

From twitter to GDP: Estimating economic activity from social media[☆]

Agustín Indaco



Carnegie Mellon University in Qatar, Education City, P.O. Box 24866, Doha, Qatar

ARTICLE INFO

JEL classification:

C53
C55
E01
O11

Keywords:

Social media data
Big data
Cities
Satellite images
National accounts

ABSTRACT

Using all geo-located image tweets shared on Twitter in 2012–2013, I find that the volume of tweets is a valid proxy for estimating GDP at the country level, explaining 78 percent of cross-country variations. I also exploit the geographic granularity of social media posts to estimate and predict GDP at the sub-national level. I find that tweets alone can explain 52 percent of the variation in GDP across cities in the US. Estimates using Twitter data perform on par with the more common night-lights proxy. Furthermore, both indicators seem to capture different aspects of economic activity and thus complement each other.

1. Introduction

Despite incessant debate about its ability to accurately measure the state of the economy, the gross domestic product (GDP) is still the most widely used indicator to gauge countries' economic performance (Masood (2014)). One of the many difficulties with estimating GDP is that its measurement is often complicated and expensive. This could lead to measurement error, particularly in developing countries, which in turn mislead policy makers and businesses. Another concern is that given the importance surrounding official GDP estimates both in terms of market fluctuations as well as public perception of politicians' performances, governments can find short-term benefits in manipulating these estimates. Additional shortfalls of GDP estimates are the delay with which they are released and their geographic rigidity. In light of this, much research has been focused on alternative ways of measuring GDP other than the traditional sample survey method, both to corroborate as well as a control mechanism.

In this paper, I argue for the use of data from social media posts as a proxy for measuring GDP. By locating and analyzing the volume of hundreds of million social media posts, I show that one can accurately estimate GDP at the country level. For this exercise, I collect all geo-located

image tweets shared on Twitter for the years 2012 and 2013. In total the dataset includes roughly 140 million tweets, with their corresponding text, time of posting, precise latitude and longitude and unique user identifier. I then aggregate the number of tweets sent from each country in each year to estimate annual GDP at the country level. I then show that the predictive power of tweets compares well with satellite night-light images, which are widely used as a proxy for economic activity. Furthermore, I show that ideally, both proxies should be included together in order to substantially increase the models' predictive power.

This paper has five main findings. First, I show that the volume of tweets can be used as a proxy for estimating current GDP at the country level; my preferred model can explain 87 percent of the cross-country variation in GDP. Second, I compare the strength of social media as a proxy for estimating GDP relative to the more commonly used satellite night-lights data. I find that Twitter data explains slightly more of the variance in cross-countries' GDP estimates. More importantly, I find that when both proxies are included, the predictive power increases and both coefficients remain statistically significant, indicating that they may not be capturing the same aspects of economic activity. Several different exercises in this paper show that the two proxies ideally should be used together for more accurate predictions. Third, I find that social

[☆] I am grateful to David Jaeger, Francesc Ortega, and Lev Manovich for their useful comments and suggestions. I also thank Tahir Butt and Agustín Mario for immensely helpful conversations, and Damon Crockett for putting together the maps that are shown in Figs. 1–5. This paper has also benefited from helpful comments from David Weil and seminar participants at The Federal Reserve Board, European Meeting of the Urban Economics Association, Queens College, Carnegie Mellon University in Qatar, Conference on Advanced Research Methods and Analytics and CUNY Graduate Center.

E-mail address: aindaco@andrew.cmu.edu.

<https://doi.org/10.1016/j.regsciurbeco.2020.103591>

Received 21 December 2019; Received in revised form 24 August 2020; Accepted 26 August 2020

Available online 15 September 2020

0166-0462/© 2020 The Author. Published by Elsevier B.V. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

media data can be used to estimate annual variation in GDP at the country level. Despite not being very precise, the relationship between annual changes in social media posts and changes in economic activity are informative. Fourth, I present and study a hypothesis of the underlying relationship between image tweets and economic activity in which image tweets are a byproduct of consumption. It seems that social media is used as a medium to share conspicuous consumption among its users. Finally, I exploit the continuous geographic granularity aspect of social media posts to study the local economic effect of a large oil reserve finding in an underdeveloped region in Argentina. Despite widespread enthusiasm regarding the local economic effect the discovery would have, no official estimates had been reported: for the first time, I estimate a 38 percent annual increase in economic activity in the region.

Although this paper focuses solely on data gathered from Twitter, this is just one of several social media platforms that could potentially serve as a proxy for different economic measurements. A priori, it is hard to assess which platform is better suited for a particular indicator. Each social media platform has its own purpose, application and user demographics that needs to be carefully considered before making any claims. This paper does not propose that image tweets are useful as a proxy for estimating all economic variables. It may also be argued that there are other social media platforms that are better suited to serve as a proxy for GDP. The goal of this paper is to show an example of social media data being used for such purpose. Throughout this paper, I detail the advantages and shortfalls of social media data for this purpose, with a focus on image tweets. I also carry out some exercises that explore the boundaries and capabilities of such tools for measuring economic activity in different instances.

Social media can contribute greatly to economic research. Measures taken from social media, as well as from other alternative data sources, can serve both as substitutes as well as complements to traditional survey data. In particular, social media data has several properties that are beneficial when estimating economic measures. First, social media data is available in real time, which allows to nowcast economic activity and thus provide companies and individuals updated information when making economic decisions.¹ Second, given that geo-tagged social media posts can be geographically assigned to a precise location within approximately a 10 m radius, one can aggregate social media posts at any sub-national geographical level one deems interesting. This includes aggregating data between areas that are not bound together inside political borders and thus fabricate meaningful areas of study that are not possible with official datasets. Third, unlike survey data that are costly to recollect, social media data is organically being generated by users from all over the world and available to data agencies, international organizations and non-governmental organizations at a relatively low cost. Finally, even though estimates in this paper are based solely on the volume of tweets from each location for a given period of time, more information can be potentially extracted from the content of each post. Topic modeling, sentiment analysis, use of language, spelling, share of posts coming from locals and tourists and other information can be extracted from posts to improve estimates.

There are several papers which explore the use of social media data to measure economic outcomes. [Antenucci et al. \(2014\)](#) analyze tweets to build a model that accurately predicts unemployment insurance claims in the US. In a similar manner, [Llorente et al. \(2015\)](#) use Twitter data to estimate regional unemployment rates in Spain. [Glaeser et al. \(2017\)](#) use Yelp data to predict changes in the number of overall establishments and restaurants across zip codes in the US. As far as I

¹ A few papers have looked at proxies to nowcast economic activity at the local level. [Askitas and Zimmermann \(2013\)](#) estimate the German business cycle at a monthly level by measuring toll activity on important highways by heavy transport vehicles. Whereas [Glaeser et al. \(2017\)](#) use data from Yelp to estimate business establishment entry and exit across zip codes in the US.

am aware, this is the first paper to explore the use of social media data to estimate economic activity across the globe.

This paper also contributes to a growing literature of papers that use novel datasets (or tools) to study topics in urban economics using micro-geographic spatial data. For example, [Henderson et al. \(2019\)](#) use LandScan data to estimate urban agglomeration effects to explain household income and wage differences across cities. [Davis et al. \(2017\)](#) on their part use granular property-level data to study the divergence of land and house values during boom-bust cycles. [Dauth and Haller \(2020\)](#) use geo-referenced information to study the relation between commuting distances and wages and whether individuals are loss-averse to changes in commuting times. Other papers have focused on developing tools better suited to work with these granular geo-coded data. For example, [Åvald Sommervoll and Sommervoll \(2019\)](#) use machine learning to circumvent the trade-off in the level of aggregation when using fixed effects models. They use a genetic algorithm to identify areas with similar location premiums which can be used to spatially aggregate areas that are both refined and robust.

2. Estimating GDP (and its problems)

GDP is the most widely used measure of a country's economic performance. Given its importance, national income measurements are governed by the United Nations System of National Accounts (SNA), which sets a global standard that allows for international comparisons of economic activity across countries. However, adherence to these standards is entirely voluntary, and cannot be rigidly enforced. While some countries adhere to the most recent standards set by the revised SNA 2008, many developing countries have still not adopted the previous 1993 SNA standards and are still using 1968 SNA methodologies.²

As an illustration of the degree of measurement error, [Johnson et al. \(2013\)](#) study the revision in the Penn World Tables (PWT), a popular dataset frequently used for making comparisons of GDP across countries. The authors find that the standard deviation of the change in countries' average growth over the period 1970–1999 was 1.1 percent per year when comparing data in version 6.1 to version 6.2.³ Given that the average growth rate is 1.56 percent, this is a relatively large discrepancy. In fact, [Dawson et al. \(2001\)](#) claim that some results in the economic literature based on PWT are purely a product of measurement error in the data.⁴

While these examples show the complications involved when measuring GDP, these problems tend to be accentuated in developing countries where statistical offices tend to have fewer resources for constructing this immensely intricate task. [Jerven \(2013\)](#) explains that in African countries the informal sector is so large that it cannot be left out of GDP estimates, and thus a variety of innovative accounting practices are implemented to take them into consideration. This leads to a series of pragmatic decisions to be made within statistical offices, which are subject to the availability of trustworthy data, financial resources and political instructions. In turn, [Jerven \(2013\)](#) believes that GDP statistics from African countries are “best guesses of aggregate production”.

A second concern is that given the importance surrounding official GDP estimates, both in terms of market fluctuations as well as public perception of politicians' performances, governments can find short-term benefits in manipulating these estimates. In less developed countries, higher growth estimates of real GDP per-capita which have

² According to [Jerven \(2013\)](#), as of 2010, in Africa only Cameroon and Lesotho were upgrading their system to incorporate SNA 2008 standards, while three African countries still used the 1968 SNA methodologies.

³ Version 6.1 of the PWT was released in 2002, while version 6.2 was released in 2006.

⁴ In particular, they find that the empirical link between output and volatility and the cross-country test in the Permanent Income Hypothesis are driven by measurement error in the data.

later been revised downward have been associated with a higher probability of reelection ([Brender and Drazen \(2005\)](#)). Even though as we have seen GDP estimates are subject to measurement error, [Kerner and Crabtree \(2018\)](#) find that there are non-random variations in official macroeconomic estimates. Furthermore, given that GDP estimates are produced by countries' statistics agencies via surveys that are not publicly available, it is difficult for non-governmental organizations, international organizations and the public at large to corroborate these estimates.

Official GDP estimates also have limitations in terms of lags and geographical aggregation. At best, advanced GDP estimates are available 30–45 days after the end of the reference quarter. After this release, revised estimates which incorporate information from additional surveys are disclosed up to three years later. Finally, GDP estimates are geographically bound to political and geographical borders. Even in the case of many developed countries which produce estimates for metropolitan areas, these estimates are geographically rigid and cannot be combined across country borders. For many developing countries, there are no sub-national estimates of economic activity.

These concerns and limitations surrounding official GDP estimates have motivated efforts in finding proxies that may be able to estimate economic activity. The bulk of this literature has centered around efforts to use satellite night-light images to estimate GDP at the country level ([Donaldson and Storeygard \(2016\)](#)). I contrast and compare Twitter and night-lights as proxies for GDP in section 5. This paper adds to this literature by exploring whether social media can be used as a proxy for GDP both at the country level as well as for sub-national regions. On the one hand, these proxies could be used by countries' statistic agencies as supplements that improve the accuracy of their estimates. On the other hand, these proxies could also be used as a tool for non-governmental agencies and international organizations to corroborate official GDP estimates.

3. Twitter and twitter data

Twitter is a social media application which allows users to post short messages of any subject of their choosing. These messages are known as *tweets*. Twitter started in 2006 and by 2012 had 140 million global users and 340 million tweets per day. Initially, tweets were limited to 140 characters, because Twitter was originally designed as an SMS mobile phone-based platform. In its early days, 140 characters were the limit that mobile carriers imposed with the SMS protocol standard so Twitter was simply creatively constrained. As Twitter eventually grew into a web platform, the 140-character limit remained as a matter of branding.⁵ Unless restricted by the user, tweets are publicly available and can be read via the application or on a web browser.

Twitter initially did not allow users to share images, videos or other sorts of media in their tweets. This changed in August 2011, when Twitter rolled out a platform that allowed users to add images to their tweets. Until September 2013, image tweets did not present the image, but instead included a link where your image could be viewed. In September 2013, images could be previewed directly on the tweet.

Detailed data from tweets sent by public users are offered directly from the company at a cost. Although this data is not publicly available at scale⁶ it is available for countries' data agencies, international organizations and non-governmental organizations at a relatively low cost.

The dataset used in this paper contains all geo-tagged image tweets

posted on Twitter between January 1, 2012 and December 31, 2013.⁷ This dataset was provided directly by Twitter, through a Twitter Data Grant submission in 2014.⁸ The total dataset contains 140 million tweets from all around the world. Each tweet contains information on: i) a unique identifier for each individual Twitter user; ii) the latitude and longitude (with 5 decimal points, which gives a precision of 1.1 m) of where the tweet was sent from⁹; iii) the date and time in which the tweet was sent; iv) the image tweeted; and v) any accompanying text¹⁰.

The underlying mechanism which relates the volume of tweets and economic activity is not clear, but there may be several factors at play. On one hand, part of the cross-country differences seems to be explained by smart-phone penetration, which is in essence an indicator of a economic development. Even though one can post tweets from a desktop or other technological devices, it is safe to assume that the large majority of tweets are sent from smart-phones.

For within country differences, particularly among developed countries where smart-phone penetration is saturated, the underlying mechanism is different. Even though there are several reasons why social media users post on social networks, a substantial amount of posts seem to display the consumption of luxury goods and services (i.e.: food from a renowned restaurant, the latest trend in shoes, exclusive sporting events). In a way, social media has become the ideal medium to exhibit what [Veblen \(1899\)](#) coined as conspicuous consumption: when the utility consumers attain stems more from revealing their wealth and income to others, rather than from the direct utility derived from the good itself. To the conspicuous consumer, the public display of wealth or income is a means to attaining or sustaining a given social status. In order to preserve this status, they must constantly exhibit these consumption patterns online. I examine these hypothesis in section 7.

[Table 1](#) summarizes this Twitter data by year and by country income groups (using the World Bank's classification). We see that the average number of tweets per country was roughly 100,000 in 2012 and increased to 500,000 in 2013. This is an indication of the growth of image tweets over this period. The breakdown of average tweets per income group shows that countries in higher income groups have more tweets, although the growth rates from 2012 to 2013 are larger among lower income countries.

[Table 2](#) provides information on the distribution of tweeted images across the 10 countries with most tweets over this time period. The US tops the list of tweets in 2012 and 2013, with 31.7 percent of the share of total tweets over this period. The US is also the largest country in the world in terms of GDP, with roughly 22.4 percent of the share of global GDP. While some countries, like Japan, have a smaller share of tweets than of global GDP, others, like the UK, have a larger share of tweets than GDP.

⁷ According to [Weidemann and Swift \(2013\)](#) roughly 20 percent of tweets are geo-located. [Lee \(2015\)](#) finds that 42 percent of tweets contain an image, but this study was based on 1 million tweets sent solely from US West Coast users, which likely biases the results. In December 2018, I queried the Twitter API on multiple occasions, collecting a dataset of 10,000 random Tweets. Among this dataset, I find that 4.9 percent of tweets are geo-located and 22.8 percent have images.

⁸ The grant was awarded to the Cultural Analytics Lab, directed by Lev Manovich.

⁹ The latitude and longitude from these posts are generated automatically via very reliable hardware and software and thus I am reasonably certain of their accuracy. Nonetheless, throughout this paper, I hardly need such level of precision, except in the exercise carried out in [subsection 8.1](#).

¹⁰ Given the concern of bots (i.e.: an autonomous program that can interact with computer systems or users) biasing the dataset, I do identify bots as users which at some point sent more than five tweets in a span of 1 min. Removing these users from the data does not significantly change the estimates presented in section 4 and are thus left in our sample.

⁵ In 2017 however, Twitter expanded this limit to 280 characters per tweet.

⁶ Individuals can freely and publicly access data on a relatively small batch of Tweets by querying the Twitter API, but these queries are limited so that users are not able to trove large and representative samples.

Table 1
Twitter data summary statistics: Mean and S.D.

	2012	2013
Tweets	109,678.1 (354,724.1)	528,694.9 (2,397,003.8)
<i>By Income Group</i>		
High (62)	211,993.9 (541,334.7)	1,126,937.3 (4,048,721.5)
Upper-middle (51)	89,123 (201,367.9)	406,004.9 (836,887.4)
Lower-middle (43)	37,394.5 (106,230.9)	209,938.9 (929,686.9)
Low (28)	735.9 (898.9)	3602.5 (4370.9)

Notes: Top row shows the mean number of tweets per country and standard deviation in brackets. The bottom half of the table shows the mean number of tweets and the standard deviation sent by each country divided by income groups, where the number of countries per income group is shown in brackets.

Table 2
Top 10 countries in tweets 2012–2013.

	2012 Tweets	2013 Tweets	2012–2013 Share total Tweets (%)	2012–2013 Share world GDP (%)
Top 10				
U.S.A.	8,260,789	29,376,135	29.3	22.4
U.K.	3,093,930	8,775,427	8.7	3.6
Indonesia	5,44,329	6,401,877	6.4	1.3
Spain	2,121,122	6,210,364	6.2	1.9
Japan	1,825,886	5,128,888	5.5	8.3
Saudi Arabia	1,121,424	5,198,632	3.1	0.9
Brazil	707,881	3,064,005	3.1	3.4
Turkey	630,344	3,112,159	3.0	1.2
Mexico	956,809	3,028,057	2.6	1.6
Russia	110,5047	2,560,921	2.5	2.4

Notes: Top 10 countries in terms of most tweets shared in total between January 2012 and December 2013. Columns 1–2 shows the total number of tweets sent out from each country for each year, Column 4 shows each country's share of total tweets over this time period and Column 5 shows each country's share of world GDP over 2012–2013.

3.1. Visual examples of what twitter data reflect

Fig. 1 shows that there are some clear visual patterns to the location and distribution of tweets worldwide that seem to represent economic activity and population density. The location from where each image tweet was sent is represented by a small light blue point. There are clusters of these light blue points both in areas that are more densely populated as well as areas with higher levels of per capita income. For example, in the United States, the largest concentration of image tweets seem to be centered along the coastal areas, but not so in the less-densely-populated South West and Rocky Mountain States. South America has a cluster of tweets mainly surrounding big cities in Ecuador, Colombia and Venezuela in the north and Brazil, Argentina, Uruguay and Chile further south. In Africa, image tweets tend to be concentrated in richer countries: Morocco, Algeria and Egypt, and in Sub-Saharan Africa in South Africa, Nigeria and Kenya. Western Europe seems to be mostly lit up and the concentration of tweets becomes sparser as we move east into Ukraine, Belarus, Latvia, Estonia and ultimately into Russia. The case of Australia is also telling: tweets are concentrated around large cities off the west coast like Melbourne, Sydney and Canberra, as well as Perth on the east coast.

Fig. 2 shows a more detailed view of East Asia that depicts the clear cutoff in the number of tweets being sent from South Korea and North Korea, respectively. The cluster of points in the north-west part of South Korea corresponds to tweets sent from Seoul. Hardly any tweets are sent

from the other side of the border, which corresponds to North Korea. Internet access is strictly limited in North Korea, and primarily used by the government and foreigners.¹¹

A closer look at image tweets sent from Europe in Fig. 3 lets us see in closer detail how tweets are in fact clustered around capitals and big cities, where population densities and economic activity tend to be higher. We can also depict the roads and highways that connect these large cities. Interestingly, several research papers (Banerjee et al. (2012) and Ghani et al. (2016)) have shown that economic development spreads along network routes; this same pattern seems to be reflected in the location of tweets as well.

Similar patterns emerge in Fig. 4 for the US West Coast. First of all we notice that large cities like Los Angeles and San Francisco are easily visible by the large and extended cluster of tweets that surrounds them. But we can also identify major highways connecting these cities easily depicted.

Moving to the Middle East and North-East Africa, Fig. 5 shows that the number of tweets is smaller and solely concentrated in a few big cities. Interestingly, in Egypt we can see a cluster of tweets flowing down the Nile River: this pattern resembles the concentration of population and economic development surrounding the river.

¹¹ As of 2016, the use of Twitter and several other social media applications has been banned in North Korea.

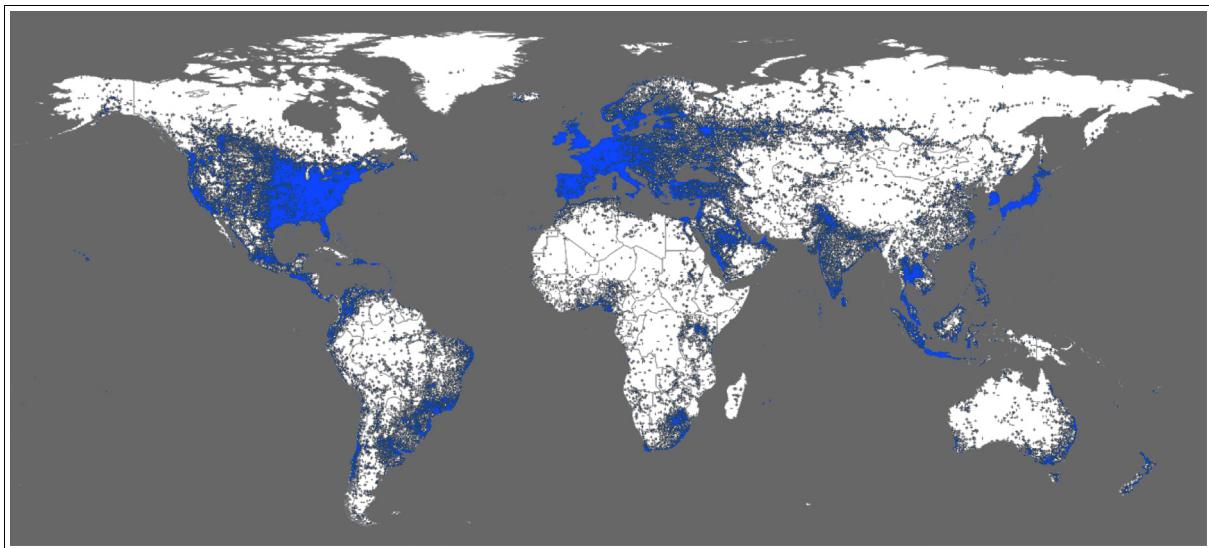


Fig. 1. Worldwide Map of Location of Image Tweets.

Notes: Each light blue dot represents an image tweet sent from that precise location using information on the latitude-longitude. This is a subset of 100 million random tweets from the complete sample for Jan. 2012–Dec. 2013 (For interpretation of the references to colour in this figure legend, the reader is referred to the Web version of this article.)

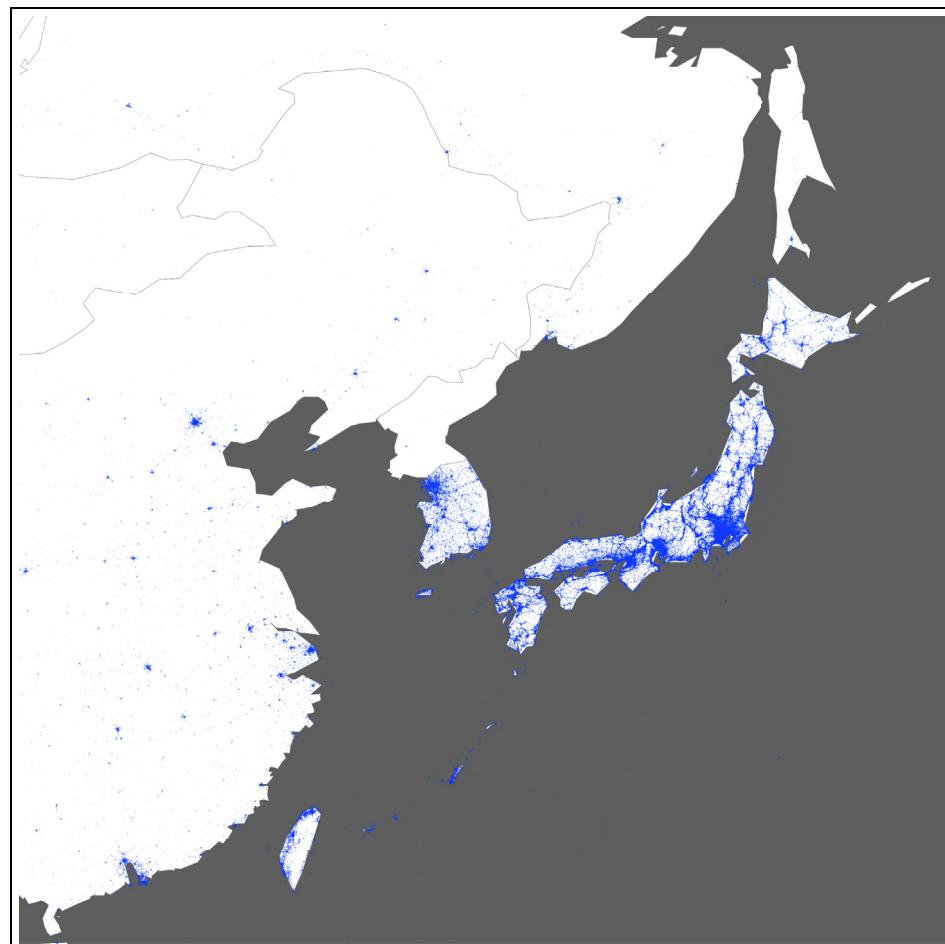


Fig. 2. Map of Image Tweets sent from Japan, South Korea and North Korea.

Notes: Each light blue dot represents an image tweet sent from that precise location using information on the latitude-longitude. This is a subset of 100 million random tweets from the complete sample for Jan. 2012–Dec. 2013. (For interpretation of the references to colour in this figure legend, the reader is referred to the Web version of this article.)

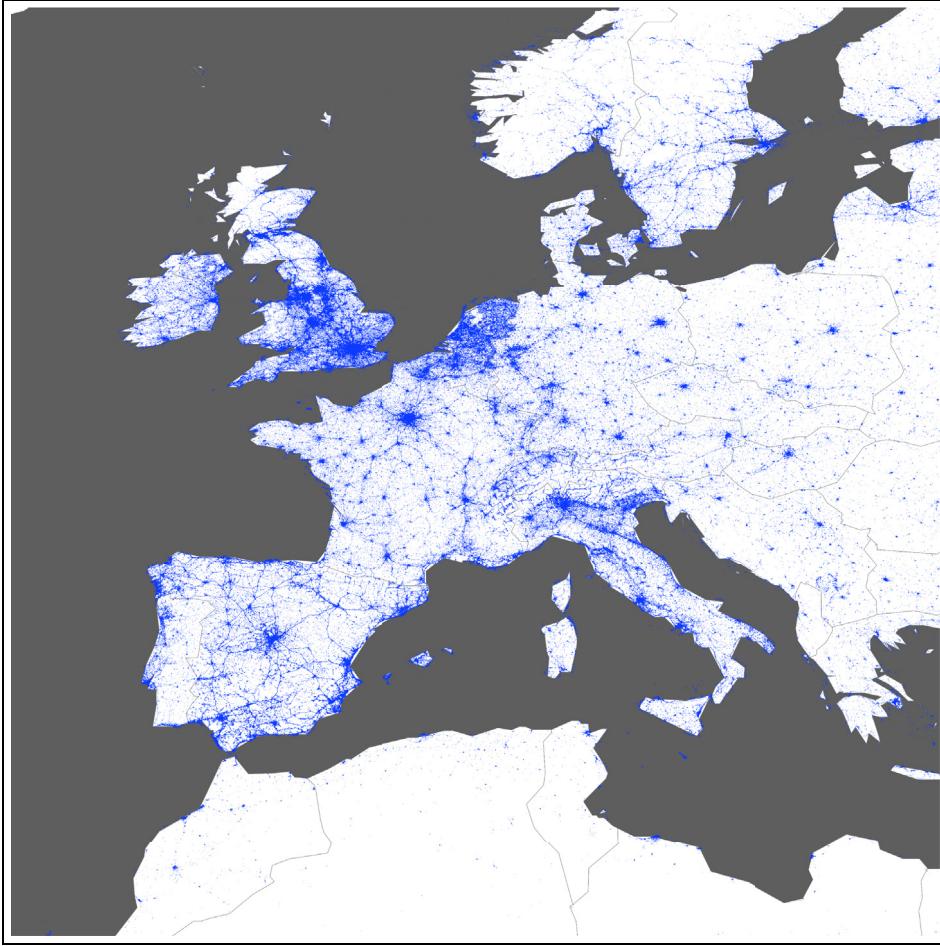


Fig. 3. Map of Image Tweets sent from Europe.
Notes: Each light blue dot represents an image tweet sent from that precise location using information on the latitude-longitude. This is a subset of 100 million random tweets from the complete sample for Jan. 2012–Dec 2013. (For interpretation of the references to colour in this figure legend, the reader is referred to the Web version of this article.)

4. Empirical estimates and methods

4.1. Using twitter to estimate GDP

The first goal of this paper is to study whether Twitter data is a valid proxy for estimating current GDP in USD at the country level. First, I aggregate the number of tweets by country and year. I use the precise location (latitude-longitude) to geocode the country of origin where each post was sent from. I then aggregate the volume of tweets by country per year.

In order to assess the validity of tweets as a proxy for GDP at the country level, I estimate:

$$\ln GDP_{i,t} = \beta_0 + \alpha_t + X'_{i,t}\gamma + \beta_1 \ln Tweets_{i,t} + \varepsilon_{i,t}, \quad (1)$$

where the explained variable is the natural log of GDP of country i in year t . The vector $X_{i,t}$ is composed of country characteristics including population, the share of the population with access to the internet and continent to which it belongs. The coefficient we are most interested in is β_1 which shows the relevance of the number of image tweets taken from that country in each of those years for estimating GDP. In Equation (1), year fixed effects (α_t) control for any differences in the use of Twitter from one year to the other as well as changes in global economic conditions.

The corresponding estimates are reported in Table 3. There are 184 countries in the dataset with data on GDP, Twitter and population for both years.¹² In column 1, I regress the natural log of GDP solely on

the number of image tweets sent from each country. This is the baseline regression. The coefficient of interest on $\ln(Tweets)$ is positive and highly significant and the R^2 is 0.78.

Columns 2–4 explore potential factors that might help explain what is driving the relationship between GDP and tweets. When the population of the country is included in column 2, the coefficient on $\ln(Tweets)$ is reduced, but remains statistically significant at the 1 percent level. This is an important result given that Fig. 1 shows that the clusters of tweets are concentrated in areas with large populations. Thus, the results in column 2 indicate that even when controlling for population, the volume of tweets still holds valuable information for estimating economic activity. Including the population increases the R^2 to 0.87 and the partial R^2 of $\ln(Tweets)$ is 0.66. Furthermore, the root-mean-square error (RMSE) is reduced to 0.85. Column 3 includes categorical dummies for the continents in which each country is situated. This captures economic development as well as cultural differences in image sharing on social media platforms that exist between regions. Neither the coefficient of interest or the goodness of fit greatly change.

Column 4 adds the share of the population with access to the internet. The number of observations are reduced to 180 countries per year because the World Bank does not have data on the share of the population with access to internet for six countries.¹³ Is it possible that the predictive power of the volume of tweets, once we control for population, comes solely as a representation of the share of people who have access to the internet; which in turn is a measure of economic

¹² I also remove countries in which Twitter was banned for a period of time during any of these years: China and Iran.

¹³ These are: Libya, Kosovo, Curacao, Palau, South Sudan and San Marino.

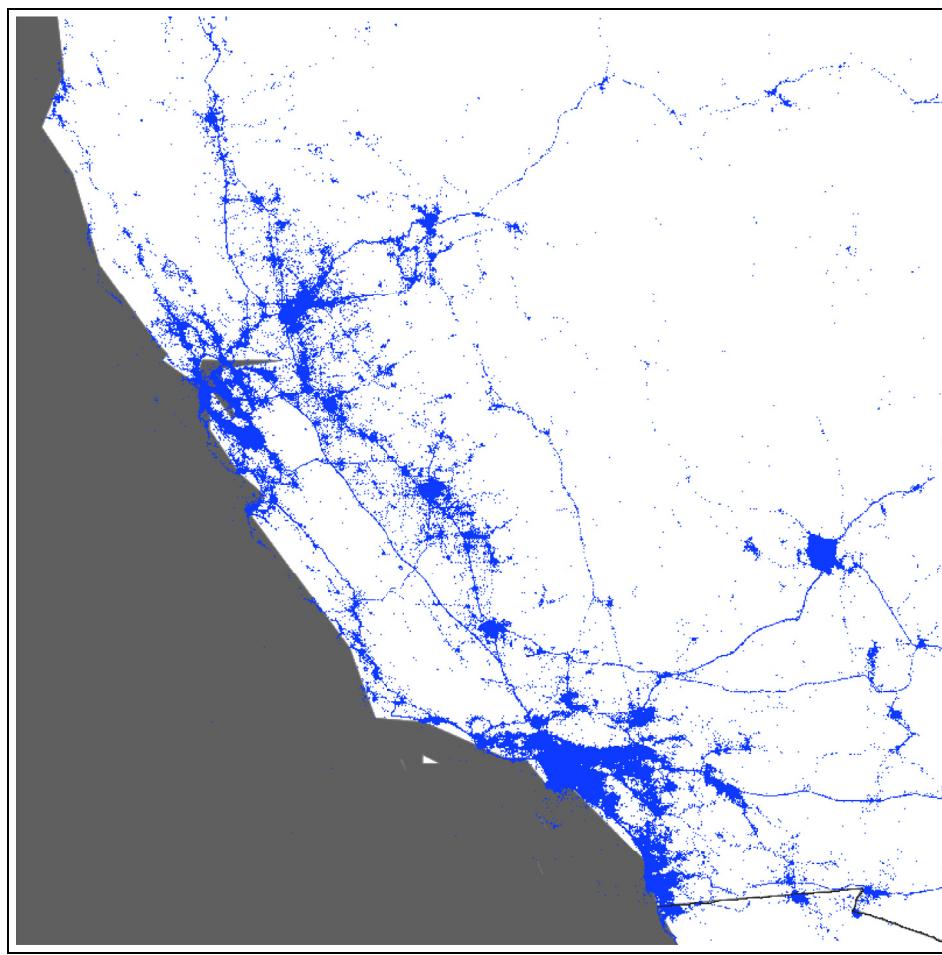


Fig. 4. Map of Image Tweets in the West Coast of the US.

Notes: Each light blue dot represents an image tweet sent from that precise location using information on the latitude-longitude. This is a subset of 100 million random tweets from the complete sample for Jan. 2012–Dec. 2013. (For interpretation of the references to colour in this figure legend, the reader is referred to the Web version of this article.)

development?¹⁴ Column 4 shows that when both measures are included together (and controlling for population and Continent) the coefficient of interest on $\ln(\text{Tweets})$ is reduced, but remains statistically significant at the 1 percent level. This indicates that the volume of tweets is capturing something more than just the share of population with access to internet once we control for population.¹⁵ The R^2 is 0.94 and the partial R^2 of $\ln(\text{Tweets})$ is 0.29. The R^2 without including tweets would be 0.78, similar to the R^2 in Column 1 with only $\ln(\text{Tweets})$ as independent variable. The RMSE is further reduced to 0.56.

Table 3 shows that the number of image tweets sent in a year is a good measure for estimating GDP at the country level, being able to explain 78 percent of the cross-country variation in GDP on its own. In all specifications, the coefficient on the number of image tweets is statistically significant. **Fig. 6** is a visual representation of these estimates for our baseline model in Column 1: the estimates lay pretty closely around the 45° line. There are a few exceptions that stand out; most notably Cuba where internet is restricted. **Fig. 7** plots the residuals of Equation (1) against the fitted values, allowing us to study the distribution of the residuals; which seems to be randomly distributed around zero (i.e.: no clear pattern emerges).

As a robustness check, I run Equation (1) independently for each income group, as per the World Bank's classification.¹⁶ I run the most complete model which includes the number of tweets, the population and the percentage of population with access to internet (i.e.: the model that corresponds to Column 4 in **Table 3**). The corresponding estimates are reported in **Table 4** and show that running the model separately for each income group has little effect on the goodness of fit and relevant coefficients of the model. The measures of R^2 in each of the different specifications are similar and vary between 0.87 for Low-income countries and 0.92 for Lower-middle income countries. The coefficients on $\ln(\text{Tweets})$ vary somewhat for each income group, but are statistically significant at the 1 percent level for all groups. More importantly, the partial R^2 of $\ln(\text{Tweets})$ is between 0.46 and 0.52.

Given that these groups are based on income levels, **Table 4** shows that the number of tweets is able to accurately estimate relatively small differences in GDP levels across countries. Furthermore, **Table 4** shows that tweets can be used to measure economic activity in developing countries where estimating GDP is particularly complicated.

4.2. Using twitter to predict GDP

Section 4.1 showed that we can use tweets to accurately estimate GDP at the country level. In this section, I extend that exercise and

¹⁴ Regressing the log of volume of tweets and population on the share of the population with access to internet gives an R^2 of 0.63. Both coefficients are statistically significant at the 1 percent level.

¹⁵ This will be explored in more detail in section 7.

¹⁶ As of 1 July 2013, low-income economies are defined as those with a GNI per capita, calculated using the World Bank Atlas method, of \$1045 or less in 2012; middle-income economies are those with a GNI per capita of more than \$1045 but less than \$12,746; high-income economies are those with a GNI per capita of \$12,746 or more. Lower-middle-income and upper-middle-income economies are separated at a GNI per capita of \$4125.

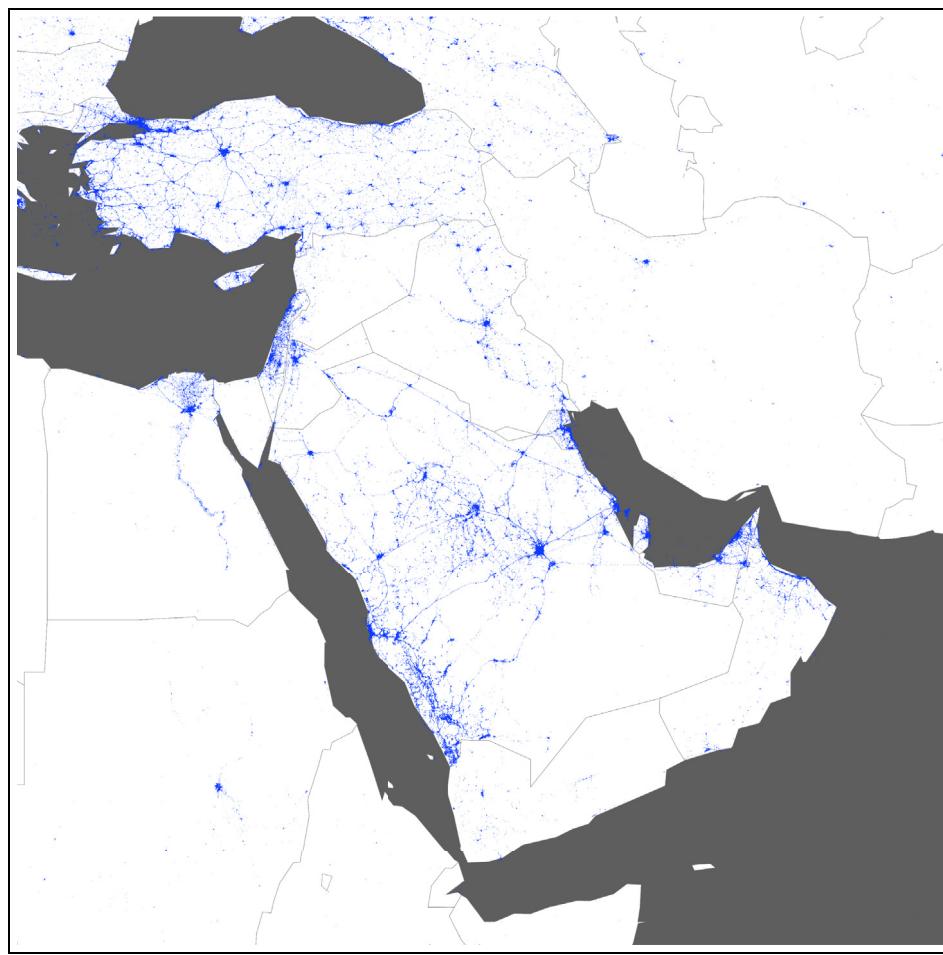


Fig. 5. Map of Image Tweets in Middle East and North-East Africa.

Notes: Each light blue dot represents an image tweet sent from that precise location using information on the latitude-longitude. This is a subset of 100 million random tweets from the complete sample for Jan. 2012–Dec. 2013. (For interpretation of the references to colour in this figure legend, the reader is referred to the Web version of this article.)

Table 3
Estimating country GDP.

Dep. var.:ln(GDP)	(1)	(2)	(3)	(4)
ln(Tweets)	0.67*** (0.02)	0.49*** (0.02)	0.46*** (0.02)	0.17*** (0.02)
ln(Population)		✓	✓	✓
Continent			✓	✓
Internet				✓
R ²	0.78	0.87	0.89	0.94
Adj. R ²	0.78	0.87	0.88	0.94
R ² w/o Tweets	0.00	0.60	0.75	0.78
Partial R ² Tweets	0.78	0.66	0.56	0.29
Fixed-effects	Year	Year	Year	Year
Num. obs.	368	368	368	356
RMSE	1.08	0.85	0.78	0.56

*** $p < 0.01$, ** $p < 0.05$, * $p < 0.10$.

Notes: The dependent variable in all columns 1–4 is the log of GDP. This is using data from 2012 to 2013. All specifications have year fixed effects. Column (4) has fewer observations because the World Bank data does not have data on the percent of population that uses the internet for four countries.

study whether we can use Twitter data to precisely predict GDP at the country level.

For this exercise, I run the preferred model, which includes tweets and population in the specifications outlined in Equation (1). In order to perform out-of-sample prediction, I do not include the complete sample of data. Instead, I carry out a both a leave-one-out as well as a k-fold cross-validation resampling method. Table 5 shows the results for both

of these methods.¹⁷

Fig. 8 shows the results for this exercise for the k-Fold method with 10,000 repetitions. The bulk of predictions are clustered closely around the 45-degree line, indicating that the majority of the out-of-sample predictions do not fall far from the official estimates. This seems to be the case across the income distribution.

¹⁷ The output presented on Table 5 is for $k = 10$ and for 10,000 repetitions. As a robustness check, both the number of folds (k) as well as the number of repetitions were varied without important differences in the results.

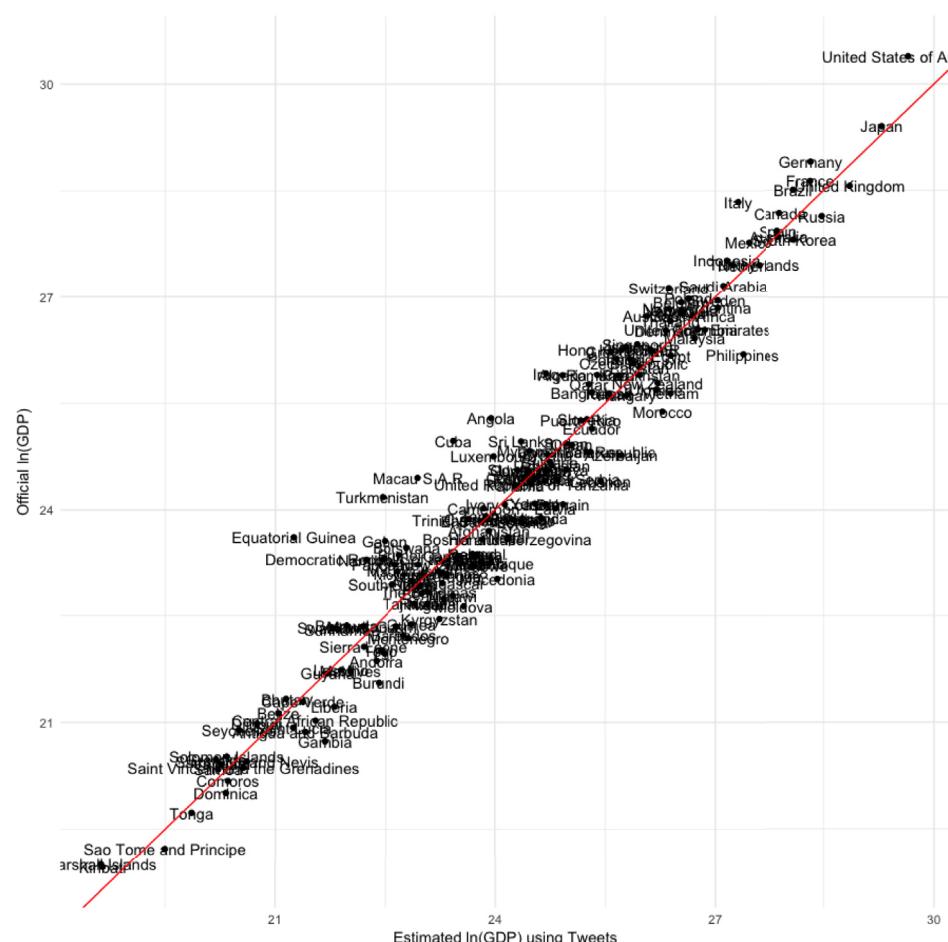


Fig. 6. Estimated vs Official GDP for 2013.
Notes: Red line represents 45 degree line, which indicates where GDP estimates based on Tweets per country and official estimates for GDP are equal. This is for year 2013 and 178 countries included in the sample. (For interpretation of the references to colour in this figure legend, the reader is referred to the Web version of this article.)

Furthermore, Table 5 shows goodness-of-fit results for both methods. The results do not vary too much between methods, overcoming any bias-variance trade-off issue that may arise between the two estimates (James et al. (2013)). More importantly though, the results show that the out-of-sample predictions are quite accurate and the root mean-squared errors (RMSE) are quite similar to those on Table 3.¹⁸

5. Twitter and Night-lights

The use of visible light emanating from earth as captured by weather satellite images has been widely suggested as a good proxy for measuring economic activity. Different studies have shown that night lights can be used to measure GDP estimates at the country level ([Pinkovskiy and Sala-i-Martin \(2016\)](#)), GDP growth at the country level ([Henderson et al. \(2012\)](#)) and GDP for sub-national regions ([Doll et al. \(2006\)](#), [Ghosh et al. \(2010\)](#), [Henderson et al. \(2012\)](#) and [Sutton et al. \(2007\)](#)). According to these studies, the intensity of artificial night-lights highly correlates with GDP and thus can be used to estimate economic activity for different geographic regions. This paper clearly draws a lot from this literature.

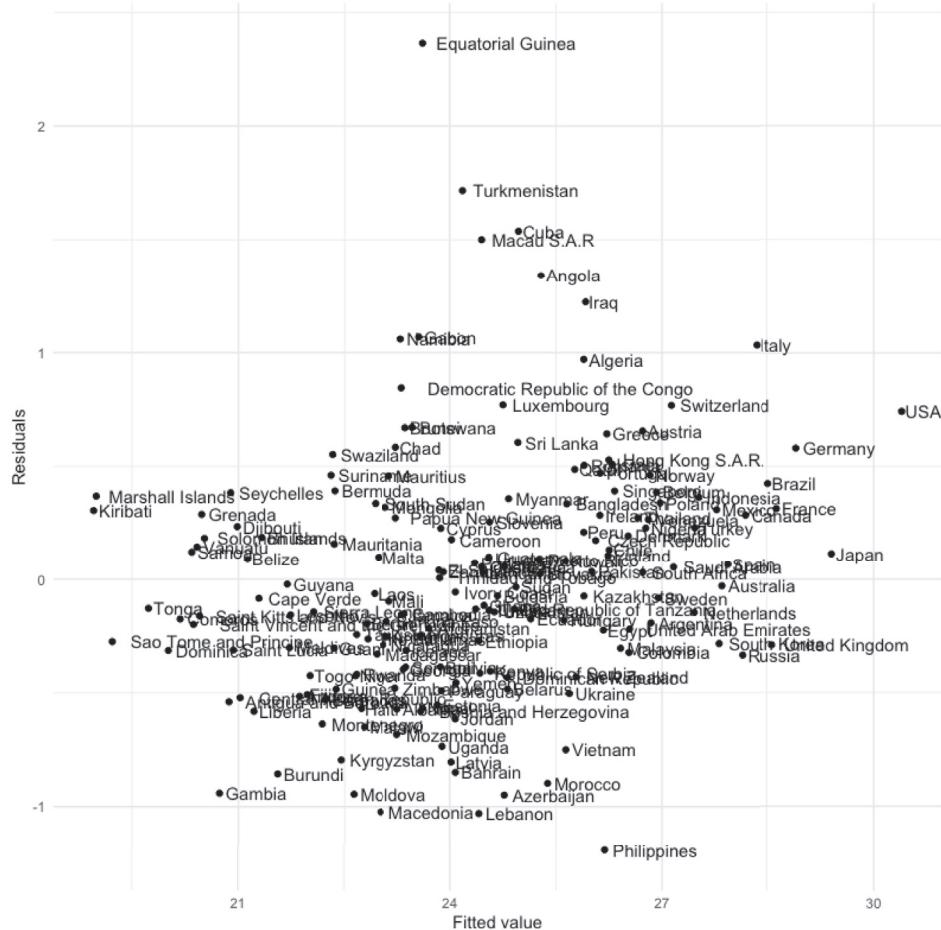
Satellite night-light images come from the US Air Force. Several of their weather satellites circle the earth 14 times per day, recording the intensity of earth-based lights. Each satellite observes every location on the planet (between 65° S latitude and 65° N latitude) every night at some time between 20:30 hs and 22:00 hs. The intensity of lights is

measured for every 30-s output pixel and is averaged across all valid evenings in a year. The raw data are heavily processed to correct for several atmospheric and visual distortions. These tasks include the identification of clouds, removal of glare, identification of intense natural light during summer months, etc. The objective is to leave only man-made light visible. They then average all valid images over the year and report the intensity of light for approximately every 0.86 square kilometer. Intensity of night lights reflects outdoor and some indoor use of lights.

Among the main advantages of using night-lights to estimate economic activity is that satellite data are available more readily than official GDP estimates, that they can be more reliable for some developing countries and that they can be used to estimate economic activity at sub-national regions for which official statistics are not available.

Social media data have these advantages as well. In fact, social media data have a higher-frequency and more granular geographic scope. While luminosity data from satellite images must be averaged over time in order to have a more accurate measure free from cloud and glare, the stream of social media posts is continuous. Furthermore, while intensity of night-lights is captured in 0.86 square km pixels, each social media post has its corresponding latitude and longitude coordinate, precise to the nearest 1.1 m; this allows for unconstrained and

¹⁸ Another important thing to point out that is not presented on the table, is that the average coefficients are similar to those presented in Table 3 for the entire dataset. This would indicate that the model is stable.

**Fig. 7.** Residual and Fitted Value for 2013 GDP Estimates.

Notes: Plot shows the distribution of the residuals of preferred model against the fitted values. This is for year 2013 and 178 countries included in the sample.

Table 4
Estimating country GDP: By income group.

Dep. var.:ln(GDP)	Low	Low-middle	Upper-middle	High
ln(Tweets)	0.44*** (0.09)	0.42*** (0.04)	0.37*** (0.05)	0.47*** (0.04)
ln(Population)	✓	✓	✓	✓
Continent	✓	✓	✓	✓
R ²	0.87	0.92	0.89	0.90
Adj. R ²	0.84	0.91	0.88	0.89
R ² w/o Tweets	0.53	0.80	0.85	0.82
Partial R ² Tweets	0.46	0.48	0.53	0.52
Fixed-effects	Year	Year	Year	Year
Num. obs.	56	78	104	132
RMSE	0.95	0.72	0.78	0.79

***p < 0.01, **p < 0.05, *p < 0.10.

Notes: The dependent variable in all columns is the log of GDP. These estimates correspond to estimating Column (4) of Table 3 independently for each income group. The first column is for the subset of Low-income countries, second column for Low-middle income countries, third column for Upper-middle income countries, and the fourth column for High-income countries. These categories are based on The World Bank classifications. This is using data from 2012 to 2013. All specifications have year fixed effects.

Table 5
Predicting country GDP out-of-sample:
Cross-validation results.

Dep. var.:ln(GDP)	LOO	k-Fold
R ²	0.863	0.866
Adj. R ²	0.863	0.866
RMSE	0.865	0.863

Notes: The table presents results for regressing the log of Tweets and population on the log of GDP. Column 1 presents results for the leave-one-out (LOO) cross-validation estimate. Column 2 presents the results for the k-Fold cross-validation method, with k = 10 and 10,000 repetitions.

more flexible geographical aggregations.¹⁹

However, as with any proposed proxy, there are shortfalls with night-light data. Chen and Nordhaus (2010) show that satellite night-light data suffers from substantial measurement error, both in the time-series and the cross-section aspect of the data. The authors regress the logarithm of luminosity for the same year across images from different satellites and find standard errors in the range of 0.20. They also find slightly smaller standard errors when regressing year-to-year variations in luminosity from images taken by the same satellite. This means that the measurement error for individual grid cells is in the order of 20 logarithmic percent. They also find that images from specific satellites are consistently dimmer than other satellites. This leads the authors to question the reliability of the night-lights data and conclude that it “is at best a noisy indicator” (Chen and Nordhaus (2010)).

Doll (2008) finds that satellite recording of nighttime light density has a tendency to overestimate the true extent of lit area on the ground. This occurs because these images tend to attribute light generated at a particular site to nearby sites as well. This effect is referred to as overglow. In part, this is explained by the procedure in which satellites capture night-lights, but it is also due to the large overlap that exists between pixels.

Furthermore, luminosity at the pixel or grid level is bottom coded at 0 and top coded at 63. As a consequence, in certain geographically small and rich countries, more than 3 percent of the pixels are topped off at 63 (i.e.: Netherlands and Belgium). Thus inhibiting the estimator to parse out differences at the top end of the distribution. When using social media data there is no top coding since we can have an infinite number of posts coming from the same location. This could prove useful when trying to detect small differences at the top end of the income distribution.

On the other hand, there are several advantages for the use of night-light images over social media data. First of all, satellite images are available since 1992, which allows for research on the evolution of GDP and medium-term growth paths. Henderson et al. (2012) use this to estimate economic growth of countries and regions over the last 30 years. Currently, this sort of analysis could not be done using social media data. Since social media became widely popular in the mid-2000s, it is not suitable to perform historical research. Furthermore, the availability of a relatively long time series data has demonstrated that night-lights have been a valid proxy for measuring economic activity for several decades. This shows that the underlying mechanism which relates night-lights to economic activity has remained fairly stable over time. This has yet to be proven with social media data. Since these platforms are fairly new and people are still adapting to these technologies, it may

¹⁹ Although it is not clear if this temporal and geographical continuum truly represent a useful advantage in terms of estimating economic activity. To take an extreme example: it is not clear that one can use the volume of tweets in every second to represent economic activity per-second, nor can one differentiate economic activity generated in every meter of a region.

be the case that our use of social media changes dramatically in some years time, which may invalidate their usefulness as a proxy for GDP²⁰.

Given that the Twitter data used in this paper correspond to 2012 and 2013, I collect Operational Linescan System data on night lights intensity from the United States Air Force Defense Meteorological Satellite Program for those same years. These two datasets will be used to compare the potential use of each of them as proxies for estimating GDP at the country level and, more interestingly, to see whether both proxies can be used together as proxies. Given that the underlying mechanism driving the relationship between each of these proxies and GDP is believed to be somewhat different (more on this in section 7) and that their measurement errors are not correlated, we can exploit that fact that several error-prone measures are better than one (Rao (1992)).

As a first approach, I map the distribution and intensity of both night-lights and Twitter data.²¹ Fig. 9 presents Twitter posts on the top and satellite night-light images on the bottom. The two maps show a striking resemblance, as the bulk of tweets and night-light intensity clusters around large cities and capitals around the world. Fig. 10 focuses on Europe and again shows the similarity between the location and volume of tweets relative to the intensity of night-lights emanating from earth.

A more formal inspection on the relationship between night-lights and tweets is shown in Fig. 11. The figure plots night-lights against the volume of Tweets in each country for the year 2013 and includes a local non-parametric regression line. A few countries are labeled for reference. Among countries with higher night-light intensity relative to tweets are Cuba, Iran and China. Internet is restricted in Cuba and Twitter was banned in both China and Iran during this period. On the other end are island countries, such ad Fiji, Palau and Dominica, which have a larger number of tweets than lights according to the overall relationship. Another interesting outlier in this regard is Macau S.A.R, which is an autonomous region on the south coast of China, where Twitter and other social media platforms are not banned. The correlation coefficient between these two variables taking both 2012 and 2013 into consideration is 0.828.

In order to compare and contrast the validity of night-lights and tweets as a proxy for GDP at the country level and see the extent to which they can be used together, I estimate:

$$\begin{aligned} \ln GDP_{i,t} = & \beta_0 + \alpha_t + \beta_1 \ln Lights_{i,t} + \beta_2 \ln Tweets_{i,t} \\ & + \beta_3 \ln Population_{i,t} + \varepsilon_{i,t}, \end{aligned} \quad (2)$$

where the explained variable is the natural log of GDP of country i in year t . The two coefficients of interest are β_1 and β_2 , which respectively show the relevance of the intensity in night-lights and the number of image tweets for estimating GDP. β_3 controls for the population. Year fixed effects (α_t) control for any differences in the use of Twitter from one year to the other and changes in global economic conditions.

Column 1 of Table 6 includes only night-light intensity and the model explains 85 percent of the variation in GDP at the country level. Column 2 includes all variables in $X'_{i,t}$, which comprises population, continent dummies and percent of population with access to the internet, and night-light data explains roughly 88 percent of the variation in GDP. Including both Twitter data as well as night-lights in Column 3 of Table 6, the coefficient on both of these variables are statistically significant, reflecting that they both provide distinct information for esti-

²⁰ It is important to stress that even though this paper studies solely the attributes of Twitter data, this is just one of several social media platforms that could potentially serve as a proxy for different economic measurements. There are several examples of social media platforms whose popularity has quickly decreased over time. This is obviously a big concern for the use of any one platform as a proxy. But less so about the overall use of social media platforms in general for these tasks.

²¹ This exercise was in fact the way I became truly excited about the idea that social media data could be used as a proxy for economic activity.

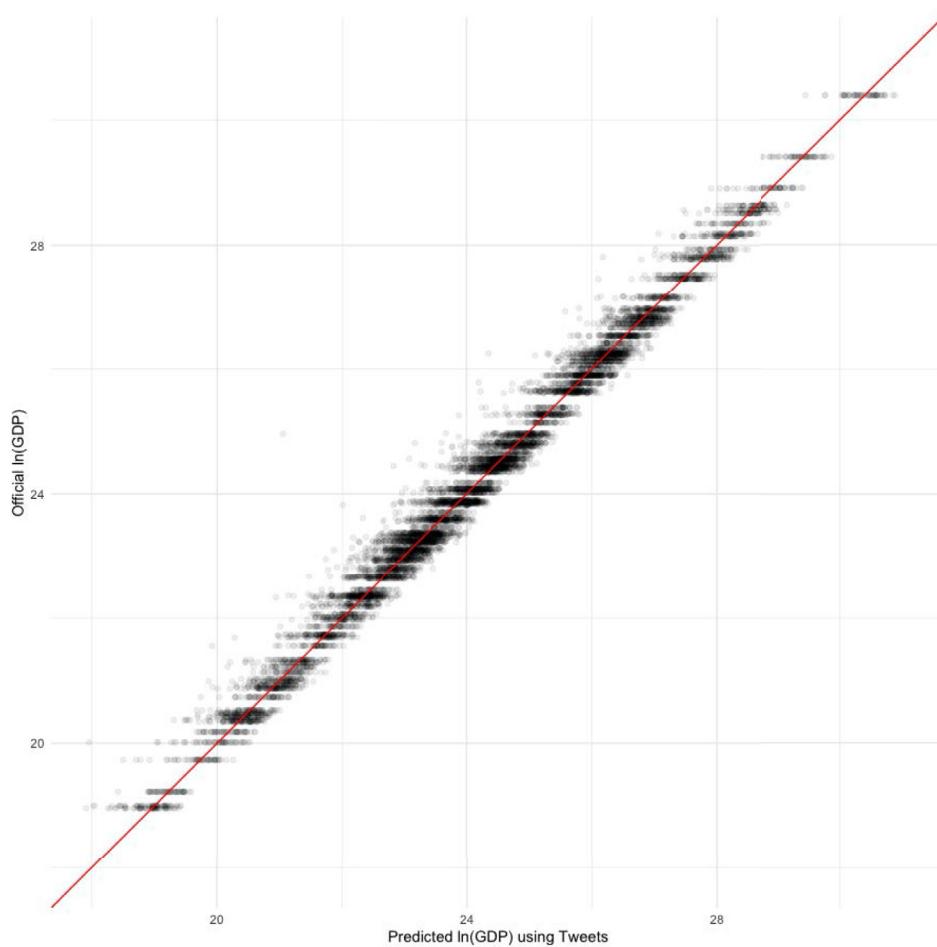


Fig. 8. Predicting Country GDP Out-of-Sample.
Notes: Red line represents 45 degree line, which indicates where GDP predictions based on Tweets and population per country and official estimates for GDP are equal. This is done using a k-Fold cross-validation with $k = 10$ and based on 10,000 resamples. (For interpretation of the references to colour in this figure legend, the reader is referred to the Web version of this article.)

mating GDP and suggesting that ideally both measures could be used together to make a more accurate proxy of GDP. R^2 increase to 0.95 and RMSE decreases to 0.53, indicating that estimates when both proxies are used together are quite accurate. An important result showing the capabilities of both proxies used together is that the partial R^2 of the two are quite evenly distributed: partial R^2 for lights is 0.27 and 0.32 for tweets.

Another way to compare the use of night-lights and Twitter data as a proxy for estimating economic activity is to analyze the estimated GDP using each of the two proxies. Fig. 12 plots both of these estimates for year 2013. We see that the bulk of the estimates lie closely around the 45-degree line which shows that the two proxies not only estimate a significant portion of the variation in GDP, but that the individual estimates for each country are similar using each of these two proxies. This validates the use of Twitter lights as a proxy.²²

In sum, this section shows that Twitter and night-light data could be used together as proxies in order to more accurately estimate economic activity. Given that the two proxies seem to be capturing different aspects of economic activity, that their measurement errors are not correlated, and that the explanatory power is markedly improved when both proxies are included, it seems these two variables have salient features that make them useful complementary proxies.

6. Estimating changes in GDP

The previous sections have established the validity of social media data as a proxy for estimating the level of economic activity. An obvious extension is to see whether social media data can be used to estimate within country changes in economic activity from one period to the next.

This exercise has already been carried out for satellite night-lights with diverse success. Henderson et al. (2012) show that night-light luminosity does a reasonable job at predicting annual fluctuations in GDP growth for countries over a 25 year period. Although they find that luminosity has some predictive power on estimating short-term annual changes in economic activity, the proxy is more accurate at estimating long-term growth between 1992 and 2007.

Due to the weakness in predicting annual changes, the authors study whether light intensity is prone to a ratchet issue: given that light growth seems to be reflecting installation of new capacity, light intensity is non-decreasing and thus does not capture decreases in economic activity. Although the authors do not find such a ratchet effect, the volume of social media posts seem to be more sensitive to short-term fluctuations in economic activity (see section 7) and thus could serve as a more accurate proxy for estimating changes in economic activity. In this sense, I argue that social media data is better equipped as a proxy for estimating short-run variations in economic activity, particularly in recessionary periods.

In order to study whether tweets have predictive power in estimating annual fluctuations in GDP and how it compares to night-lights, I estimate:

$$\ln GDP_{i,t} = \beta_0 + \alpha_i + \alpha_t + X'_{i,t} \gamma + \beta_1 \ln Lights_{i,t} + \beta_2 \ln Tweets_{i,t} + \varepsilon_{i,t}, \quad (3)$$

²² Section 6 extends this analysis and shows that both Twitter data as well as night-lights can also be used to estimate within-country variation in GDP from one year to the other. By including country fixed effects to Equation (2) the model estimates the within country changes in GDP to 2012 and 2013 and Table 7 shows that the coefficient on both tweets and night-lights is positive and statistically significant.

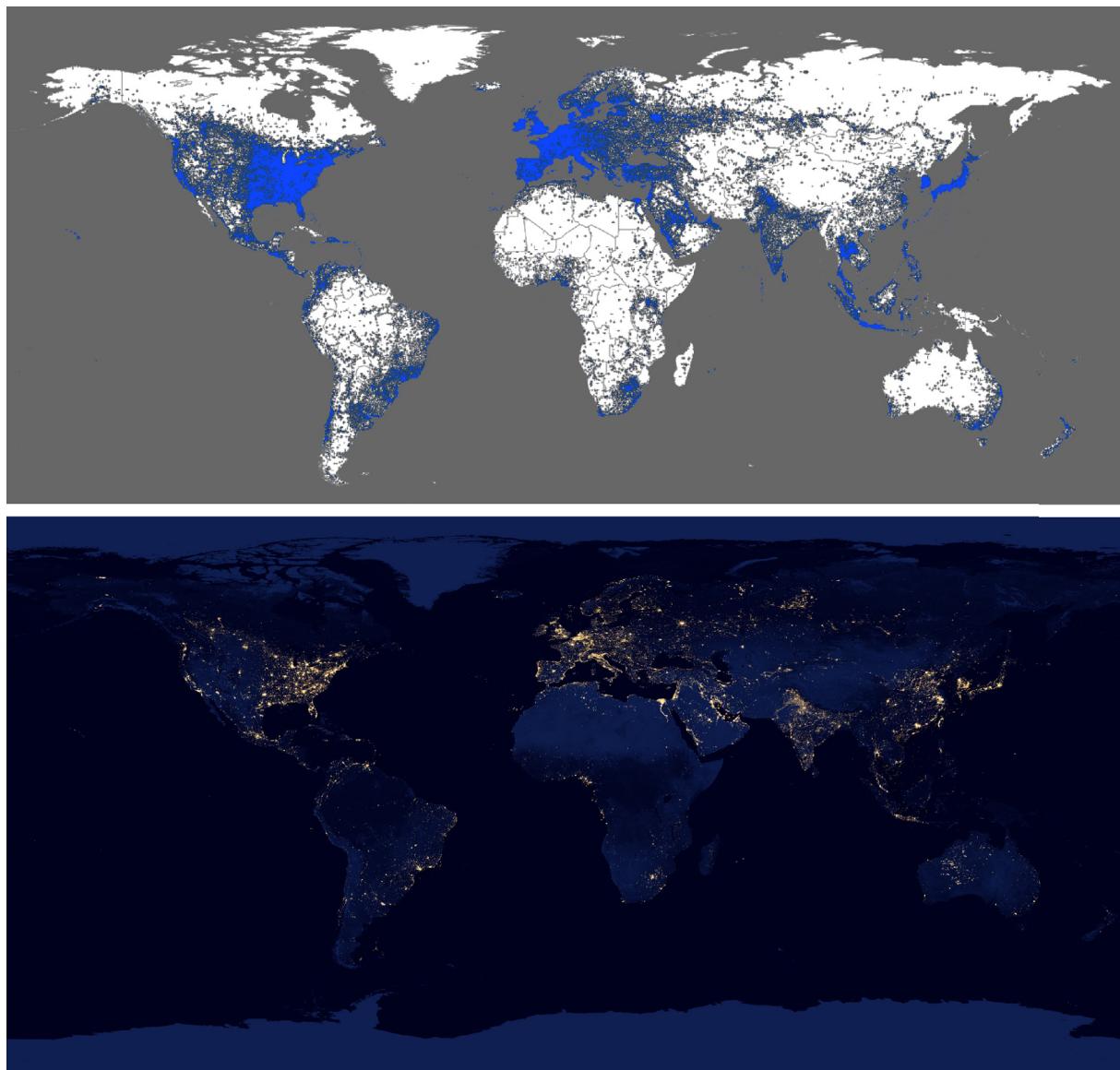


Fig. 9. Twitter and Night-Light Maps for the World

Top: Twitter posts. Each light blue dot represents an image tweet sent from that precise location using information on the latitude-longitude. This is a subset of 100 million random tweets from the complete sample for Jan. 2012–Dec. 2013.

Bottom: Satellite night-light image for 2012. (For interpretation of the references to colour in this figure legend, the reader is referred to the Web version of this article.)

Source:<https://earthobservatory.nasa.gov/images/79765/night-lights-2012-flat-map>

which is equivalent to Equation (2) with country fixed effects (α_i). In essence this measures the annual variation in GDP for country i between years 2012 and 2013. The coefficients of interest are β_1 , which indicates the relationship between changes in GDP and changes in light intensity and, more importantly, β_2 , which shows the relationship between changes in GDP and changes in the volume of tweets. The corresponding estimates are reported in Table 7.

Column 1 shows results when Equation (3) includes only tweets as explanatory variable; the coefficient is positive and statistically significant, with an R^2 of 0.23. Column 2 includes several country characteristics. Similar to section 4, the motivation behind including these variables is to study whether Twitter data is capturing more than a combination of population and the share of that population with access to the internet. The coefficient on Tweets remains positive and statistically significant, at the 5 percent level, indicating that it seems it does capture something that the other measures do not pick up. Column 3 mimics

this specifications solely for night-lights: the coefficient on night-lights is positive and statistically significant (at the 5 percent level), with an R^2 of 0.33. Finally columns 4 and 5 include both tweets and night-lights together. The coefficient on both proxies are positive and statistically significant in both specifications (at the 5 percent level for night-lights when full set of country characteristics is included). The R^2 increases to 0.37 when all covariates are included. Once again, this indicates that these two measures may not be capturing exactly the same aspects of economic activity and the explanatory power improves with the inclusion of both.

Albeit not being very precise, the relationship between annual changes in social media posts and annual changes in GDP are informative. Furthermore, it seems that tweets and night-light luminosity could be used together to produce a more accurate estimate of annual variations in economic activity.

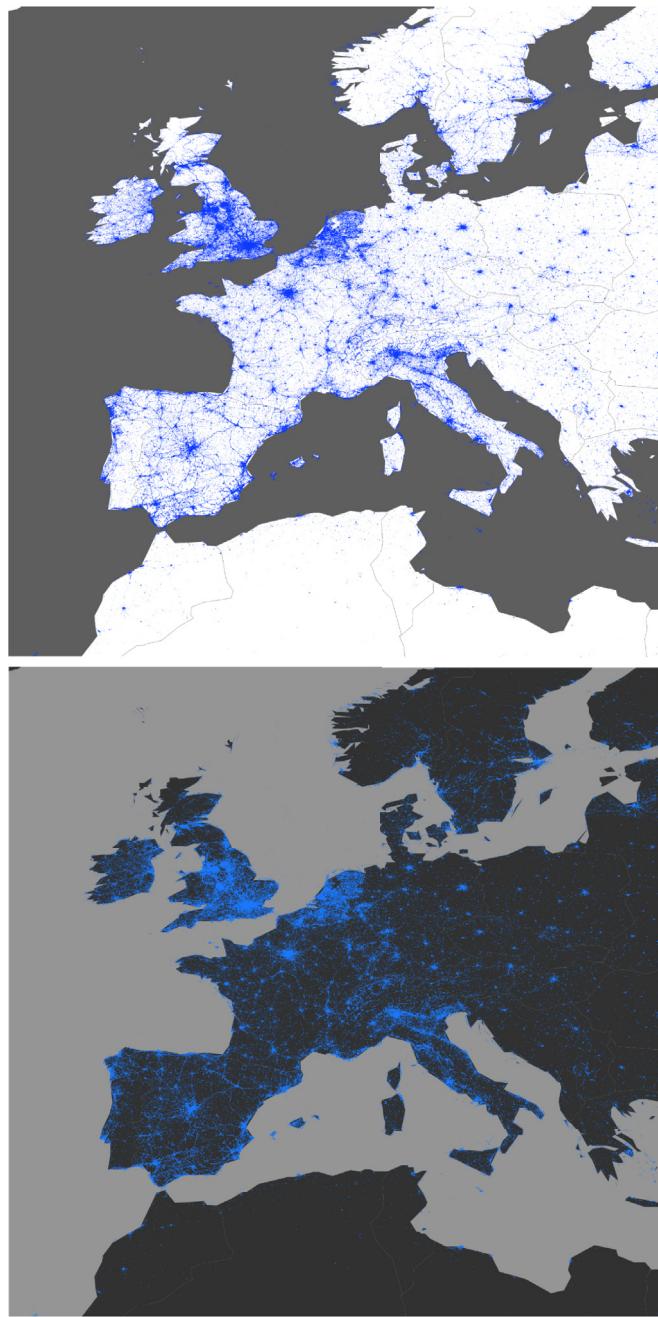


Fig. 10. Twitter and Night-Light Maps for Europe

Top: Twitter posts. Each light blue dot represents an image tweet sent from that precise location using information on the latitude-longitude. This is a subset of 100 million random tweets from the complete sample for Jan. 2012-Dec. 2013.

Bottom: Satellite night-light image for 2012. (For interpretation of the references to colour in this figure legend, the reader is referred to the Web version of this article.)

Source: <https://earthobservatory.nasa.gov/images/79765/night-lights-2012-flat-map>

7. Conspicuous consumption, leisure and tweets

The previous sections have shown diverse ways in which the volume of image tweets is related to economic activity. Nonetheless, it is still not clear what the underlying mechanism that relates social media posts

and economic activity.²³

One possible explanation is related to smart-phone penetration, which is essentially an indicator of economic development. Under this hypothesis, more developed countries have more smart-phones per individual, which in essence increases the ratio of the population with social media accounts (which are free of charge to the user), and thus increases the number of social media posts per capita. Notice that under this hypothesis, the number of posts per user does not vary across individuals, but the aggregate number of posts per country increases as a result of more individuals with smart-phones and hence social media accounts. This explanation works pretty well to explain cross-country variations, especially between developed and developing countries.

Nevertheless, this hypothesis does not explain why the volume of social media posts explains within-country differences (as shown in subsection 8.1). Notably in developed countries, where smart-phone penetration has been saturated for some time, there must be another underlying mechanism that drives more tweets in one region than another. Although there are likely several factors at play, one possible hypothesis is that social media applications serve as a medium to showcase consumption of goods and services among the network of users.

This behavior falls in line with what Veblen (1899) coined as conspicuous consumption: when the utility consumers attain stems more from revealing their wealth and income to others, rather than from the direct utility derived from the good itself. Under this hypothesis, individuals display wealth to their network as a means to attaining a given social status which is what they ultimately desire. In order to preserve this status, they most constantly exhibit the consumption patterns online. In this scenario, social media represents the ideal medium through which individuals may showcase a stream of consumption of goods and services to friends and acquaintances.²⁴

According to Becker (1965), consumption and leisure are complementary goods; thus we would expect both of them to occur at the same time. To test whether leisure tweets are better predictors of economic activity than tweets sent during work hours, I divide all image tweets sent from the US in 2012 in two groups: (i) those sent during working hours (Monday to Friday between 9:00–17:00 hs); and (ii) those sent during leisure time, defined as non-working weekday hours (Monday to Friday after 17:00 hs and before 9:00 hs) and weekends (Saturday and Sunday any time).

Under this hypothesis, I expect the coefficient on leisure tweets to be positive and statistically significant, while the coefficient on tweets during working hours to be either negative or not statistically significant. In order to directly study the hypothesis that tweets are a byproduct of consumption, I will see if the volume of tweets are a valid proxy for personal consumption expenditure in the US. In order to test this, I aggregate the number of tweets in each state which were sent during working and leisure times and run the following regression:

$$\ln \text{Consumption}_i = \beta_0 + \beta_1 \ln \text{Population}_i + \beta_2 \ln \text{Leisure Tweets}_i + \beta_3 \ln \text{Work Tweets}_i + \varepsilon_i \quad (4)$$

where I estimate personal consumption expenditure for state i in year 2012. The coefficients of interest are β_2 and β_3 which show the relevance

²³ Given that this paper is studying the validity of social media as a proxy for economic activity and not trying to establish any causal effect between the two, a thorough understanding of their relationship is not necessary. Nevertheless, a more clear understanding of the underlying mechanism that drives the two is helpful for understanding why the proxy works and more importantly, when it might stop being accurate.

²⁴ Han et al. (2019) discuss extensively the implications this representation of consumption has at an aggregate level. They present a theoretical model in which consumption is more salient than nonconsumption in the social transmission process which affects consumption behavior along the network. This visibility bias causes people to perceive that others are consuming heavily and thus increases aggregate consumption.

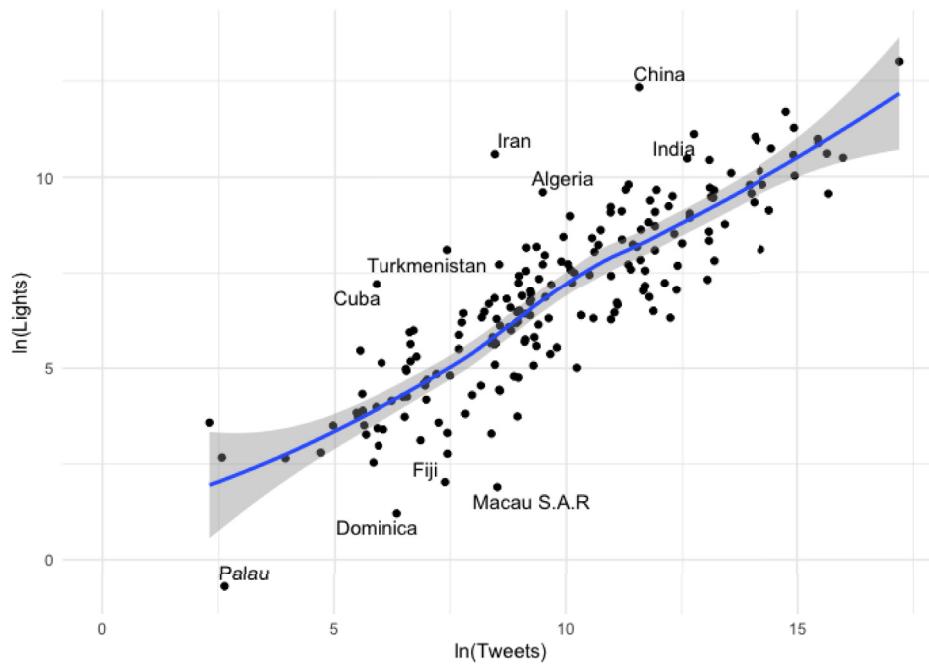


Fig. 11. Relationship Between Night-Lights and Twitter for 2013.

Notes: The figure plots natural log of night-lights against the natural log of number of Tweets in each country. Data included in this figure is for year 2013 only. A local non-parametric regression line is included. Several countries are labelled for reference. This data includes 174 countries.

Table 6
Lights vs Twitter for Estimating GDP.

Dep. var.:ln(GDP)	(1)	(2)	(3)
ln(Lights)	0.88*** (0.02)	0.31*** (0.03)	0.26*** (0.03)
ln(Tweets)			0.08*** (0.03)
ln(Population)		✓	✓
R ²	0.85	0.88	0.95
Adj. R ²	0.85	0.88	0.95
Partial R ² Lights	0.85	0.38	0.27
Partial R ² Tweets	–	–	0.32
Num. obs.	350	346	346
RMSE	0.84	0.63	0.53

***p < 0.01, **p < 0.05, *p < 0.10.

Notes: The dependent variable in all columns 1–3 is the log of GDP. This is using data from 2012 to 2013. All specifications have year fixed effects.

vance of the number of image tweets taken in each of those time slots for estimating consumption at the state level.

The corresponding estimates are reported in Table 10. Columns 1 and 2 introduce the volume of all tweets, irrespective of when they were taken. The coefficient is positive and statistically significant (at the 10 percent level) even when controlling for population in column 2. Once again, population is included as a control to test whether the volume of image tweets is capturing something more than just population. Given that the coefficient remains positive and statistically significant indicates that social media posts captures something else besides population, and the explanatory power is greatly improved by the inclusion of both. Column 3 includes only leisure hour tweets, and the coefficient is positive and statistically significant even when controlling for population. Columns 5 and 6 introduce both leisure and work tweets: while leisure tweets are positive and statistically significant, work tweets are negative (although not statistically significant).

These results seem to support the hypothesis that Twitter serves as a medium which enables users to share their consumption habits throughout their network. This hypothesis could explain one of the underlying

mechanisms that relates the volume of tweets and economic activity, with image tweets being a byproduct of consumption. Given that consumption represents roughly 2/3 of overall GDP in the US since 2000, this explains why image tweets are a good proxy for estimating economic activity at large²⁵

²⁵ The concept of using solely those tweets sent during non-working hours was not utilized during the main sections of this paper due to the fact that different countries have different workdays and weekdays schedule. Many countries observe a Monday-Friday workweek, others a Sunday-Thursday and some countries have other specifications. Moreover, the typical workday hours varies widely between countries, and these are particularly hard to find for each country. I did perform the main results on Table 3 using solely non-working hour tweets assuming a baseline 9:00–17:00 hs workday for all countries and my best attempt, using several sources, at identifying proper weekend days for each of them. The results are quantitatively similar and qualitatively the same.

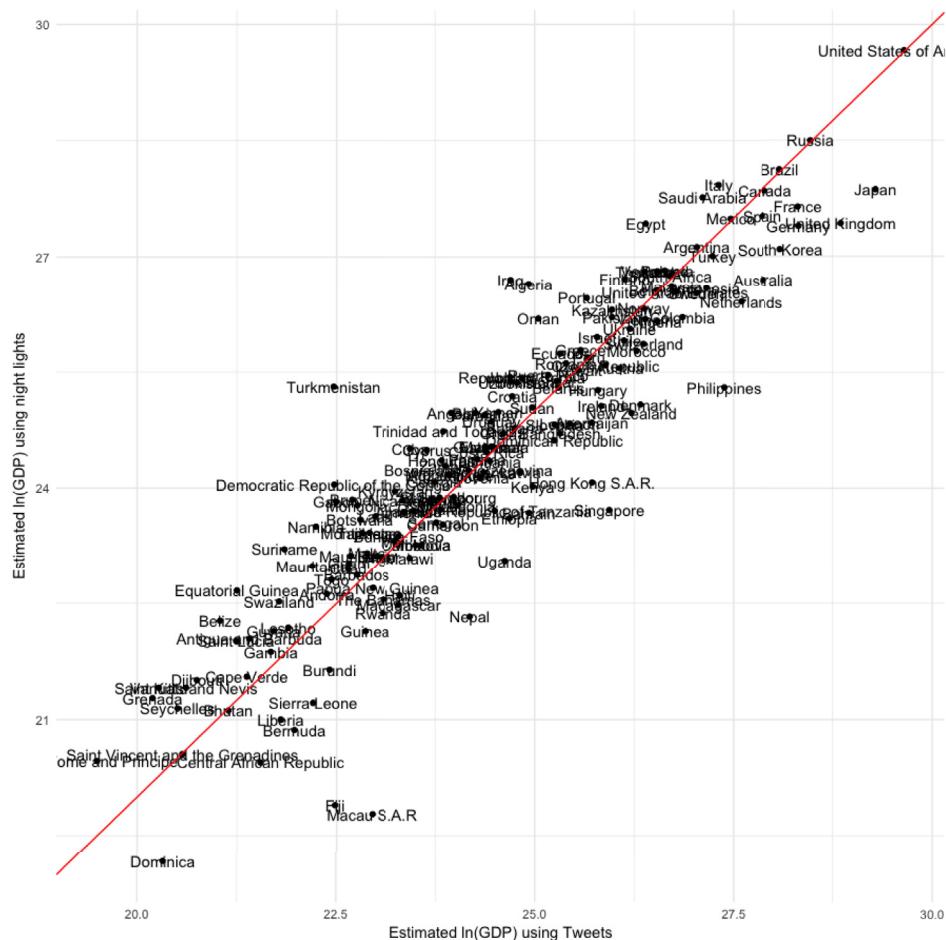


Fig. 12. Night-Lights vs Twitter GDP Estimates for 2013.

Notes: Red line represents 45 degree line, which indicates where GDP estimates using Twitter and night-lights are equal. This is for year 2013 and 171 countries included in the sample. (For interpretation of the references to colour in this figure legend, the reader is referred to the Web version of this article.)

Table 7
Country fixed effects.

Dep. var.:ln(GDP)	(1)	(2)	(3)	(4)	(5)
ln(Tweets)	0.02*** (< 0.01)	0.01** (< 0.01)		0.02*** (< 0.01)	0.02*** (< 0.01)
ln(Lights)			0.06** (0.03)	0.09*** (0.02)	0.07** (0.07)
ln(Population)		✓			✓
Continent		✓			✓
Internet		✓			✓
R ²	0.23	0.33	0.33	0.36	0.37
Adj. R ²	0.17	0.28	0.18	0.32	0.30
Fixed-effects	country	country	country	country	country
Num. obs.	364	304	356	352	298

***p < 0.01, **p < 0.05, *p < 0.1.

Notes: The dependent variable in all columns is the log of GDP. These estimates come from including country fixed-effects to Equation (2). This is using data from 2012 to 2013. Columns 2 and 5 have fewer observations because the World Bank does not have internet access data for some countries.

8. Exploiting geographical granularity of social media data

8.1. Estimating GDP for US states and cities

One of the main advantages of social media data in estimating GDP lies in the geographic detail of social media posts. Exploiting the accuracy of the location from where social media posts are sent allows one

to aggregate posts at any sub-national geographical level one deems interesting. This includes aggregating data between areas that are not officially bound together and thus fabricate meaningful areas of study that are not possible with official datasets.

Given that national accounts data tends to be aggregated at the country-level it is not suitable for this type of analysis. Even in the cases when they do offer regional level estimates, it is not possible to use esti-

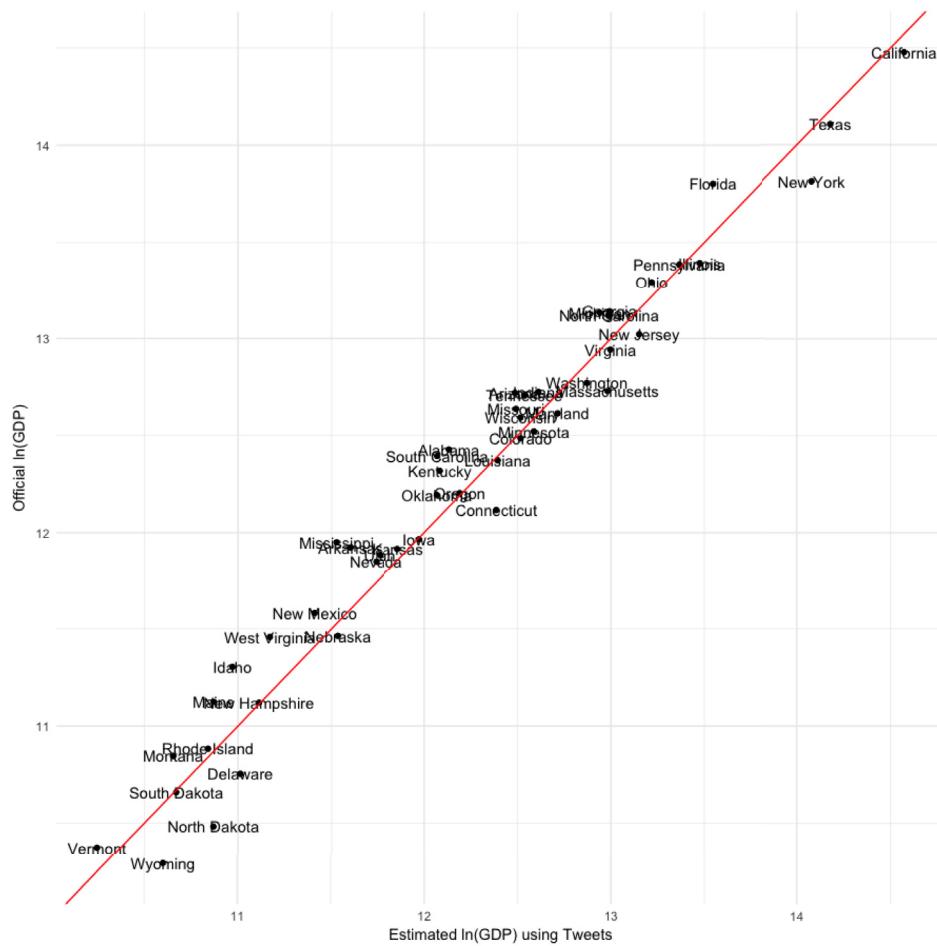


Fig. 13. Estimated vs Official GDP at the State Level for 2012.

Notes: Red line represents 45 degree line, which indicates where estimates based on volume of tweets and official estimates for GDP at the US state level are equal. This is for year 2012. (For interpretation of the references to colour in this figure legend, the reader is referred to the Web version of this article.)

Table 8
Estimating GDP for metropolitan statistical areas.

Dep. var.:ln(GDP)	(1)	(2)	(3)	(4)
ln(Tweets)	0.23*** (0.04)	0.12*** (0.03)		0.06** (0.03)
ln(Lights)			0.16*** (0.02)	0.10*** (0.02)
ln(Population)		✓		✓
R ²	0.52	0.59	0.48	0.68
Adj. R ²	0.52	0.58	0.48	0.67
Partial R ² Tweets	0.52	0.35	–	0.31
Partial R ² Lights	–	–	0.48	0.29
Fixed-effects	Year	Year	Year	Year
Num. obs.	464	464	464	464
RMSE	1.15	0.93	1.17	0.67

*** $p < 0.01$, ** $p < 0.05$, * $p < 0.10$.

Notes: The dependent variable in all columns is the log of GDP. This is using data from 2012 to 2013 for MSA in the US.

mates from different countries in the same analysis because each survey is conducted independently with different questions and methodologies and thus cannot be taken together.

Several papers successfully carried out this sort of analysis using night-light data to estimate economic activity over various regions for which no official GDP estimate exists (Henderson et al. (2012), Pinkovskiy (2017) and Michalopoulos and Papaioannou (2018)).

In order to corroborate the validity of Twitter data as a proxy for GDP at the sub-national level, I first estimate GDP for individual US

states and compare that to official state level estimates provided by the Bureau of Economic Activity (BEA). For this exercise, I collect the subset of all tweets from 2012 taken in the US and use the latitude and longitude from each post to assign them to the state from where they originate. Using the total number of tweets sent from each state, I use the coefficients from Equation (1) and collect GDP estimates for 2012 for each US state. Fig. 13 plots both the BEA as well as the estimates from the Twitter data. The bulk of the estimates lie closely around the 45-degree line which shows that Twitter data is a valid proxy for esti-

Table 9
Predicting GDP for US Metropolitan Statistical Areas: Cross-validation results.

Dep. var.:ln(GDP)	LOO	k-Fold
R ²	0.676	0.676
Adj. R ²	0.675	0.676
RMSE	0.671	0.670

Notes: The table presents results for regressing the log of Tweets and population on the log of GDP. Column 1 presents results for the leave-one-out (LOO) cross-validation estimate. Column 2 presents the results for the k-Fold cross-validation method, with k = 10 and 10,000 repetitions.

mating GDP at the state level. Furthermore, the correlation coefficient between these two estimates is 0.985.

The geographic granularity of social media data allows us to go beyond estimating economic activity at the state level and estimate economic activity at the city level. For this, I study the relationship between tweets, night-lights and economic activity at the metropolitan statistical area (MSA) level. By definition, MSA's are geographical regions which have a relatively high population density at its core and have close economic ties throughout the area. Using coordinates from tweets, I assign them to their corresponding MSA.²⁶ I then capture the night-light intensity from each MSA and the official GDP estimates for the BEA.

[Table 8](#) shows the results for Equation (2) on MSA's for 2012 and 2013. Column 1 shows the results for estimating MSA economic activity using tweets; column 2 controls for population. The coefficient of interest is statistically significant in both cases. The elasticity between tweets and economic activity when controlling for population is roughly 0.12. The R² is 0.52 and 0.58, respectively. Column 3 is estimated using solely night-lights. The coefficient is statistically significant with an R² of 0.48. Finally, column 4 combines all covariates. The coefficient on both tweets and night-lights are reduced, showing that there is some overlap in terms of what they explain in terms of economic activity. The coefficient on both variables is statistically significant (at the 5 percent level for tweets). The R² is 0.67 with an RMSE of 0.67.

In sum, both tweets and night-lights can be used to estimate economic activity at the city level. This could be extremely useful for countries where GDP data at the sub-national level (i.e.: regions, provinces and metropolitan areas) are not available.

8.2. Predicting GDP for cities

It is important and relevant that Twitter, together with night-lights, can be effectively used to estimate GDP at the sub-national level. But it is even more useful if these proxies could help predict GDP in these areas. For this, I see whether Twitter and night-light data can accurately predict out-of-sample estimates for US MSA's.

The first part of this exercise is similar to what was done in [subsection 4.2](#). I use MSA-level data and carry out both a leave-one-out as well as a k-fold cross-validation resampling method. Results are presented in [Table 9](#). The average coefficients for both these methods are quantitatively similar to those presented in [Table 8](#) for the entire dataset; indicating that the model is stable. These out-of-sample predictions are done using twitter, population and night-light data. In the leave-one-out method, the R² is 0.676 and the RMSE is 0.671. Similar results are shown for the k-fold cross-validation. This show that both these proxies

²⁶ Note that MSA's do not span the entirety of the US territory. Thus a small percentage of tweets (3.8 percent) are posted from locations that do not correspond to any MSA.

do a good job at predicting out-of-sample GDP estimates for MSA: the model explains roughly two-thirds of variation in GDP across MSA.

[Fig. 14](#) shows the results for this exercise for the k-Fold method with 10,000 repetitions. The bulk of the out-of-sample predictions are clustered closely around the 45-degree line, indicating that the majority of these do not fall far from the official estimates. This seems to be the case across the income distribution.

8.3. Economic activity across national borders

Another exercise that can be done using Twitter data is to estimate economic activity in areas around national borders. It is not possible to estimate economic activity in these areas across national borders using official GDP estimates because even though some countries do publish regional GDP estimates at the sub-national level, each survey is conducted independently with different questions and methodologies and thus cannot be taken together. Social media data, in this case Twitter, can help us around this issue: we can aggregate all the Twitter posts from a specific geographic area and take them together to use as a proxy for economic activity.

The border between Mexico and the US represents an interesting example to carry out this exercise as they are important trade partners and there is a great deal of economic synergy between the two countries. In 2012, US exports to Mexico totaled \$215,875 million, which represents 14 percent of total exports in the US and 49 percent of total imports for Mexico. On the other hand, Mexico exports to the US totaled \$323,026 million, which represents 74 percent of total exports for Mexico and 14 percent of total imports in the US.

[Fig. 15](#) shows a density map of economic activity in the US-Mexican border based on tweets. There are clear economic regions that go beyond the international border, particularly in the region close to the Gulf Coast. This shows that the area around Houston, San Antonio and Austin in Texas concentrate the largest level of economic activity which extends south of the border and into important cities in Mexico, like Monterrey and San Luis Potosí. The same phenomenon is visible in the Pacific Ocean side of the border, with a cluster of economic activity integrating US cities close to the border like Los Angeles and San Diego, with border cities in Mexico like Tijuana. These exercises are useful not only for businesses trying to make decisions in the region, but also for scholars who are interested in studying cross-border spillover effects in border-cities.²⁷

8.4. Local economic shock: Vaca Muerta, Argentina

The previous exercises have showcased the advantages of using Twitter data to estimate economic activity in terms of the continuous time and the geographic granularity of social media posts. In some scenarios, both characteristics can be combined to produce a timely and geographically detailed estimate of economic activity. This is extremely useful for studying local economic shocks, such as factory openings ([Greenstone et al. \(2010\)](#)) and closings, discovery of natural resources or natural disasters ([Kocornik-Mina et al. \(2015\)](#)). For example, nowcasting the impact of such events could be extremely useful for the government in estimating how much recovery aid they need to send to a certain affected area. Official data is not well equipped for such tasks because of delays in estimates and because many times they are not geographically granular enough to perceive local economic shocks.

One recent such case of a discovery of natural resources that has transformed a local economy occurred in the province of Neuquén, Argentina. The Vaca Muerta Formation is a geologic formation which has had oil production since 1918. But in July 2011, a large oil discovery was made, with resources estimated at 16 billion barrels of oil and

²⁷ For example, see [Hanson \(2001\)](#) and [Coronado et al. \(2015\)](#).

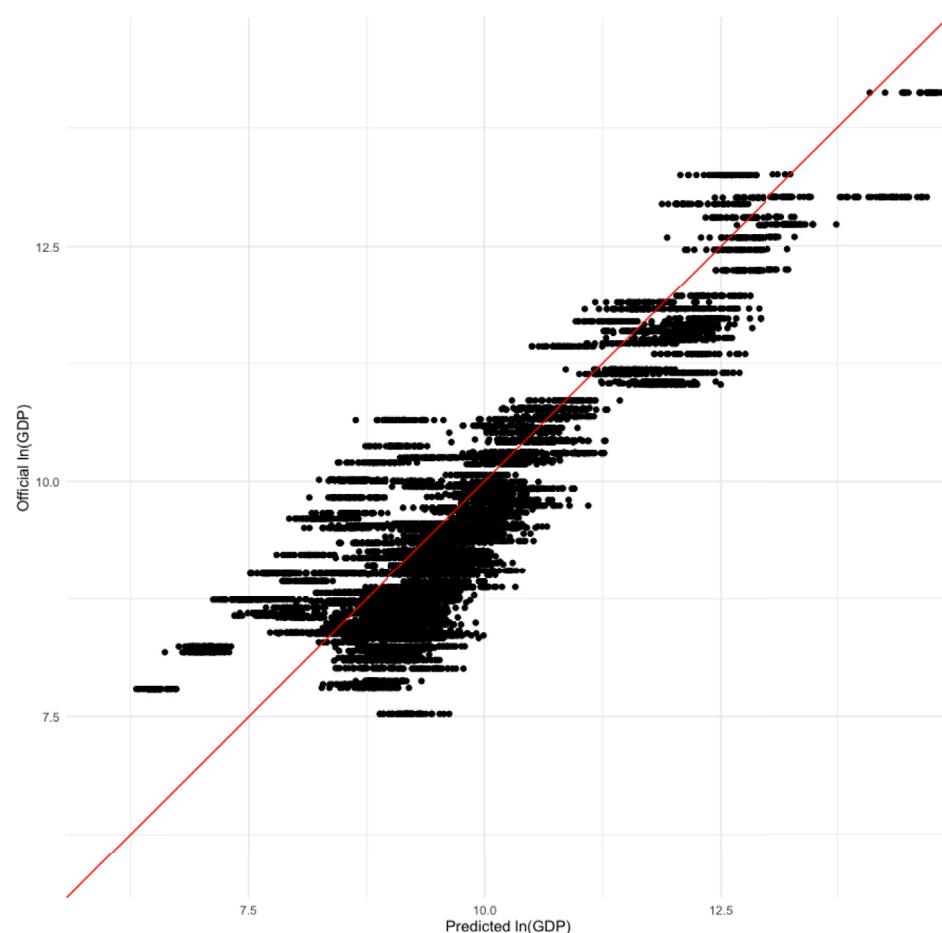


Fig. 14. Predicting GDP for Metropolitan Statistical Areas in the US.

Notes: Red line represents 45 degree line, which indicates where GDP predictions based on tweets, population and night-lights per MSA and official estimates for GDP are equal. This is done using a k-Fold cross-validation with $k = 10$ and based on 10,000 resamples. (For interpretation of the references to colour in this figure legend, the reader is referred to the Web version of this article.)

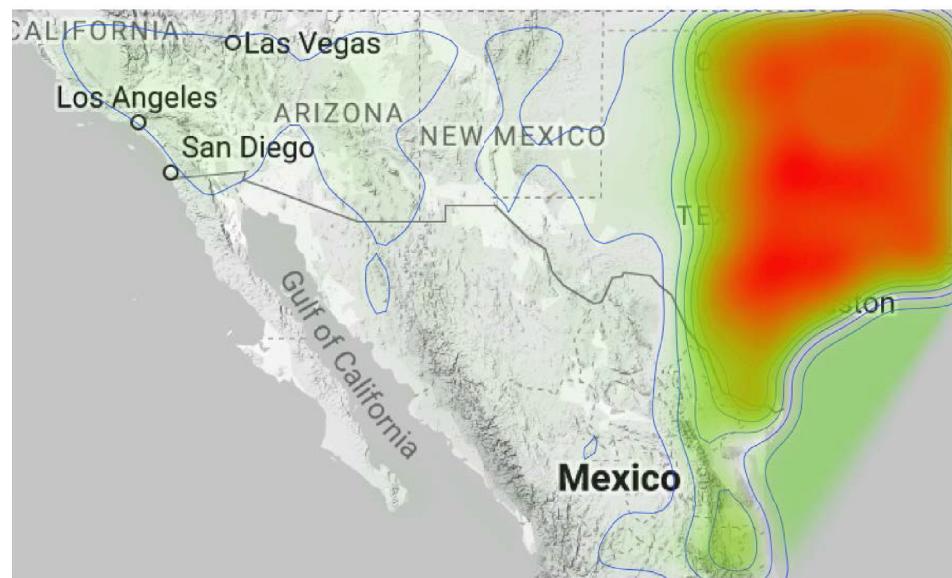


Fig. 15. Density Map for Economic Activity in the US-Mexico Border in 2012.

Notes: Contours contain areas with more economic activity with areas shaded in red representing the highest concentration of economic activity. (For interpretation of the references to colour in this figure legend, the reader is referred to the Web version of this article.)

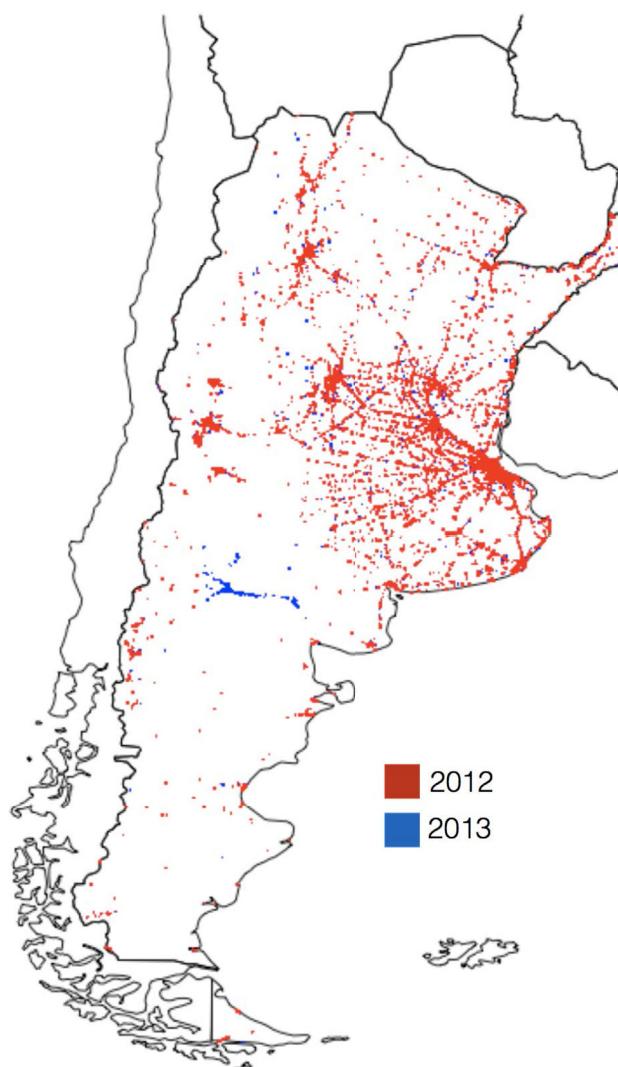


Fig. 16. Tweets in Argentina in 2012-13.

Notes: Figure shows distribution of tweets in Argentina for 2012 and 2013. Given that there was a large increase in the number of Tweets from 2012 to 2013 nationwide, I randomly selected a number of tweets for 2013 equivalent to that of 2012. Tweets from 2012 plotted appear on top to show the difference in distribution of tweets between the two years that is present in the area around Vaca Muerta.

308 trillion cubic feet of gas. If exploited, the proven reserves of the country would increase more than eight times.

Beyond the impact this caused in the energy market of the country at large, the local region was particularly disrupted by the migration, investment, changes in infrastructure and overall economic shock that such a discovery had on the small local economy. The closest town to Vaca Muerta is called Añelo, which in 2010 had 2500 inhabitants. Since then, informal estimates have reported the population to more than double and rent prices to skyrocket. There are newspaper articles claiming that prices of sq/ft in Añelo are comparable to the high end neighborhoods in Buenos Aires. But none of this has been corroborated or certified by official statistics. Argentina does produce provincial GDP estimates, but the local economic effects Vaca Muerta has on a small town like Añelo do get damped by changes in economic activity in larger cities when looking at provincial estimates. This is where we can exploit the advantage that social media offer granular geographic data, which allows us to aggregate tweets at any region of interest, even particularly small ones.

In order to estimate the growth of economic activity in the region surrounding Vaca Muerta, I identify all social media posts being sent from that region in 2012 and 2013. Fig. 16 shows the distribution of tweets in Argentina for each of the two years.²⁸ This figure highlights how the distribution of tweets in Argentina from one year to the other solely differ in the region surrounding Vaca Muerta. In order to get an estimate of the change in economic activity in this region between 2012 and 2013, I use the elasticity between tweets and GDP from Table 7. Based on the difference in volume of tweets originating from the region, I find that economic activity grew by 38.1 percent from one year to the other.²⁹ As far as I am aware, this is the first time the local economic effect of Vaca Muerta on the region has been quantified.

9. Conclusion

The main goal of this paper is to study whether social media data, in this particular case from Twitter, could be used as a proxy for estimating GDP. Using 140 million geo-located tweets from across the world, I find first that the volume of tweets is an accurate predictor of economic activity at the country level. I then extend this exercise to accurately calculate economic activity for different sub-national regions to exploit the precise geographical information retrieved from social media posts. A supplementary goal of this paper is to link the use of social media data to the widely-studied night-light proxy. All exercises carried out in this paper have consistently shown that when both proxies are used together the predictive power of the model is greatly increased. Furthermore, given that the coefficients of each of the proxies remain statistically significant when used together, they seem to be capturing different aspects of economic activity.

An important concern with the use of proxies is that they typically lack an underlying theoretical model that demonstrates the relationship between the proxy and the variable of interest. In this paper, I study the underlying relationship between tweets and economic activity and find that image tweets are a byproduct of consumption, with social media platforms being the ideal medium to showcase consumption among the network of users. Nonetheless, more research needs to be done in addressing this relationship before such a model is incorporated to estimate economic activity. Moreover, there are other social media platforms that might be better suited for estimating GDP or other relevant economic metrics. Hopefully this paper will spur such research.

There are several institutions and stakeholders that can benefit from such endeavours. For many developing countries, GDP estimates are produced annually and solely at the country level. Being able to produce sub-national estimates of economic activity at a relatively low cost could be tremendously helpful for country officials and international organizations to evaluate certain policies. For example, measuring the economic impact of *maquiladoras* in special economic zones, or studying the effects of unanticipated discoveries of natural resources in a particular region. Even in developed countries, having a reliable proxy to measure economic activity in a timely manner at a more granular level and unbound by political borders, is useful to study local economic activity. These estimates could help identify areas that are gentrifying as well as evaluate the recovery from local economic shocks, such as

²⁸ Given that the number of tweets grew substantially between 2012 and 2013 (see Table 1), I restrict the 2013 tweets to a random sample of tweets equivalent to the number of tweets sent out in 2012.

²⁹ To corroborate this estimate, I carry out the same analysis using the change in night-light intensity over the region between 2012 and 2013. Using the coefficients from Table 6, I find that economic activity grew 27.5 percent over this time period. Finally, using data and elasticities from both Twitter and night-lights, I find an increase in economic activity in the region of 31 percent. Without an official estimate to use as a benchmark it is hard to say which estimate is more accurate (or if either of them are), but given that they show a significant increase in economic activity, it expands on the notion that both these measures could be used together to represent a more accurate proxy.

Table 10
Mechanism between tweets and economic activity.

Dep. var.:ln(Consumption)	(1)	(2)	(3)	(4)	(5)	(6)
ln(All Tweets)	0.51*** (0.07)	0.15* (0.02)				
ln(Leisure Tweets)			0.51*** (0.07)	0.22** (0.02)	0.93** (0.44)	0.56*** (0.08)
ln(Work Tweets)					-0.42 (0.43)	-0.09 (0.08)
ln(Population)		1.04*** (0.03)		1.04*** (0.03)		1.03*** (0.03)
R ²	0.54	0.99	0.55	0.99	0.56	0.99
Adj. R ²	0.53	0.99	0.54	0.99	0.54	0.99
Partial R ² Leisure Tweets	-	-	0.55	0.42	0.55	0.44
Num. obs.	48	48	48	48	48	48
RMSE	0.68	0.12	0.67	0.12	0.67	0.12

***p < 0.01, **p < 0.05, *p < 0.10.

Notes: This table uses all image tweets sent from the US in 2012. Tweets are aggregated at the state level and categorized into one of two labels depending on the time in which they were posted. Work tweets includes all tweets sent during weekdays between 8am and 5pm. Leisure tweets includes tweets sent during the weekday before 9am and after 5pm as well as tweets sent on weekends, irrespective of time.

natural disasters. The data needed to evaluate these crucial questions might already be available.

Declaration of competing interest

None.

A. Supplemental data

A.0.1. Socio-economic data

The World Bank provides freely and publicly available data on various relevant socio-economic indicators at the country level. Given that one of the main objectives of this paper is to provide an effective proxy for estimating GDP that allows for more transparency in official statistics, it is important that all the data used in this paper is publicly available and thus could be replicated by individuals and institutions. Besides from current GDP in USD, I also obtain total population for each country from the World Bank database.

Another indicator I obtain from the World Bank is the percent of the population that uses the internet. Given that Twitter requires internet service access to establish a connection, the penetration of the internet in a given country is a useful variable to include in our baseline regression.

As described in section 2 official GDP estimates suffer several issues, particularly in developing countries. If this is in fact the case, I could be estimating a reported GDP that is not in fact the *true* GDP. Thus if the estimates using Twitter are imprecise, it could be in part because of measurement error in the GDP calculations I am trying to estimate in the first place. The World Bank produces a composite score assessing the capacity of a country's statistical office. In particular they focus on three specific areas: methodology, data sources, and periodicity and timeliness. The overall score is a simple average of all three area scores

on a scale of 0–100, where higher values indicate higher quality data. In section ?? I use these data quality scores to see if the discrepancies in our estimates are larger for countries with inferior data quality, as assessed by the World Bank.

A.0.2. Night-time light satellite data

Given that much of the previous literature on alternative ways of estimating GDP has come from night-lights from satellite data, I contrast and compare the performance of satellite and Twitter data as a proxy for GDP (see section 5). The United States Air Force Defense Meteorological Satellite Program satellite orbits the earth roughly 14 times a day, taking images that record the luminous intensity radiated from the earth. Although the main purpose of this task is to detect moonlit clouds, a useful byproduct is that they also capture lights emitted from human settlements.

Scientists then process these images and perform a series of tasks (i.e.: remove intense sources of natural light during summer months, auroral activity, days where cloud cover obstructs the earth's surface, etc.) that leave only man-made light visible. They then average all valid images over the year and report the intensity of light for approximately every 0.86 square kilometer. The intensity figure is an integer between 0 and 63, where higher values indicate more light. These datasets are made publicly available by the National Oceanic and Atmospheric Administration's National Geophysical Data Center. These are also the datasets used in the majority of the papers referred to in section 5 that use night-light as a proxy for GDP. I downloaded these datasets for 2012 and 2013, the same years for which I have complete Twitter data.

References

- Antenucci, Dolan, Cafarella, Michael, Levenstein, Margaret, Ré, Christopher, Shapiro, Matthew D., March 2014. Using Social Media to Measure Labor Market Flows. NBER Working Papers 20010. National Bureau of Economic Research, Inc.
 Askatas, Nikolaos, Zimmermann, Klaus F., July 2013. Nowcasting business cycles using toll data. J. Forecast. 32 (4), 299–306.

- Banerjee, Abhijit, Duflo, Esther, Qian, Nancy, March 2012. On the Road: Access to Transportation Infrastructure and Economic Growth in China. Working Paper 17897. National Bureau of Economic Research.
- Becker, Gary S., 1965. A theory of the allocation of time. *Econ. J.* 75 (299), 493–517.
- Brender, Adi, Drazen, Allan, 2005. Political budget cycles in new versus established democracies. *J. Monetary Econ.* 52 (7), 1271–1295.
- Chen, Xi, Nordhaus, William D., August 2010. The value of luminosity data as a proxy for economic statistics. In: Cowles Foundation Discussion Papers 1766, Cowles Foundation for Research in Economics. Yale University.
- Coronado, Roberto, De León, Marycruz, Saucedo, Eduardo, 2015. So close to Mexico: economic spillovers along the Texas-Mexico border. In: Ten-Gallon Economy. Palgrave Macmillan, New York, pp. 183–198.
- Dauth, Wolfgang, Haller, Peter, 2020. Is there loss aversion in the trade-off between wages and commuting distances? *Reg. Sci. Urban Econ.* 83, 103527.
- Davis, Morris A., Oliner, Stephen D., Pinto, Edward J., Bokka, Sankar, 2017. Residential land values in the Washington, DC metro area: new insights from big data. *Reg. Sci. Urban Econ.* 66, 224–246.
- Dawson, John W., Dejuan, Joseph P., Seater, John J., Frank Stephenson, E., 2001. Economic information versus quality variation in cross-country data. *Can. J. Econ. Revue canadienne d'économique* 34 (4), 988–1009.
- Doll, Christopher, 2008. CIESIN Thematic Guide to Night-Time Light Remote Sensing and its Applications. 01. .
- Doll, Christopher N.H., Muller, Jan-Peter, Morley, Jeremy G., April 2006. Mapping regional economic activity from night-time light satellite imagery. *Ecol. Econ.* 57 (1), 75–92.
- Donaldson, Dave, Storeygard, Adam, November 2016. The view from above: applications of satellite data in economics. *J. Econ. Perspect.* 30 (4), 171–198.
- Ghani, Ejaz, Goswami, Arti Grover, Kerr, William R., 2016. Highway to success: the impact of the golden quadrilateral project for the location and performance of Indian manufacturing. *Econ. J.* 126 (591), 317–357.
- Ghosh, Tilottama, Powell, Rebecca L., Elvidg, Christopher D., Baugh, Kimberly E., Sutton, Paul C., Anderson, Sharolyn, December 2010. Shedding light on the global distribution of economic activity. *Open Geogr. J.* 3 (1).
- Glaeser, Edward L., Kim, Hyunjin, Luca, Michael, November 2017. Nowcasting the Local Economy: Using Yelp Data to Measure Economic Activity. Working Paper 24010. National Bureau of Economic Research.
- Greenstone, Michael, Hornbeck, Richard, Moretti, Enrico, 2010. Identifying agglomeration spillovers: evidence from winners and losers of large plant openings. *J. Polit. Econ.* 118 (3), 536–598.
- Han, Bing, Hirshleifer, David, Walden, Johan, February 2019. Visibility Bias in the Transmission of Consumption Beliefs and Undersaving. Working Paper 25566. National Bureau of Economic Research.
- Hanson, Gordon H., 2001. U.S.–Mexico integration and regional economies: evidence from border-city pairs. *J. Urban Econ.* 50 (2), 259–287.
- Henderson, J., Vernon, Storeygard, Adam, Weil, David N., April 2012. Measuring economic growth from outer space. *Am. Econ. Rev.* 102 (2), 994–1028.
- Henderson, J., Vernon, Nigmatulina, Dzhamilya, Kriticos, Sebastian, 2019. Measuring urban economic density. *J. Urban Econ.* 103188.
- James, Gareth, Witten, Daniela, Hastie, Trevor, Tibshirani, Robert, 2013. An Introduction to Statistical Learning: with Applications in R. Springer.
- Jerven, Morten, 2013. Poor Numbers: How We Are Misled by African Development Statistics and what to Do about it. Cornell University Press.
- Johnson, Simon, Larson, William, Papageorgiou, Chris, Subramanian, Arvind, 2013. Is newer better? Penn World Table Revisions and their impact on growth estimates. *J. Monetary Econ.* 60 (2), 255–274.
- Kerner, Andrew, Crabtree, Charles, January 2018. The IMF and the Political Economy of Data Production. *SocArXiv*.
- Kocornik-Mina, Adriana, McDermott, Thomas K.J., Michaels, Guy, Rauch, Ferdinand, December 2015. Flooded Cities. GRI Working Papers 221. Grantham Research Institute on Climate Change and the Environment.
- Lee, Kevan, November 2015. What Analyzing 1 Million Tweets Taught Us. Retrieved from <https://thenextweb.com/socialmedia/2015/11/03/what-analyzing-1-million-tweets-taught-us/>. (Accessed 2 October 2018).
- Llorente, Alejandro, García-Herranz, Manuel, Cebrán, Manuel, Moro, Esteban, 2015. Social media fingerprints of unemployment. *PloS One* 10 (5), 1–13.
- Masood, Ehsan, September 2014. The Great Invention: the Story of GDP and the Making (And Unmaking) of the Modern World. Saqi Books.
- Michalopoulos, Stelios, Papaioannou, Elias, 2018. National institutions and subnational development in Africa. *Q. J. Econ.* 129 (1), 151–213.
- Pinkovskiy, Maxim L., 2017. Growth discontinuities at borders. *J. Econ. Growth* 22 (2), 145–192.
- Pinkovskiy, Maxim, Sala-i-Martin, Xavier, 2016. Lights, camera ... income! Illuminating the national accounts-household surveys debate. *Q. J. Econ.* 131 (2), 579–631.
- Rao, B. L. S. Prakasa, 1992. Identifiability in Stochastic Models: Characterization of Probability Distributions. Academic Press (Google-Books-ID: nRLvAAAAMAAJ).
- Sommervoll, Ávald, Sommervoll, Dag Einar, 2019. Learning from man or machine: spatial fixed effects in urban econometrics. *Reg. Sci. Urban Econ.* 77, 239–252.
- Sutton, Paul C., Elvidge, Christopher D., Ghosh, Tilottama, May 2007. Estimation of gross domestic product at sub-national scales using nighttime satellite imagery. *Int. J. Ecol. Econ. Stat.* 8 (S07), 5–21.
- Veblen, Thorstein, 1899. The Theory of the Leisure Class. McMaster University Archive for the History of Economic Thought.
- Weidemann, C., Swift, Jennifer, 2013. Social media location intelligence: the next privacy battle - an ArcGIS add-in and analysis of geospatial data collected from Twitter.com. *Int. J. Geoinform.* 9 (2), 21–27.