

# PHILOSOPHICAL WRITING IN EARLY NEW ZEALAND NEWSPAPERS

---

Joshua Black

February 5, 2021

UC  ARTS  
DIGITAL LAB

## 1. Problem:

- to gain insight into philosophical writing in early New Zealand newspapers

## 2. Method:

- corpus construction via self-labelling and supervised learning,
- corpus analysis with, e.g., co-occurrence networks.

## 3. Results

## 4. Upshot:

- a method applicable for many humanities research questions,
- but with shortcomings to be aware of.

## PROBLEM

---

- Organisation:
  - UC Arts Digital Lab
  - ...part of the digital humanities department.
  - Digital humanities: use of computational (and data science) techniques to achieve insight into cultural products and practices.
- Myself:
  - research background in traditional humanities research (in philosophy),
  - with newly developed data science skills.

# HUMANITIES PROBLEM

- Histories of philosophy in New Zealand don't have much to say before the mid-twentieth century:
  - 'many of those who had longstanding chairs published next to nothing' (Davies and Helgeby 2014, 24).
- An explanation: excessive on focus academic publications.
  - ...and on *academic* philosophy.
- Newspapers as an alternative source:
  - 'the fundamental infrastructure for intellectual life ...
  - ...newspapers were ascendant in New Zealand because imported books were expensive and a sustainable local periodical literature was slow to emerge' (Ballantyne 2012, 57–78).
  - Promising: a venue both for the academics and the wider public?
- An example of the kind of thing we're after:

The following is a brief abstract of  
the Debate held at the Town Hall  
East Oxford, on Thursday, 9th.

*(Continued from last week.)*

6. A simple form of metaphysical argument may be briefly put as follows:—All existence are of two kinds necessary and contingent. By a necessary existence is meant one which never began to be, and can never cease to be. By a contingent existence is meant one which commenced ~~to be and will cease to be.~~ My exist-

# DATA SCIENCE METHODS

- Data source: the National Library Papers Past Newspaper Open Data Pilot.
  - A dataset containing the output of OCR for newspaper data in English up to 1900.
  - Made available in 2020 to encourage digital experimentation.
  - 'Big' by human standards: 1,471,384 pages of content. (315GB compressed)
  - Ethics:
    - all out of copyright, no living people discussed, but ...
    - some offensive material present.
  - <https://natlib.govt.nz/about-us/open-data/papers-past-metadata/papers-past-newspaper-open-data-pilot/>
- Data science methods needed to:
  - Find the relevant material (it's a small portion of the data set!)
  - Derive insight from it once found (it's still a lot of text!)
  - We will engage in 'distant reading' (Moretti 2013)

# THE PROJECT

## 1. Corpus construction

- Aim: find the relevant material in the dataset.
- Method: labelling articles and training Naive Bayes classifiers following a 'bootstrapping' pattern.

## 2. Corpus analysis

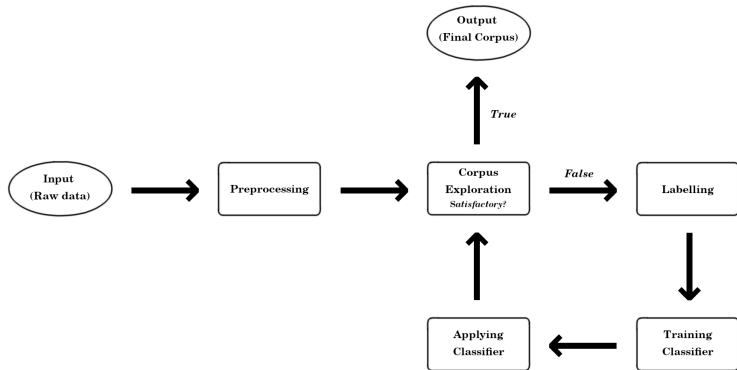
- Aim: use the corpus to learn something about philosophical writing in NZ newspapers.
- Method: many text analysis methods, including concordancing, collocations, co-occurrence networks, and topic modelling.
- This presentation will focus on co-occurrence networks.



## METHOD

---

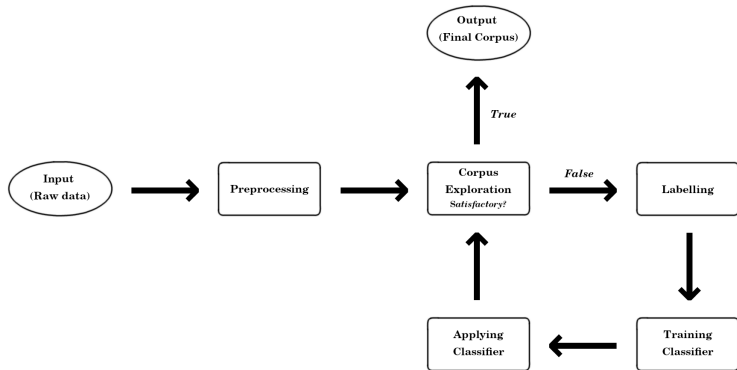
# CORPUS CONSTRUCTION FLOW DIAGRAM



# PREPROCESSING (XML → PANDAS)

- Data format for each issue:
  - a METS file with structural information and an ALTO file for each page.
- Throw away:
  - Spatial information (e.g. location of each word on page)
  - Items tagged as advertisements.
- Method:
  1. iterate through title-year tarballs,
  2. for each issue, collect list of articles and corresponding text blocks from METS file,
  3. iterate through ALTO files, collecting text blocks for each article
  4. gather all in Pandas dataframe with each row corresponding to an individual article.
- Result:
  - 7,592,619 distinct articles and their plain text,
  - ...stored as eight pickled dataframes (around 8GB total).

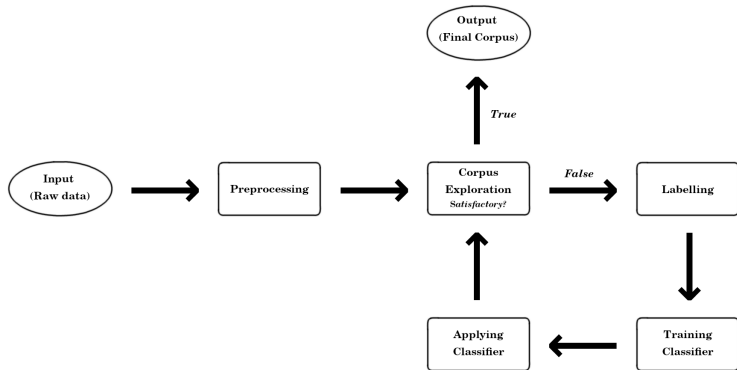
# CORPUS CONSTRUCTION FLOW DIAGRAM



# CORPUS EXPLORATION

- First stage: look at REGEX matched for 'philoso\*'
- Lots of methods used to pick out desired and non-desired articles:
  - inspecting random articles, keyword searches, word clouds, concordancing, collocations, co-occurrence networks ...
  - ...more on co-occurrence networks later.
- If the corpus contains lots of material that we are not interested in, it is not 'satisfactory'. If so, we move to the next stage.

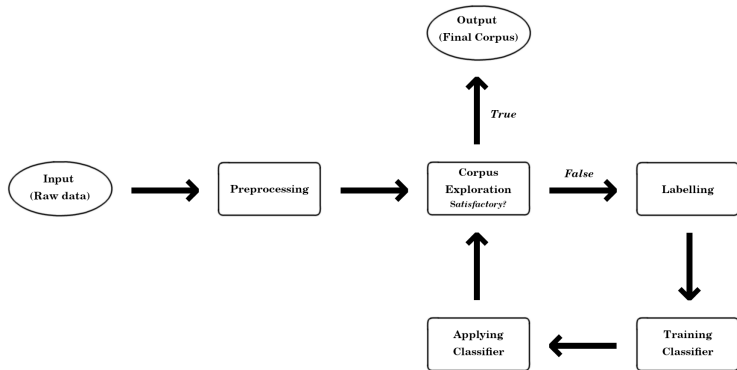
# CORPUS CONSTRUCTION FLOW DIAGRAM



# LABELLING

- Label articles in order to train a classifier to find what we are after.
- Two key labels for the project:
  - Philosophy: is the majority of the article 'philosophical discourse'?
    - A broad definition: does it develop or discuss ideas 'ultimate reality' or 'ultimate value'.
    - e.g.: is there life after death, are there multiple sources of knowledge, what is the best way to organise society and why?
    - Some reliance on my own experience in studying 19th century philosophy.
  - Philosophy type:
    - Is it about ethics, the relationship between religious belief and modern thought, metaphysics and epistemology, or other?
    - The relationship between religious belief and, e.g., evolution is a very prominent topic at this time.
  - Also attempted, but not used: 'Readable', 'Writing Type', 'NZ author'
- NB: it is important to ensure that we label a wide range of non-philosophy.

# CORPUS CONSTRUCTION FLOW DIAGRAM

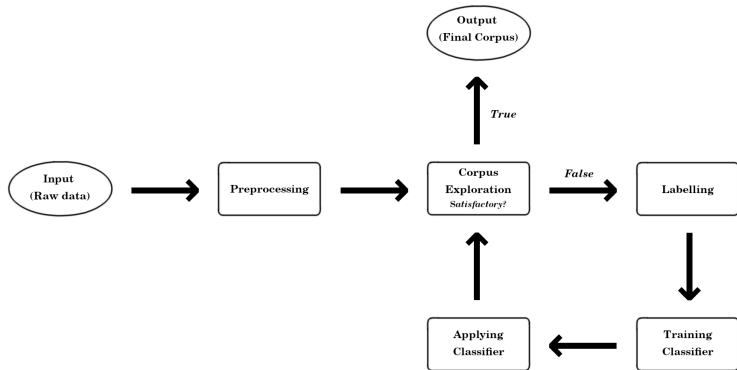




# TRAINING AND APPLYING CLASSIFIER

- Train a classifier to distinguish ‘philosophy’ articles from ‘non-philosophy’.
- Classifiers tried: Naive Bayes and Support Vector Machines
- Naive Bayes is simple and fast to train, while performing remarkably effectively.
- Training and testing data divided, and training data resampled, as appropriate given the state of the labelled collection.
- Classification algorithms implemented using Scikit-learn Pipelines.
  1. Text to bag of words (params: size of feature space)
  2. TF-IDF transformation
  3. The Naive Bayes classifier (params: prior on representativeness of labelled set.)
- Grid CV search used to select parameters for the classifier.
- ...with accuracy, recall, and precision tried as measures.
- Trained classifier applied to complete dataset, treating the ‘philosophy’ articles as a new candidate corpus.

# CORPUS CONSTRUCTION FLOW DIAGRAM



# 'BOOTSTRAPPING'

- The phrase: 'pull yourself up by the bootstraps',
- In this case:
  1. starting with nothing, we add articles to our labelled collection,
  2. having collected a good number ( $\sim 200-300$ ), with much higher representation of philosophy than the general dataset,
  3. we train and apply a classifier,
  4. we use the articles classifier as philosophy as a source of new articles to label.
- NB: after the first classifier has been applied, the new non-philosophy articles added to the labelled collection will have been classifier as philosophy by the previous classifier.
  - ...this means that subsequently trained classifiers can be more selective.
  - ...we need a picky classifier.
  - Stop when satisfied that the corpus does not contain too much 'non-philosophy'

## AIM 2: CORPUS ANALYSIS

- Corpus analysis ~ the corpus exploration step.
- One method used: co-occurrence networks:
  1. create bag of words or TF-IDF representation of documents;
  2. compute a document-term matrix (num terms x num documents) for either the BOW or TF-IDF representation;
  3. compute a term-term matrix (num terms x num terms), representing co-occurrences in documents of each pair of terms; and
  4. given a search term, pass these matrices to a statistical function to return the most closely related words.
    - Statistics implemented: log Dice and mutual information.
- Networks displayed using Dash cytoscapes:
  - see project dashboard:  
**`nz-newspaper-philosophy.herokuapp.com`**
- NB: term-term matrices can get very large. Various methods to control the size of the dictionary were employed.

## RESULTS

---

# CORPUS CONSTRUCTION (SIZE REDUCTION)

Classifiers become more selective:

Corpus	Article Count
Processed dataset	7592619
(Step 0) 'philoso*' Corpus	29647
(Step 1) Naive Bayes 1	239649
(Step 2) Naive Bayes 2	31131

**Table:** Article counts for processed dataset and general philosophy corpora.

## WORD CLOUD: SAMPLE OF PROCESSED DATASET

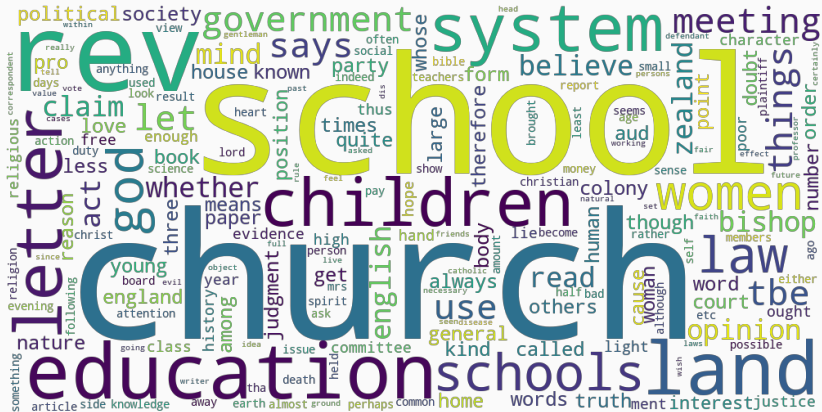


## WORD CLOUD: STEP 0





## WORD CLOUD: STEP 1





# LABELLING

Label	Step 1		Step 2	
	Value	Count	Value	Count
Readable	True	247	True	918
	False	26	False	41
Philosophy	True	101	True	299
	False	147	False	620
Philosophy Type	Religion-Science	58	Religion-Science	140
	Ethics-Politics	25	Ethics-Politics	94
	Epistemology-Metaphysics	3	Epistemology-Metaphysics	13
	Other	15	Other	52
Writing Type	Public Event	40	Public Event	97
	Letter to the Editor	23	Letter to the Editor	69
	First-order Writing	36	First-order Writing	111
	Review	2	Other	22
NZ	True	77	True	178
	False	12	False	41

**Table:** Label counts at Step 1 and Step 2.

## CLASSIFIER PERFORMANCE (STEP 2)

		Predicted	
		False	True
Actual	False	181	14
	True	15	62

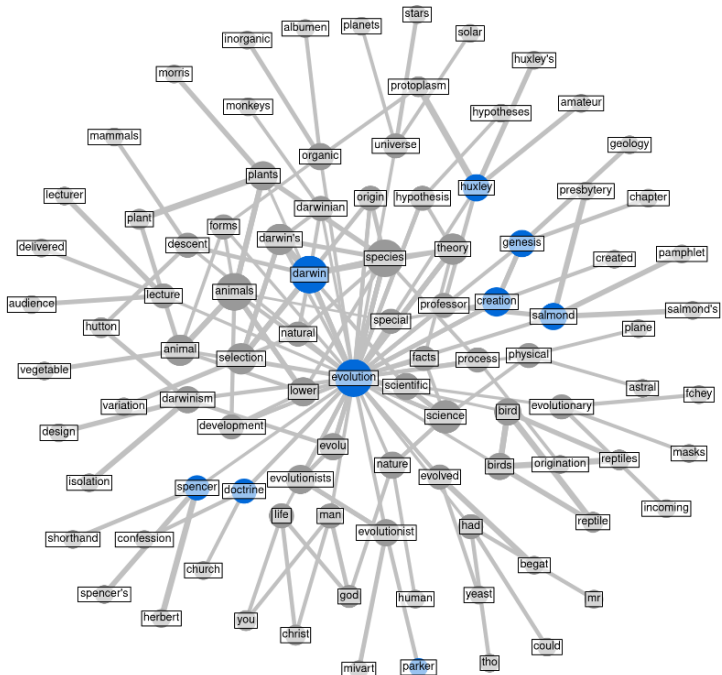
**Table:** Confusion Matrix for Second Naive Bayes Classifier

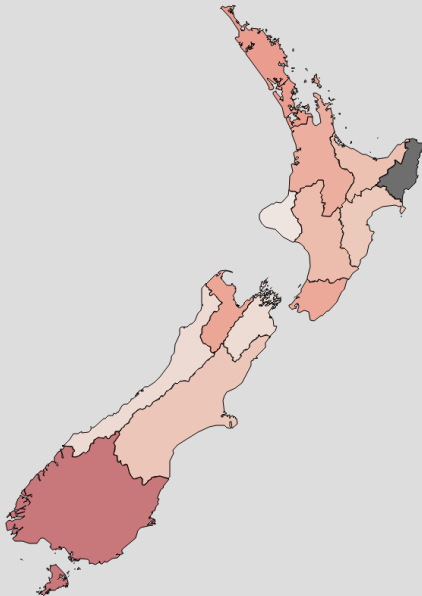
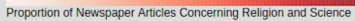
- Accuracy: 0.89
- Precision: 0.81
- Recall: 0.80

# CORPUS PERFORMANCE

- Metrics are not bad, but not great either.
- Inspection of the false positives and negatives reveals:
  - prevalence of 'composite articles' in the false negatives, where only one or two bits are philosophical; and
  - prevalence of 'edge cases' and mistaken labelling in the false positives.
- Conclusion: the performance of the classifiers is being limited by the quality of the labels.

- Many methods were used at this stage (see project report)
- The following two slides contain:
  1. an example of a telling co-occurrence network, and
  2. a choropleth revealing prominence of religion and science discourse in different regions.







UPSHOT

---

1. Single investigator, self-labelled corpus production using METS/ALTO digitised newspaper files is feasible.
  - Generalisability: METS/ALTO is the standard for newspaper digitisation, so the same methods could be applied in other countries.
2. The corpus produced at the corpus construction stage shows potential for research into early New Zealand philosophy.

## 1. Problem 1:

- Many articles at the time were made up of lots of distinct bits (especially editorials).
- Since the classifier loses many of these articles, the resulting corpus is not fully representative of philosophical discourse in early NZ newspapers.
- Possible solution: label text blocks rather than articles.

## 2. Problem 2:

- Labelling criteria were insufficiently clear.
- Better labelling might improve classifier performance.
- Possible alternative: start with easier distinctions (e.g. is the article a report of a public lecture?), then move to subject matter distinctions.

## 1. Problem:

- to gain insight into philosophical writing in early New Zealand newspapers

## 2. Method:

- corpus construction via self-labelling and supervised learning,
- corpus analysis with, e.g., co-occurrence networks.

## 3. Results:

- An interesting collection of articles for digital humanities research,
- with indications that labelling could be improved,
- and ability to reveal features of philosophical discourse about relationship between religious belief and then-new scientific ideas.

## 4. Upshot:

- a method applicable for many humanities research questions,
- but with shortcomings to be aware of.

- Dashboard: `nz-newspaper-philosophy.herokuapp.com`
- GitHub (full project):  
`github.com/JoshuaDavidBlack/NPOD_Philosophy`
- GitHub (dashboard): `github.com/JoshuaDavidBlack/NPOD_Philosophy_Heroku`