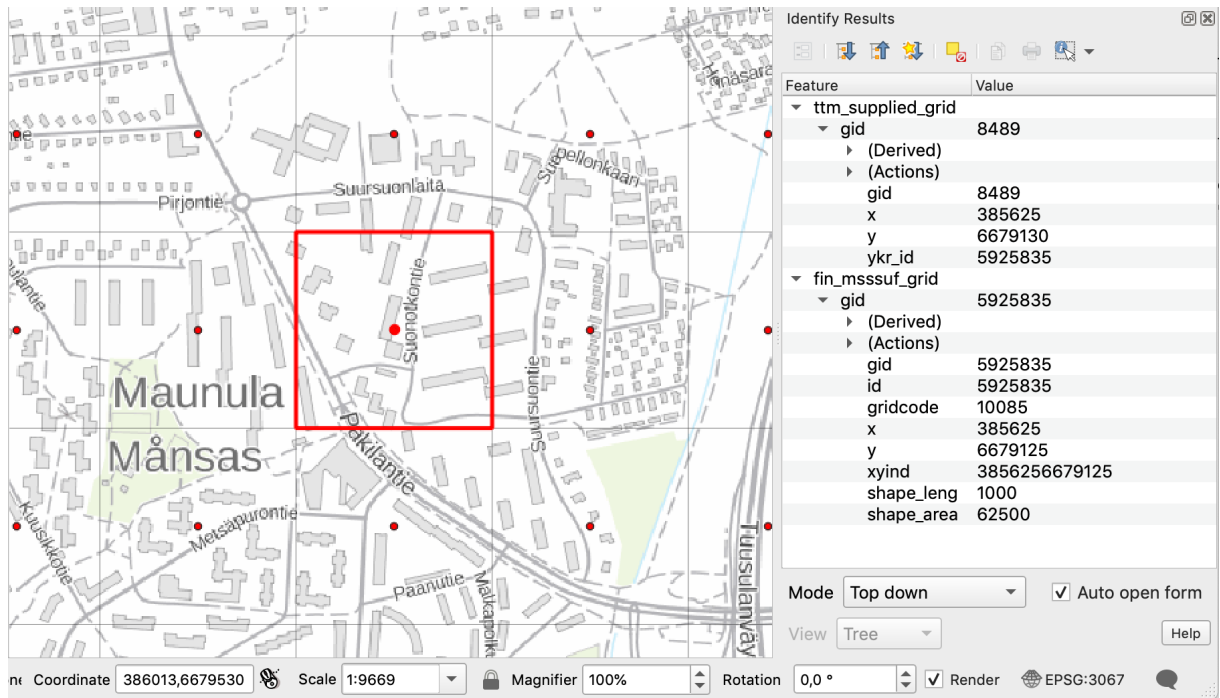


Figure 5. A 250 × 250 m cell of the SSUF grid located in the Maunula district of the HCR.

The commuting dataset includes industry classification (IC) data (table 2) for the commutes between those grid cell pairs that have ten or more single commutes between them; this limitation is for privacy reasons. The dataset is longitudinal including data for several (but not all) years between 1998 and 2016; however, there is a longitudinal break in the industry classification, so that data from 2007 onwards is

Table 1. All the subdataset classes of the complete SSUF dataset (SYKE, 2019).

Class	Description
Habitations	Information about all dwellings of 7m ² or larger with cooking facilities and a dedicated private entrance..
Households and car ownership	Information about persons and their cars who are permanently registered in a single habitation.
Wholesale and retail locations	Information about wholesale and retail stores.
Summer cottages	Information about cottages that are not registered for permanent habitation.
Land use and land cover	CORINE land cover classification of the SSUF grid cells
Buildings	Information about all buildings
Commutes	Information about commutes (the primary source dataset of this research).
Workplaces	Persons working in an SSUF grid cell. Classified by industry.
Labour force	Persons dwelling in an area that participate in the labour force. Classified by industry.
Population	Population statistics.

classified differently. I utilise only the part of the dataset that uses the newer classification; with this restriction, the journey data for the years 2007, 2009, 2010, 2012, 2014, 2015 and 2016 is included. The data collection day for each of the years is the last day of the year. The individual commutes were included in the dataset thrice: once for the total and once for each binary gender, as registered in the census-based population register data. (SYKE, 2019)

3.3 Travel time data

While the SSUF data enables straightforward calculation of the euclidian distances of the commutes (figure 7), human beings are rather obviously not capable of unsupported flight, making euclidian distances an extremely poor estimate of our realised mobility. Consequently, they are known to be a poor measure of accessibility as well (Salonen et al., 2012), thus requiring an ancillary dataset to complement the SSUF data.

Given the challenges in obtaining suitable datasets for mobility research (see p. 6) and the need to be able to analyse potential mobility changes induced by land use and transportation policies before they are actually enacted, I aimed to use a modern accessibility dataset created by the Accessibility Research Group at the Univer-

Table 2. Industry classification of the SSUF dataset.

- Agriculture, forestry and fishing
- Mining and quarrying
- Manufacturing
- Electricity, gas, steam and air conditioning supply
- Water supply; sewerage, waste management and remediation activities
- Construction
- Wholesale and retail trade; repair of motor vehicles and motorcycles
- Transportation and storage
- Accommodation and food service activities
- Information and communication
- Financial and insurance activities
- Real estate activities
- Professional, scientific and technical activities (speciality professions, e.g. law, architecture, accounting, marketing, research, vets)
- Administrative and support service activities (e.g. leasing, office services, employment agencies, security, travel agencies)
- Public administration and defence; compulsory social security
- Education
- Human health and social work activities
- Arts, entertainment and recreation
- Other service activities (e.g. NGOs, appliance repair services, laundry services, spas)
- Activities of households as employers; undifferentiated goods- and services-producing (activities of households for their own use)
- Activities of extraterritorial organisations and bodies
- Unknown industry

sity of Helsinki (Tenkanen & Toivonen, 2020) in lieu of a mobility dataset to estimate the mobility changes induced by changes in land use; more specifically, by changes in the transport networks.

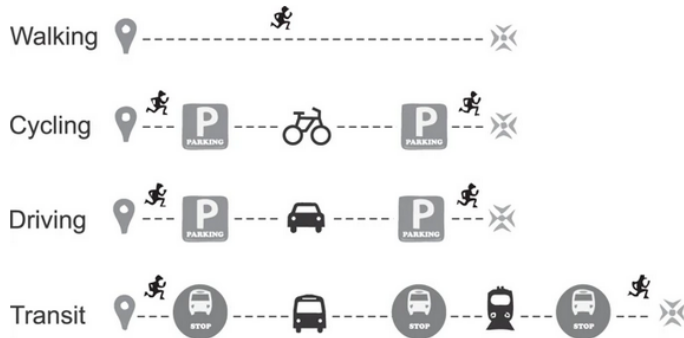


Figure 6. A door-to-door approach to travel times. Adapted from Tenkanen & Toivonen (2020).

The dataset – known as the Helsinki Region Travel Time Matrices (TTMs) – is spatially limited to the HCR, and its limited spatial extent is the primary reason I chose the HCR as my research area (see chapter 3.1). In principle, this dataset includes door-to-door travel (figure 6) time information for three distinct years – 2013, 2015 and 2018 – between any two 250×250 m grid cells in the HCR by private car, public transport (PT), bicycle and walk, including car and bike parking times and PT waiting and changing times, both during the rush-hours and during midday. In practice, however, the dimensions of the dataset matrix are not complete: the bicycle data is not longitudinal at all, as it is only included in the 2018 TTM; moreover, the 2013 TTM lacks the rush-hour data, restricting the multitemporal dimension to the 2015 and 2018 TTMs (see table 3). Thus, the full longitudinal extent of the dataset is only obtainable for midday car, PT and walk data. (Tenkanen & Toivonen, 2020)

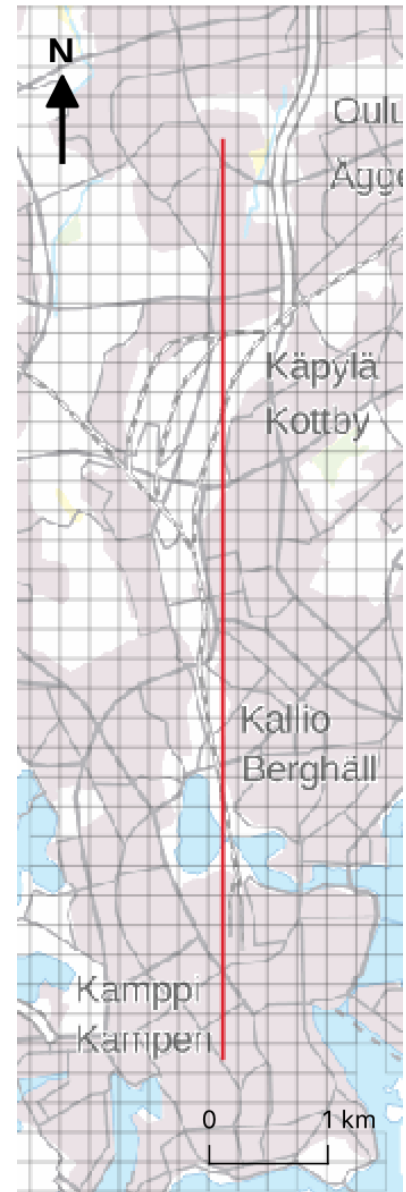


Figure 7. Euclidian route between two endpoints of a commute.

Table 3. Travel modes and mean, median and standard travel times according to the Helsinki Region Travel Time Matrices. Adopted from Tenkanen & Toivonen (2020).

Travel mode	Mean t	Median t	Std t
Public transport			
<i>Rush hour</i>			
2018	84.3 minutes	81 minutes	34.4 minutes
2015	82.9 minutes	79 minutes	34.2 minutes
2013	—	—	—
<i>Midday</i>			
2018	86.2 minutes	83 minutes	34.4 minutes
2015	84.3 minutes	81 minutes	33.6 minutes
2013	81.3 minutes	78 minutes	32.4 minutes
Private car			
<i>Rush hour</i>			
2018	44.5 minutes	44 minutes	16.5 minutes
2015	41.9 minutes	41 minutes	14.2 minutes
2013	—	—	—
<i>Midday</i>			
2018	38.9 minutes	38 minutes	14.3 minutes
2015	37.2 minutes	36 minutes	12.5 minutes
2013	37.7 minutes	37 minutes	12.6 minutes
<i>Free flow</i>			
2018	25.6 minutes	25 minutes	8.9 minutes
Cycling 2018			
<i>Fast biker</i>	59.1 minutes	57 minutes	27.6 minutes
<i>Slow biker</i>	93.2 minutes	89 minutes	43.6 minutes
<i>Walking</i>			
2018	281.9 minutes	271 minutes	133.2 minutes
2015	282.0 minutes	271 minutes	133.3 minutes
2013	281.0 minutes	269 minutes	135.6 minutes

This dataset is publicly available and utilises exactly the same grid structure as the SSUF dataset, being thus not only spatially but on the data record key level fully compatible between the SSUF commuting dataset. The PT timetable data collection days for the TTMs were 08/04/13, 28/09/15 and 29/01/18 for the 2013, 2015 and 2018 matrices, respectively (Accessibility Research Group & University of Helsinki, 2016a, 2016b, 2018). Henceforth, I refer to this dataset as the TTM data, and to the individual TTM years as the TTM 2013, TTM 2015 and TTM 2018, respectively.

3.4 Travel survey data

To utilise a multimodal accessibility dataset instead of a mobility one, having some information about the most likely travel mode of each individual is a prerequisite, as it is very unlikely that we would all use the same form of transport. For this informa-

tion, I utilised a travel survey report published by the Helsinki Region Transport (HSL). The latest such survey, conducted in 2018, was done by drawing a random stratified sample of 38,720 resident individuals from the population register. The report included the results of the previous survey, conducted in 2012, for comparison. (Brandt et al., 2019)

I used this travel survey to estimate the share of the travel mode for the commutes conducted from each grid cell. As the HCR consists of four different municipalities – Espoo, Helsinki, Kauniainen and Vantaa – and this was the highest level of granularity of the travel mode share data in the survey (Brandt et al., 2019, pp. 59–69), I also did my estimates based on the municipality that each grid cell was located in (see figure 4).

3.5 Other ancillary datasets

To spatially divide the HCR for further subregional analysis, I used a WFS map layer of HCR statistical districts provided by the City of Helsinki (2020). To weigh data by population changes, I also used a population dataset jointly maintained by the cities of the HCR and Statistics Finland (2019).

4 Data processing and analysis methods

4.1 Tools

For data processing, I utilised the computational services of CSC – IT Center for Science Ltd, a state-owned company providing computational services for universities and other higher education institutes in Finland. As a platform, I used an I/O-optimised virtual Linux server running the Ubuntu 18.04 operating system; for the data backend and for primary preprocessing, I used a PostgreSQL 10 database (DB) with PostGIS 2.4 spatial extensions. For aggregating and analysing the data further, I used Python and R scripts.

For visually analysing the data, I used QGIS 3.4 as a GIS application to create maps manually for the purpose of analysing the complete dataset without consideration for the industry classification. For the IC part of the dataset, I created per-in-

dustry maps programmatically via a Python script utilising Bokeh visualisation library. For statistical analysis, I used R to analyse the complete dataset and Python for the IC part. In R, I used the default statistical package Stats for all statistical tests, except for Levene tests, for which I used the MatrixTests package. In Python, I used the Scipy package for all tests. All of the aforementioned scripts are publicly available from my GitHub repositories (see appendix 1). For coding and writing, I used a laptop.

4.2 Workflow

The workflow of my research study is graphically described in the figure 8. It can be generally divided to organising the study, reviewing the relevant literature, preparing the data, analysing the complete dataset and the IC part of it separately and answering the research questions.

4.3 Data preprocessing

In general, I chose to eliminate any parts of the datasets not required for the analysis process as late in my workflow as possible, so as to make sure that I would have the widest possible dataset available later, should I then discover a need to widen the scope of my research in some aspect. However, this approach somewhat slowed down the data import process, especially with respect to importing the SSUF data.

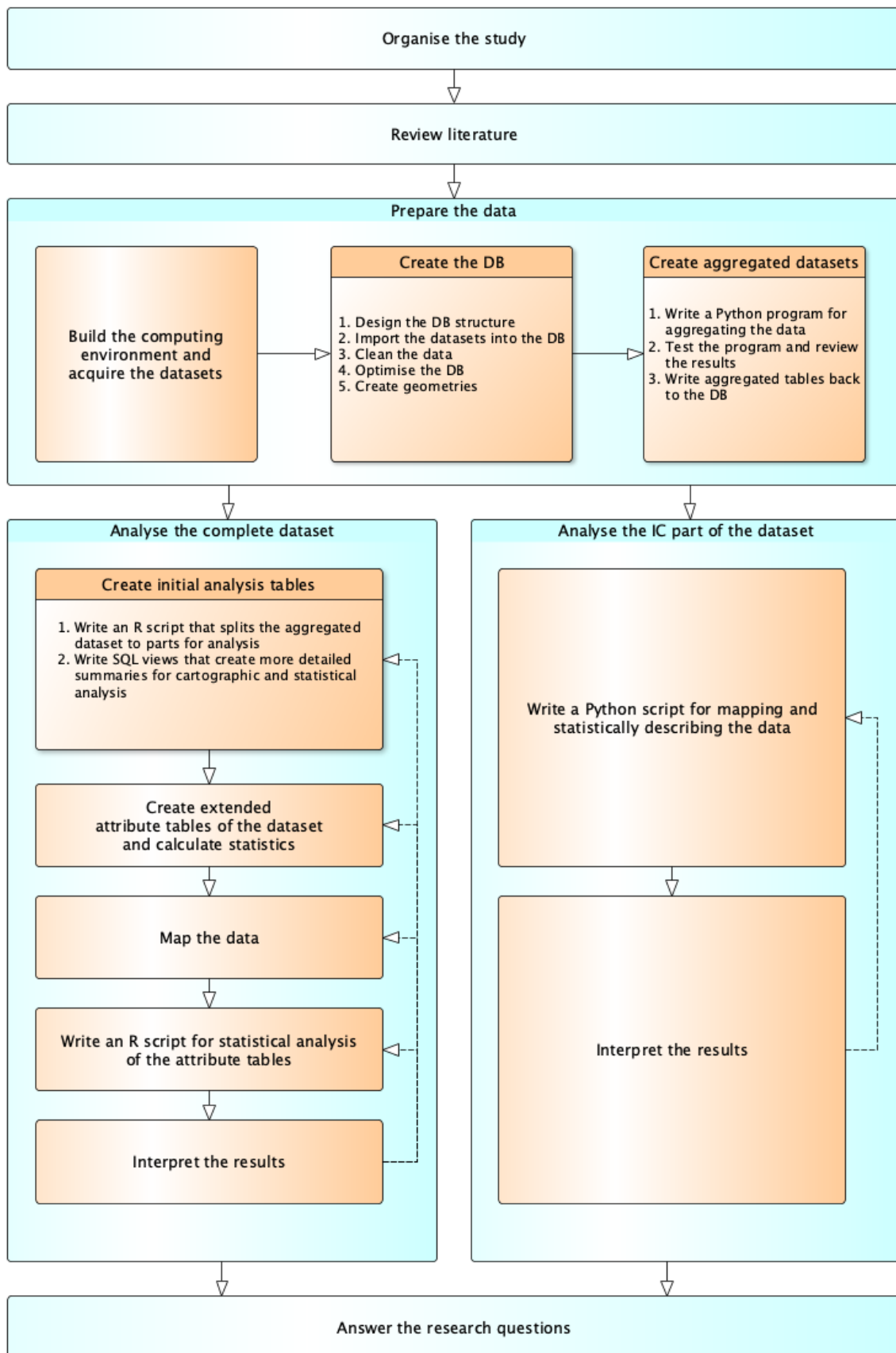


Figure 8. The workflow of my research study.

4.3.1 Importing the datasets

To import the SSUF data, which was originally provided as a set of year-specific Microsoft Access Database (MDB) files with a shared table structure, I created a DB table for them, and then used an open-source tool called MDB-export and a one-liner Bash shell script (figure 9) that imported each of the files to the same DB table; I then cleaned the DB further by converting all the -1 values to real null (empty) values. I did this for the whole national dataset and for all the years from the year 2007 onwards (that used the newer industry classification); importing the data took over 24 hours, and the full table scans required by the null conversions took nearly two days. I also created geometry columns for the start and end points of the commutes and the euclidian lines between the two, in case that those were later needed, cast some text-type field to integers, defined a primary key (id, year) and created some indexes to improve query performance.

```
for i in `find . -name '*2008.mdb' |sort |xargs`; do echo $i;
mdb-export -Q -I postgres $i `echo -n $i |sed 's/^\./[0-9]\{4\}\.mdb/\1/'` |sed 's/\T06_tma_e_[0-9]\{2\}_TOL2008/T06_tma_e_TOL2008/' |psql -q tt; done
```

Figure 9. A one-liner Bash shell script that I used to import the year-specific commute MDB files to my research DB. The files include over 26 million million records in total; the script took over 24 hours to finish.

```
for i in `find ./ -name 'travel_times_to*.txt'`; do echo $i;
cat $i |psql tt -q -c "COPY ttm2013 FROM STDIN (FORMAT csv,
DELIMITER ';', NULL '-1', HEADER)"; done
```

Figure 10. A one-liner Bash shell script to import the TTMs to my research DB.

To import the TTMs, I simply created DB tables for each of the TTMs, and then used a Bash one-liner script to find all the loose TTM files from the directory structure and to write them to the DB (figure 10). I later on discovered that the 2013 and 2015 TTMs included rows having an empty target cell (to_id), and in a hindsight, it would have made sense to drop those rows at the point of the import; as I did not do that, I deleted those rows afterwards, validating the deletion result by searching all such rows in the source data and checking that their counts matched the deletion counts (436,561 and 2,399,761 rows in the 2013 and 2015 TTMs respectively; the 2018

TTM did not include such rows). I then defined primary keys for the tables (from_id, to_id).

4.3.2 Joining the different datasets to each other in the database

To join the SSUF and TTM datasets, I first joined the helper grids supplied with both datasets; those grids had a shared ID number, which enabled direct SQL joins using the shared ID, without utilising spatial functions. This ability to join the SSUF and TTM datasets directly with each other also made sure that the observed, otherwise unexplained five-metre North-South discrepancy between the two grids (figure 11) did not warrant further investigation. Utilising this join, I created a new helper grid of my research area, using a so-called anti-join to eliminate all the rows from the SSUF-supplied national grid that did not have an equivalent ID in the TTM helper grid. This way, I got a grid, framed to my research area, that included both the grid cell numbers of the TTMs and the x/y-coordinates required to join the SSUF data, without the unexplained five-metre discrepancy between the original helper grids.

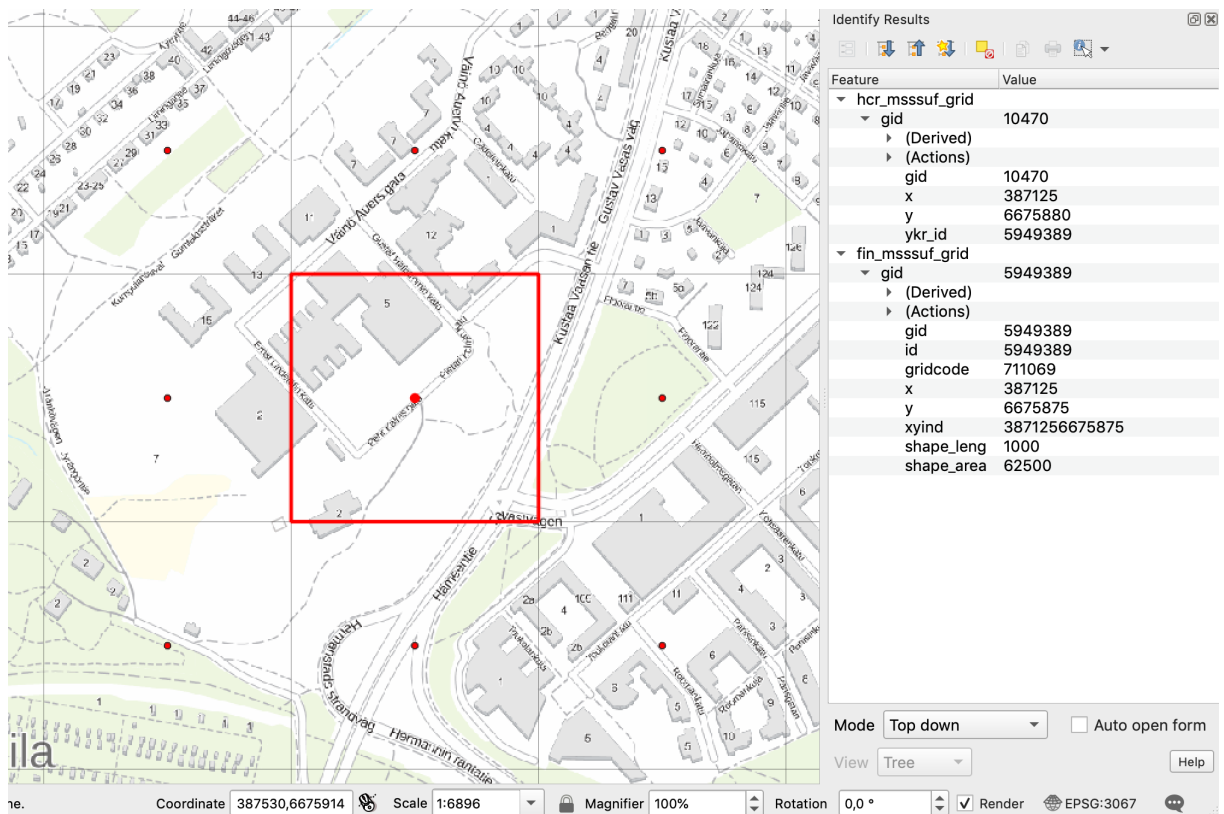


Figure 11. A 250 × 250 m grid cell from the Kumpula Campus area of the University of Helsinki, demonstrating the five-metre North-South discrepancy – the y coordinate value between the grids supplied with the TTM (top) and SSUF (bottom) datasets.

With the help of this new grid, I then created a copy of the original SSUF dataset that included only the journeys that had their origin or destination in my research area. This was a looser definition than what was ultimately necessary for my research, as I still wanted to take any steps narrowing the dataset as late in my workflow as possible. I then joined all the SSUF data records with each of the TTMs, creating for each TTM a new table that included fields for the equivalent SSUF id-year pairs (essentially a foreign key, although I did not technically define them as such) enabling further querying of the commuting data of any given year and binary gender against the travel time information of that specific TTM in a more simplified fashion. For the specific SQL queries creating these tables, see appendix 1 for my code repositories.

For the purposes of the subregional analysis, I added fields for start and end point of the commutes to the SSUF dataset and updated the fields by spatially joining statistical district layer with the SSUF data.

I did not utilise the travel survey data during preprocessing.

4.4 Data aggregation

4.4.1 *Initial aggregation of the joined datasets*

For the initial aggregation of the joined datasets, I wrote a Python script that joins the SSUF data with the TTMs, counts the amount of journeys, aggregates – based on the TTM data – the journey times and distances and finally writes the aggregated results back to the DB as new tables. (For references to this code, see appendix 1). I used this script to create the following four DB tables for further analysis of the data:

- All journeys: counts and aggregated times and distances of all commute journeys within the whole HCR for each journey data year–TTM pair.
- All journeys by region: counts and aggregated times and distances of commute journeys by statistical subregion.
- IC journeys: counts and aggregated times and distances of those commute journeys, for which industry classification was available, within the whole HCR for each TTM–journey year pair, aggregated by industry.

- IC journeys by region: counts and aggregated times and distances of those commute journeys, for which industry classification was available, within the whole HCR for each TTM–journey year pair, aggregated both by industry and by statistical subregion.

At this point of the workflow, I also hardcoded the gender variable to the Python script, thus narrowing the dataset further. As a result, all of the aforementioned tables are gender-neutral. Furthermore, I also omitted any journey years older than 2012 from the results.

4.4.2 *Final aggregation results*

For the actual analysis, I narrowed the datasets further by only selecting those journey data year–TTM pairs that temporarily best matched each other (see p. 11 for the journey data acquisition days and p. 13 for the TTM data acquisition days); thus, I used the 2012 journey data with the 2013 TTM, the 2015 journey data with the 2015 TTM, and the 2016 journey data with the 2018 TTM. For the 2012–2013 journey data–TTM pair and the 2015–2015 pair the data acquisition dates were within about four and about three months of each other, respectively, but for the 2016–2018 pair the dates were nearly thirteen months apart, due to unavailability of any temporarily better matching datasets. In addition, due to aforementioned longitudinal limitations of the TTM dataset (see chapter 3.3), only midday times were considered.

To further analyse the data, I wrote an R script that read in each of the initial four DB tables, narrowed the dataset further to the aforementioned journey data–TTM pairs, integrated the travel survey data to enable calculation of relative shares of the data by travel mode, region or industry and then finally wrote 26 separate aggregated result tables back to the database. List of these result tables with their descriptions is included in table 4. An example result of the aggregation of the IC part of the dataset is presented in table 5.

Table 4. List of the result tables created by the preparatory R script that split the aggregated tables to parts and further prepared them for the actual analysis. The tables that I actually used in the analysis are marked with an asterisk().*

Table name in the DB	Table description
res_agg_j_all_car	All car data as aggregated
res_agg_j_all_car_freq	All car data, relative share of each data year pair of the total for all years
res_agg_j_all_pt	All PT data as aggregated
res_agg_j_all_pt_freq	All PT data, relative share of each data year pair of the total for all years
res_agg_j_all_reg_car	All car data as aggregated, by region
res_agg_j_all_reg_car_mun*	All car data as aggregated, by region, multiplied by municipal travel mode factors
res_agg_j_all_reg_car_ttyfreq	All car data as aggregated, by region, relative shares grouped by TTM year (one TTM year: 100%)
res_agg_j_all_reg_car_ttyfreq_mun*	All car data as aggregated, by region, relative shares grouped by TTM year and multiplied by municipal travel mode factors
res_agg_j_all_reg_pt	All PT data as aggregated, by region
res_agg_j_all_reg_pt_mun*	All PT data as aggregated, by region, multiplied by municipal travel mode factors
res_agg_j_all_reg_pt_ttyfreq	All PT data as aggregated, by region, relative shares grouped by TTM year (one TTM year: 100%)
res_agg_j_all_reg_pt_ttyfreq_mun*	All PT data as aggregated, by region, relative shares grouped by TTM year and multiplied by municipal travel mode factors
res_agg_j_ic_car	IC car data as aggregated
res_agg_j_ic_car_freq	IC car data, relative share of each data year pair of the total for all years
res_agg_j_ic_pt	IC PT data as aggregated
res_agg_j_ic_pt_freq	IC PT data, relative share of each data year pair of the total for all years
res_agg_j_ic_reg_car	IC car data as aggregated by region
res_agg_j_ic_reg_car_icfreq	IC car data as aggregated by region, relative shares of each industry (all industries: 100%)
res_agg_j_ic_reg_car_icfreq_mun*	IC car data as aggregated by region, relative shares of each industry, multiplied by municipal travel mode factors
res_agg_j_ic_reg_car_mun	IC car data as aggregated by region, multiplied by municipal travel mode factors
res_agg_j_ic_reg_car_ttyfreq	IC car data as aggregated by region, relative shares grouped by the TTM year (one TTM year: 100%)
res_agg_j_ic_reg_pt	IC PT data as aggregated by region
res_agg_j_ic_reg_pt_icfreq	IC PT data as aggregated by region, relative shares of each industry (all industries: 100%)
res_agg_j_ic_reg_pt_icfreq_mun*	IC PT data as aggregated by region, relative shares of each industry, multiplied by municipal travel mode factors
res_agg_j_ic_reg_pt_mun	IC PT data as aggregated by region, multiplied by municipal travel mode factors
res_agg_j_ic_reg_pt_ttyfreq	IC PT data as aggregated by region relative shares grouped by the TTM year (one TTM year: 100%)

Table 5. An example result of the aggregation of the IC part of the dataset. The region in question is Otaniemi, in Espoo. I chose the region as an example in a fairly accidental fashion, without any careful consideration. Only IC fields with non-zero values are shown here: Manufacturing ('Manuf'), Wholesale and retail trade ('Trade'), Education and Arts, entertainment and recreation ('ArtsEnt'). The applied travel mode factors are 20.3% for PT and 45.9% for car.

Data as if all commuters (506 in total) were using PT (DB table: res_agg_j_ic_reg_pt)										
Measure	TTM	JourneyYear	RegID	Count	Total	Manuf	Trade	Education	ArtsEnt	
pt_m_tt	2013	2012	0492022000	506	8193	1025	913	6199	56	

Data as if all commuters (506 in total) were using car (DB table: res_agg_j_ic_reg_car)										
Measure	TTM	JourneyYear	RegID	Count	Total	Manuf	Trade	Education	ArtsEnt	
car_m_t	2013	2012	0492022000	506	5410	483	402	4483	42	

PT data as weighted by municipal travel mode factors (DB table: res_agg_j_ic_reg_pt_mun)										
Measure	TTM	JourneyYear	RegID	Count	Total	Manuf	Trade	Education	ArtsEnt	
pt_m_tt	2013	2012	0492022000	103	1663.2	208.1	185.3	1258.4	11.4	

Car data as weighted by municipal travel mode factors (DB table: res_agg_j_ic_reg_car_mun)										
Measure	TTM	JourneyYear	RegID	Count	Total	Manuf	Trade	Education	ArtsEnt	
car_m_t	2013	2012	0492022000	232	2483.2	221.7	184.5	2057.7	19.3	

Equivalent PT data relative shares (DB table: res_agg_j_ic_reg_pt_icfreq_mun)										
Measure	TTM	JourneyYear	RegID	Count	Total	Manuf	Trade	Education	ArtsEnt	
pt_m_tt	2013	2012	0492022000	103	100	12.5	11.1	75.7	0.7	

Equivalent car data relative shares (DB table: res_agg_j_ic_reg_car_icfreq_mun)										
Measure	TTM	JourneyYear	RegID	Count	Total	Manuf	Trade	Education	ArtsEnt	
car_m_t	2013	2012	0492022000	232	100	8.9	7.4	82.9	0.8	

4.5 Visual map analysis

4.5.1 Initial considerations for visualisation

Before analysing the data statistically, I considered it practical to create visualisations of it to be able to observe any changes easier. To prepare the data for visualisation I first considered which of the aforementioned (table 4) 26 aggregated tables are actually truly relevant for the final analysis; I had created the tables in a relatively mechanical fashion, and did not expect to utilise them all. The tables that I actually used as a base for further analysis are marked in table 4.

4.5.2 Visualising all journeys

To analyse the data visually, I started creating maps that demonstrate spatial differences in the data. To enable creating these tables, I created views into the DB to use as attribute tables for the maps. An example of an SQL query for creating such a view is shown in the figure 12. I also integrated a population dataset (City of Helsinki et al., 2019) to the rest of my dataset at this point to be able to weigh the data by population changes. I then used the field calculator tool in QGIS to weigh the total ag-

gregated changes by population changes for every statistical district, except for the Kauniainen area, for which the population data was not available at the same statistical level as the rest of the data.

```

1 CREATE VIEW res_agg_j_all_reg_pt_changes_2015_2018 AS (
2   » SELECT * FROM (
3   » » SELECT r.gid,
4   » » CASE WHEN r.nimi<>' ' THEN r.nimi ELSE 'Kauniainen' END AS nimi,
5   » » r.geom,
6   » » j1."MunID",
7   » » j1."AreaID",
8   » » j1."DistID",
9   » » j1."Count" AS "C1",
10  » » j1."Total" AS "F1",
11  » » j2."Count" AS "C2",
12  » » j2."Total" AS "F2",
13  » » ROUND(CAST(j2."Total"-j1."Total" AS NUMERIC), 2) AS "FChange",
14  » » j3."Total" AS "T1",
15  » » j4."Total" AS "T2",
16  » » ROUND(CAST(j4."Total"-j3."Total" AS NUMERIC)) AS "TChange",
17  » » "Pop2016",
18  » » "Pop2018",
19  » » "Pop2018"- "Pop2016" AS "PopChange"
20  » » FROM hcr_subregions r
21  » » LEFT JOIN res_agg_j_all_reg_pt_ttyfreq_mun j1
22  » » » ON r.kokotun=j1."RegID" AND j1."TTM" = 2015
23  » » LEFT JOIN res_agg_j_all_reg_pt_ttyfreq_mun j2
24  » » » ON r.kokotun=j2."RegID" AND j2."TTM" = 2018
25  » » LEFT JOIN res_agg_j_all_reg_pt_mun j3
26  » » » ON r.kokotun=j3."RegID" AND j3."TTM" = 2015
27  » » LEFT JOIN res_agg_j_all_reg_pt_mun j4
28  » » » ON r.kokotun=j4."RegID" AND j4."TTM" = 2018
29  » » LEFT JOIN hcr_population p
30  » » » ON r.kokotun=p."RegID"
31  » » ) AS Q
32  » WHERE "FChange" IS NOT NULL
33  » AND "TChange" IS NOT NULL
34  » );

```

Figure 12. An example of an SQL query used to create a view for the all-dataset analysis; this particular query creates a view that calculates the PT changes between 2015 and 2018.

The actual maps are presented in figures 14, 16, 18 and 20 for the unweighted changes and in figures 15, 17, 19 and 21 (see chapter 5.1.1) for the population-weighted to visualise the differences in the data. The classification used on the maps

is based on standard deviation to demonstrate any outliers in the data. For the population-weighted maps, the map labels are also based on z -scores, as the changes are neither absolute nor on ratio scale. I also exported the attribute tables of the maps to spreadsheet files to perform further statistical analyses in R.

4.5.3 *Visualising the journeys classified by industry*

Due to the relatively large dimensions of the IC part of the dataset, it would have been considerably time-consuming – and thus quite impractical and error-prone – to visualise these changes manually, even if this had allowed me to create maps with better classification. Thus, I programmatically created four maps for each of the 22 industries: two car maps and two PT maps – one for the change between 2013 and 2015 and one for the change between 2015 and 2018, respectively – to depict the changes; 88 maps in total. An example of these maps for the Human health and social work industry class is shown in figure 29; all of the maps are also available in appendix 2.

Changes shown on the maps are per-district percentage point changes of all journeys within a single industry for two journey year-TTM pairs. Due to technical limitations of the Bokeh library and due to the need to make the maps as comparable with each other as possible, I classified the data on each map to only three classes and manually chose class boundaries for all the maps, loosely based on the variations of the data, but also on maintaining pretty breaks; my chosen class breaks are $\pm 10\%$ -points. Despite the classification limitations, the areas of all districts on the map are labeled by their respective %-point numbers.

Unlike when visualising all journeys (see chapter 4.5.2), I did not create any separate views into the database for analysing the IC journeys. Instead, I integrated an SQL query template into the Python script that I used to create the maps. An example of a single SQL query as formed by that code is presented in the figure 13.

```

1 SELECT * FROM (
2   » SELECT CASE WHEN r.nimi<>' THEN r.nimi ELSE 'Kauniainen' END AS nimi,
3   » rf1."MunID",
4   » rf1."AreaID",
5   » rf1."DistID",
6   » rf1."Count" AS "C1",
7   » rf1."Education" AS "RF_T1",
8   » abs1."Education" AS "T1",
9   » CASE WHEN
10  » » abs1."Education" <> 0
11  » » THEN ROUND(CAST(abs1."Education"/(abs1."Education"/abs1."Total")/
    » abs1."Count" AS NUMERIC), 2)
12  » » ELSE 0
13  » END AS "M1",
14  » rf2."Count" AS "C2",
15  » rf2."Education" AS "RF_T2",
16  » abs2."Education" AS "T2",
17  » CASE WHEN abs2."Education" <> 0
18  » » THEN ROUND(CAST(abs2."Education"/(abs2."Education"/abs2."Total")/
    » abs2."Count" AS NUMERIC), 2)
19  » » ELSE 0
20  » END AS "M2",
21  » ROUND(CAST(rf2."Education"-rf1."Education" AS NUMERIC), 2) AS "RFChange",
22  » ROUND(CAST(abs2."Education"-abs1."Education" AS NUMERIC), 2) AS "AbsChange"
23  » FROM hcr_subregions r
24  » LEFT JOIN res_agg_j_ic_reg_car_icfreq_mun rf1
25  » » ON r.kokotun=rf1."RegID" AND rf1."TTM" = 2013
26  » LEFT JOIN res_agg_j_ic_reg_car_icfreq_mun rf2
27  » » ON r.kokotun=rf2."RegID" AND rf2."TTM" = 2015
28  » LEFT JOIN res_agg_j_ic_reg_car_mun abs1
29  » » ON r.kokotun=abs1."RegID" AND abs1."TTM" = 2013
30  » LEFT JOIN res_agg_j_ic_reg_car_mun abs2
31  » » ON r.kokotun=abs2."RegID" AND abs2."TTM" = 2015
32  » ) AS Q
33 WHERE "AbsChange" IS NOT NULL

```

Figure 13. An example of an SQL query – executed by the map-plotting Python script – that creates the IC maps and calculates statistics for that data. This particular query creates the attribute table for the Education industry maps, when the travel mode is car.

4.6 Statistical analysis

The relative longitudinal shallowness of the TTM dataset does not enable utilising some of the most typical statistical methods such as ordinary least squares (OLS) linear regression to analyse the changes. While the results of the visual map analysis could be used as such to identify notable outliers and to answer my research questions, I also used Student's (1908) two-samples *t*-test to assess the statistical signific-

ance of the change between the journey data-TTM pairs. The statistic can be calculated as follows:

$$t = \frac{m_A - m_B}{\sqrt{\frac{S_A^2}{n_A} + \frac{S_B^2}{n_B}}}$$

where m_A and m_B represent the means of each data group, n_A and n_B represent the group sizes and S^2 is an estimate of the pooled variance of the groups. This estimate can be calculated as follows:

$$S^2 = \frac{\sum (x - m_A)^2 + \sum (x - m_B)^2}{n_A + n_B - 2}$$

with $n_A + n_B - 2$ degrees of freedom. The null hypothesis is equality of means between the data groups. The assumptions of the test are equal variances between the input data groups and normal distribution of residuals between the groups.

To assess the equality of variances I used a Levene-type statistic (Levene, 1960), but anticipating some asymmetry in the distributions, I used median instead of mean values as the central location of each input group, as proposed by Brown & Forsythe (1974) for skewed distributions. For this test, the null hypothesis is homogeneity of variances; i.e. the test should not be significant for the result of t -test to be valid. This statistic is calculated as follows:

$$W = \frac{(N - k)}{(k - 1)} \frac{\sum_{i=1}^k N_i (\bar{Z}_i - \bar{Z}_{..})^2}{\sum_{i=1}^k \sum_{j=1}^{N_i} N_i (Z_{ij} - \bar{Z}_i)^2}$$

where $Z_{ij} = |Y_{ij} - \bar{Y}_i|$ and \bar{Y}_i is the group median of i -th subgroup – that is, the group medians of Z_{ij} . $\bar{Z}_{..}$ is the overall median of Z_{ij} (NIST, 2012).

To test for the assumed normal distribution of residuals between input data pairs used in the t -test, I used numerical Shapiro-Wilk (1965) test for normality, but also made visual checks by quantile-quantile (Q-Q) and histogram plots. As for the Levene test, the null hypothesis for the Shapiro-Wilk test should be upheld for the results of the t -test to be valid; the null hypothesis is that the values are normally distributed. The statistic of the Shapiro-Wilk test is calculated as follows:

$$W = \frac{(\sum_{i=1}^n a_i x_{(i)})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

where the $x_{(i)}$ are the sample values ordered from smallest to largest and a_i are tabulated coefficients generated from an n -sized sample drawn from a normal distribution (Shapiro & Wilk, 1965; NIST, 2012). However, the Shapiro-Wilk test is known to be oversensitive on larger sample sizes (see e.g. Ghasemi & Zahediasl (2012)); thus, the residuals should be plotted and investigated visually, even if the results of the Shapiro-Wilk test suggest rejecting the null hypothesis.

4.6.1 Statistical analyses of all journeys

For the complete, non-IC joined dataset, I performed unpaired two-samples t -tests for each of the journey data-TTM pairs, but also performed t -tests to further compare the two comparisons, i.e. I compared the change between the years 2013 and 2015 to the change between 2015 and 2018. I did this comparison of comparisons for both unweighted and population-weighted data. I also assessed the validity of the t -test assumptions – variance homogeneity and normal distribution of residuals – with Levene-type statistics, Shapiro-Wilk analysis and residual plots. In addition, I fitted simple OLS regression models – $y = \alpha + X\beta + \epsilon$, with the target variable y being the change in the aggregated travel time and the population change X being the only explanatory variable – to the data to check whether population changes statistically explain some part of the observed changes. The results of all statistical analysis are shown in tables 7, 8 and 9. Weighing the dataset by population required me to intersect the Kauniainen area from the results, as for that area no sub-region population data was available. However, less than 1% of the HCR population live in Kauniainen (Aluesarjat, 2020), so I considered making this change to my research area reasonable.

4.6.2 *Statistical analysis of journeys classified by industry*

For the IC part of the dataset, I did not consider it necessary to weigh the data for population changes, as my research question regarding them was to find any significant per-industry changes, regardless of their underlying reasons. Also, this would have been complicated in practice due to the larger dimensions of the IC part. In addition, I did not do any comparison of comparisons to this part of the dataset, as that would have been a very complex operation as well. Thus, for this part of the dataset, I integrated the counting of descriptive statistics, t -tests, Levene-type tests and Shapiro-Wilk tests to my map-producing Python code.

5 Results

5.1 Results of the analysis of all journeys

5.1.1 *Results of the visual map analysis*

The most observable trends appeared to be a continuing decrease in aggregated commuting times in the Northeastern and Eastern districts of Malmi, Mellunkylä and Vuosaari; this trend was observable regardless of the transport mode, but was more pronounced on PT. Another observable trend was a general increase of commuting times in the Kivistö district. All the changes depicted on the maps and their equivalent z -scores are listed on table 6.

5.1.1.1 *Public transport maps*

The single clear trend observable from the maps was a decrease in PT commuting times in the districts of Vuosaari, Malmi and Mellunkylä. These three districts showed a decrease of at least one standard deviation on all maps, whether adjusted for population changes or not. Vuosaari was the most notable outlier, demonstrating a decrease of more than three standard deviations in all cases.

Between 2013 and 2015 (figure 14), there was a decrease of more than three standard deviations in Kaarela, Laajasalo, Malmi and Vuosaari, and in one district, at Kampinmalmi, an increase of more than three standard deviations could be ob-

served. When adjusted for changes in population only Kampinmalmi and Vuosaari still demonstrated a change of more than three standard deviations (figure 15).

Between 2015 and 2018, there were reductions in absolute aggregated times in many districts (figure 16). These changes were most noticeable in Lauttasaari and Vuosaari, with reductions equivalent to over three standard deviations from the all-HCR mean. In no area could an increase of equal proportions be observed; the largest increase was shown in the Kivistö district.

Weighing the data for population slightly changed the number of districts where changes deviated more than one standard deviation from the all-HCR mean (figure 17), and led to some districts seeing more and others seeing less pronounced increases. Lauttasaari and Vuosaari were still among the most extreme outliers with reductions equivalent to more than three standard deviations from the mean, with Mellunkylä now rising to this category as well, and the district with largest increase was still observed at Kivistö, with 2.8 standard deviations above mean.

5.1.1.2 Private car maps

According to my maps, between the years 2013 and 2015 the aggregated commuting times by car were decreasing most notably in several western districts of the HCR, in the municipality of Espoo (figure 18). Between 2015 and 2018 this effect was partially reversed, as the aggregated times in the same districts were increasing again (figure 20), even if not quite so strongly. Weighing the data for population somewhat reduced this effect (figures 19 and 21, respectively), but did not eliminate it entirely.

In the Eastern parts of the area, in certain districts of Eastern Helsinki, the travel times by car seemed to have a slightly decreasing trend between both 2013–2015 and 2015–2018 journey data–TTM pairs. Weighing for population led to smaller z -scores on some of these districts, but larger on others. Vuosaari, the most consistent outlier on PT, did see some decrease by car as well.

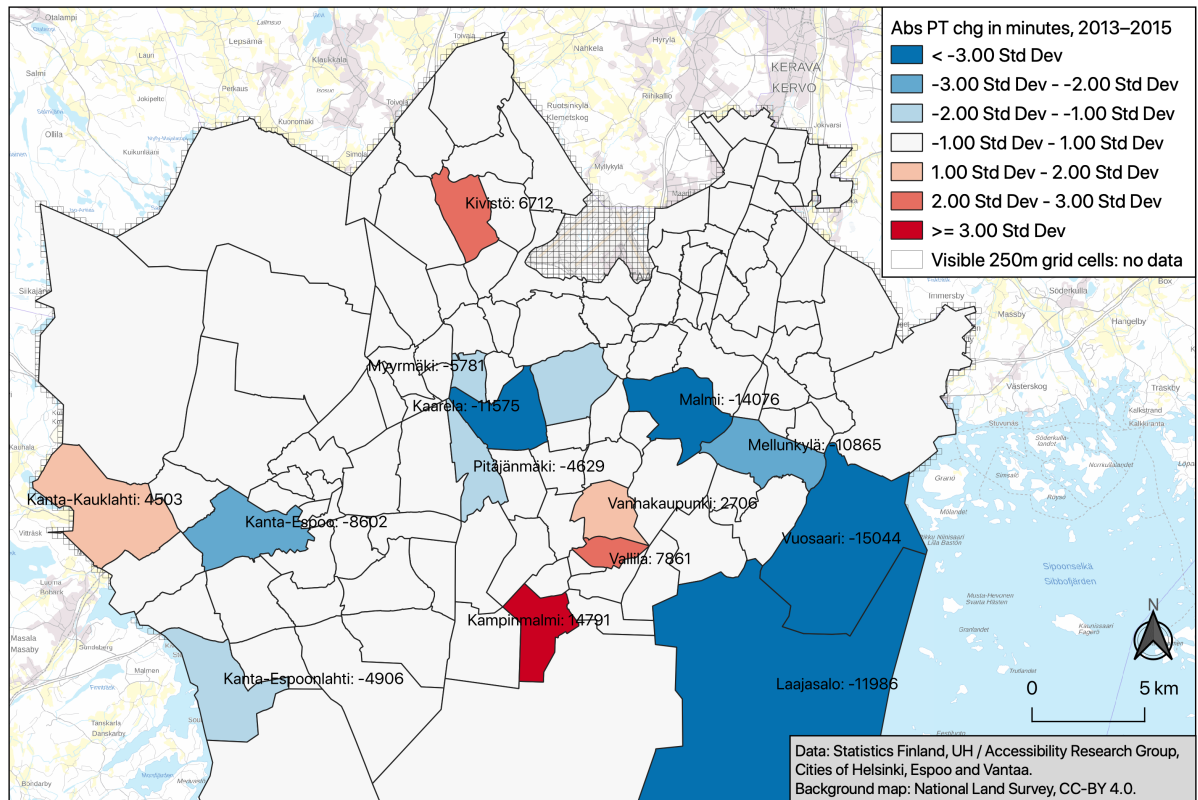


Figure 14. Map of the changes in aggregated travel times by PT between 2013 and 2015.

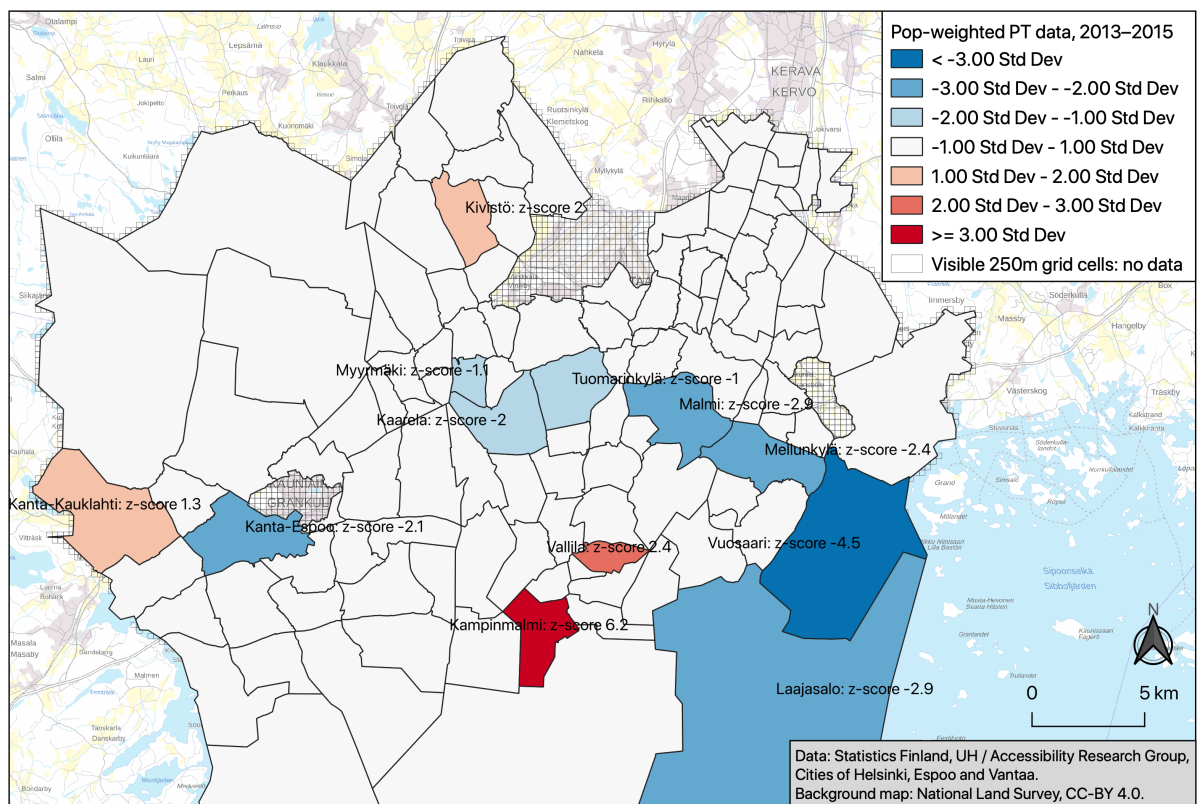


Figure 15. Map of the changes in aggregated travel times by PT between 2013 and 2015, weighted for population changes.

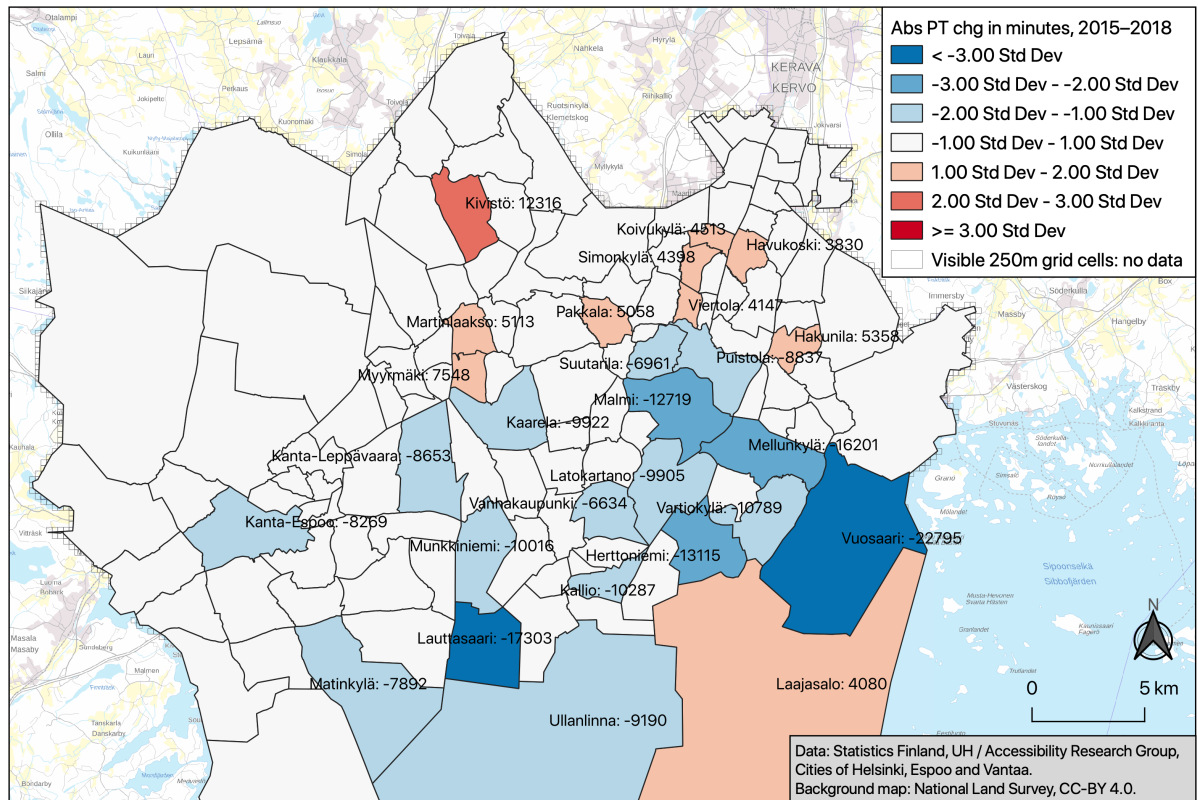


Figure 16. Map of the changes in aggregated travel times by PT between 2015 and 2018.

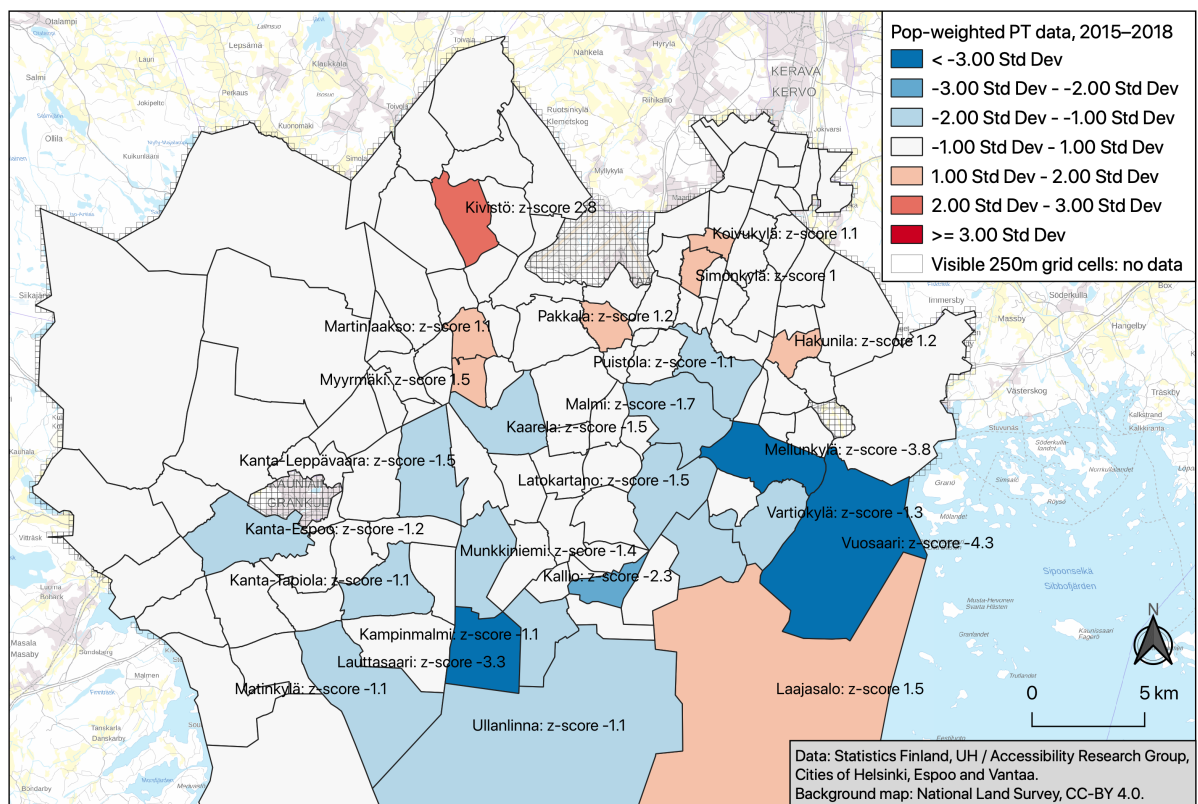


Figure 17. Map of the changes in aggregated travel times by PT between 2015 and 2018, weighted for population changes.

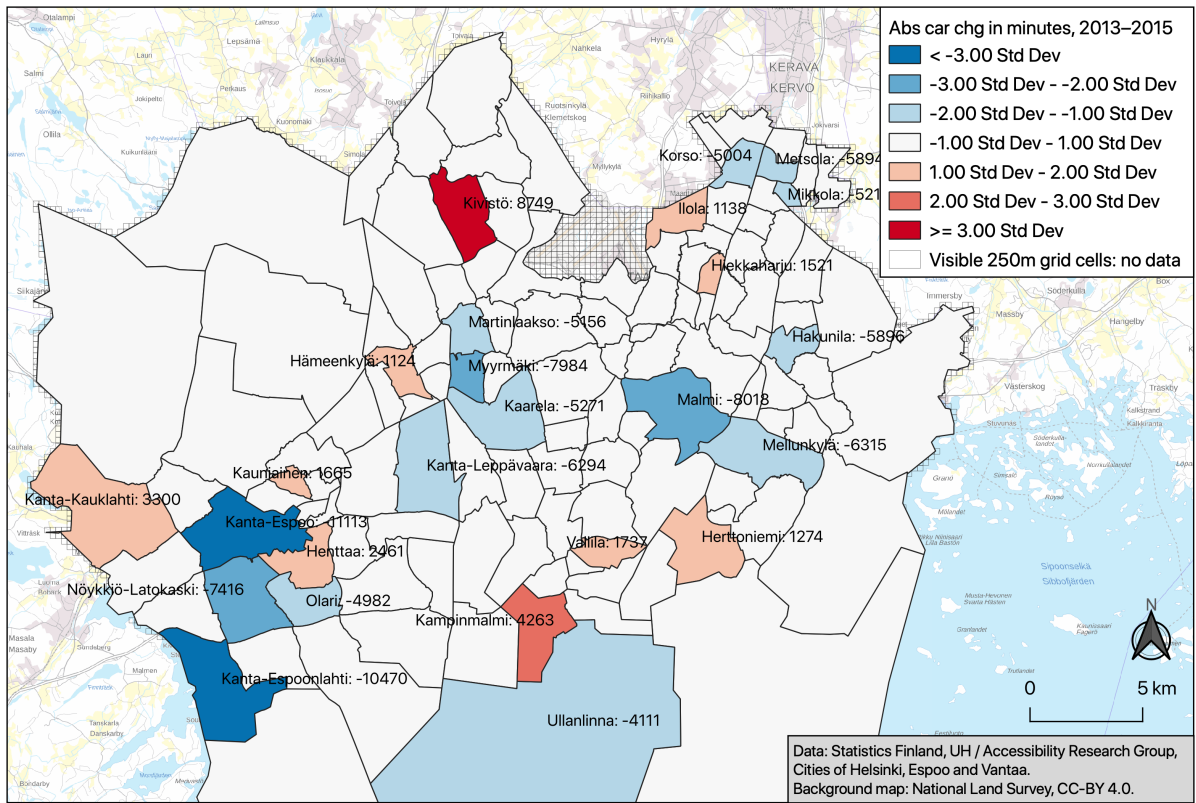


Figure 18. Map of the changes in aggregated travel times by car between 2013 and 2015.

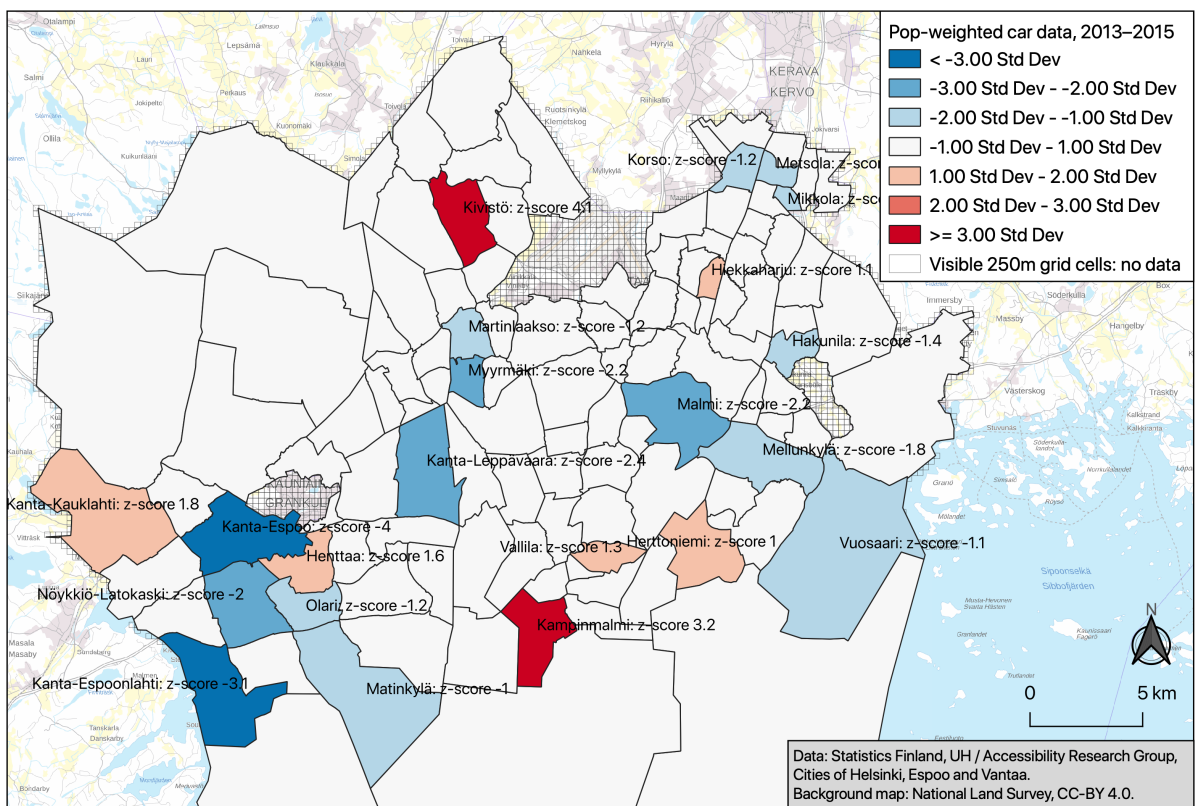


Figure 19. Map of the changes in aggregated travel times by car between 2013 and 2015, weighted for population changes.

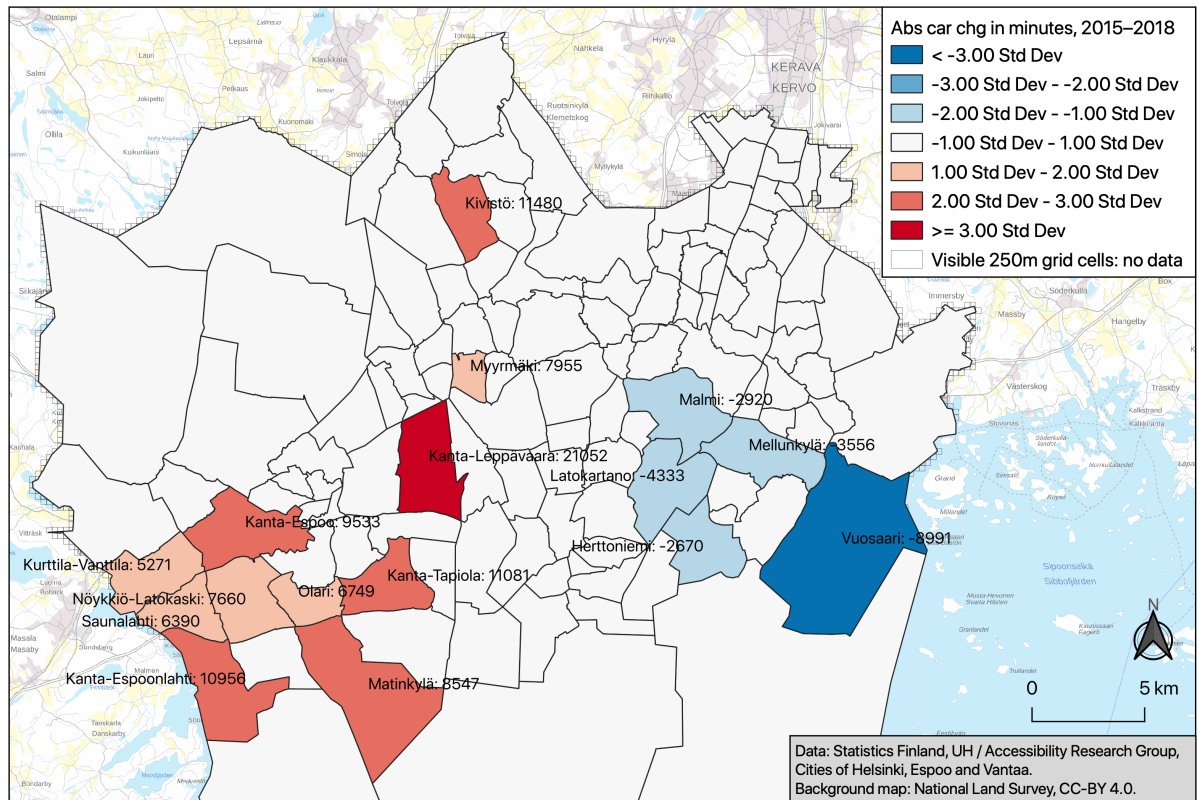


Figure 20. Map of the changes in aggregated travel times by car between 2015 and 2018.

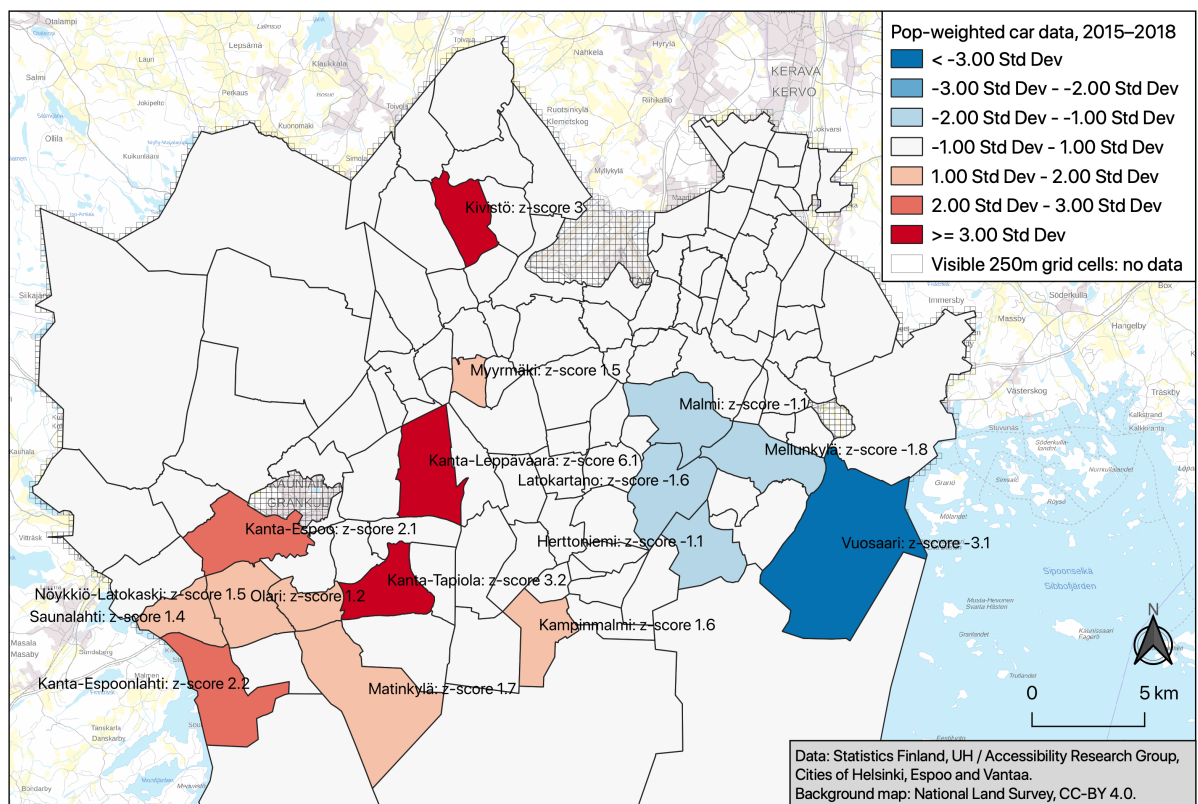


Figure 21. Map of the changes in aggregated travel times by car between 2015 and 2018, weighted for population changes.

Between 2013 and 2015 the districts with an absolute change of more than three standard deviations were Kanta-Espoo, Kanta-Espoonlahti and Kivistö; for the first two, the change meant a decrease, but for Kivistö, an increase. When weighted for population changes, the equivalent districts with commuting time changes of more than three standard deviations were Kampinmalmi, Kanta-Espoo, Kanta-Espoonlahti and Kivistö; for the first and last of these, the change marked increase, while for the two middle ones it marked decrease. Notably, one of the unnamed subdistricts of Kauniainen demonstrated a change of more than one standard deviation in absolute numbers, but could not be weighted for population changes due to missing district-level population data.

Between 2015 and 2018, the districts with a change magnitude of three standard deviations or more were Kanta-Leppävaara, and Vuosaari; for the former, the change marked an increase and for the latter, a decrease. When weighted for population, Kanta-Leppävaara had been joined by Kanta-Tapiola and Kivistö as districts demonstrating an increase of over three standard deviations, whereas Vuosaari was still shown to see an equivalent decrease.

5.1.2 Results of the statistical analysis

5.1.2.1 Comparisons between individual journey data–TTM pairs

For the complete dataset without IC, all the results of the statistical tests between individual journey data–TTM pairs are listed in the table 7. Kernel densities of the distributions of unweighted changes are shown in figure 22 and of population-weighted changes in figure 23, respectively. The density plots suggest that kurtosis of the weighted distributions are higher than that of unweighted distributions. They also suggest that the distribution of PT changes between 2015–2018 has notably lower kurtosis and is notably more skewed than the other distributions.

The assumptions of the *t*-test require equality of variances between the input data groups and normal distribution of residuals between the groups. With regard to the variance equality requirement, only the Levene test for the car data comparison between the years 2015 and 2018 is significant; the others are not, and thus the null hypotheses cannot be rejected for those.

Table 7. Descriptive statistics and t-test results of the individual journey data–TTM pairs.

Travel Mode, Pair	PT, 2013–2015	PT, 2015–2018	Car, 2013–2015	Car, 2015–2018
1 st year: Minimum	166	24	158	17
1 st year: 1 st Quantile	9,018	8,754	9,164	7,934
1 st year: Median	19,720	19,397	20,412	18,545
1 st year: Mean	40,094	38,910	27,846	26,162
1 st year: 3 rd Quantile	57,176	54,443	38,177	36,551
1 st year: Maximum	233,530	218,486	116,086	109,792
1 st year: Standard Dev.	46,743	45,206	25,418	23,955
2 nd year: Minimum	168	27	153	19
2 nd year: 1 st Quantile	8,845	9,773	8,539	8,831
2 nd year: Median	19,521	20,217	18,769	20,523
2 nd year: Mean	39,223	37,465	26,373	27,761
2 nd year: 3 rd Quantile	54,993	50,064	36,609	38,246
2 nd year: Maximum	218,486	195,691	109,792	130,844
2 nd year: Standard Dev.	45,253	41,077	23,936	25,479
Degrees of freedom	246	248	246	232
Levene (median), W	0.0273	0.2647	0.1383	12.1604
Levene (median), p	0.8690	0.6074	0.7103	0.0006
Shapiro, W	0.7612	0.9110	0.8913	0.8273
Shapiro, p	6.36E-13	4.79E-07	4.86E-08	2.19E-10
Two-sample t-test	0.1490	0.2644	0.4700	-0.5112
Two-sample t-test, p	0.8817	0.7917	0.6387	0.6097

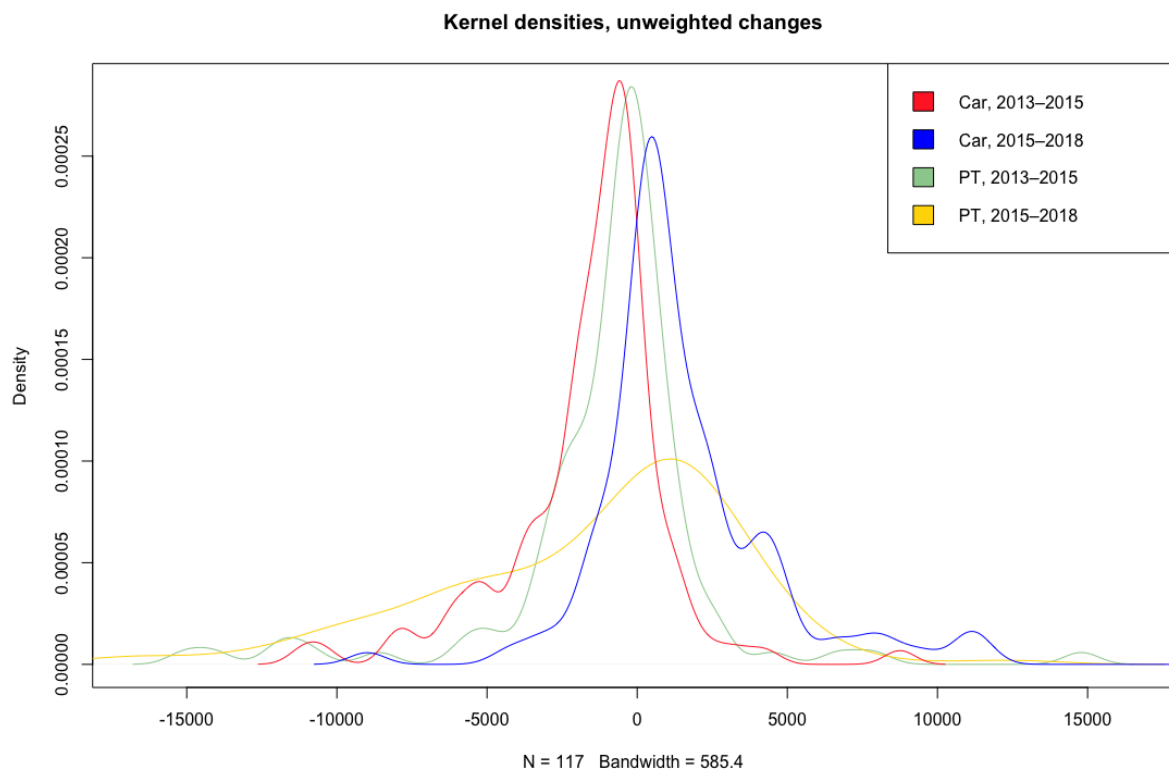


Figure 22. Kernel densities of the distributions of the unweighted changes.

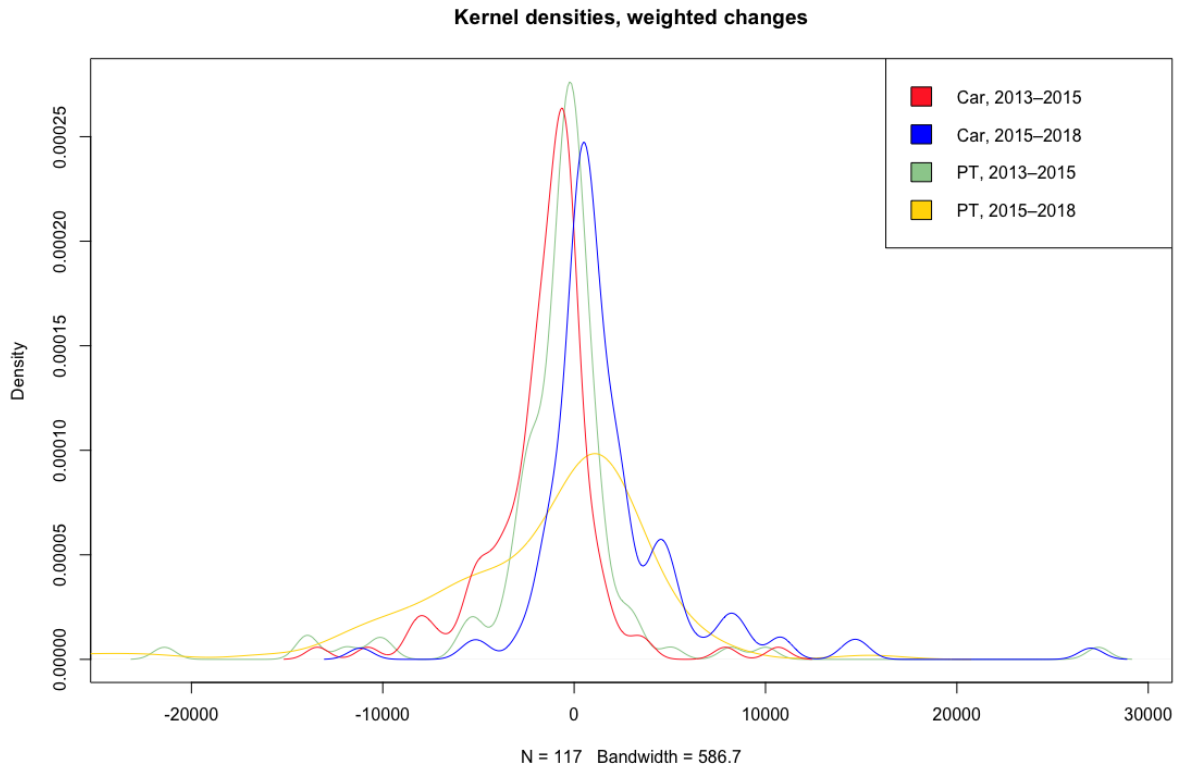


Figure 23. Kernel densities of the distributions of the population-weighted changes.

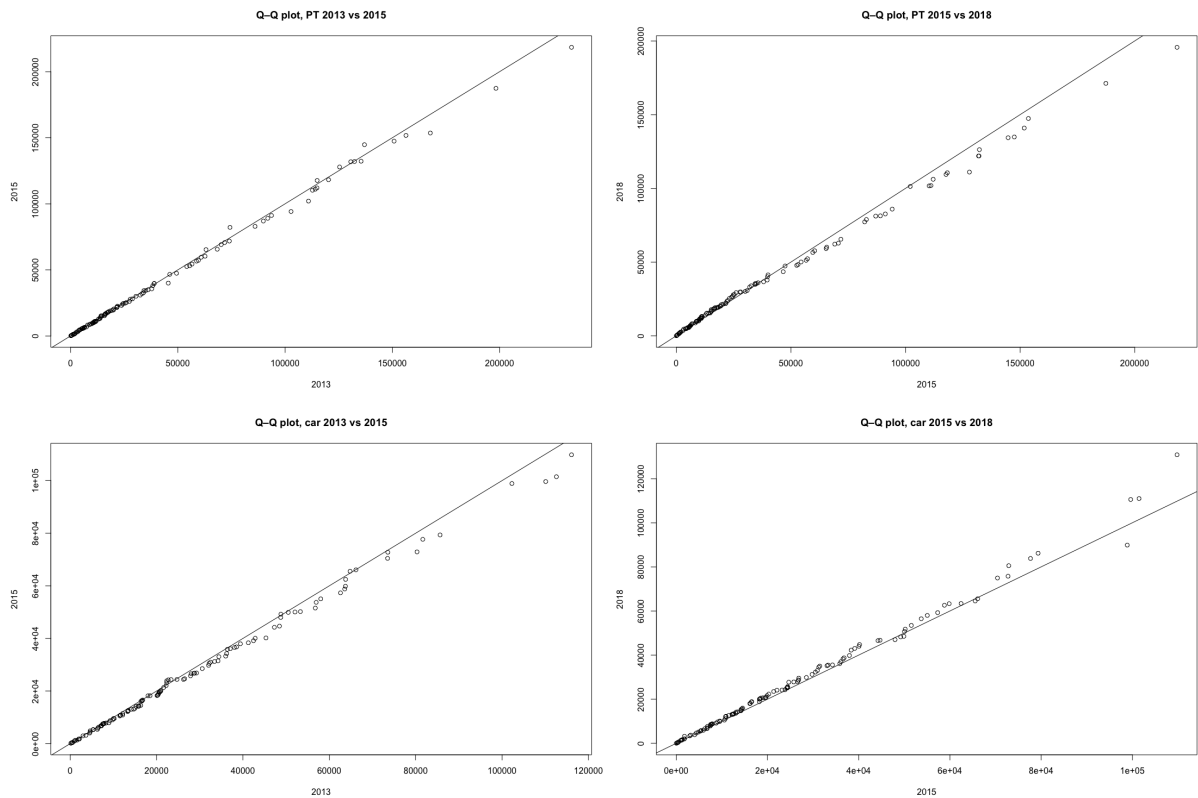


Figure 24. Q-Q plots of the individual journey data-TTM pairs.

The Shapiro-Wilk tests are highly significant for all the data pairs, which suggests that the null hypothesis should be rejected. However, due to the known oversensitivity of Shapiro-Wilk (see chapter 4.6), a visual investigation of the Q-Q plots is recommended; plots for all the individual journey data-TTM pairs are shown in the figure 24. The interpretation of the plots suggests that in reality, residuals are distributed fairly normally and the null hypotheses can be upheld.

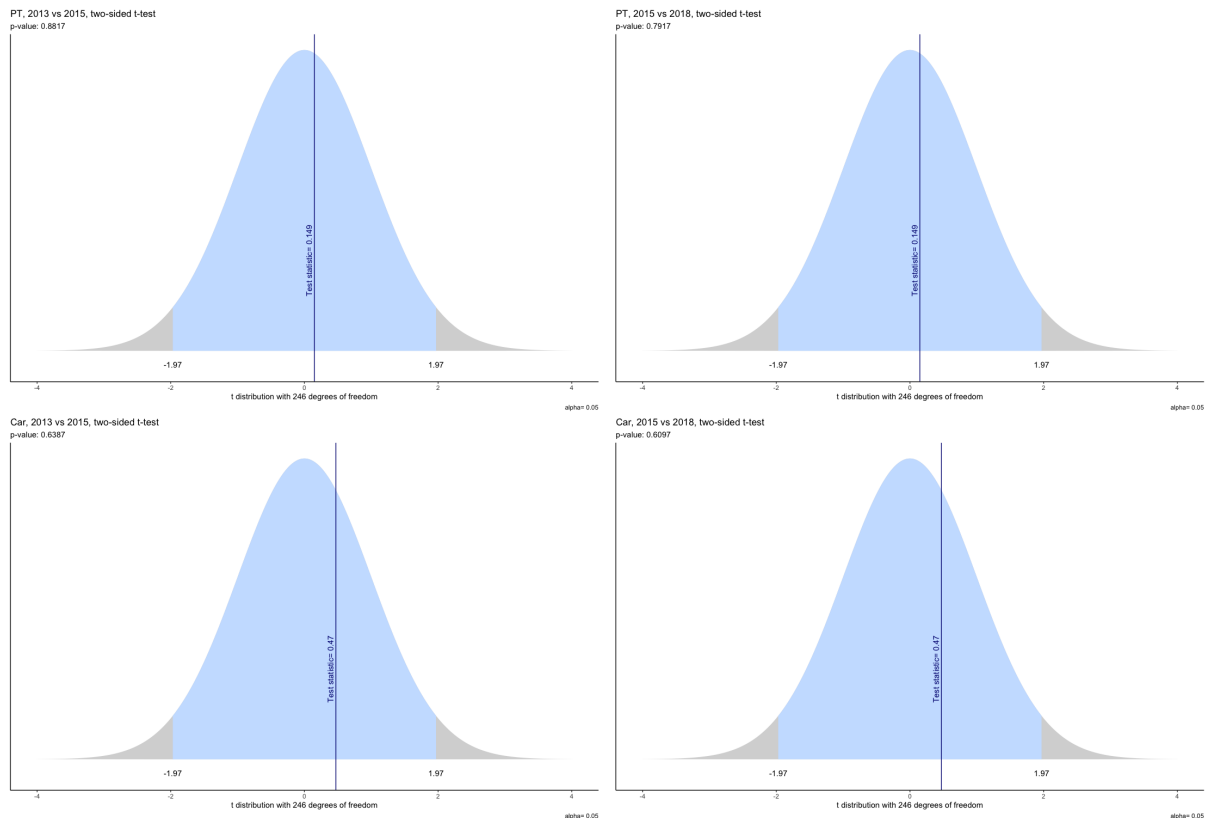


Figure 25. Density plots of the *t*-tests of the individual journey data-TTM pairs

However, even though the data of three of the four data pairs, PT 2013–2015, PT 2015–2018 and car 2013–2015 is suitable for *t*-tests, none of the test results are significant. Thus, none of the mean values between the compared data pairs have statistically significant differences. Density plots of the *t*-tests showing the distribution vs test statistic are shown in figure 25.

Table 8. Results of the linear models for assessing the effects of population changes.

Target variable: total change in travel time				Regression coefficient: population change				
Travel Mode	Year Pair	Est. Coeff.	Std. Error	<i>t</i> -value	Adj. R ²	Signif.	Signif. Code	Missing Obs.
PT	2013–2015	1.2622	0.3778	3.3410	0.0806	0.0011	**	7
PT	2015–2018	-1.8870	1.0000	-1.8870	0.0212	0.0617	.	6
Car	2013–2015	0.6387	0.2877	2.2200	0.0328	0.0284	*	7
Car	2015–2018	1.4758	0.6553	2.2520	0.0334	0.0262	*	6

Significance codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.'

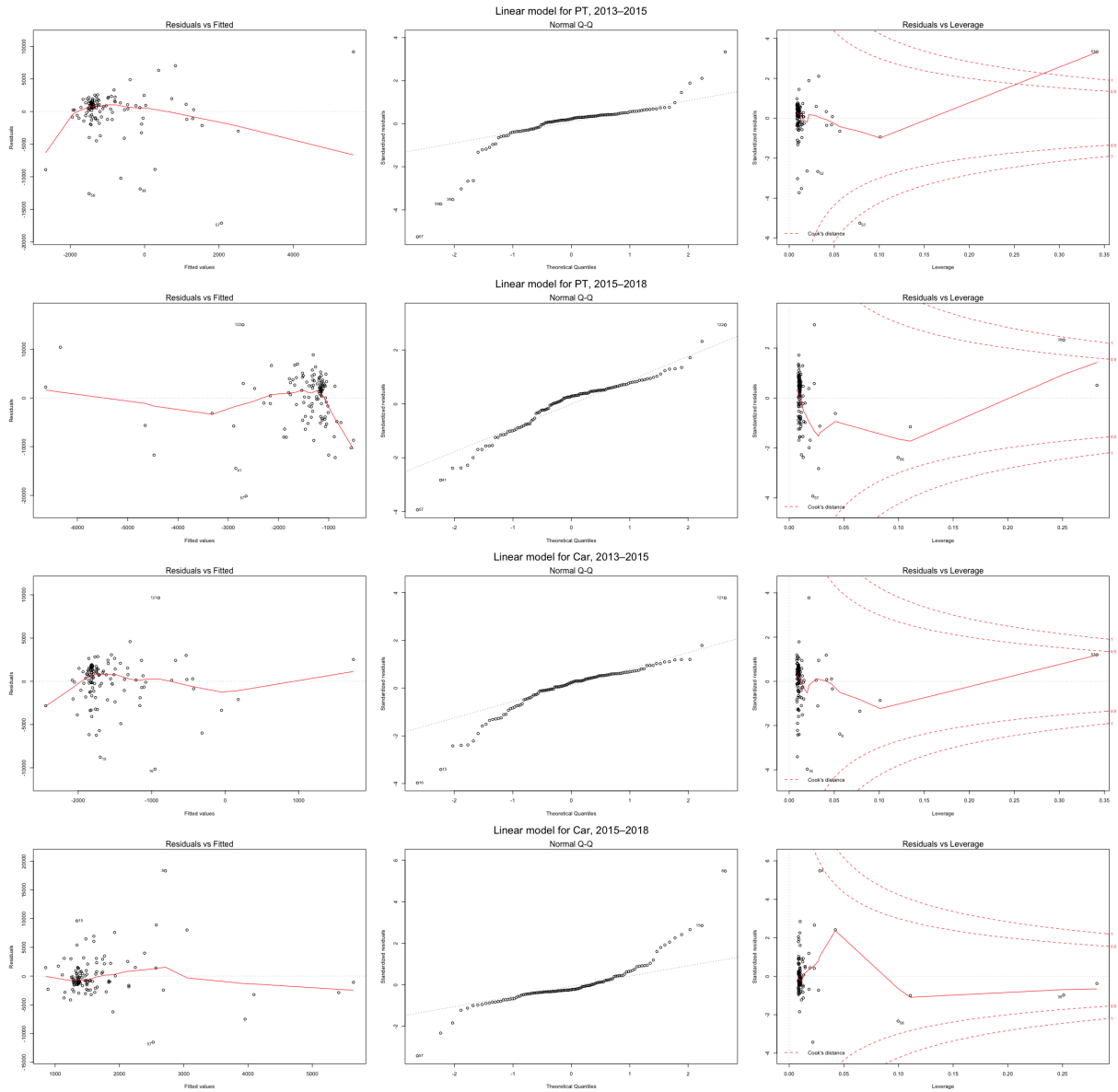


Figure 26. Plots of the residuals of the linear models.

5.1.2.2 Assessing the effects of population changes

To assess the potential effects of population changes, I constructed simple linear models for each of the individual journey data-TTM pairs. The null hypothesis for these models is that the commuting time changes are not dependent on population changes. The results of these tests are listed in the table 8. However, regardless of their perceived significance and their adjusted coefficients of determination, the errors of all these models appear to have variances that are not constant; see the first plot column in figure 26 for residuals vs fitted plots for the data pairs. This implies that the errors are dependent on the explanatory variable; thus, the basic assump-

tions of the OLS models are not upheld, and the models have to be rejected. Due to some extreme outliers, normal distribution of residuals can be questioned as well; see second column in figure 26 for Q–Q plots of the models. Also, certain single observations clearly have outside influence on the models; see the third column in figure 26 for residuals vs leverage plots.

Table 9. Descriptive statistics and t-test results of the 2013–2015 and 2015–2018 changes compared to each other.

<i>Travel Mode, Weighting</i>	<i>PT, Uweighted</i>	<i>PT, Pop-weighted</i>	<i>Car, Unweighted</i>	<i>Car, Pop-weighted</i>
<i>1st pair: Minimum</i>	-15,044	-28,327	-11,113	-13,416
<i>1st pair: 1st Quantile</i>	-1,924	-4,962	-2,243	-2,561
<i>1st pair: Median</i>	-320	153	-1,032	-1,036
<i>1st pair: Mean</i>	-931	-1,774	-1,567	-1,595
<i>1st pair: 3rd Quantile</i>	335	1,880	-293	-291
<i>1st pair: Maximum</i>	14,791	15,430	8,749	10,667
<i>1st pair: Standard Dev.</i>	3,539	6,186	2,628	2,960
<i>2nd pair: Minimum</i>	-22,795	-28,327	-8,991	-11,173
<i>2nd pair: 1st Quantile</i>	-4,962	-4,962	159	156
<i>2nd pair: Median</i>	153	153	742	815
<i>2nd pair: Mean</i>	-1,549	-1,774	1,674	1,877
<i>2nd pair: 3rd Quantile</i>	1,817	1,880	2,446	2,596
<i>2nd pair: Maximum</i>	12,316	15,430	21,052	27,005
<i>2nd pair: Standard Dev.</i>	5,277	6,186	3,474	4,167
<i>Degrees of freedom</i>	232	232	232	232
<i>Levene (median), W</i>	17.0748	12.1604	1.3339	1.6289
<i>Levene (median), p</i>	0.0001	0.0006	0.2493	0.2031
<i>Shapiro, W</i>	0.9160	0.8273	0.7706	0.7214
<i>Shapiro, p</i>	1.849E-06	2.189E-10	3.084E-12	1.329E-13
<i>Two-sample t-test</i>	1.0512	1.2947	-8.0473	-7.3486
<i>Two-sample t-test, p</i>	0.2943	0.1967	4.373E-14	3.395E-12

5.1.2.3 Comparison of comparisons

To further assess the longitudinal changes in the dataset, I also compared the 2013–2015 and 2015–2018 comparisons with each other. The results for these tests are listed in table 9.

For these results, the car changes have lower standard deviations than the PT changes. The Levene tests for both weighted and unweighted PT comparisons are significant; thus the data null hypothesis of variance equality should be rejected and the results of the *t*-tests ignored. In any case, the results of the *t*-tests for the PT comparisons are not significant.

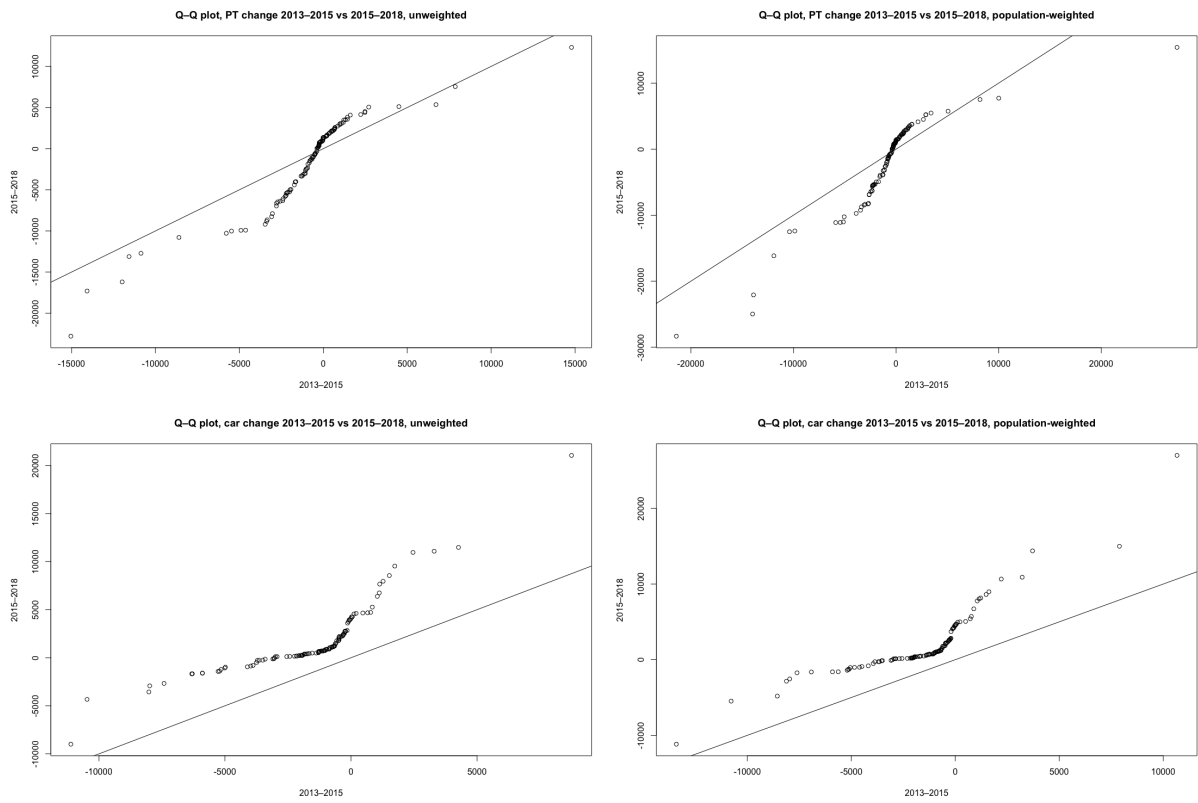


Figure 27. Q-Q plots of the 2013–2015 and 2015–2018 changes compared to each other.

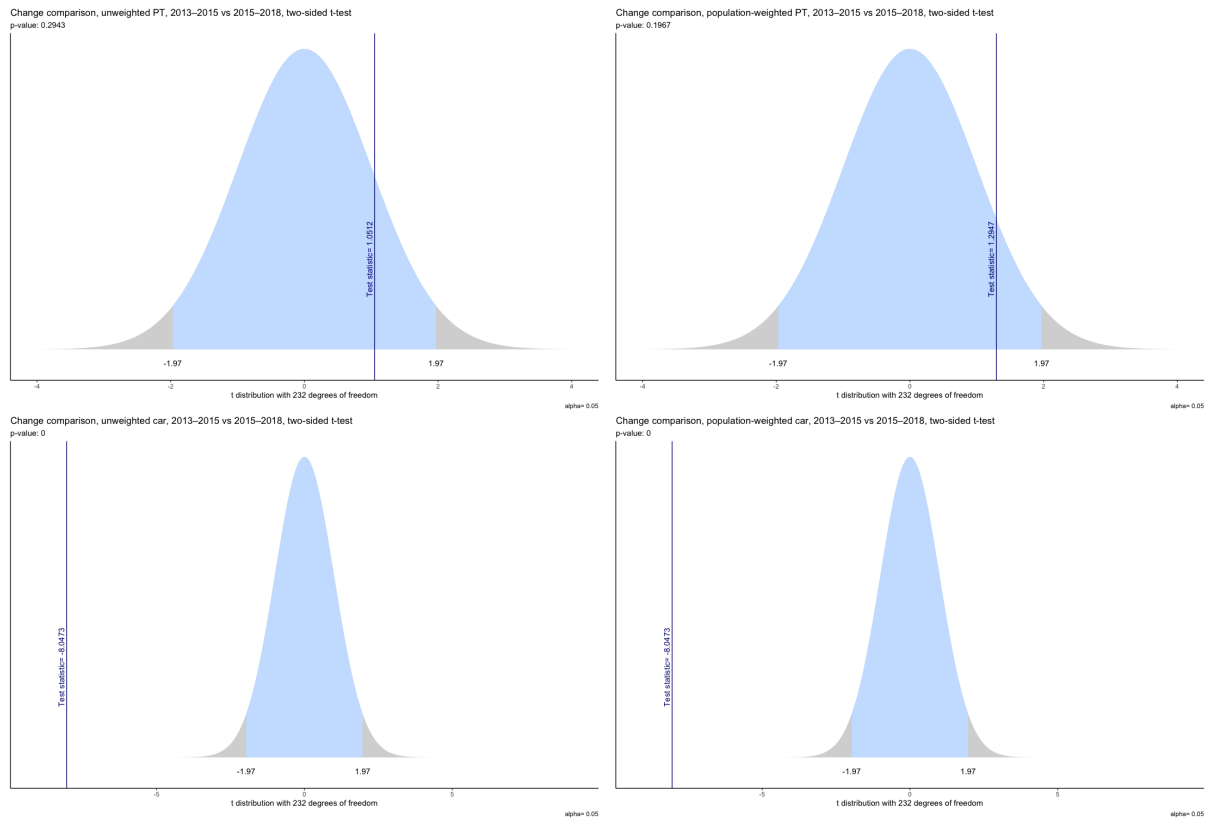


Figure 28. Density plots of the *t*-test distributions of the 2013–2015 and 2015–2018 changes compared to each other.

For car comparisons, the Levene tests are not significant, and thus the null hypothesis of variance equality cannot be rejected. The Shapiro-Wilk tests are significant, however, which suggests rejecting the null hypothesis of normal distribution of residuals. An inspection of the Q-Q plots (figure 27) suggests that the distribution is quite skewed and the normality hypothesis is somewhat questionable, and should be taken into account while interpreting the results of the actual t -tests. The results of the t -tests themselves are statistically significant, suggesting that the null hypothesis of the equality of means between the data groups should be rejected; thus, the yearly comparisons of 2013–2015 and 2015–2018 appear to have different means, for both the unweighted and weighted comparisons. Density plots of the t -tests showing the distribution vs test statistic are shown in figure 28; they suggest that the means of the comparisons of car times are relatively far apart from each other.

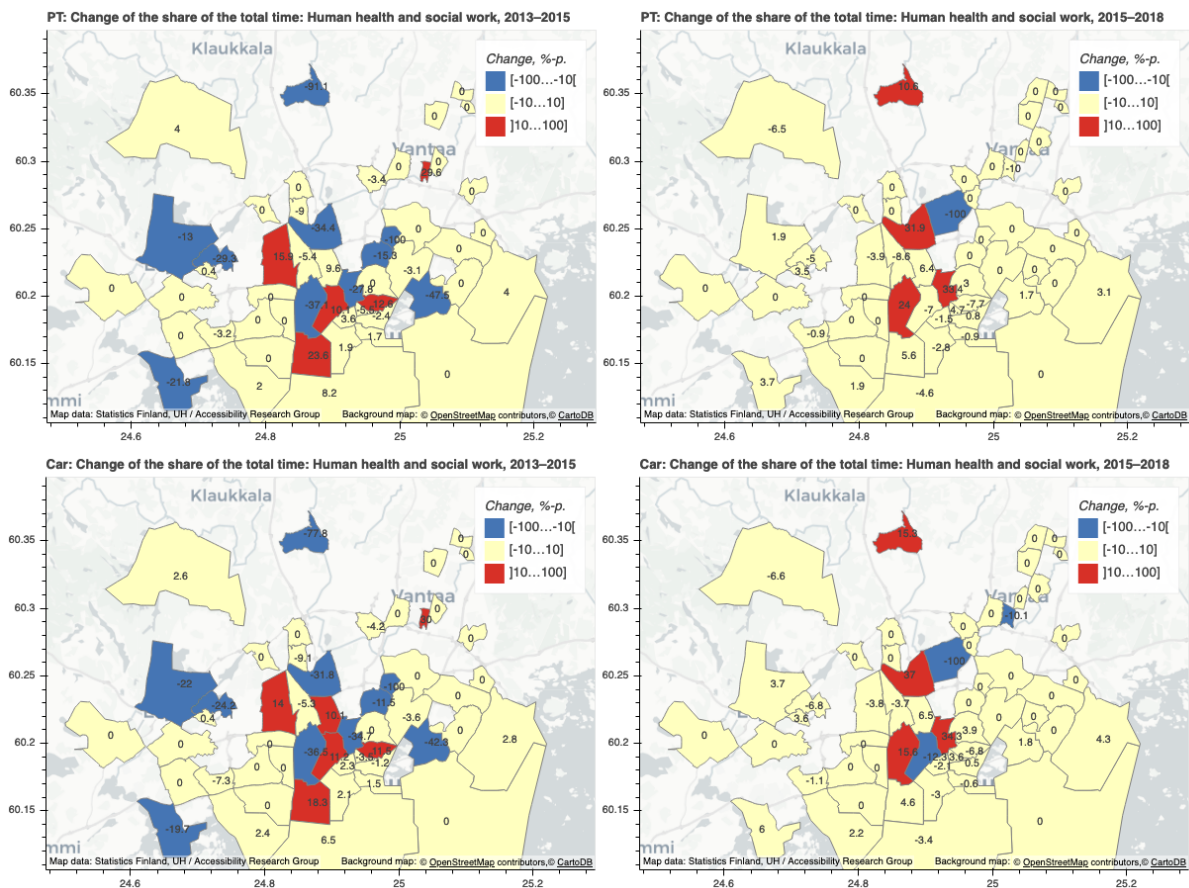


Figure 29. An example of industry-classified maps (example industry: Human health and social work). The changes depicted on maps are percentage point changes between the all-industries total of 100%; i.e. they demonstrate, whether the share of journeys of a particular industry in a specific district has changed significantly, when compared to other industries. The maps of all 22 industries are available as appendix 2.

5.2 Results of the analysis of the journeys classified by industry

5.2.1 Results of the visual map analysis

The mapped changes were relative changes between different industries; not all of the industries were subjected to changes large enough that they would have been readily visible on the maps. The industries where either the PT or the car changes exceeded ten percentage points in absolute numbers on at least one area are listed on table 10. An example of the maps for Health and

Table 10. Industries where relative change between industries in at least one district exceeded ten percentage points.

Industry	Districts with abs chg > 10 %-p.
Administrative and support services	Both car and PT
Arts, entertainment and recreation	Car only
Construction	Both car and PT
Education	Both car and PT
Financial and insurance activities	Both car and PT
Hotels, restaurants and catering	Both car and PT
Human health and social work	Both car and PT
Information and communication	Both car and PT
Manufacturing	Both car and PT
Other service activities and NGOs	Both car and PT
Public administration and defence	Both car and PT
Speciality professions	Both car and PT
Transportation and storage	Both car and PT
Wholesale and retail trade	Both car and PT

human services is shown in figure 29; all 88 maps for all the industries are available in appendix 2. Some districts on the maps are not mapped at all; if the district had no commuting data due to privacy reasons (see chapter 3.2), it is not shown on the map at all. There were also some industries where per-district changes on some districts reached 100%.

5.2.2 Results of the statistical analysis

All results of the statistical analysis are presented in table 11. None of the Levene tests are significant; thus, the null hypothesis of variance equality cannot be rejected. All Shapiro-Wilk tests are significant, suggesting that the null hypothesis of normal distribution of residuals should be rejected. Q-Q plots suggest that the changes in most areas are nonexistent and few extreme outliers define the distribution; this is further suggested by the observation that the median change of every single industry is zero (see table 11). An example of the Q-Q plots for Human health and social work is shown in figure 30; the rest of the Q-Q plots are available in appendix 3. The histogram plots further suggest that the kurtosis of the residual distributions are high. An example of the histogram plots for Human health and social work is shown in figure 31; the rest of the histograms are available in appendix 4. As the visual ana-

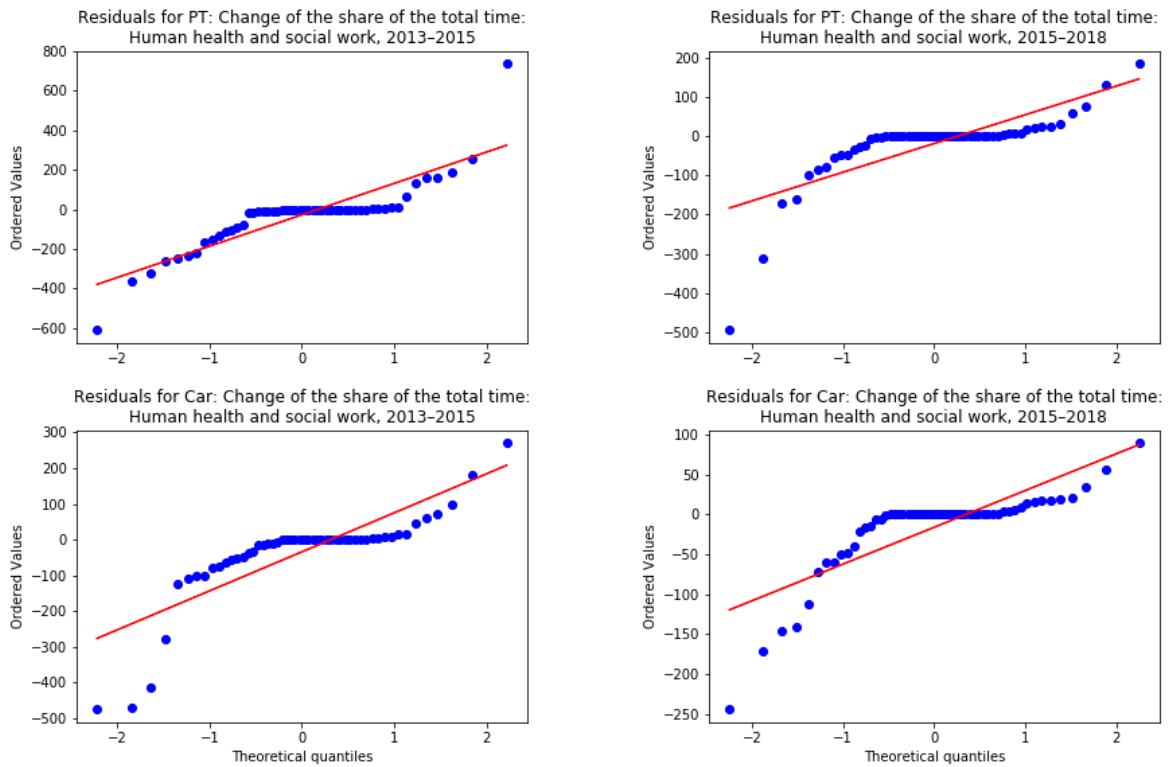


Figure 30. An example of the Q-Q plots for IC data (example industry: Human health and social work). Q-Q plots for all 22 industries are available in appendix 3.

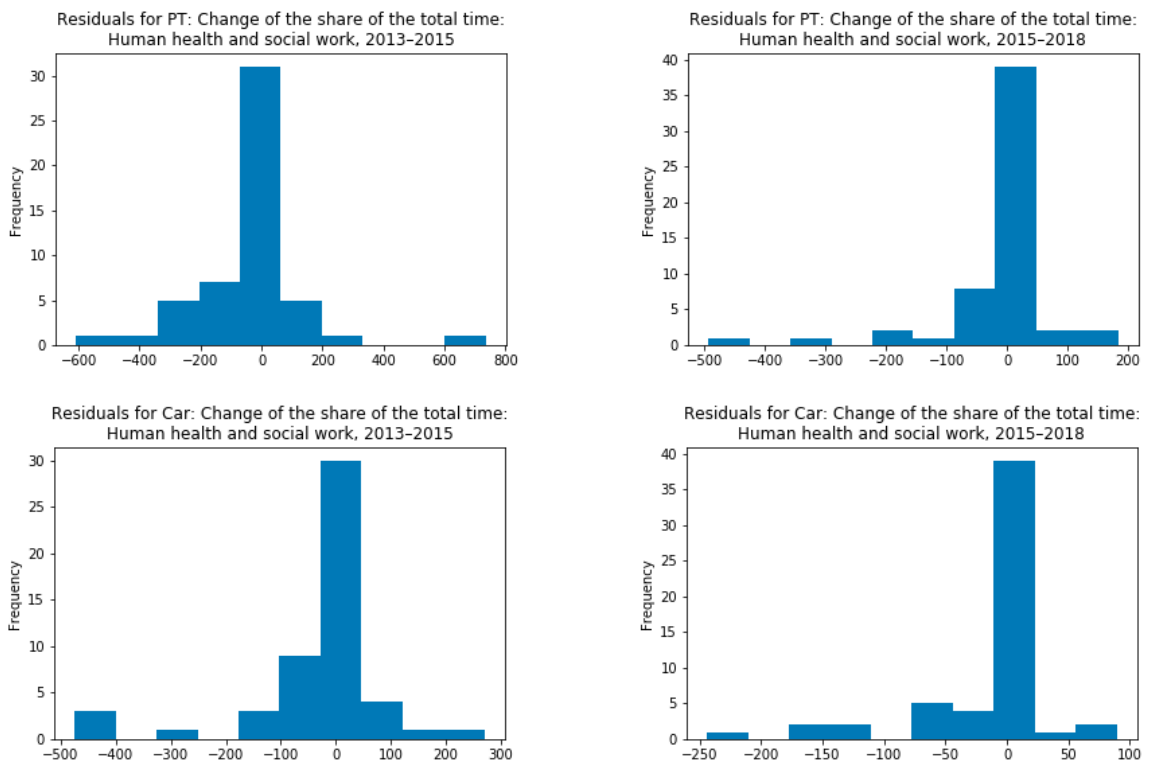


Figure 31. An example of the histogram plots for IC data (example industry: Human health and social work). Histogram plots for all 22 industries are available in appendix 4.

While the median change – and thus the smallest absolute change – of every single industry is zero, a further exploration of the descriptives of the different data pairs shows that between 2013 and 2015, the industries with greatest observed absolute changes in aggregated times for both car and PT journeys were Wholesale and retail trade, Education and Human health and social work, in this respective order. Between 2015 and 2018, the greatest absolute changes for car were in Education, Human health and social work and Public administration and defence, in that order. For PT, they were in Public administration and defence, Human health and social work and Education, in the same order.

6 Discussion

6.1 Factors affecting the changes in aggregated times

6.1.1 *Effects of population changes*

While adjusting the results for population changes appeared to slightly increase the kurtosis of the distributions (see figures 22 and 23 for the distributions of unweighted and weighted changes, respectively), it had no effect on the significance of any results. Thus, changes in population did not appear to have any big role. By visual map analysis, significant increases in travel times were observable in the areas of districts Jätkäsaari and Kivistö. The former is a new city subdistrict in an area that was the former location of the main container harbour of the city; the latter an entirely new district built to the west of the airport, along the route of the new Ring Rail rail connection. Both areas saw significant increases in their aggregated commuting times, even if adjusted for their population increases.

However, the population adjustment did have some effect on several other districts: In both 2013–2015 and 2015–2018, several regions in the Eastern Helsinki saw smaller decreases in their PT commuting times, when adjusted for population. Likewise, the increase in Kivistö at Vantaa was less dramatic, with the area being transferred to a lower category on the map – if only barely. On the other hand, weighing for population made Kampinmalmi to an even more extreme outlier, even if the map classification defies this, and the change is only shown by looking at the table 6;

its z -score increased from 4.5 to 6.2, suggesting that the transport network failed to keep up with the demands of the increased population.

By car, the effects of population adjustments did not appear to have a big role. In some areas, e.g. in Kivistö, the effects appear to be relatively dramatic on map, but in reality both the unweighted and weighted values fall close to the same class border. The notable exception is Jätkäsaari, where, as with the PT times, also car-based times have increased more than the increase in population would suggest. The congestion situation of the Jätkäsaari subdistrict has been seen as problematic both in public discussion and by the City of Helsinki itself (see e.g. Nervola, 2016) and new transport network solutions – even radical and controversial ones – have been proposed to alleviate the situation. (see e.g. Paananen, 2020).

6.1.2 Plausible effects of transport network changes

My first research question was ‘*Have the recent changes in the transport infrastructure had any significant effect to aggregated commuting times in any particular subregion of the HCR?*’ (see the first chapter). For the PT network changes, this is plausible. At least the following changes are observable:

- The Eastern districts that saw the most significant reductions in their commuting times are served by a new, relatively densely trafficked East–West trunk bus route No 560, which route runs through those districts (see figure 32). This line started operating on 10th August 2015 (HSL, 2015); thus, its schedule was included in the 2015 TTM.
- Between 2013 and 2015, the commuting times in Kivistö district were not increased as much as the increase in population would suggest; this suggests that the new Ring Rail – that started operations on 1st July 2015 – took off part of the pressure of the population increase.
- Vuosaari and Lauttasaari districts saw very significant decreases in their commuting times between 2015 and 2018; other districts along the metro line saw decreases as well. As the new, western extension of the metro line started operating on 18th November 2017, these reductions are plausibly network change induced as well. However, this effect has the potential to be seen as a

more pronounced one in my results as it really is due to unavailability of fresh commuting data (see the first sub-point in the chapter 6.2.1 below).

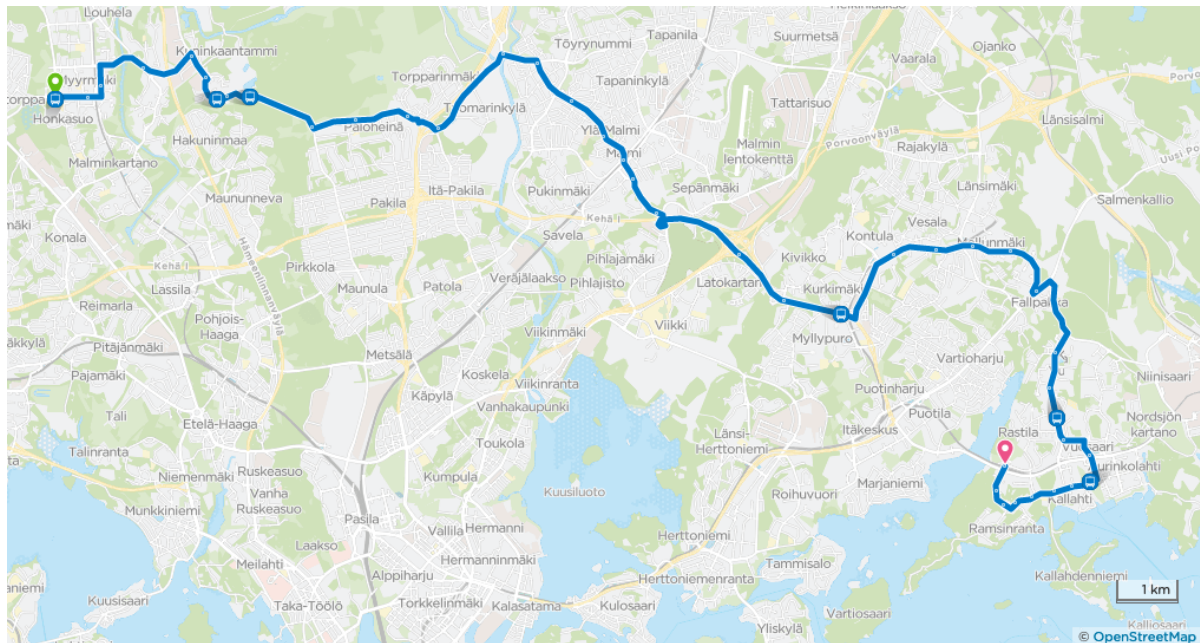


Figure 32. The route of the new, East-West trunk bus route 560, passing through Melunmäki, Malmi and Vuosaari. (HSL, 2020; here provided as a screenshot from the HSL journey planner taken on 30th May 2015, due to an unexpected availability of the line geometry through the HSL API.)

For the car, I am not aware of any road network infrastructure projects of similar scale as the aforementioned PT projects were, with the possible exception of large construction projects and their related traffic arrangements at the central parts of the Tapiola district (see e.g. Tapiola toimii ry, 2020) that potentially help to explain the significant increases in that area. The seasonal variability of speed limits on main roads (see chapter 6.2.2) may also play a part; the changes seen on the data affecting individual areas are possibly more closely related to population changes in those areas (see the chapter 6.1.1) than transport network changes, even if the population changes in the whole area did not statistically have any significant effect on the changes. Thus, despite the statistically more credibly significant difference in means between the car changes, the effects are not spatially as obviously concentrated to any particular districts the way that the PT changes are. Taking into account the lower standard deviations of the car changes (see chapter 5.1.2.3 and table 9), it is plausible to assume that with the exception of relatively minor changes

such as adjusted speed limits, no significant transport network induced changes affected the car travel times.

6.1.3 Regional changes affecting particular industries

My second research question (see the first chapter) was ‘*Are there any significant regional commuting time changes between industries?*’

Between the 2013 and 2015, there was a notable shift between industries, with the share of health and human services decreasing over six percentage points of the total and the share of construction seeing an increase of nearly seven percentage points. However, the changes involve too small groups of people for any reasonable conclusions to be made, with the 100 % increases in the share of construction in certain districts being based on the fact that due to privacy reasons only few industries have been represented in many areas; see the fifth sub-point in the chapter 6.2.1 for further discussion about this error. Thus, based on my research the answer to this question appears to be no. While some industries saw notable in-industry changes in absolute terms, my research did not attempt to explore the spatial variation of those changes.

6.2 Potential sources of errors and missing observations

Quoting Foote & Huebner (2014), I divide the potential sources of error to three categories: obvious error sources, measurement errors, and processing errors.

6.2.1 Obvious error sources

Following the classification of Foote & Huebner (2014), I considered the following obvious error sources potentially relevant in my research:

1. Age of data and temporal density of observations:
 - Ideally, the data acquisition days of the TTMs would have matched those of the commuting data. However, especially with regard to the 2018 TTM, the commuting data was almost 13 months older than the supposedly equivalent TTM data. This leads to a commute-transport network mismatch that reduces the accuracy of my accessibility-based mobility estimates (see chapter 3.3).

- It could also be argued, that ideally the longitudinal resolution of the TTM dataset would match the commuting dataset. Now the availability of commuting data varies; occasionally, it has been updated yearly, but sometimes only every other year. The longitudinal resolution of the TTM dataset has thus far been roughly thirty months.
- Travel surveys were only conducted every sixth year, which meant that I had to utilise the 2012 factors with the 2015 dataset as well. This is also an obvious source of error.

2. Spatial and temporal data coverage

- The relative shallowness of the TTM dataset is an obvious source of error. As the year 2013 in the dataset includes only midday travel times, the dataset lacks longitudinal temporal coverage. However, it is obvious that most commutes would be better represented by rush-hour travel times. Also, integrating cycling data as a third travel mode would make the analysis more extensive, but with regard to cycling data, the dataset has no longitudinal depth at all: it is only included in the 2018 TTM.
- The population dataset lacks subdistrict information in the area of the municipality of Kauniainen; however, at the point of my workflow where I realised that I should integrate a population dataset and weigh the data for population, I had already aggregated the rest of the dataset to a district-level, and adjusting this for Kauniainen would have meant that a tiny part of the dataset would have had to be separately aggregated to municipal level. Due to the complexity of this procedural exception, I chose not to do it, which led to unavailability of population-weighted results in the Kauniainen area. However, as there were no notable outliers in the unweighted results in any district of Kauniainen, this omission most likely did not significantly affect the results.

3. Spatial density of observations

- The most obvious omission in observation density is that the travel mode factors were available only on municipal level, even though it is entirely

plausible to hypothesise that the regional variation between travel modes is spatially much more fine-grained; obvious factors such as car ownership are generally known to vary on sub-regional level as well (see e.g. Brandt & Lindeqvist, 2016; and Laitinen & Vuorio, 2019 for a newspaper article). At least district-level information of travel mode preferences would likely improve accuracy.

- Due to privacy concerns, the commuting data is aggregated to 250 × 250 cells. From the point of view of data accuracy, the availability of precise endpoints of the single commutes would be desirable.
4. Relevance. This is the most important source of error: in lieu of a comprehensive mobility dataset based on e.g. mobile phone data (see the end of the chapter 2.2), I am using an accessibility dataset to estimate the mobility; this is likely to be a relatively crude estimate at best, and to my best knowledge, no research exists about how accurately the travel times of my travel time dataset match reality.
 5. Data accessibility. This is a source of error closely linked to relevance; if the most relevant data is not available due to accessibility reasons, less relevant sources are often used instead. Moreover – As Foote & Huebner (2014) state – ‘accessibility to data is not equal’.

- In the case of my research, my commuting dataset is not freely available, even if the Finnish academic institutions are free to use it in their research. While this is not a source of error for me *per se*, it hinders the repeatability of my research results.

Table 12. Counts of all journeys and IC journeys.

Journey type	Count	% of total*
2013, no IC	442,656	96.9%
2013, with IC	14,223	3.2%
2015, no IC	430,722	97.1%
2015, with IC	13,064	3.0%
2018, no IC	437,954	97.1%
2018, with IC	13,245	3.0%

**) Individually rounded values*

- Also, as stated above, the availability of precise journey endpoints and a precise mobility dataset would have been desirable.
- Due to privacy reasons, only about three percent of all journeys of the commuting dataset are classified by industry (for the exact numbers, see

table 12). It is important to understand that this minority is not a random sample of the total: it is biased towards larger housing units and workplaces, and thus might also create demographic bias; e.g. it may be hypothesised that poorer people are more likely to live in larger housing units and work within large industry complexes. Also, the lack of availability of the IC data left many districts entirely without data.

6.2.2 *Potential errors in original measurements*

There are several potential sources of significant errors in the original measurements:

- Erroneous speed limits affecting the 2015 TTM data. According to the Accessibility Research Group – the research group responsible for maintaining the TTM dataset – there were erroneous speed limit data on certain major streets thoroughfares in the 2015 TTM (Accessibility Research Group & University of Helsinki, 2018).
- Seasonal variability of speed limits not taken into account. In Finland, highway speed limits are lowered during winter months; the 2018 TTM is based on wintertime data, and thus, the calculated times are based on speed limits that are generally lower than on previous TTMs, even if policy changes are ruled out (Accessibility Research Group & University of Helsinki, 2018).
- Uncertainty regarding workplace locations. In some industries, e.g. construction, the actual workplace location generally differs from the registered location of employer’s premises. However, the metadata of the commuting dataset did not clarify how the workplace locations are determined.
- Input data quality of the TTM dataset, not only regarding speed limits, but perhaps regarding some other variables as well. As Tenkanen & Toivonen (2020) wrote: *‘Data used to estimate the congestion levels, or the time that it takes to find a parking space in different parts of the city, or how much time it takes to get and unlock a bike when departing from home, all affect the results.’*

- Potential errors in the commuting dataset. Looking at the individual commutes in the commuting dataset whose endpoints were located at the squares of my own home and workplace, I saw no commutes between the squares on some of the data years, despite being completely sure that both my own and my employer's registered addresses were at those squares; I cannot rule out this being an error in the dataset and cannot rule out other errors of the same type.
- As my work is heavily based aggregation of individual journeys and estimates of their duration, any errors in the source data may propagate to the final results of my analysis.

6.2.3 *Processing errors*

There might be errors arising from processing the data; I have done my utmost to avoid any, but cannot guarantee their existence. During my work, I have made programming and other logical errors, but also some related to statistical analysis. While I believe that I have caught any that would seriously affect the results, it cannot be ruled out that some still exist; indeed, it is ultimately possible that even my fundamental logic of aggregating the journeys should have been different. It is also possible that I have made classification and generalisation errors that have ultimately caused me to misinterpret the results. Ultimately, any certainty against processing errors is only obtainable through the scientific process, when my work is reviewed by others.

7 Conclusions

Big infrastructure projects are a fact of life for those of us who live in larger cities, as the cities evolve and trends of urban planning change. We all have to adjust our daily rhythm to them, and as we do so, our mobility patterns change. When that is aggregated to the level of population of several city districts, the potential effects are large.

My research demonstrates that a modern accessibility dataset such as that described by Tenkanen & Toivonen (2020) can be utilised to evaluate these effects in

advance; even if no dramatic patterns were revealed, the results do reveal changes that demonstrate likely causality between PT infrastructure projects and aggregated commutes.

To this date, the TTM dataset has been upgraded roughly by every two and a half years, with every upgrade increasing its longitudinal depth; the data collection day of hitherto newest upgrade was in January 2018. Should the Accessibility Research Group release a new upgrade during the second half of 2020, the longitudinal depth of rush-hour measurements will grow again, creating a new opportunity to refine my research by comparing the changes between rush-hour data and midday data. Another possibility for improving the accuracy of the results would be utilising travel mode preferences on sub-municipal level, if any are available. An obvious subject for further research would also be to obtain a modern, extensive mobility dataset based on mobile phone data to validate the usage of accessibility datasets in estimating mobility.

With regard to the subset of the commuting dataset that was classified by industry, the small relative size of the dataset combined with its many dimensions and the fact that the subset is not a stratified sample but a highly biased one (see p. 51 at the end of the chapter 6.2.1) makes analysing the dataset more difficult. My research question regarding it was *'are there any significant regional commuting time changes between industries?'*, but in retrospect *are there any significant regional changes affecting the employees of a particular industry?* might had been a better question, and one that might had been more readily answerable by this limited set of data. Taking into account the temporal constraints of my work, I did not consider it feasible to rephrase my research question and rerun my analysis to try to answer that question as well. However, given that no statistically significant patterns were identified by my *t*-tests, the answer to the rephrased question might had been equally inconclusive.

8 Acknowledgements

I want to express my deepest thanks and gratitude to the Associate professor in geoinformatics at the University of Helsinki, Tuuli Toivonen, for supporting me in my motivation crisis and suggesting me this thesis topic; without her great patience, encouragement and invaluable advice back in 2019, I would likely never have even started this study, let alone completed it. I also want to thank my actual supervisor, Petteri Muukkonen and my co-supervisor, Olle Järv, for all the support that they were able to give me despite the pandemic of 2020 upending all normal activities at the university; I should probably have asked them a lot more questions. Special thanks to Petteri for his great patience in this challenging situation.

In addition, I want to thank my employer, the Greens in Finland, for enabling me to have a four-month break from my day-to-day job to complete this thesis and generally having enabled my studies by always being a very flexible employer, even though my chosen academic discipline offers few obvious benefits for them. I am also grateful for the National Geographic Society for the grant they assigned me for my thesis work, even though in the end my actual thesis topic turned out to be completely unrelated to my original field work.

Finally, I want to thank my family and friends for all their emotional support throughout these difficult times, and express my special gratitude to those family members who helped me by proofreading this thesis.

9 References

- Accessibility Research Group, & University of Helsinki. (2016a, January 20).
Pääkaupunkiseudun matka-aikamatriisi 2013. *Saavutettavuuden maantiedettä*.
<https://blogs.helsinki.fi/saavutettavuus/paakaupunkiseudun-matka-aikamatriisi-2013/>
- Accessibility Research Group, & University of Helsinki. (2016b, January 20).
Pääkaupunkiseudun matka-aikamatriisi 2015. *Saavutettavuuden maantiedettä*.
<https://blogs.helsinki.fi/saavutettavuus/paakaupunkiseudun-matka-aikamatriisi-2015/>
- Accessibility Research Group, & University of Helsinki. (2018, June 13).
Pääkaupunkiseudun matka-aikamatriisi 2018. *Saavutettavuuden maantiedettä*.
<https://blogs.helsinki.fi/saavutettavuus/paakaupunkiseudun-matka-aikamatriisi-2018/>
- Ahas, R., Aasa, A., Roose, A., Mark, Ü., & Silm, S. (2008). Evaluating passive mobile positioning data for tourism surveys: An Estonian case study. *Tourism Management*, 29(3), 469–486. <https://doi.org/10.1016/j.tourman.2007.05.014>
- Aluesarjat. (2020). *Helsingin seudun aluesarjat tilastokanta ja Tilastokeskus*. Aluesarjat.
<http://www.aluesarjat.fi/>
- Barbosa, H., Barthelemy, M., Ghoshal, G., James, C. R., Lenormand, M., Louail, T., Menezes, R., Ramasco, J. J., Simini, F., & Tomasini, M. (2018). Human mobility: Models and applications. *Physics Reports*, 734, 1–74. <https://doi.org/10.1016/j.physrep.2018.01.001>
- Bertolini, L., le Clercq, F., & Kapoen, L. (2005). Sustainable accessibility: A conceptual framework to integrate transport and land use plan-making. Two test-applications in the Netherlands and a reflection on the way forward. *Transport Policy*, 12(3), 207–220. <https://doi.org/10.1016/j.tranpol.2005.01.006>

- Brandt, E., Kantele, S., & Rätty, P. (2019). *Liikkumistottumukset Helsingin seudulla 2018* (9/2019; HSL:n julkaisu, p. 172). HSL Helsingin seudun liikenne. https://www.hsl.fi/sites/default/files/hsl_julkaisu_9_2019_netti.pdf
- Brandt, E., & Lindeqvist, M. (2016). *Auton omistus Helsingin seudulla – katsausmenneeseen kehitykseen ja pohdintoja tulevasta* (16/2016; HSL:n julkaisu, p. 90). HSL Helsingin seudun liikenne. https://www.hsl.fi/sites/default/files/19_2016_auton_omistus_helsingin_seudulla.pdf
- Brown, M. B., & Forsythe, A. B. (1974). Robust Tests for the Equality of Variances. *Journal of the American Statistical Association*, 69(346), 364–367. JSTOR. <https://doi.org/10.2307/2285659>
- Burns, L. D., & Golob, T. F. (1976). The role of accessibility in basic transportation choice behavior. *Transportation*, 5(2), 175–198. <https://doi.org/10.1007/BF00167272>
- Candia, J., González, M. C., Wang, P., Schoenharl, T., Madey, G., & Barabási, A.-L. (2008). Uncovering individual and collective human dynamics from mobile phone records. *Journal of Physics A: Mathematical and Theoretical*, 41(22), 224015. <https://doi.org/10.1088/1751-8113/41/22/224015>
- City of Helsinki. (2018). *Statistical Yearbook of Helsinki 2018*. City of Helsinki, Executive Office, Urban Research and Statistics.
- City of Helsinki. (2020, February 25). *Avoimet paikkatiedot*. Helsingin kaupunki. <https://www.hel.fi/helsinki/fi/kartat-ja-liikenne/kartat-ja-paikkatieto/Paikkatiedot+ja+aineistot/avoimet+paikkatiedot/>
- City of Helsinki, City of Espoo, City of Vantaa, & Statistics Finland. (2019, October 30). *Population projection in the Helsinki metropolitan area by district*. https://hri.fi/data/en_GB/dataset/paakaupunkiseudun-vaestoennuste
- Cui, B., Boisjoly, G., El-Geneidy, A., & Levinson, D. (2019). Accessibility and the journey to work through the lens of equity. *Journal of Transport Geography*, 74, 269–277. <https://doi.org/10.1016/j.jtrangeo.2018.12.003>

- Dalvi, M. Q., & Martin, K. M. (1976). The measurement of accessibility: Some preliminary results. *Transportation*, 5(1), 17–42. <https://doi.org/10.1007/BF00165245>
- Delafontaine, M., Neutens, T., Schwanen, T., & Weghe, N. V. de. (2011). The impact of opening hours on the equity of individual space–time accessibility. *Computers, Environment and Urban Systems*, 35(4), 276–288. <https://doi.org/10.1016/j.compenvurbsys.2011.02.005>
- Evans, G. W., Wener, R. E., & Phillips, D. (2002). The Morning Rush Hour: Predictability and Commuter Stress. *Environment and Behavior*, 34(4), 521–530. <https://doi.org/10.1177/00116502034004007>
- Foote, K. E., & Huebner, D. J. (2014). *Error, Accuracy, and Precision*. The Geographer’s Craft Project, Department of Geography, The University of Colorado; archive.org. https://web.archive.org/web/20171009234017fw_/http://www.colorado.edu/geography/gcraft/notes/error/error.html
- Geurs, K. T., & Östh, J. (2016). Advances in the Measurement of Transport Impedance in Accessibility Modelling. *European Journal of Transport and Infrastructure Research*, 16(2), Article 2. <https://doi.org/10.18757/ejtir.2016.16.2.3138>
- Geurs, K. T., & van Wee, B. (2004). Accessibility evaluation of land-use and transport strategies: Review and research directions. *Journal of Transport Geography*, 12(2), 127–140. <https://doi.org/10.1016/j.jtrangeo.2003.10.005>
- Ghasemi, A., & Zahediasl, S. (2012). Normality Tests for Statistical Analysis: A Guide for Non-Statisticians. *International Journal of Endocrinology and Metabolism*, 10(2), 486–489. <https://doi.org/10.5812/ijem.3505>
- Giuliano, G. (1991). *Is Jobs-Housing Balance a Transportation Issue?* <https://escholarship.org/uc/item/4874r4hg>
- Giuliano, G., & Small, K. A. (1993). Is the Journey to Work Explained by Urban Structure?: *Urban Studies*. <https://doi.org/10.1080/00420989320081461>
- González, M. C., Hidalgo, C. A., & Barabási, A.-L. (2008). Understanding individual human mobility patterns. *Nature*, 453(7196), 779–782. <https://doi.org/10.1038/nature06958>

- Gordon, P., Richardson, H. W., & Jun, M.-J. (1991). The Commuting Paradox Evidence from the Top Twenty. *Journal of the American Planning Association*, 57(4), 416–420. <https://doi.org/10.1080/01944369108975516>
- Gould, P. R. (1969). *Spatial Diffusion, Resource Paper No. 4*. <http://eric.ed.gov/?q=resource+paper+no.+4&id=ED120029>
- Hägerstrand, T. (1970). What about people in Regional Science? *Papers of the Regional Science Association*, 24(1), 6–21. <https://doi.org/10.1007/BF01936872>
- Hamilton, B. W., & Röell, A. (1982). Wasteful Commuting. *Journal of Political Economy*, 90(5), 1035–1053. JSTOR.
- Hansen, W. G. (1959). How Accessibility Shapes Land Use. *Journal of the American Institute of Planners*, 25(2), 73–76. <https://doi.org/10.1080/01944365908978307>
- Hawelka, B., Sitko, I., Beinat, E., Sobolevsky, S., Kazakopoulos, P., & Ratti, C. (2014). Geo-located Twitter as proxy for global mobility patterns. *Cartography and Geographic Information Science*, 41(3), 260–271. <https://doi.org/10.1080/15230406.2014.890072>
- Hedberg, C. (2005). *Geografiska perspektiv på arbetsmarknadsrörlighet*. Arbetslivsinstitutet & CIND. <http://urn.kb.se/resolve?urn=urn:nbn:se:uu:diva-78422>
- Heikinheimo, V., Minin, E. D., Tenkanen, H., Hausmann, A., Erkkonen, J., & Toivonen, T. (2017). User-Generated Geographic Information for Visitor Monitoring in a National Park: A Comparison of Social Media Data and Visitor Survey. *ISPRS International Journal of Geo-Information*, 6(3), 85. <https://doi.org/10.3390/ijgi6030085>
- Helminen, V., & Ristimäki, M. (2007). Relationships between commuting distance, frequency and telework in Finland. *Journal of Transport Geography*, 15(5), 331–342. <https://doi.org/10.1016/j.jtrangeo.2006.12.004>
- Helminen, V., Ristimäki, M., & Oinonen, K. (2003). *Etätyö ja työmatkat Suomessa*. <https://helda.helsinki.fi/handle/10138/40503>

- HSL. (2015, August 6). *Runkolinja 560 aloittaa liikennöintinsä 10.8.—Reitillä myös uusi 1,2 kilometrin joukkoliikennetunneli*. HSL. <https://www.hsl.fi/uutiset/2015/runkolinja-560-aloittaa-liikennointinsa-108-reitilla-myos-uusi-12-kilometrin>
- HSL. (2019). *Metro services from 18 November 2017*. HSL. <https://www.hsl.fi/en/west-metro/metro>
- HSL. (2020). *A screenshot of the journey planner*. [Map]. HSL. <https://reittiopas.hsl.fi/>. Accessed 30/05/2020.
- Ingram, D. R. (1971). The concept of accessibility: A search for an operational form. *Regional Studies*, 5(2), 101–107. <https://doi.org/10.1080/09595237100185131>
- Järv, O., Ahas, R., & Witlox, F. (2014). Understanding monthly variability in human activity spaces: A twelve-month study using mobile phone call detail records. *Transportation Research Part C: Emerging Technologies*, 38, 122–135. <https://doi.org/10.1016/j.trc.2013.11.003>
- Järv, O., Tenkanen, H., & Toivonen, T. (2017). Enhancing spatial accuracy of mobile phone data using multi-temporal dasymetric interpolation. *International Journal of Geographical Information Science*, 31(8), 1630–1651. <https://doi.org/10.1080/13658816.2017.1287369>
- Kaufmann, V., Bergman, M. M., & Joye, D. (2004). Motility: Mobility as capital. *International Journal of Urban and Regional Research*, 28(4), 745–756. <https://doi.org/10.1111/j.0309-1317.2004.00549.x>
- Kenworthy, J. R., & Laube, F. B. (1996). Automobile dependence in cities: An international comparison of urban transport and land use patterns with implications for sustainability. *Environmental Impact Assessment Review*, 16(4), 279–308. [https://doi.org/10.1016/S0195-9255\(96\)00023-6](https://doi.org/10.1016/S0195-9255(96)00023-6)
- Kiiskilä, K., Tuominen, J., Frösén, N., Valli, R., & Herneoja, A. (2017). *Kehäradan liikenteelliset vaikutukset* (6/2017; HSL:n julkaisuja, p. 62). https://www.hsl.fi/sites/default/files/6_2017_keharadan_liikenteelliset_vaikutukset.pdf

- Kluger, A. N. (1998). Commute variability and strain. *Journal of Organizational Behavior*, 19(2), 147–165. [https://doi.org/10.1002/\(SICI\)1099-1379\(199803\)19:2<147::AID-JOB830>3.0.CO;2-Y](https://doi.org/10.1002/(SICI)1099-1379(199803)19:2<147::AID-JOB830>3.0.CO;2-Y)
- Laitinen, J., & Vuorio, J. (2019, March 18). Kruununhaassa asuva Pekka Hietaniemi karsastaa joukkoliikennettä ja omistaa niin BMW:n katumaasturin kuin avomallisen Audinkin – Helsingissä se on jo kapinallista. *Helsingin Sanomat*. <https://www.hs.fi/kaupunki/art-2000006038774.html>
- Levene, H. (1960). Robust tests for equality of variances. In I. Olkin (Ed.), *Contributions to Probability and Statistics* (pp. 278–292). Scopus.
- Levinson, D. M. (1998). Accessibility and the journey to work. *Journal of Transport Geography*, 6(1), 11–21. [https://doi.org/10.1016/S0966-6923\(97\)00036-7](https://doi.org/10.1016/S0966-6923(97)00036-7)
- Levinson, D. M., & Kumar, A. (1994). *The Rational Locator: Why Travel Times Have Remained Stable*. <http://dx.doi.org/10.1080/01944369408975590>
- Levinson, D., & Wu, Y. (2005). The rational locator reexamined: Are travel times still stable? *Transportation*, 32(2), 187–202. <https://doi.org/10.1007/s11116-004-5507-4>
- Lintunen, P. (2000). *Työmatkat ja työpaikkaomavaraisuus*. LYYLI-raporttisarja. Liikenneministeriö.
- Lyons, G., & Chatterjee, K. (2008). A Human Perspective on the Daily Commute: Costs, Benefits and Trade-offs. *Transport Reviews*, 28(2), 181–198. <https://doi.org/10.1080/01441640701559484>
- Mills, E. S. (1972). *Studies in the Structure of the Urban Economy*. The Johns Hopkins Press, Baltimore, Maryland 21218 (\$7. <https://trid.trb.org/view/128838>
- Muth, R. F. (1969). *Cities and Housing: The Spatial Pattern of Urban Residential Land Use*. Chicago U.P.
- Nervola, A. (2016). *Jätkäsaaren ajoaikamittaukset*. Kaupunkisuunnitteluvirasto, Helsingin kaupunki. https://www.uuttahelsinki.fi/sites/default/files/inline-attachments/2017-04/jatkasaaren_ajoajat_muistio.pdf

- Neutens, T. (2015). Accessibility, equity and health care: Review and research directions for transport geographers. *Journal of Transport Geography*, 43, 14–27. <https://doi.org/10.1016/j.jtrangeo.2014.12.006>
- Neutens, T., Schwanen, T., Witlox, F., & Maeyer, P. D. (2010). Equity of Urban Service Delivery: A Comparison of Different Accessibility Measures: *Environment and Planning A*. <https://doi.org/10.1068/a4230>
- NIST. (2012, April). *NIST/SEMATECH e-Handbook of Statistical Methods*. <https://www.itl.nist.gov/div898/handbook/index.htm>
- Paananen, V. (2020, January 14). Ruoholahteen suunniteltu rekkaramppi sai hylkäystuomion kaupunkiympäristölautakunnalta. *Helsingin Sanomat*. <https://www.hs.fi/kaupunki/art-2000006372827.html>
- Palmer, J. R. B., Espenshade, T. J., Bartumeus, F., Chung, C. Y., Ozgencil, N. E., & Li, K. (2013). New Approaches to Human Mobility: Using Mobile Phones for Demographic Research. *Demography*, 50(3), 1105–1128. <https://doi.org/10.1007/s13524-012-0175-z>
- Priemus, H., Nijkamp, P., & Banister, D. (2001). Mobility and spatial dynamics: An uneasy relationship. *Journal of Transport Geography*, 9(3), 167–171. [https://doi.org/10.1016/S0966-6923\(01\)00007-2](https://doi.org/10.1016/S0966-6923(01)00007-2)
- Ruths, D., & Pfeffer, J. (2014). Social media for large studies of behavior. *Science*, 346(6213), 1063–1064. <https://doi.org/10.1126/science.346.6213.1063>
- Salonen, M. (2014). *Analysing spatial accessibility patterns with travel time and distance measures: Novel approaches for rural and urban contexts* [Thesis]. <https://helda.helsinki.fi/handle/10138/135999>
- Salonen, M., & Toivonen, T. (2013). Modelling travel time in urban networks: Comparable measures for private car and public transport. *Journal of Transport Geography*, 31, 143–153. <https://doi.org/10.1016/j.jtrangeo.2013.06.011>
- Salonen, M., Toivonen, T., Cohalan, J.-M., & Coomes, O. T. (2012). Critical distances: Comparing measures of spatial accessibility in the riverine landscapes of Per-

- uvian Amazonia. *Applied Geography*, 32(2), 501–513. <https://doi.org/10.1016/j.apgeog.2011.06.017>
- Sadow, E., & Westin, K. (2010). The persevering commuter – Duration of long-distance commuting. *Transportation Research Part A: Policy and Practice*, 44(6), 433–445. <https://doi.org/10.1016/j.tra.2010.03.017>
- Shapiro, S. S., & Wilk, M. B. (1965). An Analysis of Variance Test for Normality (Complete Samples). *Biometrika*, 52(3/4), 591–611. JSTOR. <https://doi.org/10.2307/2333709>
- Sheller, M., & Urry, J. (2006). The New Mobilities Paradigm: *Environment and Planning A*, 38(2), 207–226. <https://doi.org/10.1068/a37268>
- Small, K. A., & Song, S. (1992). ‘Wasteful’ Commuting: A Resolution. *Journal of Political Economy*, 100(4), 888–898. JSTOR.
- SOU. (2003). *Geografisk rörlighet för sysselsättning och tillväxt* (Text No. 37; Statens offentliga utredningar från Arbetsmarknadsdepartementet, p. 150). Regeringskansliet. <https://www.regeringen.se/rattsliga-dokument/statens-offentliga-utredningar/2003/04/sou-200337/>
- Student. (1908). The Probable Error of a Mean. *Biometrika*, 6(1), 1–25. JSTOR. <https://doi.org/10.2307/2331554>
- SYKE. (2016a, July 21). *Työmatkan keskipituus—Ymparisto.fi*. [https://www.ymparisto.fi/fi-FI/Kartat_ja_tilastot/Ympariston_tilan_indikaattorit/Yhdyskuntarakenne/Työmatkan_keskipituus_kasvanut_14_kilome\(28635\)](https://www.ymparisto.fi/fi-FI/Kartat_ja_tilastot/Ympariston_tilan_indikaattorit/Yhdyskuntarakenne/Työmatkan_keskipituus_kasvanut_14_kilome(28635))
- SYKE. (2016b, December 7). *Ymparisto > Yhdyskuntarakenteen seurannan aineistot*. https://www.ymparisto.fi/fi-FI/Elinymparisto_ja_kaavoitus/Yhdyskuntarakenne/Tietoa_yhdyskuntarakenteesta/Yhdyskuntarakenteen_seurannan_aineistot
- SYKE. (2019). *Yhdyskuntarakenteen seurannan aineistot (YKR)*. Suomen ympäristökeskus.
- Tapiola toimii ry. (2020). *AINOA -kokonaisuus*. <https://www.tapiolankeskus.fi/fi/Tapiola-uudistuu/AINOA-kokonaisuus>. Accessed 30/05/2020.

- Tenkanen, H., Minin, E. D., Heikinheimo, V., Hausmann, A., Herbst, M., Kajala, L., & Toivonen, T. (2017). Instagram, Flickr, or Twitter: Assessing the usability of social media data for visitor monitoring in protected areas. *Scientific Reports*, 7(1), 1–11. <https://doi.org/10.1038/s41598-017-18007-4>
- Tenkanen, H., Saarsalmi, P., Järvi, O., Salonen, M., & Toivonen, T. (2016). Health research needs more comprehensive accessibility measures: Integrating time and transport modes from open data. *International Journal of Health Geographics*, 15(1), 23. <https://doi.org/10.1186/s12942-016-0052-x>
- Tenkanen, H., & Toivonen, T. (2020). Longitudinal spatial dataset on travel times and distances by different travel modes in Helsinki Region. *Scientific Data*, 7(1), 1–15. <https://doi.org/10.1038/s41597-020-0413-y>
- United Nations. (2014). *United Nations, Department of Economic and Social Affairs, Population Division (2014). World Urbanization Prospects: The 2014 Revision, CD-ROM Edition*. United Nations, Department of Economic and social affairs, Population Division.
- University of Oxford. (2014). *Population density by city*. Our World in Data. <https://ourworldindata.org/grapher/population-density-by-city>
- Wegener, M., & Fuerst, F. (2004). *Land-Use Transport Interaction: State of the Art* (SSRN Scholarly Paper ID 1434678). Social Science Research Network. <https://doi.org/10.2139/ssrn.1434678>

Appendix 1: Description of the research database and code repositories

To prepare the datasets for the analysis I designed a PostgreSQL database with spatial extensions and imported the source data to that database. Then, I wrote a Python application that aggregates the data. The code repository that includes the application code and instructions for recreating my research DB are publicly available at GitHub, at address <https://github.com/pinjaliina/CommuteAggregator>.

Furthermore, I also created another public code repository,

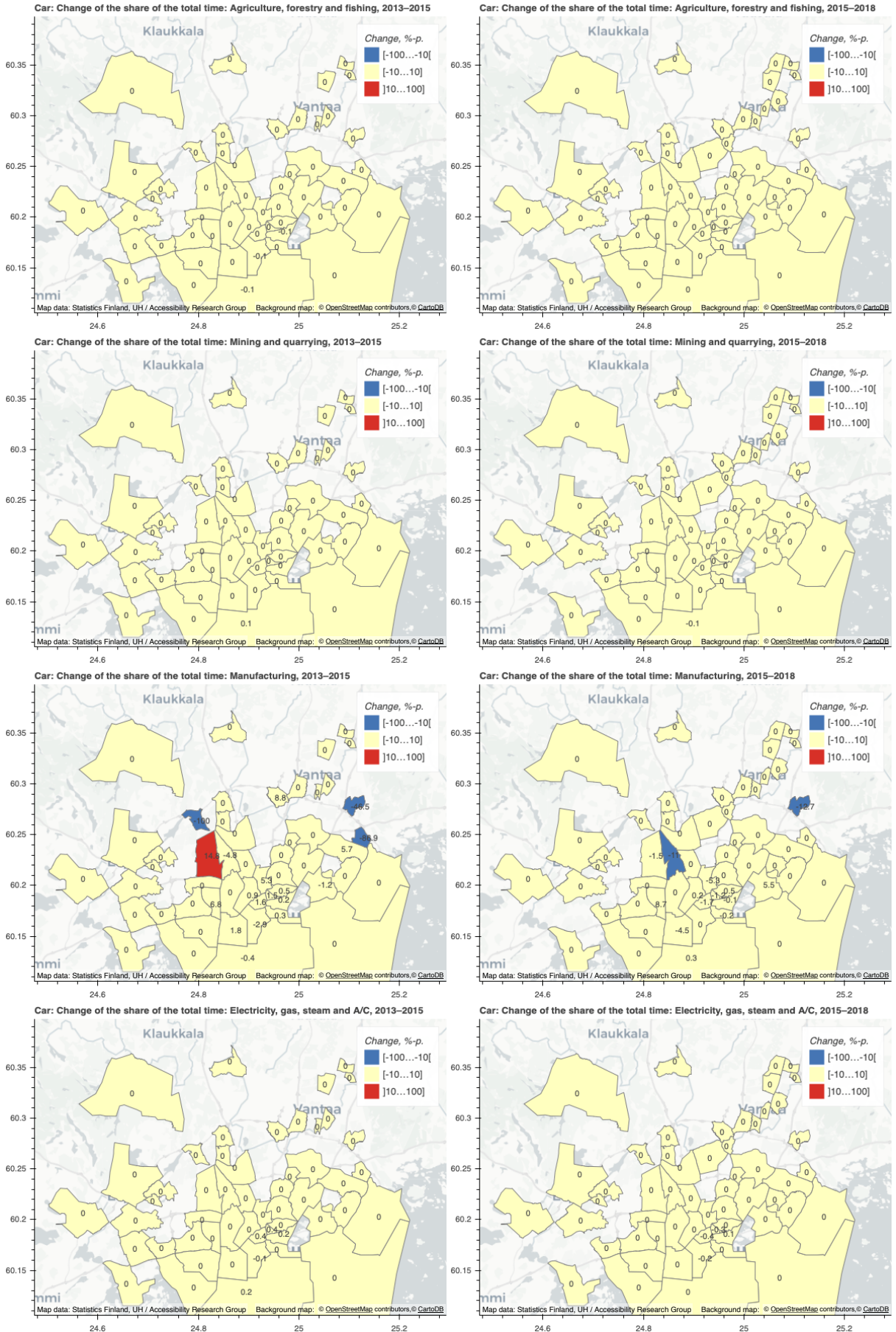
<https://github.com/pinjaliina/AnalyseCommutes>, which includes scripts for:

- further preparing the data for analysis,
- creating attribute tables of the data for QGIS use,
- statistically assessing the QGIS attribute tables and
- plotting and statistically assessing the IC data.

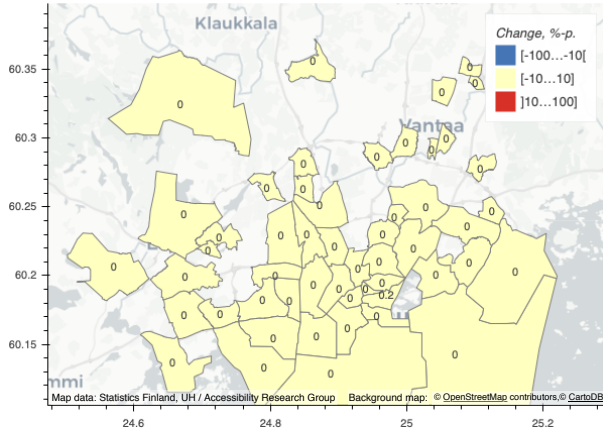
In addition, this repository also contains the attribute tables exported from QGIS for further analysis and previously created IC data maps and statistical output files (the contents of the table 11 of this thesis plus the contents of the other appendices). In addition, the maps of appendix 2 are available interactively through this repository.

With the help of these repositories, a person with a copy of the non-public SSUF dataset should be able to reproduce my work.

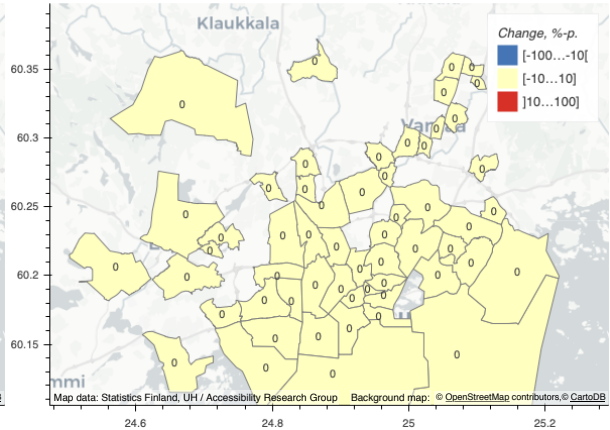
Appendix 2: IC frequency maps



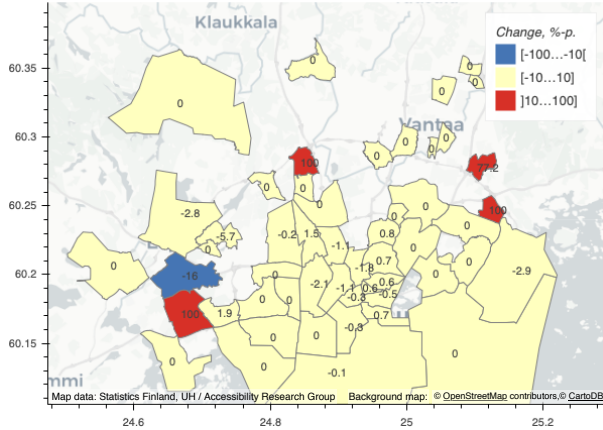
Car: Change of the share of the total time: Water supply, sewage and environment, 2013-21



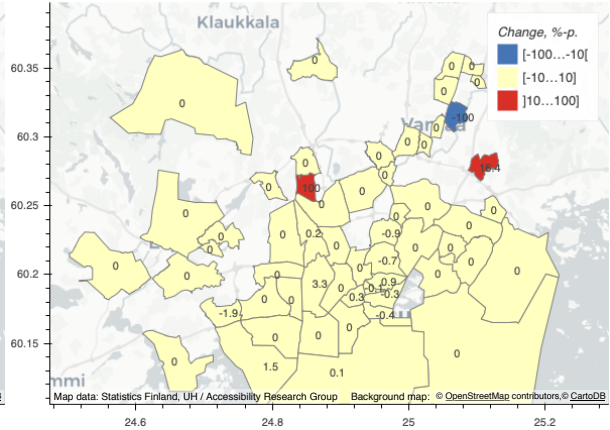
Car: Change of the share of the total time: Water supply, sewage and environment, 2015-21



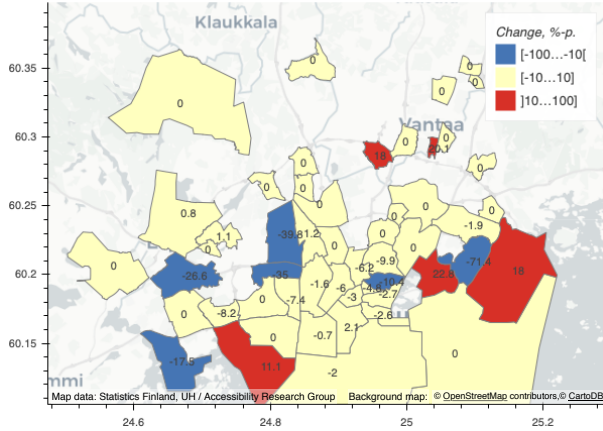
Car: Change of the share of the total time: Construction, 2013-2015



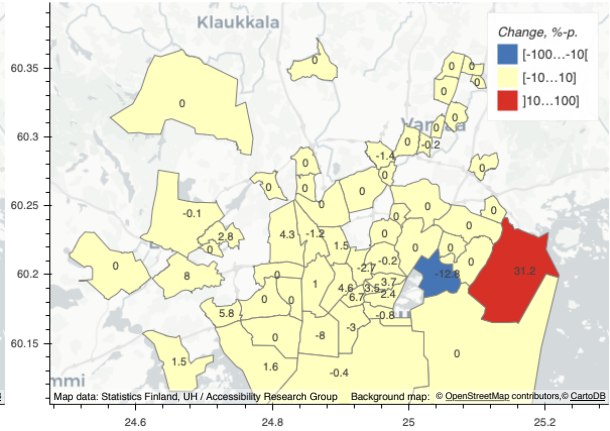
Car: Change of the share of the total time: Construction, 2015-2018



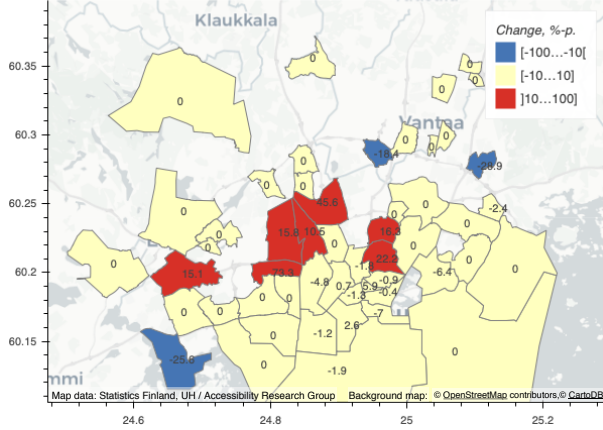
Car: Change of the share of the total time: Wholesale and retail trade, 2013-2015



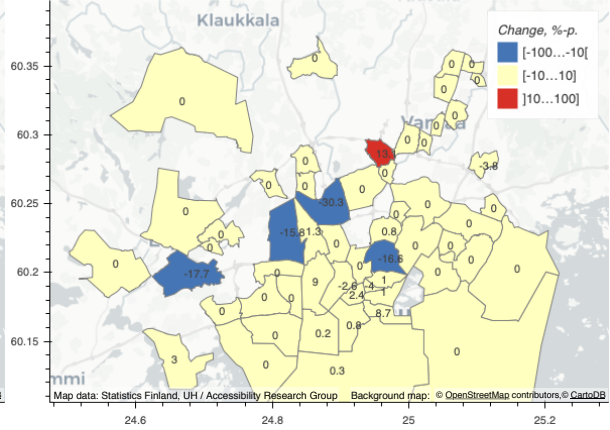
Car: Change of the share of the total time: Wholesale and retail trade, 2015-2018

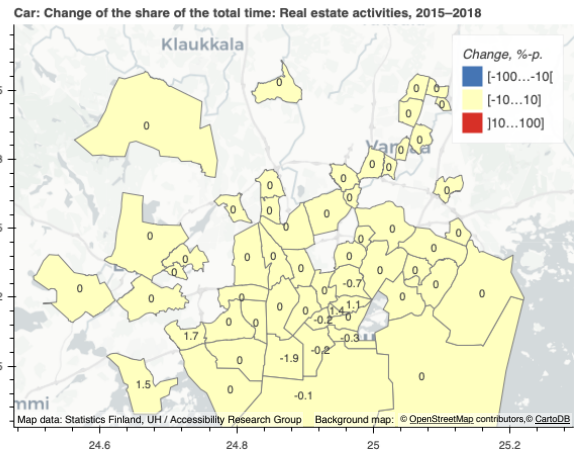
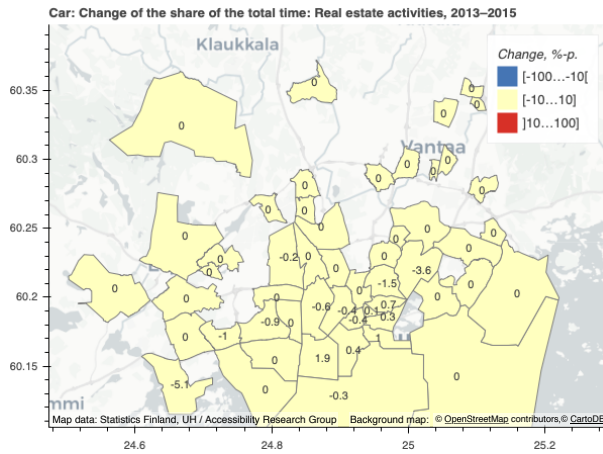
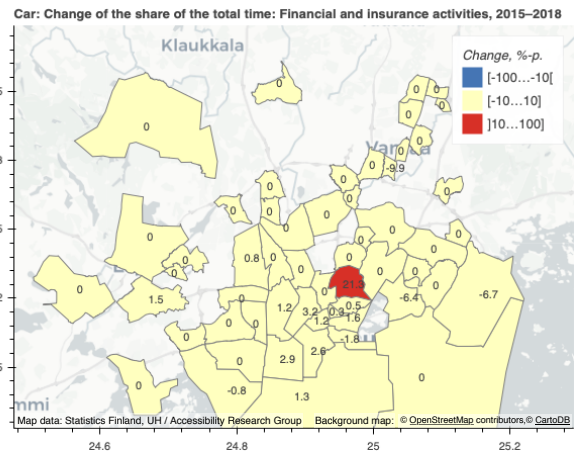
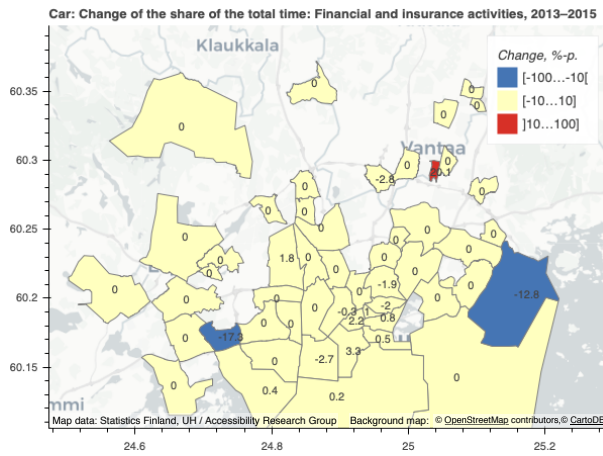
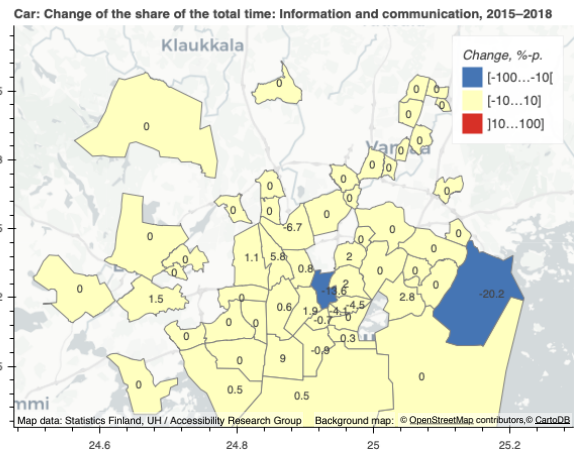
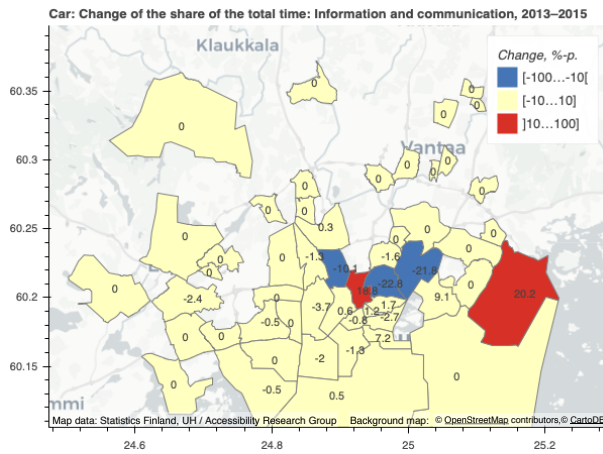
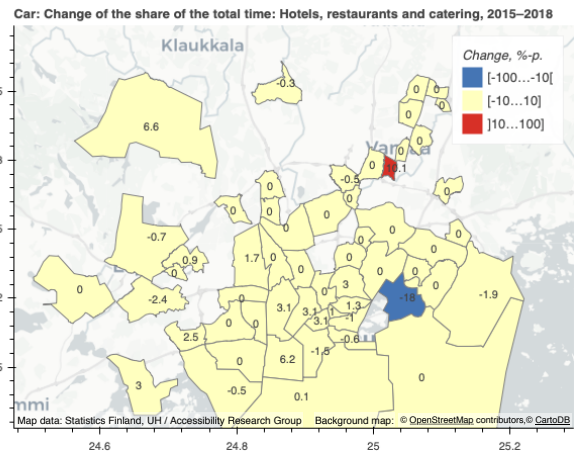
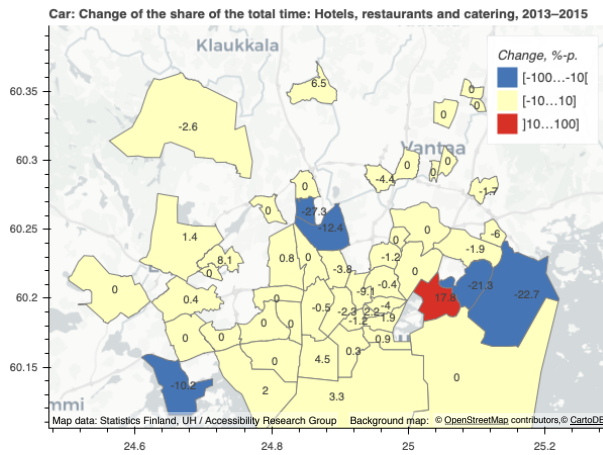


Car: Change of the share of the total time: Transportation and storage, 2013-2015

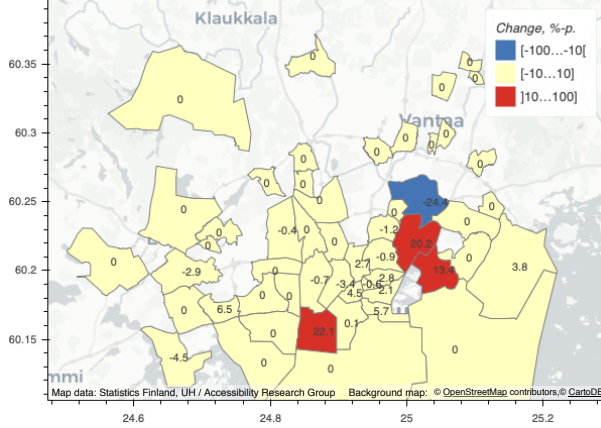


Car: Change of the share of the total time: Transportation and storage, 2015-2018

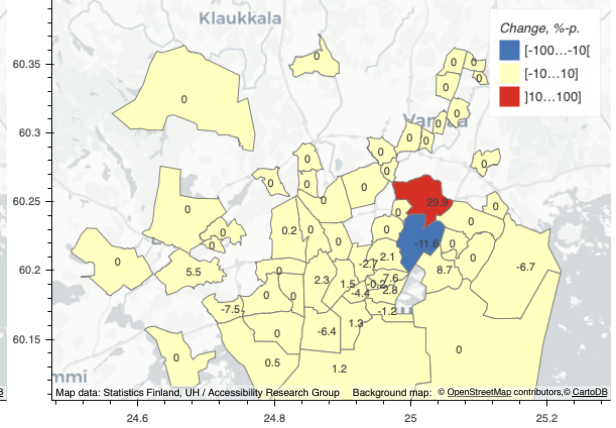




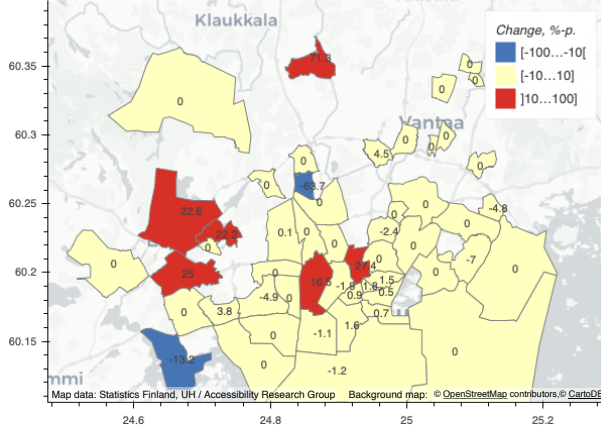
Car: Change of the share of the total time: Speciality professions, 2013–2015



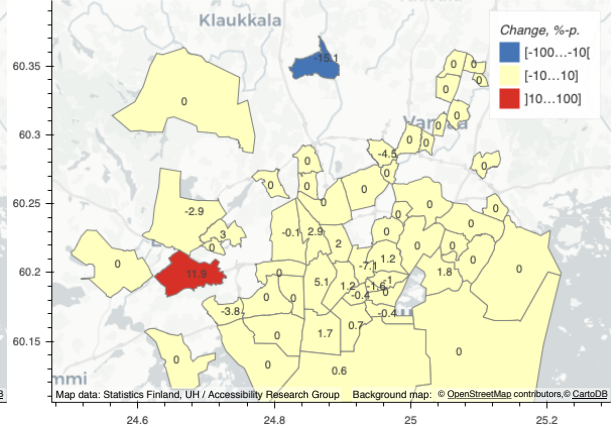
Car: Change of the share of the total time: Speciality professions, 2015–2018



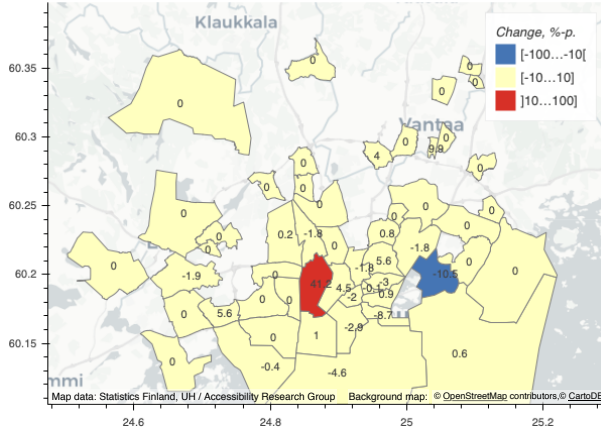
Car: Change of the share of the total time: Administrative and support services, 2013–2015



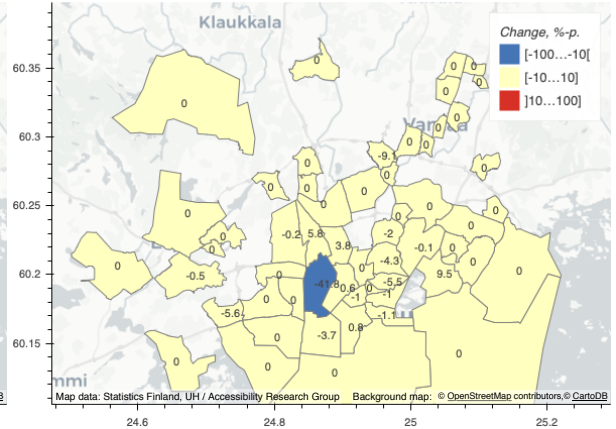
Car: Change of the share of the total time: Administrative and support services, 2015–2018



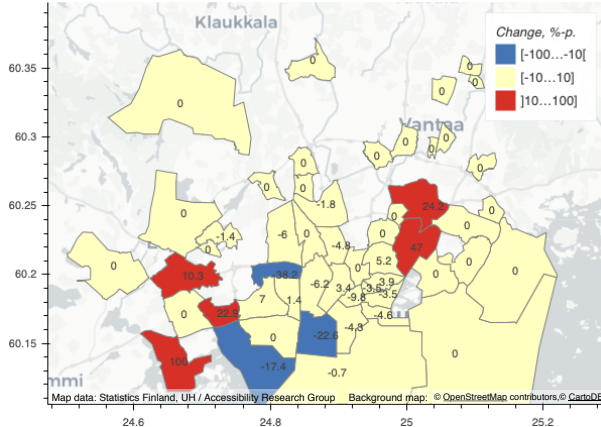
Car: Change of the share of the total time: Public administration and defence, 2013–2015



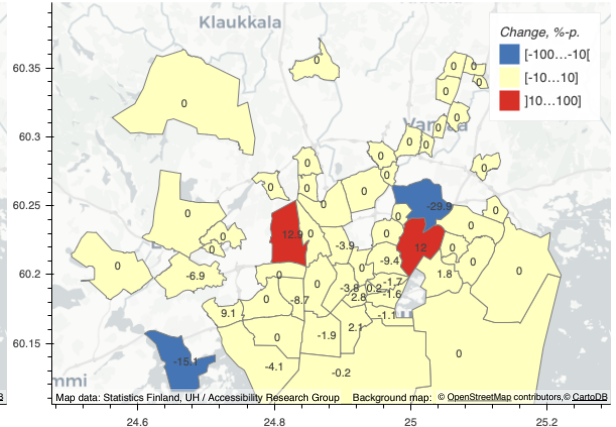
Car: Change of the share of the total time: Public administration and defence, 2015–2018

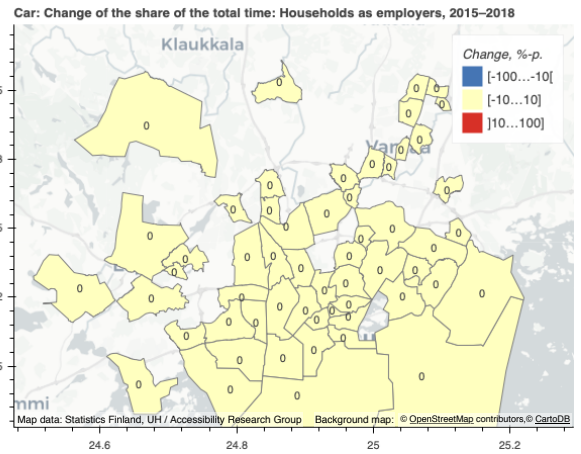
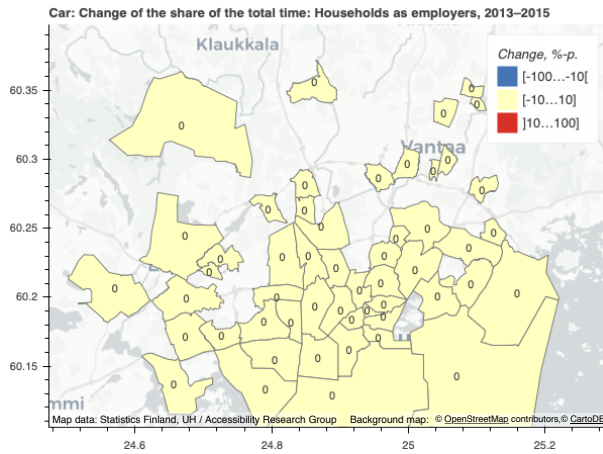
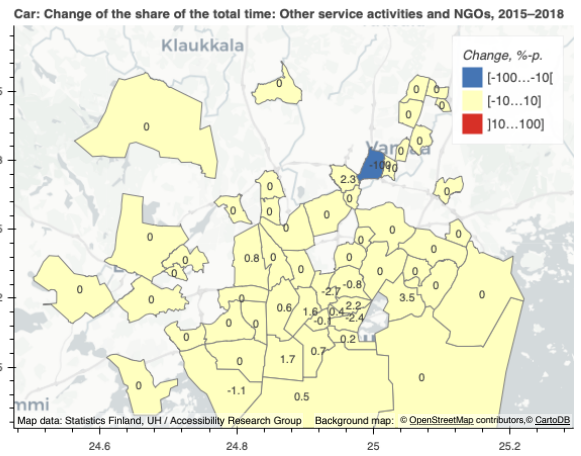
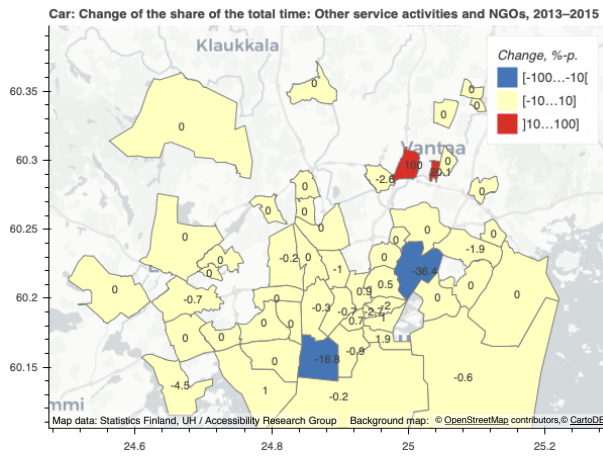
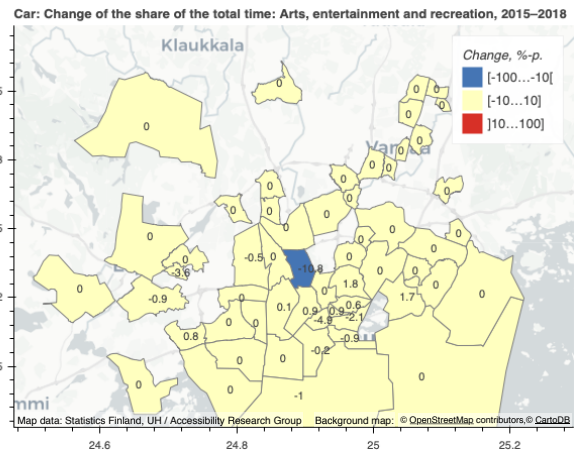
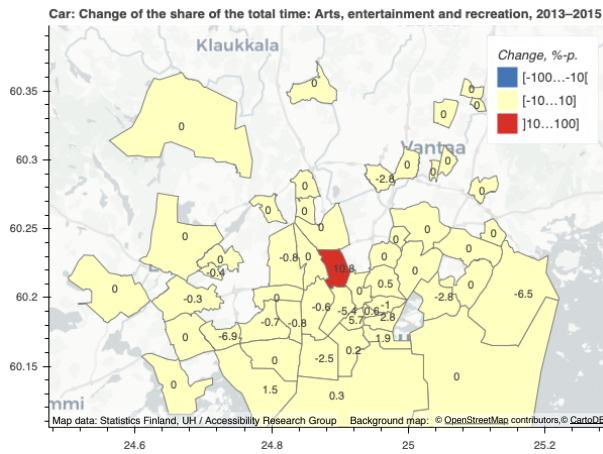
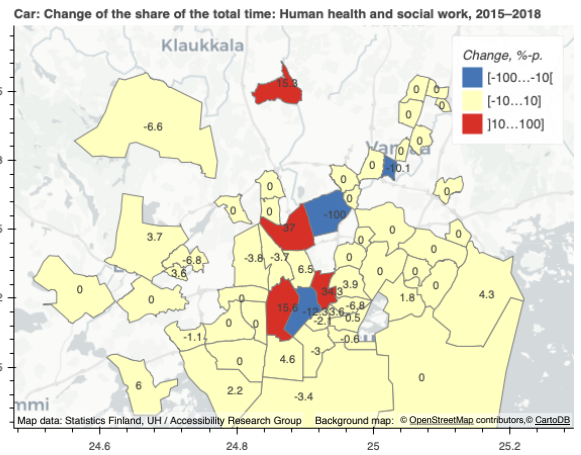
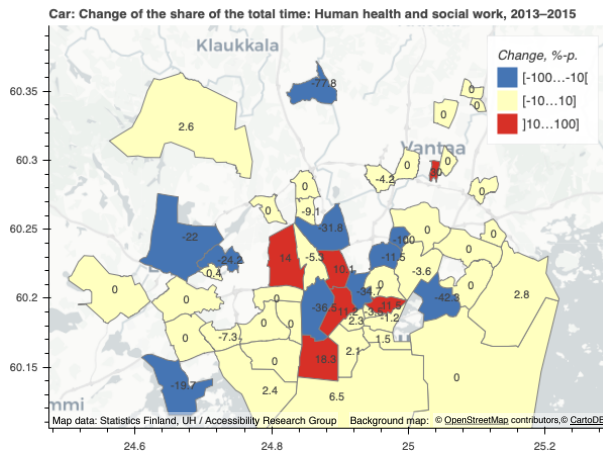


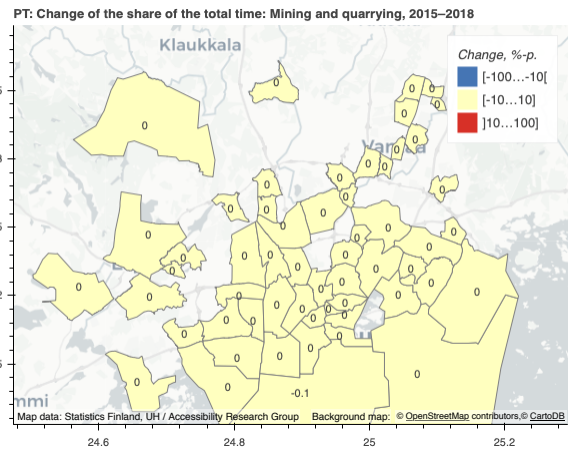
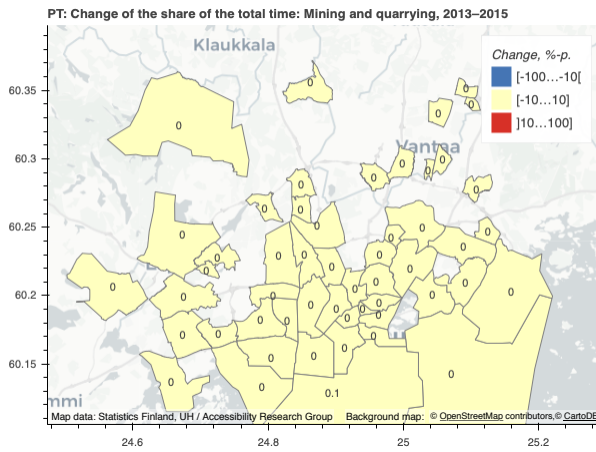
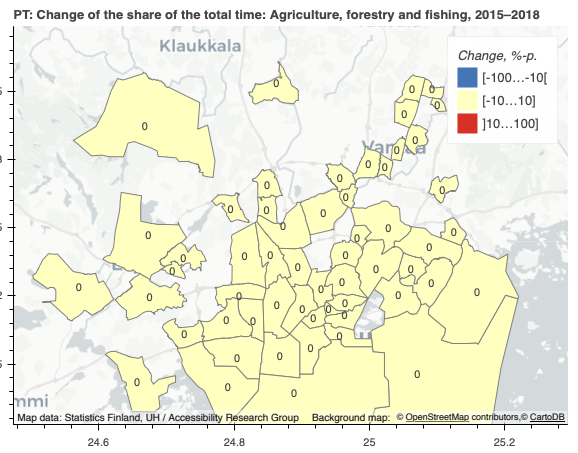
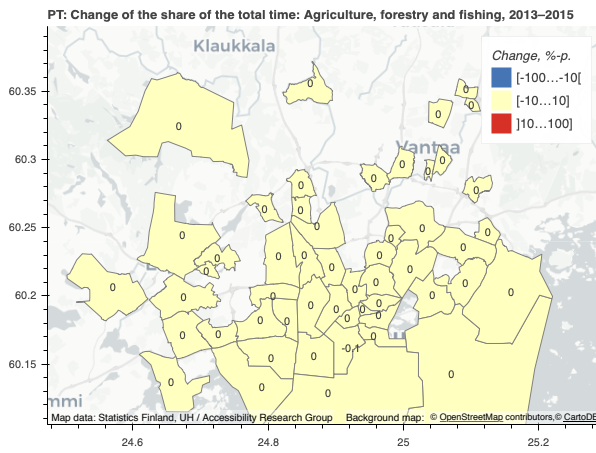
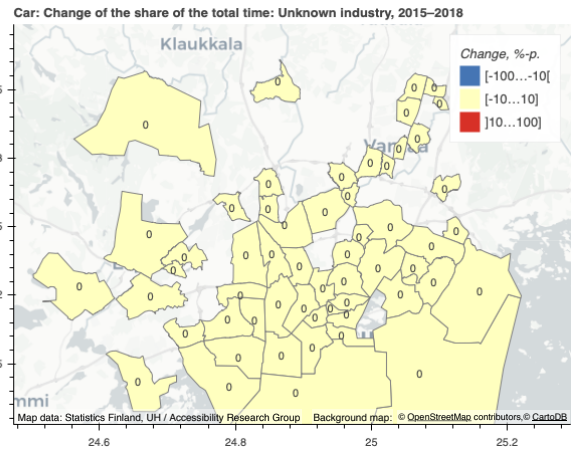
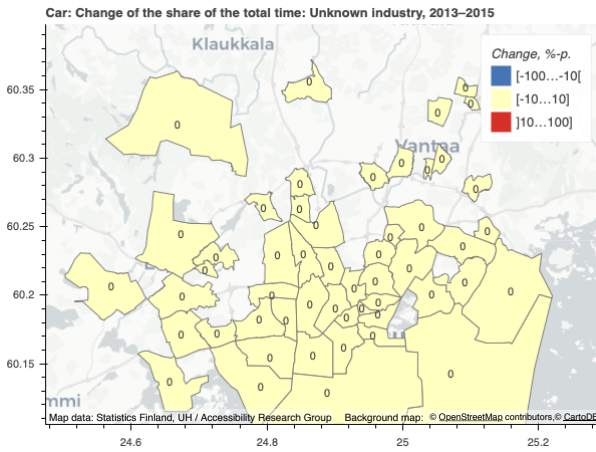
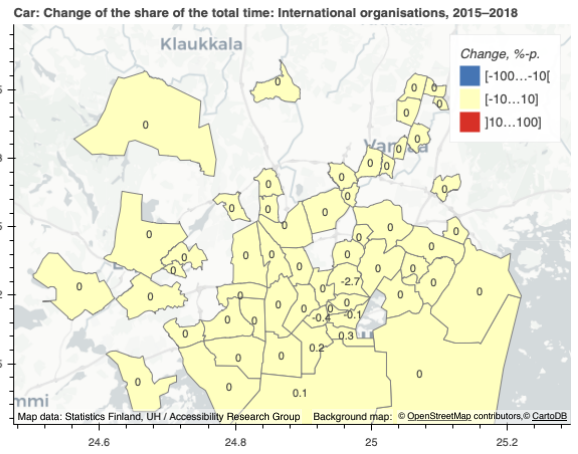
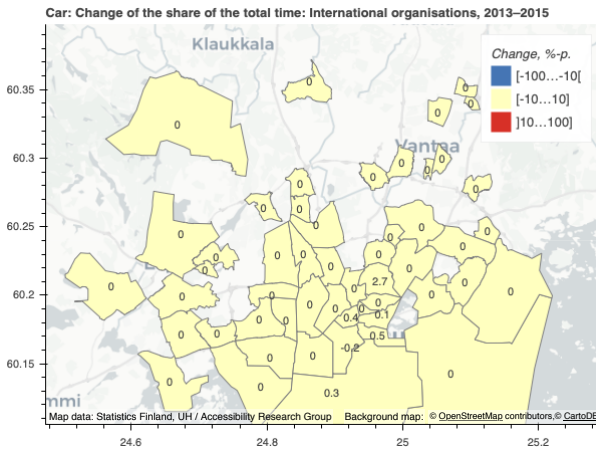
Car: Change of the share of the total time: Education, 2013–2015

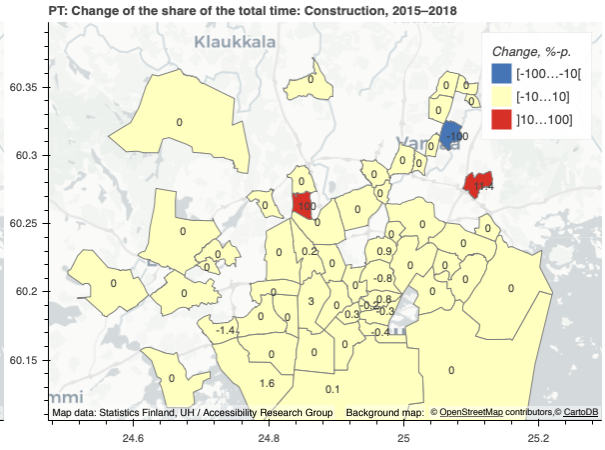
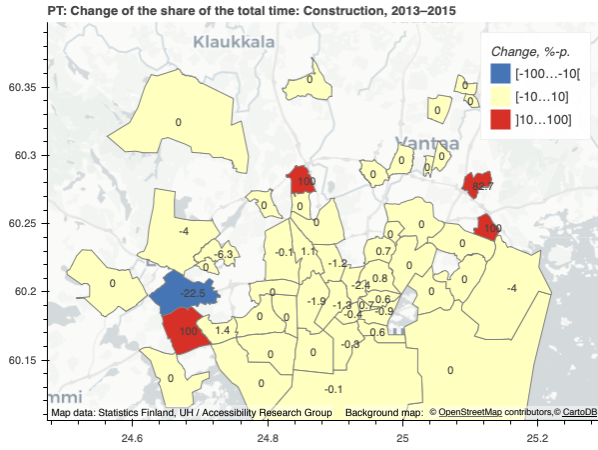
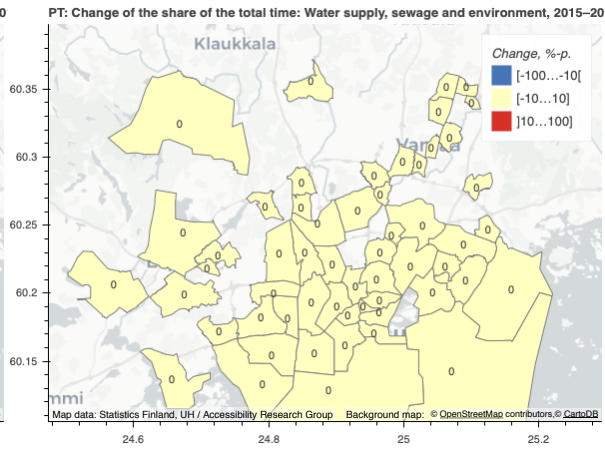
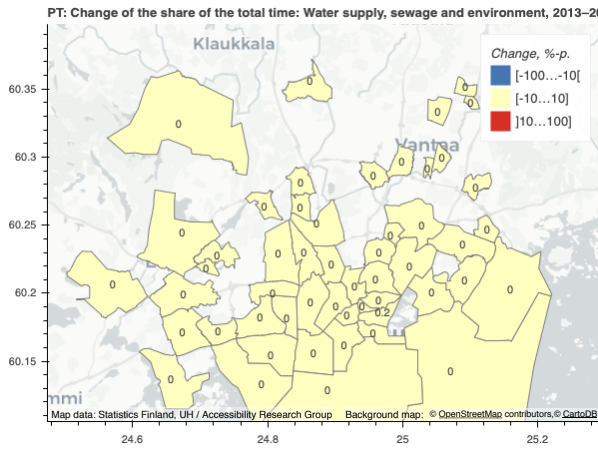
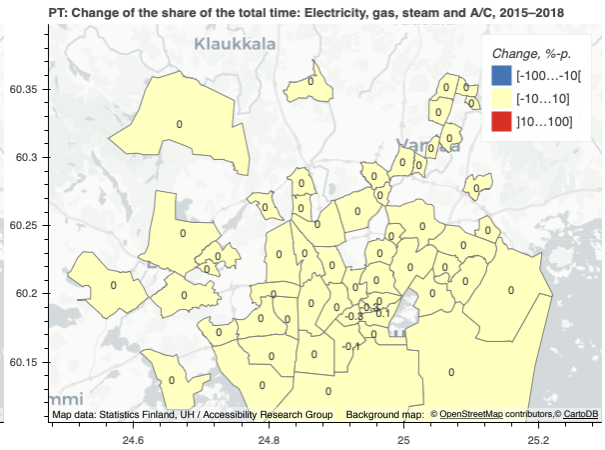
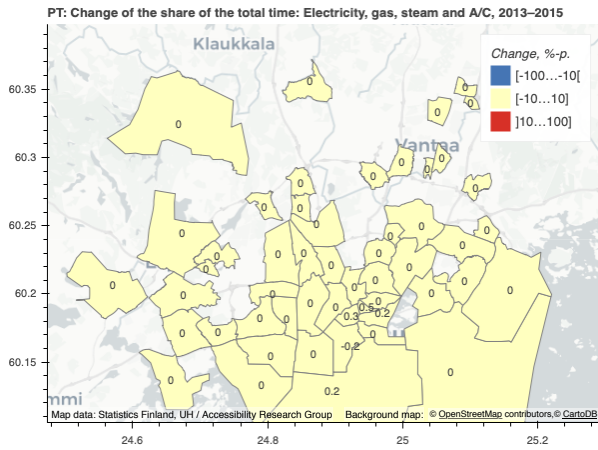
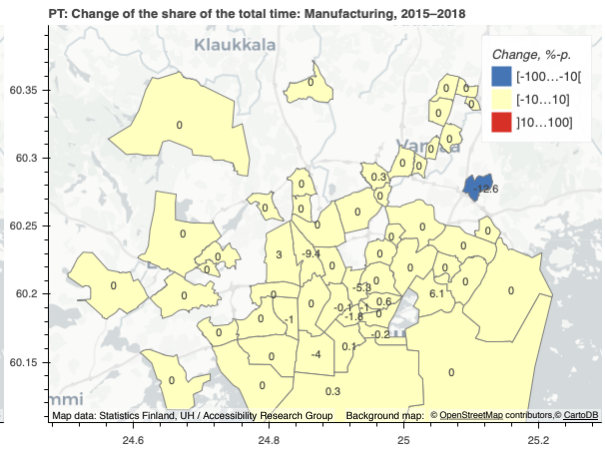
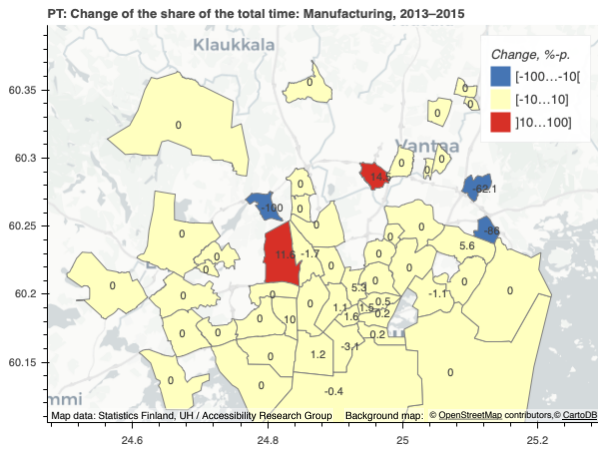


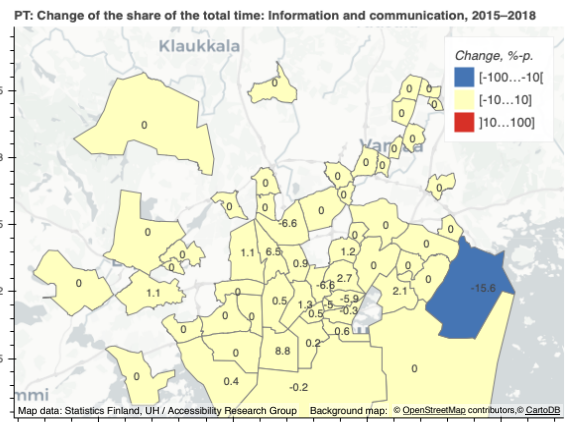
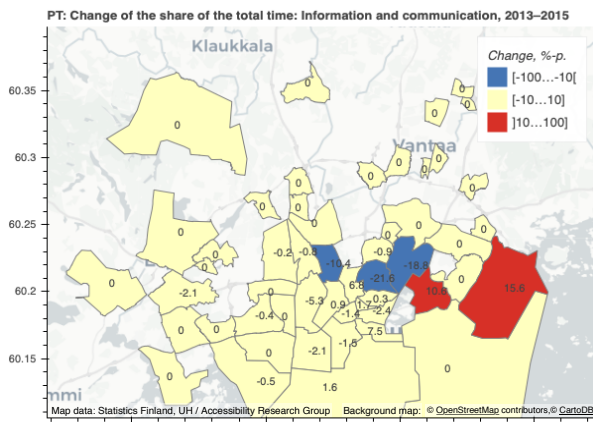
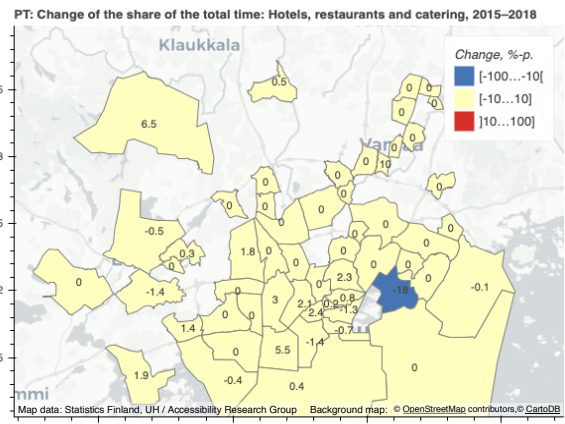
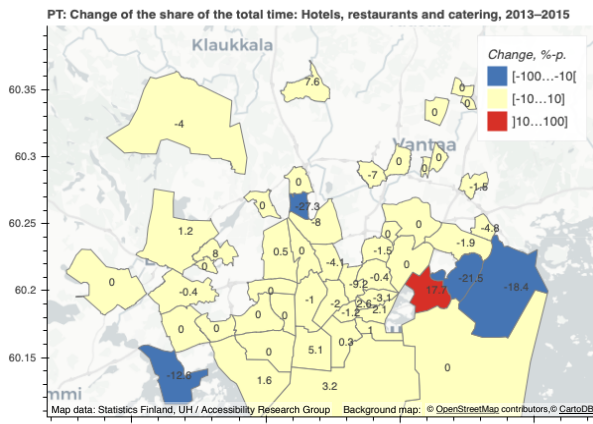
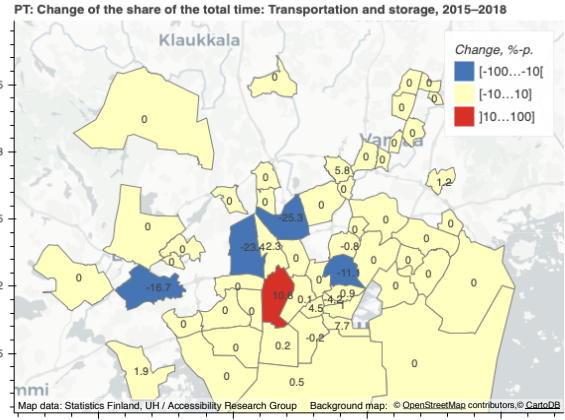
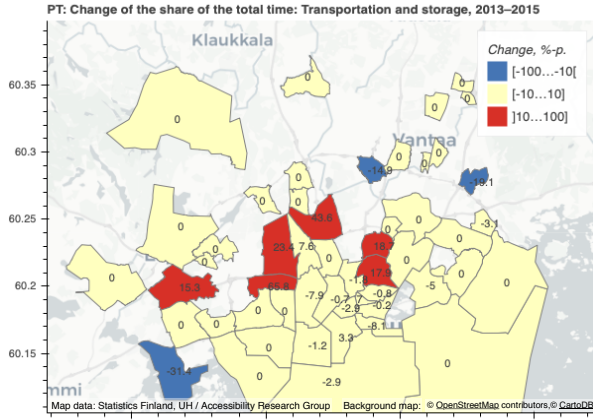
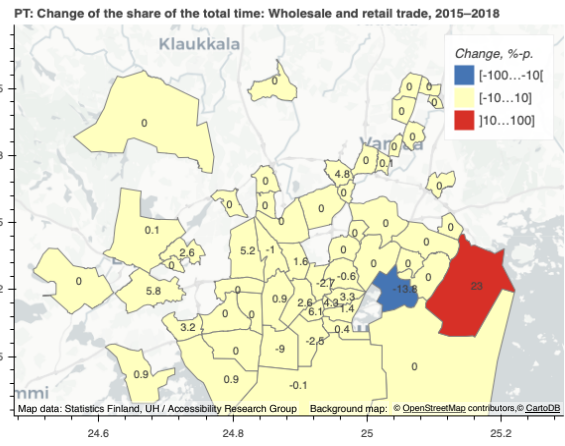
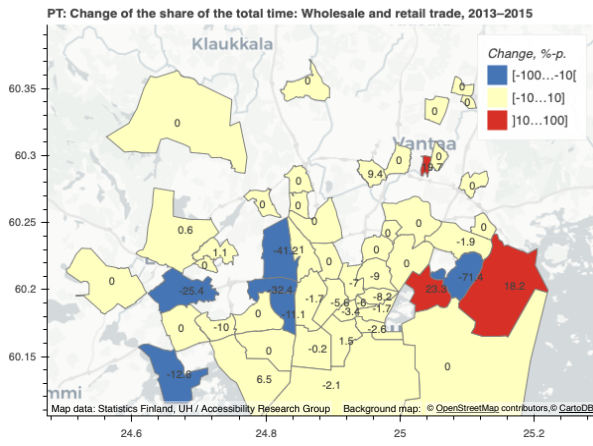
Car: Change of the share of the total time: Education, 2015–2018

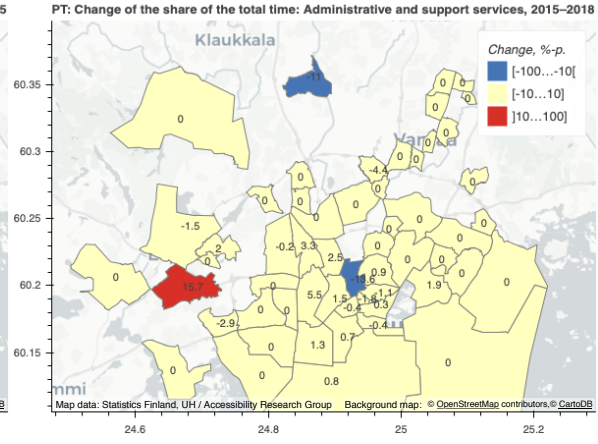
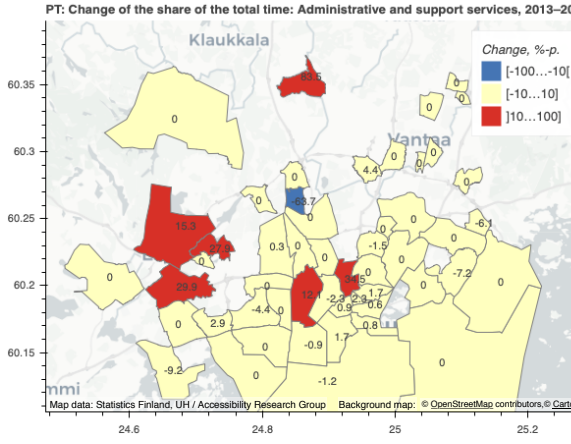
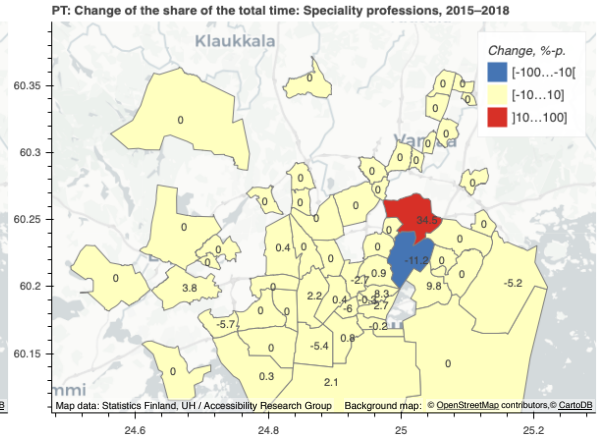
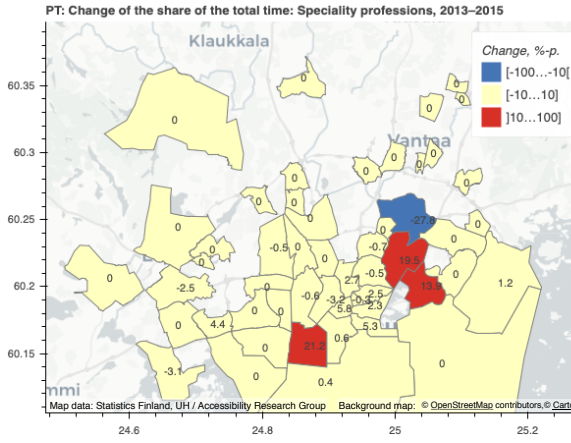
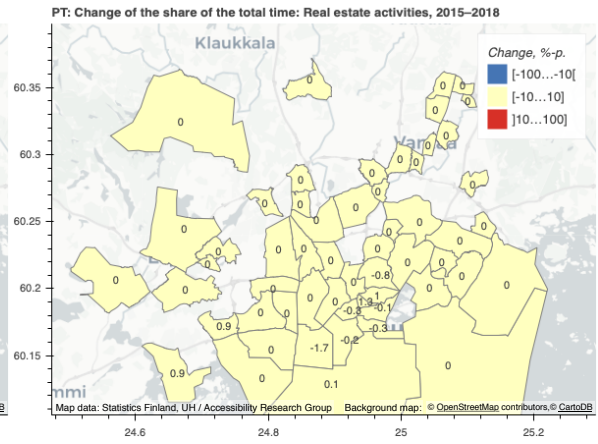
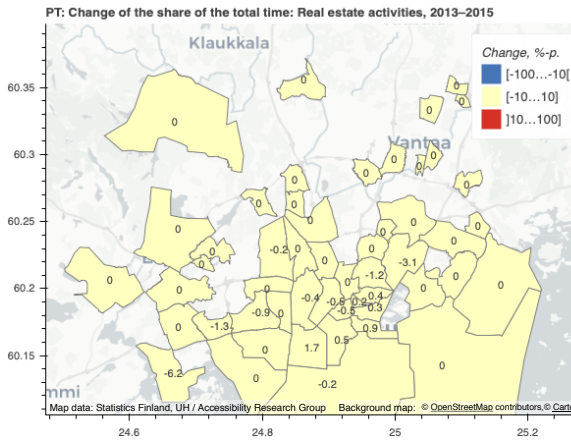
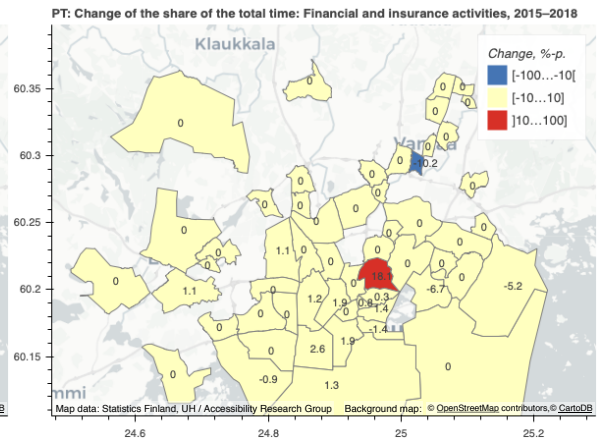
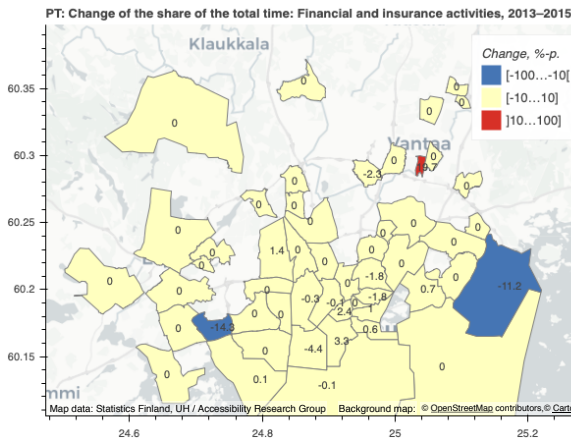


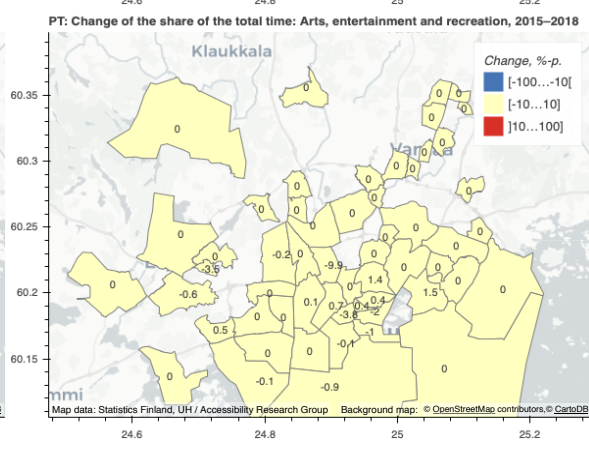
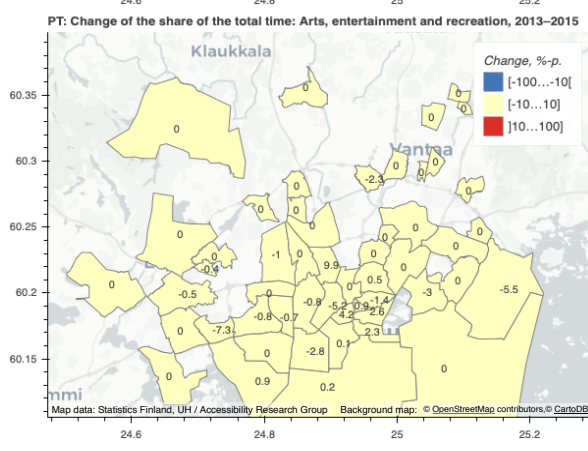
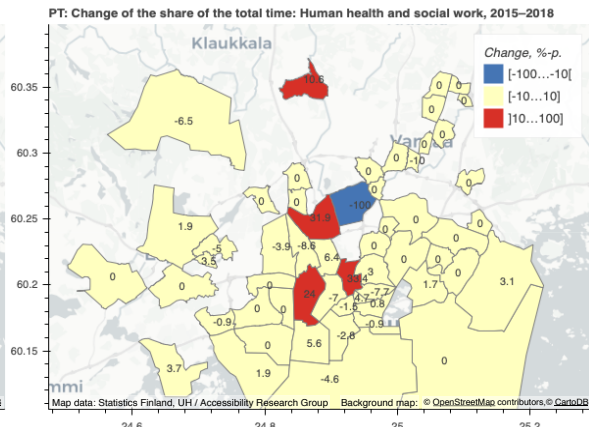
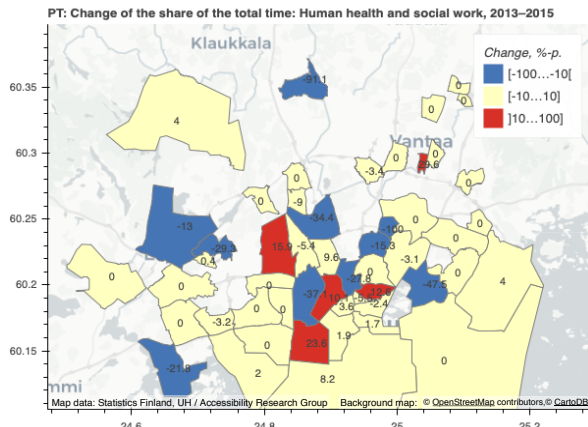
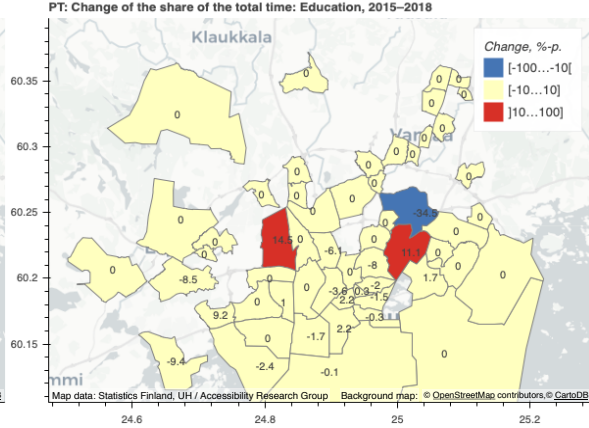
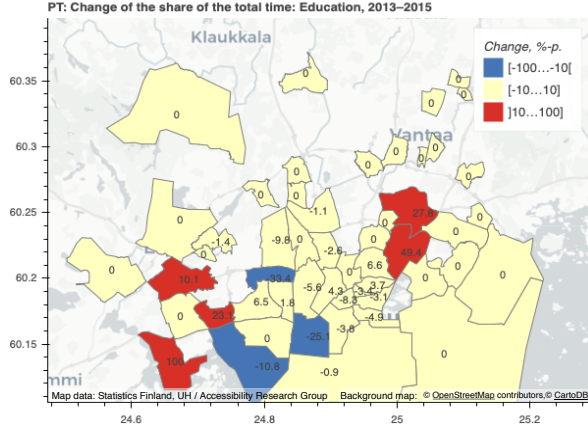
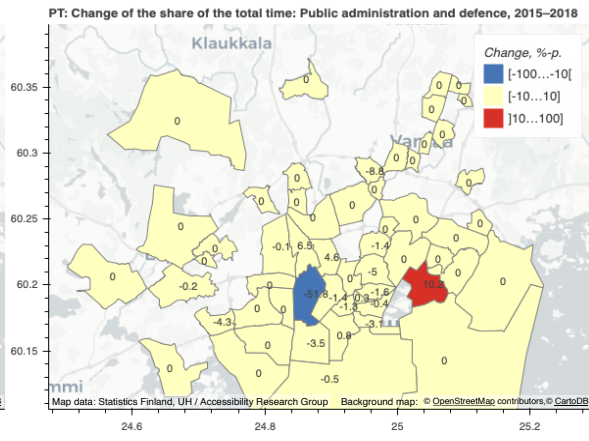
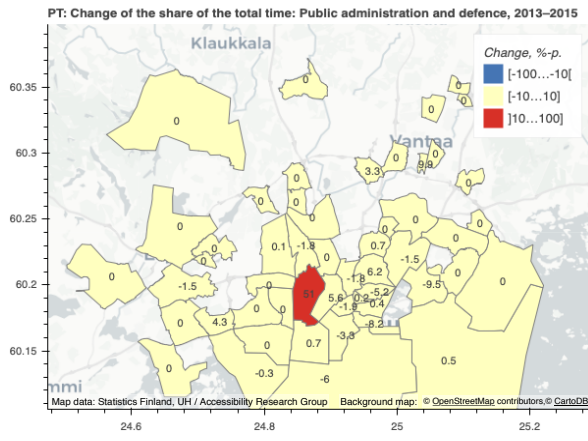


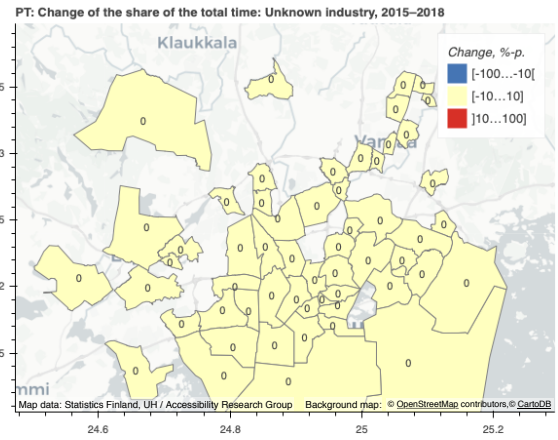
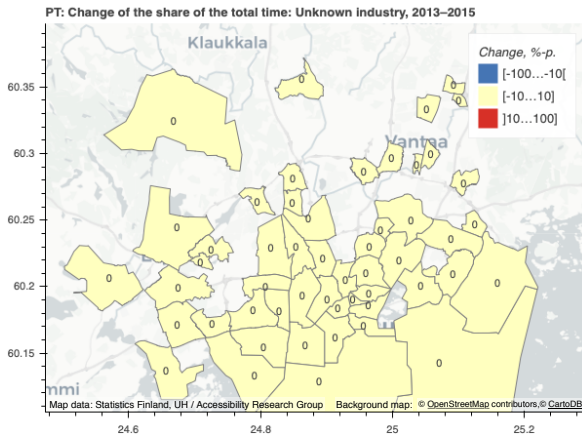
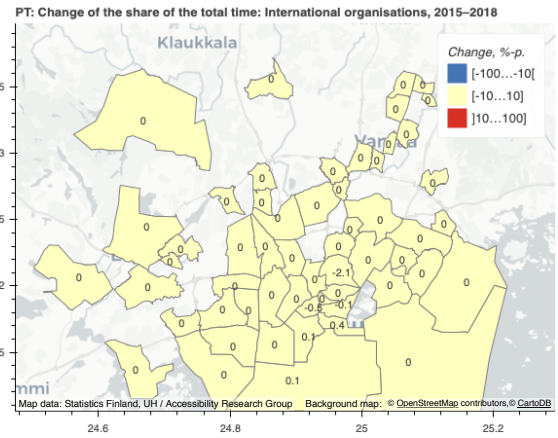
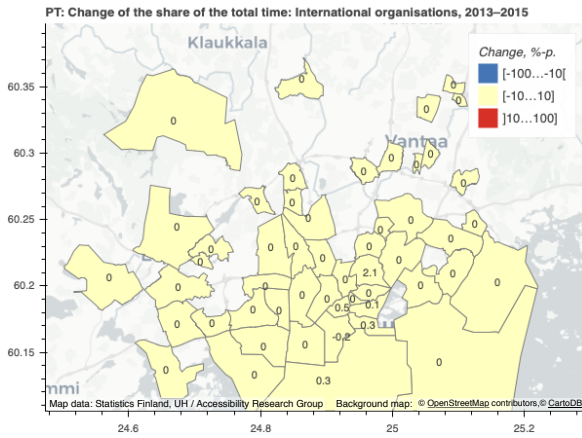
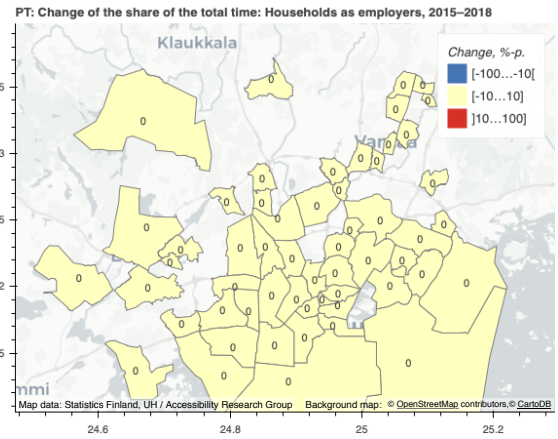
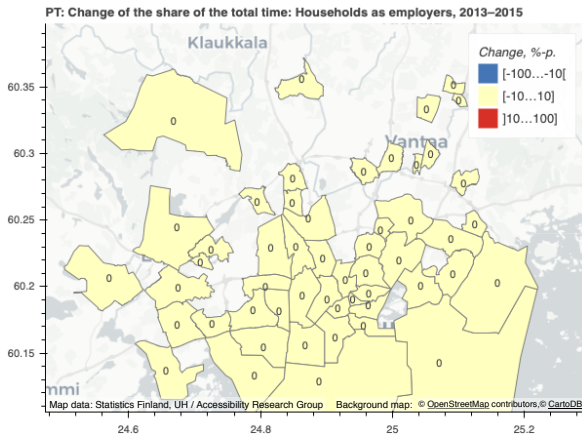
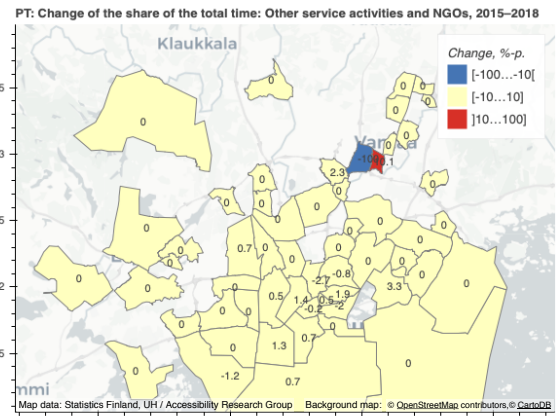
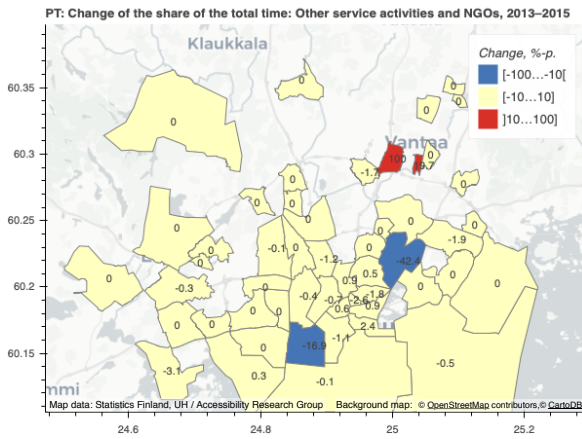




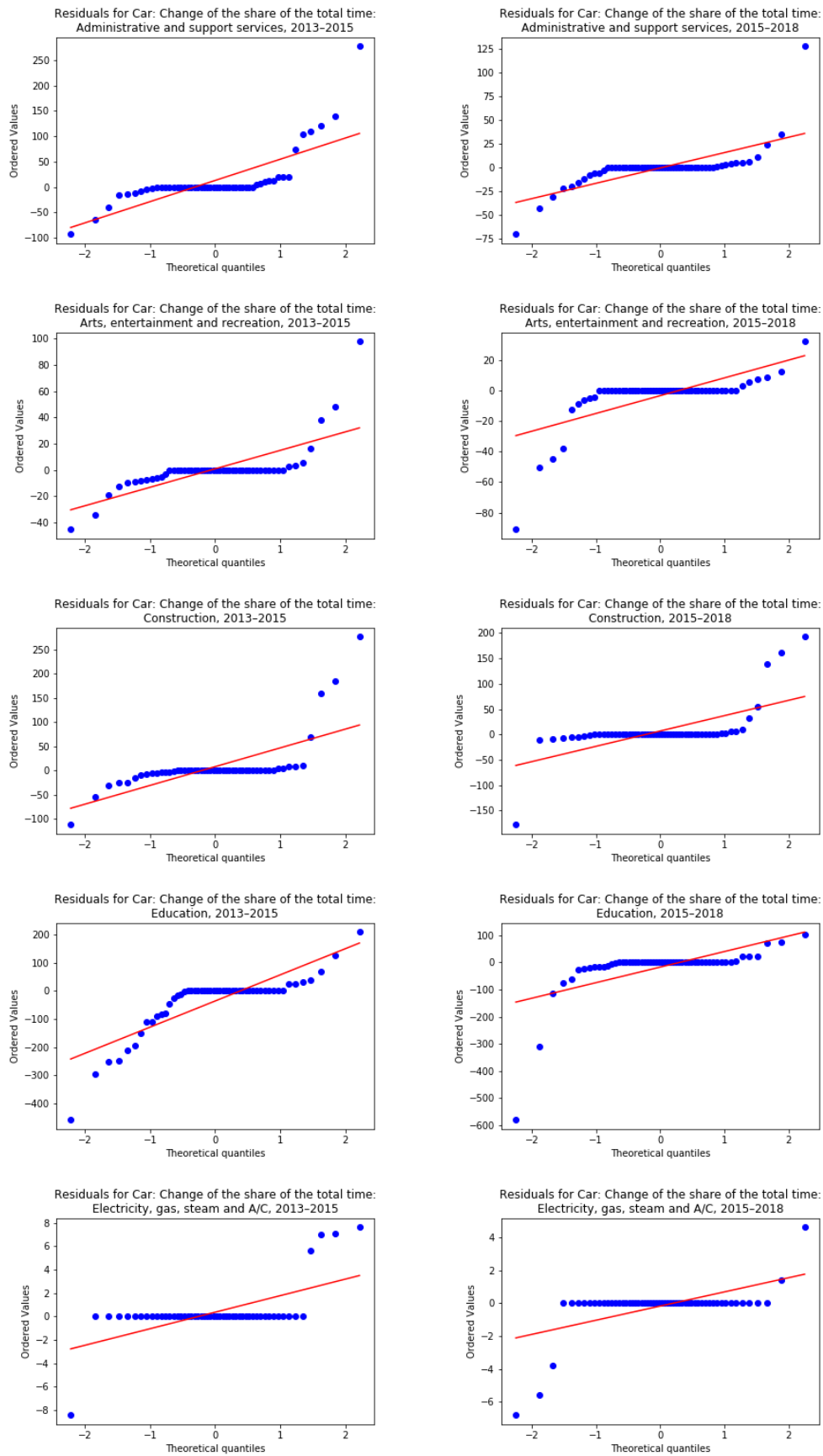


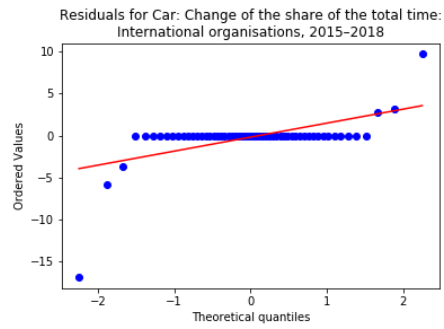
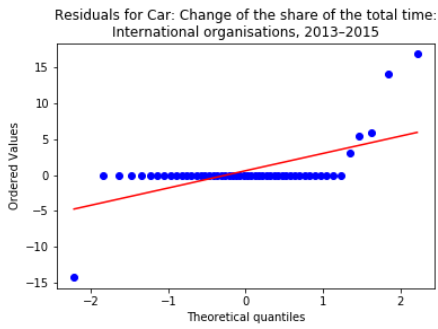
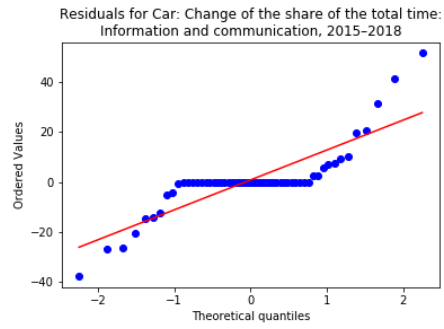
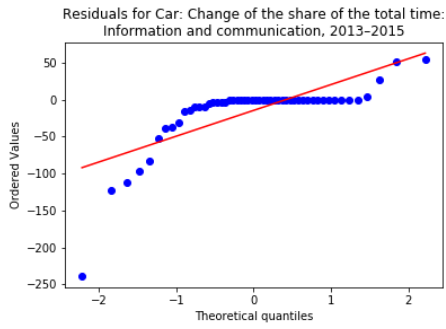
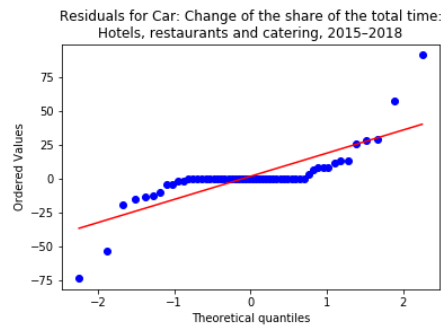
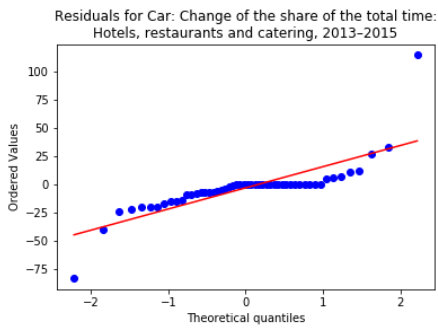
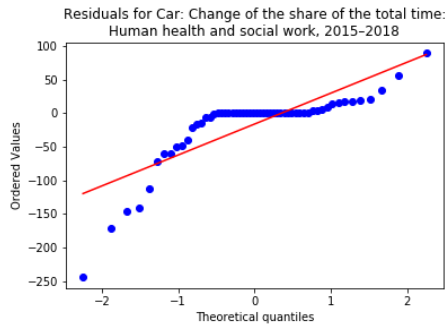
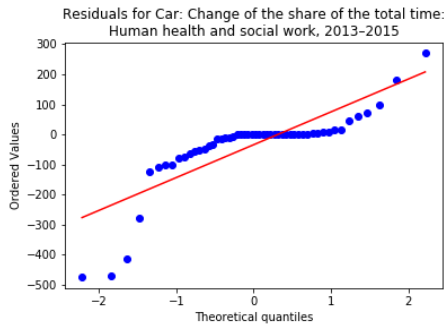
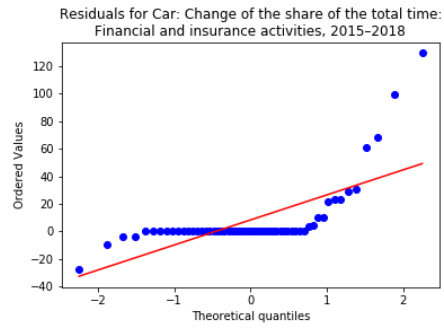
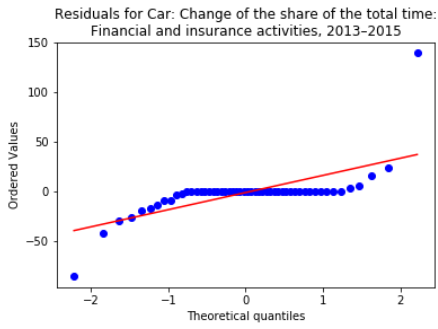


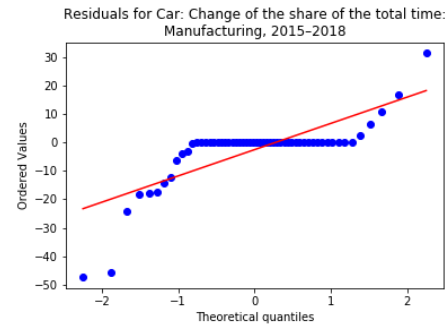
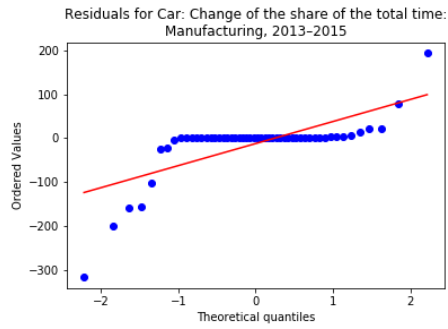
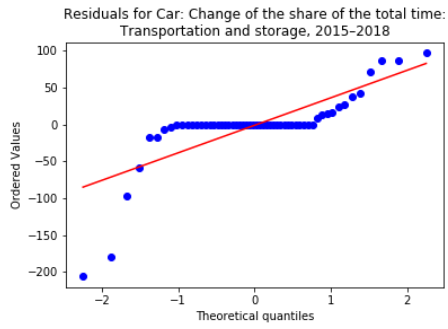
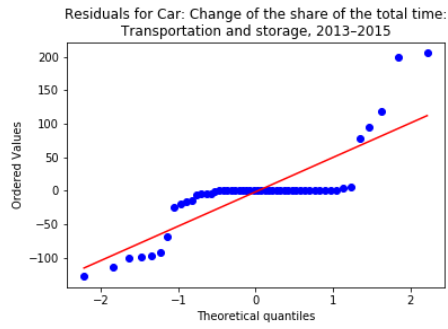




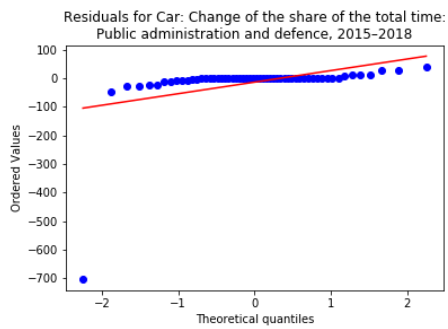
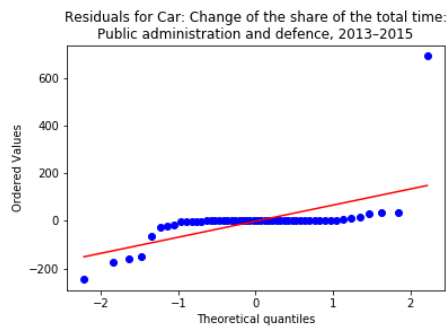
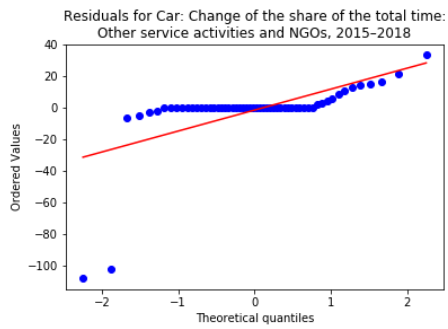
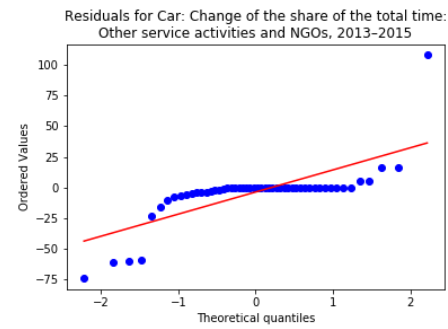
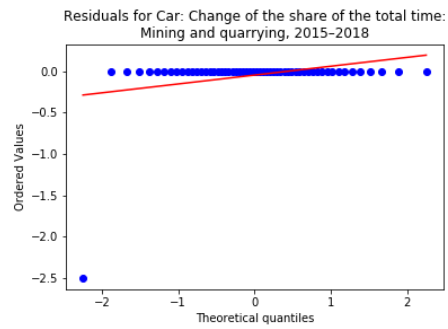
Appendix 3: Q-Q plots of IC residuals

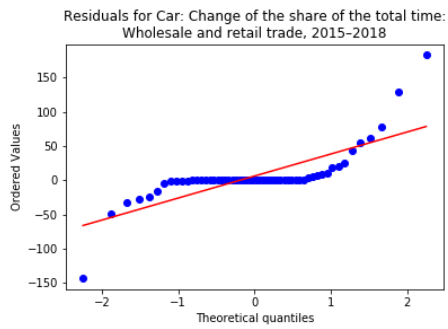
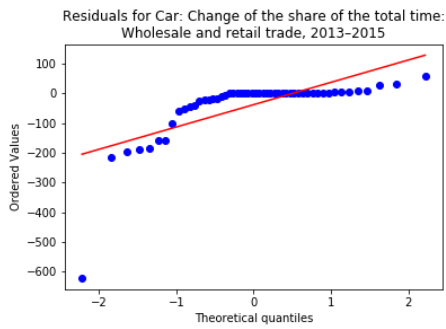
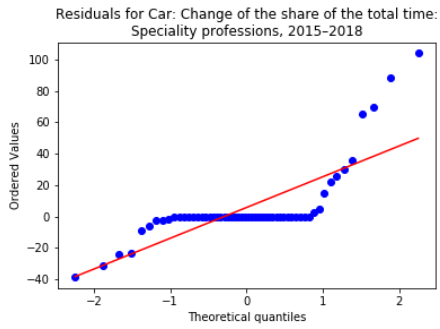
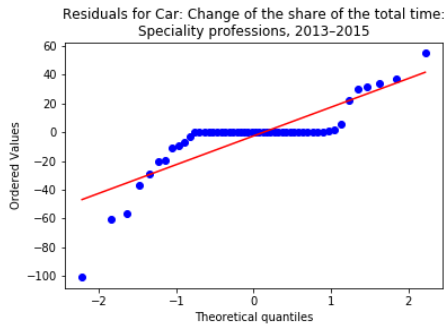
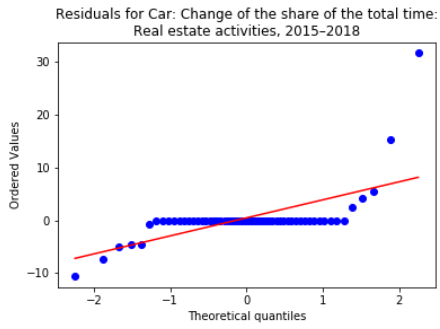
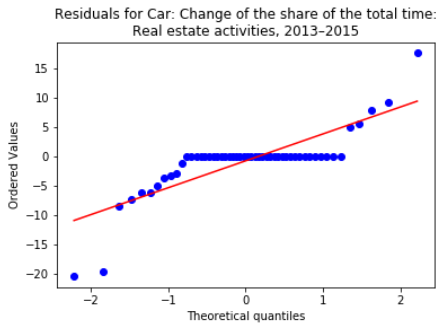




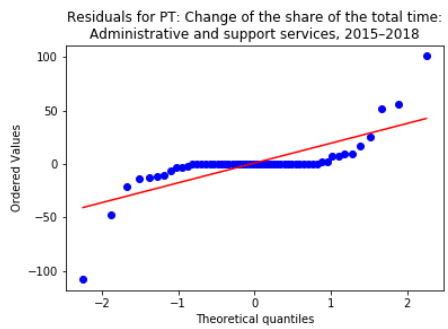
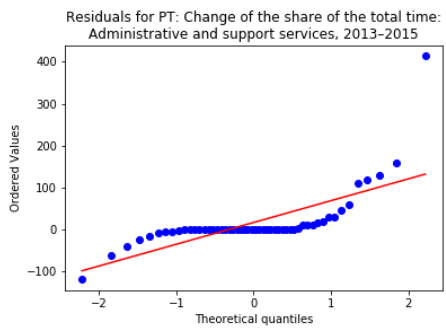
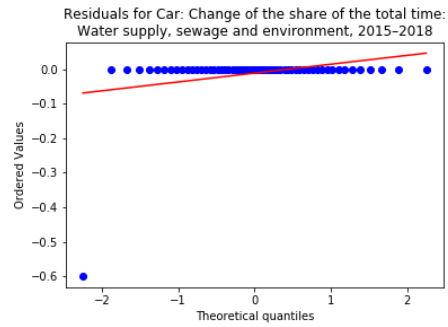


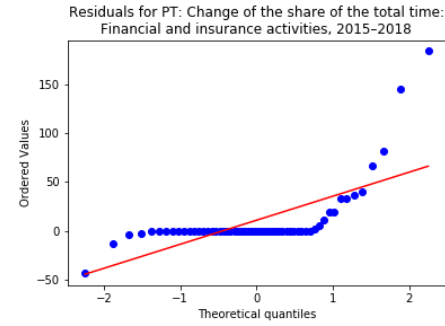
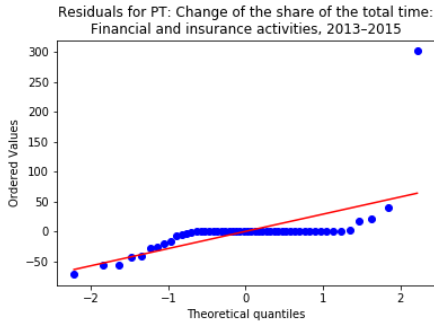
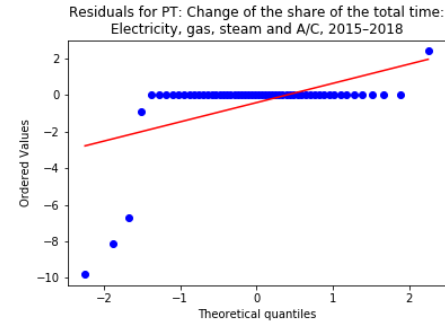
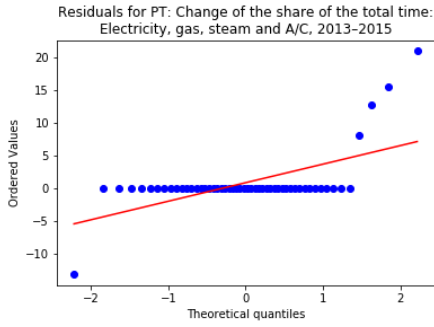
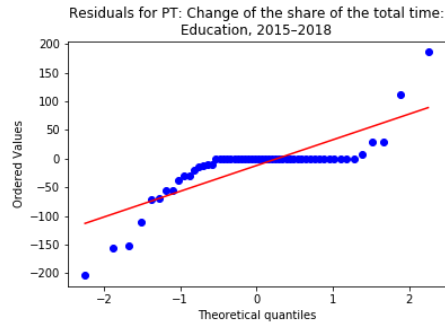
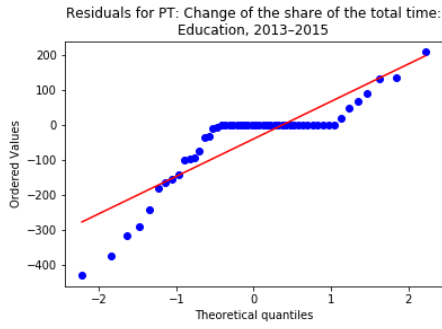
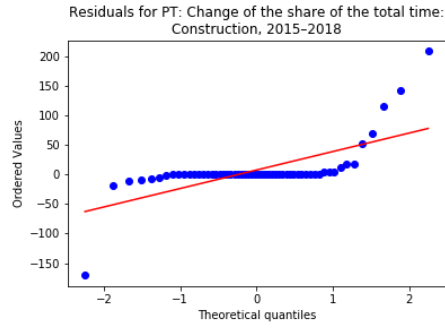
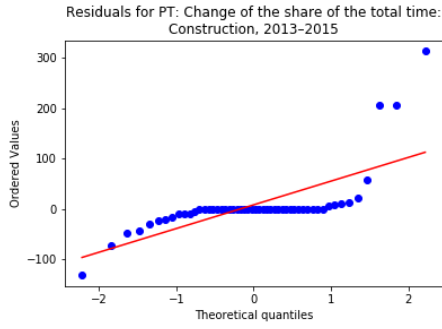
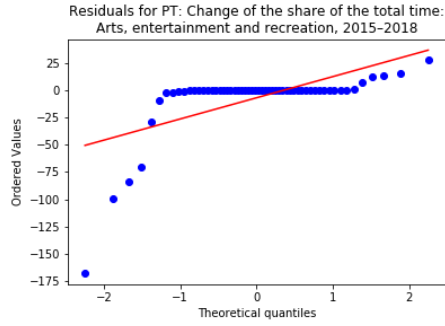
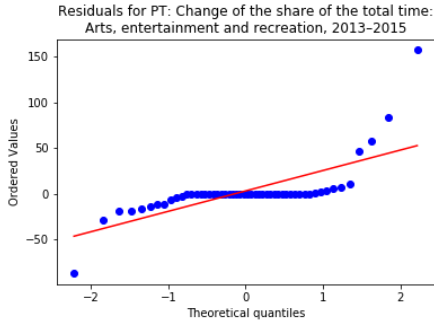
No data.

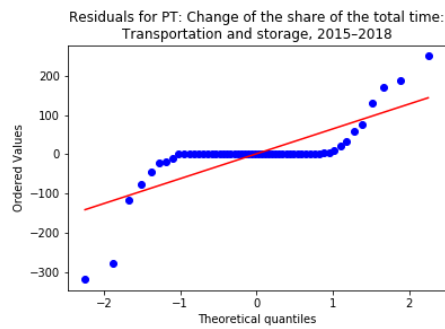
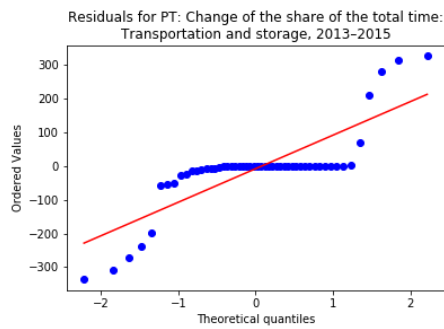
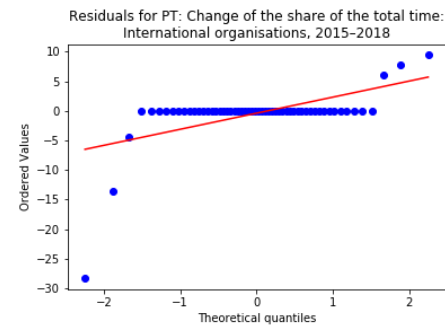
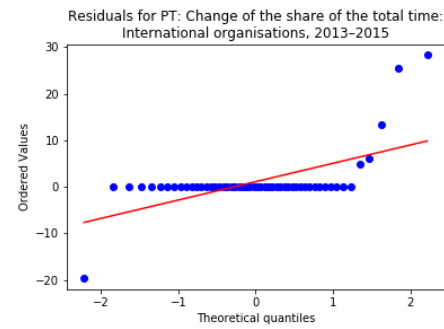
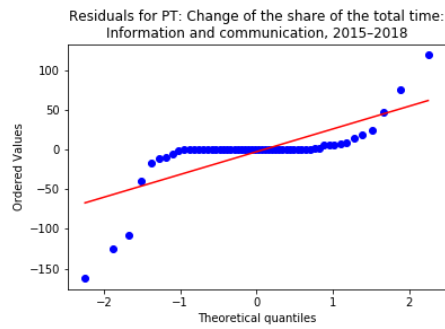
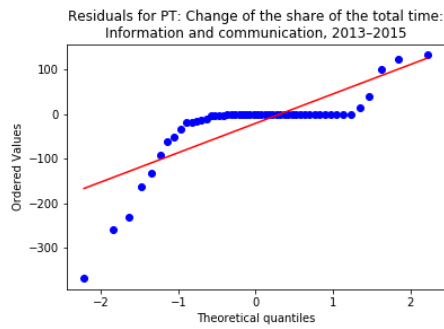
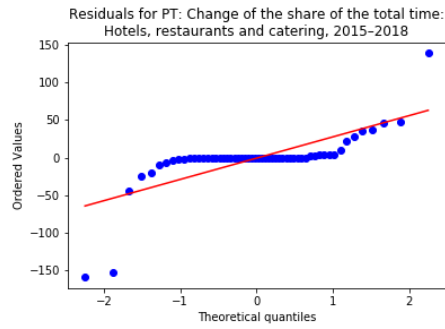
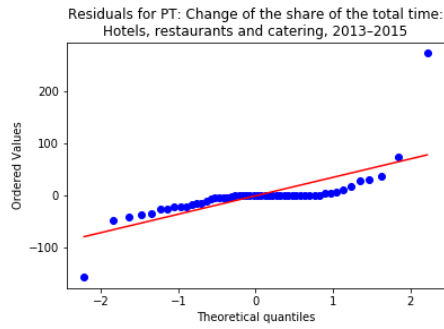
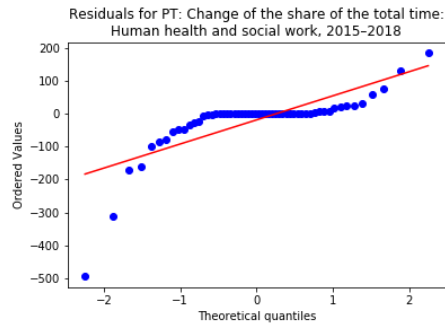
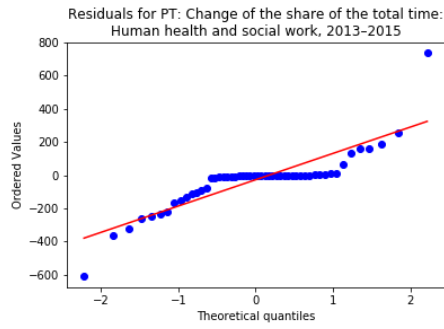


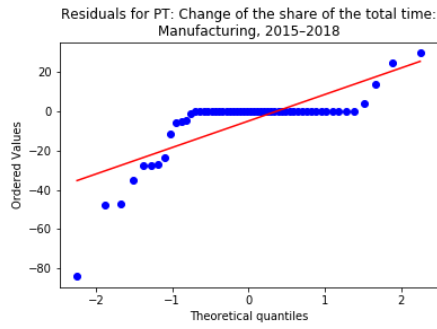
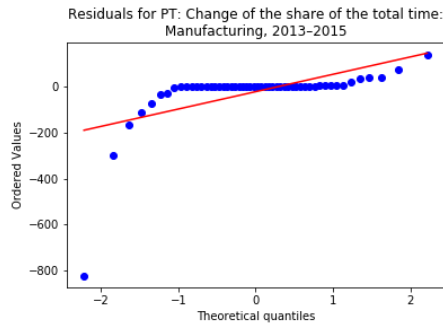


No data.

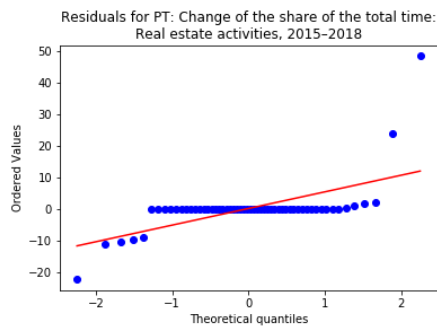
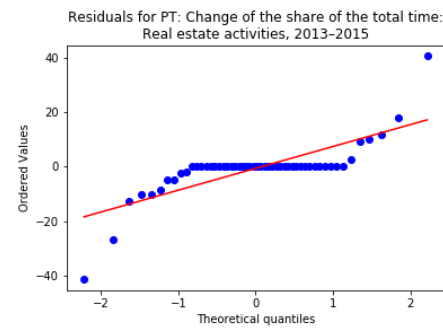
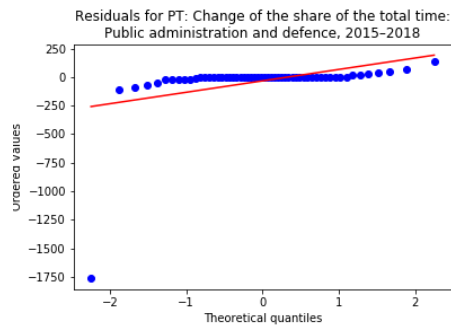
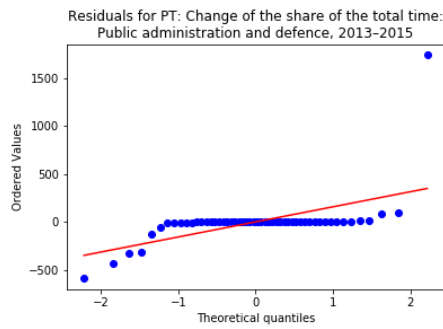
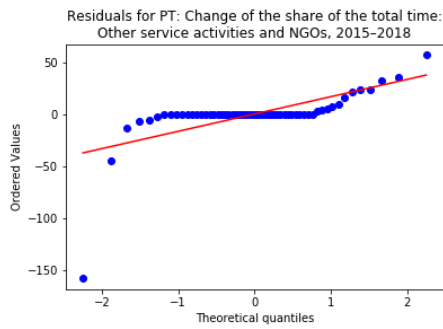
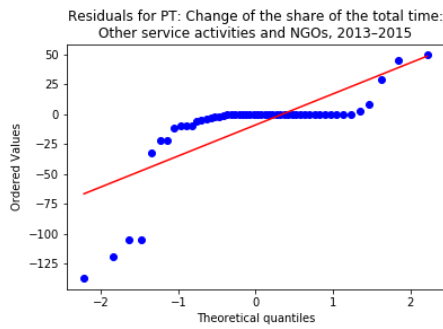
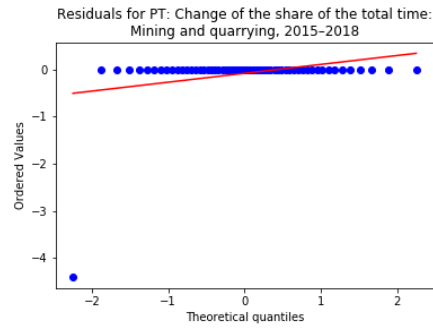


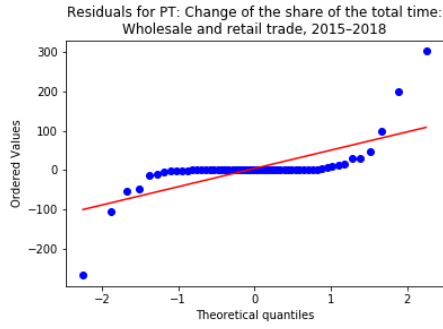
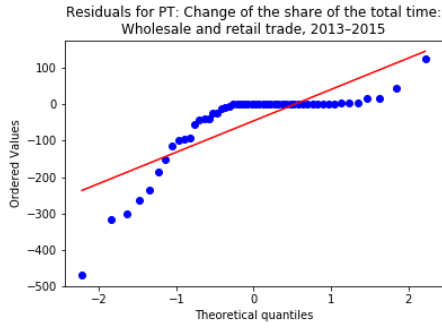
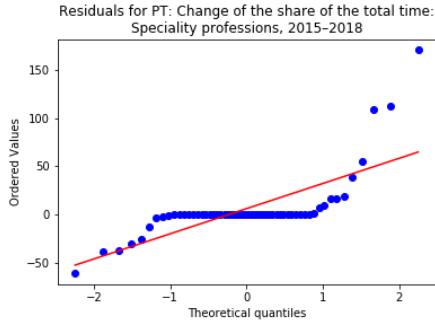
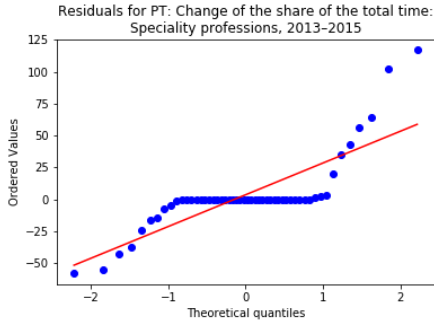




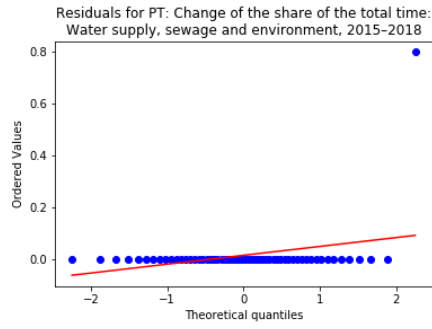


No data.

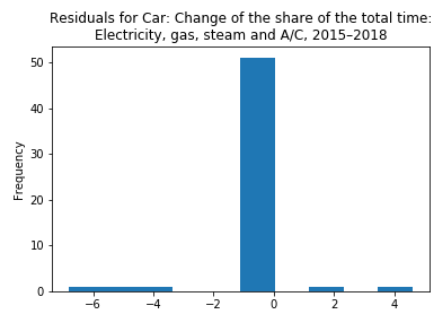
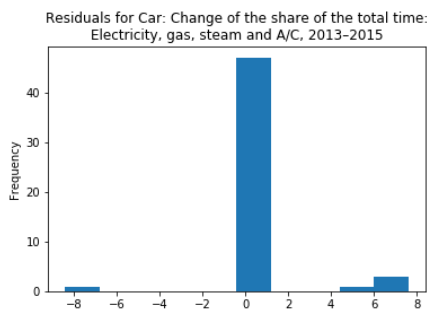
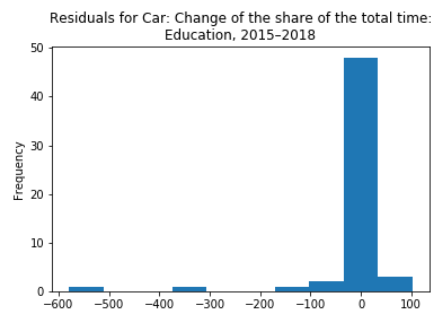
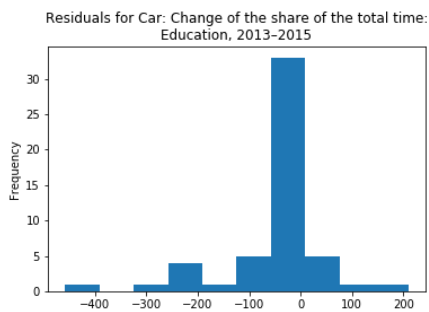
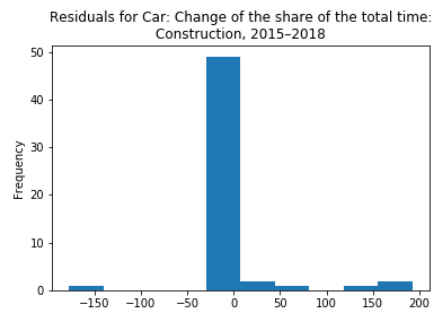
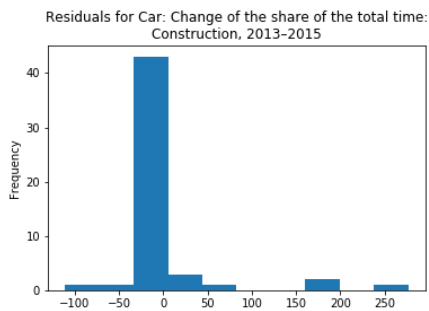
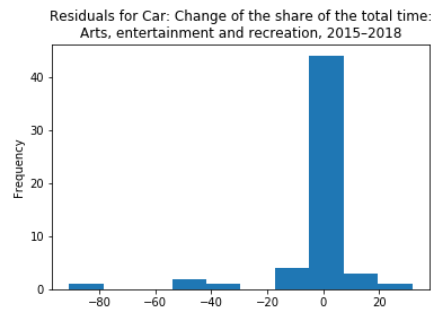
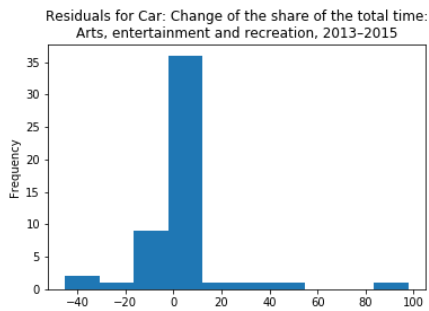
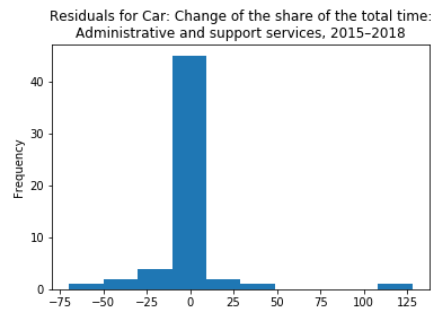
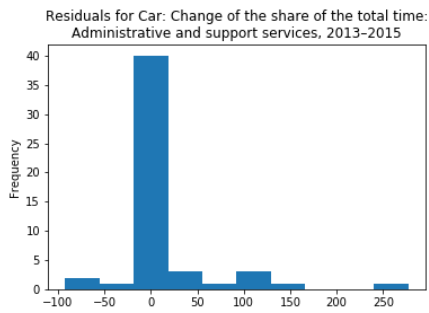


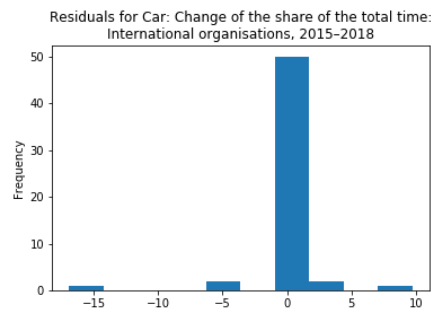
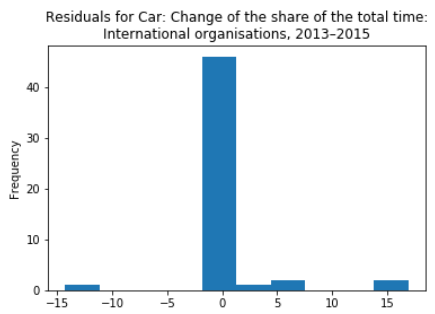
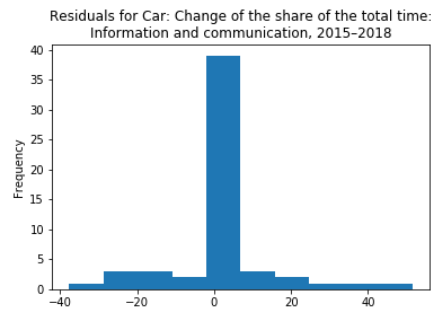
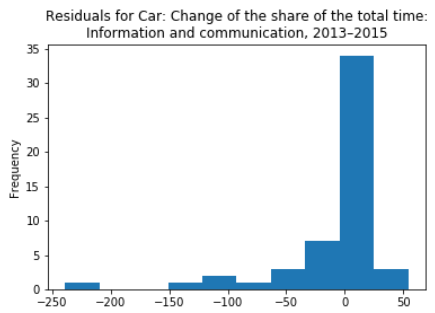
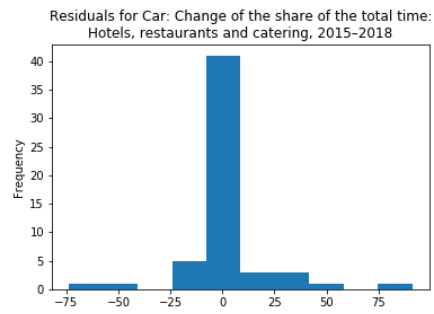
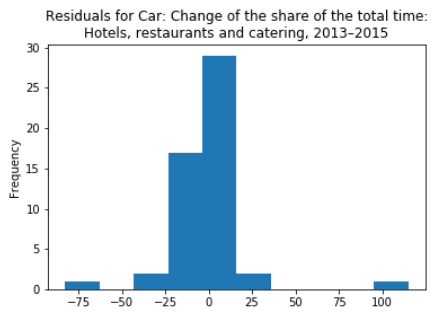
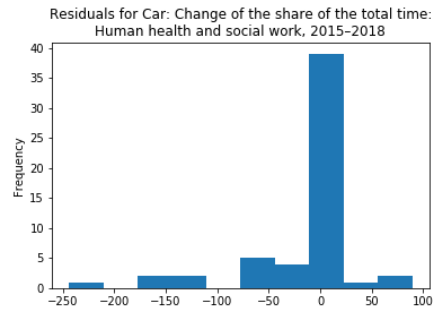
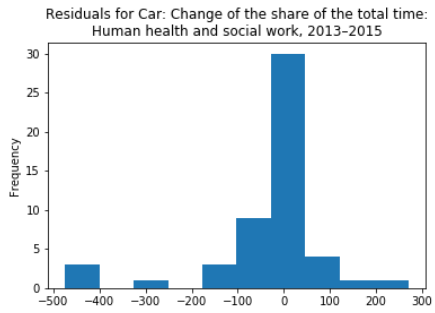
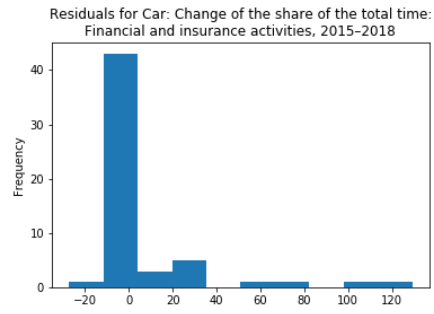
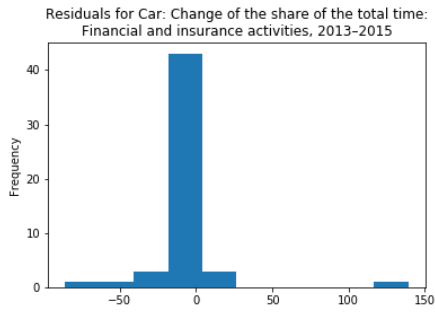


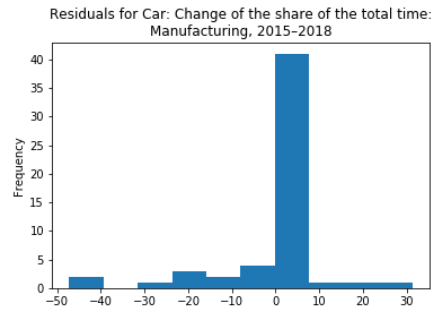
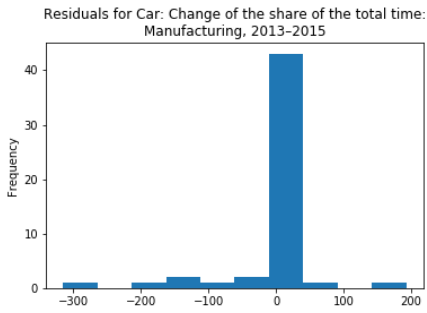
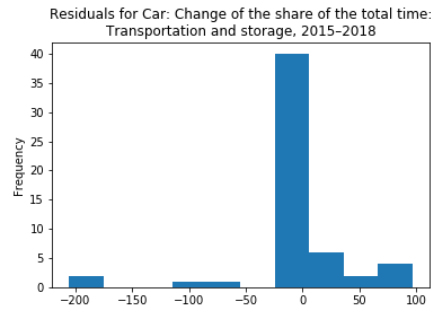
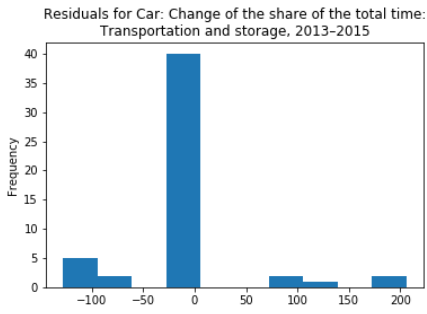
No data.



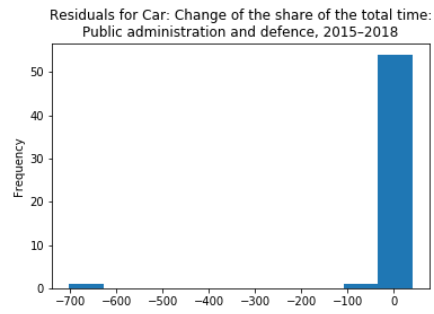
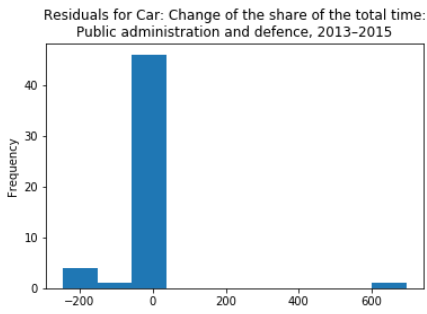
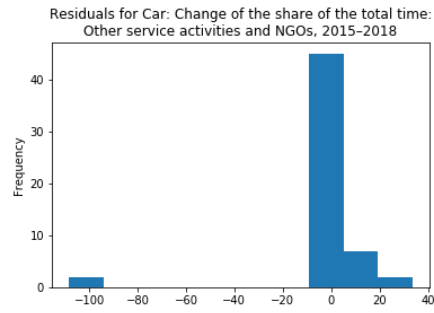
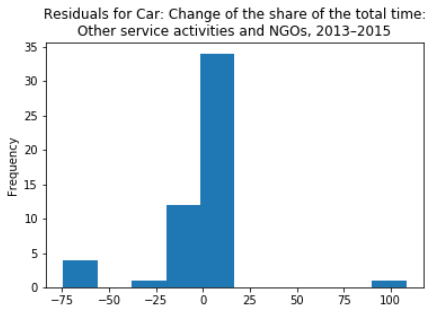
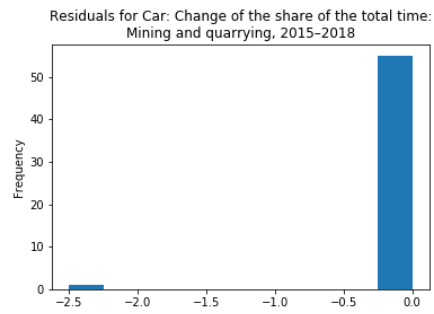
Appendix 4: Histogram plots of IC residuals

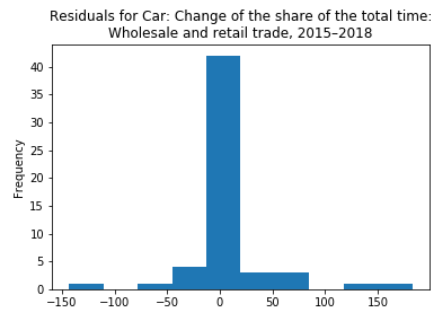
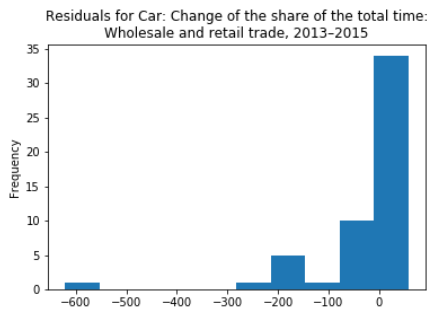
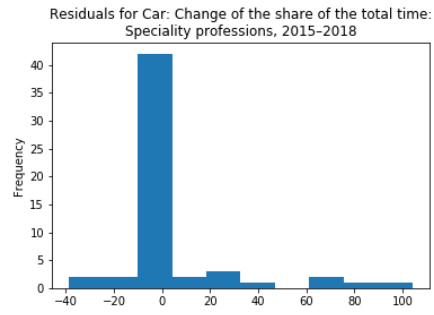
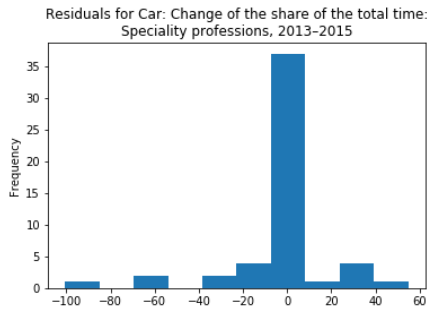
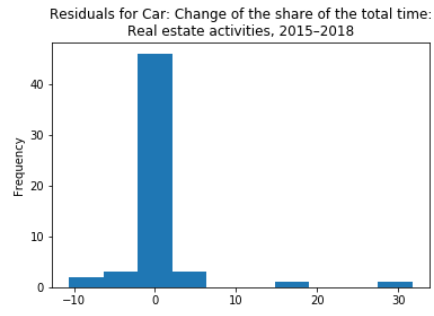
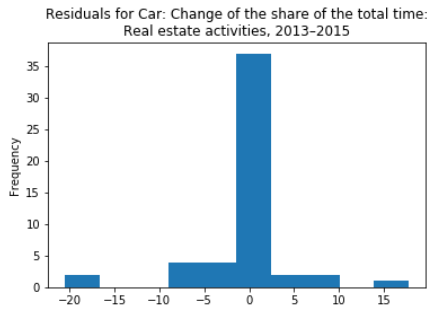




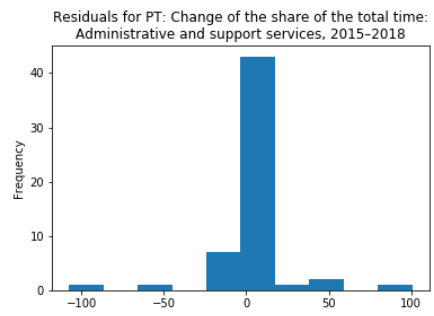
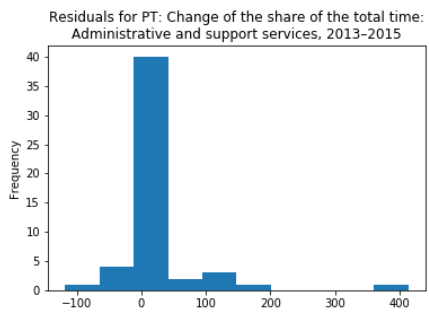
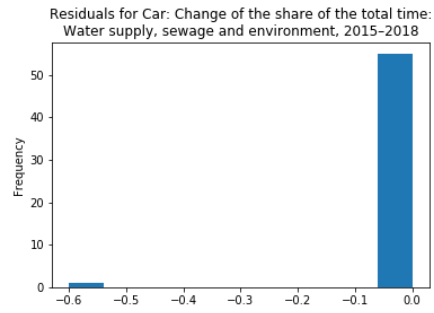


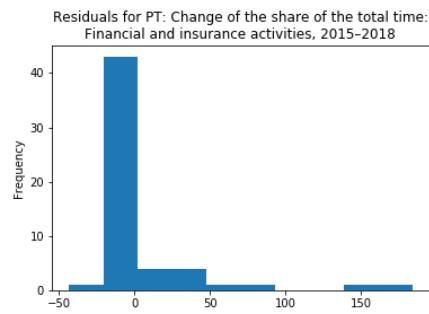
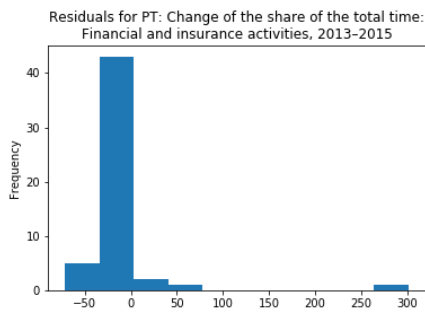
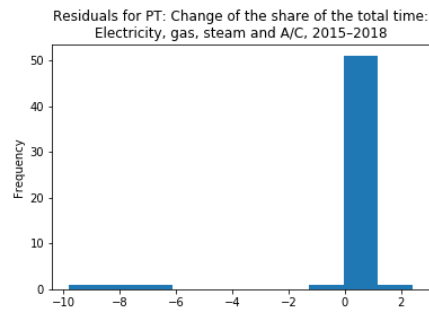
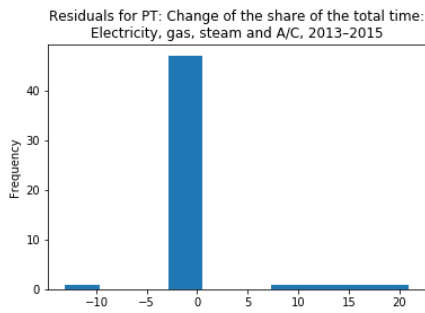
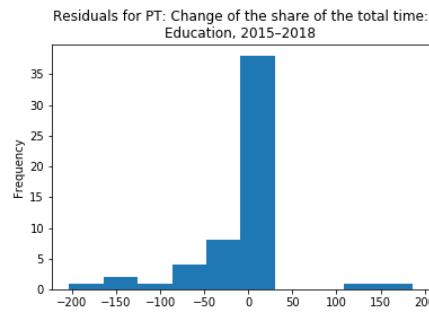
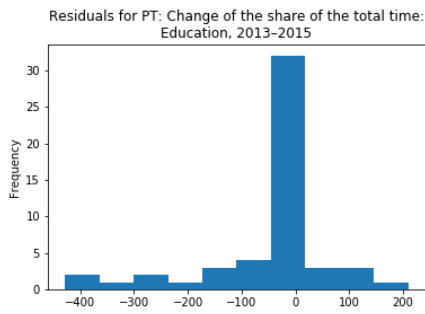
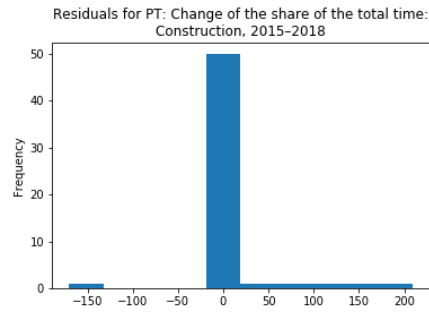
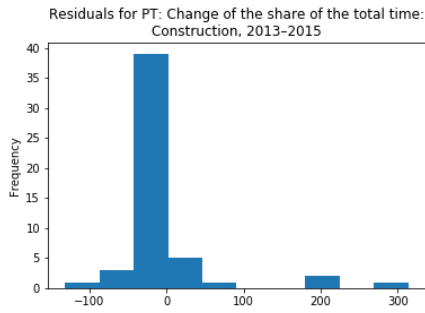
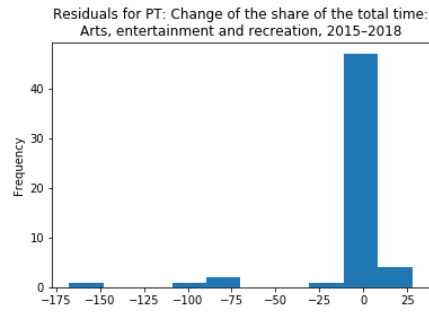
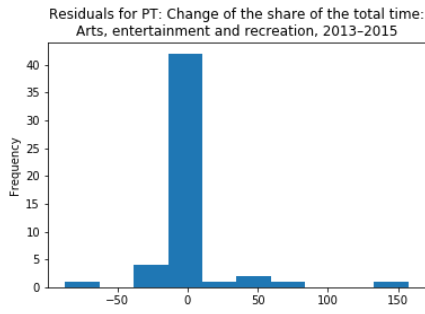
No data.

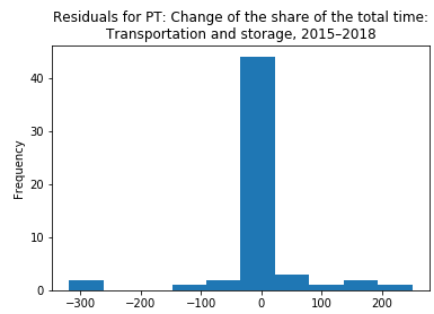
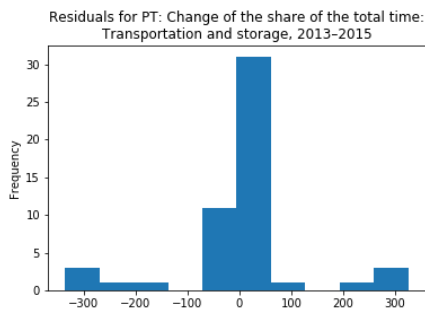
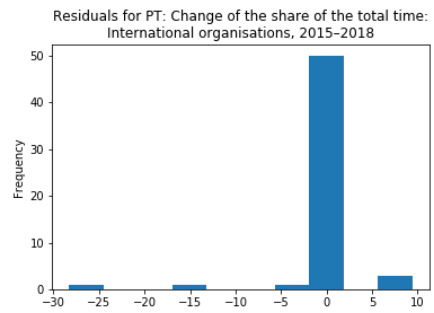
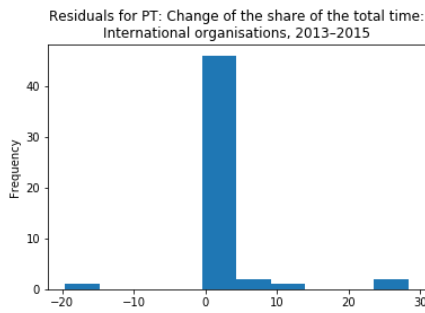
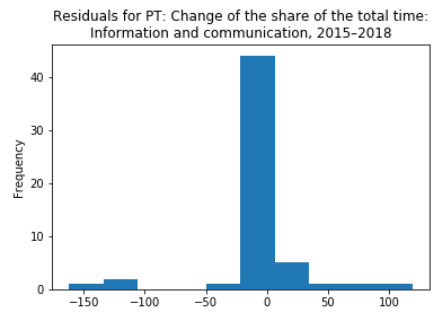
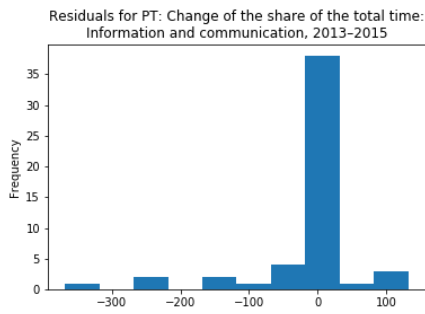
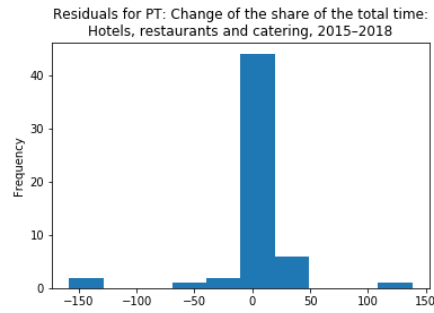
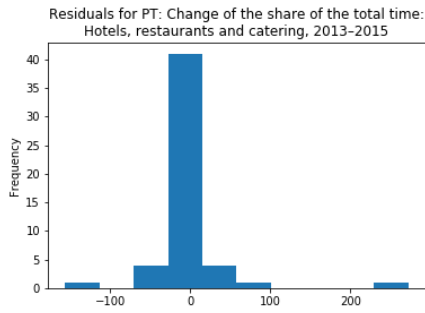
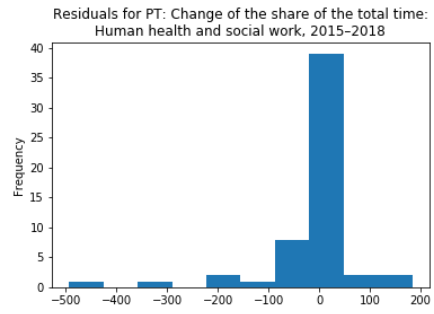
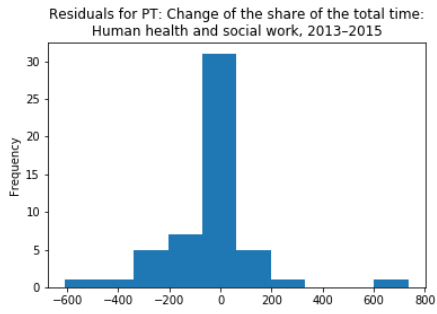


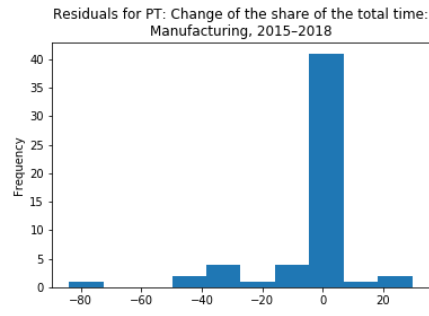
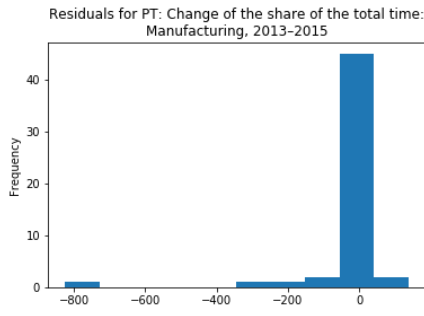


No data.

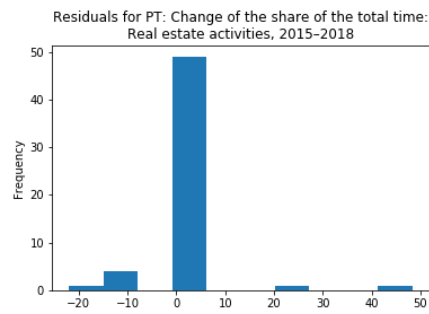
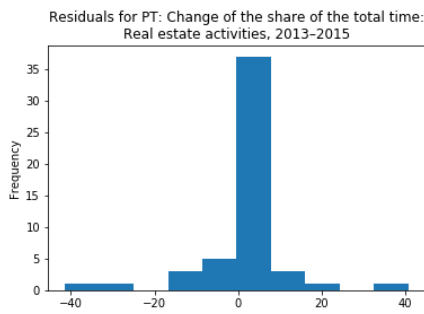
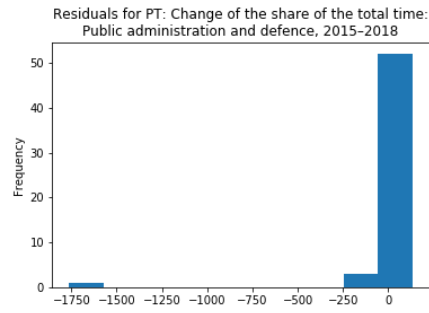
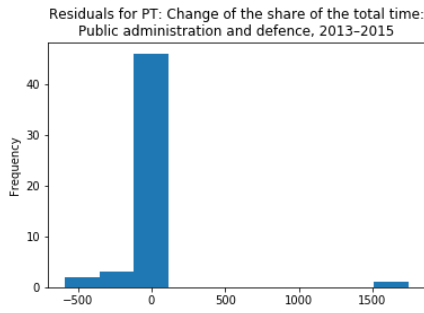
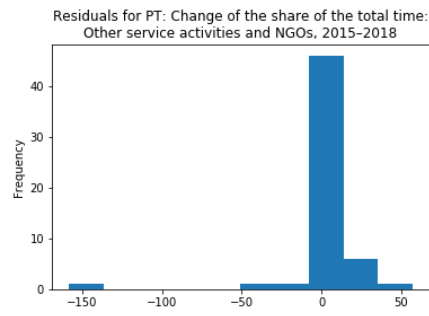
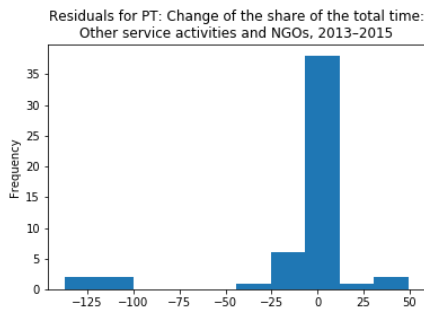
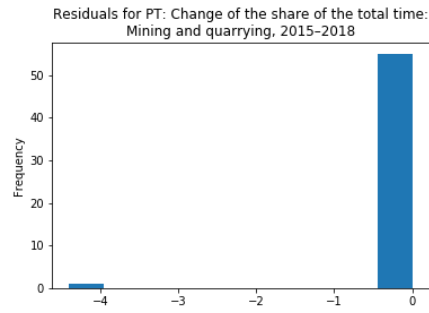


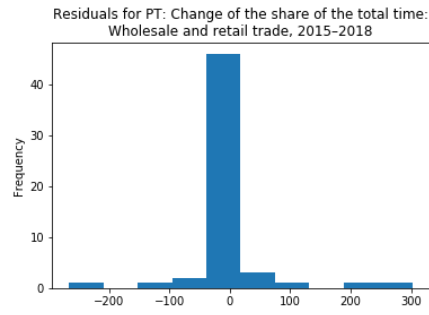
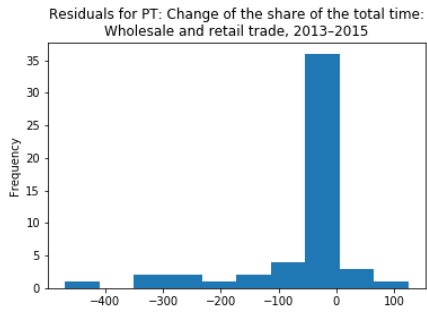
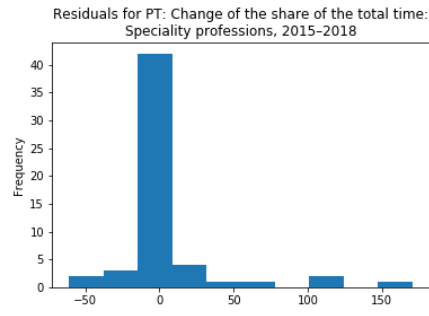
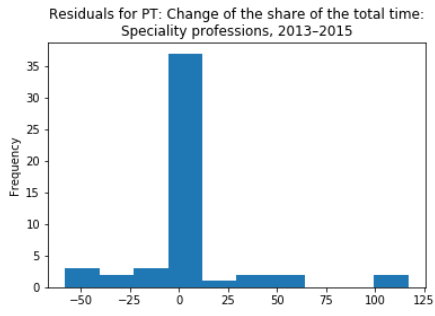






No data.





No data.

