

Rossouw, Stephanie; Greyling, Talita

Working Paper

## Big Data and Happiness

GLO Discussion Paper, No. 634

**Provided in Cooperation with:**  
Global Labor Organization (GLO)

Suggested Citation: Rossouw, Stephanie; Greyling, Talita (2020) : Big Data and Happiness, GLO Discussion Paper, No. 634, Global Labor Organization (GLO), Essen

This Version is available at:  
<http://hdl.handle.net/10419/223012>

**Standard-Nutzungsbedingungen:**

Die Dokumente auf EconStor dürfen zu eigenen wissenschaftlichen Zwecken und zum Privatgebrauch gespeichert und kopiert werden.

Sie dürfen die Dokumente nicht für öffentliche oder kommerzielle Zwecke vervielfältigen, öffentlich ausstellen, öffentlich zugänglich machen, vertreiben oder anderweitig nutzen.

Sofern die Verfasser die Dokumente unter Open-Content-Lizenzen (insbesondere CC-Lizenzen) zur Verfügung gestellt haben sollten, gelten abweichend von diesen Nutzungsbedingungen die in der dort genannten Lizenz gewährten Nutzungsrechte.

**Terms of use:**

*Documents in EconStor may be saved and copied for your personal and scholarly purposes.*

*You are not to copy documents for public or commercial purposes, to exhibit the documents publicly, to make them publicly available on the internet, or to distribute or otherwise use the documents in public.*

*If the documents have been made available under an Open Content Licence (especially Creative Commons Licences), you may exercise further usage rights as specified in the indicated licence.*

# Big Data and Happiness

Stephanié Rossouw<sup>1</sup> and Talita Greyling<sup>2</sup>

## Abstract

The pursuit of happiness. What does that mean? Perhaps a more prominent question to ask is, *'how does one know whether people have succeeded in their pursuit'*? Survey data, thus far, has served us well in determining where people see themselves on their journey. However, in an everchanging world, one needs high-frequency data instead of data released with significant time-lags. High-frequency data, which stems from Big Data, allows policymakers access to virtually *real-time* information that can assist in effective decision-making to increase the quality of life for all. Additionally, Big Data collected from, for example, social media platforms give researchers unprecedented insight into human behaviour, allowing significant future predictive powers.

**Keywords:** Happiness; Big Data; Sentiment analysis

**JEL classification codes:** C88, I31, I39, J18

## 1. Introduction

*"Data! Data! Data!" he cried impatiently. 'I can't make bricks without clay.'"*

– *Sherlock Holmes*

Yes, even this fictional private detective believed in the necessity of data! Sherlock Holmes knew that you couldn't make robust conclusions without having the most updated and relevant data. Without such data, you cannot make swift and decisive decisions, the same way that you cannot build solid bricks without clay—Elementary, my dear Watson, indeed.

This paper discusses changes in the measurement of subjective well-being brought about by the use of Big Data in happiness research. The usual way to depict subjective well-being is by just asking respondents how they feel about their daily lives. Based on this basic idea, survey analysis is one of the most used methods to measure people's perceived well-being.

---

<sup>1</sup> Corresponding author: School of Economics, Auckland University of Technology, Email: [stephanie.rossouwe@aut.ac.nz](mailto:stephanie.rossouwe@aut.ac.nz)

<sup>2</sup> School of Economics, University of Johannesburg, Email: [talitag@uj.ac.za](mailto:talitag@uj.ac.za)

However, in an era of Big Data, with information readily available via social media, the possibilities of alternative measurements are evolving rapidly. For example, one can use sentiment and emotion analysis of social media data as an alternative method to measure how people feel or think about their lives. Twitter, Facebook, Instagram, Google trends and other platforms give us a vast data source. However, to what extent is that data also useful to describe and analyze happiness?

Throughout, the paper will deal with several topics and conclude with a survey on happiness indices constructed with the by-product of Big Data, namely organic data. The topics include, for example, lessons learned from earlier studies in the field of sentiment analysis in the domain of subjective well-being. Questions arise, such as whether it could be possible to use any lessons learned by applying them on a large scale to say something about the happiness of an entire country? Or whether the coverage is limited to specific target groups in society? And, of course, what other applications there are concerning Big Data as alternatives for the classical methods in the field of happiness and subjective well-being?

## **2. What is Big Data?**

In layman's terms, Big Data is a phrase used to describe a massive *volume* of both structured (for example stock information) and unstructured data (for example social media postings) generated through information and communication technologies (ITC) such as the Internet. The volume of Big Data is so large that it is difficult to process using traditional database and software techniques, and it is also unique in the *velocity* at which it is created and collected as well as the *variety* of the data points being covered (Sarangi & Sharma 2020). Since Big Data cannot be processed using traditional database and software techniques, it is vitally important that increased data storage capacity combined with faster and more efficient computing capabilities continues to be developed.

### ***2.1. Big data vs organic data***

Let's get this debate sorted out first. What exactly is the difference between Big Data and organic data? Well, the answer is quite simple. Big Data represents data which is collected for a specific purpose. For example, movements on stock markets are collected as high-frequency intraday data to predict future stock market movements. However, organic data represents data which is seen as a by-product from its original intent. Here, the social media platform Twitter is an excellent example. Every country in the world that has active Twitter users provides a 'live' stream of tweets. The purpose of a tweet is to express oneself among those you believe to be your peers. Researchers now use this 'live' stream of tweets to determine the happiness or the mood of a nation (see section 5). The original intent of tweets was never to be used for this kind of research, but this data in its by-product form is now a rich data resource informing social scientists on human behaviour.

## ***2.2. Benefits of organic data***

People live in a fluid world where circumstances can change almost instantaneously and, as such, there is a need for data that represent this new reality. Whereas survey data is comprehensive and mostly representative of society, there is usually a significant time-lag associated with the release of information gathered by this form of data. Additionally, survey data are extremely costly, which is one of the reasons why countries or organizations, such as Gallup Analytics, collect such data annually or cyclically. In census data (the most representative of all survey data), for example, the United States only takes stock once every ten years. In contrast, countries like New Zealand and South Africa take a census every five years. Information from census data is essential for future planning. Still, many things happen in five years, and you need to be able to measure the economic and socio-economic status of your citizens more readily.

The timeliness of Big Data is the primary benefit of this kind of data as it offers an immediate source of information to policymakers, who are often confronted with short-term horizons and imperfect information during the decision-making process. A second benefit of using Big Data is the ability to allow governments and social scientists to 'listen' and capture the needs and well-being concerns of their citizens, rather than relying on answers to pre-defined explicit questions. Organic data, such as the social media platforms of Twitter and Facebook, reveal what people care about most. The third benefit of this kind of data is that it offers social scientists the possibility to observe people's behaviour and not just opinions. This approach of revealed preferences unveils a reflexive picture of society because it allows the main concerns of citizens (and the priority ranking of those concerns) to emerge spontaneously and, as such, complements the information captured by gross domestic product. Lastly, Big Data, unlike survey data, do not suffer from non-response bias (Callegaro & Yang 2018).

## ***2.3. The evolution from survey data to organic data***

How did researchers get from traditional survey data to what they now call organic data? Well, when one thinks about survey data, the first thing that comes to mind is a census collection. Most governments, as far back as when Caesar Augustus called for a census of the entire Roman Empire, collect data on its citizens. Census collection is nothing more than an extensive survey completed by every individual in a household for a country. Through this survey data, governments collect information on their citizens' economic and socio-economic statuses. Interestingly enough, Bellet and Frijters (2019) argue that the primary reason governments collected survey data on its citizens in the past was to control them and to force the payment of taxes. Governments have come a long way since then and use both survey and organic data to measure, amongst other things, the well-being of people to increase life satisfaction for all.

#### ***2.4. The usual suspects: World Happiness Report and the World Value Survey***

Survey data is still instrumental, and social scientists have been able to learn a lot about how people perceive their lives to be relative to others. One of the best-known survey datasets currently used in all areas of social science is the World Happiness Report (WHR) (Helliwell et al. 2020). The WHR uses data collected by Gallup Analytics, which include data from a telephone survey of at least 1000 people aged 18 and older. Telephone surveys are used in countries where telephone coverage represents at least 80 per cent of the population or is the customary survey methodology. In countries where telephone interviewing is employed, Random-Digit-Dial (RDD) or a nationally representative list of phone numbers is used. Telephone methodology is typical in the United States, Canada, Western Europe, Japan, Australia and the like. In the developing world, including much of Latin America, the former Soviet Union countries, nearly all of Asia, the Middle East, and Africa, an area frame design is used for face-to-face interviewing. Face-to-face interviews take approximately 1 hour, while telephone interviews take about 30 minutes.

The Gallup operated survey (Helliwell et al. 2020) includes 6 measures of self-reported positive emotions (happiness, learning, life evaluations today and in 5 years, laughing, being respected) as well as 4 measures of negative emotions (anger, sadness, stress, worry).

The World Values Survey (WVS) (Inglehart et al. 2014) is also a dataset extensively used by social scientists studying changing values and their impact on social and political life. The survey, which started in 1981, seeks to use the most rigorous, high-quality research designs in each country. The WVS consists of nationally representative surveys conducted in almost 100 countries which contain almost 90 per cent of the world's population, using a standard questionnaire. The WVS is the largest non-commercial, cross-national, time series investigation of human beliefs and values ever executed, currently including interviews with almost 400,000 respondents. Moreover, the WVS is the only academic study covering the full range of global variations, from very poor to affluent countries, in all of the world's major cultural zones.

The WVS has demonstrated over the years that people's beliefs play a key role in economic development, the emergence and flourishing of democratic institutions, the rise of gender equality, and the extent to which societies have effective government.

Despite the achievements from the WHR and WVS, the subjective well-being measure raises several challenges and concerns among the broader scientific community. Firstly, Deaton (2013) points out that subjective well-being measures are affected by the limitation characteristic of all self-reporting measures. Therefore, it is not based on revealed behaviour and choices. Additionally, the works of Kahneman and Deaton (2010) and Steptoe et al. (2014) highlight the significant challenge of measuring aspects of multidimensional

subjective well-being such as pain or chronic disability, which are considered to have a very real and long-lasting impact. Secondly, survey questions on subjective well-being are limited in coverage, space and time. They cannot be measured at frequencies most useful for decisive government decision-making, nor reflect policy decisions made at a local level. Thirdly, the wording and order of survey questions have a real effect on measured outcomes. For example, if you follow a question regarding a person's income with a question on life satisfaction, there will be a marked difference in the assigned score (lower) (Deaton 2013).

## ***2.5. The unlikely heroes: Twitter, Facebook, Google and Instagram***

Now, that the foundation is laid by discussing some of the best-known survey data to date, let's turn the spotlight onto the new up-and-coming organic data that, through no fault of their own, are underpinning current and future research agendas.

### *2.5.1 Twitter*

People frequently use Twitter to broadcast their activities and observations and share their opinions about a wide variety of topics and events. The dynamic nature of this micro-blog and the amount of user-generated content that is publicly available and openly shared has attracted the interest of the research community in the field of smart cities. This rich, valuable and reliable content generates numerous possibilities for researchers in different domains. It has been used in large-scale temporal and geographical analyses and studies to complement traditional physical sensors.

In recent years, the popularity of social media platforms, such as Twitter have experienced tremendous global growth. In 2019, Twitter had 3.2 billion social media users worldwide, which accounts for 42 per cent of the total world population (Mohsin 2019). Various fields of study, such as psychology, sociology and economics (among others) can now analyze the impact of specific events or shocks with the help of Twitter data. Over the same period, Twitter has also expanded significantly, with 7 million new users each year (Statista 2019) and a total number of 330 million active daily users in 2019. The total number of tweets sent per day is an incredible 500 million, and the top three countries by user count outside the United States are Japan (35.65 million users), Russia (13.9 million), and the United Kingdom (13.7 million) (Omnicores 2020).

### *2.5.2 Facebook*

The sheer online ubiquity of Facebook is astounding. As of March 2020, Facebook had over 2.5 billion monthly active users who watched over 100 million hours of video every day. *Every 60 seconds*, there are 317,000 status updates, 400 new users, 147,000 photos uploaded and 54,000 shared links (Omnicores 2020). In short, since its creation in February 2004, Facebook has become the embodiment of a modern-day Cinderella

story creating 'virtual' coffee shops, sitting rooms or parks allowing millions of social interactions to play out every day. This growing new way of interacting and displayed social behaviour is inherently fascinating. Still, it also provides social scientists with an unprecedented opportunity to observe human behaviour in a naturalistic setting, test hypotheses in a novel domain, and recruit participants efficiently from many countries and demographic groups.

Wilson et al. (2012) argue that there are three broad reasons why Facebook is of relevance to social scientists. Firstly, it provides an unprecedented amount of data, through activities performed by its members (such as status updates, likes and dislikes expressing different preferential patterns and connecting to others). Therefore, this social media platform offers many new opportunities for studying human behaviour that previously had to rely on behaviours that were difficult to ascertain. Secondly, the tremendous popularity of Facebook makes it a topic worthy of study in its own right. Facebook and other online social networks (OSNs) are attractive to social scientists because, in addition to reflecting existing social protocols, they are also spawning new ones by changing the way 2.5 billion people relate to one another and share information. Thirdly, the rise of OSNs brings both unique benefits and dangers to society, which warrant careful consideration. The benefits associated with Facebook, such as the strengthening of social ties, are tempered by concerns about the time people spend on social media instead of interacting face-to-face with their loved ones. Additionally, the idea that information shared by people on Facebook is forever in the 'cloud' strengthens calls for better privacy and information disclosures.

### *2.5.3 Google Trends*

Researchers have used Google Trends data to investigate several questions, from exploring the course of influenza outbreaks to forecasting economic indicators. Google Trends is a website by Google that analyses the popularity of top search queries in Google Search across various regions and languages. Google handles over 75,000 queries per second which are an incredible 2.5 trillion searches per year. Additionally, it has a 95.65 per cent share of mobile search traffic worldwide, and although it already handles trillions of queries each year, the search volume grows by roughly 10–15 per cent annually (99 firms 2020).

The website uses graphs to compare the search volume of different queries over time. There are a number of different tools to explore Google search data: the public Google Trends website, the Google Trends Application Programming Interface (API) and the Google Health API. An API is nothing more than a software intermediary that allows two applications to talk to each other. For example, each time you use an app like Twitter, send a WhatsApp message or check the weather on your phone, you're using an API. All of these allow researchers to perform the same basic functions, such as offering results based on geographic locations and distinct time periods, but they do have some significant differences. The Trends API is available to journalists and academic researchers, while the Health API is only available to academic researchers.

#### *2.5.4 Instagram*

Instagram is an American photo and video-sharing social networking service owned by Facebook, Inc. It was created by Kevin Systrom and Mike Krieger and launched in October 2010. Instagram is the 6<sup>th</sup> most popular social network worldwide and had approximately 26.9 million new users in 2020, which is almost double the projected growth of other social media platforms. In total, Instagram has 1 billion monthly active users, and more than 500 million of them use the platform every day. A whopping 100 million photos and videos are uploaded daily, and the top three countries by user count outside the United States are India (80 million users), Brazil (77 million), and Indonesia (63 million) (Omnicores 2020). Two fun facts about Instagram that not many people know are the most popular hashtags being used: #Love, #Instagood, #Photooftheday, #Fashion, and #Beautiful, and the most Instagrammed food globally: Pizza, followed by Sushi.

Interestingly, in a study conducted by Lee et al. (2015), they found that Instagram users have five primary social and psychological motives: social interaction, archiving, self-expression, escapism, and peeking.

#### *2.6 Challenges and limitations of organic data.*

Alas, just like survey data, organic data is not free from limitations. The most significant criticism that researchers face when using organic data collected from social media sites is that the results of their studies, unlike those using survey data, are not representative of a country's population. Whether it is posted on Facebook or tweets from Twitter, the data represents a biased non-uniform subsample. Additionally, social media platform users are not representative of the overall population in that older people and children are vastly underrepresented. However, the demographics of Facebook and Twitter are known, and corrections can be made in the style of stratified sampling. In studies conducted by O'Connor et al. (2010) and Schwartz et al. (2013), they successfully fit their biased samples to representative data.

Another criticism is that people may post in a socially desirable fashion and thereby not present their true emotions or opinions. Although these criticisms do have merit, they are less problematic than is often assumed. Many psychology and sociology studies (see section 4) employ non-random population samples (such as college students), and many surveys suffer from desirability biases. Non-representative data is still valuable for understanding large populations in the same way that survey research has been valuable for understanding various populations.



### 3. Is organic data relevant to happiness research?

#### 3.1.A villain not of its own making: gross domestic product

Before discussing the relationship between organic data and happiness research, let's first take a step back and see where things went wrong. Between 1654 and 1676 Sir William Petty, an English economist, physician, scientist and philosopher, came up with the basic concept of gross domestic product (GDP) to attack landlords concerning unfair taxation during warfare between the Dutch and the English. Today, there is still debate as to who is primarily responsible for the evolution of the concept of GDP as it is currently understood, Simon Kuznets or John Maynard Keynes. Regardless of who it was, Kuznets was the first person to warn the United States government not to equate income with welfare: "*The welfare of a nation can, therefore, scarcely be inferred from a measurement of national income*" (United States 2020: page 7). More specifically, as argued by Algan et al. (2015), GDP does not measure non-market social interactions such as friendship, family happiness, moral values or a sense of purpose in life. This motivates the recourse to subjective self-reported measures of well-being, such as life satisfaction or measures of happiness, which economists use increasingly as a direct measure of utility. Nevertheless, even with Kuznets' warning, GDP played a pivotal role in the evolution of how researchers understand and measure human well-being (Table 1).

**Table 1: Evolution of the concept of human well-being (1950-2015).**

Period	Meaning of human well-being	Measurement of human well-being
1950s	Economic well-being	GDP growth
1960s	Economic well-being	GDP per capita growth
1970s	Basic needs	GDP per capita growth + basic goods
1980s	Economic well-being	GDP per capita but also the rise of non-monetary factors
1990s	Human development/capabilities	Human development and sustainability captured by social indicators <sup>3</sup>
2000s	Universal rights, livelihoods, freedom	The Millennium Development Goals and 'new' areas: risk and empowerment
2015	Universal rights, livelihoods, freedom and sustainability	Sustainable Development Goals

Source: Adapted from Sumner (2003).

<sup>3</sup> It falls outside of the scope of the paper to discuss the social indicators movement, but we refer readers to the work done by Verlet and Devos (2009) and Diener and Suh (1997) for additional information.

As table 1 shows, the evolution of the meaning and measurement of human well-being over time has been influenced by the position of development economics within development studies, and the tension between economic imperialism and multidisciplinary. Since development economists have moved away from a pure economic pursuit towards multidisciplinary approaches, so the concept of human well-being has been broadened from a concern about income towards a multidimensional understanding of well-being. Thereby, it recognizes that economic well-being, as measured by the GDP per capita, cannot explain the broader quality of life in a country on its own. The discussion in the 21<sup>st</sup> century has now moved entirely from empirical studies that aim to deepen society's understanding of human well-being to happiness as a policy objective. This has been witnessed by the works done by Stiglitz et al. (2009) when they wrote the report on 'Measurement of Economic Performance and Social Progress' commissioned by the French President Nicholas Sarkozy. President Sarkozy was not satisfied with the present state of their statistical information about the economy and society (focusing on GDP) and wanted a more in-depth measurement of people's perceptions of their well-being.

### ***3.2. Affect vs evaluative happiness***

Now that it has been established that one should measure happiness to influence policy, how do you define happiness? Many wiser (and long dead by now) and future Aristotleses will attempt to define happiness, but in reality there is no point, as this word has too many meanings for too many people. According to Haybron (2013), most scholarly work under the rubric of happiness centres on two senses of the term. The first usage treats happiness as basically a *synonym for the normative concept of well-being* and as such, *the notion that I know what is good for me or what will make my life go well*. In this well-being sense, one can assign happiness to a person when they have passed judgement (evaluated) over their lives and concluded that their lives are *going well for them*. This is the most natural reading of talk about *leading a happy life*, as opposed to simply *being happy*.

Simply being happy is a pure psychological state (affect happiness), denoting some broad and typically lasting aspect of the individual's state of mind. In considering the multitude of subjective well-being literature, this is the standard usage of happiness, and the prevalent norm in the language used. As argued by Veenhoven (2009), affect happiness is primarily determined by the gratification of needs in the first place and, as such, roots in human nature. Instead, the cognitive component is determined by the realization of wants which root in human culture. Therefore, the dominant views of happiness in this sense are hedonistic theories, which roughly equate happiness to pleasure. The life satisfaction theories equate happiness to an attitude of being satisfied with your life as a whole (this typically involves a global judgment about your life, as opposed to merely having a pleasant experience).

### ***3.3. Timing is everything!***

Traditional economic theory suggests that people display rational behaviour and that this rational behaviour drives their decision-making. However, as argued by psychologists such as Diener et al. (2009) and Kahneman (2011), people are not rational decision-makers, but emotive entities and the emotional state of a person influences their decision-making, as well as the mistakes they make, in no small degree.

In today's world, government and society make decisions at unpredictable times, be it daily or even intradaily, and they need to have the latest, most relevant information at their disposal to make informed decisions. Thus, organic data comes into play! The ability of social scientists to predict the happiness of a country's citizens by making use of, for example, a 'live' stream of tweets from Twitter, opens up the enormous potential for organic data to significantly influence countries' direction. If policymakers and social scientists could develop better predictive power by uncovering previously unknown insights into how people's emotions dictate their decision-making, many beneficial outcomes await society. Additionally, decision-makers seek to implement policies and regulations that increase the quality of life. Measuring happiness by using organic data in almost real-time is one useful way to assess the impact of policies as soon as they are implemented. Also, this will allow time-sensitive debates about potential policies that address specific current societal issues such as legalizing marijuana or removing abortions from the criminals' act.

## **4. Sentiment analysis using organic data**

Sentiment analysis, as applied to social media platforms, has received increasing interest from the research community. The reason is that an increasing number of people express their feelings, opinions and attitudes in social media, thereby generating organic data. Sentiment analysis is important, due to its applicability to a wide range of fields, such as economics, social science and the political economy. Several works claim that social phenomena such as stock prices, movie box-office revenues, and political elections, are reflected by social media data (Gayo-Avello 2013, Bollen et al. 2011, Asur & Huberman 2010) and that opinions expressed in those platforms can be used to assess the public opinion indirectly (O'Connor et al. 2010).

Sentiment analysis is an automated process of determining whether a text expresses a positive, negative, or neutral opinion about a topic (Hailong et al. 2014). Natural Language Processing (NLP) methods and algorithms are used to perform sentiment analysis. Some of these methods and algorithms include (Wolff n.d):

- i) Rules-based Systems: using, for example, lexicons.
- ii) Automatic systems: relying on machine learning techniques to learn from data.
- iii) Hybrid systems: combining both rule-based and automatic approaches.

#### ***4.1 Rules-based sentiment analysis***

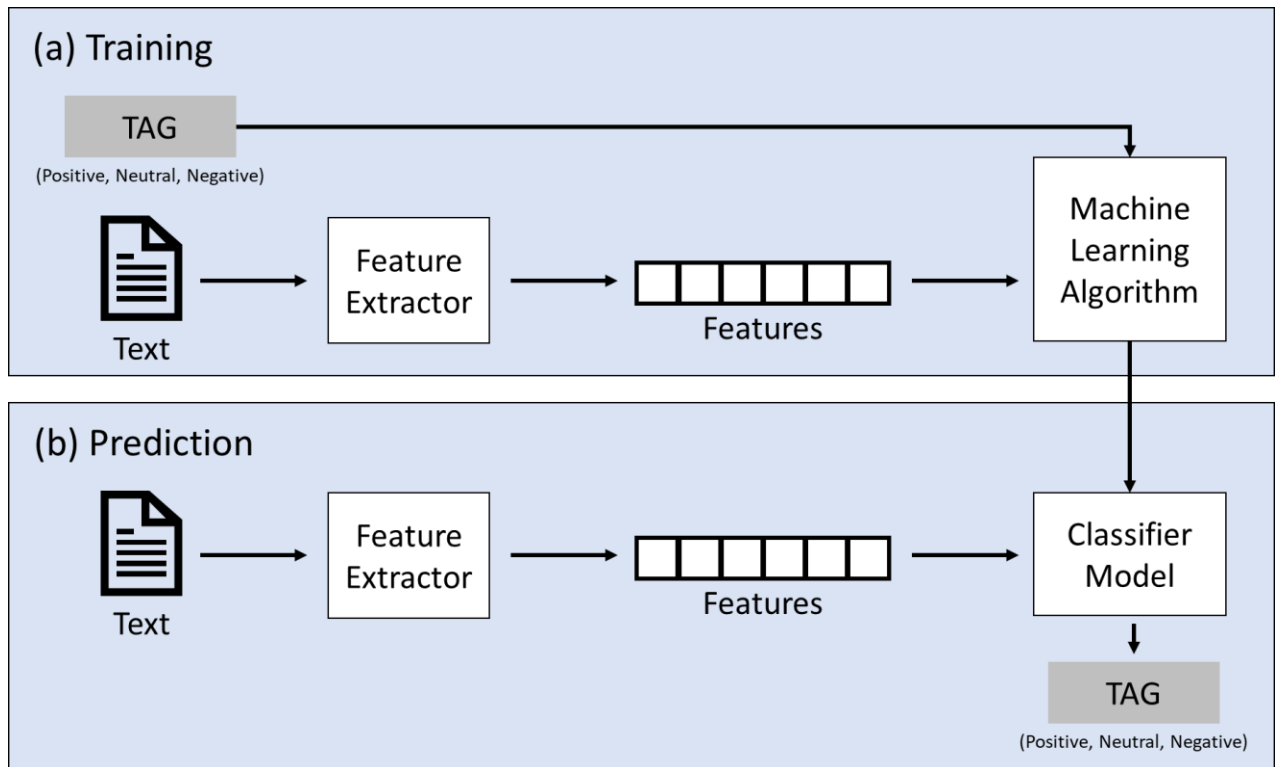
Rules-based sentiment analysis systems adopt a lexicon to perform sentiment analysis. This is done by counting and weighting sentiment-related words that have been evaluated and tagged. To collect the word list mentioned before, three main approaches are used: the manual, dictionary-based and corpus-based approach.

- i) Manual approach: This approach is very time consuming and is, thus, not usually used alone, but combined with automated approaches. The manual approach is then employed as a final validation method because automated methods make mistakes.
- ii) Dictionary-based approach: This approach uses a pre-defined dictionary of words, where each word is associated with a specific sentiment polarity strength. Feelings of people such as being happy, sad or depressed can be determined by comparing words against lexicons from dictionaries.
- iii) Corpus-based approach: This approach tries to find co-occurrence patterns of words to determine their sentiment. This approach is based on the seeding list of opinion words and then finds another opinion word which has a similar context. This method is used to assign a happiness factor to words, depending on the frequency of their occurrences in "happy" or "sad" blog posts.

#### ***4.2 Automatic sentiment analysis systems***

Different from lexicon or crafted manual rules, automatic system analysis relies on machine learning techniques, such as classification to determine the polarization of words.

To perform the automatic classification of words, the system must first be fed sample text before it can return a category, such as positive, negative, or neutral. Following this, there are two stages involved in the implementation of automatic systems, namely the training and the prediction stage. In the training stage, the sentiment analysis model is trained to correctly tag text as *negative*, *neutral* or *positive*, using sample data (see figure 1). The feature extractor then transforms the text into a feature vector, which creates pairs of feature vectors and tags, for example positive, negative, or neutral, that are fed into the machine learning algorithm to generate a model. In the prediction stage, the feature extractor is used to transform the unseen text into feature vectors, which are fed to the model, enabling it to make sentiment predictions.



Source: (Wolff n.d)

Figure 1: Automatic sentiment analysis system

#### 4.3 Hybrid Sentiment Analysis Systems

Using the hybrid systems, they combine rule-based and automatic techniques to analyze the sentiment of the text. This method often delivers more accurate results than using the systems individually.

### 5. Latest research using organic data

This section will provide an overview of the latest or most influential studies relying on organic data.

#### 5.1. Organic data in psychology

Where does psychology fit into the field of Big Data or related areas such as computational social science? There are several areas in which psychology can and has begun to weigh in, such as well-being, mental health, depression, substance use, behavioural health, behaviour change, social media, workplace well-being and effectiveness, student learning and adjustment, and behavioural genetics. While there are undoubtedly limits to that which may eventually be quantified regarding human behaviour, recent studies in psychology have demonstrated many thriving and diverse methodologies that were all impossible before the Internet age.

Alharthi et al. (2017) desired to shed light on the importance of understanding the causes of an individual's emotions and the underlying motivations behind their actions, as explained by the Human Needs Theories (HNT). They developed an annotation framework to analyze social media textual data, from Twitter, for basic needs concepts from a psychological well-being perspective. Alharthi et al. (2017) succeeded in introducing a basic psychological needs corpus consisting of 6,334 tweets that capture and analyze explicit and implicit human needs expressions. The corpus was annotated with emotion categories, psychological needs, levels of satisfaction of the needs, types of social contexts and life domains.

In a fascinating study by Hills et al. (2019), they produced a workable proxy for subjective well-being going back to 1776. This allowed the authors to compare their proxy to GDP over the same period. The period examined was post-1820 when readily accessible data became available. Additionally, the authors were able to assess changes in life expectancy, together with economic and social events, throughout this period for those countries in the study. Their analysis allowed them to gain retrospective insight, by using language corpora, into the rise and fall of subjective well-being as derived from historical language use. The availability of large corpora allowed them to do so for six languages (English (British), English (American), German, Italian, Spanish, and French).

Hills et al. (2019)'s key source was another unlikely hero (just like Twitter), namely, Google Books corpus. This corpus is a collection of *word frequency data* for over 8 million digitized versions of physically published books (6 per cent of all books published up to 2019). To assess the specific emotions underpinned by individual words, they used the largest available sets of existing emotive words rating norms for each of the six languages.

Well-being, which encompasses much more than emotion and mood, is linked with good mental and physical health. Psychologists such as Seligman (2011) and Ryff and Keyes (1995) break well-being into separate domains. In such "*dashboard approaches*" well-being is best measured as separate, correlated dimensions. In his well-being theory, Seligman (2011) suggests five major pillars that together contribute to a person's sense of well-being: Positive Emotions, Engagement, Relationships, Meaning, and Accomplishment (PERMA). Other "*dashboard approaches*", as stated by Ryff and Keyes (1995), seek to capture subtle psychological notions such as "autonomy" or "self-acceptance". Promoting whole-body health will be helped significantly by the results obtained from a combination of national surveys and the ability to quickly and accurately assess "*dashboard approaches*". Schwartz et al. (2016) predicted individual well-being, as measured on a life satisfaction scale, through the language people use on social media. Their findings highlight the important role social media play and suggest that analyses of language go beyond mere prediction by identifying the language that characterizes well-being.

Ford et al. (2018) explored the potential for internet search data to serve as indicators of subjective well-being and predictors of health at the state and metro area levels in the United States. They argued that there was a direct relationship between searches for positive and negative affect-related words and the actual emotions being experienced at the time of the search. Apart from 15 affect words collected from Google's Trends website, they also used data on health, self-reported emotions, psychological well-being, personality, and Twitter postings at the state and metro area levels. The researchers found several internet search scores correlated with indicators of cardiovascular health and depression. Some search term scores also correlated strongly with self-reported emotions, well-being metrics, neuroticism, per capita income, and Twitter postings at the state or metro area level. After additional analyses, they found that affect word searches do predict depression rates at the metro area level, more so than the effects of income and other well-being measures. This reiterates the theory stating that internet search data can be useful in ascertaining physical and mental health at the aggregate level.

## ***5.2. Organic data in economics and finance***

Many studies have used Twitter to investigate the potential of using this organic data for the prediction of stock market returns or movement (Broadstock & Zhang 2019, Tabari et al. 2018, McKay 2018, Renault 2017, Rao & Srivastava 2012, Bollen et al. 2011). Apart from the success that daily data has shown (in the studies mentioned above) in accurately predicting stock market movements, high-frequency intraday data also make a significant contribution. Investment decisions are made continuously throughout the day and not merely at one specific point of time during the day. Therefore, the ability of Twitter to provide information on the sentiment and emotions of investors, in almost 'real-time' in both developed and emerging stock markets are gathering support. To this end, Steyn et al. (2020) analyzed eight stock markets, including six developed countries (France, Germany, the U.K., the USA, Japan and Spain) and two emerging markets (Poland and India) by using high-frequency intraday data. The purpose of this study was to investigate the ability of sentiment and emotions extracted from tweets to predict stock market movements. Their study was the first one to use a combination of three different machine learning classification algorithms, various evaluation metrics, and validation techniques to ascertain the predictive ability of their models. Their results indicate significant predictive ability over stock market movements in both developed and emerging markets when one considers investor sentiment and emotions extracted from tweets.

The potential of using Twitter to provide additional information about the quality of life that can be included in a planning context and effectively used to improve the decision-making process of government officials was demonstrated by Zivanovic et al. (2020). In a study focusing on the city of Bristol, in the United Kingdom, the author analyzed 1,374,706 geo-tagged tweets by using a combination of manual coding of messages, automated classification and spatial analysis. Through this analysis, the author was able to capture people's perception of

their quality of life and found that the domains of health, transport and the environment were important to Bristol residents. Upon further analyses, Zivanovic et al. (2020) found a difference between the wards in Bristol, in the number and type of quality of life perceptions in every domain, the spatial distribution of positive and negative perceptions, and differences between the domains. Furthermore, based on people's opinions, there is a difference in the quality of life between Bristol neighbourhoods. The author concluded that Twitter data could indeed be used to evaluate the quality of life and be used to complement the official quality of life surveys.

In terms of research focusing on pure economic indicators and Big Data, a revolutionary study conducted by Richardson et al. (2018) is most enlightening. In the world, as people know it, policymakers make real-time decisions using incomplete information on current economic conditions. It is not only survey data that is released with a time lag but also most key statistics. In their paper, Richardson et al. (2018) analyze the real-time nowcasting performance of machine learning algorithms estimated on New Zealand data. Using a large set of real-time quarterly macroeconomic indicators, they train a range of popular machine learning algorithms and nowcast real gross domestic product growth for each quarter over the 2009-2018 period.

Furthermore, they compare the predictive accuracy of these nowcasts with that of other traditional univariate and multivariate statistical models. Their results find that the machine learning algorithms outperform the traditional statistical models. Moreover, combining the individual machine learning nowcasts, further improves performance relative to individual nowcasts.

### ***5.3. Organic data in sociology***

It seems that sociology has embraced the use of organic data the most thus far, with a significant number of studies using social media to analyze human behaviour and its consequences. Therefore, this section will focus the discussion on only a select few.

#### *5.3.1 Health*

Analyzing user messages in social media can measure different population characteristics, including public health measures. For example, Twitter messages have been used in tracking the influenza rates in the United Kingdom and United States (Lampos and Cristianini 2010, Culotta 2010), but this has primarily been the extent of mining Twitter for public health. Paul and Dredze (2011) considered a broader range of public health applications for Twitter. They applied the Ailment Topic Aspect Model to over one and a half million health-related tweets and discovered mentions of over a dozen ailments, including allergies, obesity and insomnia. They introduced extensions to incorporate prior knowledge into this model and apply it to several tasks: tracking illnesses over times (syndromic surveillance), measuring behavioral risk factors, localizing illnesses



by geographic region, and analyzing symptoms and medication usage. Paul and Dredze (2011) showed quantitative correlations with public health data and qualitative evaluations of model output. Their results suggest that Twitter has broad applicability for public health research.

Excessive alcohol use is a global public health challenge that causes wide-spread problems in both social-capital (violence, crime and imprisonment) and health itself. In the United States (U.S.) over 88,000 deaths per year are attributed to excessive alcohol use, which costs the health system over \$250 billion annually. Using 38 billion tweets, Giorgi et al. (2020) characterized how drinking-specific language varies across regions and cultures in the U.S. They identified 3.3 million geolocated "drunk" tweets from their initial sample and correlated their language with the prevalence of self-reported excessive alcohol consumption. Giorgi et al. (2020)'s study shows that Twitter can be used to explore the specific sociocultural contexts in which excessive alcohol use occurs within particular regions and communities (religious communities had a high frequency of anti-drunk driving tweets, Hispanic communities discussed family members drinking, and college towns discussed sexual behaviour). These findings confirm that Twitter can be used to deliver targeted public health messages for specific high-risk groups, as well as understanding those cultural determinants of substance abuse better.

### 5.3.2 Religion

Religion is "*the opium of the people*", or so Karl Marx thought in 1843. He did concede that religion could provide comfort in difficult circumstances, though Marx argued this comfort was nothing more than *people deceiving themselves*. One hundred and sixty-three years later, Dawkins (2006) reiterates this negative sentiment expressed by Marx, arguing that the world would be a happier place without religion. Despite their objections, there is evidence of a positive correlation between all four major world religions: Christianity, Islam, Buddhism, Hinduism and happiness (see, for example, Villani et al. 2019, Diener et al. 2011).

With 19 per cent of the world's population, Islam adherents, Weber et al. (2013) use public data in English and Arabic from Twitter to study the phenomenon of secular versus Islamist polarization. Their research focuses on the number of retweets and hashtags from both camps to ascertain the extended network of users and to quantify how polarized society as a whole is at a given point in time. Their analysis followed the manual approach, as discussed in section 3.1 and focused on (i) religious terms, (ii) derogatory terms referring to other religions, and (iii) references to charitable acts. The authors found strong indications that a measure of global hashtag polarization, related to the overlap between hashtags used by the two political sides, works as a "barometer for tension" with high values coinciding with periods of violent outbreaks. Given their results, they argue that Twitter could possibly be used as a forecasting tool.

To add to the debate on whether the relationship between religion and happiness is positive, Ritter et al. (2014) used nearly 2 million tweets to examine differences between Christians and atheists in natural language. The sample of Christians and atheists came from those who chose to follow the Twitter feed of five Christian public figures or five atheist public figures. The five Christian public figures were Pope Benedict XVI (@PopeBXVI), Dinesh D'Souza (@DineshDSouza), Joyce Meyer (@JoyceMeyer), Joel Osteen (@JoelOsteen), and Rick Warren (@RickWarren). The five atheist public figures were Richard Dawkins (@RichardDawkins), Sam Harris (@SamHarrisOrg), Christopher Hitchens (@ChrisHitchens), Monica Salcedo (@Monicks), and Michael Shermer (@MichaelShermer). Contradicting the very notion that happiness will be higher without religion, the results showed that Christians use less negative and more positive emotion words than their atheist counterparts.

Moreover, Ritter et al. (2014)'s study was the first one of its kind to highlight that the relationship between happiness and religion is partially mediated by *thinking style*. This result reinforces previous laboratory studies and self-reported data, suggesting that social connection *partially mediates* the relationship between happiness and religiosity.

### 5.3.3 *Crime and violence*

Traditional crime prediction systems make extensive use of historical incident patterns as well as layers of information provided by geographic information systems and demographic information repositories. Wang et al. (2012) decided to investigate Twitter's potential to be used as a predictive tool for future criminal activities. Their approach is based on automatic sentiment analysis, as discussed in section 3.2 and involved running predictions using linear modelling. Interestingly, they tested their model on predicting future hit-and-run crimes in Charlottesville, Virginia, in the United States for an eight-month period. Their model outperformed the baseline model in predicting future hit-and-run incidents uniformly across the entire period.

Similarly, Gerber (2014) investigated crime prediction using geo-tagged tweets. He used Twitter-specific linguistic analysis and statistical topic modelling to automatically identify discussion topics across Chicago, Illinois in the United States. After incorporating the aforementioned discussion topics into his crime prediction model, the author showed that for 76 per cent of all crime types studied, Twitter data significantly improved the predictive power of future crimes relative to the standard approach based on kernel density estimation (a technique to estimate the unknown probability distribution of a random variable, based on a sample of points taken from that distribution). Gerber's (2014) research has implications specifically for criminal justice decision-makers in charge of resource allocation for crime prevention. More generally, this research has implications for decision-makers concerned with geographic spaces occupied by Twitter-using individuals.

De Choudhury et al. (2014) used Twitter data in conjunction with country-level homicide data from the Mexican government to examine whether residents of four major cities in Mexico displayed desensitization to protracted violence resulting from the Mexican Drug War. Over a 2-year period they found that negative affect expressed in tweets declined despite increases in homicides.

Experiences of community-wide traumas, especially of mass violence, are associated with psychological distress and adverse physical health outcomes. However, studying communities impacted by mass violence is often costly, requires swift action to enter the field when disaster strikes and may be invasive for some traumatized respondents. Typically, individuals are studied after the traumatic event with no baseline data against which to compare their post-trauma responses. Given these challenges, Jones et al. (2016) used longitudinal Twitter data across three case studies to examine the impact of violence near or on college campuses. They studied the communities of Isla Vista, California, Flagstaff, Arizona, and Roseburg, Oregon, in the United States and compared their results with control communities, for the period 2014-2015. In Isla Vista, they observed an increase in post-trauma negative emotion expression among sampled followers after mass violence. They showed how patterns of response appear differently, based on the timeframe under scrutiny. In Flagstaff, they replicate the pattern of results among social media users in the control group from Isla Vista, after a campus shooting in that community killed one student. In Roseburg, they replicate this pattern in another group of Twitter users likely to live in a community affected by a mass shooting. They conclude that Twitter is a creative tool that allows health officials to ascertain where the biggest needs lie within certain affected communities. This will support any decision about deploying frontline staff to deal with psychological distress in the aftermath of violence.

#### *5.3.4 Politics*

Even though researchers do not know a lot about happiness and political behaviour, the issues are nevertheless of pressing significance for policy. Governments and policies affect individual's happiness, which means happiness is seen as an outcome (the end). Government institutions and policies set the stage on which people live their lives. Unfortunately, much less is known about the effects of societal happiness on political behaviour and outcomes.

A strong positive relationship exists in the U.S., U.K. and China, between the level of happiness and voter turnout (Ward 2019). People who are depressed are less likely to vote. Moreover, also demonstrated in section 5.2 for Australia, apart from economic voting, a clear and significant positive relationship exists between national happiness in the run-up to general elections and the subsequent electoral success. Therefore, it is not surprising that there is a growing interest in monitoring Twitter, as it is a space where people talk a lot about government policy. Whether governments like it or not, social media is changing the way that people participate in democracy. It's engaging some of the most disengaged. Young people, who are least likely to

vote, say that they feel more engaged politically when they participate in social media for political purposes and that they'd be more likely to vote because of it.

In a groundbreaking study by O'Connor et al. (2010) they connected measures of public opinion measured from polls with sentiment measured from tweets. Their study is one of a handful that was able to successfully fit their biased samples (from Twitter) to representative data. The authors analyzed several surveys on consumer confidence and political opinion over the 2008 to 2009 period and found that they correlate to sentiment word frequencies in contemporaneous tweets. They do recognize that their results vary across datasets, but stress that in several cases the correlations were as high as 80 per cent and captured important large-scale trends. Their results set the stage for others in highlighting the potential of using non-representative subsamples as a substitute and supplement for traditional polling.

In a related study, Tumasjan et al. (2010) used the context of the German federal election to investigate whether Twitter is used as a forum for political deliberation and whether online messages on Twitter validly mirror offline political sentiment. Using linguistic inquiry and word count text analysis software, they conducted a content analysis of over 100,000 tweets containing a reference to either a political party or a politician. Results show that Twitter is indeed used extensively for political deliberation. They find that the mere number of tweets mentioning a party reflects the election result. Moreover, joint mentions of two parties are in line with real-world political ties and coalitions. An analysis of the tweets' political sentiment demonstrates close correspondence to the parties' and politicians' political positions indicating that the content of Twitter messages plausibly reflects the offline political landscape.

Rill et al. (2014) presented a system called *PoliTwi*, which was designed to detect emerging political topics (Top Topics) in Twitter sooner than other standard information channels. The recognized Top Topics are shared via different channels with the broader public. For the analysis, they collected about 4 million tweets before and during the parliamentary election in 2013 in Germany for a six-month period. Their results indicated not only that new topics appeared almost instantaneously on Twitter but that their appearance was also faster relative to Google Trends. The ability of the German government to interact and influence voters before elections using Twitter was an additional insight into their study.

Riotta et al. (2014) presented an analysis of the behaviour of Italian Twitter users during national political elections. They monitored the volumes of the tweets related to the leaders of the various political parties and compared them to the results of the elections. Furthermore, they studied the topics that were associated with the co-occurrence of two politicians in the same tweet. They could not conclude, from a simple statistical analysis of tweet volume and their time evolution, that it is possible to predict the election outcome precisely. In saying that, they did find that the volume of tweets and their change in time provided a very good proxy of the final results.

## 6. Happiness indices and organic data

### 6.1 Hedonometer

Hedonometer is one of the first measures of happiness created using organic data. The research of Dodds and Danforth (2010) at the University of Vermont Complex Systems Center is the basis of the Hedonometer, and Brian Tivnan, Matt McMahon and their team from the MITRE Corporation do the technological support. The project started at the end of 2008 and measures happiness levels continuously per day, thus resulting in a time series from the end of 2008 to the present (for more information read the foundational paper Dodds et al. (2011)).

#### 6.1.1 Methodology

To construct the Hedonometer, the authors used different sources, namely: Google Books, New York Times articles, Music Lyrics, and Twitter to determine the 5000 words most often used in each of these corpora. They merged these words, resulting in a composite set of 10 222 unique words. Making use of Amazon's Mechanical Turk service, they had the words scored on a nine-point scale of happiness: from 1 (sad) to 9 (happy) (see Table 2). The word that has the highest happiness score is laughter, namely 8.5 out of 9, and it is ranked number one out of the 10222 words. The word with the lowest score (saddest) is terrorist, it has a score of 1.3 out of 9, and it is ranked number 10222 out of 10222 words.

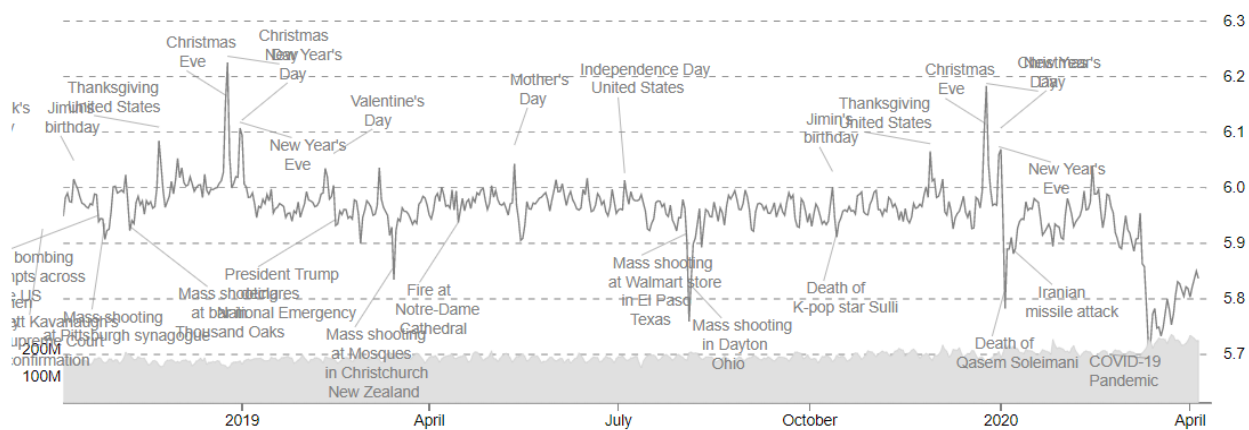
Once they have derived the happiness values of each word, they use these values to construct the Hedonometer. They use the Twitter Decahose API feed, which is a Streaming API feed that continuously sends a sample of roughly 10 per cent of all tweets to Dodds and his team.

**Table 2: Hedonometer, word scores on a 9-point scale of happiness.**

Word	Happiness rank	Score
Laughter	1	8.5
Love	3	8.42
Rainbow	13	8.06
Congratulations	25	8.00
Enjoy	112	7.66
Luxury	312	7.30
Journey	469	7.14
Horny	3313	5.82
Ministry	4410	5.58
Arrested	10209	1.64
Deaths	10219	1.64
Terrorist	10222	1.3

Source: Dodds et al. (2011).

To construct the Hedonometer, they bin all the tweets extracted on a day; however, they include only those words that are recognized to be English. On average the bin includes 200 million words, extracted across the world, per day. Each word receives the predetermined happiness score, to derive an average happiness score per day, see figure 2 for a representation of the Hedonometer for 2019 up to the day of writing this paper.



Source: Hedonometer (2020).

Figure 2: Hedonometer - average happiness using Twitter data.

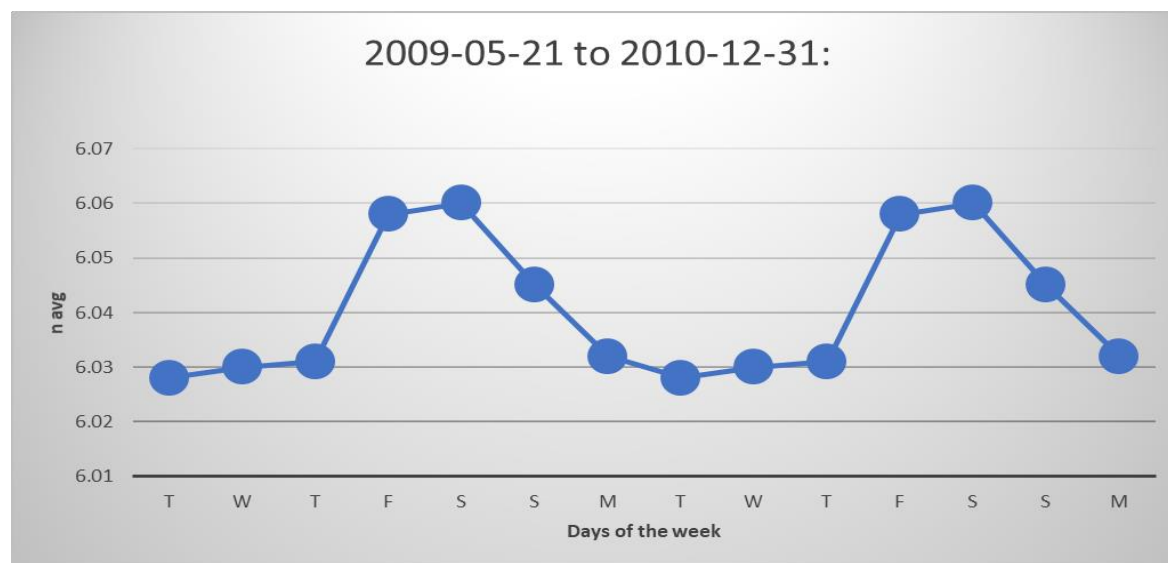
Apart from the Hedonometer time series, there have also been analyses of subsamples per geographical region. For example, for all the States of the USA from 2011 to 2013 and the Cities of the USA from 2012 to 2014. Mitchell et al. (2013) conducted a detailed investigation of correlations between the Hedonometer and a wide range of emotional, geographic, demographic, and health characteristics. They achieved this by combining their real-time data with that of annually surveyed data. The real-time data is in the form of their massive, geo-tagged dataset comprising over 80 million words generated in the year 2011 extracted from the Twitter Application Programming Interface (API). The annually surveyed data comprise characteristics of all 50 U.S. states and close to 400 urban populations. Based on the previous works of Dodds and Danforth (2010) and Dodds et al. (2011), they were able to generate taxonomies of states and cities based on their similarities in word use. Additionally, the team estimated the happiness levels of each state and city. The results demonstrated the correlation which exists between highly resolved demographic characteristics and the happiness level of each state or city's Hedonometric level. Once again, their results show how social media, and in particular Twitter, may potentially be used to estimate real-time happiness levels.

Other studies undertaken by the Hedonometer team include the measurement of the happiness levels of specific sources, such as the happiness of News as reflected by the New York Times, or the Happiness of Stories according to a ranked list of movie scripts.

### 6.1.2 Results

Some interesting results from studying the Hedonometer include: (i) the average happiness levels per day of the week, and (ii) the events related to the happiest and saddest days in the time series.

Regarding the average happiness levels per day of the week, the research team found that Tuesdays, against expectations as the general opinion is Mondays, have the lowest average happiness scores. After that the happiness levels increase to peak on Saturdays, slightly above the score for Fridays (see figure 3).



Source: Dodds et al. (2011). Note: The scale is from 1 extremely negative to 9 extremely positive.

Figure 3: Daily average happiness scores according to the Hedonometer.

The Hedonometer time series also shows those days that strongly deviate from the norm; the happiest (most positive) and saddest (most negative) days. The positive days usually occur on annual religious, cultural, and national event days, for example, Christmas Day, New Year's Day, Valentine's Day, Thanksgiving, Fourth of July, Mother's Day, and Father's Day. It is understandable that these days measure among the happiest days of the year and reflects a substantial degree of social synchrony.

The negative days are typically when unexpected societal trauma, due to, for example, natural disasters or the death of a celebrity, occurs. In 2008 the financial Bailout Bill of the U.S. induced a multi-week depression, with the lowest point on Monday, September 29, 2008. Another decrease in the Hedonometer was at the onset of 2009 with the swine flu or H1N1 pandemic. The most significant single-day decrease in the index was on Michael Jackson's death, July 7, 2009. Natural disasters also contribute to relatively low levels of happiness, for example, the Chilean earthquake in February 2010.

### 6.1.3 *Limitations of the Hedonometer*

Although the Hedonometer has many advantages and has contributed to a better understanding of our human moods, it also has certain limitations. The first limitation deals with using Twitter itself, which refers to the representativeness of the data.

Additionally, the Hedonometer cannot deal with the context in which words are used, as words in itself are evaluated and not the sentiment of the construct. For example, a phrase such as "I did not enjoy the holiday", will attract a score of 7.66 for "enjoy" and 7.96 for "holiday", thus reflecting an overwhelmingly positive sentiment, when actually the sentiment is negative.

Currently, the Hedonometer calculates a happiness index on a scale of 1 (sad) to 9 (happy), but it cannot detect the emotions underpinning the words or the tweets. Thus, it cannot determine if the changes in the levels of happiness are due to negative emotions such as fear or anger, or positive emotions such as joy and trust.

The Hedonometer only analyses tweets that are made in English and does not consider tweets made in other languages. Approximately only 34 per cent of all tweets are in English. The most used other languages include Korean, Thai, Turkish, French, Japanese, Spanish, and Portuguese. Therefore, the Hedonometer only reveals the "mood" of the tweets of the English-speaking world, excluding the "mood" of any other tongue speakers (Kalev et al. 2013).

## 6.2 *The Gross National Happiness Index (South Africa, New Zealand and Australia)*

In 1979, the King of Bhutan, Jigme Singye Wangchuck, famously made the following statement to a Financial Times journalist at the Bombay airport: "*Gross National Happiness is more important than Gross National Product*". This notion that the government should care more about its peoples' happiness (the end) than GDP (the means) has not been lost throughout the last four decades.

Forty years after this profound statement of the King of Bhutan, two well-being economists, put it to work when the team launched a Gross National Happiness Index (GNH) for three Commonwealth member states; South Africa, New Zealand and Australia (Greyling, Rossouw & Afstereo 2019). This project aimed to measure the mood of countries' citizens during different economic, social and political events.

Since February 2020, the researchers extended the project that initially analyzed the sentiment of tweets to incorporate the analysis of the emotions underpinning tweets. The team did this to determine which emotions are most prominent on specific days or events. These analyses are especially insightful in cases where there are shocks, such as COVID-19, to determine the emotions of a nation under challenging circumstances and in events where one expects changes in emotions.



## 6.2.1 Methodology

This section discusses the methodology used to (i) construct the GNH index, and (ii) evaluate the emotions underpinning tweets.

### 6.2.1.1 GNH index

To construct the GNH index for the different nations, the researchers use Big Data methods and extract tweets from the voluntary information sharing structure of Twitter. Thereafter, the team applies sentiment analysis to a live Twitter-feed and labels every tweet as having either a positive, neutral or negative sentiment. This sentiment classification is then applied to a sentiment-balance algorithm to derive a happiness score. The happiness scores range between 0 and 10, with 5 being neutral, thus neither happy nor unhappy.

The team extracts all tweets made in each of the three countries per day and calculates a happiness score per hour. The index is available live on the GNH website (Greyling et al. 2019). In South Africa, the average number of tweets extracted is 63 000 per day. South Africa has approximately 11 million Twitter users, representing almost 18 per cent of the population. In Australia, the average number of tweets extracted per day is 34 300, with 5.3 million Twitter users representing 22 per cent of the population (David 2020). In New Zealand, the average number of tweets extracted per day is 6 400, with 330 000 Twitter users, representing 8 per cent of the population (Omnicores 2020). Although the number of tweets is extensive and represents significant proportions of the populations of the countries, it is not representative. However, Twitter accommodates individuals, groups of individuals, organizations and media outlets, representing a kind of disaggregated sample, thus giving access to the moods of a vast blend of Twitter users, not found in survey data. Furthermore, purely based on the vast numbers of the datasets, it seems that the GNH index gives a remarkably robust reflection of the mood of a nation.

### 6.2.1.2 Emotions

To analyze the emotions rather than the sentiment of tweets, the team analyses the "words" of a tweet to determine the emotion underpinning the specific word. The researchers differentiate between eight emotions, namely *joy, anticipation, trust, disgust, anger, surprise, fear and sadness*.

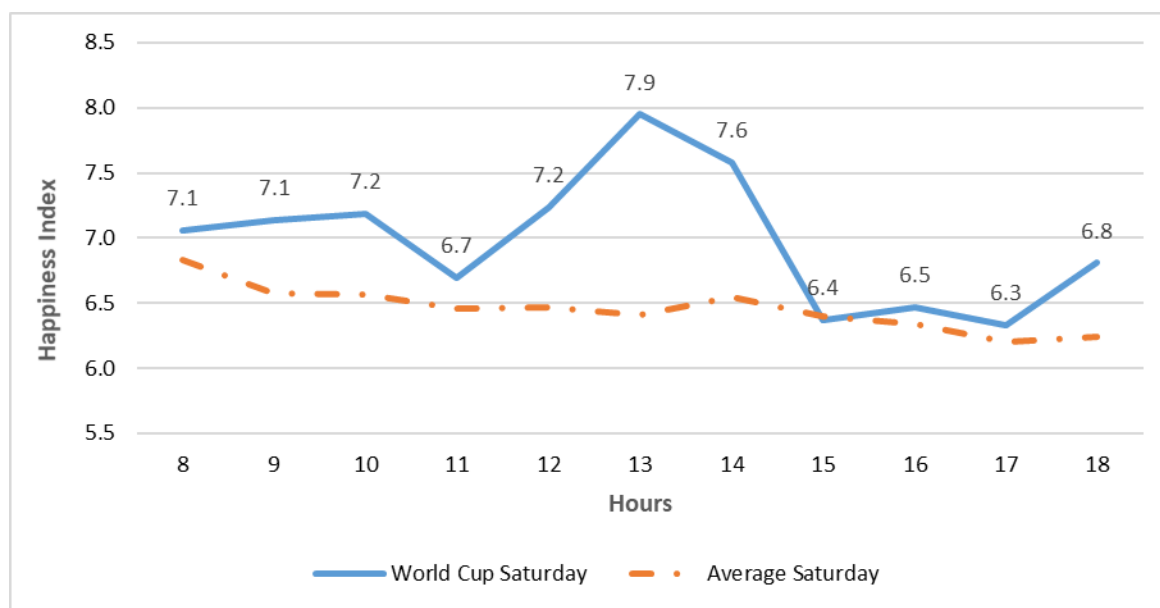
## 6.2.2 Results

### 6.2.2.1 Sentiment analysis

The researchers' primary interest in the index is to follow the mood of the nations, included in the study during specific political, economic, and social events.

To consider the mood of a nation during a sporting event, they used the example of the South African Springboks winning the Rugby World Cup on November 2, 2019. The rugby game commenced at 11:00 am (GMT+2) and

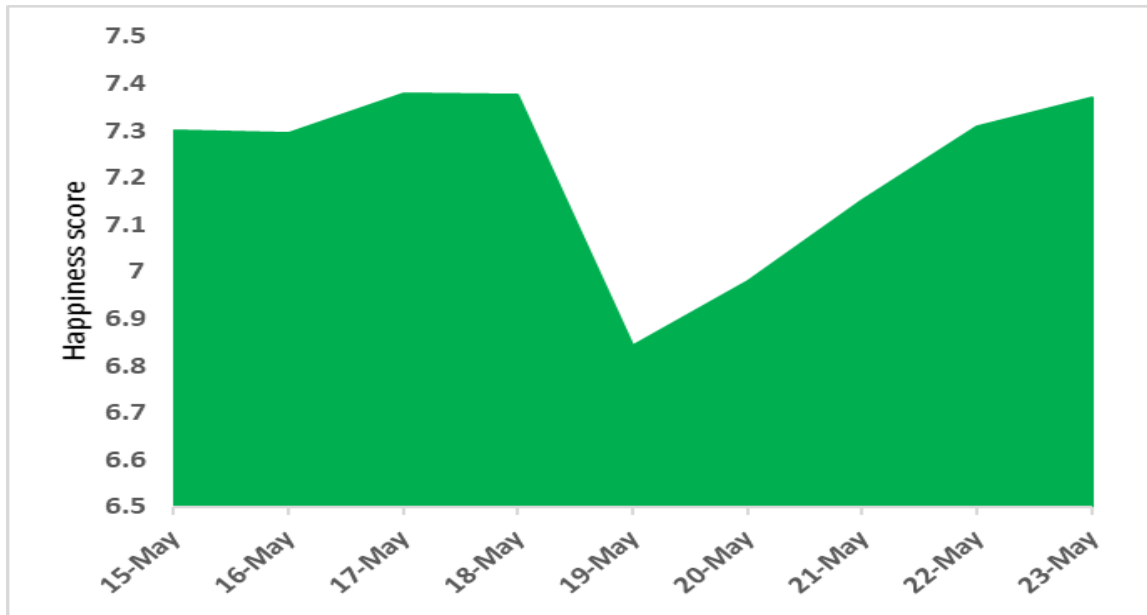
lasted until 1:00 pm. Figure 4 shows the intraday per hour measurement of happiness levels. As can be seen, since waking South Africans' happiness levels were higher than the norm of 6.3, caused by the higher expectations expressed for the day. At 11:00 am there was a slight dip, but as the game proceeded and it became clear that the Springboks had the upper hand over their opponents, England, the happiness score increased to reach an hourly high of 7.9 at 1:00 pm (the end of the game). Coincidentally, this happiness level of 7.9 is the highest in South Africa, since launching the index in April 2019. This illustrates the power of sporting events to influence the mood of a whole nation.



Source: Greyling et al. (2019)

Figure 4: South Africa's intraday (hourly) happiness levels during the Rugby World Cup 2019.

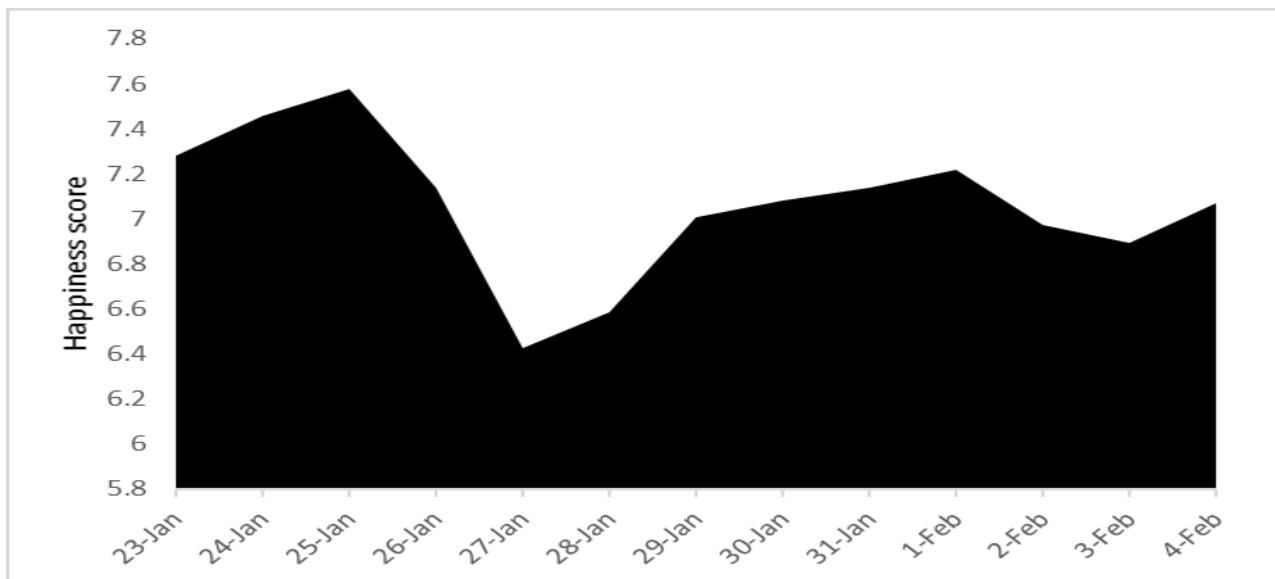
A political event also affects the mood of a nation, for example in Australia the federal elections were held on May 18, 2019. Figure 5 shows the happiness levels of Australians, before and after the election day (May 18 2019). In the period leading up to the elections, the general mood of the country was higher than the norm of 7.3. However, the happiness index showed a significant drop in the mood from 18 to May 19, 2019 (see figure 5), as initially polls and TV channels reported that results favoured the Labor party. Still, then the results turned and the Coalition party took the lead to win the elections, against expectations. These events show that the happiness index reveals the mood of a nation during a political event.



Source: Greyling et al. (2019).

Figure 5: Happiness levels during the Australian elections.

In terms of a societal trauma, the researchers considered New Zealand's reaction to the American professional basketball player, Kobe Bryant's, death on January 27, 2020. Figure 6 shows that on January 27, the happiness levels of New Zealanders dropped significantly below the norm of Monday's, showing the empathy of New Zealanders with the death of a celebrity.

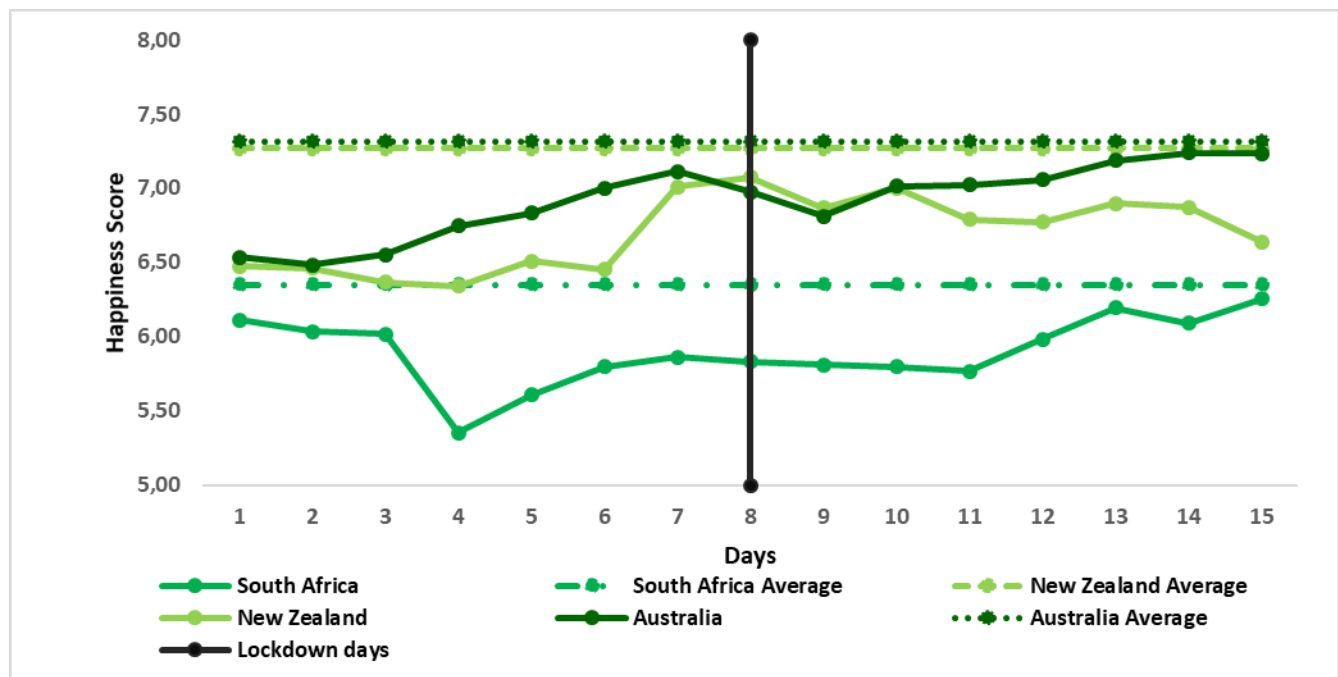


Source: Greyling et al. (2019).

Figure 6: Happiness levels in New Zealand, death of Kobe Bryant.

Another societal trauma that affected the entire world was the COVID-19 pandemic. In terms of measuring the happiness levels, the GNH was uniquely placed to record the effect of the pandemic. When the first cases of the Coronavirus (COVID-19) were confirmed in these countries (January 25 in Australia, February 28 in New Zealand and March 6 in South Africa), the research team saw that society did not react to this 'faceless, impersonal' virus at first. Still, as more cases were announced, and the public realized the threat of the disease, the happiness levels dropped far below the daily average levels (see figure 7). These lower levels of happiness remained until the day this paper was written.

The "Lockdown day" indicated by the vertical line on figure 6 shows the day on which the regulations to curb the Coronavirus were introduced in each country, respectively. Interestingly, the levels of happiness were lower in the days leading up to the implementation of the regulations than after that. This is an indication that fear and negative expectations of the unknown can drive happiness levels down. Once the regulations were implemented the happiness levels were still lower than the norm, but slightly more positive than before, as people adjusted to these new circumstances. At the time of writing this paper (April 2020), the authors reported that Australia and South Africa were nearly back to their normal happiness levels, suggesting the incredible ability of people to adjust to their new way of living after a major shock.



Source: Greyling et al. (2019).

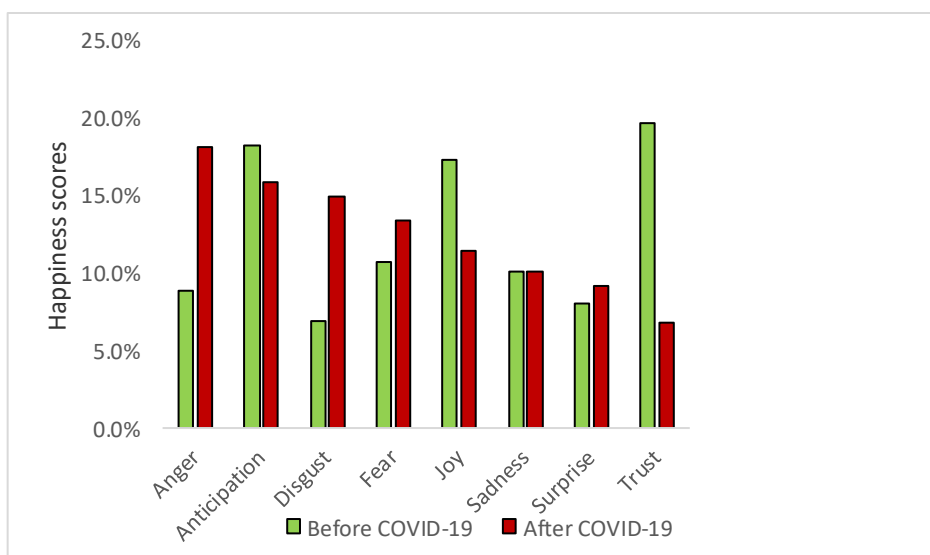
Notes: "Lockdown day" is the first day that regulations came into force for people to isolate socially in different countries. The "average" indicates the average happiness levels in the respective countries before COVID-19.

Figure 7: Happiness levels before and after lockdown regulations.

### 6.2.2.3 Emotion analysis

The team started analyzing the different emotions, namely *joy, anticipation, trust, disgust, anger, surprise, fear and sadness*, as mentioned previously. The extended research commenced just in time to study the changes in emotions before and after the global COVID-19 pandemic of 2020, a once in a 100-years' event. The team found that there were marked differences in the emotions of people before and after the regulations to curb COVID-19 were introduced. The main regulations were to enforce social isolation, although the stringency of these restrictions varied between countries. For example, in Australia, that had the most lenient regulations of the countries, people were allowed to go about their normal business, but on March 29 social gatherings were restricted to only two people. In New Zealand, the lockdown and social isolation were for four weeks from March 26 to April 23; however, New Zealanders were allowed out of their houses to buy essential goods and to exercise. In South Africa, the lockdown regulations were very severe and enforced by the police and the military. People were not allowed to leave their homes; the only time they could leave their houses was to purchase essential goods. In all three countries, there was a marked change in the emotions of the people from before to after the lockdown (see figure 7).

The emotions most noted across all three countries before COVID-19 were joy, trust and anticipation, but after the announcement of the increased number of COVID-19 victims and regulations to curb the spread of the virus, the emotions changed drastically, though differing in levels from one country to the other. The prominent emotions became fear, anger and distrust. The most severe changes in emotions were noticed in South Africa that had the strictest lockdown regulations. Anger increased significantly, whereas joy and trust decreased significantly, see figure 8.



Source: Greyling et al. (2019).

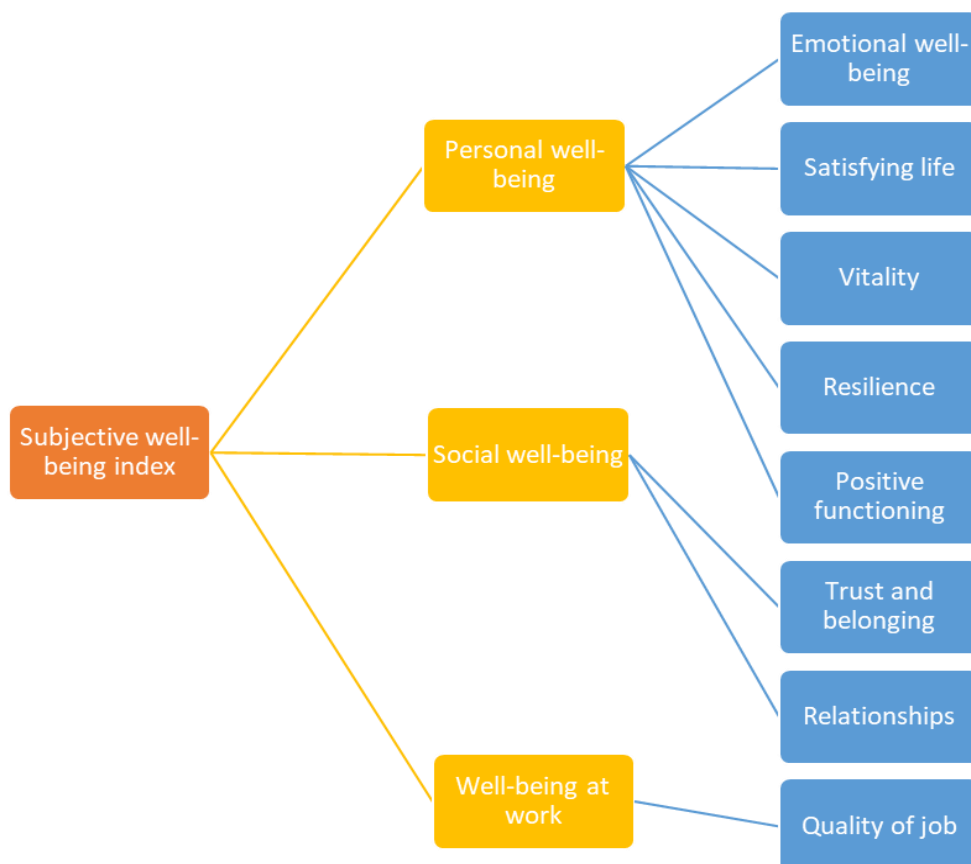
Figure 8: Emotions of South Africans before and after the enforcement of lockdown regulations.

### 6.2.3 Limitations of the GNH index

Although the GNH Index addressed some of the limitations mentioned regarding the Hedonometer, such as considering the context in which a word is used in a tweet, analyzing tweets in different languages and determining the emotions embedded in tweets, it still suffers from the problem of the representativeness of the Twitter data, as discussed before.

### 6.3. A composite index of subjective well-being (Italy and Japan)

Since 2012, Ceron, Curini and Iacus (2016) used an Integrated Sentiment Analysis (a human supervised machine learning method) on Big Data extracted from Twitter, for both Italy and Japan, to obtain a composite index of subjective and perceived well-being that captures various aspects and dimensions of individual and collective life (Iacus et al. 2020). Inspired by the Happy Planet Index (Jeffrey et al. 2016) and keeping to the philosophy that subjective well-being is multidimensional, their Subjective Well-being Index (SWBI) consists of eight dimensions. These dimensions encapsulate three different well-being areas: personal well-being, social well-being and well-being at work (figure 9).



Source: Adapted from Iacus et al. (2020).

Figure 9: Dimensions of SWBI

### 6.3.1 Methodology

Up until 2017, the researchers extracted and classified 240 million tweets over 24 quarters. To analyze the sentiment, they applied a new humanly supervised sentiment analysis, different from what was explained in section 3.2, and did not rely on lexicons or special semantic rules.

As discussed in section 1.6, one of the biggest criticisms against using social media data such as Twitter, is that it is a biased, non-uniform subsample and not representative of the overall population. To overcome this sampling bias, they use a weighted method developed by Iacus et al. (forthcoming). In essence, the researchers use a hierarchical aggregation, where Italian provincial-level data, weighted by the characteristics of provincial macro-variables, are used to estimate the composition of sentiment throughout regional society. Thus, they integrate Twitter data with official or administrative data.

For each indicator, they apply  $\hat{y}_{dt}^w$  as the regional sampling mean, where the regional units are the weighted means of province-level units, to overcome the non-uniform sampling structure of the data:

$$\hat{y}_{dt}^w = \frac{1}{\sum_{i=1}^{n_{dt}} w_{idt}} \sum_{i=1}^{n_{dt}} y_{iat} w_{iat} \quad (1)$$

where  $n_{dt}$  is the number of provinces in region  $d$  at time  $t$ , and  $w_{idt}$  are the weights used.

They calculate the weights by making use of the Twitter rate and internet coverage. The Twitter rate is the ratio between the number of tweets and the population size in the region ( $d$ ), analyzed in time  $t$ . Iacus et al. (forthcoming) calls  $w_{1;idt}$  the Twitter rate and  $w_{2;idt}$  the broadband coverage, to apply to the weighting procedure for  $\hat{y}_{dt}^w$  in equation (1) and argues that this is a good proxy of the use of Twitter for Italians.

### 6.3.2 Results

As argued by the researchers, traditional well-being scores, such as that reported in the World Happiness Report, which are usually considered on an annual basis, have a situation of temporal stationarity, which is not realistic in current life. What they uncovered, after considering the data for Italy's 19 regions, was the important role that time differences play. Additionally, they were able to show that their social well-being index components of well-being at work, having a satisfying life and vitality correlate with macroeconomic level traditional survey data.

### 6.3.3 Limitations

The index, although addressing some of the limitations of the other happiness indices using Twitter, especially the representativeness of Twitter data, has a limited-time series from January 2012 to December 2017; thus, it is not

available for analyses of current events such as the global pandemic. Furthermore, it does not analyze the emotions underpinning tweets and only considers tweets made in Italian and Japanese, thereby neglecting to analyze tweets made in English.

## 7. Summary

This paper presented a broad overview of happiness, Big Data (organic data) and the current research utilizing these two evolving concepts. As was highlighted throughout, psychologists, economists and sociologists (among others) have learned a significant amount about human behaviour that they could not have fathomed before the era of Big Data.

Alas, it would be considered biased if there is no word of caution about Big Data. With all of the 'good' aspects of Big Data, there are also a few 'bad' aspects that need to be considered. Firstly, just like the gravitational forces of the sun or moon that create tidal waves, there is no escaping the gravitational forces that create the Big Data tidal wave. It will have a permanent impact on society at large, and everyone will need to adapt. Secondly, although social media platforms provide organic data that allow researchers to investigate behaviour, there are serious concerns about the amount of time people spend on social media. Psychologists and social scientists have expressed fears that people are losing the ability to communicate face-to-face and therefore becoming antisocial. Additionally, social media platforms provide a stage for individuals to utter hate speech against specific ethnicities, genders or countries, unchecked. Finally, the debate regarding privacy and individuals' rights will only get more heated as time goes by. For some, Big Data could equate to Big Brother.

One thing is certain, for governments and policymakers to make better-informed decisions regarding the human well-being of their people, they need better and the most up to date data available. Big Data allows them the opportunity to not blindly follow prior theories about what is 'good' for their people, but to actually hear from the people what they think is 'good' for them.

*"It is a capital mistake to theorize in advance of the facts. Insensibly one begins to twist facts to suit theories instead of theories to suit facts".*

*– Sherlock Holmes*



## References

- 99firms. (2020). Agencies and software you can trust. Available at <https://99firms.com/> Accessed on March 2 2020.
- Algan, Y., Beasley, E., Guyot, F., Higad, K., Murtin, F., & Senik, C. (2015). Big Data Measures of Well-Being: Evidence from a Google Well-Being Index in the U.S. SciencesPo Working Paper 2015.
- Alharthi, R., Guthier, B., Guertin, C., & El Saddik, A. (2017). A dataset for psychological human needs detection from social networks. *IEEE Access*, 5: 9109-9117.
- Asur, S. & Huberman, B. A. (2010). Predicting the Future with Social Media. *2010 IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent Technology*, Toronto, ON, 2010, pp. 492-499.
- Bellet, C. & Frijters, P. (2019). Big data and well-being. *In World Happiness Report 2019*, chapter 6. Edited by John F. Helliwell, Richard Layard, and Jeffrey D. Sachs. Associate Editors, Jan-Emmanuel De Neve, Haifang Huang, Shun Wang, and Lara B. Aknin.
- Bollen, J., Mao, H. & Zeng, X. (2011). Twitter mood predicts the stock market. *Journal of Computational Science*, 2(1): 1–8.
- Broadstock, D. & Zhang, D. (2019). Social-media and intraday stock returns: The pricing power of sentiment. *Finance Research Letters*, 30: 116-123.
- Callegaro, M. & Yang, Y. (2018). The Role of Surveys in the Era of "Big Data. *In Vannette D., Krosnick J. (eds) The Palgrave Handbook of Survey Research*. Palgrave Macmillan, Cham.
- Ceron, A., Curini, L., & Iacus, S. M. (2016). iSA: A fast, scalable and accurate algorithm for sentiment analysis of social media content. *Information Sciences*, 367: 105–124.
- Culotta, A. (2010). Towards detecting influenza epidemics by analyzing Twitter messages. *In KDD Workshop on Social Media Analytics*.
- David, N. (2020). Social Media Statistics, Australia – January 2020. Available at <https://www.socialmedianews.com.au/social-media-statistics-australia-january-2020/> Accessed on April 23 2020.
- Dawkins, R. (2006). *The God Delusion*. Black Swan.
- De Choudhury, M., Monroy-Hernandez, A., & Mark, G. (2014). Narco emotions: Affect and desensitization in social media during the Mexican drug war. *In Proceedings of the SIGCHI Conference on Human Factors in Computing Systems* (pp. 3563–3572). New York, NY: ACM.

- Deaton, A. (2013). *The Great Escape: Health, Wealth and the Origins of Inequality*. Princeton University Press.
- Diener, E., Lucas, R. E., Schimmack, U. & Helliwell, J. F. (2009). *Well-being for Public Policy*. Published by Oxford University Press.
- Diener, E. & Suh, E. (1997). Measuring quality of life: economic, social, and subjective indicators. *Social Indicators Research*, 40(1-2): 189-216.
- Diener, E., Tay, L., & Myers, D. G. (2011). The religion paradox: If religion makes people happy, why are so many dropping out? *Journal of Personality and Social Psychology*, 101(6): 1278–1290.
- Dodds, P. S. & Danforth, C. M. (2010). Measuring the Happiness of Large-Scale Written Expression: Songs, Blogs, and Presidents. *Journal of Happiness Studies*, 11: 441–456.
- Dodds, P. S., Harris, K. D., Kloumann, I. M., Bliss, C. A. & Danforth, C. M. (2011). Temporal patterns of happiness and information in a global social network: Hedonometrics and twitter. *PLoS ONE*. DOI:10.1371/journal.pone.0026752.
- Ford, M.T., Jebb, A. T., Tay, L. & Diener, E. (2018). Internet Searches for Affect-Related Terms: An Indicator of Subjective Well-Being and Predictor of Health Outcomes across U.S. States and Metro Areas. *Applied Psychology: Health and Well-being*, 10(1): 3–29.
- Gayo-Avello, D. (2013). A Meta-Analysis of State-of-the-Art Electoral Prediction from Twitter Data. *Social Science Computer Review*, 31(6), 649–679.
- Gerber, M. S. (2014). Predicting crime using Twitter and kernel density estimation. *Decision Support Systems*, 61:115–25.
- Giorgi, S., Yaden, D. B., Eichstaedt, J. C., Ashford, R. D, Buffone, A. E. K., Schwartz, H. A., Ungar, L. H. & Curtis, B. (2020). Cultural Differences in Tweeting about Drinking Across the U.S. *International Journal of Environmental Research and Public Health*, 17(1125).
- Greyling, T., Rossouw, S. & Afstereo. (2019). Gross National Happiness Index. The University of Johannesburg and Afstereo [producers]. Available at <http://gnh.today>. Accessed on April 23 2020.
- Hailong, Z., Wenyan, G., & Bo, J. (2014). Machine learning and lexicon-based methods for sentiment classification: A survey. In *2014 11th Web Information System and Application Conference* (pp. 262-265). IEEE.
- Haybron, D. M. (2013). *The Nature and Significance of Happiness*. Oxford Handbook of Happiness, edited by Ilona Boniwell, Susan A. David, and Amanda Conley Ayers.

- Hedonometer. (2020). Average Happiness for Twitter Time-series. Vermont Complex Systems Center. Computational Story Lab. Available at [https://hedonometer.org/timeseries/en\\_all/](https://hedonometer.org/timeseries/en_all/). Accessed on April 23 2020.
- Helliwell, J. F., Layard, R., Sachs, J. & De Neve, J. (eds.). (2020). World Happiness Report 2020. New York: Sustainable Development Solutions Network. ISBN 978-1-7348080-0-1.
- Hills, T. T., Proto, E., Sgroi, D., & Seresinhe, C. I. (2019). Historical analysis of national subjective well-being using millions of digitized books. *Nature human behaviour*, 1-5: 1271-1275.
- Iacus, S. M., Porro, G., Salini, S. & Siletti, E. (2020). An Italian Composite Subjective Well-Being Index: The Voice of Twitter Users from 2012 to 2017. *Social Indicators Research*. <https://doi.org/10.1007/s11205-020-02319-6>
- Iacus, S. M., Porro, G., Salini, S., & Siletti, E. (forthcoming). Controlling for selection bias in social media indicators through official statistics: A proposal. *Journal of Official Statistics*.
- Inglehart, R., C. Haerpfer, A. Moreno, C. Welzel, K. Kizilova, J. Diez-Medrano, M. Lagos, P. Norris, E. Ponarin & B. Puranen (eds.). (2014). World Values Survey: All Rounds - Country-Pooled Datafile Version: <http://www.worldvaluessurvey.org/WVSDocumentationWVL.jsp>. Madrid: J.D. Systems Institute.
- Jeffrey, K., Wheatley, H., Abdallah, S. (2016). *The Happy Planet Index: 2016. A global index of sustainable well-being*. London: New Economics Foundation. Available at <http://happyplanetindex.org/about>
- Jones, N. M., Wojcik, S. P., Sweeting, J., & Silver, R. C. (2016). Tweeting negative emotion: An investigation of Twitter data in the aftermath of violence on college campuses. *Psychological Methods*, 21: 526–541.
- Kahneman, D. (2011). *Thinking, fast and slow*. Farrar, Straus and Giroux Publishers.
- Kahneman D. & Deaton, A. (2010). High income improves evaluation of life but not emotional wellbeing, *Proceedings of the national academy of sciences*, 107(38): 16489-16493.
- Kalev H. L., Wang, S., Cao, G., Padmanabhan, A. & Shook, E. (2013). Mapping the global Twitter heartbeat: The geography of Twitter. *First Monday*, 18(5-6). Available at <https://firstmonday.org/ojs/index.php/fm/article/download/4366/3654> Accessed on April 23 2020.
- Lamos, V., & Cristianini, N. (2010). Tracking the flu pandemic by monitoring the social web. In IAPR 2nd Workshop on Cognitive Information Processing (CIP 2010).
- Lee, E., Lee, J. A., Moon, J.H. & Sung, Y. (2015). Pictures Speak Louder than Words: Motivations for Using Instagram. *Cyberpsychology, Behaviour and Social Networking*, 18(9): 552-556.

- McKay, D. (2018). *Investigating the Effect of Sentiment in High-Frequency Financial Markets*, Dublin: The University of Dublin.
- Mitchell, L., Frank, M. R., Harris, K. D., Dodds, P. S. & Danforth, C. M. (2013). The Geography of Happiness: Connecting Twitter sentiment and expression, demographics, and objective characteristics of place. *PLoS ONE*, 8(5): e6441.
- O'Connor, B., Balasubramanian, R., Routledge, B. R. & Smith, N. A. (2010). From Tweets to Polls: Linking Text Sentiment to Public Opinion Time Series. Proceedings of the Fourth International Conference on Weblogs and Social Media, ICWSM 2010, Washington, DC, USA, May 23-26, 2010. Available from <https://dblp.uni-trier.de/db/conf/icwsm/icwsm2010.html>
- Omnicores. (2020). Omnicores Agency. Available at <https://www.omnicoreagency.com/> Accessed on February 26 2020.
- Paul, M. J. & Dredze, M. (2011). You are what you tweet: Analysing Twitter for public health. In Proceedings of the Fifth International AAAI Conference on Weblogs and Social Media. Available from <http://www.aaai.org>.
- Rao, T. & Srivastava, S. (2012). Analyzing stock market movements using Twitter sentiment analysis. In Proceedings of the 2012 international conference on advances in social networks analysis and mining (ASONAM 2012) (pp. 119-123). IEEE Computer Society.
- Renault, T. (2017). Intraday online investor sentiment and return patterns in the U.S. stock market. *Journal of Banking & Finance*, 84: 25-40.
- Richardson, A., van Florenstein Mulder, T. & Vehbi, T. (2018). Nowcasting New Zealand GDP Using Machine Learning Algorithms. Centre for Applied Macroeconomic Analysis, Working Paper 47/2018.
- Rill, S., Reinel, D., Scheidt, J., & Zicari, R. V. (2014). Politwi: Early detection of emerging political topics on twitter and the impact on concept-level sentiment analysis. *Knowledge-Based Systems*, 69: 24–33.
- Riotta, G., Riccaboni, M., Pammolli, F., Caldarelli, G., Chessa, A. & Puliga, M. (2014). A Multi-Level Geographical Study of Italian Political Elections from Twitter Data. *PLoS ONE*, 9(5): e95809.
- Ritter, R. S., Preston, J. L. & Hernandez, I. (2014). Happy tweets: Christians are happier, more socially connected, and less analytical than atheists on Twitter. *Social Psychological & Personality Science*, 5: 243–249.
- Ryff, C. D. & Keyes, C. L. M. (1995). The structure of psychological well-being revisited. *Journal of Personality and Social Psychology*, 69(4): 719-727.

Saranghi, S. & Sharma, P. (2020). *Big Data: A Beginner's Introduction*. Taylor & Francis Group. ProQuest Ebook Central.

Schwartz, H. A., Eichstaedt, J. C., Kern, M. L., Dziurzynski, L., Agrawal, M., Park, G. J., Lakshmikanth, S. K., Martin, S., Seligman, E. P., Ungar, L. & Lucas, R. E. (2013). Characterizing Geographic Variation in Well-Being Using Tweets. Proceedings of the Seventh International AAAI Conference on Weblogs and Social Media.

Schwartz, H. A., Sap, M., Kern, M. L., Eichstaedt, J. C., Kapelner, A., Agrawal, M., Blanco, E., Dziurzynski, L., Park, G. & Stillwell, D. (2016). Predicting individual well-being through the language of social media. In *Biocomputing 2016: Proceedings of the Pacific Symposium*, pp. 516–527. World Scientific.

Seligman, M. E. P. (2011). *Flourish: A Visionary New Understanding of Happiness and Well-being*. Free Press.

Steptoe, A., Deaton, A. & Stone, A. A. (2014). Subjective Wellbeing, Health, and Ageing. *The Lancet*, 385(9968): 640-648.

Steyn, D. H. W., Greyling, T., Rossouw, S. & Mwamba, J. M. (2020). Sentiment, emotions and stock market predictability in developed and emerging markets. *GLO Discussion Paper*, No. 502, Global Labor Organization (GLO), Essen.

Stiglitz, J. E., Sen, A. & Fitoussi, J. (2009). *Report by the Commission on the Measurement of Economic Performance and Social Progress*. Available at [https://www.economie.gouv.fr/files/finances/presse/dossiers\\_de\\_presse/090914mesure\\_perf\\_eco\\_progres\\_social/synthese\\_ang.pdf](https://www.economie.gouv.fr/files/finances/presse/dossiers_de_presse/090914mesure_perf_eco_progres_social/synthese_ang.pdf)

Sumner, A. (2003). Economic and non-economic well-being: a review of progress on the meaning and measurement of poverty. Helsinki, Finland. Paper prepared for WIDER international conference on inequality, poverty and human well-being, May 30-31, 2003.

Tabari, N., Praneeth, B., Biswas, P., Seyeditabari, A., Hadzikadic, M. & Zadrozny, W. (2018). Causality Analysis of Twitter Sentiments and Stock Market Returns. Proceedings of the First Workshop on Economics and Natural Language Processing pages 11–19 Melbourne, Australia, July 20, 2018.

Tumasjan, A., Sprenger, T. O., Sandner, P. G. & Welpe, I. M. (2010). Predicting elections with twitter: What 140 characters reveal about political sentiment. In Proc. Int. AAAI Conf. Web Social Media (ICWSM), 10(1): 178–185.

United States (2020). Bureau of Foreign and Domestic Commerce. Seventy-Third Congress, 1933-1935 and Kuznets, Simon, 1901-1985. *National Income, 1929-1932*, Washington: U.S. Government Printing Office, 1934, <https://fraser.stlouisfed.org/title/971> Accessed on March 16 2020.

Veenhoven, R. (2009). How do we assess how happy we are? *In* Dutt, A. K. & Radcliff, B. (eds.) 'Happiness, Economics and Politics: Towards a multidisciplinary approach', Edward Elger Publishers, Cheltenham U.K., ISBN 978 1 84844 093 7, Chapter 3, pp. 45-69.

Verlet, D. & Devos, C. (2009) The Main Determinants for Subjective Well-Being: A Quest for the Holy Grail? *In*: Møller V., Huschka D. (eds) Quality of Life and the Millennium Challenge. *Social Indicators Research Series*, 35(4): 193–219. Springer.

Villani, D., Sorgente, A., Paola, I. & Alessandro, A. (2019). The Role of Spirituality and Religiosity in Subjective Well-Being of Individuals With Different Religious Status. *Frontiers in Psychology*, 10: 1525.

Wang, X., Gerber, M. S. & Brown, D. E. (2012). Automatic crime prediction using events extracted from twitter posts. *In* Proceedings for International Conference on Social Computing, Behavioral-Cultural Modeling Predictions, September, pp.231–238.

Ward, G. (2019). Happiness and voting behaviour. *In* World Happiness Report 2019, chapter 3. Edited by John F. Helliwell, Richard Layard, and Jeffrey D. Sachs. Associate Editors, Jan-Emmanuel De Neve, Haifang Huang, Shun Wang, and Lara B. Aknin.

Weber, I., Venkata, R., Garimella, K. & Batayneh, A. (2013). Secular vs Islamist polarisation in Egypt on Twitter. *In* Proceedings of the 2013 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining, 290–97. New York, NY: ACM.

Wilson, R. E., Gosling, S. D., & Graham, L. T. (2012). A review of Facebook research in the social sciences. *Perspectives on Psychological Science*, 7(3): 203-220.

Wolff, R. (n.d). Quick introduction to sentiment analysis. Available at <https://towardsdatascience.com/quick-introduction-to-sentiment-analysis-74bd3dfb536c> Accessed on February 12 2020.

Zivanovic, S., Martinez, J. & Verplanke, J. (2020). Capturing and mapping quality of life using Twitter data. *GeoJournal* 85: 237–255.