

# rivista di statistica ufficiale

REVIEW OF OFFICIAL STATISTICS

**In this issue:**

n. 2  
2021

Putting People First: Beyond COVID-19

*Jean-Paul Fitoussi, Khalid Malik, Jill Rubery  
and Robert Skidelsky*

Optimal sampling design for household finance surveys  
using administrative income data

*Giulio Barcaroli, Giuseppe Ilardi, Andrea Neri, Tiziana Tuoto*

An analysis of the influence of tunnel length and road type  
on road accident variables

*Antonella Pireddu, Silvia Bruzzone*



# rivista di statistica ufficiale

REVIEW OF OFFICIAL STATISTICS

n. 2  
2021

**In this issue:**

- Putting People First: Beyond COVID-19  
*Jean-Paul Fitoussi, Khalid Malik, Jill Rubery  
and Robert Skidelsky* 7
- Optimal sampling design for household finance surveys  
using administrative income data  
*Giulio Barcaroli, Giuseppe Ilardi, Andrea Neri, Tiziana Tuoto* 29
- An analysis of the influence of tunnel length and road type  
on road accident variables  
*Antonella Pireddu, Silvia Bruzzone* 71

**Editor:**

Patrizia Cacioli

**Scientific committee****President:**

Gian Carlo Blangiardo

**Members:**

Corrado Bonifazi	Vittoria Buratta	Ray Chambers	Francesco Maria Chelli
Daniela Cocchi	Giovanni Corrao	Sandro Cruciani	Luca De Benedictis
Gustavo De Santis	Luigi Fabbris	Piero Demetrio Falorsi	Patrizia Farina
Jean-Paul Fitoussi	Maurizio Franzini	Saverio Gazzelloni	Giorgia Giovannetti
Maurizio Lenzerini	Vincenzo Lo Moro	Stefano Menghinello	Roberto Monducci
Gian Paolo Oneto	Roberta Pace	Alessandra Petrucci	Monica Pratesi
Michele Raitano	Maria Giovanna Ranalli	Aldo Rosano	Laura Terzera
Li-Chun Zhang			

**Editorial board****Coordinator:**

Nadia Mignolli

**Members:**

Ciro Baldi	Patrizia Balzano	Federico Benassi	Giancarlo Bruno
Tania Cappadozzi	Anna Maria Cecchini	Annalisa Cicerchia	Patrizia Collesi
Roberto Colotti	Stefano Costa	Valeria De Martino	Roberta De Santis
Alessandro Faramondi	Francesca Ferrante	Maria Teresa Fiocca	Romina Fraboni
Luisa Franconi	Antonella Guarneri	Anita Guelfi	Fabio Lipizzi
Filippo Moauro	Filippo Oropallo	Alessandro Pallara	Laura Peci
Federica Pintaldi	Maria Rosaria Prisco	Francesca Scambia	Mauro Scanu
Isabella Siciliani	Marina Signore	Francesca Tiero	Angelica Tudini
Francesca Vannucchi	Claudio Vicarelli	Anna Villa	

**rivista di statistica ufficiale**

n. 2/2021

Four-monthly Journal: registered at the Court of Rome, Italy (N. 339/2007 of 19th July 2007).

e-ISSN 1972-4829

p-ISSN 1828-1982

© 2022

Istituto nazionale di statistica

Via Cesare Balbo, 16 – Roma



Unless otherwise stated, content on this website is licensed under a Creative Commons License - Attribution - 3.0.

<https://creativecommons.org/licenses/by/3.0/it/>

Data and analysis from the Italian National Institute of Statistics can be copied, distributed, transmitted and freely adapted, even for commercial purposes, provided that the source is acknowledged.

No permission is necessary to hyperlink to pages on this website. Images, logos (including Istat logo), trademarks and other content owned by third parties belong to their respective owners and cannot be reproduced without their consent.

*The views and opinions expressed are those of the authors and do not necessarily reflect the official policy or position of the Italian National Institute of Statistics - Istat.*

*The Scientific Committee, the Editorial Board and the authors would like to thank the anonymous reviewers (at least two for each article, on a voluntary basis and free of charge, with a double-anonymised approach) for their comments and suggestions, which enhanced the quality of this issue of the Rivista di statistica ufficiale.*

## Editorial Preface

The present issue N. 2/2021 of the *Rivista di statistica ufficiale* opens with a challenging article on political economy written by Jean-Paul Fitoussi, Khalid Malik, Jill Rubery and Robert Skidelsky, in the context of the COVID-19 global pandemic.

The paper strongly supports the apparently simple concept that all policies need to be driven by the notion of “*putting people first*”. Consequently, both the beneficiaries and the disadvantaged concerning the state actions have to become an essential part of the social dialogue and policy-making.

The role of the state is undoubtedly decisive and needs to be transformative, in order to promote justice and remove barriers, advancing people’s interests and ensuring that rights and capabilities of both present and future generations are taken into account in defining and implementing policies.

The paper focusses on the notion that markets are a social construct and that their efficacy and social value largely depend on the purposes of the states themselves.

It is necessary that post-COVID-19 economic and social policies are based on a mix of recovery and reform, which cannot longer be separated and have to be addressed to the several aspects of sustainability, with the utmost attention to enhancing social protection systems.

The second article, by Giulio Barcaroli, Giuseppe Ilardi, Andrea Neri and Tiziana Tuoto, represents the result of a synergistic collaboration between *Banca d’Italia* (the Central Bank of Italy) and the Italian National Institute of Statistics - Istat.

It deals with the integration of administrative sources with sample surveys, focussing on the ensuing methodological challenges. In more detail, register data on personal income are used as auxiliary information in the sampling design of the Italian Survey on Household Income and Wealth (SHIW), part of the Eurosystem’s Household Finance and Consumption Survey (HFCS).

The aim is to further improve the information on household income and wealth that are increasingly used for policy-making, so as to best represent their full distribution and to provide a more accurate picture of the economic situation of all households.

There are, indeed, critical issues in obtaining the sufficient number of observations, due to low response rates from both richer and poorer

households. Consequently, exclusively survey-based estimators risk to result biased. In addition, given their most severe impact on estimates, it becomes imperative to develop strategies orientated towards oversampling richer households, which concentrate a large share of total income and wealth.

For this purpose, the authors illustrate a method and the related application for an optimal stratification and sample allocation, creating two unique archives, which proved to play an essential role both in the strategy developed and in the encouraging results obtained.

Antonella Pireddu and Silvia Bruzzone close this issue by signing an article on the several critical aspects related to the road accidents occurred inside tunnels.

Similarly, to the previous work, it is a positive example of exploitation of different sources and synergies, in this case between the Italian National Institute for Insurance against Accidents at Work - INAIL, and the Italian National Institute of Statistics - Istat.

The illustrated in-depth study examines the relationship between the class of vehicle involved, the time of the accident, work-related and non-work-related trips, the circumstances observed and the infrastructure characteristics.

More specifically, the recourse to geoprocessing outcomes in an integrated way with the other available data allowed to insert information also on the length of the tunnels involved in the road accidents detected. In addition, the combination of descriptive statistics and multidimensional methods based on the application of the Principal Component Analysis proved to be very suitable for maximising and enhancing the results obtained, emphasising the association between accident parameters and tunnel types.

Overall, on the one hand this article stands as methodological reference for those who are required to assess road safety, to manage risks and to evaluate the effectiveness of prevention strategies and policies, preparing the ground for future research that can be focussed on a more homogeneous and limited subset of tunnels.

On the other hand, it highlights also some limits and criticalities, by identifying their solution in further integrating different sources and new archives of data and information, in order to be able to take into account a larger number of variables than those used in the current analysis.

Patrizia Cacioli  
*Editor*

Nadia Mignolli  
*Coordinator of the Editorial board*

## Putting People First: Beyond COVID-19

Jean-Paul Fitoussi <sup>1</sup>, Khalid Malik <sup>2</sup>, Jill Rubery <sup>3</sup> and Robert Skidelsky <sup>4</sup>

### Abstract

*“The political problem of mankind is to combine three things: Economic Efficiency, Social Justice and Individual Liberty”.*

John Maynard Keynes

*The global pandemic is bringing all the issues that torment the modern world to boiling-point: inequality, job scarcity, ecological transformation, mass migration, artificial intelligence and the uncontrolled development of public goods and culture<sup>5</sup>. This paper advances the simple idea that all policies must be driven by the notion of ‘putting people first’. And, who benefits and who loses from state actions has to become an essential part of societal dialogue and policy-making.*

*The emergency of COVID-19 has forced the state into a commanding economic position.*

- 1 Jean-Paul Fitoussi is Professor Emeritus at the Institut d’Etudes Politiques de Paris (SciencesPo), and Professor at LUISS Guido Carli University, Rome. Jean-Paul Fitoussi is also a member of the Centre for Capitalism and Society at Columbia University. His last books are: *Measuring What Counts: The Global Movement for Well-Being*, with Joseph E. Stiglitz and Martine Durand (the New Press, 2019), and *Comme on nous parle: L’emprise de la novlangue sur nos sociétés* (Les Liens qui libèrent, 2020).
- 2 Khalid Malik is a former Director of the UNDP Human Development Report Office and lead author of the *Human Development Report 2013. The Rise of the South: Human Progress in a Diverse World*, and the *Human Development Report 2014. Sustaining Human Progress: Reducing vulnerabilities and Building Resilience*. His last book *Why Has China Grown So Fast For So Long* (Oxford University Press, 2012) is now available in Chinese (Renmin University Press). His current research interests include a focus on ‘development as transformation’.
- 3 Jill Rubery is Professor of Comparative Employment Systems and Director of the Work and Equalities Institute at Alliance Manchester Business School. She is an elected fellow of the British Academy and the Academy of Social Sciences. Her research focusses on the inter-disciplinary comparative analysis of employment systems, with particular interests in wage structures, employment regulation, minimum wage systems, working time and welfare systems. She has worked extensively for the European Commission and the International Labour Organisation.
- 4 Robert Skidelsky is Professor Emeritus of Political Economy, Economics Department, Warwick University. He is author of: a three-volume biography of the economist John Maynard Keynes (1983, 1992, 2001); *How Much is Enough?: Money and the Good Life* (with Edward Skidelsky, 2012); *What’s Wrong with Economics?: A Primer for the Perplexed* (2020); etc. He was elevated to the UK House of Lords as Lord Skidelsky of Tilton in 1991 and was elected Fellow of the British Academy in 1994. He is chair of the Centre for Global Studies.

*The views and opinions expressed are those of the authors and do not necessarily reflect the official policy or position of the Italian National Institute of Statistics - Istat.*

- 5 Adapted from Laura Penacchi, Social Europe, 19<sup>th</sup> April 2021.

*This presents an opportunity to recast economic and social policy in a form suitable for our times, specifically to reinstate the state's duty to protect the health, security, and well-being of its citizens and enhance their capabilities and human agency.*

*The paper advances the notion that markets are a social construct and that their efficacy and social value depend in large part on the purposes of the states themselves. The purposes of the state cannot be outsourced to the market.*

*Since the 1980s the language of economics has been that of neoclassical economics, with neoliberalism as its political expression. Neoclassical macroeconomics limited the role of the state to the 'fight against inflation', which was delegated to autonomous central banks. By means of this 'newspeak' the Keynesian system was deprived of virtue and neoliberalism came to stand for rigour and even morality. But a policy regime based on the assumption of automatic tendencies towards full employment and cyclical stability is unable to prescribe for a world normally marked by the absence of both. To do that, we need an updated macroeconomic framework indebted, but not confined, to Keynes.*

*Neoclassical microeconomics validated the market-determined allocation of capital and distribution of income. This ignored the truth that markets are necessarily embedded in a system of political economy, whose aims are broader than those of market efficiency. Exclusive attention to market efficiency omits the requirements of morality and justice that bind society: the equal right of all citizens to health, education, decent jobs, political voice, personal and economic security; an approach to well-being rooted in human development and in Amartya Sen's notion of capability.*

*Economic and social policy post-COVID-19 must, therefore, be a mix of recovery and reform. Unlike in the past these cannot be separated, since recovery must be directed to securing the sustainability of the recovered economy, including social sustainability. So, the state's role has to be transformative.*

*It is now widely recognised that recovery requires a 'stimulus' and that the stimulus will need to be primarily fiscal and not monetary. In principle, the balancing of the economy should be a fiscal responsibility, with monetary policy as a constraint on fiscal excess. Success in meeting this challenge will determine the survival of the European single currency.*

*One important fiscal idea is to make the state 'employer of the last resort' via a decent public job guarantee. This would create a more powerful automatic stabiliser*



*than unemployment insurance as well as keeping a vital link between the worker and work.*

*Transformative policy will need to address itself to the task of ensuring economic and social sustainability. It must ensure that market activity does not result in a depletion of precious natural and human resources and the social systems of trust and cooperation on which the well-being of people depend.*

*Three main lessons or conclusions arise from the COVID-19 experience:*

*1. The State matters an obvious but necessary restatement. Much can be achieved by defining clearly the purposes of the state and by taking on a 'whole of policy-making' approach that eliminates institutional lines separating macro policies from micro level action and structural reform.*

*2. Reform and protect. The imperative has to be to build on and reform, rather than strip down and fragment, the social relations of employment, to maintain social cohesion and stability at a macro level and to promote a fairer and more equal society. Employers have to play their part in 'putting people first' by providing decent jobs, contributing to the costs of social protection and providing opportunities for training and progression.*

*3. Influence technological change. There is much worry that the pace and direction of technological change may result in large swaths of society losing their livelihoods, their sense of meaning, and their freedom from control and surveillance. Putting people first means reclaiming technology for human flourishing.*

**Keywords:** Role of the state, newspeak, dominance of fiscal policy in the policy mix, social protection.

## 1. What's in a word?

The re-emergence of the state has not been as a result of conviction, but of necessity. The established theory of economic policy under 'normal' conditions has not changed. Words define concepts. Using the wrong word deprives us of a concept that may capture well one's thoughts and experiences. Concepts in turn drive policy. The wrong concept can have a profoundly negative impact on the lives of people<sup>6</sup>.

The SARS-CoV-2 has exposed the fragility of the conceptual underpinnings of policy. In the same breath we seem to have gone from the 'absolute necessity' of budgetary discipline to the 'absolute necessity' of deficits. And, there is still a lingering desire to think of the 'COVID-19 crisis' as transitory, to be followed by a return to 'normality', as was the case after the financial crisis of 2008-2009.

Economics is illustrative of the pitfalls of language. The terms we use – Keynesian, neoclassical, new classical, new Keynesian – by flattening time, obscures the chronology of theoretical contributions. We are led to suppose that a 'primitive' Keynesian school was succeeded by the 'superior' new classical or neoclassical school, whereas in fact, today's new classical economics is a reversion to primitive pre-Keynesian economics, decked out with a mathematical apparatus designed to give it authority.

Of all the schools, only the Keynesian and Marxian schools attempt to provide a theoretical explanation of unemployment. Today's new classical or neoclassical orthodoxy denies its existence or makes it the consequence of 'frictions' that prevent the market from functioning freely. Thus it remains in the pre-Keynesian period, even if its form has reached a high technical level. From this chronological perspective, Keynesian theory remains post-neoclassical.

The word Keynesian has taken on a pejorative connotation. Gregory Mankiw, President of the US Council of Economic Advisors from 2003 to 2005, highlighted the tension between the Scientist and the Engineer, arguing that the Engineer – that is Keynesianism – remained dominant in the direction of US economic policy. But despite that, it is increasingly clear that it is the

<sup>6</sup> Jean-Paul Fitoussi: *Comme on nous parle: L'emprise de la novlangue sur nos sociétés*, LLL, 2020.

Scientist who has won the intellectual battle, that is, the new classical school whose founding father, Robert Lucas, famously stated “people will no longer take Keynesian theory seriously in the future”. The New Keynesians have responded to the challenge by using the very language of the new classical school. The result is a watered-down version of neoclassical theory, one that accepts neoclassical ‘science’ but leaves a narrow policy space for Keynesian ‘engineering’.

In Europe, the scientist has won both battles – intellectual and policy. Words like full employment, fiscal stimulus, industrial policy, public investment, have been replaced with terms like competitiveness, structural reform, fiscal compact, and public debt. There is a close correspondence between the precepts of the new classical school, and European institutions, and European policies in particular.

It is therefore with an impoverished language, stripped of its diversity, that we describe the Anglo-American/European universe. But as the American experience teaches us, this language does not solve any relevant problem and that is why perversely it is well suited to Europe in particular where rules freeze the handling of economic policy instruments, and prevent their proper application.

This new language tries to influence us in two ways: first, to convince us that everything has been done to solve the haunting problems we face: unemployment, precariousness, inequalities.

By dint of repeating this, it appears that nothing more can be done about them.

The second direction is more concrete and refers to the efficacy of policy measures. The newly empowered language of neoclassical economics takes on, even promotes, ‘structural reform’ of pensions, unemployment compensation, and labour law in order to maximise efficiency, competitiveness, and save money, indifferent to the social cost of such policies.

Consider the curious term “the end of work”. On the one hand, it reflects the ancestral fear of technical progress robbing workers of their employment and remuneration. But it could just as well herald an economy of abundance, and, as Keynes said the end of the ‘economic’ problem. Which of the two hypotheses prevails depends entirely on policy.

The idea that technology is necessarily detrimental to work is invalidated by past technological revolutions. What is called the Fourth (or digital) Industrial Revolution discloses two possibilities: either the increase in productivity resulting from digitalisation makes it possible to increase remuneration and reduce working time, or the increase in productivity is captured by a tiny minority of people, which means that the end of work happens as well, but only represents a Pyrrhic victory. Who will in fact buy the products if most individuals have no incomes? There is indeed some confusion between technology and distribution: if an increase in productivity happens in a world of moderate inequality it can lead to an increase in well-being. If it happens in a highly unequal world, it may lead to a catastrophe – the bigger the degree of inequality, the deeper its catastrophic consequences.

Globalisation has also entered the ‘newspeak’ dictionary. Globalisation is ‘good’ because it widens the market and allows us to take advantage of new opportunities. At the same time it causes disruptions and losses, against which nation-states exist to protect their citizens. What is curious is that Europe seems to prefer to disarm itself. It does not want to be a federation, let alone a power. Defining itself as a federation of nation-states increases the ambiguity of its identity, which reduces its weight in the concert of nations. Depriving itself of many of the instruments of power – fiscal policy, exchange rate policy, industrial policy – it cannot devise a strategy to deal with globalisation, which is fast becoming its soft underbelly.

After the crisis of sovereign debt in 2008-2009 the ‘migraines’ of orthodox economists were budget deficits and public debt. Now they have to admit that both were not as important as previously thought, and that even more, their enlargement might be good policy. But a kind of conditioned reflex renews them: “how are we going to repay the pileup of COVID-19 debt?” they ask? Yet the two situations are radically different. In 2008-2009 the constraints on spending were self-imposed (at the expense of society). Today, in the age of COVID-19, the elites have little choice, since ‘virtuous frugality’ has become politically unacceptable.

The crisis calls governments to order by reminding them that their primary mission is to protect their populations – to put people first – not to unravel the systems of social protection. They can no longer pursue the race for competitiveness by bidding against economic security. Still, they itch for

action and reputation by pursuing reforms in labour law, health insurance, unemployment insurance, and pension systems under the umbrella of ‘structural reform’.

Policies such as deficit spending that are now being taken up as a matter of necessity are precisely those that were earlier dismissed as irrational. Wasn’t the dominance of austerity policies part of the logic of competitiveness, designed to weaken monopoly power in labour markets? This logic of the market is belied by historical experience. From 1945 to the end of the 1970s, growth and inflation kept the public debt under control, in spite of the oil shocks. In France, President Giscard D’Estaing’s seven-year term of office ended in 1981 with balanced public finances. The UK, which at the end of the WWII had a debt of more than 200% of GDP, found itself at the end of the seventies with a debt of about 40%. Even Italy in this period had a low public debt. In addition, more generous systems of social protection and a strengthening of the bargaining power of workers (the opposite of what is happening today) had favourable consequences for economic growth. The problem public debt only emerged with the restrictive policies and the conservative revolution of Thatcher and Reagan, aimed at defeating inflation and reining in the state.

The stubbornness of governments did the rest: the fight against inflation had to be continued at all times, whatever the rate of unemployment and in spite of high interest rates. To build a good public debt, you only need three ingredients: a sluggish economy, low inflation and abnormally high interest rates. This is why the soaring public debt coincided everywhere in the world not with Keynesian policies but with their abandonment, and why a system designed to secure monetary stability led to the greatest financial instability that the world has known since the 1930s.

By means of the ‘newspeak’ the Keynesian system was deprived of virtue, and neoliberalism came to stand for rigour and morality. Policies seeking to advance well-being were said to be inherently ineffective. Today, as is only right, there is much renewed reflection given the COVID-19 crisis, but whether tomorrow there will a stronger fight for social progress depends on the language used to understand what has happened and is happening. What follows is a policy discussion conducted in an updated language much indebted to Keynes but rooted in a human development framework.

## 2. Rethinking the basics and addressing structural deficits

Today, neoliberalism appears to be in retreat, given the scale of the COVID-19 crisis, with much interest in broadening the scope of state action to improve the lives of people. Even the US Federal Reserve in its conduct of monetary policy is being asked to address long-standing deficits of income and market access for minorities and the poor.

But progress is neither automatic nor inevitable. It is the product of enhanced capabilities of individuals coming together in a functioning society. With the right policies and institutional and societal support much can be achieved. In trying to understand where we are currently and where we need to go, we have to re-examine the essential building blocks of our analysis: markets, the role of government, and the intent of policies.

Markets. Markets are an essential safeguard of liberty and variety. But their efficacy has always depended on the purposes of the state. Karl Polanyi in his classic book<sup>7</sup> traced the dynamics between markets and state in 17<sup>th</sup> and 19<sup>th</sup> century Europe. The arc from belief in free markets to a regulated market system arose as a reaction to rising poverty, unemployment and insecurity. Today, COVID-19 presents a similar challenge in redefining the relationship between states and markets.

Markets of course have never existed in a vacuum of institutions, and institutions have never existed without the support of a state. To take just one example, trade cannot be sustained without contracts and trust. This implies a network of legal arrangements (including constitutional) and the taxes to finance it, but also a network of trust that gives assurance that promises will be honoured. Equally a market system cannot generate sufficient social consent if it is routinely associated with heavy unemployment and rising inequality of wealth and income.

The State. Putting people first should be obvious for a government. After all, improving the lives of people is its main responsibility. This should be reflected in the objectives of all policies, not only the economic one. It implies first and foremost that equal attention should be given to all citizens. The UN's Sustainable Development Goals represent a transformational agenda,

---

<sup>7</sup> Karl Polyani, *The Great Transformation*, 1944.

seeking to promote a more just world. Policy frameworks and markets in turn have to be revamped to deliver on that transformational intent.

Principles are important. The basic idea of human development is about promoting equal life chances for all, based on the Kantian principle that all people are of equal value, as also enshrined in the UN Charter. It is based on the universalism of life claims. And it promotes the notion that humans need to be empowered to live lives they value<sup>8</sup>.

The Developmental Perspective requires us to take a further step. John Rawls's influential work on justice<sup>9</sup> sought to develop principles for a just society and to address the problems associated with redistributive justice. While taking issue with the Rawlsian search for the ideal, Amartya Sen<sup>10</sup> proposes a focus on the actual behaviour of people and the need to remedy the injustices that are 'here and now'<sup>11</sup>.

This suggests a second principle to guide state policy in which social context matters. It posits that equal consideration for all (the 'ideal state') may well demand unequal treatment in favour of the poor and the disadvantaged – to enable the most excluded to exercise their human agency more equally (the 'here and now').

Both economic and social policies influence people's life chances and capabilities. Pursuing the broader goal of equity and justice also reinforces social competences and deepens social cohesion.

Since all policies are a means to an end, their justification should lie in their influence on the lives of people. Considerations of equity have to become embedded in the debate about policy. Who gains and who loses when policies are designed or implemented has to become part of the national conversation.

These broader concerns seem to have been forgotten since the end of the 1970s when the conservative revolution took place. Without being fully aware of it, there was a shift from a frame of reference for economic and social policy that favoured final objectives to one in which only intermediate objectives seemed to matter. In other words, there was a shift from political economy to a "technocratic" economy. Commitment to full employment was

---

8 2014 Human Development Report.

9 John Rawls, *A Theory of Justice*, 1971, and *Justice as Fairness: Political not Metaphysical*, 1985.

10 Amartya Sen, *The Idea of Justice*, 2009.

11 Khalid Malik, *SDGs and Justice*, Background Paper for the HLPF, 2019.

replaced by obsession with balancing the budget; obsession with efficiency replaced the conversation about ends and values; people came to be regarded as parts of machines that had to run without ‘frictions’.

The primary aim of government (and society more generally) therefore is to raise the level of well-being of the population: as listed in the 17 SDGs of the UN, or on specific determinants of well-being such as Health, Education, Decent Jobs, Voice, Social Connections, Environmental conditions, Personal security, Economic security, and Material Conditions (Income and Wealth). COVID-19 serves as an illustration of the importance of these determinants. Had the objectives of policies been addressed to these determinants, a lot of human, social, economic suffering would have been avoided.

Some of these capabilities (Sen, HDRs) may be quite elementary such as being adequately nourished, while others may be more complex such as having the literacy required to participate actively in political life<sup>12</sup>. All are essential for a person who wants to choose the life she/he would like to live.

The challenge of globalisation. Globalisation is the process of integrating markets in goods, capital, labour, and information across national frontiers. A globally integrated market economy is the ultimate goal of neoliberalism. It is thus a direct challenge to any attempt by the state to regulate and limit markets. At the same time it is an important cause of economic insecurity, precariousness, and the rise in inequality.

Neoliberal globalisation<sup>13</sup> is very different from the ‘internationalism’ which emerged after World War II. This contained a clear realisation of the mission of the state, to contain the adverse effects of internationalisation: to protect the population through a robust Welfare State, and domestic full employment policy. In addition, national states created, in the Bretton Woods system, the international institutions to foster economic development and to prevent ‘beggar my neighbour’ policies like currency and tariff wars.

Economic integration brings openness. Openness triggers volatility. Volatility fuels insecurity. Insecurity requires protection. The central problem of globalisation, now and then, is thus how the demand for protection against economic, social and environmental insecurity can be met within an international framework.

12 J.E. Stiglitz, A. Sen and J.-P. Fitoussi, *Mismeasuring Our Lives: Why GDP Doesn't Add Up*, 2010.

13 Jean-Paul Fitoussi, *Globalization and the twin protections*, Working Paper OFCE, 2007.



There is a need to disentangle rhetoric from reality, and recognise from the outset that the phenomenon of globalisation is happening in a world populated by nation-states. The COVID-19 pandemic dramatically brought out the central purpose of the nation state as protection of its population. More than ever the nation-states of the world are alive and well: the hyper power of the United-States, the super power of Europe, Russia, and China.

The rhetoric of globalisation clashes with the reality that power and protection are putting strict limits on the interplay of free markets. For example, the selling of a nuclear plant by a country to another (in a context where such a trade is allowed) depends much more on the interplay of power than on economic considerations. The same can be said about the trade in energy, airplanes, vaccines and the like. Trade between countries often obeys geopolitical considerations rather than just economic ones. There are political externalities to economic trade. Most of the time, trade between countries stems at the boundary between economics and diplomacy. However obvious this assessment is, it is necessary to belabour the point to shake the certainties of the free market believers. Here again, as the COVID-19 crisis highlights, globalisation has reduced the sovereignty of many states, so confident in international trade that they had left to other countries to produce goods strategic to their own survival, medicine among others.

Because of this new consciousness, governments now understand that they should control the markets rather than be controlled by them. Markets are a social construct. ‘We the peoples’ are the ones who shape them with rules and regulations and determine their scope. Their ultimate outcomes are a function of our design. Markets are a means not an end.

This implies that governments also have the mission to promote the production of strategic goods and services at home. And this is happening. In several countries, ‘outmoded’ Institutions and ideas are being rejuvenated, as in France where a Planning Bureau was re-created in September 2020. What other governments are saying is similar in substance:

*“The government will have to protect workers, all workers, but it would be a mistake to protect all economic activities indifferently; some will have to change, even radically, and the choice of which activities to protect and which to accompany in the change is the difficult task that economic policy will have to face in the coming months”.*

Mario Draghi, 17<sup>th</sup> February 2021.

### 3. Towards a framework for economic and social policy

The sluggish and increasingly uneven economic performance of the neoliberal era challenges the orthodox dictionary. COVID-19 has gone further and laid bare the structural deficits of the existing system.

Traditionally policy challenges are divided into short-term and long-term. The short-term challenge is to macro-policy, the long-term to meso and micro-policy. We can stick with this division, provided we remember that the way we handle short-run cyclical disturbances is bound to shape the long-run trajectory of the overall economy. But we need to go further.

It is generally accepted that our economies require a ‘stimulus’ after the profound shock of COVID-19. President Biden has given the lead here with a stimulus package of \$1.9 trillion, amounting to about 10% of US GDP<sup>14</sup>. This is grounded on the recognition that market-based economies lack an automatic recovery mechanism strong enough to bring them back to full employment. The consensus of forecasts is that in the post-COVID period the UK and European economies will, in the absence of stimulus measures, be up to 10% smaller than they were in 2019, with a corresponding rise in unemployment. There is a debate about how much ‘spare capacity’ this actually represents, because of hysteresis and scarring effects. Existing capacity has been destroyed, not just shut down, so new capacity will have to be created.

It is further right to state that recovery policy will have to address itself to issues of supply and not just demand, *i.e.* that the Keynesian remedy of digging up holes and filling them up again is inadequate. That approach whilst useful does not create the desirable long-term improvement in people’s life chances. True enough, any direct boost to demand, by increasing national income, is also an indirect boost to supply. However, there is a danger that any serious lag in the supply response will cause inflation. So for this reason an investment policy must pay attention to questions of quality and not just quantity. This directs attention to the nature of the supply required by economies of the future.

The big current debate in macro-policy is between the respective parts to be played by fiscal and monetary expansion. The orthodox view has long

---

<sup>14</sup> The first part of the stimulus provided under the Trump administration amounted to 3 trillion dollars.

been that any ‘heavy lifting’ an economy requires should be done by monetary policy. This goes back to Friedman’s restatement of the QTM, which posits a direct link between money and nominal income. The effects of quantitative easing (QE) will be divided between prices and output. In one theory of the transmission, a boost to (say) asset prices is expected to spill over into the real economy via wealth and/or confidence effects.

However, experience of 2008-2009 shows that monetary policy is much the weaker of the two recovery instruments. The best that QE did – and this is a considerable plus – was to prevent the banking crisis from sliding over into another Great Depression. But it provided only a marginal stimulus for recovery. The reason lay in low profit expectations by investors and the quantity of defaulted loans in banking balance sheets. So there is now almost a consensus that fiscal expansion is the more effective of the two instruments in present circumstances. It guards against ‘hoarding’ of new money, and it is more effective in directing money to the real economy.

In this context, it is now timely for governments to adopt once again full employment as a key objective of the state. And perhaps go even further by becoming ‘employer of the last resort’, offering a public sector decent job guarantee to any person of working age willing to work but unable to secure a job in the private sector. A buffer stock of decent public sector jobs which waxes and wanes with the business cycle would be a much more powerful economic stabiliser than the complex and increasingly punitive system of unemployment benefits, while its automatic character would guard against the dangers of a ‘political business cycle’. Having a job is critical to individual self-esteem. It reinforces family cohesion and enhances social solidarity.

Against this acceptance of a loosening of fiscal chains, the old language deplores the lack of any apparent limit to government spending (this is why the ECB was forbidden to finance the debt of member governments). If there is no limit to central bank financing of government spending (monetising the debt), the Central Bank, it is claimed, simply becomes an agent of the Treasury. In the UK, since March 2020, the expansion of the QE programme has exactly tracked the increase in the budget deficit. Can the idea of Central Bank independence survive the perception that the main aim of monetary policy is to enable the government to finance its deficit? Modern Monetary Theory argues that the only constraint on government spending is inflation,

but this begs the question why the government should treat inflation as a constraint, in view of its beneficial effects in stimulating activity by, among other things, reducing the real cost of borrowing.

This critique of pure ‘fiscalism’ has led the argument for retaining central bank independence to set interest rates as a curb on excessive government spending. This recognises, though, that the government should be fiscally active rather than passive.

The question of new rules to secure necessary coordination of fiscal and monetary policy remains under discussion, even after old reminiscences like Ricardian equivalence have started to yield to the brute facts of persisting under-utilisation of capacity punctuated by recurrent financial crises.

Cutting across these questions is the realisation that any COVID-19 recovery policy should be directed to securing the sustainability of the economy and not just its cyclical stability. This concern is in the light of the two long-term challenges of automation and climate change.

The first challenge may be put thus. Offsetting the policy of job recovery will be the growth of ‘technological unemployment’, as automation routinises an ever-increasing number of jobs and professions and platform work further fragments many of the remaining jobs. Such an outcome cannot be addressed by counter-cyclical policy alone. It requires interventions at the micro-level.

One can think of three kinds of intervention, separately or together. The first would be to slow down automation to enable full employment at present or reduced hours of work to continue into the medium future. A second would be to impose statutory limits on hours worked per week (as happened in France’s ‘*Reduction du temps de travail*’, 1980). Here the benefit would be also to facilitate moves towards greater gender equality through more equal sharing of both paid and unpaid work. The third intervention would be to boost training and ‘upskilling’ programmes to enable re-employment of displaced workers (though that may not necessarily minimise job losses).

We should not reject the first intervention entirely: at the very least it draws attention to differences between ‘good’ and ‘bad’ technology and between means and ends. However, the main medium term challenge will be to ensure that fruits of productivity gains do not accrue wholly to the owners of capital. So, redistribution emerges as a key issue in dealing with

technological unemployment. This is only likely to be effective if the current wave of technology does not undermine the system of secure employment with guaranteed hours and income. One can then think of different kinds of redistributionary policy – wealth tax, profit tax, progressive income tax.

In this context schemes of universal basic income (UBI) come into play. They can be thought of as a national dividend, which grows with the productivity growth of the economy. A problem with UBI is its start-up cost, since it is by definition untargeted, and the danger that it may replace other necessary social support programmes like health insurance or unemployment benefits, as well as facilitating and promoting precarious forms of employment, particularly if paid at a relatively low level. The benefits and costs of universal basic income need to be carefully assessed and trialed.

Climate change raises a somewhat different problem of sustainability. The challenge it poses is to ensure a zero-carbon emissions economy. This is separate from both cyclical stability and productivity growth. The need is to ensure a reorientation of the economy to the kind of output which safeguards not future jobs or incomes but future resources. How far can the required rebalancing of economic life be left to market forces (consumer demand) or ‘nudged’ by tax incentives? To what extent does it require direct government investment in the ‘green economy’ via institutions like Green Investment Banks – the last of which also raises the issue of public debt sustainability? Such questions should be at the forefront of any re-thinking of the role of the state.

In short, the state has a deep responsibility to advance people’s interests and to ensure that rights and capabilities of both the present and future generations are taken into account in defining and implementing policy. For that to occur, the state’s purpose has to be seen as a transformational one, promoting justice and removing barriers that hold back human agency.

## 4. The State's Investment Function

Keynesian stabilisation policy was concerned with the quantity of investment; microeconomic interventions are concerned with the allocation of investment between different uses. Why should not this be left to the financial system in line with investor choice? There is a growing consumer led demand for 'ethical' investment.

However, offsetting this movement from below is the degenerated and predatory character of the financial system. Much of it is socially useless and worse – a Ricardian rent extractor rather than a sustainable wealth enabler. The question then is by what mix of regulation, direction, and prohibition can finance be reattached to the general welfare? Does the political will exist in the EU or (now) the UK for anything beyond cosmetic tinkering?

This takes us back to fiscal policy. The case for the primacy of fiscal policy is not just that it is a more powerful macroeconomic stabiliser than monetary policy, but that government is the only entity apart from the financial system capable of allocating capital. If we are not willing to allow investment in technology and infrastructure to be shaped by a purely financial logic, the role of what Mariana Mazzucato calls 'mission-oriented' public investment becomes inescapable<sup>15</sup>. The last stand of neo-liberal orthodoxy is to assert the superior efficiency of private investment, especially as compared to public investment. The logic is impeccable, but it has lost contact with the reality and the necessities of our time.

---

15 Mariana Mazzucato, *Mission Economy: A Moonshot Guide to Changing Capitalism*, 2021.

## 5. Rethinking Social Protection and Innovation

The COVID-19 pandemic has coincided with a period of intense debate about the future of work and employment that has been fuelled by a new wave of technological innovation, presaging the demise of employment of the quantity and form to which the developed world has been accustomed. The pandemic has, however, revealed the very high risks presented by such radical changes and has added to the urgency to strengthen systems of social protection.

To deliver on the larger purpose of improved lives and overall well being, governments have to take on a transformative role, that puts people first, women as well as men, and addresses the structural impediments to minorities, the poor and the disadvantaged, from exercising their agency fully. A ‘whole of policy-making’ approach necessarily reduces the institutional lines separating macro policies from micro-level actions and structural change.

Going forward, a holistic policy approach should be governed by four guiding ideas.

### A. The State matters

It is imperative to reverse the creeping privatisation of public services in the name of ‘efficiency’ and ‘saving money’. Examples abound of ‘outsourcing’ public services, mainly to the detriment of their users, ranging from the well documented failures of state run prison systems in the US to the recent costly and ultimately ineffective UK COVID-19 tracing programmes which Britain’s Public Health Executive sub-contracted to private entities. The focus on efficiency and the mythical belief in the efficacy of the private sector has eroded common sense approaches to the provision of social services leading to a severe underfunding of state institutions.

### B. Improving the Social Relations of Production

As the pandemic hit all forms of employment, the need to maintain the linkages between employers and their workforces and to maintain wage income became apparent to most governments in developed countries.

The pandemic could be considered an experiment in the kinds of dramatic changes to employment predicted in the future of work debates but, which when contemplated as a mass change, forced governments to recognise the unsustainability of a dislocated and fragmented employment system in which citizens have no security of income or employment. Moreover, employers have showed more loyalty to their workforces than might be expected under predictions of the inexorable rise of transactional, platform-based work. Even in the United States where there was no job retention or short time work scheme, most of those made unemployed returned to work with their previous employer and most furloughed workers in the UK have been taken back into employment with their current employer.

This suggests the need to maintain and strengthen secure employment relationships. To address the growing fragmentation and insecurity of employment a dual approach is needed that universalises social protections (for example basic incomes for children and citizens pensions), while also extending the umbrella of the standard employment relationship to incorporate more flexible forms of working and more flexible careers. This also requires removing incentives to employers to create precarious jobs – by raising minimum legal standards (for example living wages and guaranteed hours) and taxing employers at the same rate for all forms of work, both direct employment and fees to freelancers or to outsourcing companies.

### C. State obligations and sharing of responsibilities

While the essential obligation to look after its citizens lies with the state, there is need for clearer thinking about how to ensure that employers share the burden.

If citizens are to retain some stability and security in their lives there is a need for employers to share the responsibilities and costs of maintaining citizens and the workforce in periods when working is not possible due to sickness, maternity, unemployment or old age.

Traditionally the employer contributes through both social contributions and by providing open-ended long duration contracts that minimise moves between employment and unemployment and maintain some stability of income both during periods when workers cannot work due *e.g.* to sickness



and in periods when there are minor demand fluctuations. If work is allowed to become more precarious, employers evade more of these costs while at the same time workers' dependency on support from the state increases. The sharing of these costs needs to be maintained by a mixture of methods.

The state needs to take action to block exit routes used by employers to evade responsibilities and to minimise opportunities for voice at work. This means reducing incentives to use outsourcing, casual work, platforms, zero hour contracts and bogus self-employed contracts. The establishment of joint responsibilities between, for example, the client and the outsourcing company, the client and the agency, the client and the platform, all to ensure fair treatment of the workers, including the right contractual status and opportunities for voice and due process within the employment relationship. These initiatives also need to be developed in conjunction with trade union and employee representatives.

There is also the need for employers to help reduce the impact of both COVID-19 and technological change on employment. In the initial recovery period, job sharing could help reduce mental health problems as well as providing a more orderly move towards new careers and restructuring. This is only likely to be achieved through cooperation, with the state providing incentives and support for work sharing and employers facilitating changes to work practices. This could have longer-term benefits as the increasing productivity stemming from technological advance could restart the move towards shorter working hours, a movement that has stalled over recent decades. A move towards a lower standard of full-time hours, more equivalent to 30 hours or a four day maximum working week could also enable a move towards more equal sharing of both paid and unpaid work between men and women. The mass homeworking experiment under COVID-19 has provided the opportunity to normalise flexible working, facilitating a more equal organisation of both work and family life, though there will be a need to guard against the danger of increased employer surveillance of home-based work.

The disruption to employment careers and life chances threatened by both COVID-19 and the application of new technologies also highlights the need for employers and the state to engage cooperatively on major programmes of training and retraining. Action is needed not only to upskill within organisations but also to help those displaced to develop new types

of skills and to change careers. This effort needs active involvement of the state, employers and trade unions and employee representatives. The state is needed to fund retraining programmes but employers are needed to provide opportunities to fully develop these skills within a work environment. Trade unions and employee representatives need to be involved to champion the capacities of workers to be retrained and redeployed, to avoid the costly waste of human potential caused by redundancy and long-term unemployment.

#### D. Influencing technological change

A UK study<sup>16</sup> refers to a potential forty-five percent reduction in middle class professional jobs and a sharp rise in low-end jobs (the precariat) over the next decades. But innovation and technological change processes are embedded in society. They are influenced by incentives and should not be seen as an exogenous variable to which people ‘must adjust’. If society adopts a full employment target, then the system is forced to find policy solutions that make that objective possible. Thus instead of taking action to remedy the potential effects of technological change – as outlined above – a better way may be to build in the objective of putting people first into the design and implementation of technology.

The link between technological choices and employment needs further review. Can we influence investment decisions to take into account broader sustainability objectives – both in relation to the planet and sustainable employment/inclusive growth? Some have argued that the most advanced technology available may not be the best way to meet sustainability objectives (even if in theory AI could be repurposed). This then leads us back to the old debates about intermediate technologies, especially in small-scale local production in order to reduce carbon emissions. Further, can more ecologically focussed technologies be developed? This may require combining public commitments to investment with strong market incentives that reward products of such technologies. Putting people first means reclaiming the focus of employment change to provide opportunities for human flourishing.

---

<sup>16</sup> Frey and Osborne, *The Future of Employment*, 2013.

## 6. Conclusion

We cannot accept the claim of the neoliberal paradigm to be morally neutral, leaving ethics a matter of personal preference. Economists' values are embedded in the models they use, the research topics they select, their choice of relevant variables, and such like. A crisis of values cannot be confronted by the so-called moral neutrality of the market system.

The task facing us is to find the moral and practical equivalents of Roosevelt's New Deal, the Keynesian full employment commitment, the Bretton Woods system, and the postwar welfare state. But, perhaps go beyond that and address long standing deficits in equity, gender and social position.

These are the conditions for any successful and inclusive market system. Their achievement will require a shift of power from finance to production, and of morals from means to ends. And, to go beyond servicing the top tier of society to meeting the needs of people across the board.

The 'bottom line' is an economic and moral system fit for purpose and that which puts people first.



# Optimal sampling design for household finance surveys using administrative income data

Giulio Barcaroli <sup>1</sup>, Giuseppe Ilardi <sup>2</sup>, Andrea Neri <sup>2</sup>, Tiziana Tuoto <sup>3</sup>

## Abstract

*Household finance surveys, which collect detailed information on household income and wealth, are increasingly used for policy-making. They should provide an accurate picture of the economic situation of all households. Unfortunately, the upper parts of the wealth distribution are often missing in household surveys. Since rich households concentrate a large share of total income and wealth, survey-based estimators may be biased. The ideal situation would be to have access to auxiliary information on household finances at the design stage. This is rarely the case. In this paper we present an application that uses tax records in the design of a major survey on household finances. We discuss the methodological challenges of using administrative information for designing the sample. We propose a method for an optimal stratification and sample allocation.*

**Keywords:** Household finance surveys, Household Finance and Consumption Survey - HFCS, register data, tax records, income, sampling design, optimal stratification, calibration.

- 
- 1 Independent expert, formerly at the Italian National Institute of Statistics - Istat ([gbarcaroli@gmail.com](mailto:gbarcaroli@gmail.com)).
  - 2 Central Bank of Italy/Banca d'Italia ([Giuseppe.Ilardi@bancaditalia.it](mailto:Giuseppe.Ilardi@bancaditalia.it); [Andrea.Neri@bancaditalia.it](mailto:Andrea.Neri@bancaditalia.it)).
  - 3 Italian National Institute of Statistics - Istat ([tuoto@istat.it](mailto:tuoto@istat.it)).

*The views and opinions expressed are those of the authors and do not necessarily reflect the official policy or position of the Italian National Institute of Statistics - Istat.*

*The authors would like to thank the anonymous reviewers for their comments and suggestions, which enhanced the quality of this article.*

## 1. Introduction

The measurement of households' economic conditions is high on the political and economic research agenda. In recent years, this topic is becoming increasingly important also for National Central Banks, as it has been recognised to interact with their functions (Eurosystem Household Finance and Consumption Network, 2009).

One of their main targets is to guarantee price stability through monetary policy. To this purpose, they need to have a good knowledge of how households make their spending decisions and how they respond to changes in their finances. Central Banks also have to supervise the risks for financial stability arising from the household sector. For this reason, they need to monitor the household's ability to face their levels of indebtedness if some shock occurs (such as the loss of a job of some member of the household) (Michelangeli and Rampazzi, 2016). Moreover, Central Banks are also increasingly interested in understanding the effects of their policies on the household's economic conditions and in particular on income and wealth inequality (Casiraghi *et al.*, 2018; Colciago *et al.*, 2019; Dobbs *et al.*, 2013; Dossche *et al.*, 2021).

Sample surveys are the main tool used to collect granular information on these aspects. In the Euro area, the European Central Bank has established a network of survey specialists, statisticians and economists to collect harmonised microdata on household income and wealth through the Household Finance and Consumption Survey (HFCS). Because of the range of purposes for which these data are used, it is particularly important that the survey adequately represent the full distribution of income and wealth. In practice, the greatest difficulties are in obtaining a sufficient number of observations in the two extremes of the distributions. Households with very poor finances may see little relevance in participating in a survey about finances. Moreover, they could live in areas that could be dangerous for the interviewers. Under-representation of these households is likely to have little impact on estimates of mean, but it would affect many other statistics such as those related to the income distribution or poverty. At the other end of the spectrum, research has shown that very affluent households are likely to be under-represented: see for example, Eckerstorfer *et al.*, 2016; Neri and Ranalli, 2011; D'Alessio and Neri, 2015; Kennickell, 2019; Vermeulen, 2018; Chakraborty *et al.*, 2019. Indeed, wealthy respondents are generally a hard-to-reach population

since they may live in multiple locations, which, also, may have security measures that make it difficult for the interviewer to contact the household to negotiate the interview. Moreover, rich persons may be difficult to persuade to participate since they are generally busy or less willing to declare their finances. Although such households are small in number, they own a large share of total income or wealth. Thus, the under-representation of these households would have negative effects on many estimates.

The availability of auxiliary information at the design stage (such as administrative records relating to household finances) would prove extremely effective in addressing these issues. Such information would enable survey agencies to identify correctly this rare population, also making it possible to oversampling it to compensate for the difficulties in enrolling it in the survey. Unfortunately, such information is rarely available, mainly because of confidentiality issues that prevent the exchange of personal data among the owner and other institutions. Moreover, even if this information is available, generally it is not consistent with the definitions and the concepts used in the survey.

This study discusses the use of register data on personal income in the sampling design of the Italian HFCS survey. It draws on a collaboration between *Banca d'Italia* (the Central Bank of Italy) which runs the survey, and the Italian National Statistical Institute – Istat which has access to the administrative records. Thanks to this collaboration, we have been able to create two unique archives that are essential for our strategy.

The HFCS survey is a two-stage sample with municipalities selected as primary sampling units (PSUs) and households selected as second-stage units (SSUs). In this paper, we discuss the first time that the information from the personal income register is used to optimise the sample design, focussing on the second stage, while treating the first-stage sample as fixed. A more general and even complex optimal sampling design, which also considers first stage units, is possible and desirable. However, the impact of such a design would have on the organisational procedures that support the survey would certainly be heavy. Therefore, to introduce and manage innovations gradually, it was a survey requirement to deal with the second stage treating the first stage as fixed.

The paper is organised as follows. The following Section will provide a brief overview of the different use of administrative records in the main household finance surveys and the main contributions of our article. Sections 3 and 4 will introduce the survey and register data we use for our application, while Sections 5 and 6 describe the methods used in our sample design. The results are presented in Section 7. The article concludes with a summary and discussion of the main results in Section 8.



## 2. The use of register data in household finance surveys

Administrative records are increasingly used for statistical purposes. Some countries already used them in the design of their household finance surveys.

The US survey of Consumer Finances employs a dual-frame design, including an area-probability (AP) and a list component. The list sample is used to oversample households that are likely to be relatively wealthy. The basis of the sample is a set of specially edited individual income tax returns developed by the Statistics of Income Division (SOI) of the Internal Revenue Service (Kennickell, 2008). The list sample is stratified using a “wealth index” computed using income data to predict a rank ordering of people by wealth. After defining the stratifying variable in terms of the whole population, the list is reduced for the actual selection to include only cases that filed returns from a municipality included in the PSUs underlying the AP sample. Within each stratum, cases are oversampled by a progressively larger proportion in richer strata (Kennickell, 2017).

In Canada, the design of the Survey of Financial Security foresees that each province is stratified into rural and urban areas and different design is used in each. In rural areas, a multi-stage sample is selected using the Labour Force Survey area frame. In urban areas, information from the administrative records at the family level, such as age and income, is used to stratify the Address Register into groups of dwellings having similar well-being.

In the 2017 wave of the HFCS, seventeen out of twenty-two countries used different strategies to oversample richer households (Household Finance and Consumption Network, 2020). Italy was one of the five countries in the HFCS which had no access to auxiliary information that could be used in the sampling design. The oversampling strategies varied significantly between countries, and are heavily dependent on the available data.

The Spanish Survey of Household Finances (EFF) has used, at least for some waves, individual wealth tax files. The sampling is achieved thanks to the collaboration of the INE (Spain’s statistical institute) and the Tax Authorities (TA), through a complex coordination mechanism (for confidentiality reasons). The population frame contains information on fiscal wealth and income for each household. The choice of defining the wealth strata is based on the households’ percentile distribution of the wealth tax for

Spain. Cases in richer strata are over-sampled progressively at higher rates (Bover *et al.*, 2014).

The French Wealth survey uses tax registers on personal wealth data to identify four strata: wealthy city dwellers, equity-based wealth, real estate-based wealth, lower wealth. Richer strata are sampled at higher rates.

Tax registers on personal income are used in Estonia, Finland, Latvia, and Luxembourg, while in Cyprus the sampling is based on the Customer register of the electricity authority.

The main limitation to the use of administrative records is the legal restrictions to protect the privacy of households. Depending on the country, the limitations may relate to the use of the data (for instance, restricting the use to detect tax-evasion purposes) or the transfer of the microdata to any institution outside the producing agency.

Other countries adopt different sampling strategies to compensate for the unavailability of register data at the individual level. Greece, Ireland, Hungary, Poland, and Slovenia use the information at area level (such as average income and real estate) as proxies of households' economic conditions).

Despite the use of register data is not a novelty, to the best of our knowledge, there are not many studies in the literature discussing the benefits and the challenges in the use of register data in the design of a household finance survey. Indeed, administrative records are not built for statistical use and therefore they generally adopt different concepts and definitions from the ones used in the survey. They may also suffer from quality issues such as under-coverage, lack of timeliness, and errors. These issues should be taken into account when using them for sampling purposes. Still, in the literature or the methodological notes of the surveys, many choices are not documented. For example, it is not always clear how the strata boundaries are chosen, how the allocation is defined, or how the above-mentioned differences are taken into account.

The few studies available are mainly focussed on the benefits of using register data. For the US survey on consumer finances, Kennickell (2008) shows that the availability of a list of individuals based on income tax returns produces far more precise estimates of wealth than would be possible with a less-structured sample of the same size, and it provides a framework for

correcting for non-response, which is higher among the wealthy. Similar results are found by Bover (2010) as far as the Spanish survey on household finances is concerned. Other research evaluates the effectiveness of the different strategies in obtaining samples that represent adequately the whole distributions of income and wealth (see for instance Household Finance and Consumption Network, 2016).

We contribute to the existing literature in two ways. The first one is that we present a discussion on the challenges and the (expected) benefits of using personal income tax data, drawing on the data of a real survey. In particular, we present a way to address the issue of biased variance estimates based on administrative records. The second contribution of our paper is to present an optimal stratification and sample allocation strategy to be used for multivariate populations. This solution enables us to jointly identify the optimal stratification based on the tax data and the optimal sample size in each stratum. The method presented in the paper has been applied in the 2020 Italian HFCS. Hopefully, our application may contribute to give insights for other data producers.

### 3. The Italian Survey on Household Income and Wealth

*Banca d'Italia* conducts the Survey on Household Income and Wealth (SHIW) since the 1960s. Starting from 2010, the survey is part of the Eurosystem's Household Finance and Consumption Survey (HFCS), coordinated by the European Central Bank.

The target population of the survey is all individuals that are officially resident in Italy. People living in institutions (convents, hospitals, prisons, *etc.*) or those who are in the country illegally are out of the scope of the survey. The survey is used to collect granular information on many aspects ranging from the socio-demographic characteristics of the household and of its members, to the different sources of income, to the household's assets and liabilities to the consumption and saving behaviours. A household is defined as a person living alone or a group of people who live together in the same private dwelling and share expenditures, including the joint provision of the essentials of living. Persons usually resident, but temporarily absent from the dwelling for less than six months (for reasons of holiday travel, work, education, or similar) are included as household members. On the contrary, possible other persons with usual residence in the dwelling but not sharing expenditures (*e.g.* lodgers, tenants, *etc.*) are treated as separate households.

The sample consists of about 8,000 households. The sample size is chosen to produce estimates at the national level. Since 1989 about half of the sample has included households interviewed in previous surveys (panel households). Data collection is entrusted to a specialised company using professional interviewers and CAPI methodology.

The sample is drawn in two stages, with municipalities and households as, respectively, the primary and secondary sampling units. In the first stage, a stratified sample of about 400 municipalities is selected. The variables used for stratification are the region and population size. In the second stage, a simple random sample of households to be interviewed is then selected from the population registers. Participation in the survey is not mandatory. In case a household refuses to participate in the survey, it is replaced by another one living in the same municipality, randomly selected from population registers.

At present, no auxiliary information relating to the household's finances is available at the design stage. This implies that in the final sample only a few rich households are selected. For instance, just by chance, only 80 households belonging to the top 1 percent will be selected. Moreover, once such a household refuses to participate, the available information does not allow replacing them with another with similar finances. Starting from the 2014 wave, *Banca d'Italia* has progressively taken all the legal steps necessary to have access to the fiscal ids of the persons in the sample to make data linkage with register data possible.

## 4. Register data

In Italy, several public administrations (including the Tax authority) are committed by law to provide their administrative data to the Italian National Statistical Institute - Istat to reduce the cost of data collection and the burden on the citizens. The two registers (held by Istat) exploited in this work are the Italian Population Register (PR) and the Italian Tax Register (TR).

The PR contains individual records for citizens enrolled in the Italian municipality registers, grouped in their administrative declared households. These registers are regularly updated by municipalities based on the declarations they receive from citizens. Whenever there is a change in the household composition, such as people getting married or moving to another city, individuals are supposed to communicate this change to the offices in charge of the population register. In most instances some incentives bring people to keep their official records updated: for example, some taxes are lower for houses that are officially primary residences, so in case of purchase of the main residence people immediately update the official records. The PR is used as a sampling frame of all the household surveys in Italy. It is also used to draw the sample of the Italian HFCS for a long time. In this study, we use the version available at the end of 2018.

The second register we use is the Italian Tax Register held by the tax authority. The latest available version of this register has a 2-year time lag, so, the reference time of the TR is 2016 when writing this paper. The TR contains all the records corresponding to the yearly taxable income of people afferent to the Italian Tax System. It is worthwhile noting that in Italy, people with an income below certain thresholds do not have to provide a tax declaration. Yet, the TR is based on multiple sources which enable to recover the information also for those who are below these thresholds. The main limit of the TR is that it does not include the income for financial assets (interest and dividends) that generally are taxed with a different system and that are not reported in fiscal declarations (according to national accounts, interests and dividends account for about 15 percent of household disposable income). This data gap may limit the utility of the personal tax data to target wealthy households, which usually concentrate a large share of financial wealth.

The income variables used in this study are “Total income”, “Dependent employment income”, “Self-employment income”, “Pension income” and “Rent”. This information is available at the individual level.

In Italy, the tax agency provides individuals from birth with a unique code, foreigners are provided with the code when they enter the country and ask for permission to stay. The two registers have been linked using these identifiers.

The final data frame contains both demographic information (including household composition) and fiscal incomes at the individual level. The new archive has been created only for the persons living in the municipalities selected as primary sampling units in the survey (around 27.5 million individuals). Individual incomes have then been aggregated at the household level using the official PR definition of household.

Households with members with an income higher than a given threshold (1 million euros) have been included in a separate self-representing stratum. It accounts for 0.01 percent of the total population and 0.6 percent of total income. Since the households in this stratum represent a very hard-to-reach population which may require different *ad hoc* strategies, we exclude them from the present analysis. The final sampling list consists of about 12 million households.

Register data are not built for statistical use and therefore they adopt concepts and definitions that may be different from those used in the survey. The first one relates to the definition of household composition. SHIW and surveys in general use “economic household” concept, *i.e.* those actually living together and sharing the essentials of living (Jäntti *et al.* 2013).

Population registers collect information on all the individuals that are officially resident in the same household, while the target of the survey is the “*de facto*” household composition in the reference year (irrespective of the official residency). The two concepts may differ because of changes that may occur between the selection of data from the registry (September of the reference year) and the time of the interview (from January and June of the year following the year of reference). Moreover, in some instances, people may not have an incentive to update their official status, such as immigrants coming back to their native countries for good. Finally, the official composition of the household may be affected by the taxation system. For example, a household

could be fictitiously divided into two groups for saving taxes linked to the different taxation of the main residence compared to secondary dwellings.

The second difference between register and survey data relates to the definitions of the income sources. In the survey, incomes are collected net of taxes and social contributions, while in the TR each income source is recorded gross and only the total amount of taxes paid by each person is available. Moreover, in the case of self-employed taxable incomes are affected by fiscal rules (such as the possibility of deducting operating losses or investments made in previous years) that do not apply in the survey. Another important incoherence is due to the difference in the methodology for assessing the incomes from non-rented dwellings, that is the amount of income a property owner would get by renting her/his own house: in SHIW is adopted the self-assessment method which consists in asking directly the respondents to provide their best estimate, while in the TR the cadastral income (*rendite catastali*) is used for evaluating the stream of these incomes. The cadastral income is a figurative income that can be obtained by multiplying the surface of the property by a specific coefficient, calculated by the Italian Tax Agency according to the municipality, the census zone, the type of dwellings, and its quality. Given that the coefficients are not regularly updated, these incomes significantly underestimate the true value of market rents.

Besides the two differences above mentioned, it worth noting that tax data have quality issues due for instance to tax evasion (Neri and Zizza, 2010; Fiorio and D'Amuri, 2006) and depending on the method used to estimate under-reporting, the magnitude of the problem varies between 7 and 14 percent (Albarea *et al.*, 2018). Moreover, tax data are available with a two-year time lag and therefore may no longer reflect the real situation of the household (especially in the case of self-employed).

One of the main consequences of the above-mentioned issues is that using administrative records for variance estimation in the sample design stage is likely to produce biased results which, in turn, may lead to a sub-optimal selection of the sample.



## 5. Optimal stratification and sample allocation methodology

Stratification is one of the most widely used techniques in sample survey design, serving the twofold purpose of providing samples that are representative of major subgroups of the population and of improving the precision of estimators.

In SHIW/HFCS, the particular aim of the stratification should be to increase precision in the top of the wealth distribution. So far, this has not been implemented in Italy.

The design of stratification involves a sequence of decisions relating the choice of the stratification variables, the choice of the number of strata to be formed, the mode in which strata boundaries are determined, the choice of sample size to be taken from each stratum (allocation of the sample) and the choice of sampling design within strata.

Studies have provided procedures for the determination of the strata boundaries under a given sample allocation, which are mainly applicable to univariate cases (see for instance Kareem and Adejumo, 2015; Horgan, 2006). On the other hand, there are studies proposing methods to solve the problem of optimum allocation for multivariate populations when the strata are already decided (see for instance Khan, 2008). To the best of our knowledge, in the literature, there are no studies proposing methods to deal simultaneously with the issue of strata boundaries definition and sample allocation for multivariate populations.

In this paper, we propose the use of a genetic algorithm (Schmitt, 2001) that can explore the universe of all the possible stratifications looking for the one that minimises the total cost of the sample required to satisfy the precision constraints. This algorithm is implemented in the *R* package *SamplingStrata* (Barcaroli *et al.*, 2020). This package, of current use in the Italian National Statistical Institute for various sampling surveys, has been used in the New Zealand Statistical Institute, tested at Statistics Denmark, and considered for evaluation at Statistics Canada. Eurostat used *SamplingStrata* for designing its 2018 *LUCAS* survey (Ballin *et al.*, 2018). In addition, the World Bank adopted *Sampling Strata* and embedded it in its *Survey Solutions Sampling Tools* integrated application.

Unlike other similar packages (as the package *stratification* Baillargeon and Rivest, 2012), *SamplingStrata* is applicable to the multivariate (more than a target variable) and multidomain (more than a domain of estimation) case, that is exactly the Italian HFCS case. The methodology is fully described in Ballin and Barcaroli, 2013; Barcaroli, 2014; Ballin and Barcaroli, 2016.

In the following, we recall its fundamentals before illustrating the application to the SHIW sampling design. It is worth recalling that the optimal stratification proposed in this paper is only related to the second stage units (the households) of the overall sampling design, since the first stage units (the municipalities) are treated as fixed due to survey requirements related to organisational and fieldwork aspects.

As the aim of the optimisation performed through the genetic algorithm is to find a stratification that minimises the variance inside the strata with respect to all the survey target variables, an important step of the method is to estimate consistently the population variance in all the strata. As already mentioned, register data use different concepts and measures compared to survey data. Moreover, they are likely to suffer from quality issues such as tax evasion and tax elusion and delays. As a consequence, they should not be used as such for the allocation of the sample. In our study, we consider the variables from tax records as proxies of the variables we want to measure. We then estimate measures of goodness-of-fit of these proxies. Finally, we use such measures to inflate our population estimates of the variance in the strata (the higher the goodness-of-fit the lower the inflating factor).

## 5.1 Optimal stratification with the *R* package *SamplingStrata*

In a stratified sampling design with one or more stages, a sample is selected from a frame containing the units of the population of interest, stratified according to the values of one or more auxiliary variables ( $X$ ) available for all units in the sampling frame. For a given stratification, the overall size of the sample and the allocation in the different strata can be determined on the basis of constraints placed on the expected accuracy of the various estimates regarding the survey target variables ( $Y$ ). If the target survey variables are more than one the optimisation problem is said to be multivariate; otherwise it is univariate. For a given stratification, in the univariate case the optimisation of

the allocation is in general based on the Neyman allocation (Cochran, 1977). In the multivariate case it is possible to make use of the Bethel algorithm (Bethel, 1989). The criteria according to which stratification is defined are crucial for the efficiency of the sample. With the same precision constraints, the overall size of the sample required to satisfy them may be significantly affected by the particular stratification chosen for the population of interest. Given  $G$  survey target variables  $Y$ , their sampling variance is:

$$\text{Var}(\hat{Y}_g) = \sum_{h=1}^H N_h^2 \left(1 - \frac{n_h}{N_h}\right) \frac{S_{h,g}^2}{n_h} \quad g = 1, \dots, G$$

where:

$$\hat{Y}_g = \sum_{h=1}^H \sum_{i=1}^{n_h} y_{g,i} w_i \quad \text{: target estimate}$$

$H$  : number of strata

$N_h$  : population in stratum  $h$

$n_h$  : sampling units in stratum  $h$

$S_{h,g}^2 = \frac{1}{n_h - 1} \sum_{i=1}^{n_h} (y_{g,i} - \bar{y}_{g,h})^2$  : sampling estimate of the variance of the  $g$ -th

variable in stratum  $h$

It should be noted that  $\text{Var}(\hat{Y}_g)$  is the design-variance of the estimator of the population total for the  $g$ -th target variable when the sample design is stratified simple random sampling. A more general variance formula would be needed if the optimisation were to find an allocation of PSUs to their strata and families within PSUs to their strata. The general formula would include both a stratified component due to first-stage sampling and another stratified component due to second-stage sampling; component variance formulas are presented, for instance, in Cochran (1977, pp. 308-310), Hansen, Hurwitz, and Madow (1953, ch. 6 and 7), or Valliant, Dever, and Kreuter (2018, ch. 9).

If we introduce the following cost function:

$$C(n_1, \dots, n_H) = C_0 + \sum_{h=1}^H C_h n_h$$

where  $C_0$  indicates a fixed cost (not dependent on the sample size) and  $C_h$

represents the average cost of collecting and processing data for a sampling unit in stratum  $h$ , then the optimisation problem can be formalised in this way:

$$\min C(n_1, \dots, n_H) = C_0 + \sum_{h=1}^H C_h n_h$$

under the constraints

$$\text{Var}(\hat{Y}_g) = \sum_{h=1}^H N_h^2 \left(1 - \frac{n_h}{N_h}\right) \frac{S_{h,g}^2}{n_h} \leq V_g \quad g = 1, \dots, G$$

where the  $V_g$  ( $g=1, \dots, G$ ) are the upper bounds for the expected sampling variance for  $\hat{Y}_1, \dots, \hat{Y}_G$ .

Bethel (1989) suggested that the problem can be more easily solved by considering the following function of  $n_h$ :

$$x_h = \begin{cases} 1/n_h & \text{if } n_h \geq 1 \\ \infty & \text{otherwise} \end{cases}$$

Using  $x_h$ , the cost function can be written as

$$C(x_1, \dots, x_H) = C_0 + \sum_{h=1}^H \frac{C_h}{x_h}$$

and the variances as

$$\text{Var}(\hat{Y}_g) = \sum_{h=1}^H N_h^2 \left(1 - \frac{1}{x_h N_h}\right) S_{h,g}^2 x_h = \sum_{h=1}^H N_h^2 S_{h,g}^2 x_h - N_h S_{h,g}^2 \quad g = 1, \dots, G$$

Consequently, the multivariate allocation problem can be defined as the search for the minimum (with respect to  $x_h$ ) of the convex function under a set of linear constraints

$$\sum_{h=1}^H N_h^2 S_{h,g}^2 x_h - N_h S_{h,g}^2 \leq V_g \quad g = 1, \dots, G$$

A numerical optimisation algorithm, that is proved to converge to the solution (if it exists), was provided by Bethel by applying the Lagrangian multipliers method to this problem.

It should be noted that there are also other algorithms for solving nonlinear programming problems in sample allocation, for instance, the *proc optmodel* in SAS offers alternatives like the trust region method, Newton-Raphson method with line search, conjugate gradient method, and a quasi-Newton method; the R packages *alabama* and *nloptr* use the augmented Lagrangian

algorithm and the method of moving asymptotes, respectively. In this paper, we focus on the use of Bethel's algorithm and genetic algorithm, through *SamplingStrata*, which allows performing the optimisation steps in two different ways, depending on the nature of the stratification variables  $X_s$ .

## 5.2 Optimisation with categorical stratification variables

Given a population frame with  $m$  auxiliary variables  $X_1, \dots, X_M$  we define as atomic stratification the one that can be obtained considering the cartesian product of the definition domains of the  $M$  variables. To each atomic stratum relevant information is attached:

- the values assumed by the stratification variables  $X_s$ ;
- the population  $N$  (number of units in the sampling frame belonging to the stratum);
- values of means and standard deviations associated to each target variable  $Y$ ;
- the average cost  $C$  of allocating a sampling unit in the stratum.

Starting from the initial atomic stratification, it is possible to generate, by differently aggregating the atomic strata, all the combinations that belong to the universe of stratifications. The number of possible different stratifications is exponential with respect to the number of the atomic strata. In concrete cases, it is therefore impossible to examine all the different possible alternative stratifications in order to individuate the best, *i.e.* the one of minimal cost. The genetic algorithm allows to explore the universe of stratifications in an efficient way, thus finding a solution not far from the optimal, by performing the following steps:

1. an initial set (*generation* in the terminology of the genetic algorithm) of stratifications (*individuals*) is randomly generated by aggregating the atomic strata: a given *individual* is a stratification where each atomic stratum is randomly attributed to one aggregate stratum identified by a combination of values of the stratification variables; each generated individual is characterised by a *genome*, *i.e.* a vector of integer numbers (*chromosomes*) indicating for each atomic stratum to which aggregate stratum it belongs;

2. for each aggregate stratum the information required (population, means and standard deviations of  $Y$ s, cost) is calculated and its *fitness* (total cost of the sample required to satisfy precision constraints) is determined by applying the Bethel algorithm;
3. the next set of individuals is generated by applying the usual operators of the genetic algorithm, *i.e.* *mutation*, *selection* and *crossover*.

Step 3 is repeated a given number of times. At the end, the individual with the best fitness (*i.e.* the stratification with the minimum cost of the associated sample) is retained as the best solution.

To clarify the above, let us consider a very simple example: a sampling frame with two stratification variables  $X_1$  and  $X_2$ , and related domains respectively (“A”, “B”) and (“1”, “2”, “3”). Considering the Cartesian product of the two domains, there will be six atomic strata:  $a_1=(\text{“A”},\text{“1”})$ ,  $a_2=(\text{“A”},\text{“2”})$ ,  $a_3=(\text{“A”},\text{“3”})$ ,  $a_4=(\text{“B”},\text{“1”})$ ,  $a_5=(\text{“B”},\text{“2”})$ ,  $a_6=(\text{“B”},\text{“3”})$ .

The initial step consists in randomly generating, say, 20 individuals (first generation), each one characterised by a genome.

For instance, the first individual could be  $I_1=(1,3,2,2,1,3)$ , that is a stratification characterised by three aggregated strata: according to the position of the elements in the genome, the first stratum is the aggregation of  $a_1$  and  $a_5$ , the second stratum is the aggregation of  $a_3$  and  $a_4$ , the third stratum is the aggregation of  $a_2$  and  $a_6$ .

For each one of these 20 stratifications, the corresponding fitness is calculated (applying the Bethel algorithm), as the cost of the sample necessary to be compliant with the precision constraints. The one with the best fitness (minimum cost) is retained as the optimal solution.

In order to generate the next generation of individuals, the 20 individuals are ordered by their fitness: if we set the *elitism rate* to 20%, the best 4 individuals will be retained as they are, with no change. Then, the remaining 80% individuals (16) will be generated in this way:

1. First, to each individual will be applied a *mutation* operator. Consider again the individual  $I_1$ . If we set the mutation chance equal to 5%, we scan the elements of its genome, each time generating a random number between 0 and 1: if it is less than 0.05, the element is changed

by assigning a random number, otherwise it is left unchanged. Suppose a mutation happens for the third element, that is changed from 2 to 1: now the genome of I1 is (1,3,1,2,1,3).

2. From the 16 elements, 16 couples are selected, with a selection probability proportional to their fitness (*selection operator*).
3. From each couple, a new individual is generated by applying the *crossover operator*. Consider a selected couple  $I1=(1,3,1,2,1,3)$  and  $I5=(3,2,1,1,1,1)$ . A number  $p$  is generated between 1 and the number of elements in the genome (6), for instance 2: the genome of the new individual will be given by the first 2 elements of I1 and the last 4 elements of I5, that is (1,3,1,1,1,1).

The new generation, composed by the 4 best individuals from the previous, and the new 16 obtained by mutation, selection and crossover, is now available. For each new individual will be calculated its fitness; if one of them has a better fitness than the current optimal solution, it will replace it.

The process continues until reaching the desired number of iterations.

### 5.3 Optimisation with continuous stratification variables

When all the stratification variables are continuous (or even categorical, but of the ordinal type), a variant of the above optimisation step is applicable. Instead of generating the atomic strata as a preliminary step, the algorithm provides to generate aggregate strata for each individual by operating in this way:

- for each continuous stratification variable, a predetermined number of values internal to its definition domain are randomly generated: these values (cuts) determine a segmentation of the domain that is equivalent to a categorisation of the variable;
- aggregate strata are consequently determined by cross-classifying units in the sampling frame according to their values belonging to the segments previously defined.

After this, the sequence of optimisation is identical to the one seen in the case of categorical stratification variables.

## 5.4 Anticipated variance

In real situations, the information contained in the sampling frame is not directly regarding the target variables of the survey, but proxy variables, *i.e.* variables that are correlated to the variables of interest. In our application, we know that income from self-employment collected in tax records is based on fiscal rules. In order to take into account this problem, and to limit the risk of overestimating the expected precision levels of the optimised solution, it is possible to carry out the optimisation by considering, instead of the expected coefficients of variation related to proxy variables, the anticipated coefficients of variation (ACV) that depend on the model that is possible to fit on couples of real target variables and proxy ones. In the current implementation, only models linking continuous variables can be considered. The definition and the use of these models is the same that has been implemented in the package *stratification* (Baillargeon and Rivest, 2012). In particular, the reference here is to two different models (applicable only to continuous variables):

1. the linear model with heteroscedasticity:  $Y = \beta \times X + \epsilon$ ,  
with  $\epsilon \sim N(0, \sigma^2 X \gamma)$  (where  $\gamma$  indicates the heteroscedasticity);
2. the log-linear model:  $Y = \exp(\beta \times \log(X) + \epsilon)$ , where  $\epsilon \sim N(0, \sigma^2)$ .

After fitting one model for each couple target / proxy variables, their parameters are given as an additional input to the optimisation function of *SamplingStrata*. The optimisation step will be then performed by calculating correctly the distributional values (means and standard deviations).



## 6. Application to the Italian HFCS

The method described in the previous Sections has been applied to the 2020 wave of the Italian HFCS survey. In particular, it has been used in the second stage of the design to select non-panel households, since the PSU and the panel households are considered fixed.

As already mentioned, register data use different concepts and definitions from the survey and they have also several quality issues. As a result, the information on household income coming from tax records is only a proxy of the actual economic situation.

As a first step, we estimate the goodness of these proxies. To this purpose, we use the refresh sample selected for the 2016 wave. These data have been linked to the Tax Register via individual ids. Considering respondents only, the link was successful for 4,328 households. For these units, we have information on the reported values for the five target variables (“Total income”, “Dependent employment income”, “Self-employment income”, “Pension income”, “Rents”) and the corresponding fiscal values. The associations between the two types of information are reported in Table 1.

**Table 1 - Linear regression models between observed variables and Tax Register variables (Italian HFCS, 2016 wave)**

Target variable	R2	Beta	Sigma
Total income	0.5771541	0.8417096	11945.78
Dependent employment income	0.6835152	0.8229064	12547.71
Self-employment income	0.2304688	0.5571044	18639.69
Pension income	0.6364706	0.7665643	5834.692
Rents	0.1366157	0.1653843	0.5436948

Source: Authors' own processing, 2018

There is an evident variability in the goodness of fitting: from a 68% in the case of “Dependent employment income” to a 13% in the case of “Rents”.

Using these models, we assign to each unit in the sampling frame the predicted values for each one of the variables of interest.

One may ask why we use data from the refresh sample for the 2016 wave, linked to Tax Register, only to fit the models, and we do not directly use it

to estimate the means, stratum variances, and the other quantities needed in the optimisation step. The answer is that this is necessary for two reasons: because the optimisation step with continuous stratification variables requires that their values are available for each unit in the sampling frame; and because, optimal strata values must be assigned to each unit in the frame when we select the final sample, and this can be done only knowing the values of the stratification variables.

As a second step, we chose the precision constraints in terms of the maximum expected coefficient of variation for the target mean estimates in the different domains (NUTS1 level Italian territorial units). The precision constraints are set equal to 5% in every domain and for all estimates.

We then run the optimisation step to define the stratification, the sample size, and its allocation. We use the sampling frame described in Section 4, containing 12,351,950 units (households). After removing the households with a source of income above 1 million euros (for the operational reasons previously explained), the resulting final population size is 12,334,342.

Numerous executions of this step have been attempted, varying the kind of optimisation (with categorical or continuous variables) and the maximum number of final strata. Even if stratification variables are continuous, we try the first algorithm after their categorisation (obtained by applying the univariate k-means clustering method). The comparison with the results obtained with the second algorithm (directly applied to stratification variables as they are) is in favour of the latter.

Another important decision is to fix the number of optimised strata to be expected in each one of the 5 territorial domains (NUTS1). This parameter is quite important in terms of the final results of the optimisation: in general, increasing the desired number of final strata determines a decrease of the sample size necessary to be compliant with the precision constraints, until a certain point, from which on, this number increases. Hints on which this point could be are given by using a particular function available in *SamplingStrata*, which performs a sequential application of a k-means algorithm, varying the number of the clusters (in this case coincident with the number of final strata) from a minimum (usually 2) to a maximum, for instance 20. The indication was to set this value to 10.

Another important parameter is the minimum number of units per stratum: too low, and the risk in case of high non-response is to have strata without respondents; too high, and the constraint may have a negative impact on the optimality of the solution. In our case, it was set to 50 households.

The optimisation has been carried out distinctly for the various domains. The number of iterations was set to 50, for each iteration 20 different solutions were generated, for a total of 1,000 solutions evaluated by applying the Bethel algorithm.

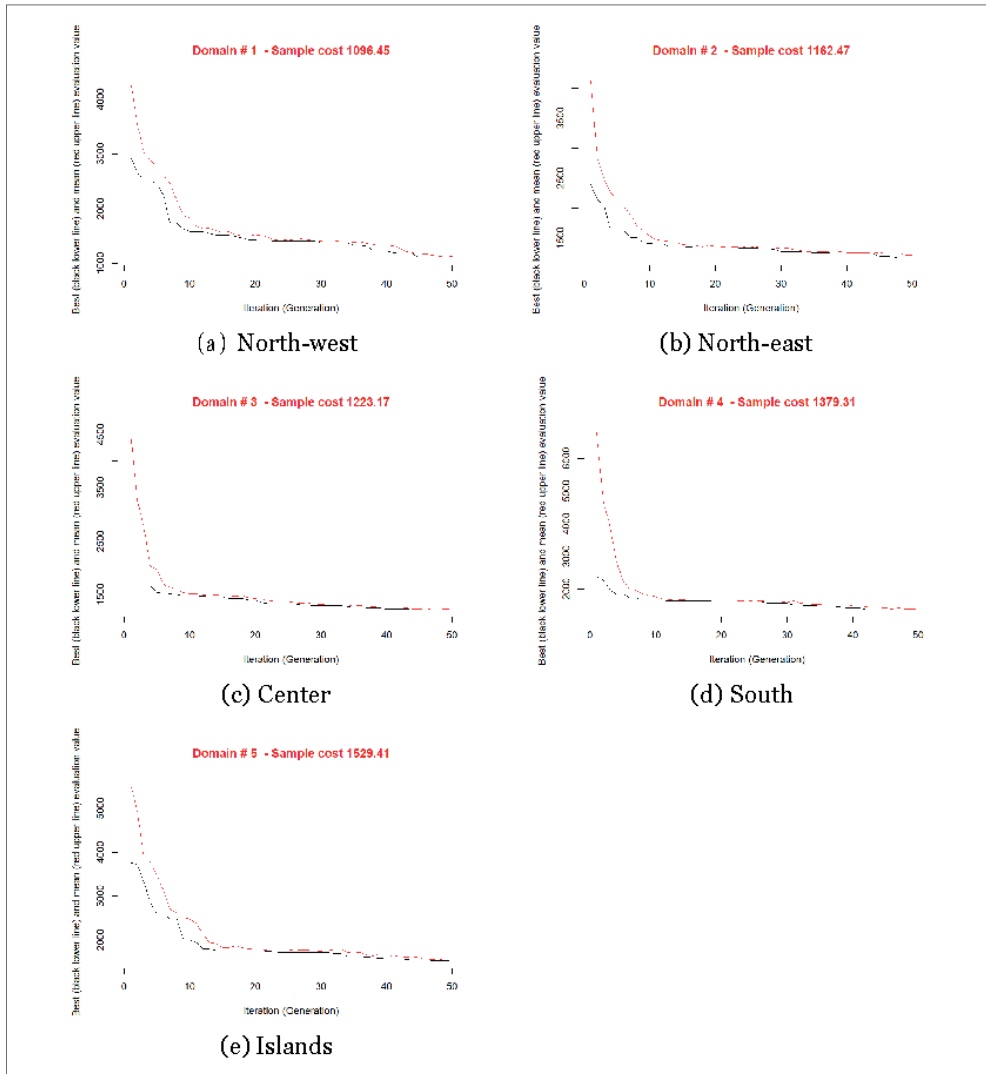
Figure 1 contains a graphical representation of the search for the optimal solution in the different domains. Each plot in this Figure can be interpreted in this way:

- in the x-axis are reported the different iterations (from 1 to 50): in correspondence to a given iteration, a set of 20 individuals have been generated, for each of them the Bethel solution in terms of sample size has been calculated;
- in the y-axis is reported the cost of the solution (in our case, the sample size);
- the red line represents the mean of the 20 Bethel solutions for each generation;
- the black line represents the cost of the best solution found so far.

Analysing these plots, a common situation for the different domains can be found: there is a smooth convergence towards the final solution of both the red and the black lines, and, more important, the lines towards the end are almost parallel to the x-axis, thus implying that adding more iterations should not increase substantially the optimality of the solution.

The overall sample size required to satisfy the precision constraints under the optimal solution is equal to 6,400.

Figure 1 - Optimisation in the different domains



Source: Authors' own processing, 2018

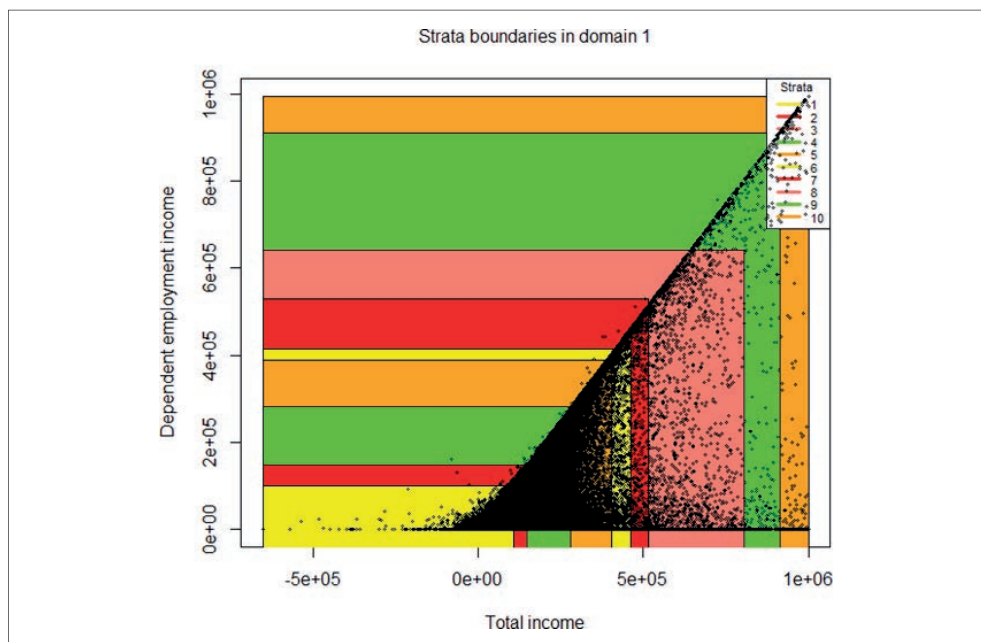
The package allows visualising in a two-dimensional graph the obtained strata, each time choosing a couple of variables. For instance, Figure 2 shows the characterisation of the strata in the first domain, by considering “Total income” and “Dependent employment income”. The points in the plot represent households in the sampling frame. Colours identify the different strata.

In Figure 3, optimised strata with population, sampling allocation, and sampling rates are reported together with the range of the two stratification variables. The intensity of the green is proportional to the values of Population and Allocation in strata, while the length of the red bar is proportional to the sampling rate.

Considering the two figures together, we can better understand this graphical representation. For instance, the first stratum (the yellow one) at the bottom left of the plot in Figure 2 includes all the households whose Total Income is less than -105,895 and Dependent Employment Income is less than 99,369. The second stratum (the red one) includes all the households whose Total Income is in the interval (-282,649; -147,433) and Dependent Employment Income is less than 162,801. The same interpretation for the other strata.

We do not report all the possible combinations of couples of variables because of their number (11 per each domain), we just report this one to show how strata appear, in their characteristic “7-shaped” format.

**Figure 2 - Strata resulting from the execution of the genetic algorithm (North-west, by *Dependent employment income* and *Total income*)**



Source: Authors' own processing, 2018

**Figure 3 - Strata population, allocation and range of stratification variables (North-west)**

Stratum	Population	Allocation	SamplingRate	Bounds Total income	Bounds Dependent employment income
1	2384770	365	0.0001529532	-652064-105895	0-99369
2	523170	272	0.0005193699	-282649-147433	0-162801
3	59346	55	0.0009336701	-421027-149039	0-149002
4	35528	50	0.0014073407	25566-280219	0-281464
5	58824	105	0.0017776216	12513-405510	0-387762
6	3846	50	0.0130005209	93801-462718	0-413450
7	2203	50	0.0226963232	161895-515268	0-530787
8	5028	50	0.0099443119	59136-804789	0-641535
9	938	50	0.0533049041	392890-912168	0-909600
10	1137	50	0.0439753738	511877-999682	0-995589

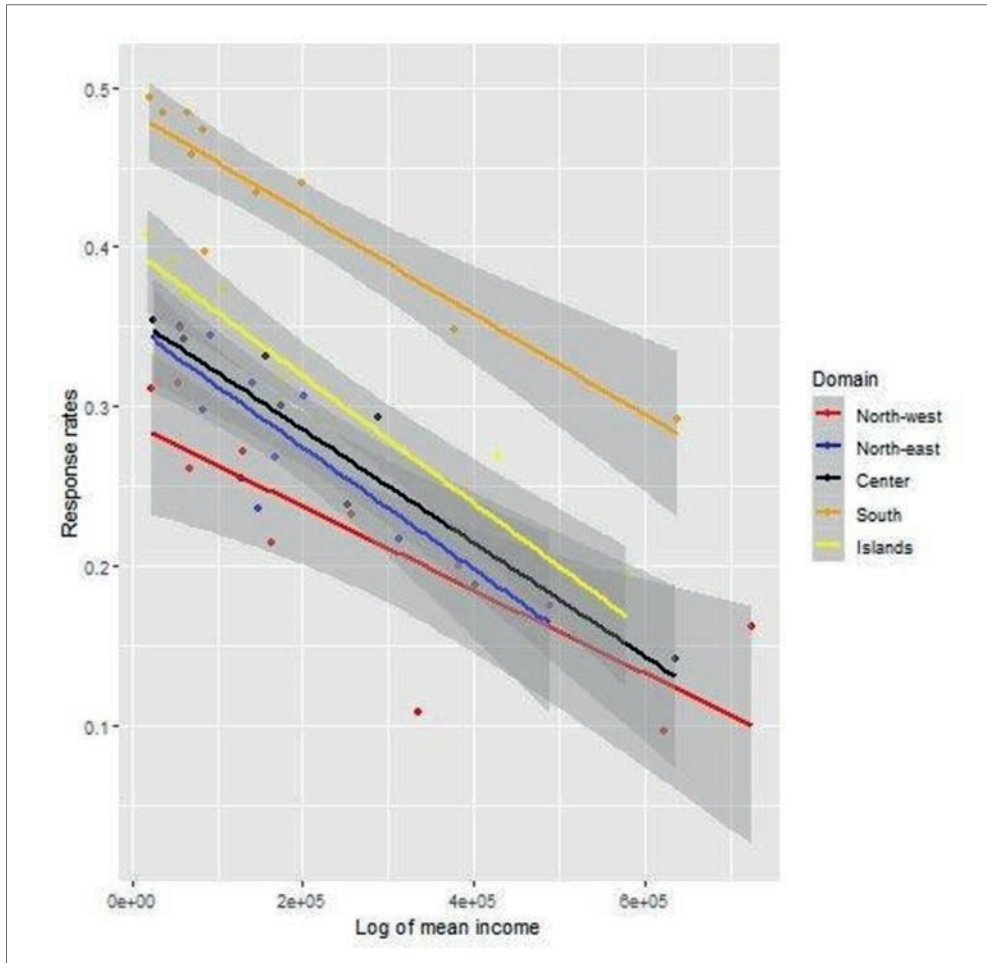
Source: Authors' own processing, 2018

The solution is characterised by a sample size equal to 6,400, and the expected coefficients of variations have been calculated assuming that all sampled units will respond to the interviewers.

In order to take the expected non-response rate into account, as a final step, we need to estimate the total sample that is required to get a final sample of around 6,400 households. Using the sample selected for the 2016 survey linked to tax records, we link both respondents and non-respondents to the Tax Register. We then estimate a model for the probability of participating in the survey using as predictors the four components of income (with the exclusion of the “Total income”) and the twenty NUTS2 Italian regions. Considering the plot in Figure 4, there is clear evidence of a linear direct inverse relationship between the log of the mean income in a stratum, and the propensity to respond. In Figure 4 we also report confidence bands around the lines, based on model standard errors.

The sample of units to be interviewed has been redefined by taking into account the propensity to non-response calculated for each unit in the sampling frame using the above model. The total number of households to be interviewed is 17,608, units that have been allocated in the optimised strata taking into account the initial allocation and the average propensity to the response in strata.

**Figure 4 - Response rate and mean income in strata**



Source: Authors' own processing, 2018

For example, in Table 2 has been reported the final solution, with the initial and final allocation, for the first domain.

**Table 2 - Optimal Stratification, initial and final allocation**

Domain	Stratum	Population	Initial Allocation	Final allocation	Sampling rate
1	1	2,384,770	365	1064	0.000446
1	2	523,170	272	784	0.001499
1	3	59,346	55	192	0.003235
1	4	35,528	50	211	0.005939
1	5	58,824	105	350	0.005950
1	6	3,846	50	195	0.050702
1	7	2,203	50	420	0.190649
1	8	5,028	50	226	0.044948
1	9	938	50	469	0.500000
1	10	1,137	50	280	0.246262

Source: Authors' own processing, 2018

Table 3 reports the coefficients of variation achievable with the selected sample (6,400 units). The solution allows meeting all the precision requirements. It can be seen that for the first variable (“Total income”) the precision is about double than prescribed.

**Table 3 - Expected coefficients of variation on estimates of the mean (%)**

Domain	Total income	Dependent emp. income	Self-employment income	Pension income	Rents
1. North-west	2.5	5.0	4.8	4.9	4.8
2. North-east	2.4	4.7	4.9	4.6	4.8
3. Centre	2.6	4.8	5.0	5.0	4.7
4. South	2.3	4.3	4.8	4.8	4.9
5. Islands	2.3	4.0	5.0	5.0	4.9

Source: Authors' own processing, 2018

These estimates of the expected CVs have been calculated (using a specific function in the package *SamplingStrata*) assuming that:

1. the survey adopts a single stage sampling process;
2. estimates are obtained by Horvitz-Thompson estimator;
3. all 6,400 units in the sample respond to the survey.

Of course, none of these assumptions hold in reality. In particular, assuming (1) and (3) leads to an under-estimation of the real values of the coefficients of variation, while the (2) might over-estimate them. For this reason, in the next Section we present the results of a simulation exercise that takes this issue into account.



## 7. Evaluation of the new sample design

In this Section, we run several simulations to have a more robust evaluation of the new design. Each simulation is based on the archive created by linking the Population and the Tax registers, and on the information coming from the 2016 SHIW survey integrated with tax records.

In the simulations, we extract 500 samples using both the new and the old design and we compute measures of precision and bias of the five income estimators. The difference between the two types of simulation is the following. In the first set, we only use the information in the Population Register for the calibration of final weights, in line with what is currently done in the SHIW survey. In the second set of simulations, we also use tax records in the weighting stage.

Each simulation is based on the following assumptions:

1. the survey uses a two-stage sampling design, so when evaluating variance of estimates, weights associated with Primary Sampling Units (the municipalities selected at the first stage) have to be taken into account;
2. estimates are obtained by calibration estimators, to handle total non-response;
3. sample size has been inflated to 17,608 households to take into account the expected non-response.

### 7.1 Simulations using Population Register for calibration

The first simulation consists of the following steps.

First, we use the models introduced in Section 6 to predict, for each unit in the sampling frame, the values of target variables.

Then, 500 samples of the required size (17,608 households) are selected from the sampling frame. For each household, we simulate the non-response mechanism using the model described in Section 6. The decision to participate is then taken by drawing a value from a Bernoulli variable with the probability of success (the propensity to respond) equal to the propensity estimated by the non-response model.

For each sample of respondents, initial weights are computed considering the probabilities of inclusion of both first and second stage, and the final weights are obtained by calibrating using the total number of households in the strata in the Population Register, as defined by the new design.

In the end, for each target estimate (means of total income and of the four components), coefficients of variation and relative bias have been calculated, averaging over the 500 replicated samples. Bias is measured as the difference between the mean value of the 500 survey-based estimates and the population means coming from administrative records.

Results are reported in Tables 4 and 5. The precision of the estimators is in line with one of the selected samples.

**Table 4 - Estimated coefficients of variation of the new sample design (%)**

Domain	Total income	Dependent emp. income	Self-employment income	Pension income	Rents
1. North-west	2.6	5.4	4.3	4.8	3.5
2. North-east	2.4	4.8	4.5	4.5	3.3
3. Centre	2.4	4.8	3.7	4.5	3.1
4. South	2.3	4.4	3.5	4.6	3.3
5. Islands	2.3	4.1	3.5	4.8	3.1

Source: Authors' own processing, 2018

The simulation shows the presence of a negative bias for incomes from employment and rents. The opposite situation holds for incomes from self-employment and pensions. The presence of bias depends on our response probability model, which is estimated using household-specific administrative information. In some strata, this model generates a high (within) variability of response propensities. Therefore, a simple calibration of the weights of respondents to the total number of households in the population is not enough to compensate for missing households.

**Table 5 - Estimated relative bias of the new sample design (%)**

Domain	Total income	Dependent emp. income	Self-employment income	Pension income	Rents
1. North-west	-3.8	-12.8	5.6	8.3	-1.5
2. North-east	-2.6	-8.8	3.0	6.4	-2.0
3. Centre	-2.7	-10.5	4.5	8.5	-1.5
4. South	-2.2	-6.5	0.6	4.4	-3.0
5. Islands	-2.1	-8.0	1.4	5.7	-1.4

Source: Authors' own processing, 2018

The old sample design is a two-stage process where the first stage is identical to the new one, with the selection of the same 454 municipalities (via PPS). The allocation of SSU units is based on the following rule: if the total population in the selected municipality is higher than 500,000 then 200 households are assigned, otherwise only 32. The total number of SSU units is 14,864. Based on this SSU stratification and allocation, we run a sample of 6,400 units for the frame. This sample represents therefore the one we have selected using the old design. The expected CVs for the selected sample are reported in Table 6.

**Table 6 - Expected coefficients of variation for the old sample design (%)**

Domain	Total income	Dependent emp. income	Self-employment income	Pension income	Rents
1. North-west	4.6	6.3	23.1	6.4	17.2
2. North-east	3.7	5.0	20.4	5.6	15.4
3. Centre	4.3	5.8	23.4	6.8	17.3
4. South	4.6	5.8	25.4	6.7	21.3
5. Islands	6.0	7.8	32.5	9.8	26.1

Source: Authors' own processing, 2018

This table has been computed using the same assumptions made for Table 3. By comparing the two, it is clear that the expected CVs for the old design are higher than those calculated for the new one. In particular, they are much higher for Self-employment income and Rents.

For comparison, we report in Tables 7 and 8 the observed CVs of the target variables computed using the 2014 and 2016 Italian HFCS. These tables are not directly comparable with the previous ones for two main reasons. First, the sample size is larger (about 8,000 households for each wave). Second, the sampling weights are calibrated in a way that is not possible for the 2020 survey since we miss some demographic information on respondents. The possibility to calibrate using other information (such as the job status) contributes to reducing the final variability of the estimators. Still, two important points can be drawn from these tables. First, the expected CVs shown in this paper are probably upper bounds for the actual ones that will be observed for the 2020 wave. Second, the advantage of the new design is also in reducing the instability of the estimators across surveys. This is particularly the case for incomes from self-employment and rents, which show significant changes in the precision from one wave to another. This is because the

available information does not allow us to have full control of the final sample composition. This situation will change thanks to the new design.

**Table 7 - Coefficients of variation estimated in the 2016 Italian HFCS wave (%)**

Domain	Total income	Dependent emp. income	Self-employment income	Pension income	Rents
1. North-west	2.1	3.8	11.0	4.1	25.4
2. North-east	3.4	4.4	12.5	4.1	14.2
3. Centre	2.3	5.3	11.5	4.8	18.9
4. South	2.5	4.6	14.1	4.5	28.5
5. Islands	3.1	4.9	22.4	5.2	34.4

Source: Authors' own processing, 2018

**Table 8 - Coefficients of variation estimated in the 2014 Italian HFCS wave (%)**

Domain	Total income	Dependent emp. income	Self-employment income	Pension income	Rents
1. North-west	2.2	2.9	8.2	3.8	11.3
2. North-east	1.9	2.7	9.6	4.4	14.4
3. Centre	2.4	3.4	18.3	4.0	10.4
4. South	2.8	4.6	21.4	3.7	22.2
5. Islands	2.7	4.1	12.7	7.3	44.4

Source: Authors' own processing, 2018

Following the same approach previously used, we then run a simulation based on the old design. In particular, we perform the following steps:

1. 500 samples have been drawn from the same sampling frame, *i.e.* the one enriched by predicted target variables;
2. for each sample, the mechanism of non-response has been simulated accordingly to the predicted non-response propensity associated with each unit in the frame;
3. for each resulting sample of respondents, calibrated estimates of interest have been calculated, where known totals are given by the number of households by strata in the Population Register as defined in the old design.

In other words, the simulation has been carried out with the same setting used for the new sample design.

In the end, coefficients of variation and relative bias for the old sample design have been calculated, averaging over the 500 replicated samples. Results are reported in Tables 9 and 10.

**Table 9 - Estimated coefficients of variation of the old sample designs (%)**

Domain	Total income	Dependent emp. income	Self-employment income	Pension income	Rents
1. North-west (Rip1)	6.18	9.50	31.73	10.34	6.61
2. North-east (Rip2)	4.96	8.17	25.77	9.34	5.48
3. Centre (Rip3)	5.51	8.68	22.22	10.00	5.68
4. South (Rip4)	4.85	7.54	19.60	8.18	5.30
5. Islands (Rip5)	7.42	11.33	29.30	14.26	7.79

Source: Authors' own processing, 2018

**Table 10 - Estimated relative bias of the old sample designs (%)**

Domain	Total income	Dependent emp. income	Self-employment income	Pension income	Rents
1. North-west (Rip1)	-5.92	-14.55	-6.78	8.12	-2.71
2. North-east (Rip2)	-4.20	-10.15	-6.26	5.97	-2.03
3. Centre (Rip3)	-4.60	-11.39	-7.39	7.14	-2.55
4. South (Rip4)	-2.85	-7.19	-3.8	4.03	-2.12
5. Islands (Rip5)	-3.18	-8.36	-4.40	4.71	-2.22

Source: Authors' own processing, 2018

**Figure 5 - Comparison of coefficients of variation obtained for the new and old sample designs**



Source: Authors' own processing, 2018

Figures 5 and 6 summarise the over-performance of the new sample compared to the old sample in terms of both coefficients of variation and bias, respectively.

It can be seen that as for the CVs, there is a clear indication of the superiority of the new design compared to the old one in terms of the sampling variance component of the Mean Squared Error (MSE).

As for the bias, the new sample design is still better, but there are 6 cases out of 25 in which the old design performs better.

**Figure 6 - Comparison of relative bias obtained for the new and old sample designs**

Source: Authors' own processing, 2018

## 7.2 Simulations using Tax Register for calibration

In the previous simulations, we did not use the known totals available from the Tax Register, *i.e.* the sum of the components of the income (Dependent Employment, Self-Employment, Pensions, Rents) by the different domains of interest (the five Italian NUTS1 geographical zones).

To fully exploit the information achievable in the administrative sources, we carried out the same simulations described before but using a different calibration model: instead of the known totals of households in the strata defined by the old and new sampling designs, we made use of both totals of households at NUTS1 level and the Tax Register incomes at stratum level.

Results in terms of CVs and bias are reported in Tables 11 and 12.

**Table 11 - Estimated coefficients of variation of the new and old sample designs (%) with calibration using Tax Register variables**

Domain	Total income		Dependent emp. income		Self-employment income		Pension income		Rents	
	New	Old	New	Old	New	Old	New	Old	New	Old
1. North-west (Rip1)	0.54	0.99	0.23	0.41	1.93	3.10	0.53	1.06	3.34	6.17
2. North-east (Rip2)	0.46	0.72	0.21	0.39	1.80	3.26	0.48	0.72	2.73	3.70
3. Centre (Rip3)	0.51	0.76	0.24	0.31	2.02	3.42	0.53	0.79	2.70	4.20
4. South (Rip4)	0.55	0.71	0.24	0.48	2.16	2.89	0.56	0.84	2.89	3.80
5. Islands (Rip5)	0.52	1.21	0.24	0.63	2.14	4.33	0.54	1.16	2.61	5.76

Source: Authors' own processing, 2018

**Table 12 - Estimated relative bias of the new and old sample designs (%) with calibration using Tax Register variables**

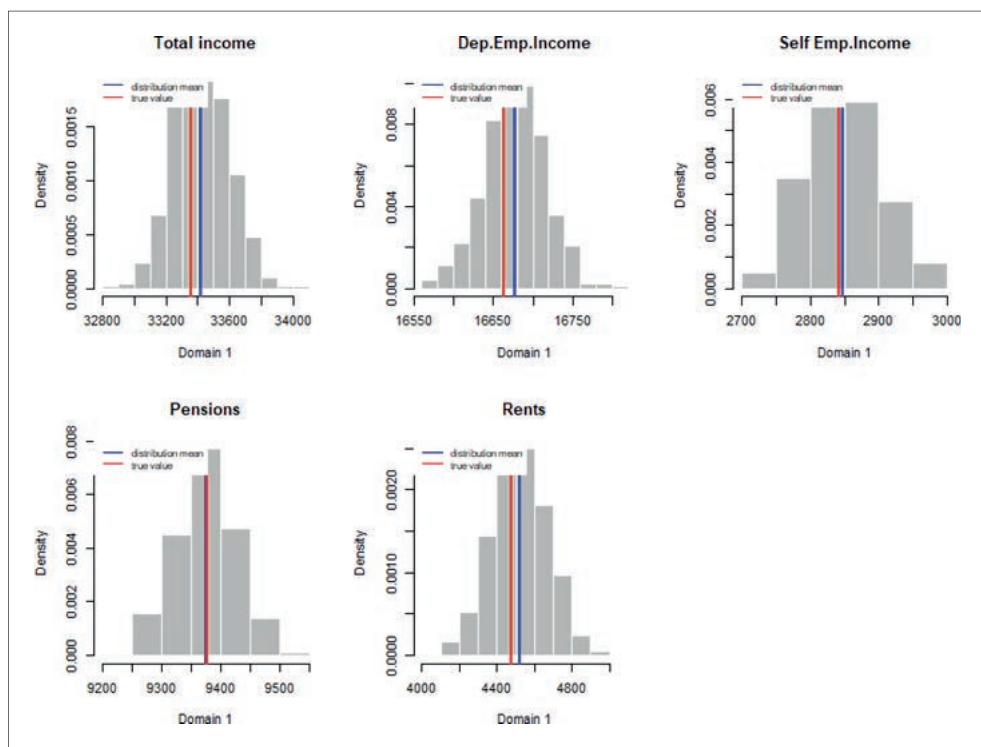
Domain	Total income		Dependent emp. income		Self-employment income		Pension income		Rents	
	New	Old	New	Old	New	Old	New	Old	New	Old
1. North-west (Rip1)	0.19	0.32	0.08	0.16	0.20	0.63	-0.02	-0.28	1.00	1.93
2. North-east (Rip2)	0.00	-1.77	0.08	-2.09	0.17	-1.20	-0.04	-2.48	-0.26	0.44
3. Centre (Rip3)	0.22	0.59	0.08	0.00	-0.02	0.94	-0.01	0.17	1.25	3.21
4. South (Rip4)	-0.21	-1.43	0.02	-2.02	-0.57	-0.58	-0.07	-1.88	-0.99	0.78
5. Islands (Rip5)	0.09	0.44	0.05	-0.17	0.22	2.26	0.03	-0.53	0.27	3.06

Source: Authors' own processing, 2018

The distribution of the 500 replicated estimates is reported in Figure 7, only for the first domain and only for the new sample design.



**Figure 7 - Distribution of the 500 replicated estimates in the first domain (new design, calibration adding Tax Register totals)**



Source: Authors' own processing, 2018

There is an evident reduction of CVs and bias for both new and old sample design, with a comparison always in favour of the new design.

This simulation is only indicative of the potential of this calibration, because results so positive depend on the fact that the target values in the frame have been generated by models that make use of the Tax Register variables as explanatory variables. Using the same Tax Register variables as known totals in the calibration model introduces a great simplification of the real situation, that may somehow compromise the full validity of these results. Nonetheless, it is expected that a model-assisted approach which also includes Tax Register variables would substantially improve the accuracy of the estimates.

## 8. Conclusions

The paper presents an empirical application of tax personal income data in the sampling design of finance surveys. Tax data are not collected for statistical purposes and therefore they use definitions and measures different from those adopted in the survey. Furthermore, they are subject to various quality problems (such as tax avoidance or evasion, the presence of thresholds below which the declaration is not necessary, and time delays before becoming available).

As a consequence, their use for statistical purposes is not straightforward. Nonetheless, this application has shown that one possible solution is to consider them as proxies for the variables of interest and to inflate the estimators of variance used for determining sample size accordingly. We are able to estimate the goodness of these proxies by linking survey data to administrative records. Our simulations show that their use enables us to take under control the expected accuracy of income estimators, despite all the limits of tax data. A second (and strictly related) advantage is that the availability of register data enables us to keep under control the fieldwork of the survey. This implies, for instance, specific households can be oversampled and those refusing to participate could be replaced with others belonging to the same stratum. This should guarantee to obtain a final sample, which is very close to the selected one, *i.e.* the most efficient one. Consequently, the expected benefits in terms of variance reduction should turn into effective advantages.

Another potential advantage is linked to the possibility of reducing bias due to non-response. Our simulation has shown that the new sample design allows not only greatly reducing the sampling variance, but also the bias component of the Mean Square Error of estimates even if we do not include Tax Register variables in the calibration model. If we also include these variables, results in terms of an overall reduction of MSE should be even greater.

## References

Albarea, A., M. Bernasconi, A. Marenzi, and D. Rizzi (eds.). 2018. “Income under reporting and tax evasion in Italy. Estimates and distributive effects” *Documento di Valutazione*, N. 8. Roma: Senato della Repubblica, Ufficio Valutazione di Impatto/*Impact Assessment Office*.

Baillargeon, S., and L.-P. Rivest. 2012. “Univariate Stratification of Survey Populations”. *R Package*, Version 2.2-3. The Comprehensive R Archive Network – CRAN.

Ballin, M., and G. Barcaroli. 2016. “Optimization of Stratified Sampling with the R Package SamplingStrata: Applications to Network Data”. In Dehmer, M., Y. Shi, and F. Emmert-Streib. *Computational Network Analysis with R: Applications in Biology, Medicine, and Chemistry. Volume 7*. Hoboken, NJ, U.S.: John Wiley & Sons.

Ballin, M., and G. Barcaroli. 2013. “Joint determination of optimal stratification and sample allocation using genetic algorithm”. *Survey Methodology*, Volume 39, N. 2: 369-393.

Ballin, M., G. Barcaroli, M. Masselli, and M. Scarnò. 2018. “Redesign sample for Land Use/Cover Area frame Survey (LUCAS) 2018”. *Statistical Working Papers*, Eurostat. Luxembourg: Publications Office of the European Union.

Barcaroli, G. 2014. “SamplingStrata: An R Package for the Optimization of Stratified Sampling”. *Journal of Statistical Software*, Volume 61, Issue 4: 1–24.

Barcaroli, G., M. Ballin, H. Odendaal, D. Pagliuca, E. Willighagen, and D. Zardetto. 2020. “SamplingStrata: Optimal Stratification of Sampling Frames for Multipurpose Sampling Surveys”. *R Package*, Version 1.5-1. The Comprehensive R Archive Network – CRAN.

Bethel, J. 1989. “Sample allocation in multivariate surveys”. *Survey Methodology*, Volume 15, N. 1: 47–57.

Bover, O. 2010. “Wealth Inequality And Household Structure: U.S. vs. Spain”. *Review of Income and Wealth*, Volume 56, Issue 2: 259–290.

Bover, O., E. Coronado, and P. Velilla. 2014. “The Spanish Survey of Household Finances (EFF): Description and Methods of the 2011 Wave”. *Documentos Ocasionales/Occasional Papers*, N. 1407. Madrid, Spain: Banco de España.

Casiraghi, M., E. Gaiotti, L. Rodano, and A. Secchi. 2018. “A “reverse Robin Hood”? The distributional implications of non-standard monetary policy for Italian households”. *Journal of International Money and Finance*, Volume 85: 215–235.

Chakraborty, R., I.K. Kavonius, S. Pérez-Duarte, and P. Vermeulen. 2019. “Is the top tail of the wealth distribution the missing link between the Household Finance and Consumption Survey and national accounts?”. *Journal of official Statistics - JOS*, Volume 35, Issue 1: 31–65.

Cochran, W.G. 1977. *Sampling Techniques. Third edition*. New York, NY, U.S.: John Wiley & Sons.

Colciago, A., A. Samarina, and J. de Haan. 2019. “Central Bank Policies and Income and Wealth Inequality: A Survey”. *Journal of Economic Surveys*, Volume 33, N. 4: 1199–1231.

D’Alessio, G., and A. Neri. 2015. “Income and wealth sample estimates consistent with macro aggregates: some experiments”. *Questioni di Economia e Finanza, Occasional Papers*, N. 272. Roma, Italy: Banca d’Italia.

Dobbs, R., S. Lund, T. Koller, and A. Shwayder. 2013. “QE and ultra-low interest rates: Distributional effects and risks”. *Discussion Paper*, McKinsey Global Institute. New York, NY, U.S.: McKinsey & Company.

Dossche, M., J. Slačálek and G. Wolswijk. 2021. “Monetary policy and inequality”. In *ECB Economic Bulletin*, Issue 2/2021. Frankfurt am Main, Germany: European Central Bank – ECB.

Eckerstorfer, P., J. Halak, J. Kapeller, B. Schütz, F. Springholz, and R. Wildauer. 2016. “Correcting for the Missing Rich: An Application to Wealth Survey Data”. *The Review of Income and Wealth*, Volume 62, Issue 4: 605–627.

Eurosystem Household Finance and Consumption Network. 2009. “Survey Data on Household Finance and Consumption. Research Summary and Policy Use”. *Occasional Paper Series*, N. 100. Frankfurt am Main, Germany: European Central Bank – ECB.

Fiorio, C.V., and F. D'Amuri. 2006. "Tax Evasion In Italy: An Analysis Using A Tax-Benefit Microsimulation Model". *The IUP Journal of Public Finance*, Volume IV, Issue 2: 19–37.

Hansen, M.H., W.N. Hurwitz, and W.G. Madow. 1953. *Sample Survey Methods and Theory: Volumes I-II*. Hoboken, NJ, U.S.: John Wiley & Sons.

Horgan, J.M. 2006. "Stratification of Skewed Populations: A review". *International Statistical Review*, Volume 74, N. 1: 67–76.

Household Finance and Consumption Network - HFCN. 2020. "The Household Finance and Consumption Survey: methodological report for the 2017 wave". *Statistics Paper Series*, N. 35. Frankfurt am Main, Germany: European Central Bank – ECB.

Jäntti, M., V.-M. Törmälehto, and E. Marlier (eds.). 2013. "The use of registers in the context of EU–SILC: challenges and opportunities". Eurostat, *Statistical working papers*. Luxembourg: Publications Office of the European Union.

Kareem, A.O., I.O. Oshungade, G.M. Oyeyemi, and A.O. Adejumo. 2015. "Moving Average Stratification Algorithm for Strata Boundary Determination in Skewed Populations". *CBN Journal of Applied Statistics*, Volume 6, N. 1: 205–217.

Kennickell, A.B. 2019. "The tail that wags: differences in effective right tail coverage and estimates of wealth inequality". *The Journal of Economic Inequality*, Volume 17, Issue 4, N. 1: 443-459.

Kennickell, A.B. 2017. "Modeling Wealth with Multiple Observations of Income: Redesign of the Sample for the 2001 Survey of Consumer Finances". *Statistical Journal of the IAOS*, Volume 33, Issue 1: 51-58.

Kennickell, A.B. 2008. "The Role of Over-sampling of the Wealthy in the Survey of Consumer Finances". In Bank for International Settlements (ed.). *The IFC's contribution to the 56<sup>th</sup> ISI Session*, Volume 28: 403-408. Lisbon, August 2007.

Khan, M.G.M., N. Nand, and N. Ahmad. 2008. "Determining the optimum strata boundary points using dynamic programming". *Survey Methodology*, 34, N. 2: 205–214.

Michelangeli, V., and C. Rampazzi. 2016. “Indicators of financial vulnerability: a household level study”. *Questioni di Economia e Finanza, Occasional Papers*, N. 369. Roma, Italy: Banca d’Italia.

Neri, A., and M.G. Ranalli. 2011. “To Misreport or not to report? The Case of the Italian Survey on Household Income and Wealth”. *Statistics in Transition - new series*, Volume 12, N. 2: 281–300.

Neri, A., and R. Zizza. 2010. “Income reporting behaviour in sample surveys”. *Temi di discussione, Working papers*, N. 777. Roma, Italy: Banca d’Italia.

Schmitt, L.M. 2001. “Fundamental Study. Theory of genetic algorithms”. *Theoretical Computer Science*, Volume 259, Issues 1-2: 1–61.

Valliant, R., J.A. Dever, and F. Kreuter. 2018. *Practical Tools for Designing and Weighting Survey Samples. Second Edition*. New York, NY, U.S.: Springer.

Vermeulen, P. 2018. “How fat is the top tail of the wealth distribution?”. *The Review of Income and Wealth*, Volume 64, Issue 2: 357–387.

# An analysis of the influence of tunnel length and road type on road accident variables

Antonella Pireddu <sup>1</sup>, Silvia Bruzzone <sup>2</sup>

## Abstract

*In 2018, the Italian National Institute of Statistics - Istat recorded about 1,150 injuries and 700 accidents in Italian tunnels where, over the past five years, there has been an increase in the accident rate. A nationwide study, conducted in order to investigate some aspects of this phenomenon, analysed the relationship between the class of vehicle involved, the time of the accident, the trip purpose (work-related or non-work-related), the circumstances observed and the infrastructure characteristics in terms of tunnel length and road type. The study offers a methodological reference for all those who, with a reactive or proactive approach, are required to assess road safety, manage risks in a specific context and evaluate the effectiveness of prevention strategies.*

**Keywords:** Accidents, accident circumstances, Principal Component Analysis, road tunnels, road type, tunnel length.

---

1 Italian National Institute for Insurance against Accidents at Work - INAIL ([an.pireddu@inail.it](mailto:an.pireddu@inail.it)).

2 Italian National Institute of Statistics - Istat ([bruzzone@istat.it](mailto:bruzzone@istat.it)).

*The views and opinions expressed are those of the authors and do not necessarily reflect the official policy or position of the Italian National Institute of Statistics - Istat.*

*The authors would like to thank the anonymous reviewers for their comments and suggestions, which enhanced the quality of this article.*

## 1. Introduction<sup>3</sup>

Road accidents in Italian tunnels have increased during the period 2013-2017. Data concerning tunnel collisions resulting in death or injury occurred in Italy in 2018 were analysed by tunnel type. A descriptive and multivariate analysis was applied to detect the relationship between tunnel type and road crash parameters. Absolute frequency was higher in shorter tunnels while relative frequency (per kilometre) was negatively correlated with length.

A positive association was found between short urban and motorway tunnels and 4-wheel vehicles, non-work-related trip and factors such as distances between vehicles, distraction and speeding, while urban and rural tunnels were positively associated with motorcycles, work-related trip and unspecified circumstances. The study points out essential findings for further targeted interventions in the tunnel road safety.

---

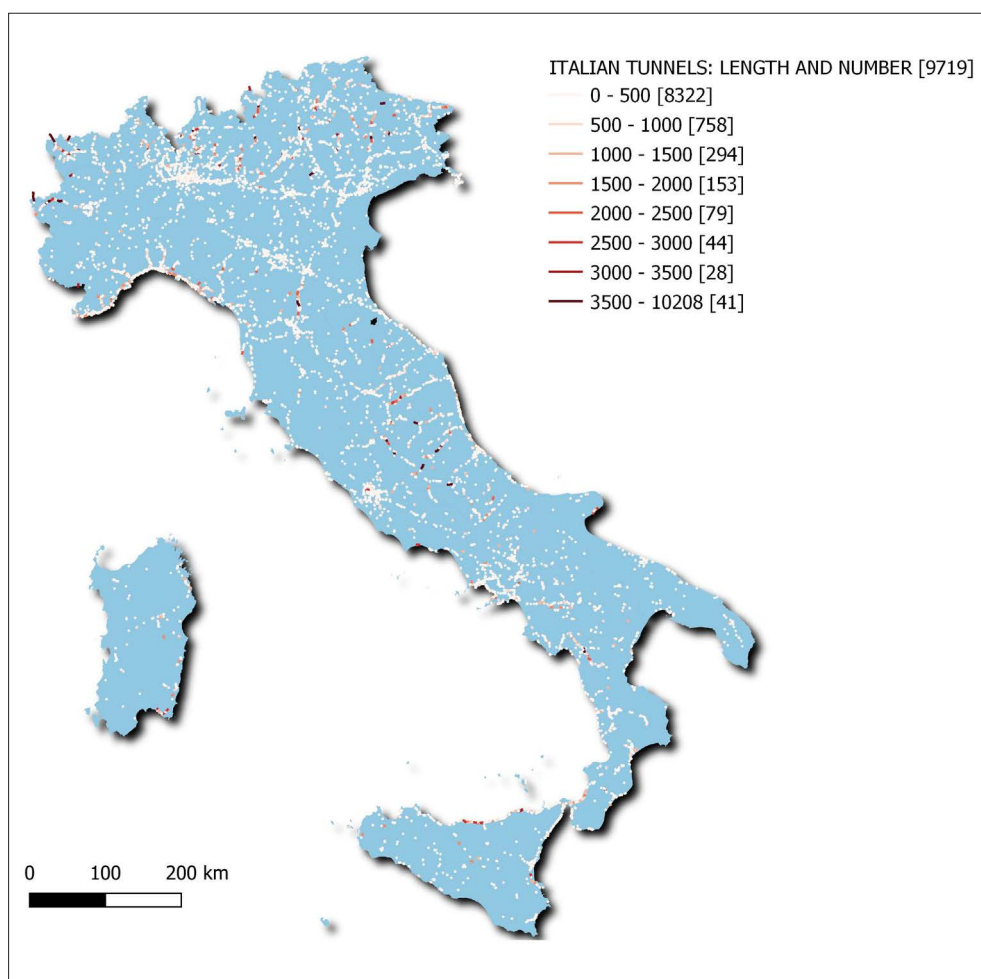
3 In producing this article: Antonella Pireddu dealt mainly with conceptualisation, data processing, methodology, writing the original draft, review of contents and results; Silvia Bruzzone dealt in particular with data processing, methodology, validation, review of content and results.



## 2. The Italian situation

In Italy, road tunnels extending for approximately 2,600 km, are distributed throughout the country: 907 km in the North West, 488 km in the North East, 410 km in the Centre, 527 km in the South and 268 km in the Islands. Underpasses, which make up 10% of the total, were also included in the analysis of the data. The authors provide an overview of the Italian road tunnel network, aggregated by classes of length (Figure 1).

**Figure 1 - Italian tunnels by length and number. Year 2018**



Source: Authors' processing on OpenStreetMap data

Figure 1 shows the distribution of tunnels classified by length within the national territory. The class of tunnels of up to 500 m is the largest with over 8,300 units, while there are approximately 750 tunnels of 500 to 1,000 m in length. The remaining classes are less numerous and decrease from about 300 (1,000-1,500 m) to 30 (3,000 - 3,500 m) units. The above illustration is the result of the geoprocessing of *OpenStreetMap* (OSM) data with *Quantum Gis* (QGIS) carried out by the authors. In the study, the tunnel type label represents a combination of the tunnel road (motorway, rural, urban) and length (Legislative Decree/*Decreto Legislativo* 30 April 1992, n. 285).

According to Istat data, in the five-year period 2013-2017, 3,175 accidents occurred in Italian road tunnels, 98 of which were fatal, and no less than 193 were work-related, involving at least one worker in route to/from work or driving as part of the work. In the same five-year period, 5,022 injuries were recorded. These involved at least 304 workers, 94 deaths within 24 hours and 11 deaths in the first 30 days following the accident. A 2019 study reported an increase in frequency for both collisions and injuries, with the latter increasing from 921 in 2013 to 1,161 in 2017 (Pireddu *et al.*, 2019). This increase also included work-related cases. Due to incomplete data for the five-year period, this study provides only a descriptive analysis without the accident parameters involved, the correlation between the latter and the type of tunnel affected. The trend observed between 2013 and 2017 was confirmed in 2018 when 716 road accidents were reported, with 1,159 injuries involving at least 77 work-related road injuries and 18 deaths. Throughout Italy, about 490 km out of a total of 2,600 km of tunnels were affected.

### 3. Data source

Statistical information on road accidents are produced by the Italian National Institute of Statistics - Istat on the basis of a survey of all road accidents occurring in Italy: the “*Survey on road accidents resulting in death or injury*”. The field of observation of the survey includes all road accidents involving deaths within 30 days or injuries that occur throughout the country over a one-year period and are recorded by a police authority. Detection refers to the time the accident occurred. Istat provides data on all reported collisions (Istat data warehouse *I.Stat.* 2018, 2019) involving at least one vehicle that occur on Italian roads.

The data considered for the purposes of this paper, as defined by national standards, are those where “at least one vehicle is involved and where at least one injured person is recorded” (Vienna Convention on Road Traffic, 1968)<sup>4</sup>.

Therefore, either accidents without injuries, or that did not occur in public roads, or without vehicles involved, or without police report were excluded (Gariazzo *et al.*, 2019).

This study data included information regarding the vehicle, the time of the accident, the trip purpose (work-related or non-work-related), the accident circumstances, the geographical coordinates and road type, and the number of deaths and injuries (Appendix F).

This work has omitted some variables such as those related to road sections, alignment and vehicular traffic on the grounds, that they should be considered in a more limited and homogeneous geographical area and over a longer period of investigation.

The availability of data and records concerning road traffic accidents is a well-established reality on the European scene. However, only recently these data have become more detailed thanks to the inclusion of information about road infrastructures, vehicles, trip purpose, time and circumstance of the crash.

---

4 The road accident is defined as “that event in which at least one vehicle is involved on the road network, occurring in the streets or squares open to traffic, which involves personal injuries (dead within 30 days and / or injured)” - (Convention of Vienna in 1968, UNECE, ITF and Eurostat 2019). For this reason, if the accident involves damage to objects only, it is then excluded from the statistics. This definition therefore reserves attention exclusively for reported accidents involving injury to people.

Lack of data was often due to difficulty in detecting accidents. Problems in accessing and using geolocation technologies in tunnels sometimes made empirical research impossible on account of unavailable or incomplete data.

The 2018 dataset was more detailed than in previous years since the information contained accident parameters and type of tunnel involved, thus enabling us to perform an analysis of the relationship between the most reliable variables.

This investigation was however limited to 716 records regarding 716 collisions involving 1,159 injured and 18 deaths that occurred in Italy in 2018 on 490 km of the approximately 2,600 km of Italian road tunnels.

For the purpose of this study, the Istat accident variables and label modes used for analysis were:

- Road tunnel location: motorway or road outside and inside urban areas, rural or road outside urban areas and not motorway, urban or road inside urban areas and not motorway (Legislative Decree/*Decreto Legislativo* 30 April 1992 n. 285);
- Vehicle type (Appendix A): Cars, Heavy goods vehicle (Heavy V), Motorcycles, Other vehicles (Other V);
- Time of occurrence (Appendix B): Night or Day defined according to a conventional interval;
- Trip purpose (Appendix C): work-related (trip/journey purpose in route to/from work or driving as part of the work) and non-work-related (trip/journey purpose non-work-related). Accidents work-related are underestimated due to the difficulty during the intervention by the police at the site of the accident, to record this information;
- Accident circumstances (Appendices D and E): Not keeping distances between vehicles (Distances), Distraction, Normal driving, Speeding, Unspecified and Other circumstances (Other C), corresponding to driver behaviour recorded when the accident occurred (Amundsen *et al.*, 2000; Gariazzo *et al.*, 2018). These classes, listed in Appendix D, are based on Istat and European coding (European Commission CARE 2016). The results of geoprocessing and query provided frequency of vehicle type, time of occurrence, trip purpose of driver, accident circumstances by tunnel length and road type.

## 4. Methods

This study used descriptive analysis combined with a multivariate approach (Bolasco, 1999; Di Franco, 2017; Johnson *et al.*, 2002; Jolliffe, 2002) based on the Principal Component Analysis (PCA). By means of geoprocessing operations (Cima *et al.*, 2014; Costabile *et al.*, 2012) the subset of Istat data was integrated with the length of the tunnel involved in the collision. The 1.3 version of the *Rstudio* and the 3.18.3 version of the *Quantum Gis* (QGIS) geographic information system were used.

For the purpose of descriptive analysis, for each variable and tunnel type, the frequency has been determined (see Appendix F). Tunnel sections were divided into 8 classes: up to 500 m (0-500); from 500 to 1,000 m (500-1000); from 1,000 to 1,500 m (1000-1500); from 1,500 to 2,000 m (1500-2000); from 2,000 to 2,500 m (2000-2500); from 2,500 to 3,000 m (2500-3000); from 3,000 to 3,500 m (3000-3500); over 3,500 m (>3500), corresponding to a range of lengths expressed in metres (Figure 1). The corresponding statistical parameters were then determined.

Based on the results of descriptive analysis, tunnels were grouped into two length classes more suitable for PCA: up to 500 metres and over 500 metres (>500) (Directive 2004/54/EC). These were then combined with tunnel location on motorway, rural and urban roads (Legislative Decree/*Decreto Legislativo* 30 April 1992 n. 285) so that indirectly different traffic conditions were also taken into consideration. Therefore, descriptive analysis was based on 24 types of tunnels, while PCA was based on 6 types (Appendix G). The two combined approaches were used to analyse the association between the accident parameters and the infrastructure in terms of length and road type.

The PCA model reproduces the information contained in the original variables, which turned out to be collinear, by concentrating all the information within the latent variables or Principal Components (PCs). The PCA enabled us to reduce redundancy (Benzécri, 1980) and provide a new representation of the coordinates of the original standardised variables that express the correlation. The combination of descriptive and multivariate analysis then revealed the association between the accident parameters and tunnel type. The frequency of vehicle type, time of accident occurrence, trip purpose and accident circumstances were then associated with the length and road type of the tunnel involved.

## 5. Results

### 5.1 Descriptive analysis

For each variable and tunnel type, defined by length and road type, frequencies (see Appendix F) and statistical parameters have been determined.

Table 1 summarises frequencies, minimum, maximum and quartiles for each variables and suggests using class aggregation more suitable for PCA.

**Table 1 - Descriptive analysis of injuries by variables. Italy, 2018**

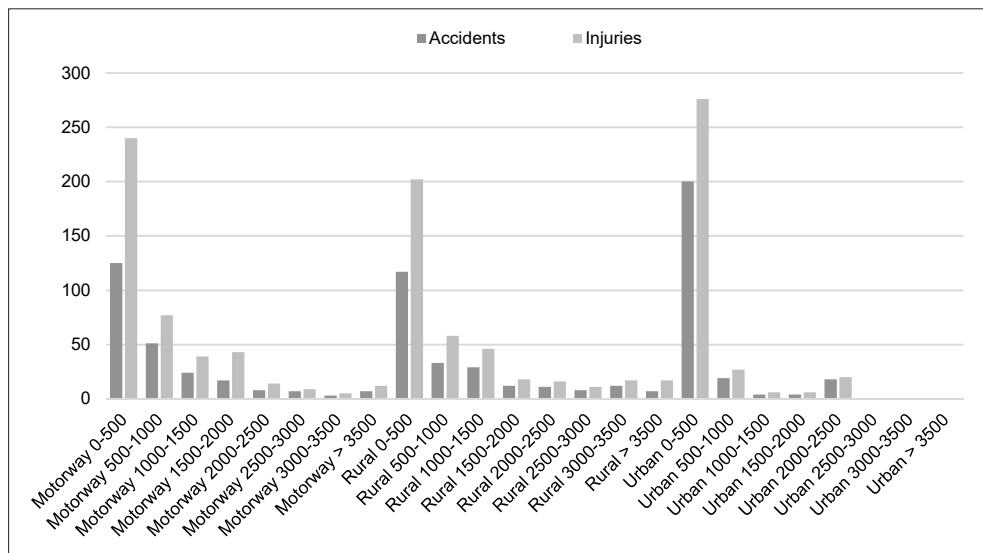
Variables	Freq.	Min.	1 <sup>st</sup> Qu.	Median	Mean	3 <sup>rd</sup> Qu.	Max.
Cars	876	-	6.5	15.0	36.5	35.7	200
Heavy goods vehicles (Heavy V)	111	-	-	2.0	5.0	4.0	27
Motorcycles	104	-	-	1.0	4.0	3.0	49
Other vehicles (Other V)	68	-	-	-	3.0	0.5	40
Night	159	-	0.7	3.0	7.0	6.0	35
Day	1000	-	6.5	12.5	41.7	40.5	248
Work-related trip	77	-	-	-	3.0	2.0	26
Non-work-related trip	1082	-	7.5	16.0	45.1	43.0	249
Not keeping distances from other vehicles (Distances)	198	-	-	2.5	8.2	6.5	62
Distraction	153	-	-	3.0	6.0	6.0	43
Other circumstances (Other C)	213	-	1.7	4.0	8.9	10.5	56
Normal driving	309	-	0.7	3.5	12.9	12.5	79
Speeding	183	-	-	3.0	8.0	8.0	47
Unspecified circumstances (Unspecified)	103	-	-	-	4.0	2.0	32

Source: Authors' processing on Istat data - Survey on road accidents resulting in death or injury

Figure 2 illustrates absolute frequencies of accidents and injuries for the 24 tunnel types, with peaks for motorway, rural, urban tunnels of up to 500 m (the most numerous in Italy) and motorways 1,500-2,000 m. By normalising the absolute data with the total length of tunnels in each class, accidents were further characterised in relation to length and road type.

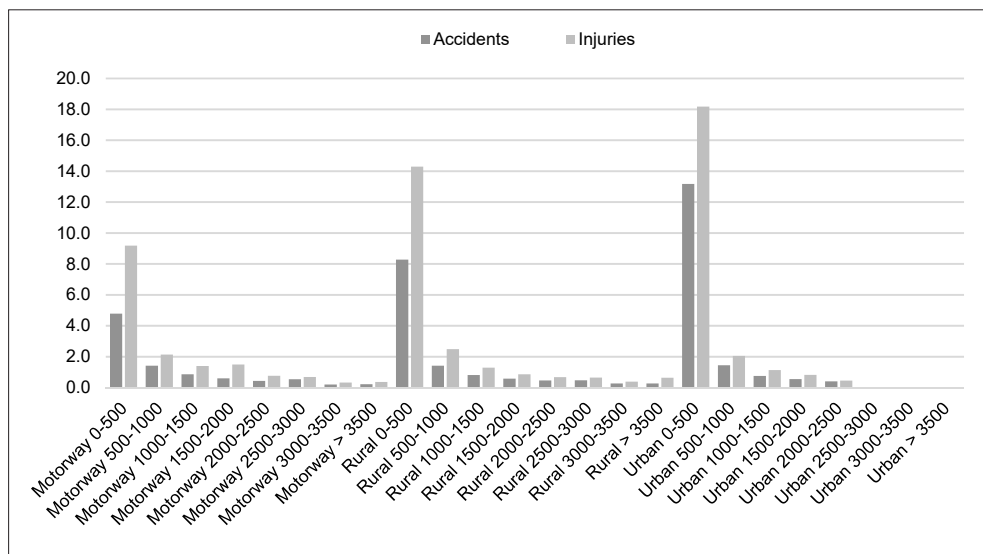
Figure 3 illustrates the ratio of the number of accidents to the total length of tunnel involved for each class, where accidents and injuries decrease with length. Further studies based on a wider range of data are needed to investigate the influence of traffic.

**Figure 2 - Injuries and accidents by tunnel type. Absolute frequencies. Italy, 2018**



Source: Authors' processing on Istat data - Survey on road accidents resulting in death or injury, and OSM

**Figure 3 - Injuries and accidents per kilometre, by tunnel type. Relative frequencies. Italy, 2018**



Source: Authors' processing on Istat data - Survey on road accidents resulting in death or injury, and OSM

## 5.2 Principal Component Analysis

A number of conditions were checked before conducting PCA (the verification of the linear relationship, the significance of the correlation between all quantitative variables, the absence of outliers and the number of observations). The Bartlett test results suggested proceeding with PCA. Once the conditions for applying PCA had been verified, analysis was performed on data matrix  $D=X_1, \dots, X_p$  (see Appendix G with Matrix D). The components were extracted from the correlation matrix and analysis of variables and cases as well as the interpretation of factors were performed on the basis of variable-component or tunnel type-component association.

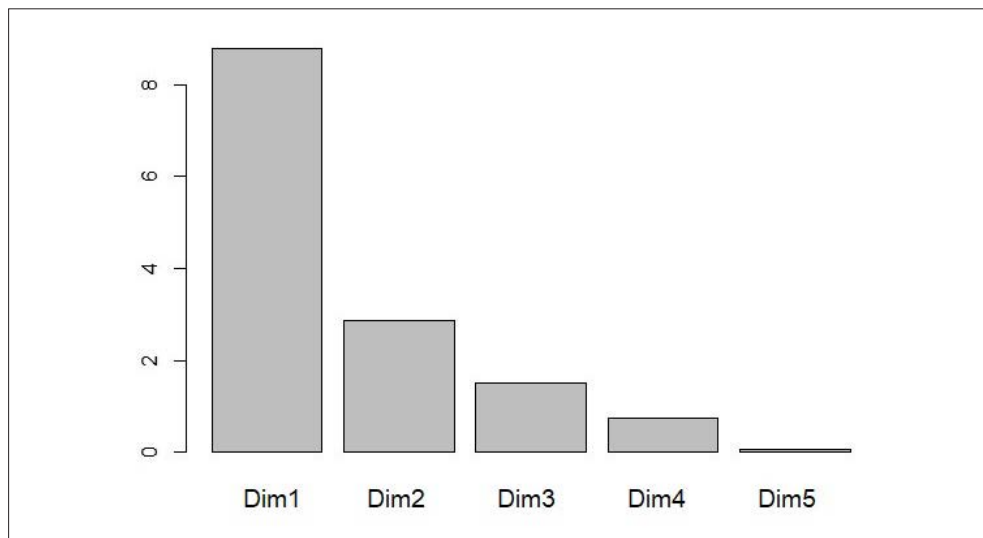
**Table 2 - Eigenvalues and explained variance of the extracted components. Year 2018**

Component	Eigenvalue	Percentage of variance	Cumulative percentage of variance
Dim 1	8.78	62.71	62.71
Dim 2	2.88	20.56	83.28
Dim 3	1.52	10.79	94.07
Dim 4	0.75	5.35	99.42
Dim 5	0.08	0.58	100.00

Source: Authors' processing on Istat data - Survey on road accidents resulting in death or injury, and OSM

The first results obtained were eigenvalues and explained variance (Table 2), which were useful for determining the number of Principal Components to use in the study. The PCs were identified by means of three criteria corresponding to eigenvalues  $>1$ , total explained variance reaching 70-90% (Table 2), and corresponding to values on the left of the inflection point on the bar plot (Figure 4).



**Figure 4 - Bar plot of eigenvalues and components. Year 2018**

Source: Authors' processing on Istat data - Survey on road accidents resulting in death or injury

The application of these three criteria yielded three PCs (Dim1, Dim2, Dim3) that summed up 94% (63%, 20% and 11%) of the total variance without a significant loss of information. These were explained by a general formula (1) as a combination of the eigenvalues (correlation matrix coefficients) and eigenvectors:

$$(1) \quad \text{Dim}_i = l_{i1}X_1 + l_{i2}X_2 + \dots + l_{ip}X_p \quad (i=1,2, \dots, p).$$

The next results were the factorial coordinates or correlations, cosines squared and contributions of variables (Table 3) and the coordinates or correlations, cosines squared and contributions of tunnel type (Table 4) and PCs (Dim1, Dim2, Dim3). Table 3 illustrates the coordinates and degree of correlation between variables and components. The output of *RStudio* in the case of standardised original variables implies that the coordinates coincide with the correlation and therefore the importance of each variable in relation to a factor can be deduced from the coordinates or correlation: the higher the coordinates of the variable, the more the latter affects the construction of the axis. Standardisation makes it possible to obtain a simplified representation of the phenomenon together with the correlation between variables in the system.

**Table 3 - Coordinates (Correlation), contributions and cosines squared by variables and PCs. Italy, 2018**

Variables	Coordinate (Correlation)			Cos <sup>2</sup>			Contributions		
	Dim1	Dim2	Dim3	Dim1	Dim2	Dim3	Dim1	Dim2	Dim3
Cars	0.9	-0.4	0.2	0.8	0.1	0.0	9.2	5.2	3.0
Heavy goods vehicles (Heavy V)	0.8	-0.5	-0.2	0.7	0.2	0.0	8.1	7.1	1.6
Motorcycles	0.5	0.7	-0.4	0.2	0.5	0.2	2.6	18.7	12.2
Other vehicles (Other V)	0.7	0.5	-0.3	0.6	0.3	0.1	6.4	9.8	6.3
Night	0.8	-0.2	0.4	0.7	0.1	0.2	7.8	1.9	13.3
Day	1.0	0.0	-0.1	1.0	0.0	0.0	11.3	0.0	0.2
Work-related trip	0.4	0.7	0.6	0.2	0.4	0.3	2.2	15.5	20.8
Non- work-related trip	1.0	-0.2	-0.1	1.0	0.0	0.0	10.8	1.1	0.8
Not keeping distances from other vehicles (Distances)	0.7	-0.4	-0.5	0.5	0.1	0.3	6.0	4.8	17.1
Distraction	0.9	0.3	0.0	0.8	0.1	0.0	9.6	4.0	0.1
Other circumstances (OtherC)	0.9	0.2	-0.1	0.8	0.0	0.0	8.6	1.5	1.1
Normal driving	0.8	-0.1	0.6	0.6	0.0	0.3	7.3	0.2	22.5
Speeding	0.8	-0.5	-0.1	0.7	0.3	0.0	7.9	9.8	0.9
Unspecified circumstances (Unspecified)	0.4	0.8	0.0	0.2	0.6	0.0	2.2	20.4	0.0

Source: Authors' processing on Istat data - Survey on road accidents resulting in death or injury

The coordinates, cosines squared and contributions referring to the type of tunnel were also calculated. The absolute contribution expresses the extent to which each variable explains the component. The cos<sup>2</sup> or relative contribution indicates the extent to which each component explains a variable (inertia or variance).

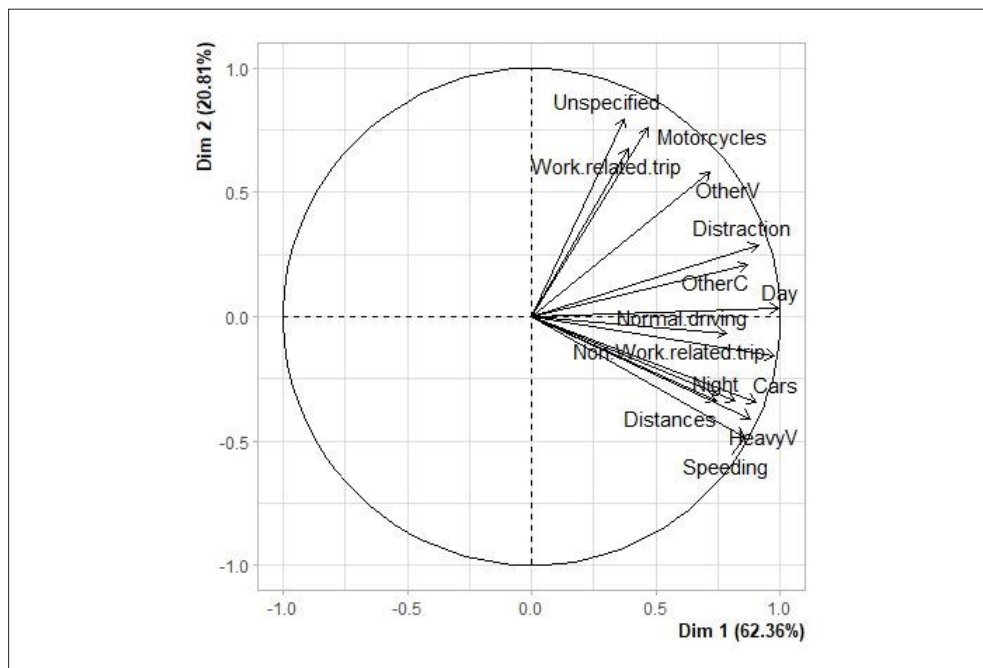
**Table 4 - Coordinates (correlation), cosines squared and contributions by tunnel type and PCs. Italy, 2018**

Tunnel type	Coordinate (Correlation)			Cos <sup>2</sup>			Contributions		
	Dim1	Dim2	Dim3	Dim1	Dim2	Dim3	Dim1	Dim2	Dim3
Motorway 0-500	1.82	-2.59	-0.3	0.29	0.59	0.0	6.29	38.92	1.2
Motorway >500	-0.09	-1.68	-1.0	0.00	0.55	0.2	0.02	16.35	10.0
Rural 0-500	0.53	0.44	2.3	0.05	0.03	0.9	0.54	1.13	56.6
Rural >500	-0.27	0.34	1.0	0.03	0.04	0.3	0.14	0.68	10.1
Urban 0-500	3.85	2.56	-1.2	0.65	0.29	0.1	28.15	37.92	16.3
Urban >500	-5.85	0.93	-0.7	0.95	0.02	0.0	64.87	5.00	5.8

Source: Authors' processing on Istat data - Survey on road accidents resulting in death or injury, and OSM

Further results of the two combined analytical approaches are summarised in Figures 5 and 6 that show the correlation between each variable and tunnel type and the Principal Components Dim1 and Dim2.

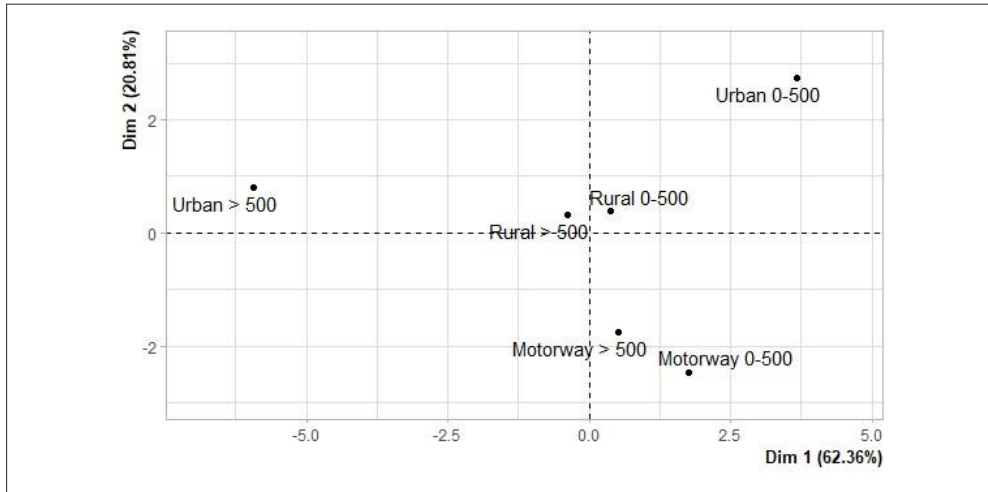
**Figure 5 - Correlation between each variable analysed and the first two PCs (Dim1 and Dim2). Year 2018**



Source: Authors' processing on Istat data - Survey on road accidents resulting in death or injury

On the basis of the results shown in Table 3 (contributions), the labels to be given to the three components were: “cars, heavy V, daytime driving, non-work-related trip, distances, distraction and speeding” (Dim1); “motorcycles, work-related trip, unspecified circumstances” (Dim2); “night-time driving, normal driving” (Dim3).

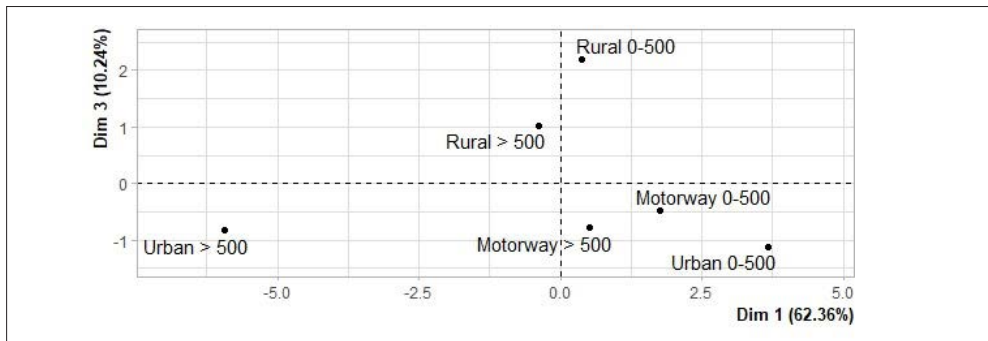
**Figure 6 - Correlation between tunnel type and the first two PCs (Dim1 and Dim2). Year 2018**



Source: Authors' processing on Istat data - Survey on road accidents resulting in death or injury, and OSM

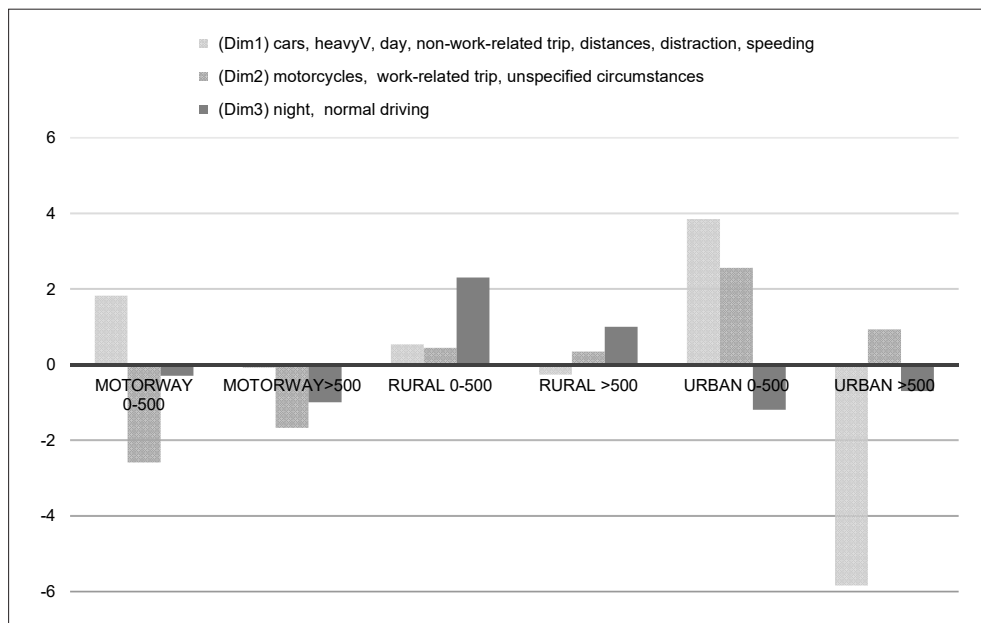
Dim1 (cars, heavy vehicles, daytime driving, non-work-related trip, distances, distraction and speeding) is seen to be positively associated with shorter urban and motorway tunnels and negatively with longer urban tunnels. Dim2 (motorcycles, work-related trip, unspecified circumstances) is negatively correlated with motorway tunnels and positively with urban ones (Figure 6). Dim 3 (night-time driving, normal driving) is seen to be positively associated with rural tunnels and shows a slight negative correlation with urban and motorway tunnels (Figure 7).

**Figure 7 - Correlation between tunnel type and Dim1 and Dim3. Year 2018**



Source: Authors' processing on Istat data - Survey on road accidents resulting in death or injury, and OSM

**Figure 8 - Positive and negative correlation between the first three PCs and the tunnel type. 2018**



Source: Authors' processing on Istat data - Survey on road accidents resulting in death or injury, and OSM

Figure 8 illustrates the correlation (ordinate) between the PCs (Dim1, Dim2 and Dim3) and the type of tunnel. Dim1 is seen to be positively associated with shorter motorway rural and urban tunnels and negatively with longer ones. Dim2 is negatively correlated with motorway tunnels and positively with urban ones. Dim3 is associated positively with rural tunnels and shows a slight negative correlation with urban and motorway tunnel (see Appendix H).

PCA highlighted some aspects that descriptive analysis alone failed to detect with regard to vehicles, the time of accident occurrence, the type of trip (work-related or non-work-related), the circumstance and the positive and negative correlations with the different types of tunnel.

A high frequency of injuries is found for the work-related trip variable in rural tunnels but has a poor correlation with the first two PCs and a slight correlation with the third component (Figures 7 and 8), which is, however, the least important. On the contrary, in shorter urban tunnels, the same variable has

higher frequencies and correlations with both Dim1 and Dim2, thus providing a more reliable variable - tunnel type association. This type of information is the value added that PCA gives to the descriptive analysis and thanks to which it was possible to ascertain more objectively the interrelationship between the descriptive parameters of the 716 accidents recorded in the Istat database and the road infrastructure involved.

## 6. Discussion

Road tunnel accident rates are influenced by many factors such as tunnel length, traffic, horizontal alignment, lane width, tunnel cross section, quality of tunnel lighting, driving speed and, last but not least, national driving habits and the technical standard of the vehicles used (PIARC, 2008). As can be seen by the amount of research that has analysed the design and safety of tunnels according to differing perspectives and objectives, issues regarding these infrastructures constitute a really vast topic (Calvi *et al.*, 2013) that involves the way drivers approach, exit from and behave inside a road tunnel (Mühlberger *et al.*, 2015).

Study findings are not consistent as regards the impact of tunnel length on safety (Bassan, 2016) or on the quality of traffic. Several factors can influence in a positive or negative way the probability or the effect of a tunnel accident. All these influencing factors make it difficult to compare collision rates at a statistical level. There is still a lack of knowledge, indeed, regarding the interaction of the various factors that influence crash rates in road tunnels. Several tunnel conditions may result in misjudgement of horizontal and vertical alignment and the perception of safe distances from other vehicles and obstacles (PIARC, 2008). Records of the number of road crashes and their victims represent essential information for road safety practitioners allowing them to analyse their spatial and temporal aspects. However, they cannot provide details on the factors causing road crashes (Hollò *et al.*, 2010). The availability of national road crash data together with information on accident road type enabled us to investigate road tunnel crashes involving death or injury throughout Italy and tunnel types nationwide. Although the study refers only to 2018 datasets, our statistical analysis (descriptive and PCA) revealed that the tunnel crash rate was associated both with tunnel length and with road type (Figures 2 and 3). Recent studies highlighted that crash rate, crash type, and contributing factors are variable in different zones of the tunnel (Amjad *et al.*, 2020).

The overall negative association with absolute frequency confirms findings obtained in Norway for road tunnel crashes, where short tunnels were associated with higher crash rates than long tunnels (Amundsen, 1994; Amundsen *et al.*, 2009). Yeung and Wong (Yeung *et al.*, 2013), who investigated expressway tunnels in Singapore, also found that Road Traffic

Crashes (RTCs) were lowest in the inner zones of tunnels and highest in the entry and exit zones. In contrast, compared with previous studies, another investigation of Chinese freeway tunnels revealed a sharp decline in the crash rate at tunnel portals and in the first 100 m of tunnels (tunnel entrance), as well as an increased number of crashes inside the tunnels. This finding was explained by better lighting in the entrance zone to the tunnels (Ma *et al.*, 2009). Caliendo showed by a negative binomial regression model for non-serious and serious accidents that crash frequency on unidirectional Italian motorway tunnel sections increases with tunnel length, in addition to other factors (Caliendo *et al.*, 2019). Further studies based on a wider range of data are needed to investigate accidents in motorway tunnels.

Our study puts forward two sets of hypothetical explanations of road tunnel accidents: the first concern all lengths linked to tunnel entrance and the impact on the driver of its confined environment that leads to reactions that increase accident risk and persist in relation to the time needed to reduce the driver's discomfort. This period of adaptation appears to be lacking in shorter tunnels and could therefore explain the higher accident rate in the latter. On the other hand, the second set of causes concern longer motorway and rural tunnels where, what counts is the gradual adaptation of the driver as he proceeds and which could be associated with the increase in speed in the more inner parts of the tunnel.

The study shows that the majority of accidents occur in tunnels of up to 500 metres in length (Figure 2). The number of accidents per kilometre, in motorway, rural and urban tunnels of up to 500 metres in length was 4.8, 8.3 and 13.2 while the same index, in longer motorway, rural and urban tunnel was 0.7, 0.6 and 0.6 (Figure 3).

The average ratio of injured to crashes was 1.9 in motorway tunnels, 1.7 in rural tunnels and 1.4 in urban tunnels of up to 500 m, while the same ratio in longer tunnels was 1.7, 1.6 and 1.3). As regard road type, the average ratio of injured to crashes was 1.8 in motorway tunnels, 1.7 in rural tunnels and 1.5 in urban ones. The higher motorway ratios were confirmed by Caliendo (*et al.*) who reported a larger number of serious accidents in motorway tunnels. The same authors found that if an accident occurs in these tunnels, the severity of injuries sustained is significantly higher than on open stretches of motorways (Caliendo *et al.*, 2012; Caliendo *et al.*, 2019).



Thanks to the PCA it was possible to ascertain more objectively the interrelationship between the descriptive parameters of the 716 accidents recorded in the Istat database and the road infrastructure involved. As far as the Dim1 is concerned: cars, heavy V, daytime driving, non-work-related trips, not respecting safe distances from other vehicles (rear-end collisions), distraction and speeding had the strongest positive association with short urban and motorway tunnels, while had the strongest negative association with long urban.

Dim2 (motorcycles, work-related drivers, unspecified circumstances) was found to be positively associated with urban and rural tunnels, while a strong negative association was observed with motorway tunnels. According to recent studies conducted over the last decade, due to traffic congestion in metropolitan areas, there has been an increase in the use of motorcycles as a mode of transport. This has led to a dramatic rise in the number of crashes involving these vehicles and a concomitant increase in the number of motorcyclists who have died or been injured (Gariazzo *et al.*, 2021). This positive association, which will need to be further studied, could be attributed to the growing number of delivery riders that operate in urban areas and drive motorcycles. With regard the Dim3 (night-time driving and normal driving) was found to be positively associated with rural tunnels.

Our findings enabled us to explain by means of an integrated interpretation of the over 700 road tunnel accidents, 1,159 injuries and 18 deaths recorded in Italy in 2018 in 491 km of tunnels, the interaction of the various parameters included in the model, influencing crash rates. For reasons of data robustness and completeness, this study focusses on a single year of observation (2018) with its number of cases collected. Therefore, the results achieved represent as much detail as possible. Variables such as traffic (marginally highlighted by location on motorway, rural and urban roads), horizontal alignment, lane width, tunnel cross section, quality of tunnel lighting and different driving habits throughout Italy were also omitted, together with the safety devices and technical standard of the vehicles in use and individual driver characteristics such as fatigue and physiological conditions.

Although these limitations need to be carefully evaluated, this study offers a methodological reference for all those who, with a reactive or proactive approach, are required to assess road safety, manage risks in a specific context

(Legislative Decree/*Decreto Legislativo* 9 April 2008 n. 81) and evaluate the effectiveness of prevention strategies (Directive 2004/54/EC). In addition, this study prepares the ground for future research that can be focussed on a more homogeneous and limited subset of tunnels. This taking into account also a larger number of variables derived from the exploitation and integration of more archives than the sources used in the current analysis.

Findings from our analysis confirm some of the results obtained by other authors and suggest that strategies designed to prevent tunnel accidents should be based on improving driver behaviour and on reducing the impact on drivers entering these infrastructures. In existing tunnels, better lighting in the areas near the entrance can mitigate impact due to a reduced road section, while a reduction in the speed limit in the inner areas may be useful for improving peripheral vision and the perception of the distances from other vehicles.

## 7. Conclusions

Our study analysed accidents resulting in death or injury in Italian road tunnels in 2018. It investigated the relationship between tunnel type (length-road type) and other variables involved in the collisions, such as vehicle type, time of occurrence, trip purpose and driving behaviour. The study shows that the majority of accidents occurs in tunnels of up to 500 metres in length. Our analysis revealed a positive association between cars, heavy vehicles, daytime driving, non-work-related drivers, not respecting safe distances between vehicles (rear-end collisions), distraction and speeding and short urban and motorway tunnels while a strongest negative association with long urban tunnels. Motorcycles, work-related drivers and unspecified circumstances were found to be positively associated with urban and rural tunnels, while a strong negative association was observed with motorway tunnels. The analysis revealed a positive association between night-time driving, normal driving and rural tunnels.

A strong positive association was found between urban tunnels of up to 500 m and accident circumstances such as careless driving, resulting in failure to observe a safe distance between vehicles and the speed limit in motorway tunnels.

Tunnel length was found to affect the frequency of crashes, while road type was found to impact the degree of severity of accident consequences. This difference based on road type provided a possible explanation which will definitely have to be studied in more detail in further research, also in relation to the at times conflicting results obtained by other studies.

Our analysis confirms some of the results obtained by other authors and proposes prevention strategies based on engineering measures designed to reduce the impact on the driver during access to tunnels or on adapting speed limits in inner tunnel areas. Despite its limits, this study provides a useful methodological reference since it highlights the influence of tunnel length on road accident variables and the relationship between various parameters that influence accident rates in road tunnels.

However, further and more thorough research based on a longer period of time will be needed to improve the interpretative capacity of this methodology.

## Appendix A - Injuries by tunnel type (road type and length) and vehicles involved. Absolute values. Italy, 2018

Tunnel type (a)	Cars	Heavy V	Motorcycles	Other V	Total
Motorway 0-500	200	27	2	11	240
Motorway 500-1,000	50	16	9	2	77
Motorway 1,000-1,500	32	4	1	2	39
Motorway 1,500-2,000	38	1	4	-	43
Motorway 2,000-2,500	14	-	-	-	14
Motorway 2,500-3,000	7	1	1	-	9
Motorway 3,000-3,500	3	2	-	-	5
Motorway > 3,500	8	2	2	-	12
Rural 0-500	167	15	10	10	202
Rural 500-1,000	42	7	9	-	58
Rural 1,000-1,500	35	5	3	3	46
Rural 1,500-2,000	16	2	-	-	18
Rural 2,000-2,500	15	-	1	-	16
Rural 2,500-3,000	9	2	-	-	11
Rural 3000-3500	13	3	1	-	17
Rural > 3,500	16	-	1	-	17
Urban 0-500	164	23	49	40	276
Urban 500-1,000	23	1	3	-	27
Urban 1,000-1,500	4	-	2	-	6
Urban 1,500-2,000	5	-	1	-	6
Urban 2,000-2,500	15	-	5	-	20
Urban 2,500-3,000	-	-	-	-	-
Urban 3,000-3,500	-	-	-	-	-
Urban >3,500	-	-	-	-	-
Total	876	111	104	68	1,159

Source: Authors' processing on Istat data - Survey on Road Accidents resulting in death or injury, and Open Street Map (a) Motorway (road outside and inside urban area), Rural (road outside urban area and not motorway), Urban (road inside urban area and not motorway).

## Appendix B - Injuries by tunnel type (road type and length) and time of occurrence. Absolute values. Italy, 2018

Tunnel type (a)	Night	Day	Total
Motorway 0-500	32	208	240
Motorway 500-1,000	9	68	77
Motorway 1,000-1,500	9	30	39
Motorway 1,500-2,000	4	39	43
Motorway 2,000-2,500	7	7	14
Motorway 2,500-3,000	-	9	9
Motorway 3,000-3,500	-	5	5
Motorway > 3,500	2	10	12
Rural 0-500	35	167	202
Rural 500-1,000	6	52	58
Rural 1,000-1,500	1	45	46
Rural 1,500-2,000	2	16	18
Rural 2,000-2,500	5	11	16
Rural 2,500-3,000	-	11	11
Rural 3000-3500	3	14	17
Rural > 3,500	6	11	17
Urban 0-500	28	248	276
Urban 500-1,000	2	25	27
Urban 1,000-1,500	1	5	6
Urban 1,500-2,000	3	3	6
Urban 2,000-2,500	-	-	-
Urban 2,500-3,000	4	16	20
Urban 3,000-3,500	-	-	-
Urban >3,500	-	-	-
<b>Total</b>	<b>159</b>	<b>1,000</b>	<b>1,159</b>

Source: Authors' processing on Istat data - Survey on Road Accidents resulting in death or injury, and Open Street Map (a) Motorway (road outside and inside urban area), Rural (road outside urban area and not motorway), Urban (road inside urban area and not motorway).

## Appendix C - Injuries by tunnel type (road type and length) and work-related or non-work-related trip. Absolute values. Italy, 2018

Tunnel type (a)	Work-related trip (b)	Non-work-related trip	Total
Motorway 0-500	-	240	240
Motorway 500-1,000	-	77	77
Motorway 1,000-1,500	-	43	43
Motorway 1,500-2,000	-	43	43
Motorway 2,000-2,500	-	14	14
Motorway 2,500-3,000	-	9	9
Motorway 3,000-3,500	-	5	5
Motorway > 3,500	-	10	10
Rural 0-500	26	176	202
Rural 500-1,000	10	48	58
Rural 1,000-1,500	7	39	46
Rural 1,500-2,000	-	18	18
Rural 2,000-2,500	3	13	16
Rural 2,500-3,000	3	8	11
Rural 3000-3500	-	17	17
Rural > 3,500	2	15	17
Urban 0-500	25	249	274
Urban 500-1,000	1	26	27
Urban 1,000-1,500	-	6	6
Urban 1,500-2,000	-	6	6
Urban 2,000-2,500	-	20	20
Urban 2,500-3,000	-	-	-
Urban 3,000-3,500	-	-	-
Urban >3,500	-	-	-
<b>Total</b>	<b>77</b>	<b>1,082</b>	<b>1,159</b>

Source: Authors' processing on Istat data - Survey on Road Accidents resulting in death or injury, and Open Street Map (a) Motorway (road outside and inside urban area), Rural (road outside urban area and not motorway), Urban (road inside urban area and not motorway).

(b) The data listed are "work-related" (trip/journey purpose in route to/from work or driving as part of the work) and "non-work-related" (trip/journey purpose not professional).

## Appendix D - Classification of accident circumstances

Circumstances (a) (codes by accident type and labels)	Circumstances (group labels)
22 Driving without maintaining a safe distance between vehicles	Distances
62 Driving without maintaining a safe distance between vehicles	
21 Driving in a careless or indecisive manner	Distraction
61 Driving in a careless or indecisive manner	
25 Straying from the right of the carriageway	Other C
26 Driving in the wrong direction	
31 Overtaking incorrectly on the right, on a bend, on a hump, with poor visibility	
33 Overtaking incorrectly on the right despite 'No Overtaking' sign	
37 Driving normally to stop or park	
36 Turning left	
35 Merging into traffic lane	
34 Reversing or U-turn manoeuvres	
45 Manoeuvring	
48 Driving off the carriageway and running down a pedestrian	
49 Failure to stop at pedestrian crossings	
51 Injuring the pedestrian with the load	
66 Proceeding despite 'No Transit' or 'No Access' signs	
70 Skidding and going off-road to avoid crashing	
71 Skidding and going off-road due to careless driving	
72 Skidding and going off-road due to excessive speeding	
73 Braking suddenly with consequences for passengers	
74 Passengers falling from vehicle when opening car door	
75 Passengers falling when alighting from vehicle	
76 Passengers falling from vehicle for not wearing seat belts	
20 Driving normally	Normal driving
40 Driving normally	
60 Driving normally	
23 Driving too fast	Speeding
24 Speeding	
41 Driving too fast	
64 Driving too fast	
65 Speeding	
00 Unspecified circumstance	Unspecified
Null	

Source: Authors' processing on Istat data - Survey on Road Accidents resulting in death or injury

(a) The same label and different code depend on the different accident type (e.g. head-on collision or vehicles travelling in the same direction).

## Appendix E - Injuries by tunnel type (road type and length) and group of accident circumstances. Absolute values. Italy, 2018

Tunnel type (a)	Distances	Distraction	Other C	Normal driving	Speeding	Unspecified	Total
Motorway 0-500	62	22	31	63	47	15	240
Motorway 500-1,000	15	9	14	12	25	2	77
Motorway 1,000-1,500	13	5	7	10	3	1	39
Motorway 1,500-2,000	4	7	10	14	8	-	43
Motorway 2,000-2,500	3	6	3	2	-	-	14
Motorway 2,500-3,000	-	2	2	5	-	-	9
Motorway 3,000-3,500	1	-	4	-	-	-	5
Motorway > 3,500	6	-	4	-	-	2	12
Rural 0-500	17	27	31	79	27	21	202
Rural 500-1,000	2	9	12	25	3	7	58
Rural 1,000-1,500	6	5	12	17	6	-	46
Rural 1,500-2,000	4	3	2	3	6	-	18
Rural 2,000-2,500	3	3	6	1	3	-	16
Rural 2,500-3,000	1	5	1	2	-	2	11
Rural 3000-3500	-	3	3	5	3	3	17
Rural > 3,500	-	-	6	1	8	2	17
Urban 0-500	51	43	56	59	35	32	276
Urban 500-1,000	8	4	2	4	9	-	27
Urban 1,000-1,500	-	-	6	-	-	-	6
Urban 1,500-2,000	1	-	-	5	-	-	6
Urban 2,000-2,500	1	-	1	2	-	16	20
Urban 2,500-3,000	-	-	-	-	-	-	-
Urban 3,000-3,500	-	-	-	-	-	-	-
Urban >3,500	-	-	-	-	-	-	-
<b>Total</b>	<b>198</b>	<b>153</b>	<b>213</b>	<b>309</b>	<b>183</b>	<b>103</b>	<b>1,159</b>

Source: Authors' processing on Istat data - Survey on Road Accidents resulting in death or injury, and Open Street Map (a) Motorway (road outside and inside urban area), Rural (road outside urban area and not motorway), Urban (road inside urban area and not motorway).



## Appendix F - Injuries by tunnel type (road type and length) and vehicles involved, time of occurrence, work-related or non-work-related trip, group of accident circumstances. Absolute values. Italy, 2018

Tunnel type (a)	Car	Heavy V	Motor-Cycles	Other V	Night	Day	Work-related trip	Non-work-related trip	Distances (b)	Distraction	Other C	Normal driving	Speeding	Unspecified
Motorway 0-500	200	27	2	11	32	208	-	240	62	22	31	63	47	15
Motorway 500-1,000	50	16	9	2	9	68	-	77	15	9	14	12	25	2
Motorway 1,000-1,500	32	4	1	2	9	30	-	43	13	5	7	10	3	1
Motorway 1,500-2,000	38	1	4	-	4	39	-	43	4	7	10	14	8	-
Motorway 2,000-2,500	14	-	-	-	7	7	-	14	3	6	3	2	-	-
Motorway 2,500-3,000	7	1	1	-	-	9	-	9	-	2	2	5	-	-
Motorway 3,000-3,500	3	2	-	-	-	5	-	5	1	-	4	-	-	-
Motorway > 3,500	8	2	2	-	2	10	-	10	6	-	4	-	-	2
Rural 0-500	167	15	10	10	35	167	26	176	17	27	31	79	27	21
Rural 500-1,000	42	7	9	-	6	52	10	48	2	9	12	25	3	7
Rural 1,000-1,500	35	5	3	3	1	45	7	39	6	5	12	17	6	-
Rural 1,500-2,000	16	2	-	-	2	16	-	18	4	3	2	3	6	-
Rural 2,000-2,500	15	-	1	-	5	11	3	13	3	3	6	1	3	-
Rural 2,500-3,000	9	2	-	-	-	11	3	8	1	5	1	2	-	2
Rural 3000-3500	13	3	1	-	3	14	-	17	-	3	3	5	3	3
Rural > 3,500	16	-	1	-	6	11	2	15	-	-	6	1	8	2
Urban 0-500	164	23	49	40	28	248	25	249	51	43	56	59	35	32
Urban 500-1,000	23	1	3	-	2	25	1	26	8	4	2	4	9	-
Urban 1,000-1,500	4	-	2	-	1	5	-	6	-	-	6	-	-	-
Urban 1,500-2,000	5	-	1	-	3	3	-	6	1	-	-	5	-	-
Urban 2,000-2,500	15	-	5	-	-	-	-	20	1	-	1	2	-	16
Urban 2,500-3,000	-	-	-	-	4	16	-	-	-	-	-	-	-	-
Urban 3,000-3,500	-	-	-	-	-	-	-	-	-	-	-	-	-	-
Urban >3,500	-	-	-	-	-	-	-	-	-	-	-	-	-	-

Source: Authors' processing on Istat data - Survey on Road Accidents resulting in death or injury, and Open Street Map (a) Motorway (road outside and inside urban area), Rural (road outside urban area and not motorway), Urban (road inside urban area and not motorway).

(b) Not keeping distances from other vehicles.

## Appendix G. Matrix D - Injuries by tunnel type (road type and length) and vehicles involved, time of occurrence, work-related or non-work-related trip, group of accident circumstances. Absolute values. Italy, 2018

Tunnel type (a)	Car	Heavy V	Motor-Cycles	Other V	Night	Day	Work-related trip	Non-work-related trip	Distances (b)	Distraction	Other C	Normal driving	Speeding	Unspecified
Motorway 0-500	200	27	2	11	32	208	-	240	62	22	31	63	47	15
Motorway > 500	152	26	17	4	31	168	-	201	42	29	44	43	36	5
Rural 0-500	167	15	10	10	35	167	26	176	17	27	31	79	27	21
Rural > 500	146	19	15	3	23	160	25	158	16	28	42	54	29	14
Urban 0-500	164	23	49	40	28	248	25	249	51	43	56	59	35	32
Urban > 500	47	1	11	-	10	49	1	58	10	4	9	11	9	16

Source: Authors' processing on Istat data - Survey on Road Accidents resulting in death or injury, and Open Street Map (a) Motorway (road outside and inside urban area), Rural (road outside urban area and not motorway), Urban (road inside urban area and not motorway).

(b) Not keeping distances from other vehicles.

## Appendix H - Accidents and injuries by tunnel type (road type and length). Italy, 2018

Tunnel type (a)	Number of Accidents	Number of Injuries	Tunnel Extension (in m)	Accidents per km	Injuries per km
Motorway 0-500	125	240	26,145.6	4.8	9.2
Motorway >500	117	199	173,294.7	0.7	1.1
Rural 0-500	117	202	14,136.1	8.3	14.3
Rural > 500	112	183	192,714.5	0.6	0.9
Urban 0-500	200	276	15,186.5	13.2	18.2
Urban > 500	45	59	70,463.0	0.6	0.8

Source: Authors' processing on Istat data - Survey on Road Accidents resulting in death or injury, and Open Street Map (a) Motorway (road outside and inside urban area), Rural (road outside urban area and not motorway), Urban (road inside urban area and not motorway).

## References

Amjad, P., H. Helai, H. Chunyang, W. Jie, L. Ye. 2020. “Revisiting freeway single tunnel crash characteristics analysis: A six-zone analytic approach”. *Accident Analysis & Prevention*, Volume 142. <https://doi.org/10.1016/j.aap.2020.105542>.

Amundsen, F.H. 1994. “Studies of driver behaviour in Norwegian road tunnels”. *Tunnelling and Underground Space Technology*, Volume 9, Issue 1: 9-15.

Amundsen, F.H., and A. Engebretsen. 2009. “Studies on Norwegian Road Tunnels II. An Analysis on Traffic Accidents in Road Tunnels 2001—2006”. *Rapport N. TS4-2009*, Roads and Traffic Department, Traffic Safety Section. Oslo, Norway: Statens vegvesen Vegdirektoratet.

Amundsen, F.H., and G. Ranes. 2000. “Studies on traffic accidents in Norwegian road tunnels”. *Tunnelling and Underground Space Technology*, Volume 15, Issue 1: 3–11.

Bassan, S. 2016. “Overview of traffic safety aspects and design in road tunnels”. *International Association of Traffic and Safety Sciences - IATSS Research*, Volume 40, Issue 1: 35-46.

Benzécri, J.-P. (ed.), F. Benzécri-Leroy, A. Birou, Laboratoire de statistique (Paris), et al. 1980. *L'analyse des données. II. L'analyse des correspondances: introduction, théorie, applications diverses notamment à l'analyse des questionnaires, programmes de calcul*. Paris, France: Dunod.

Bolasco, S. 1999. *Analisi multidimensionale dei dati. Metodi, strategie e criteri d'interpretazione*. Roma, Italia: Carocci.

Caliendo, C., and M.L. De Guglielmo. 2012. “Accident Rates in Road Tunnels and Social Cost Evaluation”. *Procedia - Social and Behavioral Sciences*, Volume 53: 166-177.

Caliendo, C., M.L. De Guglielmo, and I. Russo. 2019. “Analysis of crash frequency in motorway tunnels based on a correlated random-parameters approach”. *Tunnelling and Underground Space Technology*, Volume 85: 243-251.

Calvi, A., and F. D'Amico. 2013. "Study of the effects of road tunnel on driver behavior and road safety using driving simulator". *Advances in Transportation Studies*, Volume 30: 59-76.

Cima, V., M. Carroccio, e R. Maseroli. 2014. "Corretto utilizzo dei Sistemi Geodetici di Riferimento all'interno dei software GIS". In *Atti della 18a Conferenza Nazionale ASITA*, 14-16 ottobre 2014, Firenze: 359-363. Milano, Italia: Federazione delle Associazioni Scientifiche per le Informazioni Territoriali e Ambientali - ASITA

Costabile, S, S. Martini, L. Petriglia, G. Corrarello, e A. Avanzi. 2012. "Geoportale Nazionale. Un approfondimento sulle metodologie di conversione e trasformazione coordinate". *GEOMEDIA*, Volume 16, N. 6: 26-28.

Decreto Legislativo 30 aprile 1992, n. 285. "Testo aggiornato recante il nuovo codice della strada". *Gazzetta Ufficiale della Repubblica Italiana, Serie Generale*, n. 67 del 22 marzo 1994, *Supplemento Ordinario* n. 49.

Decreto Legislativo 9 aprile 2008, n. 81. "Attuazione dell'articolo 1 della legge 3 agosto 2007, n. 123, in materia di tutela della salute e della sicurezza nei luoghi di lavoro". *Gazzetta Ufficiale della Repubblica Italiana, Serie Generale*, n. 101 del 30 aprile 2008, *Supplemento Ordinario* n. 108.

Di Franco, G. 2017. *Tecniche e modelli di analisi multivariata*. Milano, Italia: Franco Angeli Editore.

European Commission, Community Database on Road Accidents - CARE. 2016. "Road Safety: new statistics call for fresh efforts to save lives on EU roads", *Press Release* 31 March 2016. Brussels, Belgium: European Commission.

European Parliament. 2004. *Directive 2004/54/EC of the European Parliament and of the Council of 29 April 2004 on minimum safety requirements for tunnels in the Trans-European Road Network*.

Gariazzo, C. (a cura di), A. Brusco, A. Bucciarelli, M. Bugani, C. Giliberti, A. Marinaccio, S. Massari, A. Pireddu, L. Veronico, G. Baldassarre, S. Bruzzone, M. Scortichini, M. Stafoggia, e S. Salerno. 2019. "Gli incidenti con mezzo di trasporto. Un'analisi integrata dei determinanti e dei fattori di rischio occupazionali". *Collana Ricerche*. Roma, Italia: Istituto Nazionale per l'Assicurazione contro gli Infortuni sul Lavoro - INAIL.

Gariazzo, C., S. Bruzzone, S. Finardi, M. Scortichini, L. Veronico, e A. Marinaccio. 2021. "Association between extreme ambient temperatures and general indistinct and work-related road crashes. A nationwide study in Italy". *Accident Analysis & Prevention*, Volume 155, 106110.

Gariazzo, C., M. Stafoggia, S. Bruzzone, A. Pelliccioni, and F. Forastiere. 2018. "Association between mobile phone traffic volume and road crash fatalities: A population-based case-crossover study". *Accident Analysis & Prevention*, Volume 115: 25-33.

Hollò, P., V. Eksler, and J. Zukowska. 2010. "Road safety performance indicators and their explanatory value: A critical view based on the experience of Central European countries". *Safety Science*, Volume 48, Issue 9:1142-1150.

Istituto Nazionale di Statistica - Istat. *I.Stat, Your direct access to the Italian Statistics: the complete data warehouse for experts*. Roma, Italy: Istat. <http://dati.istat.it/Index.aspx>.

Istituto Nazionale di Statistica - Istat, e Automobile Club d'Italia - ACI. "Incidenti stradali in Italia. Anno 2018". *Comunicato Stampa*. Roma, Italia: Istat. <https://www.istat.it/it/archivio/232366>.

Istituto Nazionale di Statistica - Istat. "Rilevazione degli incidenti stradali con lesioni a persone". *Informazioni sulla Rilevazione*. Roma, Italia: Istat. <https://www.istat.it/it/archivio/4609>.

Istituto Nazionale di Statistica - Istat. "Survey on Road accidents resulting in death or injury". In *SIQual - Information system on quality of statistical production processes*. Roma, Italy: Istat. <http://siqua.istat.it/SIQual/lang.do?language=UK>.

Johnson, R.A., and D.W. Wichern. 2002. *Applied Multivariate Statistical Analysis*. Hoboken, NJ, U.S.: Prentice Hall.

Jolliffe, I.T. 2002. *Principal Component Analysis, Second Edition*. Cham, Switzerland: Springer Nature, *Springer Series in Statistics*.

Ma, Z.L., C.F. Shao, S.R. Zhang. 2009. "Characteristics of traffic accidents in Chinese freeway tunnels". *Tunnelling and Underground Space Technology*, Volume 24: 350-355.

Mühlberger, A., M. Kinateder, J. Brütting, S. Eder, M. Müller, D. Gromer, and P. Pauli. 2015. “Influence of Information and Instructions on Human Behavior in Tunnel Accidents: A Virtual Reality Study”. *Journal of Virtual Reality and Broadcasting*, Volume 12, N. 3, 0009-6-42521.

PIARC, Comité technique 3.3 - Exploitation des tunnels routiers. 2008. *Human factors and road tunnel safety regarding users*. Paris, La Défense CEDEX, France: PIARC.

Pireddu, A., e S. Bruzzone. 2019. “Incidenti in gallerie stradali”. *Fact Sheet*. Roma, Italia: Istituto Nazionale per l’Assicurazione contro gli Infortuni sul Lavoro - INAIL.

Yeung, J.S., and Y.D. Wong. 2013. “Road traffic accidents in Singapore expressways tunnels”. *Tunnelling and Underground Space Technology*, Volume 38: 534-541.



The *Rivista di statistica ufficiale* publishes peer-reviewed articles dealing with cross-cutting topics: the measurement and understanding of social, demographic, economic, territorial and environmental subjects; the development of information systems and indicators for decision support; the methodological, technological and institutional issues related to the production process of statistical information, relevant to achieve official statistics purposes.

The *Rivista di statistica ufficiale* aims at promoting synergies and exchanges between and among researchers, stakeholders, policy-makers and other users who refer to official and public statistics at different levels, in order to improve data quality and enhance trust.

The *Rivista di statistica ufficiale* was born in 1992 as a series of monographs titled “*Quaderni di Ricerca Istat*”. In 1999 the series was entrusted to an external publisher, changed its name in “*Quaderni di Ricerca - Rivista di Statistica Ufficiale*” and started being published on a four-monthly basis. The current name was assumed from the Issue N. 1/2006, when the Italian National Institute of Statistics – Istat returned to be its publisher.

*La Rivista di statistica ufficiale pubblica articoli, valutati da esperti, che trattano argomenti trasversali: la misurazione e la comprensione di temi sociali, demografici, economici, territoriali e ambientali; lo sviluppo di sistemi informativi e di indicatori per il supporto alle decisioni; le questioni metodologiche, tecnologiche e istituzionali relative al processo di produzione dell'informazione statistica, rilevanti per raggiungere gli obiettivi della statistica ufficiale.*

*La Rivista di statistica ufficiale promuove sinergie e scambi tra ricercatori, stakeholder, policy-maker e altri utenti che fanno riferimento alla statistica ufficiale e pubblica a diversi livelli, al fine di migliorare la qualità dei dati e aumentare la fiducia.*

*La Rivista di statistica ufficiale nasce nel 1992 come serie di monografie dal titolo “Quaderni di Ricerca Istat”. Nel 1999 la collana viene affidata a un editore esterno, cambia nome in “Quaderni di Ricerca - Rivista di Statistica Ufficiale” e diventa quadrimestrale. Il nome attuale è stato scelto a partire dal numero 1/2006, quando l'Istituto Nazionale di Statistica - Istat è tornato a esserne l'editore.*