

Making the ‘Next Billion’ Demand Access

The Effect of Local Content: Google.co.za in Setswana

Bastiaan Quast

Abstract

This paper shows that an exogenous increase in available local content leads to a large increase in demand for internet connectivity among native speakers, even as demand as a whole is falling. Internet connectivity provides enormous improvements in quality of life as well as opportunities for the newly connected. Attempts to connect the ‘next billion’ in Africa have not met expectations, even in places where infrastructure has come online and prices have gone down. The introduction of the Setswana (Tswana) language in the South-African Google Search website was a spillover effect of this translation work being done for the Botswanan Google Search website. This exogenous event created a large increase in the number of internet-connected native speakers, as well as usage of the Setswana language online.

1 Introduction

- Setswana is also an official language of South Africa, but only a relatively small percentage of people speak it, there are also Setswana speakers in Zimbabwe and Namibia.
- Data from South Africa on 2008, 2010-2011, and 2012 (Southern Africa Labour and Development Research Unit 2008, 2012, 2013).
- Introduction in Botswana in late 2010, presumably some lag of information on non-internet users.
- Also shows increased computer ownership.

Local content is a vital means to connect new internet users. Since the term ‘Connecting the Next Billion’ was introduced in The Economist’s 2006 ‘End of Year Report’ (Standage 2006), XXX many people have been estimated to have connected to the internet. Yet despite increased range and improved affordability, many key growth markets such as sub-Saharan Africa are showing stagnation in internet connections. This paper shows that exogenous increase in accessibility of local content gave rise to a vast increase in the number of internet users among native speakers. In 2010 Google collaborated with a Botswanan team of linguists (Otlogetswe 2010) to make its Botswanan website (google.co.bw) available in the local language: ‘Setswana’. In addition to being spoken in Botswana, there is also a sizable population of Setswana directly across the border in South Africa, where it is also one of the official state languages. This led to the introduction of the Setswana language on the South African Google website (google.co.za) as spillover of the translation work for Google’s Botswanan website. This exogenous led to a vast increase in the number of native Setswana speakers reporting to have spent some amount of money in the past 30 days on internet access.

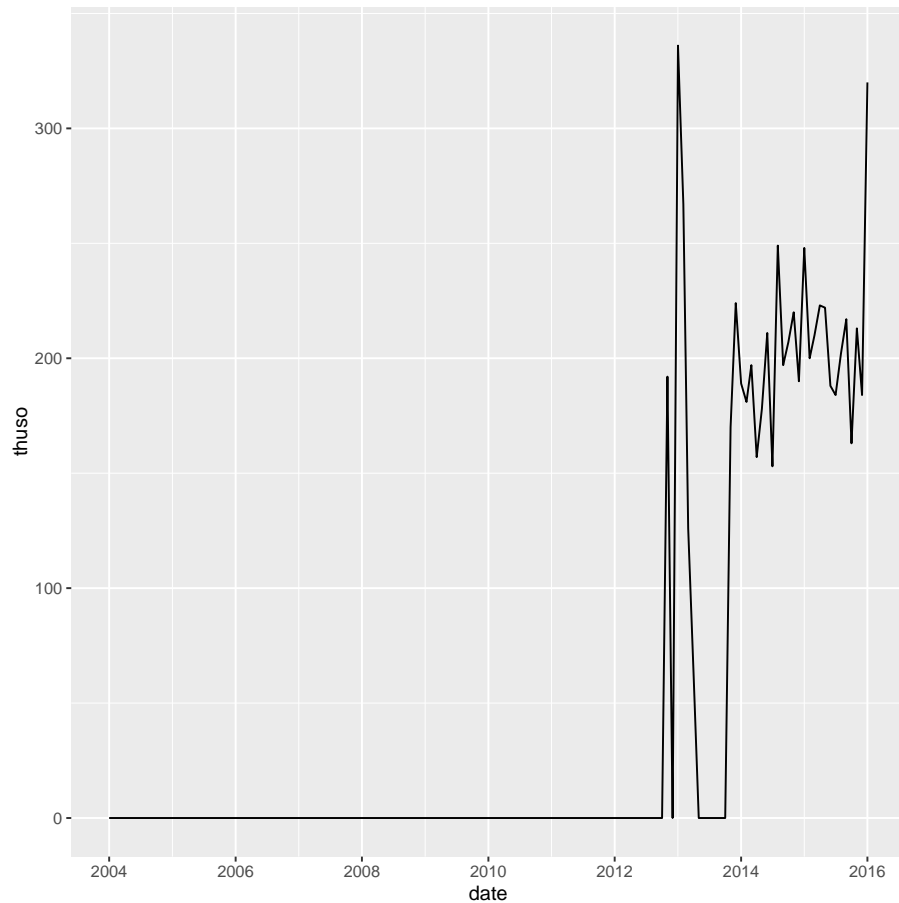
It is not required to use a certain interface language in order to search in this language. However, the search page is in many instances the first website viewed by users and the asides from being able to understand the interface, having the interface be in a certain language also encourages usage of this language, which in turn reveals more content in this language.

In short we can identify several major channels which promote further engagement.

1. Being able to read and understand the words.
2. Encouragement from the familiarity with the language on what is often the first website visited.
3. Increased likelihood to using local language (see figure 1) and thus finding more content in the native language.

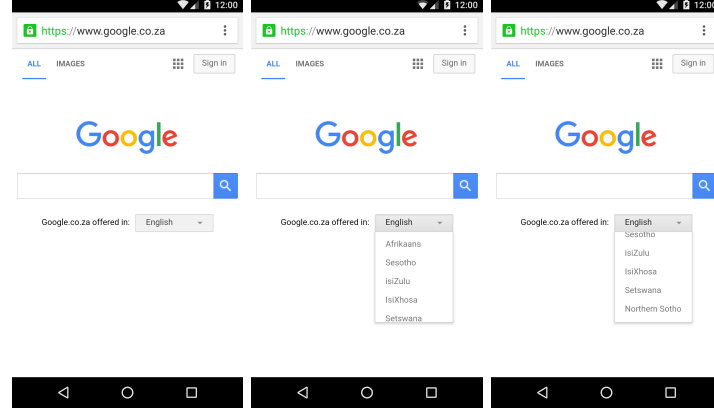
Figure 1: Usage of Setswana Words on Google.co.za

```
ggplot(thuso) + geom_line(aes(x = date, y = thuso))
```



The vast majority of internet access in developing countries is through hand-held devices such as smartphones. However, due to the limited ‘real estate’ on a mobile website, the link to changing the interface language is replaced with a dropdown menu that reveals the additional language options. Generally, the website will default to the operating system (Android / iOS) language, however, since many local African languages are not available as a system language, this is not possible there. The fact that the introduction seems to benefit desktop usage but not mobile usage is further substantiated by our results that isolate an increase in computer ownership, but no increase in cell phone ownership.

Figure 2: Changing Interface Language on Mobile



The data used for this study comes from the South African National Income Dynamics Survey, provided by Southern Africa Labour and Development Research Unit 2008, 2012, 2013, the data is further discussed in 2. Section 3 discusses the methods employed in this study, specifically, the discussion of the identification strategy can be found in 3.1 and the use of the Difference-in-Difference estimator in 3.2.

2 Data

- South African National Income Dynamics Survey (Southern Africa Labour and Development Research Unit 2008, 2012, 2013)
- descriptive stats:
 - number of adults
 - Setswana speakers
 - people using internet / cell phone etc.
 - male / female
 - income distribution

South Africa's National Income Dynamics Survey collect data on a representative set of around 10,000 households over time. The first survey took place in 2008, the second one in late 2010 and early 2011, and third one took place in 2012.

3 Methods

3.1 Identification Strategy

- Spillover from development in Botswana

This paper exploits the introduction of the Setswana interface language to Google Search in South Africa as a spillover of the development of that interface for the Botswanan Google Search website. By comparing the number of native Setswana speakers in South Africa being internet users, with the number of South Africans with a different native language around the same time, we isolate the effect of this introduction.

The Setswana language was first developed for the Botswanan Google Search website (google.co.bw). As such, the introduction of Setswana to the South African Google Search (google.co.za) was a spillover effect of that development. This allows us to rule out any possible endogeneity issues that might otherwise arise in context such as these. For instance, the Afrikaans language is almost solely spoken in South Africa. When we observe that the introduction of the Afrikaans Google Search interface occurs around the same time as a growth in the number of native Afrikaans internet users, it will be hard to isolate the effect from the introduction from its cause (since an increase in native Afrikaans internet users would be a good reason to introduce it as an interface language).

Substantial numbers of Setswana speakers exist in Botswana, South Africa, Zimbabwe, and to some extent Namibia. However the language is most important in Botswana, where it is spoken by approximately 80% of all people, and where it is the only official language other than English. As such, it is also the place where most linguistic work on the Setswana language takes place. The Setswana Google Search interface was also developed at the university of XXX by prof. OtseXXX.

It is worth noting that it is very common not to personally own a computer and ‘paying for internet access’ therefore also includes a lot of people who use the internet in other locations such as internet cafe’s.

In addition to using the propensity to spend on internet (in the last thirty days), we also use the propensity to own a computer as a dependent variable.

3.2 Estimation

As mentioned in the above section, we compare the change in the level of internet users among native Setswana speakers in South Africa, with that of native speakers of other language in South Africa around the introduction of the Setswana interface to the South-African Google Search. For this we use a Difference in Difference estimator using a native-Setswana speaker dummy variable, interacted with a event dummy variable on the introduction of the Setswana interface on google.co.za.

In addition to this estimation, use an alternative specification whereby a factor variable of native language is interacted with the event dummy variable. In a linear regression context, factor variables are estimated as dummy variables for all levels (here: all languages) except for once ‘base’ level, which is where all language dummies are false (0) and the level (native language) thus has to be the nth one (here English).

The dependent variable here is also a dummy variable, which would normally allow for the usage of an estimator such as logit. However, since we are employing

the Difference in Difference methodology

4 Results

- Base model's variable of interest (interaction of event dummy and Setswana dummy) finds strong significant result of interaction effect.
- Alternative formulation's variable of interest (interaction of event dummy and factor of categorical language variable) only significant growth only for Setswana and Venda.
- Propensity to own or spend on cellphone are not affected, presumably less obvious language change.

Table 1: Base model

```
lm(h_nfnet ~ post_event*factor(a_lng))
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.0153374	0.0069104	2.2194845	0.0264583
post_eventTRUE	-0.0153374	0.0116070	-1.3213959	0.1863755
factor(a_lng)2	-0.0119771	0.0071295	-1.6799254	0.0929782
factor(a_lng)3	-0.0101387	0.0070313	-1.4419371	0.1493265
factor(a_lng)4	-0.0118839	0.0073298	-1.6212958	0.1049606
factor(a_lng)5	-0.0017102	0.0073480	-0.2327407	0.8159637
factor(a_lng)6	-0.0100812	0.0072710	-1.3864957	0.1656019
factor(a_lng)7	-0.0106935	0.0084765	-1.2615403	0.2071202
factor(a_lng)8	0.1182366	0.0101957	11.5966572	0.0000000
factor(a_lng)9	0.0258487	0.0085674	3.0171200	0.0025532
factor(a_lng)10	0.0293054	0.0071303	4.1099620	0.0000396
factor(a_lng)11	0.0884438	0.0077638	11.3917636	0.0000000
factor(a_lng)12	0.0927707	0.0216448	4.2860579	0.0000182
post_eventTRUE:factor(a_lng)2	0.0139527	0.0119551	1.1670893	0.2431800
post_eventTRUE:factor(a_lng)3	0.0146315	0.0117940	1.2405899	0.2147632
post_eventTRUE:factor(a_lng)4	0.0169315	0.0122240	1.3851020	0.1660275
post_eventTRUE:factor(a_lng)5	0.0130134	0.0123009	1.0579188	0.2900977
post_eventTRUE:factor(a_lng)6	0.0207076	0.0121878	1.6990391	0.0893181
post_eventTRUE:factor(a_lng)7	0.0241705	0.0141697	1.7057877	0.0880539
post_eventTRUE:factor(a_lng)8	-0.1182366	0.0155235	-7.6166264	0.0000000
post_eventTRUE:factor(a_lng)9	-0.0234958	0.0140356	-1.6740145	0.0941341
post_eventTRUE:factor(a_lng)10	0.0106404	0.0119616	0.8895436	0.3737153
post_eventTRUE:factor(a_lng)11	0.0139328	0.0132554	1.0511063	0.2932149

4.1 LM4_1

```
lm(h_nfnet ~ post_event*factor(a_lng) +
      a_edlitrden +
      a_edlitwrten +
      a_edlitrldhm +
      a_edlitwrthm +
      a_woman)
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.0318595	0.0071928	4.4293548	0.0000095
post_eventTRUE	-0.0160615	0.0115788	-1.3871439	0.1654045
factor(a_lng)2	-0.0099108	0.0073031	-1.3570604	0.1747686
factor(a_lng)3	-0.0074446	0.0072117	-1.0322986	0.3019376
factor(a_lng)4	-0.0117912	0.0075144	-1.5691408	0.1166220
factor(a_lng)5	-0.0024480	0.0075095	-0.3259828	0.7444388
factor(a_lng)6	-0.0090931	0.0074556	-1.2196300	0.2226114
factor(a_lng)7	-0.0109606	0.0086869	-1.2617360	0.2070501
factor(a_lng)8	0.1112536	0.0102699	10.8329397	0.0000000
factor(a_lng)9	0.0241967	0.0088206	2.7432030	0.0060866
factor(a_lng)10	0.0327904	0.0073039	4.4894292	0.0000072
factor(a_lng)11	0.0832838	0.0079063	10.5339061	0.0000000
factor(a_lng)12	0.0929623	0.0216230	4.2992237	0.0000172
a_edlitrden	-0.0033726	0.0022402	-1.5054902	0.1322049
a_edlitwrten	-0.0068192	0.0021992	-3.1007081	0.0019317
a_edlitrldhm	0.0005543	0.0021514	0.2576646	0.7966669
a_edlitwrthm	0.0021000	0.0021548	0.9745694	0.3297790
a_womanTRUE	-0.0012292	0.0011597	-1.0599272	0.2891832
post_eventTRUE:factor(a_lng)2	0.0132434	0.0119238	1.1106677	0.2667171
post_eventTRUE:factor(a_lng)3	0.0145833	0.0117679	1.2392479	0.2152600
post_eventTRUE:factor(a_lng)4	0.0177728	0.0121993	1.4568644	0.1451606
post_eventTRUE:factor(a_lng)5	0.0129743	0.0122610	1.0581779	0.2899798
post_eventTRUE:factor(a_lng)6	0.0200726	0.0121631	1.6502852	0.0988914
post_eventTRUE:factor(a_lng)7	0.0253515	0.0141395	1.7929538	0.0729868
post_eventTRUE:factor(a_lng)8	-0.1117654	0.0154157	-7.2500872	0.0000000
post_eventTRUE:factor(a_lng)9	-0.0181736	0.0140391	-1.2944988	0.1954996
post_eventTRUE:factor(a_lng)10	0.0095912	0.0119282	0.8040750	0.4213578

4.2 LM4_5

```
lm(h_nfnet ~ post_event*setswana +
      factor(a_edlitrden) +
      factor(a_edlitwrten) +
```

```

factor(a_edlitrthm) +
factor(a_edlitwrthm) +
a_woman +
hhincome +
best_edu)

```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-0.0007802	0.0018382	-0.4244139	0.6712660
post_eventTRUE	-0.0118598	0.0012192	-9.7275870	0.0000000
setswanaTRUE	-0.0136327	0.0024304	-5.6091357	0.0000000
factor(a_edlitrden)2	0.0015400	0.0040296	0.3821762	0.7023324
factor(a_edlitrden)3	0.0002169	0.0052826	0.0410653	0.9672440
factor(a_edlitrden)4	-0.0037000	0.0068226	-0.5423219	0.5875994
factor(a_edlitwrten)2	-0.0102695	0.0040201	-2.5545159	0.0106367
factor(a_edlitwrten)3	-0.0108074	0.0052286	-2.0669776	0.0387418
factor(a_edlitwrten)4	-0.0071782	0.0067009	-1.0712332	0.2840702
factor(a_edlitrthm)2	-0.0026154	0.0036909	-0.7086027	0.4785746
factor(a_edlitrthm)3	-0.0029346	0.0051387	-0.5710741	0.5679522
factor(a_edlitrthm)4	-0.0080255	0.0069159	-1.1604527	0.2458705
factor(a_edlitwrthm)2	0.0008491	0.0037294	0.2276776	0.8198979
factor(a_edlitwrthm)3	0.0013947	0.0051466	0.2709891	0.7864006
factor(a_edlitwrthm)4	-0.0069746	0.0069204	-1.0078396	0.3135367
a_womanTRUE	-0.0014132	0.0011472	-1.2318972	0.2179937
hhincome	0.0000028	0.0000001	46.4115514	0.0000000
best_edu	0.0013633	0.0001165	11.6980667	0.0000000
post_eventTRUE:setswanaTRUE	0.0120675	0.0038748	3.1143779	0.0018445

4.3 LM2_5

```

lm(a_owncom ~ post_event*setswana +
  factor(a_edlitrden) +
  factor(a_edlitwrten) +
  factor(a_edlitrthm) +
  factor(a_edlitwrthm) +
  a_woman +
  hhincome +
  best_edu)

```


	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.0076236	0.0031357	2.4312519	0.0150504
post_eventTRUE	-0.0054318	0.0020923	-2.5960394	0.0094334
setswanaTRUE	-0.0147962	0.0041674	-3.5504730	0.0003849
factor(a_edlitrden)2	-0.0306508	0.0069013	-4.4412904	0.0000090
factor(a_edlitrden)3	-0.0309578	0.0090312	-3.4278696	0.0006089
factor(a_edlitrden)4	-0.0389988	0.0116608	-3.3444228	0.0008252
factor(a_edlitwrten)2	-0.0174611	0.0068860	-2.5357494	0.0112239
factor(a_edlitwrten)3	-0.0210480	0.0089368	-2.3552002	0.0185168
factor(a_edlitwrten)4	-0.0187070	0.0114489	-1.6339589	0.1022741
factor(a_edlitrdhm)2	-0.0018694	0.0063225	-0.2956791	0.7674765
factor(a_edlitrdhm)3	-0.0042389	0.0088029	-0.4815352	0.6301384
factor(a_edlitrdhm)4	-0.0267295	0.0118766	-2.2506064	0.0244150
factor(a_edlitwrthm)2	0.0012267	0.0063829	0.1921868	0.8475967
factor(a_edlitwrthm)3	-0.0019980	0.0088200	-0.2265280	0.8207918
factor(a_edlitwrthm)4	-0.0357502	0.0118852	-3.0079574	0.0026315
a_womanTRUE	-0.0229898	0.0019601	-11.7288993	0.0000000
hhincome	0.0000058	0.0000001	56.9305683	0.0000000
best_edu	0.0058348	0.0002002	29.1431892	0.0000000
post_eventTRUE:setswanaTRUE	0.0238541	0.0066835	3.5690970	0.0003586

5 Conclusions and Limitations

- need more local content
- need more research

References

- Otlogetswe, Thapelo J.
2010 “Setswana Google is here!”, *T.J. Otlogetswe Blog*, <http://otlogetswe.com/2010/08/13/setswana-google-here/>.
- Southern Africa Labour and Development Research Unit
2008 *National Income Dynamics Study, Wave 1*, version 5.3, <http://www.nids.uct.ac.za/home/>.
2012 *National Income Dynamics Study, Wave 2*, version 2.3, <http://www.nids.uct.ac.za/home/>.
2013 *National Income Dynamics Study, Wave 3*, version 1.3, <http://www.nids.uct.ac.za/home/>.
- Standage, T
2006 “Connecting the next billion”, *The Economist-The World in 2006*, p. 117, <http://www.economist.com/node/5134746>.