

Collegio Carlo Alberto



Wikipedia Matters

Marit Hinnosaar

Toomas Hinnosaar

Michael Kummer

Olga Slivko

No. 508

November 2017

Carlo Alberto Notebooks

www.carloalberto.org/research/working-papers

Wikipedia Matters*

Marit Hinnosaar[†] Toomas Hinnosaar[‡] Michael Kummer[§]
Olga Slivko[¶]

September 29, 2017

Abstract

We document a causal impact of online user-generated information on real-world economic outcomes. In particular, we conduct a randomized field experiment to test whether additional content on Wikipedia pages about cities affects tourists' choices of overnight visits. Our treatment of adding information to Wikipedia increases overnight visits by 9% during the tourist season. The impact comes mostly from improving the shorter and incomplete pages on Wikipedia. These findings highlight the value of content in digital public goods for informing individual choices.

JEL: C93, H41, L17, L82, L83, L86

Keywords: field experiment, user-generated content, Wikipedia, tourism industry

1 Introduction

Asymmetric information can hinder efficient economic activity. In recent decades, the Internet and new media have enabled greater access to information than ever before. However, the digital divide, language barriers, Internet censorship, and technological constraints still create inequalities in the amount of accessible information. How much does it matter for economic outcomes?

In this paper, we analyze the causal impact of online information on real-world economic outcomes. In particular, we measure the impact of information on one of the primary economic decisions—consumption. As the source of information, we focus on Wikipedia. It is one of the most important online sources of reference. It is the fifth most

*We are grateful to Irene Bertsek, Avi Goldfarb, Shane Greenstein, Tobias Kretschmer, Thomas Niebel, Marianne Saam, Greg Veramendi, Joel Waldfogel, and Michael Zhang as well as seminar audiences at the Economics of Network Industries conference in Paris, ZEW Conference on the Economics of ICT, and Advances with Field Experiments 2017 Conference at the University of Chicago for valuable comments. Ruetger Egolf, David Neseer, and Andrii Pogorielov provided outstanding research assistance. Financial support from SEEK 2014 is gratefully acknowledged.

[†]Collegio Carlo Alberto and CEPR, marit.hinnosaar@gmail.com

[‡]Collegio Carlo Alberto, toomas@hinnosaar.net

[§]Georgia Institute of Technology, michael.kummer@econ.gatech.edu

[¶]ZEW, slivko@zew.de

popular website in the world¹ and receives about 14 billion direct page views per month.^{2,3} However, the information available across Wikipedia’s 299 language editions is not the same. We analyze whether the differences in available information affect consumption choices.

We quantify the causal impact of information in Wikipedia on consumption choices, by conducting a randomized field experiment. Analyzing the impact of information using observational data would have been challenging, because of potential endogeneity. Popular products tend to attract more attention, and therefore, more information is available about them. While the amount of information on Wikipedia tends to be correlated with the products’ popularity, the information isn’t necessarily causing consumption, but may instead be its byproduct. We overcome the identification problem using randomization.

We added content to randomly chosen Wikipedia pages in randomly chosen languages. We measured the outcome using data on tourists’ overnight hotel stays in Spain. The Spanish tourism sector is important in itself by accounting for almost 5% of Spain’s GDP.⁴ It also provided a good setting for the study, since the Spanish National Statistical Institute collects information about overnight stays in Spanish hotels at the level of city, month, and tourist country of origin. Our treatment added text and photos to the Wikipedia pages of Spanish cities in different language editions of Wikipedia. The added text was translated mainly from the Spanish Wikipedia. The text was on topics relevant to tourists, such as the city’s main sights and culture. We focused our attention on cities with rather short Wikipedia pages. The randomization was done across city and language pairs. By varying the information in different language editions of Wikipedia, we can isolate the causal impact on tourists’ choices.

We find that information on Wikipedia has a sizable impact on consumption choices. Our estimates show that adding about 2000 characters (approximately two paragraphs) of text and one photo to a city’s Wikipedia page increased the number of nights spent in this city by about 9% during the tourist season compared to cities in the control group.⁵ The effect comes mostly from pages that were initially relatively incomplete. In particular, the treatment increases hotel stays by about 33% in cities which initially had very short pages in a particular language, while there was no effect on city-language combinations, where the pages were well developed.

Using data on readership from Wikipedia page views and search activity from Google Trends, we can shed some light on the mechanism that drives our findings. The added information has no significant impact on search activity outside Wikipedia but significantly increases the articles’ readership. That is, more detailed Wikipedia articles gain more attention from potential readers. The size of this effect is similar in magnitude to the effect on tourists’ choices.

¹Alexa Internet. <http://www.alexa.com/siteinfo/wikipedia.org>, accessed September 23, 2017.

²Page Views for Wikipedia. Wikimedia Statistics. <https://stats.wikimedia.org/EN/TablesPageViewsMonthlyCombined.htm>, accessed September 23, 2017.

³This does not include indirect uses such as Apple’s Siri or Google.

⁴Tourism statistics. Eurostat. http://ec.europa.eu/eurostat/statistics-explained/index.php/Tourism_statistics, accessed June 21, 2017.

⁵Our experiment doesn’t allow us to distinguish between absolute increase in demand and substitution between control and treatment. Some of the effect likely arises from rerouting tourists from other cities. The implications we highlight in this paper hold in either case.

Our results have three policy implications, which are likely to reach beyond the setting of our experiment. First, the results have implications on economic inequality and the digital divide. Language can pose barriers that hinder efficient economic activity. Language barriers have slowed innovation (Peri, 2005), decreased trade (Anderson and van Wincoop, 2004), and affected investments (Grinblatt and Keloharju, 2001). In particular, languages create a major obstacle to access to information. Large differences remain across languages in terms of information available online. Our results imply that these differences may lead to significant differences in economic behavior between various groups.

Second, on the macroeconomic level we show that online user-generated content can have a significant causal impact on economic behavior and economic outcomes. The treatment increased the number of hotel visits by 9%. If we extend this to the entire tourism industry, the impact is large. In 2015, international tourists spent 270 million nights in Spain. The same year international travel receipts equaled 51 billion euros in Spain and 116 billion in the EU.⁶ While we cannot say whether online user-generated content is changing the size of expenditures or reallocating them, it could affect the choices in the order of billions of euros.

Third, on the microeconomic level, our results highlight the importance of online presence. A 9% increase in consumption as a result of additional user-generated information is quite large, given that each international tourist spends about 101 euros per day while visiting Spain on average (García-Sánchez, Fernández-Rubio, and Collado, 2013). The findings suggest that it is beneficial to ensure that a city, firm, or product is accurately represented online in all relevant languages.

The results of this paper pose a puzzle—why is the online presence so limited? Increasing online presence is relatively inexpensive, while our results suggest a high return on investment. The online presence puzzle differs from most of the literature examining contributions to online public goods. This literature finds that contributions exceed what the economic theory would suggest. While the public goods literature assumes contributions are altruistic, we concentrate on a setting where the involved parties would benefit from making more information available.

Our paper makes three methodological contributions. First, it is among the first papers to use Wikipedia as a treatment in a field experiment for studying the impact on behavior outside Wikipedia.⁷ Wikipedia provides a good ground for this, since anyone can freely improve it⁸ and the whole process is automatically recorded in the form of revision histories. Moreover, readership of Wikipedia articles is well-recorded in the form of page views.

Second, we use a novel dataset of real-life outcomes—overnight hotel stays. Most importantly, this dataset provides a precise measure of demand of an identical product for consumers from different countries. In Spain, hotels are legally required to record guests' country of residence. We obtained the data from the Spanish National Statistical Institute aggregated to monthly level for each city and each country of origin. For example, we know how many nights German tourists spent in a particular city in July 2015. We

⁶Source: Tourism statistics. Eurostat.http://ec.europa.eu/eurostat/statistics-explained/index.php/Tourism_statistics, accessed June 21, 2017.

⁷There is a literature examining the editing behavior in Wikipedia, which we will review below.

⁸Following Wikipedia's Terms of use and policies.

use the fact that German tourists are more likely to get their information from German Wikipedia and Italian tourists from Italian Wikipedia to map consumption choices back to their potential information sources.

Finally, we make a technical contribution in analyzing Wikipedia’s revision histories. As our treatment adds information to Wikipedia pages, which can then be changed by other Wikipedia users, the first step in the analysis is to see how much of our additions are modified by other Wikipedia users over time. For this, we use a diff algorithm describing the shortest sequence of additions and deletions of characters to change the original text to the revised one.⁹ We apply this algorithm twice. First, to quantify which parts of the page our experiment added, and second, to measure how much of our additions had survived after a few months. We find that our edits are rather persistent: about 93% of our added text still existed about four months after the treatment. This could be because information on the pages we edited was relatively scarce and (hopefully) our contributions were considered sufficiently valuable by the Wikipedia community.

Our paper contributes to media economics literature studying the impact of media on economic outcomes (for an overview see DellaVigna and Ferrara (2016)). In particular, our paper adds to studies on the impact of media on consumption. Most notably, Bursztny and Cantoni (2015) use geographic variation in access to Western TV to study its long-run impact on East German consumption choices. The paper also contributes to studies on the impact of new media and online user-generated content.¹⁰ Among others Chevalier and Mayzlin (2006) and Luca (2011) study how product reviews affect sales. Enikolopov, Petrova, and Sonin (2017) analyze the impact of blog posts exposing corruption in state-controlled companies on their market returns. Xu and Zhang (2013) study the impact of Wikipedia on financial markets combining data of financial records, management disclosure records, news article coverage, and Wikipedia editing histories. Our paper adds to the literature by providing evidence of how Wikipedia informs consumers and affects their choices. It differs from these papers in terms of the research method. The above papers use either a natural experiment or detailed observational data, while we conduct a randomized field experiment which helps us to identify the effect.

Methodologically, our paper is related to a recent study by Thompson and Hanley (2017). In a work independent from ours, they also conduct a randomized field experiment in Wikipedia. They find that Wikipedia content affects scientific articles. Their work is complementary to ours—they find that Wikipedia has a significant impact on knowledge production outside Wikipedia, whereas we find that the available information affects consumption choices.

Our paper also relates to the emerging small branch of literature on information production in Wikipedia. Most of this literature analyzes contributions to Wikipedia (including Zhang and Zhu, 2011; Aaltonen and Seiler, 2015) and biases in Wikipedia (Greenstein

⁹For a description of the algorithm, see Myers (1986).

¹⁰More generally, our paper relates to the literature on how ICT affects economic outcomes by changing access to information. Among other topics, this literature has studied the impact of Internet on economic growth (Czernich, Falck, Kretschmer, and Woessmann, 2011), on labor market outcomes (Forman, Goldfarb, and Greenstein, 2012; Akerman, Gaarder, and Mogstad, 2015), on the airline industry (Dana and Orlov, 2014; Ater and Orlov, 2015), the impact of medical records on hospital costs (Dranove, Forman, Goldfarb, and Greenstein, 2014); and the impact of e-commerce on price dispersion (Overby and Forman, 2014).

and Zhu, 2012; Greenstein, Gu, and Zhu, 2016; Greenstein and Zhu, 2017). Our paper stresses the importance of understanding the Wikipedia production process and its biases by quantifying the impact of Wikipedia on offline economic behavior.

2 Background on Wikipedia

Wikipedia is a free-access Internet encyclopedia. It is the fifth most popular website in the world.¹¹ It is arguably one of the most important knowledge repositories and digital public goods. Wikipedia is written by volunteers: anyone can create Wikipedia articles or edit almost any of its existing articles.

While Wikipedia exists in 299 languages, the amount of available information differs across languages. English Wikipedia is the largest, with over five million articles. Only 13 other language editions have more than a million articles.¹²

A significant share of the population can access information only in their mother tongue. Almost half of the population in the EU does not speak any foreign language.¹³ They can only access the information from their local language Wikipedia. Figure A.1 shows local language Wikipedia sizes and the percentage of the population speaking more than one language. Language affects not only the topics covered, but also the depth of coverage. For example, among the 1000 most important articles in Wikipedia¹⁴ the median text length (relative to the corresponding page in English) varies from 5% in Latvian to 55% in French (see figure A.2). Not all topics are covered equally (see figure A.3). Overall, the worst covered topics are in categories like philosophy and religion (12%) and health and medicine (13%).

The relevant implication for this paper is that the amount of information available in each language edition of Wikipedia is not the same. It varies both in terms of the pages that exist and the depth of coverage on each topic. Figure 1 presents an example of information about a city. It describes pages about Murcia, a large Spanish city, across the different language editions of Wikipedia. This page exists in 84 different language editions of Wikipedia.¹⁵ The figure contrasts the length of the Murcia page in the 20 languages in which the page is the longest. Because it is a Spanish city, the page is the longest in Spanish Wikipedia. In all other language editions the page is at least five times shorter.

3 Experimental design

We conducted a field experiment in which we added content (text and photos) to the Wikipedia pages of Spanish cities in different language editions of Wikipedia. The ran-

¹¹Only Google, Youtube, Facebook, and Baidu are more popular than Wikipedia. The popularity is measured by the web traffic measurement company Alexa Internet (<http://www.alexa.com/siteinfo/wikipedia.org>, accessed June 19, 2017).

¹²https://meta.wikimedia.org/wiki/List_of_Wikipedias, accessed June 19, 2017.

¹³About 46% of the population speaks only their mother tongue. (cf. Eurobarometer (2012)).

¹⁴Wikipedia keeps a list of 1000 vital articles (https://en.wikipedia.org/wiki/Wikipedia:Vital_articles, accessed June 26, 2017).

¹⁵Wikipedia data on Murcia was accessed on June 20, 2017.

domization was done across city and language pairs. The outcome variable is the number of overnight hotel stays by the tourists from the countries where the population speaks one of the treated languages. The experimental design is discussed in detail below.

Sample We restricted attention to four languages and tourists from the corresponding countries: Dutch (the Netherlands), German (Germany), French (France), Italian (Italy). Altogether we had hotel data from 135 Spanish cities. However, in many of these cities, hotel data was missing for some months and some tourist countries of origin. Hence, we expected to encounter the problem of not being able to measure the effect of treatment because of missing outcome (hotel) data. We were also concerned that our fixed length treatment might not be strong enough in the case of cities which already had long Wikipedia pages.

Therefore, we restricted attention to a sample of cities that satisfied two criteria. First, the Wikipedia page for the city had to be relatively short—no more than 24,000 characters in each of the four languages. Second, there could be no missing hotel data for the city. Specifically, we required the data on hotel stays to exist for each month from May to October 2013 and for all four countries. Sixty cities satisfied these two criteria. This gave us a sample of 240 Wikipedia pages (or city-language pairs).

Randomization We randomized across 240 Wikipedia pages (60 Spanish cities in four languages). Our goal was to treat each city equally. Therefore, for each city, we treated its page in two randomly chosen language editions of Wikipedia. In each language edition of Wikipedia, we treated 30 city pages. This resulted in a design where, for each city, some languages are assigned to the treatment and some to the control group. Similarly, in each language, some cities are in the treatment and some in the control group.

To ensure balance in the treatment and control groups, we used a stratified randomization design. We ordered the 60 cities by the total number of tourists. Then we divided the cities into ten groups of six cities each. Within each group, we randomly assigned the city to one of six treatments. The six treatments were as follows: treat the city page in one of the six possible language pairs (Dutch & German; Dutch & French; Dutch & Italian; German & French; German & Italian; French & Italian). Hence, 120 city pages were treated and 120 pages remained as controls.

Treatment The pages were treated mid-August, 2014. To the pages in the treatment group, text and photos were added. The added text and photos were on topics relevant for tourists, such as the main sights and culture. Added text was translated mostly from the corresponding Spanish or English language Wikipedia pages. Typically, the photos were also from these corresponding Wikipedia pages.

Our goal was to improve the Wikipedia pages. We did not decrease the quality of Wikipedia pages, for example, by deleting existing material. On the contrary, following Wikipedia’s policies, we added material that according to our understanding was knowledge already approved by the editors of Spanish Wikipedia.

Survival of added material While editing German, French, and Italian Wikipedia was not problematic, we were not successful in editing Dutch Wikipedia. Wikipedia

allows anyone to edit it. This also means that anyone can delete all or part of an article, or undo the latest changes by reverting to a previous version. All our additions to Dutch Wikipedia were deleted in less than 24 hours. That is, all Dutch Wikipedia pages were essentially untreated from the point of view of a person reading these Wikipedia pages or accessing these indirectly, e.g. through Apple’s Siri or Google information box. Therefore, we exclude all Dutch Wikipedia articles from our analysis. Note that the results do not significantly change if we consider all Dutch articles as non-treated.

Table 1 shows that in the German, French, and Italian Wikipedias, our added text and photos survived well. (The methodology for measuring the survival of our additions is described in Section B.) Of the added text, on average 96 percent had survived by the beginning of the month following treatment and 93 percent by the beginning of the year following treatment. We interpret this in two ways. First, the edits were sufficiently persistent to provide hope that many people had seen the information our treatment added. Strictly speaking, it is not necessary that the precise wording of our treatment survives—it is to be expected that the other Wikipedia editors improve any added contributions over time in terms of wording, references, or content. However, measuring the preserved content is more difficult than measuring the actual text. Second, we hope that our treatment additions were considered useful by fellow Wikipedia editors; otherwise, they would have either reversed the edits or further revised them.

Descriptive statistics Table 2 shows that there were no significant differences in the main characteristics between the treatment and control groups.

Table A.1 shows descriptive characteristics of treatment. The median treatment added about 2000 characters of text and one photo. The treatment added relatively more to pages that were initially shorter (see Figure A.4). The initial page length by language is described in Table A.2.

Figure A.5 presents the histogram of the logarithm of the number of hotel nights. There is a large variation in the number of hotel nights. Figure A.6 presents the percentage of missing data by calendar month. It describes seasonality with slightly above ten percent missing data from May to October and up to 40 percent in December and January.

4 Results

Empirical strategy Our goal is to estimate the impact of additional information in Wikipedia on hotel stays in the corresponding city by tourists from the corresponding country. The main outcome variable is the logarithm of the number of hotel nights that tourists from country (exposed to language) j spent in city i during month t . In our main analysis, we estimate the following difference-in-differences regression:

$$\log(Nights_{ijt}) = \alpha + \beta Treatment_{ijt} + \gamma X_{ijt} + CityLanguageFE_{ij} + \varepsilon_{ijt} \quad (1)$$

The variable of interest $Treatment$ equals one for the treated city-language pairs during the months after treatment and equals zero otherwise. The regression includes fixed effects for city-language pairs $CityLanguageFE_{ij}$ and time varying control variables, X_{ijt} . The time varying control variables include: first, an indicator for period after treatment interacted with language fixed effects to take into account tourist country of origin-specific

trends; second, an indicator for period after treatment interacted with city fixed effects to take into account city-specific trend; third, logarithm of number of tourists from Spain interacted with language fixed effects to take into account events in the city which lead to an overall increase in tourism. We cluster the standard errors by city-language pair. Due to the missing data problem discussed above, in the main analysis, we restrict the sample to May–October during each year 2010–2015.

Main results Table 3 presents the main results. According to the estimates in Column 1, the treatment increases the number of hotel nights on average by 9%. Column 2, adds an interaction of the treatment variable and an indicator for Wikipedia pages that were initially relatively short. The estimates in Column 2 show that our treatment increases hotel stays by about 33% in cities where the pages were initially very short in a particular language, while there was no effect on cities with longer pages. Column 3, tries to explain the result by interacting the treatment variable and an indicator for the Wikipedia pages to which we added relatively longer text compared to the initial text length. Recall that since the length of text added was about the same, the treatment was relatively larger on pages that were initially short (Figure A.4). The results in Column 3 confirm that the effect is larger on pages where the treatment was relatively larger.

Robustness Table 4 presents our robustness checks. Columns 1–5 repeat regression in Column 1 of Table 3, so the magnitudes of the estimates are comparable.

Column 1 substitutes missing observations by zeros (only for city-year pairs, where data exists for some month and tourist country of origin). It excludes the variables that measure the number of tourists from Spain because the number of tourists from Spain is also missing. The results are very similar.

Column 2 adds observations for tourists from the Netherlands and considers these all as non-treated. The results are very similar. Recall that half of the city pages in Dutch Wikipedia were assigned to treatment, but editing Dutch Wikipedia proved impossible (24h after treatment all the pages remained untreated). We could estimate the same regression and add a separate indicator variable that equals one for months after treatment only for Dutch pages assigned to treatment. The results regarding the treatment effect remain the same.

Columns 3 and 4 add the excluded months, and Column 4 substitutes missing observations by zeros.¹⁶ In Column 4, again, the variables that measure the number of tourists from Spain are excluded. The results are similar, but in Column 3, less statistically precise.

Column 5 adds additional controls, namely, the logarithm of the number of tourists from UK interacted by language. The variables that measure the number of tourists from Spain are excluded. The results are similar.

In Column 6, the dependent variable is the number of tourists from country j divided by the number of tourists from country j plus those from Spain and the UK. Again, the variables that measure the number of tourists from Spain are excluded. While the

¹⁶We substituted missing observations only for city-year pairs, when data exists for some month and tourist country of origin

results are not comparable in magnitude, the treatment effect is positive and statistically significant.

Mechanism We analyze the mechanism by which additional information on Wikipedia changes choices. We consider three main channels. First, additional information could increase conversion rate. That is, it could lead to a larger share of readers choosing the destination. Second, the information could increase the number of readers. Third, it could increase the underlying interest in the destination via indirect effects, such as word-of-mouth. We proxy the third channel using data from Google Trends. Google Trends data measures how often a particular city is searched for on Google by the population of a particular country. We can measure the combination of the first two channels using data on the page views of Wikipedia articles. Unfortunately, we don't observe whether this reflects one person reading the page many times or many people reading it once.¹⁷ Therefore, we cannot distinguish between a higher conversion rate and a larger audience.

Table 5 presents estimates of analogous regressions as equation 1. In Columns 1–3, the outcome variable is the logarithm of the number of page views of a Wikipedia page for city i in language j during month t . In Columns 4–6, the outcome variable is the Google Trend for city i from country j during month t . The estimates in Column 1 show that the treatment increased page views by about 11 percent. Column 2 separates the effect by the length of the article (before treatment), showing that the treatment effect is larger on shorter pages. Similarly, the regression results in Column 3 show that the treatment effect is larger on pages where our treatment added a relatively larger share of text (these tended to be shorter pages). The estimates in Columns 4–6 show that our treatment had no effect on Google Trends (Google Search volume). The robustness of these estimates is studied in Table A.3.

Altogether, these results show that our treatment increases article readership, and the effect is similar in magnitude to the effect on the number of hotel nights. We find no evidence that the Google Search volume increased. We conclude that the added content on Wikipedia increased demand mostly through additional readership.

Limitations Our study faces limitations and raises questions for future research. First, our experiment was not designed to distinguish between a substitution and an overall increase. We would expect that our estimated treatment effect is at least partly explained by substitution from other possible tourist destinations. It appears unlikely that more information about interesting destinations leads to a significant increase in the entire tourism sector. The implications highlighted in the paper apply regardless of this ambiguity, though it would be interesting to distinguish these two effects.

Second, there is a question of generalizability, as the results may be specific to the types of pages and languages used in the experiment. In our sample, the Wikipedia pages were relatively short. We would expect that additional content would have less impact when the relative improvement is small. Moreover, the presence of short Wikipedia pages partly reflected the fact that these cities were not the most popular destinations. We would expect that the impact of Wikipedia is smaller in the case of major tourist

¹⁷Wikipedia did not collect unique page views prior to 2015, therefore we cannot distinguish between new and returning readers.

attractions. On the other hand, these places were notable enough to have Wikipedia pages and to receive regular tourist flows. It is unlikely that additional information could lead tourists to destinations without interesting attractions. In the languages included in the experiment, Wikipedia editions are still among the largest with relatively large readerships. The availability of information in local languages is probably less important in countries where people are used to obtaining information in English. Additionally, the countries in the experiment send large tourist flows to Spain. This means there was already preference for Spain and left room for substitution that was discussed above. The absolute level of the treatment effect is likely to be smaller in case of languages and countries where Spain was not a popular tourist destination.

On a more positive note regarding generalizability, the impact of Wikipedia is unlikely to be specific to the tourism industry. Instead, we would expect that the information on Wikipedia affects choices and behavior in many domains.

5 Discussion

We found a significant causal impact of user-generated content on Wikipedia on real-life choices. The estimated effect suggests that a well-targeted two-paragraph improvement of Wikipedia may lead to a 9% increase in tourists' overnight visits. The median monthly number of hotel nights spent by tourists from the three effectively treated countries to the cities in the control group was about 3000 (during the six months from May to October). This implies an increase of about 270 nights per month. Even if there were no tourists in the remaining 6 months, this implies about 1,600 additional hotel nights per year.

What are the implications for the local economy? According to recent estimates (García-Sánchez, Fernández-Rubio, and Collado, 2013), each international tourist visiting Spain spends about 101 EUR per day on average. Back-of-the-envelope calculations suggest that improving a city's Wikipedia page can lead to approximately 160,000 euros of additional revenue per year. This implies a considerable impact on local hotels and the overall local tourist industry.

Our results highlight the importance of online presence. Ensuring that a city, firm, or product is accurately represented in online information sources of all relevant languages is relatively cheap, i.e. almost free or a few hundred dollars in mainly one-time costs. In comparison, the 9%-increase in demand is rather large, suggesting a high return to investment.

Finally, the amount of information available in different languages varies significantly. Our results imply that this may lead to large differences in economic decisions and economic outcomes as well. This opens up a more general discussion about economic inequality and the digital divide across cultural and ethnic groups.

References

AALTONEN, A., AND S. SEILER (2015): "Cumulative Growth in User-Generated Content Production: Evidence from Wikipedia," *Management Science*, 62(7), 2054–2069.

- AKERMAN, A., I. GAARDER, AND M. MOGSTAD (2015): “The Skill Complementarity of Broadband Internet,” *The Quarterly Journal of Economics*, 130(4), 1781–1824.
- ANDERSON, J. E., AND E. VAN WINCOOP (2004): “Trade Costs,” *Journal of Economic Literature*, 42(3), 691–751.
- ATER, AND E. ORLOV (2015): “The Effect of the Internet on Performance and Quality: Evidence from the Airline Industry,” *The Review of Economics and Statistics*, 97(1), 180–194.
- BURSZTYN, L., AND D. CANTONI (2015): “A Tear in the Iron Curtain: The Impact of Western Television on Consumption Behavior,” *The Review of Economics and Statistics*, 98(1), 25–41.
- CHEVALIER, J. A., AND D. MAYZLIN (2006): “The Effect of Word of Mouth on Sales: Online Book Reviews,” *Journal of Marketing Research*, 43(3), 345–354.
- CZERNICH, N., O. FALCK, T. KRETSCHMER, AND L. WOESSMANN (2011): “Broadband Infrastructure and Economic Growth,” *The Economic Journal*, 121(552), 505–532.
- DANA, J., AND E. ORLOV (2014): “Internet Penetration and Capacity Utilization in the US Airline Industry,” *American Economic Journal: Microeconomics*, 6(4), 106–137.
- DELLAVIGNA, S., AND E. L. FERRARA (2016): “Economic and social impacts of the media,” in *Handbook of media economics*, ed. by S. Anderson, D. Stromberg, and J. Waldfoel. Elsevier, Amsterdam.
- DRANOVE, D., C. FORMAN, A. GOLDFARB, AND S. GREENSTEIN (2014): “The Trillion Dollar Conundrum: Complementarities and Health Information Technology,” *American Economic Journal: Economic Policy*, 6(4), 239–270.
- ENIKOLOPOV, R., M. PETROVA, AND K. SONIN (2017): “Social media and corruption,” *American Economic Journal: Applied Economics*, forthcoming.
- EUROBAROMETER (2012): “Europeans and their Languages Report,” Special Report 386, European Commission.
- FORMAN, C., A. GOLDFARB, AND S. GREENSTEIN (2012): “The Internet and Local Wages: A Puzzle,” *American Economic Review*, 102(1), 556–575.
- GARCÍA-SÁNCHEZ, A., E. FERNÁNDEZ-RUBIO, AND M. D. COLLADO (2013): “Daily expenses of foreign tourists, length of stay and activities: evidence from Spain,” *Tourism Economics*, 19(3), 613–630.
- GREENSTEIN, S., Y. GU, AND F. ZHU (2016): “Ideological Segregation among Online Collaborators: Evidence from Wikipedians,” Working Paper 22744, National Bureau of Economic Research.
- GREENSTEIN, S., AND F. ZHU (2012): “Is Wikipedia Biased?,” *American Economic Review: Papers and Proceedings*, 102(3), 343–348.

- (2017): “Do Experts or Crowd-based Models Produce More Bias? Evidence from Encyclopedia Britannica and Wikipedia,” *MIS Quarterly*, forthcoming.
- GRINBLATT, M., AND M. KELOHARJU (2001): “How Distance, Language, and Culture Influence Stockholdings and Trades,” *The Journal of Finance*, 56(3), 1053–1073.
- LUCA, M. (2011): “Reviews, Reputation, and Revenue: The Case of Yelp.com,” *manuscript*.
- MYERS, E. W. (1986): “AnO(ND) difference algorithm and its variations,” *Algorithmica*, 1(1-4), 251–266.
- OVERBY, E., AND C. FORMAN (2014): “The Effect of Electronic Commerce on Geographic Purchasing Patterns and Price Dispersion,” *Management Science*, 61(2), 431–453.
- PERI, G. (2005): “Determinants of Knowledge Flows and Their Effect on Innovation,” *The Review of Economics and Statistics*, 87(2), 308–322.
- THOMPSON, N., AND D. HANLEY (2017): “Science Is Shaped by Wikipedia: Evidence from a Randomized Control Trial,” SSRN Scholarly Paper ID 3039505, Social Science Research Network, Rochester, NY.
- XU, S. X., AND X. ZHANG (2013): “Impact of Wikipedia on Market Information Environment: Evidence on Management Disclosure and Investor Reaction,” *MIS Q.*, 37(4), 1043–1068.
- ZHANG, X., AND F. ZHU (2011): “Group Size and Incentives to Contribute: A Natural Experiment at Chinese Wikipedia,” *The American Economic Review*, 101(4), 1601–1615.

Figures

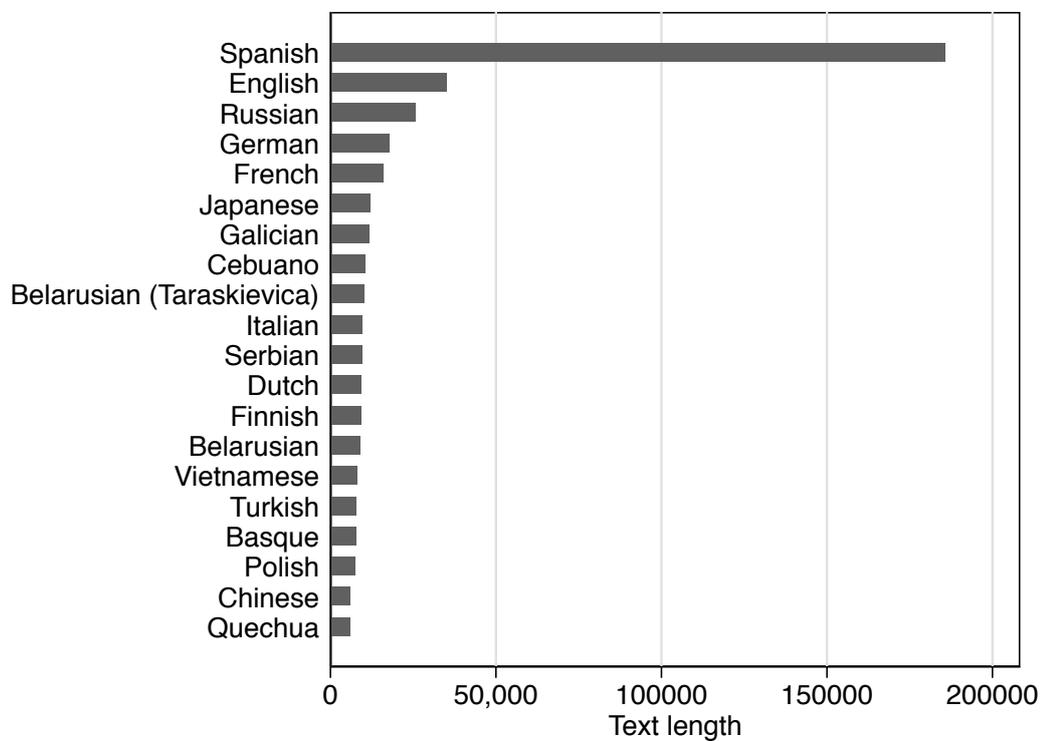


Figure 1: Length of a city page by Wikipedia language edition

Note: The page of the Spanish city exists in 84 Wikipedia language editions. Graph includes 20 languages in which the page is the longest.

Tables

Table 1: Survival over time of text and photos which we added to Wikipedia

	France	Germany	Italy	Total
% text survived: 24h	100.0	94.7	100.0	98.2
% text survived: next month	98.7	90.2	99.9	96.3
% text survived: next year	95.1	86.7	97.5	93.1
% photos survived: 24h	100.0	96.2	100.0	98.8
% photos survived: next month	100.0	92.3	96.4	96.4
% photos survived: next year	100.0	88.5	92.9	94.0
Number of observations	30	30	30	90

Note: Unit of observation is a city page in a given language Wikipedia. Percentage of text survived is calculated as described in section 3. % of text or photos survived is calculated over three time periods: 24 hours, by the beginning of the next calendar month after treatment, by the beginning of the next calendar year after treatment.

Table 2: Ability of covariates to predict treatment status

	Coef.	p-value
Log(Sum of tourists in 2013)	-0.002	0.958
Log(Number of tourists)	-0.012	0.527
Tourist data missing	0.045	0.556
Log(Initial text length)	-0.000	0.994

Note: Dependent variable is the treatment group (an indicator that equals one if a city-language pair is assigned to the treatment group and zero if it is assigned to the control group). Each row presents estimates from a separate regression of the form: $TreatmentGroup_i = Constant + \beta Variable_i + \varepsilon_i$, where $Variable$ is listed in the first column. In rows 1 and 4, a unit of observation is a city-language pair. In rows 2 and 3, a unit of observation is a city-language-month triplet and the sample covers time period until treatment.

Table 3: Dependent variable: Logarithm (number of hotel nights)

	(1)	(2)	(3)
Treatment	0.089** (0.045)	0.002 (0.038)	0.039 (0.045)
Treatment: Small page		0.332*** (0.100)	
Treatment: Large % added			0.196* (0.099)
City-Language FE	Yes	Yes	Yes
Adj. R-squared	0.245	0.248	0.246
Observations	5688	5688	5688

Note: Unit of observation is a month, city, and language (tourist country of origin) triplet. Sample includes tourists from Italy, France, and Germany to the 60 cities in Spain in May–October in 2010–2015. *Treatment* equals 1 for months after treatment for treated city-language pairs, and 0 otherwise. *Small page* equals 1 if the initial page size is below the 25th percentile, and 0 otherwise. *Large % added* equals 1 if text added to the page (as a % of the initial text in the page) is above the 75th percentile, and 0 otherwise. *Controls* include an indicator for period after treatment interacted with language fixed effects, an indicator for period after treatment interacted with city fixed effects, logarithm of number of tourists from Spain interacted with language fixed effects. Standard errors clustered by city-language pair (180 clusters).

Table 4: Robustness

	(1)	(2)	(3)	(4)	(5)	(6)
	Add missing	Add Dutch	All 12 months	12 months, add missing	Add UK	Share of tourists
Treatment	0.091** (0.045)	0.086* (0.047)	0.064 (0.041)	0.078** (0.039)	0.084* (0.043)	0.007* (0.004)
City-Language FE	Yes	Yes	Yes	Yes	Yes	Yes
Log(Tourists from Spain)	No	Yes	Yes	No	No	No
Other controls	Yes	Yes	Yes	Yes	Yes	Yes
Adj. R-squared	0.052	0.212	0.265	0.002	0.104	0.026
Observations	5724	7584	9818	11448	5688	5688

Note: Repeats the regression in column (1) in table 3. In columns 1–5, dependent variable is logarithm of number of hotel nights of tourists from a given country (Germany, France, Italy). Column 1 substitutes missing observations by zeros (only for city-year pairs, when data exists for some month and tourist country of origin). Removes variables of number of tourists from Spain. Column 2 adds observations for tourists from the Netherlands, considers these all as non-treated. Column 3 adds remaining months. Column 4 adds remaining months and substitutes missing observations by zeros (only for city-year pairs, when data exists for some month & tourist country of origin), and removes variables of number of tourists from Spain. In column 5, adds logarithm of the number of tourists from UK interacted with language. In column 6, dependent variable is the number of tourists from country x divided by the number of tourists from country x plus from Spain and UK, and it removes variables of number of tourists from Spain.

Table 5: Wikipedia page views and Google Trends

	Log(Page Views)			Google Trends		
	(1)	(2)	(3)	(4)	(5)	(6)
Treatment	0.116*** (0.030)	0.070** (0.033)	0.069** (0.032)	-0.180 (0.815)	-0.415 (0.862)	-0.317 (0.871)
Treatment: Small page		0.219*** (0.073)			0.892 (1.655)	
Treatment: Large % added			0.183*** (0.069)			0.537 (1.634)
City-Language FE	Yes	Yes	Yes	Yes	Yes	Yes
Adj. R-squared	0.581	0.566	0.582	0.231	0.231	0.231
Observations	12709	12709	12709	12709	12709	12709

Note: In columns 1-3, dependent variable is logarithm of Wikipedia page views. In columns 4-5, dependent variable is Google Trend. Unit of observation is a month, city, and language (country) triplet. Sample includes 3 languages (countries): Italian, French, and German. Sample includes 60 cities in Spain. Time period is 2010–2015 excluding August 2014 (treatment month). *Treatment* equals 1 for months after treatment for treated city-language pairs, and 0 otherwise. *Small page* equals 1 if the initial page size is below the 25th percentile, and 0 otherwise. *Large % added* equals 1 if text added to the page (as a % of the initial text in the page) is above the 75th percentile, and 0 otherwise. *Controls* in all regressions include an indicator for period after treatment interacted with language fixed effects, an indicator for period after treatment interacted with city fixed effects. In columns 1-3, *Controls* include logarithm of page views in Spanish Wikipedia interacted with language fixed effects. In columns 4-6, *Controls* include Google Trends from Spain interacted with language fixed effects. Standard errors clustered by city-language pair (179 clusters).

A Online Appendix: Additional tables and figures

Table A.1: Descriptive statistics of treatment

	mean	sd	p25	p50	p75	count
Length of text added	2047.2	697.2	1671	2082	2377	90
Number of photos added	1.2	1.1	1	1	1	90
% of text added	43.2	37.9	18	29	56	90

Note: Unit of observation is a Wikipedia page in a given language (30 pages in each of the three languages: German, French, Italian).

Table A.2: Wikipedia page length before treatment, by language

	Initial text length			
	p25	p50	p75	count
France	2435	8336	13101	30
Germany	5483	9420	13387	30
Italy	2354	4974	8534	30
Total	2824	8098	11675	90

Note: Unit of observation is a city page in a given language Wikipedia. Sample includes pages in the treatment group.

Table A.3: Robustness: Wikipedia page views and Google Trends

	Page Views		Google Trends	
	(1)	(2)	(3)	(4)
	Add English	Share of views	Add UK	Share of trend
Treatment	0.153*** (0.047)	0.011*** (0.004)	-0.147 (0.829)	0.000 (0.005)
City-Language FE	Yes	Yes	Yes	Yes
Controls: English-UK	Yes	No	Yes	No
Other controls	Yes	Yes	Yes	Yes
Adj. R-squared	0.379	0.101	0.180	0.009
Observations	12709	12575	12709	12709

Note: The table largely repeats regressions in table 5. Dependent variable, in column 1, is logarithm of Wikipedia page views, and in column 2, the number of page views of the article in language x divided by the sum of the number of page views of English, Spanish, and language x . Dependent variable, in column 3, is Google Trend, and in column 4, Google Trend from country x divided by the sum of Google trends from UK, Spain, and country x . Unit of observation is a month, city, and language (country) triplet. Sample includes 3 languages (countries): Italian, French, and German. Sample includes 60 cities in Spain. Time period is 2010–2015 excluding August 2014 (treatment month). *Treatment* equals 1 for months after treatment for treated city-language pairs, and 0 otherwise. *Controls: English-UK* include either logarithm of page views in English Wikipedia (column 1) or Google Trend from UK (column 3), all are interacted with language fixed effects. *Other controls* include an indicator for period after treatment interacted with language fixed effects, an indicator for period after treatment interacted with city fixed effects. Standard errors clustered by city-language pair (179 clusters).

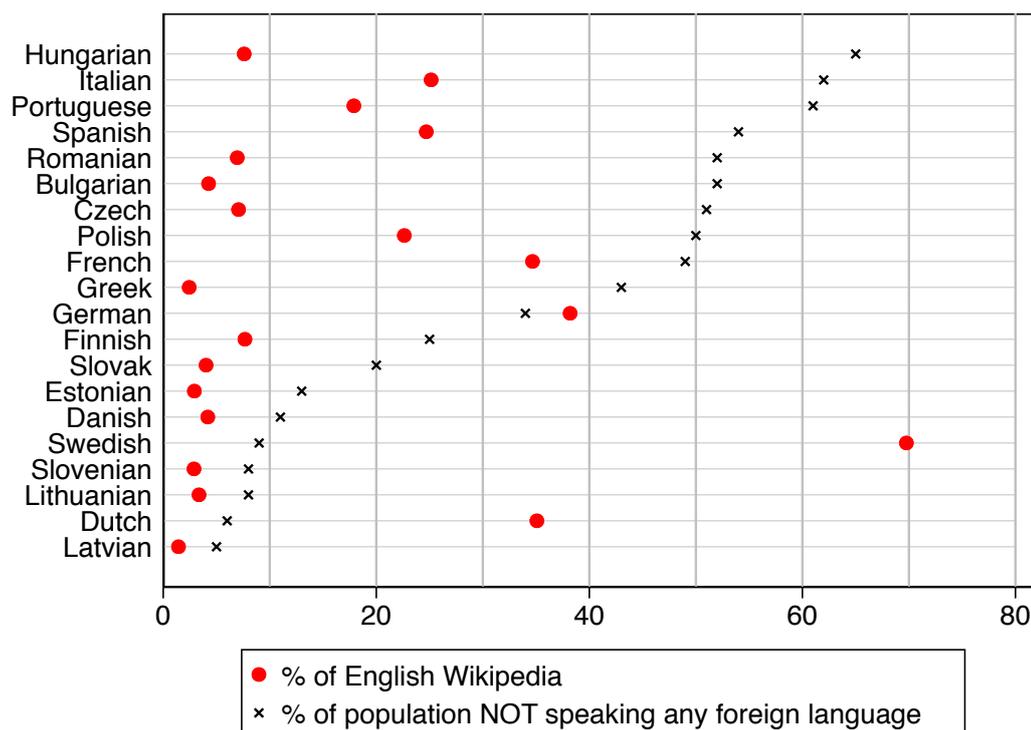


Figure A.1: Size of Wikipedia and percentage of population not speaking any foreign language

Note: The size is measured by the number of articles in the local language Wikipedia as a percentage to the number of articles in English language Wikipedia. Data source for language skills is Eurobarometer (2012).

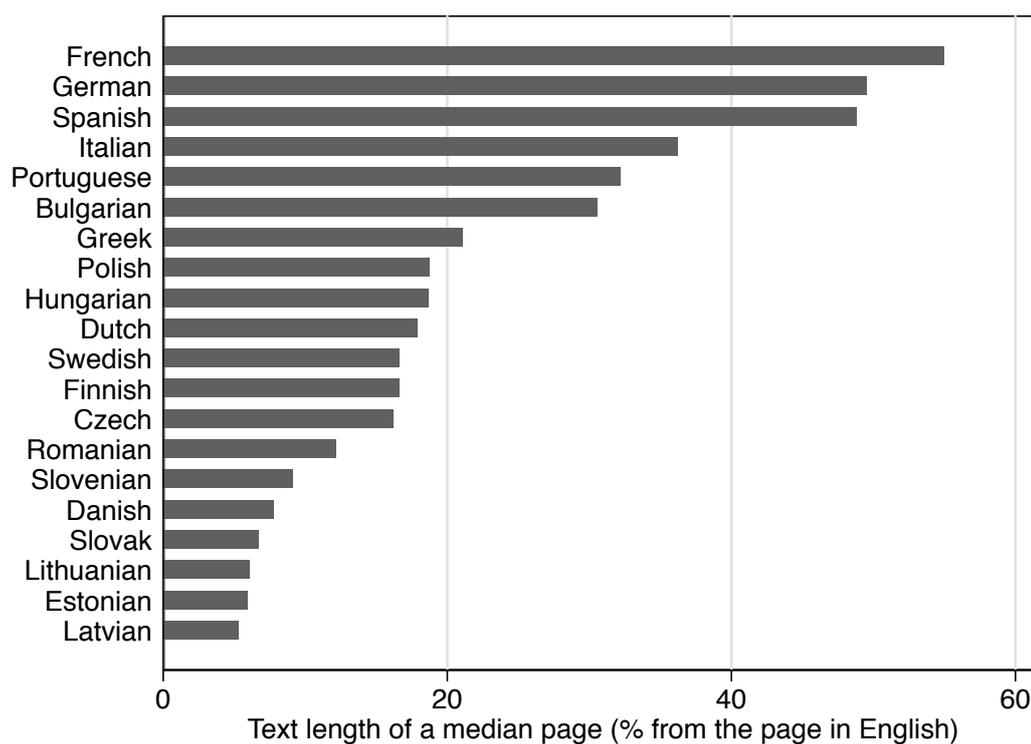


Figure A.2: Median article length by language

Note: The sample includes pages in the list of 1000 vital articles chosen by Wikipedia community. For each page, the relative text length is calculated as the percentage of of the length of text in the local language Wikipedia compared to that of the English language Wikipedia edition. The graph presents the median of the relative text lengths by language.

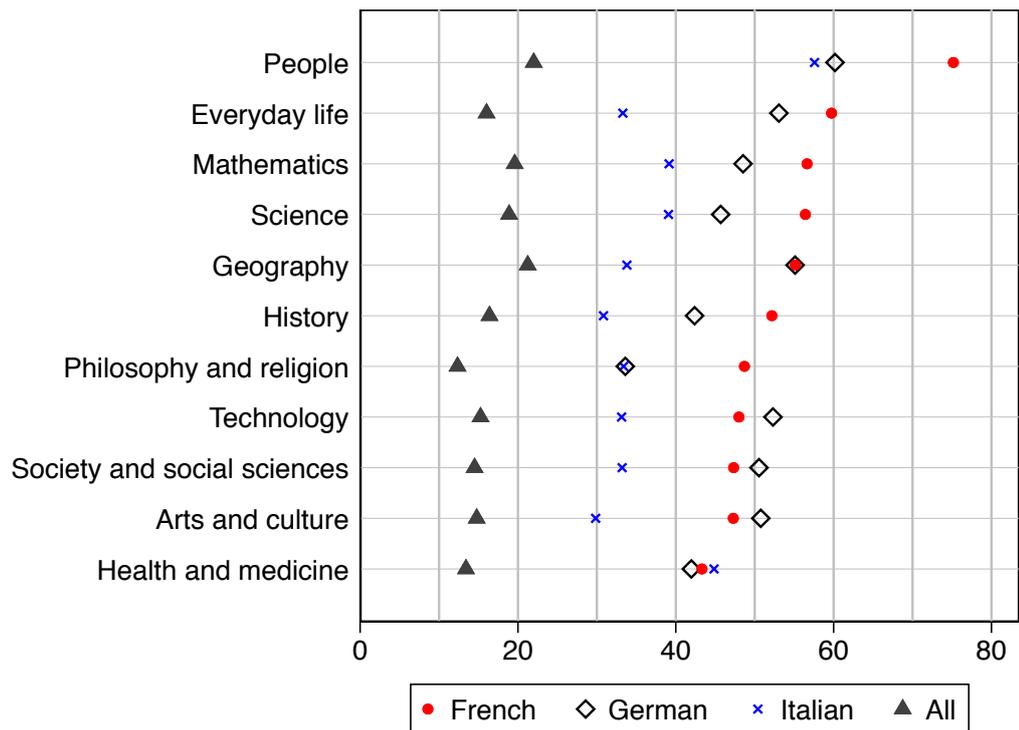


Figure A.3: Median article length by topic

Note: The sample includes pages in the list of 1000 vital articles chosen by Wikipedia community. For each page, the relative text length is calculated as the percentage of of the length of text in the local language Wikipedia compared to that of the English language Wikipedia edition. The graph presents the median of the relative text lengths by article category. For each category, it presents the overall median and median by language (French, German, Italian).

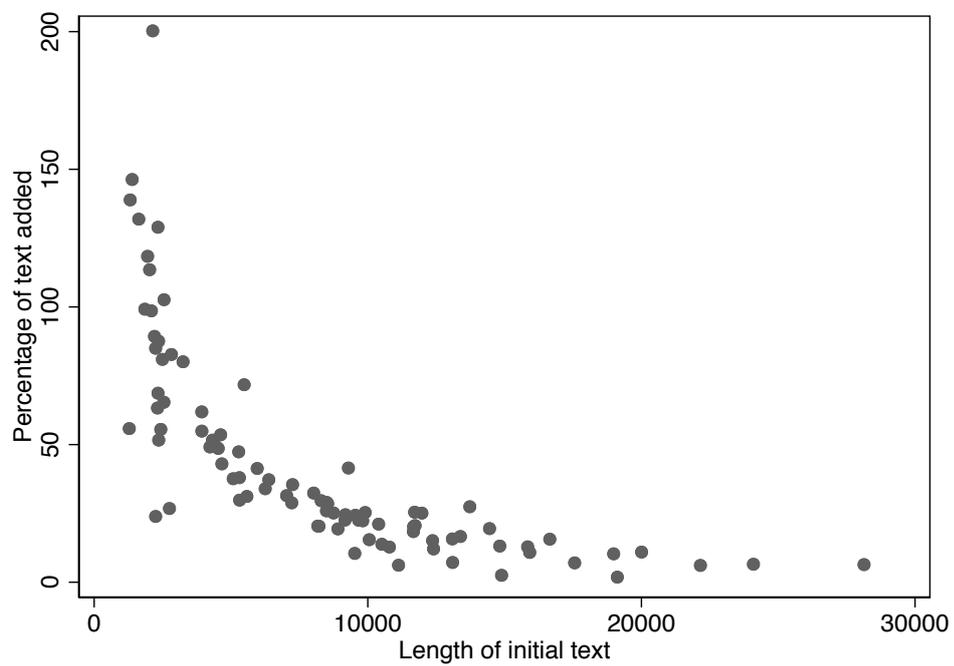


Figure A.4: Length of text added (as % of initial text) vs length of initial text

Note: Unit of observation is a Wikipedia page in a given language (30 pages in each of the three languages: German, French, Italian). Sample includes treated pages.

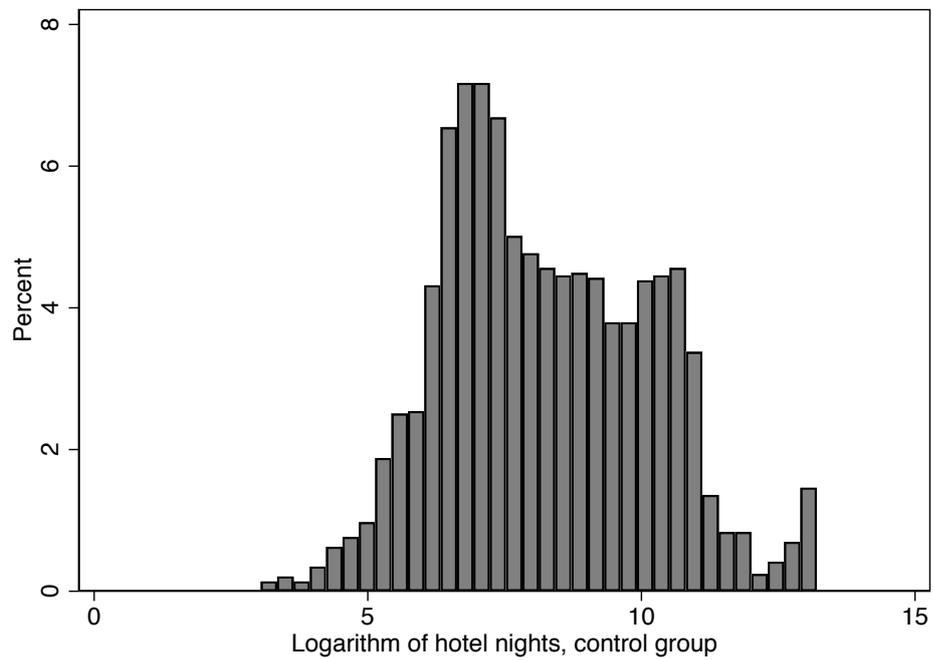


Figure A.5: Logarithm of number of hotel nights in the control group

Note: Unit of observation is a month, city, and tourist country of origin triplet. Sample includes tourists from Italy, France, Germany to the 60 cities in Spain, but only the city-country of origin pairs, which were assigned to the control group. The time period of the sample is May–October in 2010 - 2015.

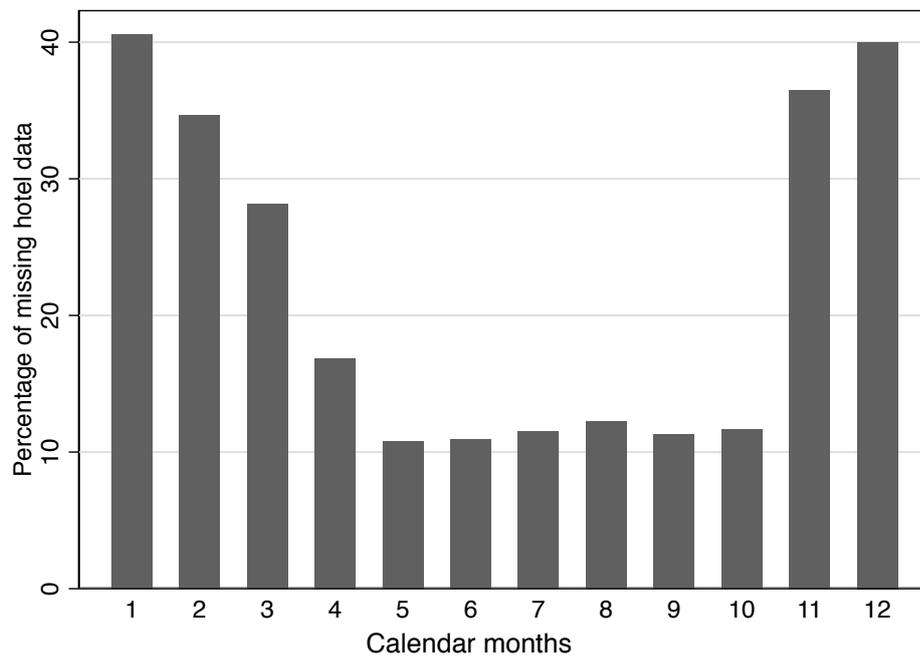


Figure A.6: Percentage of missing hotel data, over 12 calendar months (January–December)

Note: Unit of observation is a month, city, and tourist country of origin triplet. Sample includes tourists from Italy, France, Germany to the 60 cities in Spain, but only city-country of origin pairs, which were assigned to the control group. The time period of the sample is 2010 - 2015.

B Online Appendix: Measuring our treatment and its survival

We applied diff algorithm twice to quantify how much we added by our treatment and how much of it was preserved a few months later. In particular, for each page we compared three revisions that we took from the Wikipedia revision history: the last revision prior to our changes (which we call *pre-treatment* revision), the last revision created by our treatment (*post-treatment*), and version a few months later (*survived*). In the revision history, the text is always in the Wikitext format, which means that some of it is not visible for the viewer. We normalized all the three revisions as follows. We used Wikipedia’s built-in parser to get the html-version of the content, which we then converted to plain text by removing the html commands, i.e. removed all pictures, links, etc. This gave us three texts.

The length of pre-treatment is our page length measure. To quantify the content added by our treatment, we used a diff algorithm. It computes the smallest number of character additions and deletions from pre-treatment to post-treatment. The algorithm outputs which characters stayed the same, which ones were deleted, and which ones added. The total length of the added text is our measure of treatment length. Finally, to compute how much of the text survived after the editing process a few months later we computed diff from the added text to the survived text.¹⁸ See figure B1 for illustration.

Revision	Text	Difference	Length
Pre-treatment	abc		3
Post-treatment	adce	diff(abc,adce)=a d <u>ce</u>	Added 2 (<u>de</u>)
Survived	acef	diff(de,acef)= <u>a</u> c <u>ef</u>	Survived 1

Figure B1: Illustration how we used diff algorithm to quantify the additions by treatment and the survival of the additions.

¹⁸It is slightly imperfect measure, as there could be some text that was deleted, but the algorithm is unable to differentiate it from the other parts of the page (that were unrelated to our treatment), but in examples we checked by hand the results were accurate within a reasonable margin.