

Contribution

Literature Review

William Marsey
Imperial College London

June 2014

1 Introduction

1.1 The Research Question

Given a collaboratively edited document, and information about the collection of edits that constitute that document, may we measure the quality of each contribution? And may we use that to give all the collaborators a algorithmically-defined 'stake' in that final document?

Collaborative work is becoming a big deal. It is both interesting and an important trend in modern computer use. And the data is abundant.

Amongst many other things, this topic is a playground for sociology, machine learning, network studies, as well as more general studies of conflict, and personality. My work intends to focus on the algorithmic side of things - approximate string matching in particular. I look at how we may use Levenshtein distance, and the various favours, varieties and optimizations thereof, to measure contribution to a collaborative text, and how we may implement a version of this algorithm specifically tailored to our needs.

The main questions we ask are:

- What does Levenshtein distance define of a contribution in the context of massive online collaboration?
- What are the limitations and implications of defining contribution in this way?
- What else may we learn from analysing contribution?

We base our studies around data from Wikipedia. This study is defined by – and in some ways determined by – the specific context of Wikipedia, but, as we will see, is ultimately enriched by it. Due to its open-source nature, and its size, studies that touch upon Wikipedia cover a very broad range of topics. Many of them are directly related to the topic we concern ourselves with here, and many more may enrich our study tangentially.

2 Previous work

There are three sections here for the three different topics that come to bear on this subject:

- Wikipedia, studies of wikipedia, and the nature of Wikipedia
- The 'edit distance problem', Levenshtein distance
- The various pre-existing studies that apply the latter to the former

2.1 On Wikipedia

Wikipedia's pre-eminence as an online resource is self-evident to anyone who has searched the internet for a generic topic. The website is ranked 6th globally in terms of website traffic,¹ and is the highest-ranked reference website by far - most of the sites it shares the top spots with are portals, search engines, shopping mega-sites, and social media websites.[3] Despite some skepticism (particularly concern over the inherent chaos of the system: "...edits, contributed in a predominantly undirected and haphazard fashion by ... unvetted volunteers." [25]), it is widely claimed to be a success, 'the best-developed attempt thus far of the enduring quest to gather all human knowledge in one place' [19].

That Wikipedia has become a hub of research in many fields is also self-evident to anyone who has searched for articles on the subject. Mesgari et al, just quoted, has prepared a very recent 'systematic review of scholarly research on the content of Wikipedia', which gives an overview of 110 articles on the subject — an attestation to the observation that Wikipedia has been 'irresistable point of enquiry for researchers from various fields of knowledge', and a useful touching stone for this study. Mesgari et al's review finds 82 out of the 100 to concern quality in Wikipedia articles, some of these are also referenced here, and many of the others will come to bear on the study as it progresses.

Other important general sources will be WikiLit,[5] AcaWiki[1] and WikiPapers[8], all of which are online repositories of academic research into Wikipedia and other Wikis (as well as being Wikis themselves...).

2.1.1 The particularities of Wikipedia, and implications

The six 'risks' one takes when referencing Wikipedia, as defined in an early article on the subject,[12] is a good starting block for identifying the ways in which to regard the 'quality' of content in Wikipedia. We list them here, describing the implications for our work with each. Some are particular to Wikipedia, some are inherent to all Wikis.

- **Accuracy.** It is important to remember that, without severely increasing the complexity of our algorithm, we may not verify the accuracy of information. And, if accuracy of information is proportionate to value (surely it must be in a reference text), then our algorithm may misplace its reward. We may most usefully look at the problem as follows. The texts that are edited most often are those that are visited most often. The previously cited studies of Lih and Mestyan et al attest to this - they both study the peaks in activity in articles that are brought to attention in some way. We find in the work of Bongwon Suh et al that the growth of wikipedia is inverse-exponential, as the overheads of coordination and beaucroscopy temper content creation.[23][14] Content is more likely to be refined and corrected as an article grows old.[25] We can assume, then, that all articles tend towards accuracy (we may find this bore out in Giles's 2005 semi-formal comparison of Encyclopedia Britannica articles to Wikipedia articles - finding an average three mistakes in the former and 4 mistakes in the latter)[13]. We may possibly extend this to say that all edits improve a text. This is complicated by malicious, misinformed or malformed edits, but we will discuss dealing with these later.

Response:

- **Motives.** It has been found that different contributors may edit Wikipedia in various different ways, according to their proclivities.

Response:

¹According to Alexa, an Amazon-owned company. The statistics are wide-ranging based on a combined measure of Unique Visitors and Pageviews, and the data mined from around 25,000 different browser extensions, as well as sites that have installed Alexa's scripts.[2] Alexa may well be biased towards English speakers and Internet Explorer users, but this may underestimate Wikipedia.org's popularity, since 'two thirds of all Wikipedia articles are in languages other than English'[7]

- **Uncertain Expertise.** We may take this to mean malformed and misinformed editing, but we may also take it to mean malicious editing. As for malicious edits - we find a lot of useful information in Potthast et al's work on automatic detection of vandals,[21] as well as the discussions around Wikipedia's own Counter-Vandalism Unit ('CVU').[9] including
- **Volatility.**
- **Coverage.** Cite that structure is a problem.
Response: We may want to reward extra for restructuring.
- **Sources.** There has been some work explicitly taking external links to be relative to quality,[?] and seems to have been borne out by a further study which found this to be a heuristic used by actual readers.[?].
Response: We may give extra weight to the value of (working) hyperlinks, and fixing hyperlinks.

2.2 Wikipedia Self-Valuation

[Wikipedia] cannot attain the status of a true encyclopedia without more formal content-inclusion and expert review procedures.[12]

Most visited articles in an hour - correlates with (american-centric) events [10]

Denning says it cannot attain the status of a true encyclopedia without more formal content-inclusion and expert review procedures[12] this corroborates by findings in [13]?

2.2.1 On Wikipedia

'robust and remarkable growth' [14][24]

Wikipedia, at the last dump, consisted 800G of compressed data [6]

2.2.2 Evaluating Wikipedia articles

identify, analyse

after article mentioned in press [17]

compared by 'experts' to 'equivalent' Encyclopedia Britannica articles [13]

found metrics of article quality through factor analysis [22]

Analysis by conflict - revisions?[14]

WikiTrust. The most 'complete' of the many of the. Exists as firefox plugin (though it doesn't work any more) Culmination of various studies that try to QUOTE [11] and QUOTE CITE. It was assessed as recently as 2011 [18]

2.3 Edit difference

Standard: Levenshtein distance [15]

2.4 My work in context of these sources

different views emerging topics gaps and inconsistencies

Goals.

Methods. Python. Levenshtein. Optimisations of.

3 Conclusions

3.1 Intentions / Progress

The Code in appendix A

Possible extensions. Perhaps we can examine more about what we may find out about the articles. Cite other studies Lih's 2004 study of articles immediately after they have been cited in the press[17], and Metyan's 2012 use of the site to predict box-office success [20]. Lieberman worked out their locale.[16]:

1. peaks in activity (rate of levenshtein distance... may have to be in a log graph...)
- 2.

PREDIFINED / NOT-PREDEFINED ideas of quality. look for when the the article levels off? And do this by DATE rather than REVISION. We may assume that pageviews are more well-distributed than revisions

summarize major contributions (which do we care about?) evaluate your current position point out any flaw in methodology/research/contradictions are there any gaps in the area which you will cover in your research? How will you integrate sources you have mentioned into your dissertation? Point out any areas for further study

A Appendix A: Code progress

A.1 Python class for scraping a Wikipedia article's revision history

The following code is a first draft of a class which incrementally traces, parses, and stores the revision history of select articles. It chooses random articles up to a maximum amount of articles unless an article (or articles) are specified. It traces the entire discoverable² history of the, unless a specific depth is specified by the user.

The class already yields workable data, but here is some immediate further work for this code:

- Allow the user to specify timeframe
- Allow for integration with a postgres database (at the moment the code saves the data in CSV format).

This code is an extension to an existing wikipedia API class for python (which did not provide the features we need here).[4]

```
import requests
import time
import json
import csv
import wikipedia
from bs4 import BeautifulSoup
from datetime import datetime, timedelta
from decimal import Decimal

WIKI_API_URL = 'http://en.wikipedia.org/w/api.php'
WIKI_USER_AGENT = 'wikipedia (https://github.com/goldsmith/Wikipedia/)'

class WikiRevisionScrape:
    par = {
        'format': 'json',
        'action': 'query',
        'prop': 'revisions'
    }
    head = {
        'User-Agent': WIKI_USER_AGENT
    }
    rand = True
    pagelimit = 1
    historylimit = -1
    rl = False
    rl_minwait = None
    rl_lastcall = None
    pageid = 0
    parentid = 0
    childid = 0

    #atm naively assuming headers, params, titles to be in correct format
    def __init__(self, pagelimit=1, historylimit=-1, _headers=None, _params=None,
        _titles=None):
        if(_params):
            params = _params

        if(_headers):
            self.head = _headers
```

²Using the Wikipedia API, articles can either be traced back to their origin (revision parent ID = 0), or to the point at which a loop is found in the revision history - this usually happens with older articles.

```

if(_titles):
    self.params['titles'] = _titles
    self.rand = False

self.pagelimit = pagelimit
self.historylimit = historylimit

def scrape(self, indexfilename, contentsfilename):
    self._pace()
    index_f = open(indexfilename + ".csv", "ab") #HACK = needs to migrate to postgres
    contents_f = open(contentsfilename + ".csv", "ab") #HACK = needs to migrate to
        postgres
    index = csv.writer(index_f)
    contents = csv.writer(contents_f)
    index.writerow(["PAGEID", "REVISION", "USER", "USERID", "TIMESTAMP", "SIZE", "COMMENT"])
    contents.writerow(["PAGEID", "REVISION", "CONTENT"])

    for i in range(self.pagelimit):
        if 'rvprop' in self.par:
            del self.par['rvprop']
        if 'revids' in self.par:
            del self.par['revids']
        print "fetching page"
        if(self.rand):
            self.par['titles'] = wikipedia.random() #get random title
        self.childid = self._getlatest()
        r = requests.get(WIKI_API_URL, params=self.par, headers=self.head)
        self._rate()
        del self.par['titles']
        self._tracehist(index, contents)

def _getlatest(self):
    r = requests.get(WIKI_API_URL, params=self.par, headers=self.head)
    r = r.json()

    #HACK = should grab multiple pages
    for key, value in r['query']['pages'].iteritems():
        self.pageid = key
    #HACK = should grab multiple revisions (for each pageid)
    self.parentid = self.childid =
        r['query']['pages'][self.pageid]['revisions'][0]['revid']
    return self.childid

def _tracehist(self, index, contents):
    ##We store revisions we've visited
    ##loops can occur in revision histories
    visited = []
    i = self.historylimit
    j = 0

    self.par['rvprop'] =
        'userid|user|ids|flags|tags|size|comment|contentmodel|timestamp|content'

    while (self.parentid not in visited) and i is not 0 and self.parentid is not 0:
        self.par['revids'] = self.parentid

        self._pace()

        r = requests.get(WIKI_API_URL, params=self.par, headers=self.head)

```

```

r = r.json()

self._rate()

visited.append(self.childid)

#print r

self.childid = r['query']['pages'][self.pageid]['revisions'][0]['revid']
self.parentid = r['query']['pages'][self.pageid]['revisions'][0]['parentid']
user = r['query']['pages'][self.pageid]['revisions'][0]['user']
userid = r['query']['pages'][self.pageid]['revisions'][0]['userid']
size = r['query']['pages'][self.pageid]['revisions'][0]['size']
timestamp = r['query']['pages'][self.pageid]['revisions'][0]['timestamp']
comment = "" #comments sometimes don't return from old revisions...
try:
    comment = r['query']['pages'][self.pageid]['revisions'][0]['comment']
except:
    comment = ""
content = r['query']['pages'][self.pageid]['revisions'][0]['*']

index.writerow([self.pageid, self.childid, user.encode("UTF-8"), userid,
                timestamp, size, comment.encode("UTF-8")])
contents.writerow([self.pageid, self.childid, content.encode("UTF-8")])

if(self.historylimit > 0):
    print self.pageid, "fetch", j+1, "of", self.historylimit, ", revid",
          self.childid, "timestamp", str(timestamp)
    i = i - 1
else:
    print self.pageid, "fetch", j+1, ", revid", self.childid, "timestamp",
          str(timestamp)
j = j + 1
print "limit reached"

def _pace(self):
    if self.rl and self.rl_last_call and self.rl_lastcall + self.rl_minwait >
        datetime.now():
        wait_time = (self.rl_lastcall + self.rl_minwait) - datetime.now()
        time.sleep(int(wait_time.total_seconds()))

def _rate(self):
    if self.rl:
        self.rl_lastcall = datetime.now()

```

A.2 Example output

B Appendix B: Python class for basic, space-naive Levenshtein implementation

B.1 Code

```
import sys

class LevDistBasic:
    e = [] #edit operation array
    t = [] #grid array
    x = "" #string1
    y = "" #string2
    m = 0 #length string1
    n = 0 #length string2
    dist = 0 #Levenshtein distance
    ed = [] #the edit operation, calculated in _calculate()
    isFile = False

    def __init__(self, _x, _y, isFile=False):
        self.x = self._variablehandle(_x)
        self.y = self._variablehandle(_y)
        self.m = len(self.x)
        self.n = len(self.y)
        self.t = [[0]*(self.n+1) for _ in xrange(self.m+1)]
        self.e = [[" "]*(self.n+1) for _ in xrange(self.m+1)]
        self.dist = self._calculate()

    def __str__(self):
        return str(self.distance())

    def distance(self):
        return self.dist

    def strings(self):
        return self.x, self.y

    def table(self):
        return self.t

    def operation(self):
        return self.ed

    ##ADD WARNING for long strings / deal with them
    def showtable(self):
        result = ""
        for ch in self.y:
            result = result + ch + " "
        print " ", result
        for r in range(len(self.t)):
            s = ' '
            if r:
                s = self.x[r-1]
            print s, ' ', self.t[r]

    def showop(self):
        for i, op in enumerate(self.ed):
            l = str(i) + ": "
            if op[0] == 'I':
                l += "insert " + op[-1]
```



```

elif op[0] == 'K':
    l += "keep " + op[-1]
elif op[0] == 'D':
    l += "delete " + op[-1]
elif op[0] == 'S':
    l += "swap " + op[-1][0] + " for " + op[-1][-1]
else:
    return "FAIL: incorrect operation"
print l

def _ed(self):
    i, j = len(self.e)-1, len(self.e[0])-1
    self._ed_recursive(i,j)

def _ed_recursive(self,i,j):
    if self.e[i][j] == ' ':
        if i == 0 and j > 0:
            self.ed.append(('D', self.y[0]))
        if j == 0 and i > 0:
            self.ed.append(('D', self.x[0]))
        return
    if self.e[i][j] == 'K':
        self._ed_recursive(i-1, j-1)
        self.ed.append((self.e[i][j], self.x[i-1]))
    elif self.e[i][j] == 'S':
        self._ed_recursive(i-1, j-1)
        self.ed.append((self.e[i][j], (self.x[i-1] + ',' + self.y[j-1])))
    elif self.e[i][j] == 'D':
        self._ed_recursive(i-1,j)
        self.ed.append((self.e[i][j], self.x[i-1]))
    else:
        self._ed_recursive(i,j-1)
        self.ed.append((self.e[i][j], self.y[j-1]))

def _variablehandle(self,v):
    if not isinstance(v, str):
        try:
            return v.read()
        except:
            try:
                return str(v)
            except:
                print "Argument cannot be of type" + type(v)
                raise
        pass
    return v

def _calculate(self):
    for i in xrange(self.m+1):
        self.t[i][0] = i
    for j in xrange(self.n+1):
        self.t[0][j] = j
    j = 1
    while j < self.n+1:
        i = 1
        while i < self.m+1:
            c = (self.x[i-1] != self.y[j-1])
            dl = self.t[i-1][j] + 1
            ins = self.t[i][j-1] + 1
            sbs = self.t[i-1][j-1] + c
            self.t[i][j] = min(ins, dl, sbs)

```

```
        if ins < dl and ins < sbs:
            self.e[i][j] = 'I'
        elif dl <= sbs:
            self.e[i][j] = 'D'
        else:
            if(self.x[i-1] != self.y[j-1]):
                self.e[i][j] = 'S'
            else:
                self.e[i][j] = 'K'
        i += 1
    j += 1
self._ed()
return self.t[self.m][self.n]
```

B.2 Example output

```
>>> import LevDistBasic
>>> test = LevDistBasic.LevDistBasic("bank","book")
>>> test.showtable()
>>>
      b o o k
      [0, 1, 2, 3, 4]
b    [1, 0, 1, 2, 3]
a    [2, 1, 1, 2, 3]
n    [3, 2, 2, 2, 3]
k    [4, 3, 3, 3, 2]

>>> t = test.table()
>>> print t
>>>
[[0, 1, 2, 3, 4], [1, 0, 1, 2, 3], [2, 1, 1, 2, 3], [3, 2, 2, 2, 3], [4, 3, 3, 3, 2]]

>>> s = test.strings()
>>> print s
>>>
('bank', 'book')

>>> test.showop()
>>>
0: keep b
1: swap a for o
2: swap n for o
3: keep k

>>> ed = test.operation()
>>> print ed
>>>
[('K', 'b'), ('S', 'a,o'), ('S', 'n,o'), ('K', 'k')]

>>> print test
>>>
2

exit()
```

Bibliography

- [1] Acawiki. <http://acawiki.org>. Accessed: 2014-04-31.
- [2] Alexa about us. <http://www.alexa.com/about>. Accessed: 2014-04-31.
- [3] Alexa the top 500 sites on the web. <http://www.alexa.com/topsites>. Accessed: 2014-04-31 (updated daily).
- [4] Python wikipedia 1.3.0. <https://pypi.python.org/pypi/wikipedia>. “Wikipedia is a Python library that makes it easy to access and parse data from Wikipedia.” Accessed: 2014-04-31.
- [5] Wikilit. <http://wikilit.referata.com>. A temporary site, relating directly to [19]. Accessed: 2014-04-31.
- [6] Wikimedia enwiki dump progress on 20140502. <http://dumps.wikimedia.org/enwiki/20140502/>. Accessed: 2014-04-29.
- [7] Wikimedia wikipedia.org is more popular than...: Note on alexa rankings. https://meta.wikimedia.org/wiki/Wikipedia.org_is_more_popular_than...#Note_on_Alexa_rankings. Accessed: 2014-04-31.
- [8] Wikipapers. <http://wikipapers.referata.com>. Accessed: 2014-04-31.
- [9] Wikipedia counter-vandalism unit. <https://en.wikipedia.org/wiki/Wikipedia:CVU>. Accessed: 2014-04-31.
- [10] Wikipedia wikipedia signpost/2013-02-04/special report: Examining the popularity of wikipedia articles: catalysts, trends, and applications. https://en.wikipedia.org/wiki/Wikipedia:Wikipedia_Signpost/2013-02-04/Special_report. Accessed: 2014-04-31.
- [11] B. Thomas Adler and Luca de Alfaro. A Content-driven Reputation System for the Wikipedia. In *Proceedings of the 16th International Conference on World Wide Web, WWW '07*, pages 261–270, New York, NY, USA, 2007. ACM.
- [12] Peter Denning, Jim Horning, David Parnas, and Lauren Weinstein. Wikipedia risks. *Commun. ACM*, 48(12):152–152, December 2005.
- [13] Jim Giles. Internet encyclopaedias go head to head. *Nature*, 438(7070):900–901, December 2005.
- [14] Aniket Kittur, Bongwon Suh, Bryan A. Pendleton, and Ed H. Chi. He says, she says: Conflict and coordination in wikipedia. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems, CHI '07*, pages 453–462, New York, NY, USA, 2007. ACM.
- [15] VI Levenshtein. Binary Codes Capable of Correcting Deletions, Insertions and Reversals. *Soviet Physics Doklady*, 10:707, 1966.
- [16] Michael Lieberman and Jimmy Lin. You are where you edit: Locating wikipedia contributors through edit histories, 2009.
- [17] Andrew Lih. Wikipedia as participatory journalism: reliable sources? metrics for evaluating collaborative media as a news resource. In *In Proceedings of the 5th International Symposium on Online Journalism*, pages 16–17, 2004.
- [18] Teun Lucassen and Jan Schraagen. Evaluating wikitrust: A trust support tool for wikipedia. *First Monday*, 16(5), 2011.

- [19] Mostafa Mesgari, Chitu Okoli, Mohamad Mehdi, Finn Årup Nielsen, and Arto Lanamäki. “the sum of all human knowledge”: A systematic review of scholarly research on the content of wikipedia. *Journal of the American Society for Information Science and Technology*, April 2014. This is a postprint of an article accepted for publication in Journal of the American Society for Information Science and Technology copyright © 2014 (American Society for Information Science and Technology).
- [20] Mrton Mestyn, Taha Yasseri, and Jnos Kertsz. Early prediction of movie box office success based on wikipedia activity big data. *CoRR*, abs/1211.0970, 2012.
- [21] Martin Potthast, Benno Stein, and Robert Gerling. Automatic Vandalism Detection in Wikipedia. In Craig Macdonald, Iadh Ounis, Vassilis Plachouras, Ian Ruthven, and RyenW. White, editors, *Advances in Information Retrieval*, volume 4956 of *Lecture Notes in Computer Science*, page 663668. Springer Berlin Heidelberg, 2008.
- [22] B. Stvilia, M.B. Twidale, L.C. Smith, and L. Gasser. Assessing information quality of a community-based encyclopedia. In *Proceedings of the International Conference on Information Quality*, pages 442–454, 2005.
- [23] Bongwon Suh, Gregorio Convertino, Ed H. Chi, and Peter Pirolli. The singularity is not near: Slowing growth of wikipedia. In *Proceedings of the 5th International Symposium on Wikis and Open Collaboration*, WikiSym ’09, pages 8:1–8:10, New York, NY, USA, 2009. ACM.
- [24] J Voss. Measuring wikipedia. *Proceedings 10th International Conference of the ...*, January 2005.
- [25] Dennis M. Wilkinson and Bernardo A. Huberman. Assessing the value of cooperation in wikipedia. *CoRR*, abs/cs/0702140, 2007.