

# IPSW - Modelling Change of Website Archives

Group 4

---

Ian Milligan, Ian Roper, Caoimhe Rooney,  
Nathan Taback, Jessica Williams, Nich Worby

May 9, 2019

# The Problem

- Decades of web archives with lots of vital historic information
- Gets overridden from history when the website is updated
- Ever increasing amount of data
- Scale overwhelms the search for the meaningful information

# The Problem

- Decades of web archives with lots of vital historic information
- Gets overridden from history when the website is updated
- Ever increasing amount of data
- Scale overwhelms the search for the meaningful information

Can we

- find out when large changes have occurred?
- predict when a big change is going to occur?

## Aims for this week

- Find and explore ways to quantify change in a website
- Compare these quantifications
- See if we can identify big changes in an organisation from our research

# Big events in the NDP

- 28 November 2005: election called.
- 23 January 2006: federal election.
- 14 October 2008: federal election.
- 2 May 2011: federal election.
- July 2011: NDP leader announces leave of absence; replaced by interim.
- 22 August 2011: NDP leader dies.
- 24 March 2012: New NDP leader selected.
- 19 October 2015: federal election.
- 10 April 2016: NDP leader loses vote of confidence.
- 1 October 2017: New NDP leader selected.

## Attempted Approaches

- How many words on the domain change?
- How do the links out of the domain change?
- How does the way the website looks changes?
- How does the structure of the websites within the domain change?

## Four Metrics for Text

- Byte-wise comparison:
  - If any change in characters has occurred, = 1
  - If text is *exactly* the same, = 0
- TF-IDF
  - Calculates cosine distance between two different vectors of characters  $p$  and  $p'$
- Word distance
  - How many words have changed
- Edit distance
  - “Edit distance”  $\delta$  is the amount of insertion/deletion/substitution needed to turn one sequence into the other

$$1 - \frac{p \cdot p'}{\|p\|_2 \|p'\|_2}$$

$$1 - \frac{2|common\ words|}{m + n}$$

$$\frac{\delta}{m + n}$$

- Justification
  - Links to other websites are important to the website designer
  - If these change, the topic of the website has most likely changed as well
- Method
  - Compare vector of links on homepage at  $t_i$  and  $t_{i+1}$  as  $\mathbf{v}_i$  and  $\mathbf{v}_{i+1}$

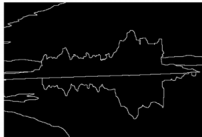
$$1 - \frac{2|\text{common links}|}{|\mathbf{v}_i| + |\mathbf{v}_{i+1}|}$$



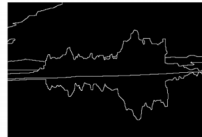
# Screenshot Comparison



(q)



(r)



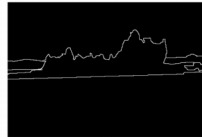
(s)



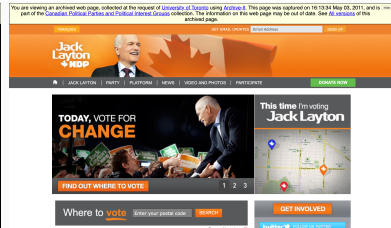
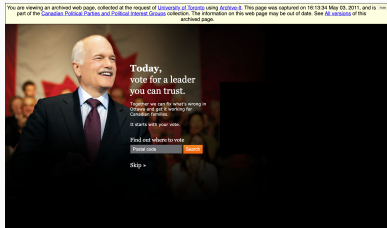
(u)



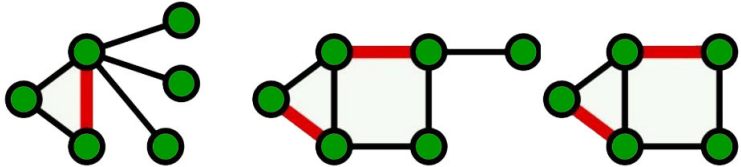
(v)



(w)



# Structure



<https://www.geeksforgeeks.org/mathematics-matching-graph-theory/>

- Data takes a long time to retrieve from servers
- Different languages
- Failed renderings
- Difficult to decipher where website fits in structure just from URL
- Internal links data showing strange behaviour