

National Security & Intelligence Applications of Text Analytics

Patrick Lam

**Lead Data Scientist, Thresher
Visiting Fellow, Harvard IQSS**

March 26, 2015

What is text analytics?

"a set of linguistic, statistical, and machine learning techniques that model and structure the information content of textual sources for business intelligence, exploratory data analysis, research, or investigation"

- Wikipedia



What are some common text analysis methods?

Natural Language Processing (NLP)

Enable computers to understand human text input.

Machine Learning

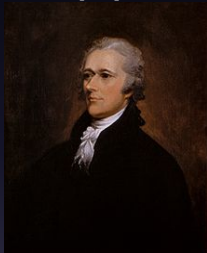
Grouping, classifying, predicting

Applications

Who Authored the Federalist Papers?

(Mosteller and Wallace, 1963)

43 papers



Hamilton

14 papers



Madison

5 papers



Jay

12 papers

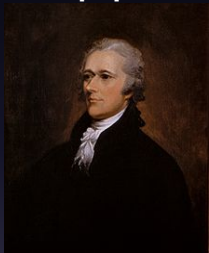


H, M, or J?

**Model the usage of high-frequency
function words for each author.**

also, and, by, of, on, there, ...

43 papers



Hamilton

14 papers



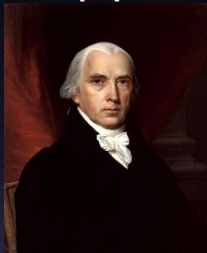
Madison

5 papers



Jay

12 papers



Madison

Authorship attribution using text analysis has a large literature and can be useful for national security & intelligence purposes.

Harvard IQSS

Institute for Quantitative Social Science



iq.harvard.edu

**Four recent projects using text
analysis by Harvard IQSS affiliates
with relevance to national security &
intelligence**

1. Reverse-Engineering Censorship In China

(King, Pan, and Roberts, 2013 and 2014)

What Gets Censored In China?

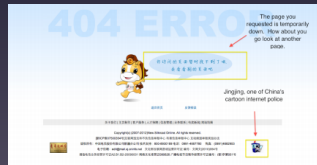
Monitor ~1,400 Chinese social media sites over 6 months across 85 content areas.

Download posts the instant they appear.

Revisit each post later to check if it was censored.

Analyze with new methods of text analysis.

Experiment with writing social media posts.



**Censorship program targets
collective action rather than criticism
of the government.**

2. Jihadi Radicalization of Muslim Clerics

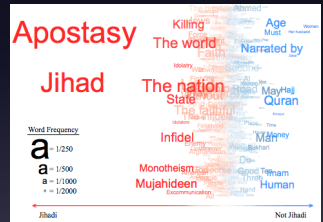
(Nielsen, 2014)

Why Do Some Clerics Preach Jihad While Others Do Not?

Download writings of a sample of Muslim clerics.

Use machine learning methods to score clerics on level of jihad by comparing writings to known Jihadi texts.

Analyze along with other data on clerics.



Clerics with *weak educational networks and connections* often use Jihadi ideology to appeal to lay audiences and further their careers.

3. Predicting Crowd Behavior

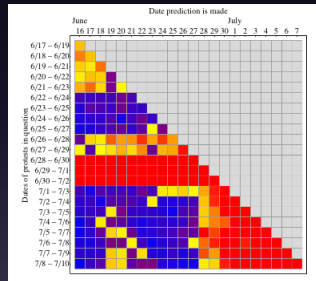
(Kallus, 2014)

Can We Predict Major Protests?

Scan over 300,000 web sources (news, blogs, forums, Twitter) in 7 languages for mentions of past, current, or future events.

Extract type of event, entities involved, and timeframe using NLP methods.

Predict on each day whether a significant protest will occur over the next three days using machine learning methods.



**Massive online public discourse data
has the potential to *predict crowd
behavior* using text analysis methods.**

4. Anti-Americanism in the Middle East

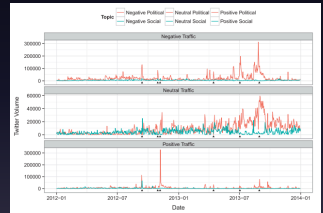
(Jamal, Keohane, Romney, and Tingley, 2015)

Anti-Americanism on Twitter

Download Arabic Twitter posts by keywords.

Consider discourse about US in general and in reaction to specific events.

Use text analysis methods to classify proportion of posts in specified categories.



Anti-Americanism in the Middle East
*is directed toward the impingement of
the US on other countries rather than*
toward American society.

**All applications of text analysis have
a common starting point.**

**How do we *define* our set of texts
that we want to analyze?**

**How do we *retrieve* the relevant texts
from the vast set of texts available?**

**How do we *follow* the relevant
conversations as they evolve?**

Define.
Retrieve.
Follow.

In some cases, the relevant set of texts is well-defined, static, and easily accessible.

Federalist Papers, complete works of Shakespeare

In most cases, we have to constantly search and retrieve the relevant texts, often using *keywords and Boolean searches*.

Twitter, news articles, blogs and forums

DHS Analyst's Desktop Binder

2.13 Key Words & Search Terms

This is a current list of terms that will be used by the NOC when monitoring social media sites to provide situational awareness and establish a common operating picture. As natural or manmade disasters occur, new search terms may be added. The new search terms will not use FBI in searching for relevant mission-related information

DHS & Other Agencies

Department of Homeland Security (DHS)
Federal Emergency Management Agency (FEMA)
Customs and Border Protection (CBP)
Border Patrol
Secret Service (USSS)
National Operations Center (NOC)
Homeland Defense

Immigration Customs Enforcement (ICE)
Agent
Task Force
Central Intelligence Agency (CIA)
Fusion Center
Drug Enforcement Agency (DEA)
Secure Border Initiative (SBI)
Federal Bureau of Investigation (FBI)

Alcohol Tobacco and Firearms (ATF)
U.S. Citizenship and Immigration Services (CIS)
Federal Air Marshall Service (FAMS)
Transportation Security Administration (TSA)
Air Marshall
Federal Aviation Administration (FAA)
National Guard
Red Cross
United Nations (UN)

20

Domestic Security

Assassination
Attack
Domestic security
Dell
Cops
Law enforcement
Authorities
Disaster assistance
Disaster management
DHS (Domestic Threatener Detection Office)
National preparations
Migration
Prevention
Response
Rescue
Dirty bomb
Domestic nuclear detection

Emergency management
Emergency response
State of emergency
Homeland security
Natural disaster evacuation (NDAI)
Natural disaster relief
SWAT
Severing
Lockdown
Bomber (suspect or threat)
Evacuation
Duel
Hanging
Explosion (suspect)
Police
Disaster medical assistance team (DMAT)
Organized crime

HAZMAT & Nuclear

Hazard
Nuclear
Chemical spill
Suspicious package/device
Toxic
National Laboratory
Nuclear facility
Nuclear threat
Chemical
Plasma

Lock
Biological infection (or agent)
Chemical
Chemical agent
Biological
Epithelial
Therapeutic
Biological material (hazard)
Biological agent

Gas
Spillover
Asbestos
Bioterror agent
Chemical agent
Explosion
Nuclear
Nerve agent
Radio

Sarin
Quarantine
HINI
Vaccine

Tsuri
New Virus
Epidemic

World Health Organization (WHO) (and components)
Viral Hemorrhagic Fever
E. Coli

Infrastructure Security

Infrastructure security
Airport
CBR (Critical Infrastructure & Key Resources)
AMTRAK
College
Computer infrastructure
Communications
Infrastructure
Telecommunications
Critical infrastructure
National infrastructure
Metro
WMATA

Airplane (and derivative)
Chemical line
Subway
BART
MARTA
Port Authority
NBC (National Business Surveillance Integration Center)
Transportation security
Grid
Power
Smart
Body scanner

Electric
Failure or outage
Black out
Brown out
Port
Dock
Bridge
Caterpillar
Delays
Service disruption
Power lines

Southwest Border Violence

Drug cartel
Violence
San Diego
Cuidad Juarez
Nogales
Sonora
Cocaine
Marijuana
Jaram
MS13 or MS-13
Bomber
Mexican army
Cartel
Cuidad de Lofio
Gulf Cartel
Sinakos
Tijuana
Tijuana
Tijuana
Yuma
Tucson
Decapitated
U.S. Consulate
Consular
El Paso
Execution

Gunfight
Trafficking
Kidnap
Calderon
Reyes
Bust
Tamaulipas
Muh Lab
Drug trade
Illegal immigrants
Smuggling (smugglers)
Mazatlan
Michoacan
Guzman
Aguilano Felix
Beltraz-Leyva
Barro Arreza
Artistic Assassins
Mexico
New Federation

Terrorism

Terrorism
AQ/Al Qaeda (all spellings)
Attack
Iraq
Afghanistan
Iran
Pakistan
Ago
Environmental terrorist
Iran terrorism
Conventional weapon
Target
Weapons grade
Dirty bomb
Hiroshima
Nuclear
Taliban
Biological weapon
Weapons cache
Arrested inmate
Improvised explosive device

Suspicious substance
AQAP (Al Qaeda Arabian Peninsula)
AQIM (Al Qaeda in the Islamic Maghreb)
TTP (Tehrik-e-Taliban Pakistan)
Yemen
Prison
Extremism
Somalia
Nigeria
Radicals
Al-Shabaab
Home grown
Plot
Nationalist
Recruitment
Fundamentalism
Islamist

Weather/Disaster/Emergency

Emergency
Ice
Hurricane
Tornado
Help
Twister
Tsunami
Earthquake
Bust
Flood
Tsunami Warning Center
Magnitude
Avalanche
Typhoon
Shelter-in-place
Terrorism
Disaster
Snow
Bilgund
Ski
Cyber Security

Mad slide or Mudslide
Erosion
Power outage
Brewer out
Warning
Watch
Lightning
Aid
Relief
Chairs
Intense
Barani
Emergency Broadcast System

Cyber security
Botnet
DDOS (dedicated denial of service)
Denial of service
Malware
Virus
Trojan
Keylogger
Cyber Command

Hacker
China
Conficker
Worms
Scammers
Social media

**Our analyses are only as good as our
keywords!**

**What is our current most commonly
used technology for defining the
relevant keywords to retrieve our
texts?**



Example: Think of keywords you would use to follow the Twitter conversation around the Boston Marathon Bombings.

Keywords about the event

#bostonbombings, explosion, terrorism, attack, ...

Keywords about the suspects

suspect, tsarnaev, dzhokhar, tamerlan, ...

Keywords about the victims

innocent, victim, collier, ...

Keywords about the reaction

tragedy, prayers, #prayforboston, ...

Keywords about the politics

obama, #tcot, #benghazi, ...

**We ran a similar experiment with 43
Harvard undergrads.**

#bostonstrong explosion marathon
brothers mondaysearch
cooker #marathonmonday
response #bostonmarathon strong
official overcome wheelchair victim
boat blue running hurt fund child safe race terrorist
#wewillneverforget what injured attack
usa intense shoes
horrrifying deaths news
shootout arrest where
swat collier guns people
restrict dead madness
fire hiding blast love runner
harvard pride boylston officer
sad drive allston again hero emt russian
diversion sean army culp brave rip rescue
legs april run find shocking rice mbta help
hospital building hunt ambulance zubeidat support radio
god qaeda dartmouth mgh lawyer somerville
pray stand #staysafe horror feet pressure lockdown
memorial detective unexpected family debris save tragedy
scary weather children leaves drunk danger
terrorism runners shutdown shooting tsarnaev
survivors chase line watertown
backpack finish

**59% of the words were suggested by
only 1 out of 43 undergrads.**

**Median number of words per
respondent was 7.**

**Humans are good at recalling a small
list of good keywords and
recognizing a good keyword when
they see it.**

Humans are bad at recalling a long list of keywords that capture different ways of representing a concept.

Some Existing Options for Automated Keyword Discovery

1. Mine search queries

Google Adwords

2. Thesaurus methods

reference books, WordNet

3. Co-occurrence methods

Enter Thresher.

(Based on King, Lam, and Roberts, 2014)

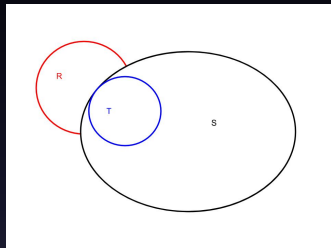
**Thresher keeps humans in the loop
and helps them find more and better
words faster.**

The Thresher Algorithm

Reference set R: texts about concept of interest

Search set S: broad set of texts

Target set T: texts in S about the same concept as in R



Goal: Estimate T and find keywords that define T.

Thresher Applications

The Boston Marathon Bombings on Twitter

R: #bostonbombings

S: boston

**Thresher separates relevant words
from irrelevant words.**

Target set words

suspect, bomb, police,
people, fbi, tsarnaev, arrest,
die, terror, attack, kill,
obama, custody, prayer,
dzhokhar, god, false,
#prayforboston, #tcot,
#bostonmarathon, picture,
identify, russia, #watertown,
tamelan, islam, jihad,...

Non-target set words

game, red sox, bruins,
celtics, back, tonight, fan,
#mlb, night, chicago, new
york, garnett, fenway, rondo,
#job, playoff, yankees,
blackhawks, stanley, pizza,
#nhl, draft,...

The Bo Xilai Scandal on Chinese blogs and forums

R: 薄熙来 (Bo Xilai)

S: 重庆 (Chongqing)

王立军	Wang Lijun
政治	government
事件	event (euphemism for the scandal)
打黑	strike corruption
犯罪	commit a crime
民主	democracy
权力	power
文革	Cultural Revolution
领导	leader
改革	reform
群众	the masses
中央中共	Central Communist Party
社会主义	socialism
唱红	sing red songs
黑社会	black society
干部	cadre
路线	party line

Writings About Suicide Bombings in Arabic

R: الاستشهادية عمليات (martyrdom operations)
from "Haqibatu'l-Mujahid" (A Mujahid's Bookbag)

S: "Pulpit of Tawhid and Jihad"
A Jihadist web library

العدو enemy

قتل killing

والنكاية to vex or spite ("vex the infidels")

يَعْلَمُهُمْ teach them

الْخَيْلَ steed

وَأَعَدُّوا fight

تُظْلَمُونَ wronged

ترهبون terrify

الغلام boy

(the story of the boy and the king, relevant to jihadis)

"And prepare against them whatever you are able of power and of steeds of war by which you may terrify the enemy of Allah and your enemy and others besides them whom you do not know [but] whom Allah knows. And whatever you spend in the cause of Allah will be fully repaid to you, and you will not be wronged."

- Quran 8:60

**Can we use Thresher to *define*
different Arabic dialects by keywords
and *retrieve* texts from them?**

Egyptian

Gulf

Levantine

MSA

Translation

ايه
معرفش
عايز
عاوز
عايزين

وش
مادري
يبي
يبغى
يبون

شو
ما بعرف
بدي
بدك
بدكم

ماذا
لا يعرف

**What
Don't know
To want/I want
To want/You want
To want/You want (pl.)**

Conversations can evolve quickly in response to certain actors and we need to be able to *follow* them.

Evading Censors in China

自由 **Freedom**

目田 **Eye field**

(homograph)

和谐 **Harmonious [Society]**
河蟹 **River crab**

(homophone)

Evading Authorities and the Distribution of Child Pornography

PatrickLam.org
ThresherVentures.com
iq.harvard.edu