Massaroetal.
RESEARCH
Leveraging insurance customer data to characterize socioeconomic indicators of Swiss municipalities
Emanuele Massaro
1\*
, Lorenzo Donadio
1
1, ClaudiaR.Binder
1 and Rossano Schifanella
14 2,3
15 \*

17 emanuele.massaro@epfl.ch

18 **1** 

19 HERUS Lab,

16 Correspondence:

20 ′

21 Ecole polytechnique

22 f'ed'erale de Lausanne, ENAC, IIE

23 GR C1 455 (B^atiment GR) -

24 Station 2, 1015 - Lausanne, CH

25 Full list of author information is

26 available at the end of the article

27 Abstract

28 In this paper, we explore the potential for insur ance customers data to predict

29 socio-economic indicators of Swiss municipalitie
s. First, we create a feature space

30 by defining various meaningful characteristics of a particular city through the

 $\ensuremath{\mathtt{31}}$  aggregation of individual data. We then evaluate those quantities by considering

32 available ocial statistics for Swiss municipaliti es for di⊶erent categories, i.e.,

33 Population, Transportation, Work, Space and Terri tory, Housing, and Economy,

34 and we investigate whether they can be predicted based on our feature space. For

35 the purpose of prediction, our study adopts two g eographical regression models

36 to explore the impact of spatial dependency e→ect s. Results show consistently a

37 correlation between customers' individual charact eristics and ocial

38 socio-economic indexes. Performance fluctuates de pending on the category,

39 reaching an adjustedR

40 2

41 =0.6in the case of per inhabitant housing area us ing a

 $\ensuremath{\mathsf{42}}$  5–fold cross validation setting. As a case study, we discuss in detail the

43 predictability of the percentage of the populatio n using public transportation and 44 2 Donadioetal.

3 RESEARCH

4 Leveraging insurance customer data to

5 characterize socioeconomic indicators of Swiss

6 municipalities

7 Lorenzo Donadio

8 1

9 , Rossano Schifanella

10 2,3

11 ,ClaudiaR.Binder

12 **1** 

13 and Emanuele Massaro

14 **1** 

15 \*

16 Correspondence:

17 ema.massaro@gmail.com

18 **1** 

19 HERUS Lab,

20 ′

21 Ecole polytechnique

22 f'ed'erale de Lausanne, ENAC, IIE

23 GR C1 455 (B^atiment GR) -

24 Station 2, 1015 - Lausanne, CH

25 Full list of author information is

26 available at the end of the article

27 Abstract

28 In this paper, we explore the potential for insurance customers data to predict

29 socio-economic indicators of Swiss municipalitie s. First, we create a feature space

30 by defining various meaningful characteristics of a particular city through the

 $\ensuremath{\mathtt{31}}$  aggregation of individual data. We then evaluate those quantities by considering

32 available ocial statistics for Swiss municipaliti es for diwerent categories, i.e.,

33 Population, Transportation, Work, Space and Terri tory, Housing, and Economy,

 $34\,$  and we investigate whether they can be predicted based on our feature space. For

35 the purpose of prediction, our study adopts two g eographical regression models

36 to explore the impact of spatial dependency exect s. Results show consistently a

37 correlation between customers' individual charact eristics and ocial

38 socio-economic indexes. Performance fluctuates de pending on the category,

39 reaching an adjustedR

40 2

41 =0.6in the case of per inhabitant housing area us ing a

 $\ensuremath{\mathsf{42}}$  5-fold cross validation setting. As a case study, we discuss in detail the

43 predictability of the percentage of the populatio n using public transportation and

for analyzing and predicting the impact of human activity on

- 46 urban well being on spatial and temporal scales t hat are no consistent with the
- 47 ocial statistics available (e.g., small communes, districts or neighborhoods).
- 48 Keywords:insurance data; geographical regression; socioeconomic indicators;
- 49 Swiss municipalities
- 50 1 Introduction
- 51 National Statistical Institutes (NSIs) play an im portant role in modern societies
- 52 to release precise information on social, environ mental or economical activities [1]
- 53 in the form of a census. For example, the census records key aspects such as how
- 54 many people live in an area, their ages and their income per capita, and it en-
- 55 ables the prediction of future population to info rm the need for schools, homes or
- 56 public services. Censuses are essential for many of the indicators that enable us to
- 57 measure progress towards the Sustainable Developm ent Goals [2]. Ocial statistics
- 58 on socioeconomic status are increasingly addressi ng a significant modernization of
- 59 their production processes, nationally and intern ationally [3]. This is also due to the
- 60 opportunities o⊷ered by the use of new data sourc es, such as mobile phone data [4],
- 61 social media [5], satellite images [6], credit ca rd transactions [7] and others [8,9,10].
- 62 The goal of NSIs is to integrate and combine this new information with the tra-
- 63 ditional sources of investigations and administra tive archives, and the increasingly
- 65 Massaroetal.Page 2 of29

64

- 66 widespread orientation towards the construction o f registers of elementary data inte-
- 67 grated. Three important challenges rise: i) the q uality and methodology of the data
- 68 collection, ii) privacy and legal issues and iii) the processing, storage and transfer
- 69 of large data sets. Data sources such as social m edia, and mobile phone records,
- 70 do not have a well-defined target population, str ucture and quality (see Section2
- 71 for a literature review) that make dicult to appl y traditional statistical methods
- 72 based on sampling theory. Privacy and legal aspec ts pose another challenge: the
- 73 prevention of the disclosure of the identity of i ndividuals is regulated and enforced
- 74 by international laws, and ensuring an appropriat e level of privacy is challenging
- 75 in case of heterogeneous, and multi-source large scale data streams [11]. Copyright

#### 46 open new directions

for analyzing and predicting the impact of human activity on  $% \begin{array}{c} \left( \left( \frac{1}{2}\right) -\frac{1}{2}\right) & \left( \frac{1}{2}\right) & \left( \frac{1}{2}\right)$ 

- 47 urban well being on spatial and temporal scales t hat are no consistent with the
- 48 ocial statistics available (e.g., small communes).
- 49 Keywords:insurance data; geographical regression; socioeconomic indicators;
- 50 Swiss municipalities
- 51 1 Introduction
- 52 National Statistical Institutes (NSIs) play an im portant role in modern societies
- 53 to release precise information on social, environ mental or economical activities [1]
- 54 in the form of a census. For example, the census records key aspects such as how
- 55 many people live in an area, their ages and their income per capita, and it en-
- 56 ables the prediction of future population to info rm the need for schools, homes or
- 57 public services. Censuses are essential for many of the indicators that enable us to
- 58 measure progress towards the Sustainable Developm ent Goals [2]. Ocial statistics
- 59 on socioeconomic status are increasingly addressi ng a significant modernization of
- 60 their production processes, nationally and intern ationally [3]. This is also due to the
- 61 opportunities o⊷ered by the use of new data sourc es, such as mobile phone data [4],
- 62 social media [5], satellite images [6], credit ca rd transactions [7] and others [8,9,10].
- 63 The goal of NSIs is to integrate and combine this new information with the tra-
- 64 ditional sources of investigations and administra tive archives, and the increasingly

## 66 Donadioetal.Page 2 of30

- 67 widespread orientation towards the construction o f registers of elementary data inte-
- 68 grated. Three important challenges rise: i) the q uality and methodology of the data
- 69 collection, ii) privacy and legal issues and iii) the processing, storage and transfer
- 70 of large data sets. Data sources such as social m edia, and mobile phone records,
- 71 do not have a well-defined target population, str ucture and quality (see Section2
- 72 for a literature review) that make dicult to appl y traditional statistical methods
- 73 based on sampling theory. Privacy and legal aspec ts pose another challenge: the
- 74 prevention of the disclosure of the identity of i ndividuals is regulated and enforced
- 75 by international laws, and ensuring an appropriat e level of privacy is challenging
- 76 in case of heterogeneous, and multi-source large scale data streams [11]. Copyright

- 78 the storage and the transfer of large amount of h eterogeneous information ensur-
- 79 ing security [13]. In this context, countries are increasingly favouring alternative
- 80 means of gathering information, instead of thetra ditionaltechniques of sending out
- 81 printed forms, interviewing people in person, or via the use of online questionnaires.
- 82 Alternatively, they are looking to indirect means of collecting data, taking advan—
- 83 tage of a wide spectrum of administrative data st reams that act as a proxy for the
- 84 variables of interest.
- 85 In this direction, customers insurance records re present a valuable input to model
- 86 the socioeconomic substrate of cities, and an opp ortunity for policy makers and
- 87 researchers to broaden the scope of their studie s. Social scientists raised the issue
- 88 of representativeness and sampling bias of large scale digital data. For example, in
- 89 [14] the authors show how age, gender, ethnicity, socioeconomic status, online ex-
- 90 periences, and Internet skills, influence the soc ial network sites that users generally
- 91 adopt. This has implications for the extent of th e conclusions that a study could
- 92 claim given a particular audience. Like census da ta, insurance customers records
- 93 share a similar size, reliability, and structural complexity [15]. However, they di⊷er
- 94 in their spatio-temporal granularity and collecti on costs. In fact, the information of
- 95 insurance customers is collected constantly by th e provider while the census runs
- 96 generally with a multi-yearly frequency due to it s organizational costs. A downside
- 97 is the proprietary nature of customers records th at could invalidate the possible
- 98 benefits for a broader community. However, we emb race the vision of initiatives
- 99 likeData Collaboratives
- 100 [1]
- 101 that propose a new form of collaboration, beyond the
- 102 public-private partnership model, in which partic ipants from di⊷erent sectors, in
- 103 particular companies, exchange their data to crea te public value. In this research,
- 104 we develop a methodology to predict socioeconomic indicators at a city level using
- 105 individual customers data from an insurance compa
  ny in Switzerland. Swiss mu-
- 106 nicipalities, sometimes also calledcommunities, a re the lowest administrative level
- 107 in the country. The responsibilities of the 2.212
   (as of 1 January 2019) Swiss mu-
- 108 nicipalities is decided by each Cantons. These ma
   y include the provision of local
- 109 public services such as education, medical and so

- 79 the storage and the transfer of large amount of h eterogeneous information ensur-
- 80 ing security [13]. In this context, countries are increasingly favouring alternative
- 81 means of gathering information, instead of thetra ditionaltechniques of sending out
- 82 printed forms, interviewing people in person, or via the use of online questionnaires.
- 83 Alternatively, they are looking to indirect means of collecting data, taking advan—
- 84 tage of a wide spectrum of administrative data st reams that act as a proxy for the
- 85 variables of interest.
- 86 In this direction, customers insurance records re present a valuable input to model
- 87 the socioeconomic substrate of cities, and an opp ortunity for policy makers and
- 88 researchers to broaden the scope of their studie s. Social scientists raised the issue
- 89 of representativeness and sampling bias of large scale digital data. For example, in
- 90 [14] the authors show how age, gender, ethnicity, socioeconomic status, online ex-
- 91 periences, and Internet skills, influence the soc ial network sites that users generally
- 92 adopt. This has implications for the extent of th e conclusions that a study could
- 93 claim given a particular audience. Like census da ta, insurance customers records
- 94 share a similar size, reliability, and structural complexity [15]. However, they di⊷er
- 95 in their spatio-temporal granularity and collecti on costs. In fact, the information of
- 96 insurance customers is collected constantly by th e provider while the census runs
- 97 generally with a multi-yearly frequency due to it s organizational costs. A downside
- 98 is the proprietary nature of customers records th at could invalidate the possible
- 99 benefits for a broader community. However, we emb race the vision of initiatives
- 100 likeData Collaboratives
- 101 [1]
- 102 that propose a new form of collaboration, beyond the
- 103 public-private partnership model, in which partic ipants from di⊷erent sectors, in
- 104 particular companies, exchange their data to crea te public value. In this research,
- 105 we develop a methodology to predict socioeconomic indicators at a city level using
- 106 individual customers data from an insurance compa
  ny in Switzerland. Swiss mu-
- 107 nicipalities, sometimes also calledcommunities, a re the lowest administrative level
- 108 in the country. The responsibilities of the 2.212 (as of 1 January 2019) Swiss mu-
- 109 nicipalities is decided by each Cantons. These ma
   y include the provision of local
- public services such as education, medical and so

- 112 ecutive power) and by the general assembly of all adult Swiss residents (legislative
- 113 [1]
- 114 **4**
- 116 Massaroetal.Page 3 of29
- 117 power). Many cantons leave the larger municipalit ies the option of opting for a city
- 118 parliament. Swiss citizenship is based on the cit izenship of a municipality. Every
- 119 Swiss citizen is, first, a citizen of a municipal ity (right of citizenship of the city or
- 120 of origin) and, then, of a canton (right of canto nal or indigenous citizenship). For
- 121 all these reasons, our analysis adopts the munici pality as a spatial unit of reference.
- 122 We propose a two-steps process to predict a wide range of socioeconomic indi-
- 123 cators. First, we compute a set of behavioral met rics using customers activity logs
- 124 concerning housing properties and vehicles insure d by "la Mobili`ere" in 2017. Sec-
- 125 ond, we use those microeconomics indicators as ex planatory variables on two spatial
- 126 regression models to predict 12 socioeconomic ind ices of 170 Swiss municipalities
- 127 for which we have reliable ocial statistical data as ground truth. In this work, we
- 128 focus on indices related to six di⊷erent categori es, i.e., Population, Transportation,
- 129 Work, Space and Territory, Housing, and Economy.
  We show that insurance data
- 130 customers can represent a valid resource to model socioeconomic indicators at scale.
- 131 The rest of the paper is organized as follows. Se ction2describes how insurance
- 132 data can benefit the urban data science research
   and it provides an overview of
- 133 the previous work in this area. Section3provides insights into the two datasets
- 134 that we use in this paper: insurance and census d ata. In Section4we describe the
- 135 methodology and the modeling framework. Section6d iscusses the results presented
- 136 in Section5and it provides a critical view on the limitations of the implemented
- 137 approach. Finally, in Section7we summarize the im portance and the findings of
- 138 this research and we provide insights for future directions.
- 139 2 Related Work
- 140 Researchers across various disciplines including sociology, demography and public
  - health have been keen on examining how society fu
- 141 nctions observing

143

## populations at

142 scale. Socio-economic indicators of cities, which were investigated be fore the digital

- 113 ecutive power) and by the general assembly of all adult Swiss residents (legislative
- l14 [**1**]

116

115 https://datacollaboratives.org

#### 117 Donadioetal.Page 3 of30

- 118 power). Many cantons leave the larger municipalit ies the option of opting for a city
- 119 parliament. Swiss citizenship is based on the cit izenship of a municipality. Every
- 120 Swiss citizen is, first, a citizen of a municipal ity (right of citizenship of the city or
- 121 of origin) and, then, of a canton (right of canto nal or indigenous citizenship). For
- 122 all these reasons, our analysis adopts the munici pality as a spatial unit of reference.
- 123 We propose a two-steps process to predict a wide range of socioeconomic indi-
- 124 cators. First, we compute a set of behavioral met rics using customers activity logs
- 125 concerning housing properties and vehicles insure
   d by "la Mobili`ere" in 2017. Sec-
- 126 ond, we use those microeconomics indicators as ex planatory variables on two spatial
- 127 regression models to predict 12 socioeconomic ind ices of 170 Swiss municipalities
- 128 for which we have reliable ocial statistical data as ground truth. In this work, we
- 129 focus on indices related to six di⊷erent categori es, i.e., Population, Transportation,
- 130 Work, Space and Territory, Housing, and Economy.
  We show that insurance data
- 131 customers can represent a valid resource to model socioeconomic indicators at scale.
- 132 The rest of the paper is organized as follows. Se ction2describes how insurance
- 133 data can benefit the urban data science research and it provides an overview of
- 134 the previous work in this area. Section3provides insights into the two datasets
- 135 that we use in this paper: insurance and census d ata. In Section4we describe the
- 136 methodology and the modeling framework. Section6d iscusses the results presented
- 137 in Section5and it provides a critical view on the limitations of the implemented
- 138 approach. Finally, in Section7we summarize the im portance and the findings of
- 139 this research and we provide insights for future directions.
- 140 2 Related Work
- 141 Researchers across various disciplines including sociology, demography and public
  - health have been keen on examining how society fu
- 142 nctions observing

# populations

143 at scale. Socioeconomic indicators of cities, which were investigated be fore the dig-

144 ital

finding correlations between demographic factors di⊶erences 145 between urban and suburban areas [16], crime rate [18], p opulation health [19], res-146 idential segregation [20] or waste production [17]. Since the aforementioned findings 147 were mostly based on survey results, they may have been adected by the fact that 148 people could have altered their answers knowing that they weremonitored. Today in the digital era in some cases information abou t people behavior is collected even without them being aware of that, let alone with their informed consent. In the era where the usage of digital technologie 151 s is so omnipresent, people every 152 day leave more digital trails than we are currently able to process. Indeed, most of 153 the recent studies focused on the use of digital andbig datato predict and study 154 socio-economic indicators of cities and countries. In extensive amount of related 155 work scholars used diwerent digital datasets for diwer ent research purposes such as 156 for studying human dynamics through cell phone data [4,21], social media posts [22], 157 vehicle GPS traces [23] or credit card data [7,24]. 158 It has been shown that social media data can predict the interplay between de-159 mographic attributes and gender gap [25], the monitoring [26] and assimilation of 160 migrants [27]), unemployment [28], and health outcomes [9]. Social media are in-161 creasingly used for demographic attribute prediction [29]. Using web search engine 163 Massaroetal.Page 4 of29 164 Figure 1: Comparison between the spatial coverage

of the insurance and the Union 165 of Swiss Cities datasets. (A) Total number of ins urance customers at zip code level. 166 The total number of customers is 1,341,328 which represents nearly 15% of the 167 Swiss population. (B) Number of inhabitants aggre

gated at municipality level. The 168 total surface covered by those municipalities is almost 3890km

170 which represents

162

169 2

171 9.35% of the country's area.

finding correlations between demographic factors di<del></del>erences in urban and suburban areas [16], crime rate [18], p 146 opulation health [19],

residential

147 segregation [20] or waste production [17]. In all these cases the results were based

on the active participation of individuals to sur 148 veys, that

may have been a⊶ected

149 by the tendency of respondents to alter their behavior knowing that they aremon-

itored. On the contrary, the digitalization of th e modern society allows to model

human behavior by means of indirect and less intr

151 usive data collection methodolo-

gies, that are often a byproduct of services desi

152 gned for di⊷erent purposes. In a scenario where the collectivity produces every da

153

more digital footprints than we

154 are able to process, the majority of the recent studies focused on the use ofbig data

155 streams to predict and study socioeconomic indicators of cities and countries. An

156 extensive body of work adopt digital traces such as cell phone records [4,21], social

157 media posts [22], vehicle GPS traces [23] or credit card transactions [7,24]tomodel

158 human dynamics at scale. For example, it has been shown that social media data

159 can predict the interplay between demographic attributes and gender gap [25], the 160 monitoring [26] and assimilation of

migrants [27]), unemployment [28], and health

161 outcomes [9]. Social media are increasingly used for demographic attribute predic-

162 tion [29]. Using web search engine datasets, Weber and Castillo [30] inferred gende

164 Donadioetal.Page 4 of30

165 Figure 1: Comparison between the spatial coverage of the insurance and the Union

166 of Swiss Cities datasets. (A) Total number of ins urance customers at zip code level.

167 The total number of customers is 1,341,328 which represents nearly 15% of the

168 Swiss population. (B) Number of inhabitants aggre gated at municipality level. The

169 total surface covered by those municipalities is almost 3890km

170 2

163

171 which represents

172 9.35% of the country's area.

173 ber and Jaimes [31] analyzed data rom di⊶erent ZIP codes enriched by US census 174 data and exploited them to highlight di⊶erences in user behaviour and search pat-175 terns among several demographic groups. Gender and age can also be inferred using 176 call detail records from smartphone devices over large populations. 177 Our work belongs to this line of work, however, i t explores for the first time, 178 to the best of our knowledge, the use of insuranc e customer records to predict 179 census variables. Insurance data have been mostly used to study the impact of 180 specific diseases [32,33], to propose models of customers' fraud detection [34,35], 181 to understand the correlation between census-based socioeconom for 182 and injury causes [36] or to evaluate disparities within hea lth care systems [37]. 183 3 Data 184 In this work, we make use of two main datasets: 1) the customer activity logs of a 185 Swiss insurance company and 2) the ocial socioeconomic indicators of all the Swiss 186 municipalities with more than 10,000 inhabitants, collected by the initiative called 187 Union des villes Suisse , i.e. Union of Swiss Cities that has published s 189 tatistics 190 on 173 Swiss municipalities every year since 2006: in this research we use data 191 related to the year 2017 [38]. The two datasets have di⊶erent spatial aggregations: 192 while the information on the insurance customers is at the zip code level; the socio-193 economic indicators have been collected at the municipality level. In aggregating the insurance data to match the spatial granulari 194 ty of the municipality data, we 195 restrict our analysis to the 170 municipalities that are prese nt in both datasets.

196 Figure1compares

- highlight di⊷erences in navigational and search patterns among several demographicgroups. Gender and age are also inferred using

call detail records from smartphone

176 devices over large populations.

- 177 Our work belongs to this line of work, however, i t explores for the first time,
- 178 to the best of our knowledge, the use of insuranc e customer records to predict
- 179 census variables. Insurance data have been mostly used to study the impact of
- specific diseases [32,33], to propose models of customers fraud detection [34,35], to
- the correlation between census-based socioeconom
  ic indicators
  and injury
- causes [36] or to evaluate disparities within hea  $182\$ lth care systems [37].
- 183 3 Data
  184 In this work, we make use of two
   datasets: 1) the customer activity logs of
  a Swiss
- insurance company and 2) the ocial socioeconomic indicators of all the Swiss
- 186 municipalities with more than 10,000 inhabitants, collected by the initiative called
- 187 Union des villes Suisse
- 188 [2]
- 189 that publishes statistics on 173 Swiss municipalities every
- 190 year since 2006. In this research, we
  focus on a cross-sectional snapshot for
  [38].
- 191 The two datasets have di⊶erent spatial aggregations: while the information on the
- insurance customers is at zip code level; the socioeconomic indicators have been
- To match the spatial granularity, we restrict our analysis to the 170 municipalities that are prese nt in both
- datasets. Figure1com-

196

195 pares
the spatial coverage in the two cases, while
Table1highlights di⊷erences and

The dataset contains the housing and vehicles ins urance policies of 1,341,328

anonymized customers of La Mobili`ere who were ac 200 tive

in 2017. La Mobili`ere is

201 [2]

202 https://uniondesvilles.ch/

203

## 204 Massaroetal.Page 5 of29

- 205 Table 1: Comparison between the two datasets.
- 206 Insurance DataUnion of Swiss Cities
- 207 FrequencyEvery yearEvery year
- 208 Spatial aggregationZip codeMunicipality
- 209 Data points3,185173
- 210 CostNot expensiveExpensive
- 211 DesignFor insurance marketingFor statistical anal
   vsis
- 212 AvailabilityPrivatePublic
- 213 a Swiss insurance group (brands: Die Mobiliar, La Mobili`ere, La Mobiliare) that is
- organized as a holding company headed by a cooper ative. The company was founded
- 215 in 1826, making it the oldest private insurance c ompany in Switzerland. Mobiliar
- is an all-insurance company operating exclusively
  in Switzerland and Liechtenstein.
- 217 With a market share of over 29%, it is the leader in the personal property insurance
- 218 market. Customers' details are aggregated at the level of the 170 municipalities for
- 219 which we have ocial census data using the postal code of provenance; this step
- 220 leaves us with 568,426 customers matching the geo graphical boundaries. For each
- 221 user, we have three types of information: i)demog
  raphic, e.g., age, gender, zip code
- of the residential area, employment status, civil status; ii)car insurance, e.g., how
- 223 many cars are insured, the type, brand and price
   of the vehicles, as well as the
- record of claims in terms of frequency and cost; iii)housing insurance, e.g., the
- number of private buildings or houses insured wit h the company, the price of the
- building and the logs of the claims. Table2summar izes the information available.
- 227 From this complete set of variables, we perform a feature engineering step in which
- 228 we select the variables of interest with the aid of a domain expert. As a result of
- 229 this process, we end up with 34 features as summa rized in Table3.Werefertothis
- 230 feature space in the modeling section of this wor k.
- 231 3.2 Swiss Census Data
- The ocial statistics for the Swiss municipalities

The dataset contains the housing and vehicles ins urance policies of 1,341,328

anonymized customers of La Mobili`ere who were ac

during 2017. La Mobili`ere

is a Swiss insurance group (brands: Die Mobiliar,

- 200 La Mobili`ere, La Mobiliare) that
- is organized as a holding company headed by a coo
- 201 perative. The company was
- 202 [2]
- 203 https://uniondesvilles.ch/

204

## 205 Donadioetal.Page 5 of30

- 206 Table 1: Comparison between the two datasets.
- 207 Insurance DataUnion of Swiss Cities
- 208 FrequencyEvery yearEvery year
- 209 Spatial aggregationZip codeMunicipality
- 210 Data points3,185173
- 211 CostNot expensiveExpensive
- 213 AvailabilityPrivatePublic

- 235 [3]
- 236
- 237 The report is published in the first quarter of t he year and it presents varied facets
- 238 of the urban life; we focus on six domains:popula tion,transportation,employment,
- 239 space and territory, housing and economy. We collect a total of 86 features for each
- 240 municipality: 11 indicators for transportation (t), 29 for population (p), 11 for em-
- 241 ployment (w), 8 for space and territory (s), 18 f
  or housing (h) and 9 for economy
- 242 (e). From the original dataset we focus on the ke y target variables that are not
- 243 redundant and can be a proxy for quality of life
  in cities, such as the unemploy-
- 244 ment rate [24], use of public transportation [39] or investment in culture [40]. The
- 245 complete list of selected target variables is lis
  ted below. As a result, we select two
- 246 characteristics for each domain: fraction of fore igners and rate of beneficiaries of
- 247 social assistance (p), number of private cars per person and the fraction of com-
- 248 muters using public transportation (t), unemploym ent rate and the unemployment
- 249 rate among women (w), percentage of areas covered by buildings or green areas (s),
- 250 [3]

- 251 https://www.bfs.admin.ch/bfs/fr/home/statistique
  s/catalogues-banques-donnees/publications/
- 252 ouvrages-synthese/statistiques-villes-suisses.htm
  1
- 254 Massaroetal.Page 6 of29
- 255 Table 2: Information for each customers in the in surance dataset
- 256 CatergoryVariable NameDescriptionVariable Type
- 258 JobStateEmployement statusString
- 259 CivilCivil StatusString
- 260 GenderGenderString
- 261 YearOfBirthYear of birthInteger
- 262 Own/RentIf own or rent an houseBoolean (Yes/No)
- 263 LangSpeaking languageString (French, German,
- 264 Italian)
- 265 NationNation of originString
- 266 ZIPZip code of residence5-digit code
- 267 Children0-26How many childrenInteger
- 268 CarsCarlCanton\*Canton where the car is registered String
- 269 CarlBrand\*Brand of the carString
- 270 Car1Price\*Price of the car (CHF)Integer
- 271 Car1ccm\*Cylinder capacityInteger
- 272 Car1ClaimsCt5Y\*Number of claims over the last 5  $\,$
- 273 years
- 274 Integer
- 275 Car1ClaimsSum5Y\*Sum of the money from the claims
- 276 over the last 5 years (CHF)
- 277 Float

- 214 Table 2: Information for each customers in the in surance dataset
- ${\tt 215} \>\>\> {\tt CatergoryVariable} \>\>\> {\tt NameDescriptionVariable} \>\>\> {\tt Type}$
- ${\tt 216~DemographicNmbrAnonymous~IDAlphanumeric}$
- 217 JobStateEmployement statusString
- 218 CivilCivil StatusString
- 219 GenderGenderString
- 220 YearOfBirthYear of birthInteger
- 221 Own/RentIf own or rent an houseBoolean (Yes/No)
- 222 LangSpeaking languageString (French, German,
- 223 Italian)
- $224\,$  NationNation of originString
- 225 ZIPZip code of residence5-digit code
- 226 Children0-26How many childrenInteger
- 227 CarsCarlCanton\*Canton where the car is registered
   String
- 228 CarlBrand\*Brand of the carString
- 229 Car1Price\*Price of the car (CHF)Integer
- 230 Car1ccm\*Cylinder capacityInteger
- 231 Car1ClaimsCt5Y\*Number of claims over the last 5
- 232 years
- 233 Integer
- 234 Car1ClaimsSum5Y\*Sum of the money from the claims
- 235 over the last 5 years (CHF)
- 236 Float

- 282 Integer
- 283 StandoffurnStandard of furnitureInteger (1-2-3-4-5)
- 284 RoomsNumber of roomsInteger
- 285 BuildZipZip Code of the insured building5-Digit c ode
- 286 BuildInsSumTotal sum of the insured values of
- 287 the building (CHF)
- 288 Integer
- 289 YearofcontrsYear of constructions of the building
  4-Digit Integer
- 290 TypeType of buildingString
- 291 HHaBClaimsCt5YNumber of claims over the last 5
- 292 years
- 293 Integer
- 294 HHaBClaimsSum5YSum of the money from the claims
- 295 over the last 5 years (CHF)
- 296 Integer
- 297 HHandBldPremPremium class of the buildingString
- vacancy rate and average area per inhabitant (h), 98 and municipal debt and fraction
- of investment in culture (e). The complete list o
- 299 f the selected features is reported
- 300 in Table4.
- Validation.As a validation step, we test the repr
- 301 esentativeness of the insurance
  - data along four dimensions: (a) total population,
- 302 (b) percentage of foreigners, (c)
- percentage of population aged 20-40, and (d) perc
- 303 entage of population aged 0-19.
- Figure2presents the Pearson's correlation between these variables extracted from
- the ocial census data and the La Mobili`ere custo 305 mers base respectively. A high
- degree of correlation ( $\rightarrow$ =0.91) can be observed for the total population, meaning
- that the insurance dataset mimics quite well the population distribution at the
  - municipality level. Focusing on the age facet, we
- 308 observe a strong relation for
  - the
- 309 case of
  - customers in the age range of 20-40 ( $\rightarrow$ =0.8) while the correlation disappears
  - (→=0.05) for customers aged 0-19. This behavior i
- 310 s expected since children and
  - teenagers are not usually the owners of insurance
- 311 policies on vehicles or houses.
- 312 Last, we observe a solid relation with the
  percentage of foreigners (→=0.6).
  - Spatial dependency.An aspect worth mentioning is
- 313 that the target variables we
  - are interested in are embedded in space. In this
- 314 direction, we compute the Moran's
- 315 I coecient to test spatial dependencies and, in p

- 241 Integer
- 242 StandoffurnStandard of furnitureInteger (1-2-3-4-5)
- 243 RoomsNumber of roomsInteger
- 244 BuildZipZip Code of the insured building5-Digit c
- 245 BuildInsSumTotal sum of the insured values of
- 246 the building (CHF)
- 247 Integer
- 248 YearofcontrsYear of constructions of the building 4-Digit Integer
- 249 TypeType of buildingString
- 250 HHaBClaimsCt5YNumber of claims over the last 5
- 251 years
- 252 Integer
- 253 HHaBClaimsSum5YSum of the money from the claims
- 254 over the last 5 years (CHF)
- 255 Integer
- 256 HHandBldPremPremium class of the buildingString founded in 1826, making it the oldest private ins
- 257 urance company in Switzerland.
- With a market share of over 29%, it is the leader
- 258 in the personal property insurance
- market. For each user, we have three classes of i
- 259 nformation: i)demographic, e.g.,
- age, gender, zip code of the residential area, em
- 260 ployment and civil status; ii)cars,
  - e.g., how many cars are insured, the brand and th
- 261 e price of the vehicles, as well
- as the record of the claims and the respective co
- 262 mpensations; iii)housing, e.g.,
  - the
- 263 number of
  - private buildings or houses insured, the price of the building and the logs
- of the claims. Customers' details are aggregated
- 264 **at** 
  - the level of the 170 municipal-
- ities for which we have ocial census data using t
- 265 he zip code as spatial reference;
- this step leaves us with 568,426 customers matchi
- 266 ng the geographical boundaries.
- 267 Table2summarizes the information available.

Table 3: Features extracted from the insurance da Table 3: Final set of features aggregated at the 319 tabase aggregated at the munici-270 municipality level. pality level. We extracted features from di⊶erent 320 customers information: i) Demographics (age, employment status, etc.), ii) Cars 321 (engine displacement - CCM, price of the vehicle, etc. ) and iii) Housing (year of 322 the building, number of claims, etc.) ] 323 CategoryNameDescription 271 CategoryNameDescription 324 Demographicf1:Unemployment rate 272 Demographicf1:Unemployment rate 325 f2:Average age in the municipality 326 f3:Fraction of owners (house) 274 f3:Fraction of owners (house) 327 f4:Fraction of foreigners 275 f4:Fraction of foreigners 328 f5:Average number of customers with at least one child child 329 f6:Market share 277 f6:Market share 330 f7:Fraction of women 278 f7: Fraction of women 331 f8:Number of customers divided by total customers 332 Carsf9:Average price of the cars 333 f10:95th percentile price of the cars 334 f11:Average year of the car 282 f11:Average year of the car 335 f12:5th percentile year of the car 336 f13:Average CCM of the car 284 f13:Average CCM of the car 337 f14:95th percentile CCM of the car 338 f15:Average number of claims per cars 339 f16:95th percentile number of claims of the car 340 f17:Average sum of claims of the car 341 f18:95th percentile number of price of the car 342 f19:Average premium of the car 290 f19:Average premium of the car 343 f20:Percent of insured cars 291 f20:Percent of insured cars 344 Housingf21:Average class of furniture 345 f22:95th percentile class of furniture 346 f23:Average number of rooms 294 f23:Average number of rooms 347 f24:95th percentile number of rooms 348 f25:Average building insured sum 349 f26:95th building insured sum 297 f26:95th building insured sum 350 f27:Average building year of Construction 351 f28:5th percentile building year of construction 352 f29:Average type of building 300 f29:Average type of building 353 f30:Average number of claims per building 354 f31:Average sum of claims per building 355 f32:95th sum of claims per building 356 f33:Average Insured Premium 304 f33:Average Insured Premium 357 f34:95th sum of insured premium 305 f34:95th sum of insured premium 358 Table 4: List of the selected 306 3.2 Swiss Census Data indicators for the 6 di⊶erent categories. 359 CategoryTarget Variable 309 Federal Statistical Oce (FSO) 310 [3] 311

273 f2:Average age in the municipality 276 f5:Average number of customers with at least one 279 f8:Number of customers divided by total customers 280 Carsf9:Average price of the cars 281 f10:95th percentile price of the cars 283 f12:5th percentile year of the car 285 f14:95th percentile CCM of the car 286 f15:Average number of claims per cars 287 f16:95th percentile number of claims of the car 288 f17:Average sum of claims of the car 289 f18:95th percentile number of price of the car 292 Housingf21:Average class of furniture 293 f22:95th percentile class of furniture 295 f24:95th percentile number of rooms 296 f25:Average building insured sum 298 f27:Average building year of Construction 299 f28:5th percentile building year of construction 301 f30:Average number of claims per building 302 f31:Average sum of claims per building 303 f32:95th sum of claims per building The ocial statistics for the Swiss municipalities 307 are collected and made available online within the initiativeStatistics of Swiss C 308 itiesthat is is the result of a collaboration between the Union of Swiss Cities and the

The report is published in the first quarter of t

the urban life; we focus on six domains:populatio

ment(w), space and territory(s), housing(h) and econ

312 he year and it presents varied facets of

313 n(p),transportation(t),employ-

the original dataset, we focus on the key target

that can be a proxy for quality of life in citie

317 variables that are not redundant and

```
318 s, such as the unemployment rate [24],
                                                                use of public transportation [39] or investment i
                                                           319 n culture [40]. As a result of this
                                                                process, we restrict the analysis to two indicato
                                                           320 rs for each domain: the fraction of
                                                                foreigners and the rate of beneficiaries of socia
                                                           321 l assistance (p), the number of private
                                                                cars per person and the fraction of commuters usi
                                                           322 ng public transportation (t), the
                                                                unemployment rate and the unemployment rate among
                                                           323 women (w), the percentage
                                                                of areas covered by buildings or green areas (s),
                                                           324 the vacancy rate and the aver-
                                                                age area per inhabitant (h), and the municipal de
                                                           325 bt and fraction of investment in
                                                                culture (e). The complete list of selected target
                                                           326 variables is summarized in Table4.
                                                           327 [3]
                                                                https://www.bfs.admin.ch/bfs/fr/home/statistique
                                                           328 s/catalogues-banques-donnees/publications/
                                                                ouvrages-synthese/statistiques-villes-suisses.htm
                                                           320 1
                                                           330
                                                           331 Donadioetal.Page 7 of30
                                                           332 Table 4: List of the target
                                                                indicators for the 6 di⇔erent domains.
                                                           333 DomainVariable
360 Populationp
                                                           334 Populationp
361 1
                                                           335 1
362 :Fractionofforeigners
                                                           336 :Fractionofforeigners
363 p
                                                           337 p
364 2
                                                           338 2
365 :Fractionofbeneficiariesofsocialassistance
                                                           339 :Fractionofbeneficiariesofsocialassistance
366 Transportationt
                                                           340 Transportationt
367 1
                                                           341 1
368 :Carsper1000inhabitants
                                                           342 :Carsper1000inhabitants
369 t
                                                           343 t
371 :Fractionofcommutersusingpublictransportation
                                                           345 :Fractionofcommutersusingpublictransportation
372 Employmentw
                                                           346 Employmentw
                                                           347 1
374 :Unemploymentrate
                                                           348 :Unemploymentrate
375 w
                                                           349 w
376 2
                                                           350 2
377 :Unemploymentrateamongwomen
                                                           351 :Unemploymentrateamongwomen
378 Space and Territorys
                                                           352 Space and Territorys
379 1
                                                           353 1
380 :Areacoveredbybuildings(%)
                                                           354 :Areacoveredbybuildings(%)
381 s
                                                           355 s
382 2
                                                           356 2
383 :Greenarea(%)
                                                           357 :Greenarea(%)
384 Housingh
                                                           358 Housingh
385 1
                                                           359 1
386 :Vacancyrate(%)
                                                           360 :Vacancyrate(%)
387 h
                                                           361 h
388 2
                                                           362 2
```

394 2 368 2 :Fractionofinvestmentinculture :Fractionofinvestmentinculture 396 measured overnspatial units and is given by: 370 3.3 Validation As a validation step, we test the representativen 371 ess of the insurance data along four dimensions: (a) total population, (b) percentage of foreigners, (c) percentage of population aged 20-40, and (d) percentage of populat ion aged 0-19. Figure2shows, for each dimension, a scatter plot and the correspond 374 ing Pearson's correlation coecient computed using the ocial census data and th 375 e La Mobili`ere customers base. A high degree of correlation (→=0.91) can be obse 376 rved for the total population variable, meaning that the insurance dataset mimi 377 cs quite well the population distribution at the municipality level. Focusing on age, we observe a strong relation for the case of customers in the age range 20-40 379 (→=0.8) while the correlation disappears (→=0.05) for customers aged 0-19. This 380 behavior is expected since children and teenagers are not usually the owners 381 of insurance policies on vehicles or houses. Last, we observe a solid relation with th 382 e percentage of foreigners (→=0.6). It is worth noting that socioeconomic processes o 383 ften manifest non-random spatial patterns that make close areas more similar than 384 distant ones. Moreover, spatial e⇔ects do not apply only to the case of neighbori 385 ng areas; on the contrary, a consistent body of literature in geography define spati 386 al relationships between aerial units as a function of distance [41]. Often this choice 387 depends on prior knowledge about the area under study or a conceptualization of th 388 e interactions between neighboring locations with regards to quantity under stud 389 y. In this work, we refer to the Moran's I statistic [42] to assess the presence o 390 f spatial autocorrelation in the census variables. Moran's I measures the global spatial 391 autocorrelation of an attributey measured overnspatial units using the following r 392 elation: 397 I= 393 I= 398 n 394 n 399 s 395 s 400 **0** 396 0 401 X 397 X 402 i 398 i 399 X 403 X 404 i 400 j 405 z 401 z 406 i 402 i

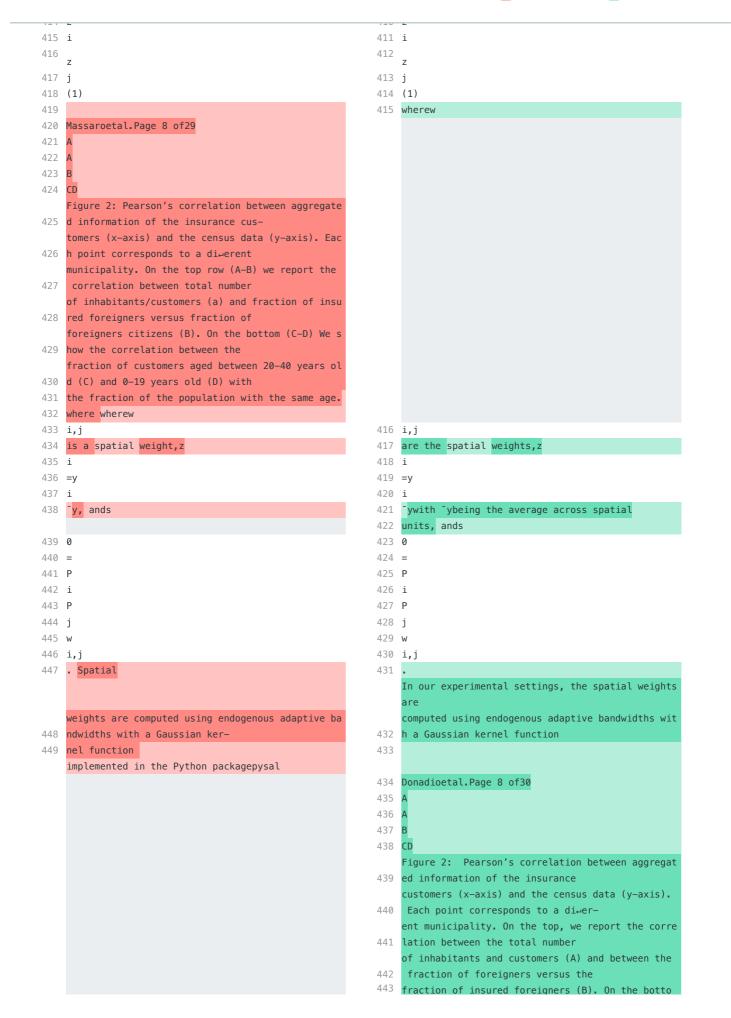
403 w

405 z

404 i,j

407 w 408 i,j

409 z



```
446 implemented in the Python packagepysal
450
                                                            447
    [4]
451
    . Table5shows how all
                                                            . Table5shows how all the selected
                                                                target variables are positively spatially autocor
452 the selected
                                                            449 related, ranging
                                                                fromI=0.56 to
    target variables are positively spatially autocor
    related, ranging
    from
                                                                I=0.8. This implies that municipalities that are
453 I=0.56 toI=0.8.
                                                            450 closer in space tend to share
                                                           451 similar socioeconomic characteristics.
454 4Methods
                                                            452 4Methods
455 In our analysis,
                                                            453 In
     we adopt multivariate linear regression to predi
                                                                this section, we describe the methodological step
                                                                s of our predictive pipeline. After
    ct
    the socio-economic
                                                            454 a features selection module,
456 indicators of
    interest. We present two spatially aware models t
                                                                 we adopt multivariate linear regression to predi
    hat we introduce to
                                                                ct
                                                                the
    capture the geographical dependencies emerging in
    our problem.
                                                            455 socioeconomic indicators of
457
                                                                interest using two spatially-aware models that ca
                                                            456 the global and local geographical dependencies.
458 [4]
459 https://pysal.readthedocs.io/en/latest/
                                                            458 https://pysal.readthedocs.io/en/latest/
460
                                                           459
461 Massaroetal.Page 9 of29
                                                           460 Donadioetal.Page 9 of30
462 Table 5: Moran's I coecients for the main target
                                                            461 Table 5: Moran's I coecients for the main census
     variables.
                                                                 variables.
463 VariableMoran's I
                                                           462 VariableMoran's I
464 p
                                                           463 n
465 1
                                                           464 1
466 :Fractionofforeigners0.7
                                                            465 :Fractionofforeigners0.7
467 p
                                                           466 p
                                                           467 2
468 2
469 :Fractionofbeneficiariesofthesocialassistance0.73
                                                            468 :Fractionofbeneficiariesofthesocialassistance0.73
470 t
                                                            469 t
471 1
                                                           470 1
472 :Carsper1000inhabitants0.56
                                                           471 :Carsper1000inhabitants0.56
                                                           473 2
475 :Fractionofcommutersusingpublictransportation0.76
                                                           474 :Fractionofcommutersusingpublictransportation0.76
476 w
                                                           475 w
477 1
                                                           476 1
478 :Unemploymentrate0.8
                                                           477 :Unemploymentrate0.8
479 w
                                                           478 w
                                                           479 2
481 :Unemploymentratebetweenwomen0.79
                                                           480 :Unemploymentratebetweenwomen0.79
                                                           481 s
482 s
483 1
484 :Buildingarea(%)0.74
                                                           483 :Buildingarea(%)0.74
485 s
                                                           484 s
486 2
                                                           485 2
487 :Greenarea(%)0.64
                                                           486 :Greenarea(%)0.64
488 h
                                                           487 h
                                                           488 1
489 1
490 : Vacancyrate(%)0.66
                                                           489 : Vacancyrate(%) 0.66
491 h
                                                            490 h
```

```
497 e
                                                            496 e
498 2
                                                            497 2
499 : Fractionofinvestmentinculture0.7
                                                            498 : Fraction of investment inculture 0.7
500 4.1 Feature selection
                                                            499 4.1 Features selection
501 The first step in constructing cost-e⊶ective pred
                                                            500 The first step in constructing cost-e⊸ective pred
    ictors is to select the features that
                                                                ictors is to select the features that
502 will best predict a given
                                                            501 will best predict a given outcome
                                                                variable. For each of the socioeconomic indica-
    variable. For each of the
    variables in Table4, asubset
    of explanatory variables was selected from all th
                                                                tors in Table4, we select a subset of explanatory
503 e variables present in Table3by
                                                            502 variables from the initial pool
     the means of di⇔erent selection algorithms. Once
                                                                of covariates summarized in Table3using theLassoL
     the features are selected for each
                                                            503 arsICmodule
                                                                 available in
    variable, the spatial dimension of the problem wi
505 ll be integrated to the model using
                                                            504 scikit-learn
    two di⊷erent approaches, as described in subsecti
506 ons4.2and4.3. The following
    three di⊶erent feature selections approaches were
    explored:
    Simple [41]:standard OLS multivariate regression
508
    model where features are se-
509 lected according to theirp
510 value
511 with a threshold ofp
512
513
    value
514 =0.05.
    Lasso [42]:to reduce model complexity and to prev
515 ent overfitting we turn to the
    Least Absolute Shrinkage and Selection Operator
516 (LASSO) model that is well-suited
    for cases showing high levels of multicollinearit
    y. In particular, we use theLasso-
    LarsICclass available in scikit-learn
519 [5]
                                                            505 [5]
                                                                . To reduce model complexity and to prevent overf
    that relies on Least Angle Regression and
                                                            506 itting.LassoLarsIC
                                                                adopts the Least Absolute Shrinkage and Selection
                                                            507 Operator [43] (LASSO) model
521 the Bayes Information Criterion for model
     selection and to find a trade—o⊣between
                                                                for fit and it relies on theLeast Angle
522 the goodness of fit
     and the complexity of the model.
                                                                Regression[44] (LARS) and theBayes Infor-
    Interactions [43]:this variant is also based on L
523 asso but it considers the interac-
                                                            509 mation Criterion[45] (BIC) for model
                                                                selection, trying to find the right trade-o⊷
    tions among the features by ensuring that the num
524 ber of predictorsn
                                                            510 between fitting performance
                                                                 and the complexity of the model.
                                                                Since variable selection methods may su⊷er from m
525
                                                            511 odel instability or potential
                                                                bias in parameter estimates and confidence interv
526 are between
                                                            512 als (especially relevant in explana-
                                                                tory modeling), we implemented the methodology an
527 defined boundariesb1 n
                                                            513 d practical suggestions pro-
                                                                posed in [46,47] to control for these e⊶ects. In
528
                                                            514 particular, we aim at estimating the
                                                                stability of the selection procedure to random pe
529 €b2, whereb1 = 2 andb2 = 15. This is because we
                                                            515 rturbations of training samples. We
need to ensure in the interactions step that the
                                                            516 implemented a subsampling without replacement rou
```

even in cases where the classical bootstrap fails 533 519 [48]. We performed 500 subsampling iterations and we computed the stability estimato 534 **b1** 520 r proposed by Nogueira et al. [49] that is a frequency-based statistics, ranging 0 t 535 521 o 1 and monotonically increasing as the stability of the feature selection grows. The <= 170. 522 536 idea is that the stability measure is a linear function of the sample variances with 523 a strictly negative slope. According to the proposed framework, stability values a 524 bove 0.75 represent an excellent agreement of the feature sets beyond chance, betw een 0.75 and 0.4 intermediate to good agreement, while values below 0.4 represent 526 a poor agreement. 537 4.2 Spatial Lag Model 527 4.2 Spatial Lag Model In a Spatial Lag Model (SLM), the first and most To characterize the influence of neighboring spat 538 straightforward way to introduce 528 ial units, we implement a Spatial space is by "spatially lagging" the dependent var Lag Model [50] (SLM) where the local e⊣ects are e iable. One must then treat it as 529 ncoded adding a term that contains a spatially lagged version of the dependent 540 an endogenous variable, this is known as 530 variable. SLM is an instance of a Spatial Autoregressive Model. Formally. Spatial Autoregressive Models where the additiona 541 This can be expressed in matrix notation, as 531 1 term follows [44]: is treated as an endoge-532 nous variable. More formally, this can be expressed in matrix notation, as follows: 533 **[5]** https://scikit-learn.org/stable/modules/generate 534 d/sklearn.linear\_model.LassoLarsIC.html 535 536 Donadioetal.Page 10 of30 542 y=←+X+W y+∞(2) 537 y= ↔ + X+W y+ ⇒ (2) 543 wherevis the vector of observations on the depend 538 wherevis the vector of observations on the depend ent variable, Xis the matrix of ent variable, Xis the matrix of 544 observations on the exogenous variables, Wis the s 539 observations on the exogenous variables, Wis the s patial weighting matrix of known patial weighting matrix of known 545 constants, is the vector of regression parameters 540 constants, is the vector of regression parameters andis the scalar autoregressive andis the scalar autoregressive 546 [5] 547 https://scikit-learn.org/stable/modules/generate d/sklearn.linear\_model.LassoLarsIC.html 548 549 Massaroetal.Page 10 of29 541 parameter. The variableWyis typically known as th 550 parameter. The variableWyis typically known as th e spatial lag ofy. The weighting e spatial lag ofy. The weighting 542 scheme will determine how the spatial dimension o 551 scheme will determine how the spatial dimension o f the problem is incorporated in f the problem is incorporated in 552 the model. A k-nearest neighbors (KNN) scheme 543 the model. We adopt a was chosen, and the number of k-nearest neighbors (KNN) scheme and we optimize the nearest neighbors to include in the scheme was op hyperparameterkwith a grid search approach that a 553 timized by maximizing 544 ims at finding the optimal the average 554 R 545 value that maximizes the averageR 555 2 over all the dependent across all the dependent variables. We use the

variables and found to be equal to ten neighbors.

- 559 coecients the PySAL package was used [45]. 560 4.3 Geographical Weighted Regression 561 Geographical Weighted Regression (GWR) is a local form of linear regression used
- 562 to model spatially varying relationships. In regr ession analysis, we try to explain
- 563 the variations of a dependent variable using a su
  - ite of uncorrelated and normally
- 564 distributed independent variables. The strength a nd direction of association is in-
- 565 dicated by the regression coecients, with one co ecient given for each variable in
- 566 the dataset. In GWR, instead of one global coecie nt for each variable, coecients
- 567 are able to vary according to space. This spatial variation in coecients can reveal
- 568 interesting patterns which otherwise would be mas ked. The general formula for a
- 569 GWR is an extension of Equation3where one regress ion is calculated for each point
- 570 using spatial weights.
- 571 y
- 572 i
- 573 =
- 574 **0**i
- 575 +
- 576 m 577 X
- 578 i=1
- 579
- 580 i,j
- 581 X
- 582 i.i
- 583 +e
- 584 i 585 . (3)
- 586 The indexiindicates the location of the city of i nterest. GWR basically fits a set
- 587 ofcoecients for each location:
- 588 589 i
- 590 = (X
- 591 T
- 592 W
- 593 i
- 594 X)
- 595 1
- 596 X
- 597 T 598 W
- 599 i
- 600 y, (4)
- 601 whereW

602 i

- 603 is the diagonal matrix of the spatial weights, un ique to locationi. There
- 604 are various schemes for calculating the weights, nearest neighbors, cubic or expo-
- 605 nential kernels. The weights are computed by the following:

- 549 4.3 Geographical Weighted Regression
- 550 Geographical Weighted Regression (GWR) is a local form of linear regression used
- 551 to model spatially varying relationships. In regr ession analysis, we try to explain
- 552 the variations of a dependent variable using a su
  - ite of uncorrelated and normally
- 553 distributed independent variables. The strength a nd direction of association is in-
- 554 dicated by the regression coecients, with one co ecient given for each variable in
- 555 the dataset. In GWR, instead of one global coecie nt for each variable, coecients
- 556 are able to vary according to space. This spatial variation in coecients can reveal
- 557 interesting patterns which otherwise would be mas ked. The general formula for a
- 558 GWR is an extension of Equation3where one regress ion is calculated for each point
- 559 using spatial weights.
- 560 v
- 561 i
- 562 =
- 563 **0**i
- 564 +
- 565 m
- 566 X 567 i=1
- 568
- 569 i.i
- 570 X
- 571 i,j
- 572 +e
- 573 i 574 . (3)
- 575 The indexiindicates the location of the city of i nterest. GWR basically fits a set
- 576 ofcoecients for each location:
- 577
- 578 i
- 579 = (X
- 580 T
- 581 W
- 582 i
- 583 X)
- 584 **1**
- 585 X
- 586 T
- 587 W
- 588 i
- 589 y, (4)
- 590 whereW
- 591 i
- 592 is the diagonal matrix of the spatial weights, un ique to locationi. There
- 593 are various schemes for calculating the weights, nearest neighbors, cubic or expo-
- 594 nential kernels. The weights are computed by the following:

```
611 i,j
                                                            600 i,j
612 b
                                                            601 b
613 ),(5)
                                                            602 ),(5)
614 whered
                                                            603 whered
615 ii
                                                            604 ii
616 is the Euclidean distance between municipalitiesi
                                                            605 is the Euclidean distance between municipalitiesi
    andjandbis the
                                                                andjandbis the
                                                                bandwidth of the kernel that has to be chosen. Fo
    bandwidth of the kernel that has to be chosen. Fo
617
                                                            606
    r each municipality we
                                                                r each municipality we
    calculate
                                                                calculate a
                                                                vector of weights and then regress using the fina
618 a
                                                            607 l formula (whether linear
                                                                or with in-
    vector of weights and then regress using the fina
    l formula (whether linear
619 with interactions),
                                                            608 teractions),
     in order to estimate all indicators for each mun
                                                                 in order to estimate all indicators for each mun
    icipality.
                                                                icipality.
                                                                The parameters
    The
620 parameters forming the GWR will be the focus of
                                                            609 forming the GWR will be the focus of the
    some analysis because of the non-
                                                                 analysis because of the non-stationarity
    stationarity of the problem. It is interesting to
621 explore how the influence of certain
    explanatory variables changes from city to city a
622 nd whether there are underlying
    tendencies. We estimated a fixed bandwith for eve
623 ry city and every model. The
    optimum bandwidth can be estimated by minimizing
    an information criterion; in
624
    practice, a corrected version of the AIC is used,
625 which unlike basic AIC is a function
626
                                                            610
627 Massaroetal.Page 11 of29
                                                            611 Donadioetal.Page 11 of30
                                                                of the problem. It is interesting to explore how
628 of sample size [46]. Thus for
                                                            the influence of certain explanatory
     a GWR model with a bandwidthb, theAIC
                                                                variables changes from city to city and whether t
                                                            613 here are underlying tendencies. We
                                                                estimate a bandwidth for each city and model. The
                                                            614 optimal bandwidth is estimated
                                                                by minimizing an information criterion; in practi
                                                            615 ce, we adopt a corrected version
                                                                of the AIC that, in contrast with the original de
                                                            616 finition, is a function of sample
                                                            617 size [51]. In more details, in
                                                                 a GWR model with a bandwidthb, theAIC
629 c
                                                            618 c
630 is given
                                                            619 is given
631 by:
                                                            620 by:
632 AIC
                                                            621 AIC
633 c
                                                            622 c
634 =2nln ^+nln 2:+n
                                                            623 =2nln ^+nln 2:+n
635 n+tr(S)
                                                            624 \text{ n+tr(S)}
636 n2tr(S)
                                                            625 n2tr(S)
637 ,(6)
                                                            626 ,(6)
638 where ^is the estimated standard deviation of the
                                                            627 where ^is the estimated standard deviation of the
    residuals; nis the number of
                                                                residuals; nis the number of
639 observations andtr(S) is the trace of the hat mat
                                                            628 observations andtr(S) is the trace of the hat mat
    rixS. The hat matrix is the
                                                                rixS. The hat matrix is the
```

629 projection matrix from the observedyto the fitted

640 projection matrix from the observedyto the fitted

644 = X633 **=X** 645 i 634 i 646 (X 647 T 636 T 648 W 637 W 649 i 638 i 650 X) 639 X) 651 **1** 640 1 652 X 641 X 653 T 642 T 654 W 643 W 655 i 644 i 656 (7) 4.4 Cross Validation 657 646 4.4 Evaluation To test the performance of the predictive pipelin The regression parameters and hyper-parameter sel ection is performed on the train-647 e we refer to an out-of-sample valiing set (80%) using cross-validation and then onc dation where the estimation of the regression par 659 e the model is fitted, the final 648 ameters and the hyper-parameter tuning are performed on a training set using a 8 performance is tested on the remaining 20% of the data. The municipalities were 0% split and cross-validation, while partitioned into three classes: communes with les the predictive performance is tested on the holdout data (remaining 20%). To cope s than 25k inhabitants, those with between 25 and 100k, and those with a population with the heterogeneity of the population distribu 662 higher than 100k. Note that 651 tion in our sample and to allow to there are five cities that belong to the last cla train and test the models with a sample that is r 652 epresentative of the entire specss: Z"urich, Gen`eve, Basel, Lausanne and Bern, which represent the five main Swiss cit trum of population size, we implement a stratific ies. In each round of the cross-653 ation approach. It is worth noting 664 validation, the procedure ensures that each fold that using a random sampling strategy instead, we 665 is a fair representation of the whole 654 could end up in the cross validadistribution. A training set with 80% of the poin tion procedure with splits that contain only high 666 ts and a validation set with the 655 ly or low populated municipalities, remaining 20% is selected. The training set is us introducing a bias in the evaluation pipeline. In ed to calibrate the parameters for this direction, we partition the muthe regression and the hyper-parameters of the GW nicipalities in three classes: communes with less than 25k inhabitants, those between 668 R (bandwidth selection) and SLM (spatial lags). The validation set is used to 25k and 100k, and those with a population higher 669 test the ability of the model to be 658 than 100k. Note that there are generalized to unknown locations. We repeat this five cities that belong to the last class: Z"uric 659 h, Gen`eve, Basel, Lausanne and Bern, procedure five times by ensuring 670 that the municipalities in the third class appear which represent the five main Swiss cities. In ea 671 only once for each training phase 660 ch round of the cross-validation, and that the same municipalities don't appear in the procedure ensures that each fold is a fair re 672 both the training and i the valida-661 presentation of the whole distribution balancing the three classes. We adopt a 5-673 tion phase. In the results section 662 cross validation accordingly. In we report the average and the standard deviation 663 the results section, 674 of the models performances we report the average and the standard due to cross-validation. deviation of the models 664 performance due to cross-validation. 675 5Results 665 5Results 676 In the first part of the section, we show 666 In the first part of this section, we present the results of the features selection process. the results of the features selection process for each of the target indicators. After In particular, in Table6we report the significant features for each target variable applying theLassoLarsICmethod 668 the number of selected features spans from using the Lasso method. The 2 (for the Housing indicatorh number of selected features varies between 2 (for the

679

housing variableh

```
we present the selected features for each model.
                                                       Overall, we observe a fair degree
                                                            of robustness to random perturbations with the me
                                                       675 asure of stability that varies
                                                       676 across dimensions. In particular,p
                                                       678 (0.82) andt
                                                       679 2
                                                       680 (0.77) show the highest stability
                                                       681 that reaches an excellent level of agreement; w
                                                       682
                                                       683 ,h
                                                       684
                                                       685
                                                       686 1
                                                       687 ,p
                                                       688 2
                                                       689 ,e
                                                       690 2
                                                       691
682 2
                                                       692 2
683 ). In the Population cate-
                                                       693 ,t
    gory, the per percentage of foreigners (p
                                                       694 1
                                                       695 ,s
                                                       696 1
                                                       697 , andh
                                                       698 1
                                                       700 Donadioetal.Page 12 of30
                                                            Table 6: Summary of the results of the Spatial La
                                                       701 g Model (SLM) for the target
                                                            indicators in the domains Population, Transportat
                                                       702 ion, and Employment.
                                                            VariableFeaturesCoecient Probability [0.025]
                                                       703 0.975]
                                                       704 plIntercept-0.0039.48E-01 -0.103 0.097
                                                            Fraction of foreignersf3-0.1191.24E-01 -0.271
                                                       705 0.034
                                                       706 f40.4601.65E-12 0.332 0.589
                                                       707 f70.2651.06E-04 0.130 0.400
                                                       708 f23-0.1091.74E-01 -0.267 0.049
                                                       709 Wdepvar0.0194.24E-01 -0.028 0.067
                                                       710 p2Intercept-0.0098.49E-01 -0.101 0.083
                                                       711 Fraction off10.2981.23E-07 0.187 0.409
                                                            beneficiaries off9-0.4171.39E-02 -0.752 -0.08
                                                       712 2
                                                            social assistancef130.1344.45E-01 -0.213
                                                                                                     0.4
                                                       713 82
                                                       714 f150.1571.51E-02 0.029 0.284
                                                       715 f23-0.3276.13E-10 -0.431 -0.223
                                                       716 Wdepvar0.0321.49E-01 -0.012 0.076
                                                       717 t1Intercept0.0088.84E-01 -0.101 0.117
                                                       718 Cars perf1-0.1372.94E-02 -0.260 -0.013
                                                       719 1000 inhabitantsf60.3811.48E-09 0.256 0.505
                                                       720 f7-0.1191.08E-01 -0.264 0.027
                                                        721 f8-0.1304.98E-02 -0.261
                                                       722 f170.0762.95E-01 -0.068 0.221
                                                       723 f190.1651.69E-02 0.029 0.301
                                                       724 f200.0683.69E-01 -0.082 0.217
                                                        725 f21-0.1384.38E-02 -0.272 -0.003
```

```
730
    -0.138
731 using publicf6-0.1011.63E-02 -0.184 -0.018
732 transportationf20-0.2611.39E-07 -0.358 -0.163
733 f220.1201.39E-02 0.024 0.217
734 f250.0913.64E-02 0.005 0.176
    Wdepvar0.1193.23E-16 0.090 0.148
736 w1Intercept-0.0088.62E-01 -0.097
                                       0.082
737 Unemployment ratef10.2201.26E-04 0.106 0.333
738 f40.1876.07E-04 0.079 0.295
739 f70.2061.18E-03 0.081 0.332
740
    f13-0.1251.25E-02 -0.225 -0.026
741 f23-0.2121.36E-04 -0.322 -0.102
    Wdepvar0.0541.26E-02 0.011 0.096
743 w2Intercept-0.0127.74E-01 -0.096
                                       0.071
744 Unemployment ratef10.1514.60E-03 0.046 0.255
745 between womenf40.1503.05E-03 0.050 0.251
746 f6-0.0876.45E-02 -0.181 0.006
747 f70.1067.38E-02 -0.011 0.224
748 f9-0.2501.39E-04 -0.379 -0.120
749 f160.0335.13E-01 -0.066 0.132
750 f190.2866.83E-05 0.144 0.428
751 f23-0.1245.06E-02 -0.250 0.001
752 f33-0.2685.80E-06 -0.384 -0.151
753 Wdepvar0.0612.60E-03 0.021 0.101
    cover a range between good (0.72) and intermediat
754 e (0.43) stability (variables are
755 listed in decreasing order), whilee
756 1
757 (0.22) shows a poor agreement. This low value
    indicates how the model characterizing the munici
758 pal debte
759 1
760 is highly dependent
    on variations of the training set to define signi
761 ficant determinants. Consistently,e
762 1
    is also the indicator with the lowest performance
763 in the predictive task, indicating
    how the insurance data is not really able to capt
764 ure its behavior.
    Switching the focus on the predictive task, Table
765 6and Table7summarize the
    results of the Spatial Lag Model for all the indi
766 cators. We present the direction and
768 Donadioetal.Page 13 of30
    Table 7: Summary of the results of the Spatial La
769 g Model (SLM) for the target
    indicators in the domains Space and Territory, Ho
770 using, and Economy.
    VariableFeaturesCoecient Probability [0.025]
771 0.975]
772 s1Intercept-0.0325.94E-01 -0.150 0.086
    Building area (%)f3-0.2862.29E-04 -0.439 -0.1
773 33
774 f40.0752.88E-01 -0.064 0.213
775 f70.1384.31E-02 0.003 0.274
776 f250.1001.16E-01 -0.026 0.226
777 f310.1347.74E-02 -0.016 0.283
778 Wdepvar0.0702.53E-02 0.008 0.132
```

0.132

0.044

784 Wdepvar0.0886.98E-05

```
785 h1Intercept0.0395.54E-01 -0.091 0.169
                                                       786 Vacancy rate (%) f30.0436.29E-01 -0.132
                                                                                                    0.217
                                                       787 f200.1672.54E-02 0.019 0.314
                                                       788 Wdepvar0.1571.19E-04 0.076
                                                                                         0.237
                                                           h2Intercept0.0167.03E-01 -0.067
                                                       790 Average areaf1-0.1611.66E-03 -0.263
                                                       792 in square metersf30.1748.51E-03  0.043  0.305
                                                       793 f4-0.0305.68E-01 -0.133 0.073
                                                       794
                                                           f60.1181.06E-02 0.027 0.210
                                                       795 f210.1372.44E-02 0.017 0.256
                                                       796 f220.0523.70E-01
                                                                            -0.062 0.165
                                                       797 f230.3133.30E-06 0.180 0.446
                                                       798 f27-0.0821.07E-01 -0.182 0.018
                                                       799 Wdepvar0.0631.14E-04 0.031 0.096
                                                       800 e1Intercept-0.0148.31E-01 -0.143
                                                                                              0.115
                                                       801 Municipal debtf9-0.1062.04E-01 -0.270
                                                                                                   0.059
                                                       802 f160.0981.84E-01 -0.047 0.243
                                                           f270.2364.07E-04 0.104 0.368
                                                           Wdepvar0.0811.03E-01 -0.017 0.180
                                                       805 e2Intercept-0.0128.23E-01 -0.114
                                                           Fraction of investmentf10.1229.03E-02 -0.020
                                                       807 in culturef40.0029.72E-01 -0.138
                                                                                              0.143
                                                       808 f70.0009.98E-01 -0.144 0.144
                                                       809 f11-0.2091.59E-02 -0.380 -0.038
                                                       810 f120.0297.49E-01 -0.148
                                                                                     0.205
                                                       811 f21-0.1001.20E-01 -0.228
                                                                                     0.027
                                                       812 f23-0.1287.58E-02 -0.271 0.014
                                                       813 f260.1079.72E-02 -0.020 0.234
                                                       814 f340.0366.35E-01 -0.113 0.184
                                                       815 Wdepvar0.1318.81E-09 0.086
                                                                                        0.176
                                                           the intensity of the relations along with confide
                                                           nce intervals; significant determi-
                                                       817 nants are marked in bold.
                                                       818 In the Population domain, the fraction
                                                           of foreigners (p
                                                       819 1
686 ) is strongly related to the featuref
                                                       820 ) is positively (=0.46)
                                                       821 linked to the demographic featuref
687 4
                                                           that represents the fraction of foreigners cus-
688
    that
    characterizes the fraction of foreign customers
689
                                                       824 tomers of La Mobili`ere
     of La Mobili`ere as expected. This
                                                           and, in the same direction, to the fraction of
                                                           womenf
690 provides an additional validation on the
                                                       825
    representativeness of the dataset used in
                                                       826 (=0.265).
   this analysis.
     Moreover, the percentage of people that receive
                                                            Moreover, the percentage of people that receive
     social assistance (p
                                                            social assistance (p
692 2
                                                       827 2
693 )
694 is strongly linked to the unemployment ratef
                                                       829 is positively linked to the unemployment ratef
                                                           (=0.298), and the average num-
    , providing another reasonable ex-
    planation, and to the average number of roomsf
                                                       832 ber of claims per carf
                                                       833 15
                                                       834
                                                           (=0.157). We observe a negative relation with ave
```

```
the social class. The most relevant features in t
700 he Transportation category are not
                                                            840 as indirect proxy for the social class.
    of immediate interpretation. In fact, for the var
                                                                 In the Transportation domain, the number of cars
                                                            841 per 1000 inhabitantst
701 iablet
                                                            842 1
                                                            844
                                                                 a negative relation with the unemployment ratef
702 1
                                                            845 1
703 (cars per 1000 inhabitants),
                                                            846 (=0.137) and the average
704 the most relevant feature is the market share (f
                                                            847 class of furnituref
                                                            848
                                                                 (=0.138). A positive link is found with the marke
                                                            849
                                                                 t share
                                                            850
                                                            851 Donadioetal.Page 14 of30
                                                            852 f
705 6
                                                            853 6
                                                            854 (=0.381), the average premium of the carsf
706 ) while for the variablet
                                                            855 19
                                                                 (=0.165) and the average
                                                            857
                                                                years of constructionf
                                                            858 27
                                                                 (=0.154). Concerning the commuters that use publi
                                                            859
                                                            860 transportationt
707 2
                                                                we observe a negative link with the market sharef
708
    muters that use public transportation) the most r
    elevant feature is the fraction of
709
                                                            863
                                                            864 (=0.101)
710
711 Massaroetal.Page 12 of29
                                                            865 and the fraction of house ownersf
    Table 6: Results of the Lasso features selection
    for the di⇔erent target variables.
                                                            866 3
    The features a are order in an ascending of the p
                                                            867 (=0.239). This could be explained by
    -value (lowest to highest).
                                                                 the observation that individuals living in rental
    CategoryVariable Feature p-valueVariable Fe
714 ature p-value
                                                            868 houses show a higher frequency of
                                                                ride-sharing use and commuting using public trans
715 Populationp
                                                            869 portation than those who own
                                                                 their houses [53]. A positive relation is found w
716
                                                            870 ith the percent of insured carsf
717
    f4
                                                                 (=0.261), the 95th percentile of the class of the
                                                            872 insured furnituref
718 к-к-к-
719 1.02E-14p
720 2
                                                            873 2
721 f23
                                                            874 \ 2(=0.120),
                                                                and the average insured sum per buildingf
722
    ----
                                                            875
                                                            876 25
723 2.01E-09
                                                            877 (=0.091).
724 f7
                                                                 Focusing on the Work category, the unemployment r
725 к-к-к-
726 4.07E-06f1
727 <del>к-к-к-</del>
728 9.94E-09
729 f31.07E-01f15
730 к−к−
731 3.40E-03
732 f232.32E-01f9
```

```
to a set of demographics features, primarily the
739 <mark>к-к-к-</mark>
                                                             881
                                                                  fraction of foreignersf
740 2.33E-08t
                                                             882 4
741 2
                                                             883 (=0.187),
742 f3
                                                             884 the fraction of womenf
743 к-к-к
                                                                 (=0.206) and, as expected, the unemployment rate
744 1.69E-07
                                                            886
745 f1
                                                             887 the Mobili`ere customersf
746
747
    2.19E-02f22
748 к-к-к
749 9.42E-07
750 f19
751
752 2.41E-02f20
753 к-к-к-
754 3.55E-06
755 f87.93E-02f6
756
757
    4.35E-03
758 f72.85E-01f25
759
760 1.84E-02
761 f179.63E-02
762 f203.74E-01
763 Workw
764 1
                                                                 (=0.220). We observe an opposite relation with th
765 f1
                                                             889 e
766 к-к-к-
                                                             890 average CCM of the carsf
767 6.77E-06w
                                                             891 13
                                                             892 (=0.125), and the average number of roomsf
                                                             893 23
                                                                 (=0.212). For the case of the women unemployment
                                                             894 ratew
768 2
                                                             895 2
769 f19
                                                             896 , the dominant
                                                                 features are related to the economic characterist
770 к-к-к-
                                                             897 ics of the items insured, being
771 6.03E-08
                                                             898 the average premium of the carsf
772 f7
                                                             899 19
773 к-к-к-
                                                             900 (=0.286) in a positive relation and the
774 1.78E-05f9
                                                             901 average price of the carsf
775 к−к−
776 4.34E-06
                                                             903 (=0.250) or average insured premiumf
777 f4
                                                            904 33
778 к-к-к-
                                                             905 (=
                                                                 0.268) linked negatively. These observations tend
779 2.59E-05f33
                                                                 to indicate gender di⇔erences in
                                                             906
780 к-к-к-
                                                             907 the insurance sector. The fraction of foreignersf
781 1.39E-05
                                                             908
782 f23
                                                             909 (=0.150) and the customers
783 ₩-₩-
                                                             910 unemployment ratef
784 7.80E-04f1
785
786 5.02E-04
787 f13
788
789 1.11E-02f4
```



```
855 1.54E-02
856 f4
857
858 2.36E-01
859 f227.89E-01
860 Economye
861 1
862 f27
863 к-к-
864 5.47E-04e
                                                           931
865
866 f1
                                                           932 (=0.142), the fraction of
867
                                                           933 house ownersf
868 1.17E-06
    f9
869
870
871 1.05E-02f21
872 ⊭
873 6.10E-03
874 f16
875
876 4.07E-02f34
877
878 1.28E-02
879 f4
880
881 2.48E-02
882 f11
883
884 3.49E-02
885 f236.60E-02
886 f71.46E-01
887 f252.12E-01
    f129.72E-01
888
    ***<0.0001,**<0.001and *<0.05.
890 customers that own an house (f
891 3
                                                           934 3
    ). For a possible explanation in this direction,
                                                           935 (=0.174), the market sharef
892 it
    has been observed that individuals living in rent
893 al houses show a higher frequency
    of ridesharing use and commuting using public tra
894 nsportation than those who own
    their houses [48]. In the Work category, the unem
895 ployment ratew
896 1
897 is primarily
    connected to a set of demographics features, e.
898 g.,f
899 6
                                                           936 6
                                                           937 (=0.118), the average class
900 andf
                                                           938 of furnituref
                                                           940 (=0.137) and the average number of roomsf
                                                           941 23
                                                           942 )(=0.313).
                                                           943 Higher values forh
                                                           944 2
                                                           945 corresponds to lower unemployment ratef
                                                           946 1
901 1
```

```
951 average year of constructions of the buildingsf
                                                           952 27
                                                           953 (=0.236) that is in accordance
                                                                with the literature where modern buildings have b
                                                           954 een considered a proxy for eco-
                                                                nomic status [54]. Moreover, the fraction of inve
                                                           955 stment in culturee
904 2
                                                           956 2
905 , the domi-
                                                           957 is negatively
906 nant features are related to the
                                                           958 connected to the average year of the carf
    economic characteristics of the
    objects insured, e.g.,
907
                                                           959 1
908
    19
                                                           960 1(=0.209).
                                                                It is worth noting that for a group of indicator
    average sum of class premium of the car orf
                                                           961 s, the corresponding predictive
                                                                models identify significant relations with expect
910
                                                           962 ed determinants: this is the case of
    average price of the cars orf
                                                           963 the pair (p
911
912
913
    average insured premium, that tend to indicate ge
914 nder di⊶erences in the insurance
    sector. Within the Space and Territory category b
915 oth variabless
916 1
                                                           964 1
917 building area
                                                           965
918 ands
                                                            966
                                                                ) where the fraction of foreigners is explained u
                                                           967 sing the information
                                                                on the nationality of La Mobili`ere customers. A
                                                           968 similar case happen for the pairs
                                                           969 (p
919 2
                                                           970 2
    green area are strongly connected to the fraction
    of house ownerf
920
                                                           971
921
922 that
    suggests a link between urban characteristics of
    neighborhoods and their average
923
     population. A similar observation applies to the
924
    Housing category especially for the
925
926 Massaroetal.Page 13 of29
    Figure 3: (A-B). Comparison between the spatial r
927 egression models (bars) and stan-
    dard multivariate linear regression (lines) for t
928 he di⊶erent features selection models.
    (A) Spatial Lag Model and (B) Geographical Weight
929 ed Regression. (C-D) Perfor-
    mance using stratified cross-validation for the f
930 ull (black triangles), the training
    and the validation sets respectively. (C) Spatial
931 Lag Model and (D) Geographical
932 Weighted Regression.
933 variableh
                                                           972 1
    vacancy rate. Finally, in the Economy category th
935 e municipal debte
                                                           973 ), (w
                                                            974 1
936 1
```

940 iterature where modern buildings have been considered a proxy for the economic sta 941 tus of a city [49]. The fraction 942 of investment in culturee 943 2 978 2 is connected to the unemployment rate (f 979 945 1 980 1 ). However, we think that these not surprising re 946 ) and other 981 lations do not undermine the validity of the experimental measures of wealth such as the average class of f 947 urnituref 982 settings for several reasons. First, the considered models identify alternative predic 948 983 tors that are complementary and 21 cross-domain to the target indicators. Second, th 949 and the sum of the 984 e observation that a variable constructed from a sample of customers of an insuran insured premiumf 985 ce company is able to predict a 950 census indicator at the national level is not tri vial. This represents another sugges-951 952 987 tion of the validity of the data collected as a proxy for socioeconomic s tatus. Third, After the first phase of variables selection, we 953 988 compare the performance of the spatially-aware models with a standard multivaria te linear regressor. Performance 989 Donadioetal.Page 15 of 30 is measured using the coecient of determination p 955 **seudoR** 990 991 BC 992 BC Figure 3: (A). Comparison of the performance betw 993 een the spatial regression models (GWR and SLM) and standard multivariate linear re 994 gression (OLS). (B-C) Performance using stratified cross-validation for the f 995 ull (black triangles), the training and the validation sets respectively. (B) Geograp 996 hical Weighted Regression and (C) 997 Spatial Lag Model. to quantitatively evaluate the impact of these no 998 t surprising variables, we compare the performance of the original models with a var 999 iation where we remove them. The 1000 accuracy in terms ofR 956 2 1001 2 1002 for both the SLM and GWR 957 for both the spacemodels remains substantially stable for all the indicators, with an average pe 958 agnostic and GWR 1003 nalty of 0.034 and 0.014 for SLM models. As shown in Figure3(A-B), both the geogra models outperform OLS across features selections and GWR, respectively. Refer to FigureA.13for a d methods and target variables, 1004 etailed comparison. After the analysis of determinants, we focus on c with a gain in performance up to 30%. It is worth 960 noting that GWR is able to 1005 omparing the performance of the achieve satisfactory results across categories wi global (SLM) and local (GWR) spatial models to a th values ranging from 0.49 fors 1006 standard multivariate linear 1007 regressor (OLS) to quantify to benefit of exploit

```
both the spatial models outperform OLS across tar
                                                           1011 get indicators, with a gain in
                                                                performance up to 30%. It is worth noting that GW
                                                           1012 R is able to achieve satisfactory
                                                                 results across categories with values ranging fro
                                                           1013 m 0.49 fors
                                                           1014 2
                                                           1015 to 0.83 in the case
                                                           1016 ofw
964 1
                                                           1017 1
965 orh
                                                           1018 orh
966
                                                           1019
967 . This provides a hint on the ability of
                                                           1020 . This provides a hint on the potential of
    insurance
                                                                 insurance customers data to
968 customers data to characterize socio-economic
                                                           1021 characterize socioeconomic processes embedded in
     processes embedded in space.
                                                                space. Allowing the relationships
969 Allowing the relationships
                                                           1022 between the independent and dependent variables
    between the independent and dependent variables
                                                                 to vary by locality to capture con-
970 vary by locality to capture contextual
                                                           1023 textual factors, GWR is useful as an exploratory
     factors, GWR is useful as an exploratory
                                                                 technique; however, its usefulness
971 technique; however, its usefulness
                                                           1024 as a prediction tool is debated when it comes to
    as a prediction tool is debated when it comes to
                                                                model generalizability. To test the
                                                           1025 ability to perform out-of-sample predictions we
972 model generalizability. To test the
    ability to perform out-of-sample predictions we
                                                                turn to stratified cross-validation as
                                                                 described in Section4.4. As shown in Figure3(B-
                                                           1026 C),
973 turn to
    cross-validation. We focus on the Lasso features
                                                                 we observe a decrease of the
     selection case that provides
    in general the best performance across categories
974 (see Figure3(A-B)). Since popu-
                                                           1027 overall performance;
                                                                 however, especially for certain target
                                                                 variables, we are still able
    lation varies broadly across municipalities, we h
975 ypothesize that training and testing
                                                           1028 to achieve a reasonable
                                                                 performance on the validation set, for example,h
    on instances of cities of very di⊷erent size coul
976 d introduce a deterioration of perfor-
977 mance. To cope with this e⇔ect, we implement a
    stratified 5 fold cross-validation in
    which the stratification process is regulated by
    population size. Figure3(C-D) shows
    the results of the regression task using stratifi
979 ed cross-validation. As expected, the
980 general performance deteriorates;
    however, especially for certain target
    variables,
981 we are still able to achieve a good
     performance on the validation set, e.g.,h
                                                           1029 2
982 2
983 =0.6,
                                                           1030 = 0.6,
984 w
                                                           1031 w
985 1
                                                           1032 1
986 =0.53,p
                                                           1033 =0.53,p
987 2
                                                           1034 2
988 = 0.49, andt
                                                           1035 = 0.49, andt
989 2
                                                           1036 2
990 =0.49.
                                                           1037 =0.49.
                                                                 The values of the performances of the models
```

1038 are also reported in the Appendix in Table??and T

to a set of baseline models in which each target 1042 indicator is predicted using the remaining variables from the census. For instanc 1043 e, let us model the fraction of 1044 foreignersp 1045 1 1046 using the explanatory variablesp 1047 2 1048 ,t 1049 **2** 1050 ,...,e 1051 2 1052 from Table4. In Figure4 we report a comparison between the performance of 1053 the census-based baseline and the insurance-based models for the cases of SML a 1054 nd GWR. We observe overall Figure 4: Comparison between the census-based and 1055 the insurance-based explanatory models for the SML and GWR cases. Positive a 1056 nd negative values mean, respectively, an increase or decrease in performance us 1057 ing La Mobili`ere data in comparison 1058 to the census baseline. a comparable performance using our approach, with 1059 the baseline having a positive delta of 0.019 on average across indicators. This 1060 is expected being the baseline based on ocial census data where cross-correlatio 1061 n e⇔ects are present. However, it is worth nothing that in our reference scenario t 1062 he census is not available and, as such, the baseline approach not feasible. The obs 1063 ervation that insurance customers records are able to achieve comparable results is 1064 an additional proof of the potential 1065 of this approach. 1066 6 Discussion  $1067\,$  In the first part of the paper, we showed how to predict a wide range of socioeco-1068 nomic indicators using insurance customers activi ty logs. In this section we shift the 1069 attention to a specific use case that has a stron g impact on urban mobility and citi-1070 zens well being: the relation between commuting a nd public transport (the variable 1071 t 1072 **2** 1073 in our settings). The use of public transportatio n is an important contributing 1074 factor to urban sustainability; it has a heavy en vironmental footprint reducing air 1075 pollution and trac congestion, among the others. It has also positive financial 1076 benefits for families and communities as a whole. higher level of security and direct 1077 positive e⊣ects on well-being and healthier habit

s. We chose transportation to exem-

1078 plify our data analysis as it is the third most i

- 993 6 Discussion
- 994 In the first part of the paper, we showed how to predict a wide range of socioeco-
- $995\,$  nomic indicators using insurance customers activi ty logs. In this section we shift the
- 996 attention to a specific use case that has a stron g impact on urban mobility and citi-
- $997\,$  zens well being: the relation between commuting a  $\,$  nd public transport (the variable
- 998 t
- 999 2
- 1000 in our settings). The use of public transportatio
   n is an important contributing
- 1001 factor to urban sustainability; it has a heavy en vironmental footprint reducing air
- 1002 pollution and trac congestion, among the others. It has also positive financial
- 1003 benefits for families and communities as a whole, higher level of security and direct
- 1004 positive e⊸ects on well-being and healthier habit s. We chose transportation to exem—
- 1005 plify our data analysis as it is the third most i

As such, the question of which variables are abl e to predict the use of 1009 public transport is a key issue. In Table7and in Table8we report the values of the predictors for the global and the GWR mode 1010 ls respectively. While for the 1011 global model, the parameters have the same values for ea ch municipality, for the 1012 **GWR** we reported the average values of the coecient, t he standard deviation, and 1013 minimum and maximum values. For the GWR, we observe a hi gh variability in the 1014 intercept, this is mainly due to the high spatial autocorrelation. The GWR adapts 1015 the intercept so it is closer to its neighbors, a nd thus achieves higher accuracy. More 1016 detailed diagnostic information on the regressio n, such as the kernel bandwidth is 1017 provided in Table9. Turning the attention to the coecients, we observe that the 1018 fraction of customers that own a house (f 1019 3 1020 ) is negatively correlated with the target 1021 variable: as expected, it has been observed that individuals living in rental houses 1022 show a higher frequency of ride sharing and publi c transportation adoption than 1023 house owners [48]. As expected, also the percentage of insured cars (f 1024 20 1025 ) is nega-1026 tive linked to the probability of commuting via p ublic transport; the more cars an 1027 individual possesses the less she turns on the pu blic system when it comes to mobil-1028 ity. Moreover, we observed a higher public transp ort adoption in major cities, e.g., 1029 Zurich, Basel, Bern and Geneva as shown in Figure4(b). This is consistent with our 1030 analysis, in fact, the fraction of house owners

(f

As such, the question of which variables are abl e to predict the use of 1082 public transport is a key issue. In Table6and in Table8we report the values of the predictors for the global and the GWR mode 1083 ls respectively (an analysis of the GWR statistics for the target variables is re 1084 ported in the Appendix from FigureA.1to FigureA.12and in TableA.3and TableA.4). 1085 While for the global 1086 1087 Donadioetal.Page 17 of30 model, the parameters have the same values for ea 1088 ch municipality, for we reported the average values of the coecient, t 1089 he standard deviation, and min-1090 imum and maximum values. For the GWR, we observe a hi gh variability in the 1091 intercept, this is mainly due to the high spatial autocorrelation. The GWR adapts 1092 the intercept so it is closer to its neighbors, a nd thus achieves higher accuracy. More 1093 detailed diagnostic information on the regressio n, such as the kernel bandwidth is 1094 provided in Table9. Turning the attention to the coecients, we observe that the 1095 fraction of customers that own a house (f 1096 3 1097 ) is negatively correlated with the target 1098 variable: as expected, it has been observed that individuals living in rental houses 1099 show a higher frequency of ride sharing and publi c transportation adoption than 1100 house owners [53]. As expected, also the percentage of insured cars (f 1101 20 1102 ) is nega-1103 tive linked to the probability of commuting via p ublic transport; the more cars an 1104 individual possesses the less she turns on the pu blic system when it comes to mobil-1105 ity. Moreover, we observed a higher public transp ort adoption in major cities, e.g.,

1106 Zurich, Basel, Bern and Geneva as shown in

(f

Figure5(b). This is consistent with our

 $\frac{1107}{}$  analysis, in fact, the fraction of house owners

1111 e (f 1112 6 1113 ). We believe that the market share feature is representative information becau 1114 se even if, the insurance company la Mobiliere is a national company, is not used e 1115 qually across the Swiss country because of the competition between di⊶erent insuran 1116 ce companies. Moreover, having an insurance contract is mandatory also for renti 1117 ng an apartment, the information of the market share of a given company can tell u 1118 s important information about a certain kind of population living in that area. 1119 One of the main characteristic of GWR is that the inferred relationships vary by lo 1120 cality, that implies each munici-1121 pality has a di⊷erent fitting performanceR 1123 and coecients. In Figure5we show spatial distribution of GWR accuracy in di⊷erent 1124 regions. Mapping the localR 1034 2 values could provide a useful tool to identify ar 1035 1126 eas where the independent variables might or might not explain the phenomenon under s 1036 VariableCoecientStd.Errorz-StatisticProbability 1127 tudy. This could be useful, for example, to identify contextual anomalies that ar 1128 e linked to specific characteristics 1037 Intercept0.1960.1151.6930.089 f3: Fraction of owners (house)-0.3290.141-2.3240. of a community. While we are able to achieve good 1038 020 1129 results in several cities, the performance for the Grisons and Ticino cantons are l 1039 f6: Market Share-0.3940.118-3.3350.000 1130 ow. These cantons are fairly small f20: Percent of insured cars-0.4270.138-3.0960.00 and isolated regions. For example, Ticino is high 1040 1 1131 ly influenced by the adjacency to f22: 95th percentile class of furniture 0.3640. Italy this influence is not captured by the mode 1132 l. Two of the main cities in Ticino; 1041 1332,7270,006 f25: Average Building Insured Sum0.2140.1221.765 Lugano and Belinzona have a very low use of publi 1042 0.078 1133 c transport as shown in Figure5. 1134 Another interesting aspect is that the local R 1135 1136 shows a clustered behavior, with adjacent areas having similar performance. Note t 1137 hat these clusters tend to match administrative boundaries and we can clearly dist 1138 inguish regions such as Lausanne, Basel and St. Gallen (light blue), central Switze 1139 rland (orange) and the Valais (dark 1043 Table 8: Summary statistics for the GWR parameter 1140 Table 8: Summary statistics for the GWR parameter s for predictingt s for predictingt 1044 2 1141 2 1045 . 1142 . 1046 VariableMeanSTDMinMax 1143 VariableMeanSTDMinMax 1047 Intercept0.0050.313-1.1690.339 1144 Intercept0.0050.313-1.1690.339 1048 f3: Fraction of owners (house)-0.3060.156-0.6490. 1145 f3: Fraction of owners (house)-0.3060.156-0.6490. 1049 f6: Market Share-0.2360.218-0.8960.202 1146 f6: Market Share-0.2360.218-0.8960.202 1050 f20: Percent of insured cars-0.3010.201-0.5860.54 1147 f20: Percent of insured cars-0.3010.201-0.5860.54

1054 Massaroetal.Page 15 of29 1055 A 1056 B 1057 Figure 4: (A) Spatial distribution of the local coecient o f determinationR 1058 2 1059 using 1060 GWR to predict the fraction of commuters using pu blic transportationt 1061 2 1062 .(B) 1063 Comparison between predicted and actual values of the percentage of commuters 1064 using public transport. 1065 prices are higher and people tend not to settle a nd start a family. Another variable 1066 that distinguishes rural from urban environments is the market share (f 1067 6 1068 ). Since it 1069 is the oldest insurance company in Switzerland, L a Mobiliere has reached customers 1070 all across the country; however, in major cities it has a lower market share, due to 1071 the fiercer competition with other companies and the higher incidence of short-term 1072 and foreigner dwellers. One of the main character istic of GWR is that the inferred 1073 relationships vary by locality, that implies each municipality has a di⊶erent fitting 1074 performanceR 1075 2 1076 and coecients. In Figure4we show spatial distribu 1077 accuracy in di⊶erent regions. Mapping the localR 1078 2 1079 values could provide a useful 1080 tool to identify areas where the independent vari ables might or might not explain 1081 the phenomenon under study. This could be useful, for example, to identify contex-1082 tual anomalies that are linked to specific charac teristics of a community. While we 1083 are able to achieve good results in several citie s, the performance for the Grisons 1084 and Ticino cantons are low. These cantons are fai rly small and isolated regions. For 1085 example, Ticino is highly influenced by the adjac ency to Italy this influence is not 1086 captured by the model. Two of the main cities in Ticino; Lugano and Belinzona have 1087 a very low use of public transport as shown in Fi gure4. Another interesting aspect 1088 is that the local R 1089 2 1090 shows a clustered behavior, with adjacent areas h aving similar 1091 performance. Note that these clusters tend to mat

ch administrative boundaries and

1151 Donadioetal.Page 18 of30

1152 A

1153 B

1154 Figure 5:

(A) Spatial distribution of the local coecient of determinationR

1155 2

1156 using

1157 GWR to predict the fraction of commuters using public transportationt

1158 2

1159 .(B)

1160 Comparison between predicted and actual values of the percentage of commuters

1161 using public transport.

ity in the communities leaving 1095 in the di⊶erent areas of the country. 1096 Limitations. The approach proposed in this paper h as few limitations that should 1097 be carefully discussed. In details: 1098 Sample bias. Even if we showed a fair level of re presentativeness along di⊷erent 1099 dimensions, the input dataset contains informatio n only on the fraction of popu-1100 lation that owns an insurance policy with a speci fic company. Several segments of 1101 1102 Massaroetal.Page 16 of29 1103 Table 9: Information on the GWR oft 1105 ; percentage of commuters using public 1106 transport. 1107 Diagnostic Information 1108 Spatial kernel:Fixed Gaussian 1109 Bandwidth used:29.030 1110 Residual sum of squares:36.026 1111 E⇔ective number of parameters (trace(S)): 41.95 1112 Degree of freedom (n - trace(S)):128.048 1113 Sigma estimate:0.53 1114 Log-likelihood:-109.337 1115 AIC:304.578 1116 ATCc: 334.533 1117 BIC:439,268 1118 R2:0.788

1119 Adj. alpha (95%):0.007

1120 Adj. critical t value (95%):2.723

- $1121\,$  the population are left out of the analysis, addi  $\,$  ng a validity bias in the results,
- 1122 especially for indicators that cover a wider spec trum of the society.
- 1123 Spatial granularity
  mismatch. Ocial statistics are available at the level of
- 1124 municipality and only for a subset of the commune s, while the insurance customers
- 1125 data provides information at the finer granularit y of zip codes. From one side, we
- 1126 have complete knowledge for a subset of the area s, while from the other side, a
- $1127\,$  partial view with a higher coverage. Our analysis is restricted to the intersection

- 1162 Table 9: Information on the GWR oft
- 1163 **2**
- 1164 ; percentage of commuters using public
- 1165 transport.
- 1166 Diagnostic Information
- 1167 Spatial kernel:Fixed Gaussian
- 1168 Bandwidth used:29.030
- 1169 Residual sum of squares:36.026
- 1170 E⊶ective number of parameters (trace(S)): 41.95
- 1171 Degree of freedom (n trace(S)):128.048
- 1172 Sigma estimate:0.53
- 1173 Log-likelihood:-109.337
- 1174 AIC:304.578
- 1175 AICc:334.533
- 1176 BIC:439.268
- 1177 R2:0.788
- 1178 Adj. alpha (95%):0.007
- 1179 Adj. critical t value (95%):2.723
- 1180 blue) in Figure5. This phenomenon might be linked to the inherent diversity in the
- 1181 communities leaving in the di⊷erent areas of the country.
- 1182 Limitations.The approach proposed in this paper h
  as few limitations that should
- 1183 be carefully discussed. In details:
- 1184 Sample bias. Even if we showed a fair level of representativeness along di-erent
- dimensions, the input dataset contains information only on the fraction of popu-
- lation that owns an insurance policy with a speci fic company. Several segments of  $% \left( 1\right) =\left( 1\right) \left( 1\right$
- $1187\,$  the population are left out of the analysis, addi  $\,$  ng a validity bias in the results,
- 1188 especially for indicators that cover a wider spec trum of the society.
- 1189 Spatial granularity
  mismatch. Ocial statistics are available at the level of
- 1190 municipality and only for a subset of the commune s, while the insurance customers
- 1191 data provides information at the finer granularit y of zip codes. From one side, we
- 1192 have complete knowledge for a subset of the area s, while from the other side, a
- 1193 partial view with a higher coverage. Our analysis is restricted to the intersection

- 1130 heterogeneity of social processes at a micro-leve l, e.g., neighborhoods in cities.
- 1131 Temporal evolution. In our analysis we currently focus on a static snapshot cov-
- 1132 ering a year of statistics. However, socioeconomi c conditions vary over time and in
- 1133 which extent and how fast this change is reflecte d into the insurance data records
- 1134 is something not explored yet.
- 1135 Data availability and privacy. The current approach is based on the assumption
- 1136 that customers data is available to the researche rs to tackle relevant challenges that
- 1137 have a broad social impact. This raises two main issues related to privacy and the
- 1138 compliance to the current legislation especially in the European framework, and
- 1139 the sharing policy. Proprietary data is usually e xploited for commercial advantages
- 1140 and profit within the organization, and not avail able to the broad scientific com-
- 1141 munity. To ground a methodology to model social p henomena on the availability of
- 1142 proprietary data that is not in control of the policy makers raises few concerns on
- 1143 the actual implementation in a real scenario.
- 1144 7 Conclusions
- 1145 In this paper we proposed 34 di⊶erent characteris tics of individual socio-economic
- 1146 behavior quantifiable through the dataset of anon ymized insurance customers, and
- 1147 then evaluated them on the example of Swiss munic ipalities. We showed that those
- 1148 quantities could be used for estimating economic performance of the regions in

the country, as proposed geographical regression 1149 models technique

## demonstrated

1150 **to** 

perform well on the validation samples for predicting major ocial

## statistical

1151 quantities

for di⊷erent categories such as Population, Tran sportation, Work,

## Space

1153

1152 and Territory, Housing and Economy on the level of Swiss municipalities. Moreover,

# 1154 Massaroetal.Page 17 of29

the same approach

is applicable in cases when ocial statistics is not available or is

inconsistent, for example when considering geographical units of a finer spatial scale

- 1197 Donadioetal.Page 19 of30
- 1198 heterogeneity of social processes at a micro-leve l, e.g., neighborhoods in cities.
- 1199 Temporal evolution. In our analysis we currently focus on a static snapshot cov-
- 1200 ering a year of statistics. However, socioeconomi c conditions vary over time and in
- 1201 which extent and how fast this change is reflecte d into the insurance data records
- 1202 is something not explored yet.
- 1203 Data availability and privacy. The current approach is based on the assumption
- 1204 that customers data is available to the researche rs to tackle relevant challenges that
- 1205 have a broad social impact. This raises two main issues related to privacy and the
- 1206 compliance to the current legislation especially in the European framework, and
- 1207 the sharing policy. Proprietary data is usually e xploited for commercial advantages
- 1208 and profit within the organization, and not avail able to the broad scientific com-
- 1209 munity. To ground a methodology to model social p henomena on the availability of
- 1210 proprietary data that is not in control of the po licy makers raises few concerns on
- $1211\,$  the actual implementation in a real scenario.
- 1212 7 Conclusions
- 1213 In this paper we proposed 34 di⊶erent characteris tics of individual socio-economic
- 1214 behavior quantifiable through the dataset of anon ymized insurance customers, and
- 1215 then evaluated them on the example of Swiss munic ipalities. We showed that those
- 1216 quantities could be used for estimating economic performance of the regions in

the country, as proposed geographical regression models technique

# demonstrated to

perform well on the validation samples for predic

1218 ting major ocial

## statistical quan-

1219 tities

for  $\text{di}_{\mbox{\tiny \it ele}}\text{erent}$  categories such as Population, Tran sportation, Work,

## Space and

- 1220 Territory, Housing and Economy at
  - the level of Swiss municipalities. This approach
- 1221 is applicable in cases when ocial statistics are not available or they are inconsis
  - tent, and the experimental pipeline demonstrated
- its ability to reach comparable performance to a scenario with complete knowledg
- 1223 e. Advantages
  - and disadvantages
- of local and global spatial regression models have been discussed extensively, high-

the next steps, we aim at applying the same appro 1159 ach also allows evaluating tem-

poral variation of socio-economic condition perfo 1160 rmance of the Swiss cities, which

is especially useful to study process of urbanization and gentrification. Finally, the proposed model can be further employed for estima

ting more specific characteristics
 of local quality of life of cities and neighbourh
oods.

1164 Author details

1165 **1** 

1166 HERUS Lab,

1167 '

1169 Lausanne, CH.

1170

1171 University of Turin,, Corso Svizzera 185, 10149 T urin, Italy.

1172

1173 ISI Foundation, Via Chisola 5, 10126

1174 Turin, Italy.

1175 References

1176 1.MacFeely, S.: The big (data) bang: Opportunitie
 s and challenges for compiling sdg indicators. Gl
 obal Policy10,

1177 121-133 (2019)

1178 2.IEAG, U.: A world that counts—mobilising the da ta revolution for sustainable development. New Yo rk: United

1179 Nations (2014)

1180 3.Struijs, P., Braaksma, B., Daas, P.J.: Ocial st atistics and big data. Big Data & Society1(1),

1181 2053951714538417 (2014)

1182 4.Pappalardo, L., Pedreschi, D., Smoreda, Z., Gia nnotti, F.: Using big data to study the link betw een human

1183 mobility and socio-economic development. In: 2015
 IEEE International Conference on Big Data (Big Data), pp.

1184 871-878 (2015). IEEE

1186 market flows. Technical report, National Bureau o f Economic Research (2014)

1187 6.Watmough, G.R., Marcinko, C.L., Sullivan, C., T
 schirhart, K., Mutuo, P.K., Palm, C.A., Svenning,
 J.-C.:

1188 Socioecologically informed use of remote sensing data to predict rural household poverty. Proceed ings of the

1189 National Academy of Sciences116(4), 1213-1218 (20

the modeling of temporal variations, which is especially useful to study processes

1228 of urbanization and gentrification.

We also aim at developing models for estimating attributes at finer geographical resolutions such 1229 as districts or neighborhoods.

1230 Author details

1231 **1** 

1232 HERUS Lab,

1233 ′

1234 Ecole polytechnique f'ed'erale de Lausanne, ENAC, IIE GR C1 455 (B'atiment GR) - Station 2, 1015 -

1235 Lausanne, CH.

1236 **2** 

1237 University of Turin,, Corso Svizzera 185, 10149 Turin, Italy.

1238

1239 ISI Foundation, Via Chisola 5, 10126

1240 Turin, Italy.

1241 References

1242 1.MacFeely, S.: The big (data) bang: Opportunitie
 s and challenges for compiling sdg indicators. Gl
 obal Policy10,

1243 121-133 (2019)

1244 2.IEAG, U.: A world that counts—mobilising the da ta revolution for sustainable development. New York: United

1245 Nations (2014)

1246 3.Struijs, P., Braaksma, B., Daas, P.J.: Ocial st atistics and big data. Big Data & Society1(1),

1247 2053951714538417 (2014)

1248 4.Pappalardo, L., Pedreschi, D., Smoreda, Z., Gia nnotti, F.: Using big data to study the link betw een human

1249 mobility and socio-economic development. In: 2015 IEEE International Conference on Big Data (Big Data), pp.

1250 871-878 (2015). IEEE

1252 market flows. Technical report, National Bureau o
 f Economic Research (2014)

1254 Donadioetal.Page 20 of30

1253

1255 6.Watmough, G.R., Marcinko, C.L., Sullivan, C., T schirhart, K., Mutuo, P.K., Palm, C.A., Svenning, J.-C.:

1256 Socioecologically informed use of remote sensing data to predict rural household poverty. Proceed ings of the

1257 National Academy of Sciences116(4), 1213—1218 (20 19)

- 1192 8.Giannone, D., Reichlin, L., Small, D.: Nowcasti ng: The real-time informational content of macroe conomic data.
- 1193 Journal of Monetary Economics55(4), 665-676 (200 8)
- 1194 9.Ginsberg, J., Mohebbi, M.H., Patel, R.S., Bramm
  er, L., Smolinski, M.S., Brilliant, L.: Detecting
  influenza
- 1195 epidemics using search engine query data. Nature4 57(7232), 1012 (2009)
- 1196 10.Zhou, K., Yang, S.: Understanding household en ergy consumption behavior: The contribution of en ergy big
- 1197 data analytics. Renewable and Sustainable Energy
   Reviews56,810-819(2016)
- 11.98 11.De Montjoye, Y.-A., Radaelli, L., Singh, V.K.,
   etal.:Uniqueintheshoppingmall:Onthereidentifiabil
   ityof
- 1199 credit card metadata. Science347(6221), 536-539 (2015)
- 1200 12.Ritter, J., Mayer, A.: Regulating data as property: a new construct for moving forward. Duke L.
  - & Tech. Rev.
- 1201 16.220(2017)
- 1202 13.Zhang, Y., Chen, X., Li, J., Wong, D.S., Li,
  H., You, I.: Ensuring attribute privacy protecti
  on and fast
- 1203 decryption for outsourced data security in mobile cloud computing. Information Sciences379,42-61(20 17)
- 1204 14.Sivarajah, U., Kamal, M.M., Irani, Z., Weerakk ody, V.: Critical analysis of big data challenges and analytical
- 1205 methods. Journal of Business Research70,263—286(2 017)
- 1206 15.Florescu, D., Karlberg, M., Reis, F., Del Cast illo, P.R., Skaliotis, M., Wirthmann, A.: Will 'b ig data'transform
- 1207 ocial statistics. In: European Conference on the QualityofOcial Statistics. Vienna, Austria, pp. 2-5 (2014)
- 1208 16.Schnore, L.F.: The socio-economic status of ci ties and suburbs. American Sociological Review, 7 6-85 (1963)
- 1209 17.Dennison, G., Dodd, V., Whelan, B.: A socio-ec onomic based survey of household waste characteri stics in the
- 1210 city of dublin, ireland. i. waste composition. Re sources, Conservation and Recycling17(3), 227—244 (1996)
- 1211 18.Cameron, S.: The economics of crime deterrenc
  e: A survey of theory and evidence. Kyklos41(2),
  301-323
- 1212 (1988)
- 1213 19.Galea, S., Freudenberg, N., Vlahov, D.: Cities and population health. Social science & medicine6 0(5),
- 1214 1017-1033 (2005)
- 1215 20.Clark, W.A.: Residential segregation in americ an cities: A review and interpretation. Populatio

- 1260 8.Giannone, D., Reichlin, L., Small, D.: Nowcasti ng: The real-time informational content of macroe conomic data.
- 1261 Journal of Monetary Economics55(4), 665-676 (200 8)
- 1262 9.Ginsberg, J., Mohebbi, M.H., Patel, R.S., Bramm
  er, L., Smolinski, M.S., Brilliant, L.: Detecting
  influenza
- 1263 epidemics using search engine query data. Nature4 57(7232), 1012 (2009)
- 1264 10.Zhou, K., Yang, S.: Understanding household en ergy consumption behavior: The contribution of en ergy big
- 1265 data analytics. Renewable and Sustainable Energy Reviews56.810-819(2016)
- 1266 11.De Montjoye, Y.-A., Radaelli, L., Singh, V.K.,
   etal.:Uniqueintheshoppingmall:Onthereidentifiabil
   ityof
- 1267 credit card metadata. Science347(6221), 536-539 (2015)
- 1268 12.Ritter, J., Mayer, A.: Regulating data as prop erty: a new construct for moving forward. Duke L.
  - & Tech. Rev.
- 1269 16.220(2017)
- 1270 13.Zhang, Y., Chen, X., Li, J., Wong, D.S., Li, H., You, I.: Ensuring attribute privacy protecti on and fast
- 1271 decryption for outsourced data security in mobile cloud computing. Information Sciences379,42—61(20 17)
- 1272 14.Sivarajah, U., Kamal, M.M., Irani, Z., Weerakk ody, V.: Critical analysis of big data challenges and analytical
- 1273 methods. Journal of Business Research70,263-286(2
- 1274 15.Florescu, D., Karlberg, M., Reis, F., Del Cast illo, P.R., Skaliotis, M., Wirthmann, A.: Will 'b ig data'transform
- 1275 ocial statistics. In: European Conference on the QualityofOcial Statistics. Vienna, Austria, pp. 2-5 (2014)
- 1276 16.Schnore, L.F.: The socio-economic status of ci ties and suburbs. American Sociological Review, 7 6-85 (1963)
- 1277 17.Dennison, G., Dodd, V., Whelan, B.: A socio-ec onomic based survey of household waste characteri stics in the
- 1278 city of dublin, ireland. i. waste composition. Re sources, Conservation and Recycling17(3), 227-244 (1996)
- 1279 18.Cameron, S.: The economics of crime deterrenc e: A survey of theory and evidence. Kyklos41(2), 301-323
- 1280 (1988)
- 1281 19.Galea, S., Freudenberg, N., Vlahov, D.: Cities and population health. Social science & medicine6 0(5),
- 1282 1017-1033 (2005)
- 1283 20.Clark, W.A.: Residential segregation in americ an cities: A review and interpretation. Populatio

- 1218 453(7196), 779 (2008)
- 1219 22.Podobnik, V., Ackermann, D., Grubisic, T., Lov
  rek, I.: Web 2.0 as a foundation for social media
  marketing:
- 1220 global perspectives and the local case of croati
  a. In: Cases on Web 2.0 in Developing Countries:
   Studies on
- 1221 Implementation, Application, and Use, pp. 342-37 9. IGI Global, ??? (2013)
- pooling with shareability networks. Proceedings o
   f the National Academy of Sciences111(37), 1329013294
- 1224 (2014)

- 1225 24.Sobolevsky, S., Massaro, E., Bojic, I., Arias,
  J.M., Ratti, C.: Predicting regional economic ind
  ices using big data
- 1226 of individual bank card transactions. In: 2017 IE EE International Conference on Big Data (Big Dat a), pp.
- 1313-1318 (2017). IEEE
- 1229 Massaroetal.Page 18 of29
- 1230 25.Fatehkia, M., Kashyap, R., Weber, I.: Using fa cebook ad data to track the global digital gender gap. World
- 1231 Development107,189-209(2018)
- 1232 26.Zagheni, E., Weber, I., Gummadi, K., et al.: L
   everaging facebook's advertising platform to moni
   tor stocks of
- 1233 migrants (2017)
- 1234 27.Dubois, A., Zagheni, E., Garimella, K., Weber,
   I.: Studying migrant assimilation through faceboo
   k interests. In:
- 1235 International Conference on Social Informatics, p
  p. 51-60 (2018). Springer
- 1236 28.Llorente, A., Garcia-Herranz, M., Cebrian, M.,
   Moro, E.: Social media fingerprints of unemployme
   nt. PloS one
- 1237 10(5), 0128692 (2015)
- 1238 29.Cesare, N., Grant, C., Hawkins, J.B., Brownste
  in, J.S., Nsoesie, E.O.: Demographics in social m
  edia data for
- 1239 public health research: does it matter? arXiv pre print arXiv:1710.11048 (2017)
- 1240 30.Weber, I., Castillo, C.: The demographics of w eb search. In: Proceedings of the 33rd Internatio nal ACM SIGIR
- 1241 Conference on Research and Development in Informa tion Retrieval, pp. 523-530 (2010). ACM
- 1242 31.Weber, C., Johnson, M., Arceneaux, K.: Genetic
  s, personality, and group identity. Social scienc
  e quarterly
- 1243 92(5), 1314-1337 (2011)
- 1244 32.Augustin, M., Reich, K., Glaeske, G., Schaefe r, I., Radtke, M.: Co-morbidity and age-related p revalence of
- 1245 psoriasis: analysis of health insurance data in g

- 1286 453(7196), 779 (2008)
- 1287 22.Podobnik, V., Ackermann, D., Grubisic, T., Lov
   rek, I.: Web 2.0 as a foundation for social media
   marketing:
- 1288 global perspectives and the local case of croati
  a. In: Cases on Web 2.0 in Developing Countries:
  Studies on
- 1289 Implementation, Application, and Use, pp. 342-37 9. IGI Global, ??? (2013)
- 1291 pooling with shareability networks. Proceedings o f the National Academy of Sciences111(37), 13290-13294
- 1292 (2014)
- 1293 24.Sobolevsky, S., Massaro, E., Bojic, I., Arias,
  J.M., Ratti, C.: Predicting regional economic ind
  ices using big data
- 1294 of individual bank card transactions. In: 2017 IE EE International Conference on Big Data (Big Dat a), pp.
- 1295 1313-1318 (2017). IEEE
- 1296 25.Fatehkia, M., Kashyap, R., Weber, I.: Using fa cebook ad data to track the global digital gender gap. World
- 1297 Development107,189-209(2018)
- 1298 26.Zagheni, E., Weber, I., Gummadi, K., et al.: L
  everaging facebook's advertising platform to moni
  tor stocks of
- 1299 migrants (2017)
- 1300 27.Dubois, A., Zagheni, E., Garimella, K., Weber,
   I.: Studying migrant assimilation through faceboo
   k interests. In:
- 1301 International Conference on Social Informatics, p
  p. 51-60 (2018). Springer
- 1303 10(5), 0128692 (2015)
- 1304 29.Cesare, N., Grant, C., Hawkins, J.B., Brownste in, J.S., Nsoesie, E.O.: Demographics in social m edia data for
- 1305 public health research: does it matter? arXiv pre print arXiv:1710.11048 (2017)
- 1306 30.Weber, I., Castillo, C.: The demographics of w
  eb search. In: Proceedings of the 33rd Internatio
  nal ACM SIGIR
- 1307 Conference on Research and Development in Informa tion Retrieval, pp. 523-530 (2010). ACM
- 1308 31.Weber, C., Johnson, M., Arceneaux, K.: Genetic
  s, personality, and group identity. Social scienc
  e quarterly
- 1309 92(5), 1314-1337 (2011)
- 1310 32.Augustin, M., Reich, K., Glaeske, G., Schaefe
   r, I., Radtke, M.: Co-morbidity and age-related p
   revalence of
- 1311 psoriasis: analysis of health insurance data in g

- 1247 Analysis of statutory health insurance data. Woun d Repair and Regeneration24(2), 434-442 (2016)
- 1248 34.Rawte, V., Anuradha, G.: Fraud detection in he alth insurance using data mining techniques. In: 2015
- 1249 International Conference on Communication, Inform ation & Computing Technology (ICCICT), pp. 1–5 (2 015).
- 1250 **IEEE**
- 1251 35.Nian, K., Zhang, H., Tayal, A., Coleman, T., L
  i, Y.: Auto insurance fraud detection using unsup
  ervised spectral
- 1252 ranking for anomaly. The Journal of Finance and D ata Science2(1), 58-75 (2016)
- 1253 36.Bell, N., Arrington, A., Adams, S.A.: Census-b ased socioeconomic indicators for monitoring inju rv causes in the
- 1254 usa: a review. Injury Prevention21(4), 278-284 (2 015)
- 1255 37.Berkowitz, S.A., Traore, C.Y., Singer, D.E., A tlas, S.J.: Evaluating area-based socioeconomic s tatus indicators
- 1256 for monitoring disparities within health care sys
   tems: Results from a primary care network. Health
   services
- 1257 research50(2), 398-417 (2015)
- 1258 38.Brechb"uhl, B., M"uller, L., Tschirren, M.: St atistiques des villes suisses 2017. Technical rep ort, Union des villes
- 1259 suisses UVS, Swiss Federal Statistical Oce (FSO) (2017)
- 1260 39.Glaeser, E.: Cities, productivity, and quality of life. Science333(6042), 592-594 (2011)
- 1261 40.Kong, L.: Introduction: culture, economy, poli cy: trends and developments. Geoforum31(4), 385 (2010)
- 41.Harrell Jr, F.E.: Regression modeling strategi 1262 es
  - 42.Tibshirani, R.: Regression shrinkage and selection via the lasso. Journal of the Royal Statisti
- 1263 cal Society: Series
- B(Methodological)58(1), 267-288 (1996)
  43.Schleich, M., Olteanu, D., Ciucanu, R.: Learni
  ng linear regression models over factorized join
- 1265 s. In: Proceedings
   of the 2016 International Conference on Managemen
- t of Data, pp. 3–18 (2016). ACM
  - 44.Anselin, L.: Spatial Econometrics vol. 310330.
- 1267 Chap. Fourteen
  45.Rey, S.J., Anselin, L.: PySAL: A Python Librar
  y of Spatial Analytical Methods. The Review of Re
- 1268 gional Studies
- 1269 37(1), 5-27 (2007)

- 1313 Analysis of statutory health insurance data. Woun d Repair and Regeneration24(2), 434-442 (2016)
- 1314 34.Rawte, V., Anuradha, G.: Fraud detection in he alth insurance using data mining techniques. In: 2015
- 1315 International Conference on Communication, Inform ation & Computing Technology (ICCICT), pp. 1-5 (2 015).
- 1316 IEEE
- 1317 35.Nian, K., Zhang, H., Tayal, A., Coleman, T., L
  i, Y.: Auto insurance fraud detection using unsup
  ervised spectral
- 1318 ranking for anomaly. The Journal of Finance and D ata Science2(1), 58-75 (2016)
- 1319 36.Bell, N., Arrington, A., Adams, S.A.: Census-b
  ased socioeconomic indicators for monitoring inju
  rv causes in the
- 1320 usa: a review. Injury Prevention21(4), 278-284 (2 015)
- 1321 37.Berkowitz, S.A., Traore, C.Y., Singer, D.E., A tlas, S.J.: Evaluating area-based socioeconomic s tatus indicators
- 1322 for monitoring disparities within health care sys
   tems: Results from a primary care network. Health
   services
- 1323 research50(2), 398-417 (2015)
- 1324 38.Brechb"uhl, B., M"uller, L., Tschirren, M.: St atistiques des villes suisses 2017. Technical rep ort, Union des villes
- 1325 suisses UVS, Swiss Federal Statistical Oce (FSO) (2017)
- 1326 39.Glaeser, E.: Cities, productivity, and quality of life. Science333(6042), 592-594 (2011)
- 1328 Donadioetal.Page 21 of30

- 40.Kong, L.: Introduction: culture, economy, poli cy: trends and developments. Geoforum31(4), 385
- 41.Khan, G., Qin, X., Noyce, D.A.: Spatial analys is of weather crash patterns. Journal of Transpor tation
  - Engineering134(5), 191-202 (2008). doi:10.1061/(A
- 1331 SCE)0733-947X(2008)134:5(191)
  - 42.MORAN, P.A.P.: NOTES ON CONTINUOUS STOCHASTIC
- 1332 PHENOMENA. Biometrika37(1-2), 17-23
- 1333 (1950). doi:10.1093/biomet/37.1-2.17.
  - https://academic.oup.com/biomet/article-pdf/37/1-
- 2/17/487420/37-1-2-17.pdf 43.Santosa, F., Symes, W.W.: Linear inversion of
- band-limited reflection seismograms. SIAM Journa
- and Statistical Computing7(4), 1307-1330 (1986).
  doi:10.1137/0907087.https://doi.org/10.1137/0907
- 1336 087
  44.Efron, B., Hastie, T., Johnstone, I., Tibshira
  ni, R.: Least angle regression. Ann. Statist.32
- 1337 (2), 407-499 (2004).

- 1272 Methodology)60(2), 271-293 (1998)
- 47.Hoaglin,

  D.C., Welsch, R.E.: The hat matrix in regression and anova. The American Statistician32(1), 17–22

  1274 (1978)
- 1275 48.Zhang,
  Y., Zhang, Y.: Exploring the relationship betwee
  n ridesharing and public transit use in the unite
- 1276 International journal of environmental research a nd public health15(8), 1763 (2018)
- 49.Horowitz,
   I.: City Politics and Planning. Routledge, ???
   (2017)
- Fan, Y., Perry, S., Kleme's, J.J., Lee, C.T.: A review on air emissions assessment: Transportati on. Journal
- 1279 of cleaner production194,673-684(2018)
- 51.Saboori,
  B., Sapri, M., bin Baba, M.: Economic growth, en
  ergy consumption and co2 emissions in oecd

- 1339 el. Ann. Statist.6(2), 461-464 (1978).
- 1340 doi:10.1214/aos/1176344136
  46.Sauerbrei, W., Buchholz, A., Boulesteix, A.L., Binder, H.: On stability issues in deriving m
- 1341 ultivariable
   regression models. Biometrical Journal57(4), 531-
- 1342 555 (2015). doi:10.1002/bimj.201300222. https://onlinelibrary.wiley.com/doi/pdf/10.1002/b
- imj.201300222
  47.Heinze, G., Wallisch, C., Dunkler, D.: Variabl
- e selection a review and recommendations for th
  1344 e practicing
- statistician. Biometrical Journal60(3), 431—449
- 1345 (2018). doi:10.1002/bimj.201700067. https://onlinelibrary.wiley.com/doi/pdf/10.1002/b
- imj.201700067
  48.Chernick, M.R.: Bootstrap Methods: A Guide for
  Practitioners and Researchers. Wiley Series in Pr
- obability and Statistics. Wiley, ??? (2011). https://books.goog
- 1348 le.it/books?id=UxDKh5Spwp8C
  - 49.Nogueira, S., Sechidis, K., Brown, G.: On the stability of feature selection algorithms. Journ
- 1349 al of Machine
- Learning Research18(174), 1-54 (2018)
  50.Anselin, L.: Spatial Econometrics vol. 310330.
- 1351 Chap. Fourteen
- 1352 51.Hurvich,

  C.M., Simono↔, J.S., Tsai, C.-L.: Smoothing para meter selection in nonparametric regression using an
- improved akaike information criterion. Journal of the Royal Statistical Society: Series B (Statisti
- 1354 Methodology)60(2), 271-293 (1998)
- 1355 52.Hoaglin,
  D.C., Welsch, R.E.: The hat matrix in regression
  and anova. The American Statistician32(1), 17-22
- 1356 (1978)
- 1357 53.Zhang,
   Y., Zhang, Y.: Exploring the relationship betwee
   n ridesharing and public transit use in the unite
   d states.
- 1358 International journal of environmental research a nd public health15(8), 1763 (2018)
- 1359 54.Horowitz,
   I.: City Politics and Planning. Routledge, ???
   (2017)
- Fan, Y., Perry, S., Kleme's, J.J., Lee, C.T.: A review on air emissions assessment: Transportati on. Journal
- 1361 of cleaner production194,673-684(2018)
- 1362 56.Saboori,
  B., Sapri, M., bin Baba, M.: Economic growth, en
  ergy consumption and co2 emissions in oecd
- 1363 (organization for economic co-operation and devel opment)'s transport sector: A fully modified bi-d irectional

- 1286 E.M. collected the data, analysed the data, devel oped the models, analysed the results, supervised the research,
- 1287 wrote the paper. L.D. analysed the data, develope d the models. R.S. analysed the results, wrote th e paper,
- supervised the research. C.R.B. wrote the paper, supervised the research. All authors read and approved the final
- 1289 manuscript.
- 1290 Acknowledgements
- 1291 The authors would like to thank the insurance com pany La Mobiliere for supporting this research an d for providing
- 1292 the dataset used for the analysis.

- 1368 E.M. collected the data, analysed the data, devel oped the models, analysed the results, supervised the research,
- 1369 wrote the paper. L.D. analysed the data, develope d the models. R.S. analysed the results, wrote th e paper,
- supervised the research. All authors read and app 1370 roved the
  - final manuscript.
- 1371 Acknowledgements
- 1372 The authors would like to thank the insurance com pany La Mobiliere for supporting this research an d for providing
- 1373 the anonymized dataset used for the analysis.