

An abstract graphic on the left side of the slide. It features several vertical bands of color: a dark teal band at the top left, a light teal band below it, a light blue band, a light purple band, and a light pink band. Overlaid on these bands are several wavy, ribbon-like shapes in various colors including blue, purple, teal, yellow, and orange, creating a dynamic, flowing effect.

Analyzing web archives

Using Python and
Solrwayback to explore
Covid-19 content on Niagara
Region municipal webpages

Introduction

- Fletcher Johnson - Research Assistant @ Brock
- B.Sc Computer Science, UofT (2010), MLIS Western University (2021)
- Previous industry experience, transition into library work

“Crisis Communication in the Niagara Region during the COVID-19 Pandemic”

- Analysis of 13 municipal webpages from April, 2020 – December, 2021
- Examination of how municipalities talked about Covid during the crisis
- Team members: Tim Ribaric, David Sharron, Cal Murgu, Dr. Karen Louise Smith, Dr. Duncan Koerber, and myself.
- My role: Technical support (Solrwayback + Notebooks)

The Data

- 300~ GB, captured using Archive-it.org service (IA commercial arm)
- WARC format
- Full datasets with API
- Derivatives accessible using Archives Research Compute Hub (ARCH)

The Tools: SolrWayback

- Project by Royal Danish Library for viewing Danish Internet (2005-onwards)
- “Search engine for WARC’s”
- Like Wayback machine, but better

SolarWayback – Why?

- Much faster than Wayback machine –
 - We host ours on an AWS t2.small (1 cpu, 2 GB ram)
- Full text search (Solr)
- Complex queries with Booleans, wildcards, and facets
- No API support, though CSV data export available.

SolrWayback Facets

content	content
content_encoding	content_encoding
content_language	url
content_type	url_path
crawl_date	title
domain	status_code
links	public_suffix
links_domains	... and more
links_images	
source_file	
type	

GPS IMAGE SEARCH TOOLBOX

RESULTS

See available export options

Previous 20

Next 20

LORD MAYOR AND COUNCIL | NIAGARA ON THE LAKE

Type: [html](#), [web page](#) @ [notl.com](#)

Url: <https://www.notl.com/content/lord-mayor-and-council>

[See all images](#)

- html (841,776)
- other (385,462)
- image (266,563)
- text (100,611)
- video (2,066)
- pdf (1,681)
- audio (47)

View data fields



SEARCH HINTS

Two colons without quote signs. When a qualified search is performed, consider quoting the value - (domain:portcolborne.ca OR domain:welland.ca) AND (content_type_norm:html AND links_domains:niagararegion.ca AND crawl_date:"[2021-01-01T00:00:00Z" TO 2022-02-01T00:00:00Z]) covid

(domain:portcolborne.ca OR domain:welland.ca) AND (content_type_norm:html AND links_domains:niagararegion.ca AND crawl_date:"[2021-01-01T00:00:00Z" TO 2022-02-01T00:00:00Z]) covid

☐ GROUPED SEARCH [?]

☐ IMAGE SEARCH [?]

☐ URL SEARCH [?]

SEARCH WITH UPLOADED FILE

GPS IMAGE SEARCH 

TOOLBOX 

FACETS

domain

- portcolborne.ca (658)
- welland.ca (21)

content_type_norm

- html (679)

type

- web page (679)

crawl_year

- 2021 (679)

status_code

- 200 (679)

public_suffix

- ca (679)

RESULTS

Showing 0 - 20 of 679 entries matching query.

[See available export options](#)

Previous 20

Next 20

COVID-19

#1

score: 39.759163

Type: html, web page @ welland.ca

Date: 27/08-2021

Url: <https://www.welland.ca/hottopics/COVID-19.asp>

Highlighted content:

"Skip to Main Content Search for: Menu Listen [COVID-19](#) information details Welland Response"

Images: showing 4 out of 4



[View data fields](#)



News / Events

[News Releases](#)[COVID-19 Media Briefings](#)[Regional Budget Info](#)[Councillor Profiles](#)[Upcoming Events](#)[Public Notices](#)[Media Contacts](#)

Preparation for COVID-19 Immunizations for Residents 80+

Niagara Region Public Health is busy preparing for the launch of our COVID-19 mass immunization clinics in the coming weeks.

Last week the 11 clinic locations were announced, and Public Health continues to work with these community partners to be

HARVEST DATE: 2021-04-30 21:43:51 HTTP status code: 200
URL: https://www.niagararegion.ca/news/article.aspx?news=1171&t=Preparation%20for%20COVID-19%20Immunizations%20for%20Residents%2080%2B&fbclid=IwAR08FQCdinQPfilwUTT8v-6SeVUcleqd0uS0GGtrMSDfSY_GnSnsbbeeius
#Harvested: 8
DOMAIN: niagararegion.ca #Harvested: 20745 #Content length harvested: 233706065
PAGE RESOURCES: #Found: 15 #Not found: 3

[Close](#) [Hide](#)

se in frequency
a's 80+ residents as



Harvest
calendar



PWID xml



Page
previews



View page
resources

First: [2021-03-12 22:18:38](#) Previous: [2021-04-23 23:26:52](#) Next: none Last: [2021-04-30 21:43:51](#)

- Drive time for large portions of our overall population
- Accessible by public transit
- Connectivity to cellular internet services
- Public safety implications, with a focus on possible impacts to major roadways

Once vaccine supply is increased, the clinics will be supplemented by local pharmacies and family doctors, so that by the time the majority of the population can receive their vaccine, they can do so even closer to home.

Clinics for residents 80 years of age and older will be starting within the next couple weeks, and residents will be able to book their appointments through the Provincial online and phone registration system that is scheduled to launch on Monday, March 15.

As soon as we have more information about it, including the website link, phone number, and any other pertinent information, we will share it with the public.

We are grateful for the large number of people who responded to our call for volunteers, and we are still looking to fill a small number of paid staff positions to help run our clinics. More information about these postings can be found on our [career portal](#).

The eagerness for vaccines rollout to ramp up over the coming weeks and months is appreciated and echoed by everyone at Public Health. We ask for continued patience as this next phase begins, as there is a lot of interest in vaccine appointments across the province, and wait times in the registration portals, both online and by phone, are possible.

Harvests for: <http://welland.ca/hottopics/covid-19.asp>

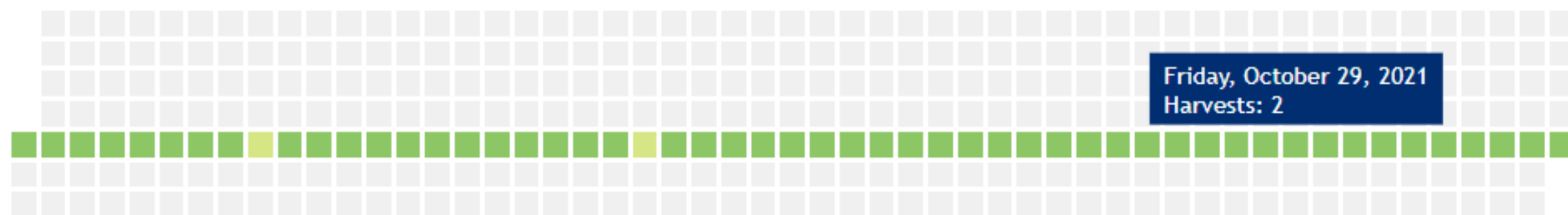
First harvest: **April 17, 2020**

Latest harvest: **December 31, 2021**

Total harvests: **179**

Show: [Months](#) - [Days](#)

2021 - [Hide details](#)



Less  More

Harvests for October 29, 2021

1. [October 29, 2021 15:54](#)
2. [October 29, 2021 15:55](#)

The Tools: Python notebooks

- Notebooks
 - Mixed media: (Python) code, rich-text, results (pictures, videos)
 - Easy to distribute and use
 - Results can be viewed without running code
 - Browser based, often run in the cloud
- Why use notebooks at all?

Answer: Sharable integrated environment with repeatable results.

Google Colab

Google Colab is an environment for running notebooks

- Available (free) to everyone, literally no config necessary
- No issues with dealing with infrastructure (GPUs offered!)
- Complete environment within browser – very portable

The Tools: Python + Google Colab

- Why is it useful for our project?
 - Distant reading (quantitative results)
 - Sharable graphs + and results
 - Drastically lowers barrier to entry for modifications
 - Reproducible results (post-publication)










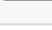






Our notebooks use ...

- Text analysis libraries used extensively. (Scikit-learn, nltk, SpaCy)
- Pandas
- Matplotlib
- Github for notebooks, S3 for data (ARCH full text derivation)

Notebook examples

- Sentiment analysis
- Text similarity (of URLs)
- Update frequency + size of updates
- First mention of vaccine by municipality

Notebook listing

ARCH Data Exploration	Notebook	 Open in Colab
COMM 4P35 Tutorial	Notebook	 Open in Colab
Hackfest notebook	Notebook	 Open in Colab
Muni Data Export	Notebook	 Open in Colab
Prep Domain Data	Notebook	 Open in Colab
Twitter Data Export	Notebook	 Open in Colab
Municipal Data Similarity	Notebook	 Open in Colab
Another example of Municipal Data Similarity using SpaCy	Notebook	 Open in Colab
Municipal Data Similarity using TF-IDF	Notebook	 Open in Colab
Content size of pages over time	Notebook	 Open in Colab
Frequency of page updates over time	Notebook	 Open in Colab
Crawl frequency visualized	Notebook	 Open in Colab
Sentiment scores of Municipal Data	Notebook	 Open in Colab
Vaccine keyword frequency	Notebook	 Open in Colab
Vaccine keyword frequency normalized	Notebook	 Open in Colab
First mention of vaccine keywords	Notebook	 Open in Colab

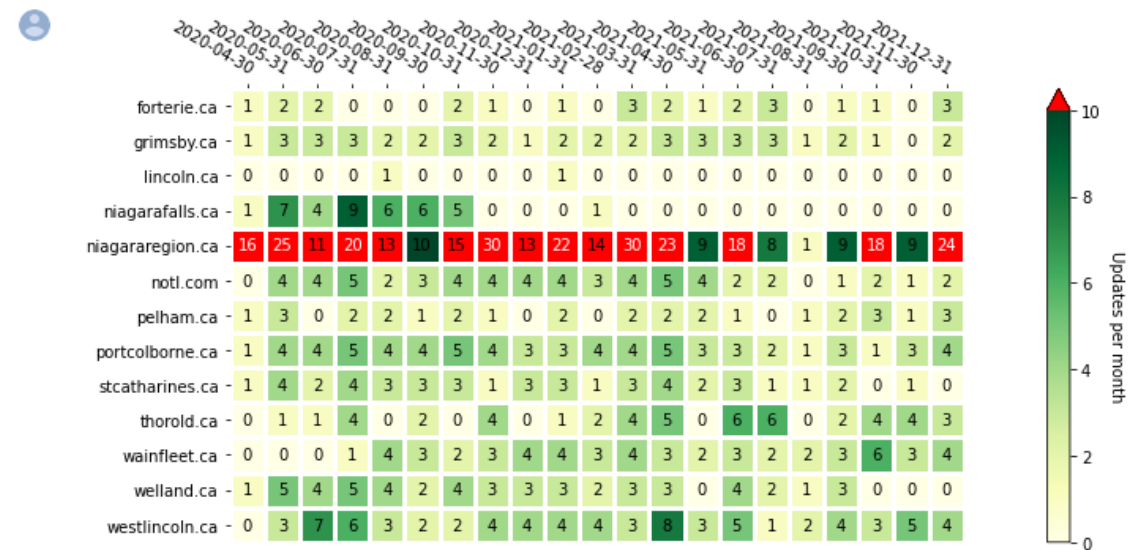


+ Code + Text Copy to Drive

Update frequency

```
[ ] 1 # 0.0 values in the delta_table mean that no updates occurred
    2 # set these values to NaNs so the groupby count() aggregator ignores them
    3 dt_cp = delta_table.copy()
    4 dt_cp[dt_cp == 0] = np.nan
    5
    6 update_table = dt_cp.groupby(pd.Grouper(freq='1M')).count().transpose().groupby(level=0).sum()
    7 update_table
```

```
1 fig, ax = plt.subplots(figsize=(15,5))
2
3 column_labels = list(map(lambda ts: ts.strftime('%Y-%m-%d'), update_table.columns))
4 im, cbar = heatmap(update_table, update_table.index, column_labels, ax=ax, cmap='YlGn',
5                   cbar_kw={'extend':'max'}, cbarlabel="Updates per month", vmin=0, vmax=10)
6 texts = annotate_heatmap(im, valfmt="{x}", threshold=15)
7
8 fig.tight_layout()
9 plt.show()
```





Brock University Library
Digital Scholarship Lab

About

This notebook creates two heatmaps based on the number of updates and size of updates released by each municipality.

The methodology for determining update size is as follows:

1. For each municipality, determine the content length - the character count - of all crawls.
2. For each crawl calculate the difference in content length compared to the previous time it was crawled.
3. For each url calculate the average of the differences over a month timeframe.
4. For each municipality calculate the average of the monthly url averages.

The way the *number of updates* is determined is very similar. The number of crawls for each municipality are determined. The difference in content length compared to a *previous crawl* is calculated. If there is a difference it is considered an update. The *updates* are then summed over a period of a month.

Limitations

- The update size heat map is quite large. Making it smaller is difficult because it becomes hard to read the numbers.
- The assumption is that the character count of each crawl is indicative of information related to covid but this is not necessarily true

Take away

- Notebooks place code, data, results all in one place
- Google Colab makes running notebooks incredibly easy
- SolrWayback is a powerful WARC search engine: useful for anyone with large collection of WARCs
- Both make analyzing web archives much easier

Links

- BrockDSL "Covid-19 in Niagara" Notebooks
https://github.com/BrockDSL/ARCH_Data_Explore
- SolrWayback
<https://github.com/netarchivesuite/solrwayback>
- Our Archive-it.org web archive
<https://archive-it.org/collections/13781>
- More about Google Colab
<https://research.google.com/colaboratory/faq.html>