

Data Challenge Workshop: Classification of business-related News Articles

Christian Ritter, Data Science Accelerator, Enterprise Statistics Division

www.statcan.gc.ca

August 20, 2018

christian.ritter@statcan.gc.ca

Code & data available at

https://github.com/chritter/Talks/tree/master/2018/Data_Challenge_UOttawa

100

STATISTICS CANADA

ONE HUNDRED YEARS AND COUNTING

STATISTIQUE CANADA

CENT ANS BIEN COMPTÉS



Statistics
Canada

Statistique
Canada



Statistique
Canada

Statistics
Canada

Data Science
Accelerator

Enabling you.

Canada

II. Our team & mandate

Sevgui Erman, PhD,
Assistant Director,
DSA Lead

Christian Ritter, PhD

Nic Denis, MSc

Monica Piccard, MSc
Unit Head

Joanne Yoon, BSc

Stan Hatko, MSc

100



Saeid Molladavoudi,
PhD



Statistics
Canada Statistique
Canada

Accelerator
Enabling you.

Canada

Motivation

The problem

- Business analysts in enterprise statistics division spend large amounts of time to search for information
- Amount of news exponentially growing
- What happened? When? Where? Who?

The solution

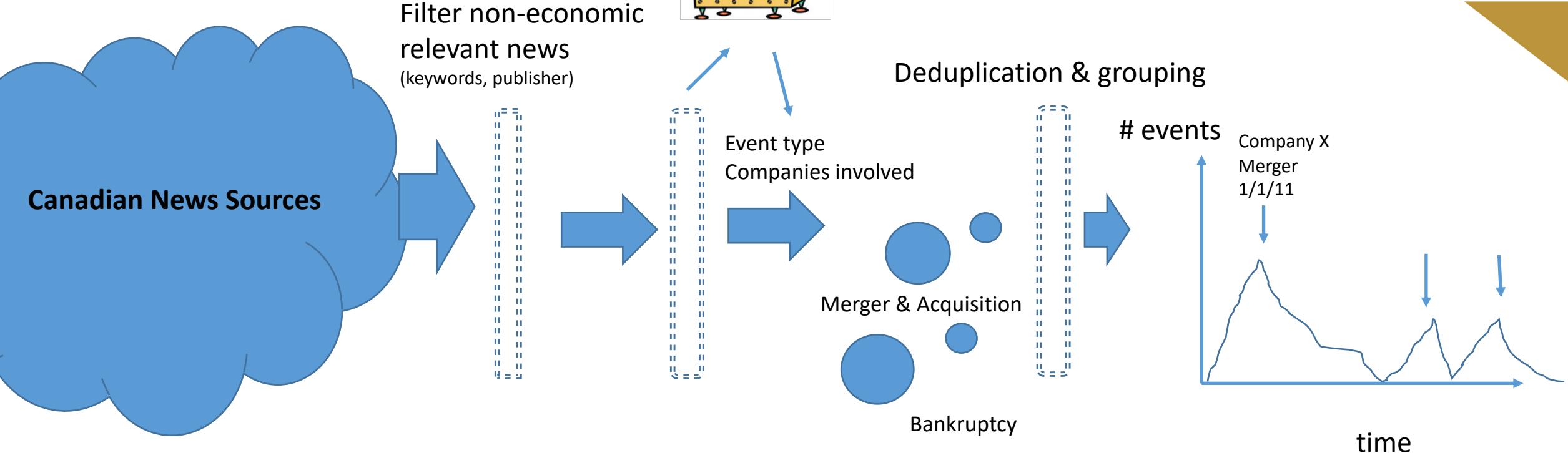
- Automatic analysis of text
- Capture more information faster
- Event and trend detection of unexpected news
- Knowledge extraction for companies and industries



Statistics
Canada Statistique
Canada



Approach

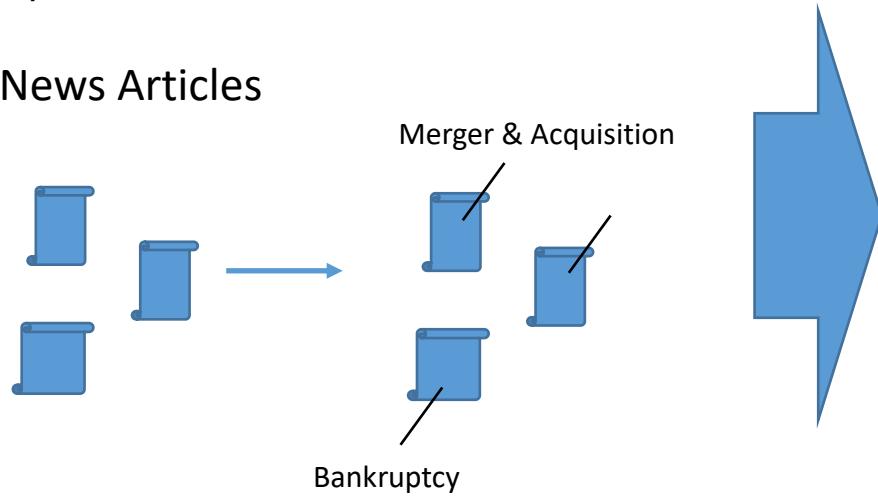


Statistics Canada
Statistique Canada

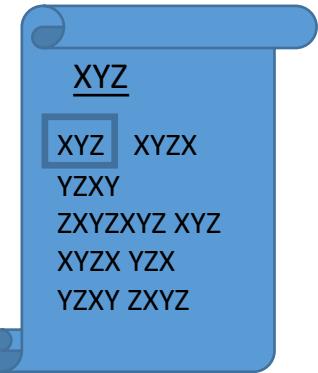
The event detector: Supervised machine learning

A) Creating a labeled dataset

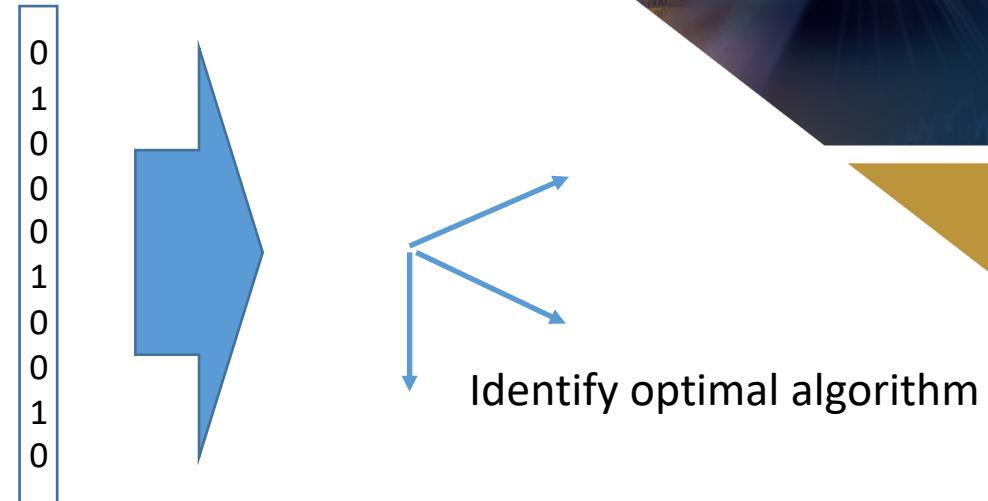
News Articles



B) Feature engineering



C) Classification



- Document-level event detection
- Name-entity recognition for deduplication and to identify companies involved



Statistics
Canada



Statistique
Canada

Statistics
Canada

Data Science
Accelerator

Enabling you.

Canada

Feature engineering



100

- Major challenge in data science workflow
- Design characteristics from which your algorithm can learn to distinguish between classes
- In our case: Lexical and syntactic features
- For example:
 - Type of word, compound (n-gram) words
 - Weight occurrence of each word: Term frequency–Inverse document frequency (tfidf)

t: term

d: document

D: all documents

$$\text{tfidf}(t, d, D) = \text{tf}(t, d) \cdot \text{idf}(t, D)$$

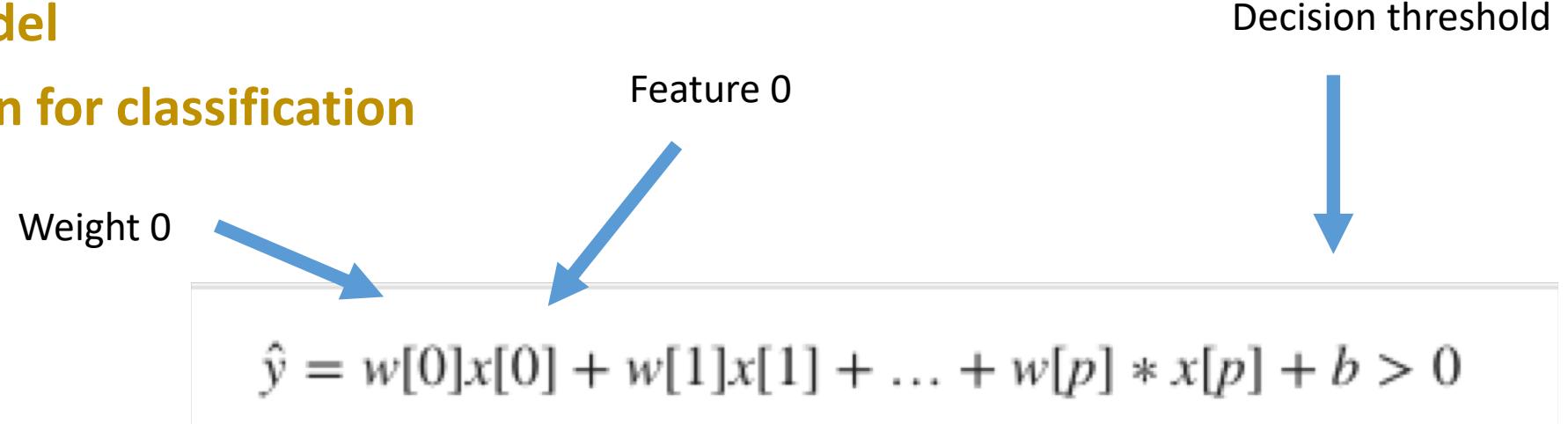


Statistics
Canada Statistique
Canada

Linear Regression

100

- Advantage to start with simple model
- Easy to interpret due to linear relationship with input x and output y
- Fast baseline model
- Logistic regression for classification



Evaluation Metrics

100

Precision/Recall

$y = 1$ in presence of rare class that we want to detect

Actual class		
Predicted class	1	
1	True positive	False positive
0	False negative	True negative

Precision

(Of all patients where we predicted $y = 1$, what fraction actually has cancer?)

$$\frac{\text{True positives}}{\#\text{predicted positive}} = \frac{\text{True positive}}{\text{True pos} + \text{False pos}}$$

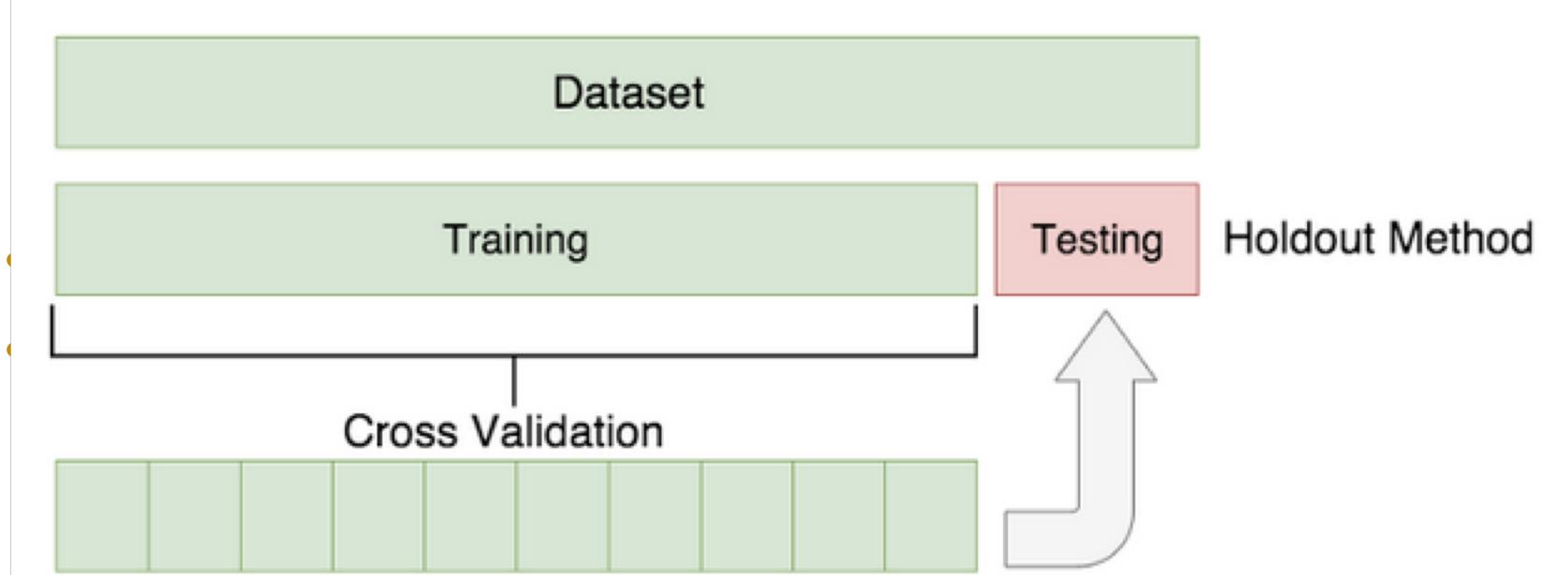
Recall

(Of all patients that actually have cancer, what fraction did we correctly detect as having cancer?)

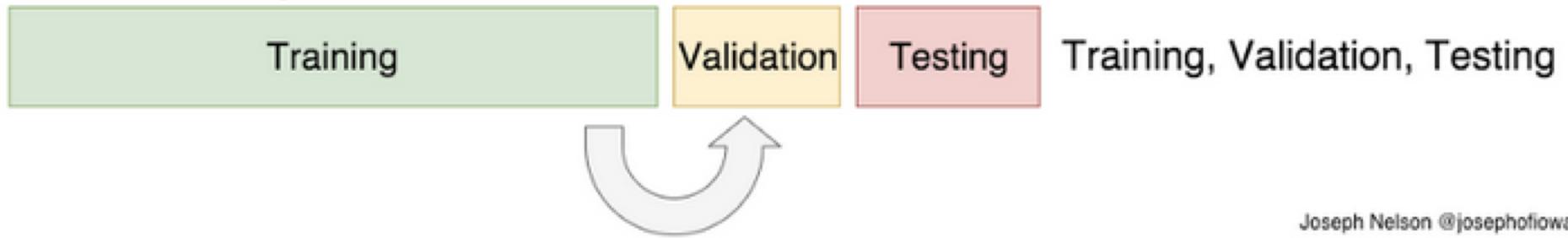
$$\frac{\text{True positives}}{\#\text{actual positives}} = \frac{\text{True positives}}{\text{True pos} + \text{False neg}}$$

Slide by Andrew Ng (Machine Learning on Coursera)

Cross validation



Data Permitting:



Joseph Nelson @josephofiowa

Implementation tips

- Do you have a regression or classification problem?
- Do you have training data and what is the quality?
- Could you leverage open data and algorithms?
- Which features could help to distinguish between the classes?
Features often matter more than algorithms
- Each model needs cross validation



Statistics
Canada Statistique
Canada

Q&A

100

- **Code & data available at
https://github.com/chritter/Talks/tree/master/2018/Data_Challenge_UOttawa**
- **Could you adopt classification in your problem case?**



Statistics
Canada Statistique
Canada



Statistique
Canada Statistics
Canada

Data Science
Accelerator

Enabling you.

Canada