

R News Digest

Devis Shehu (ds1902)



- [illegible]

6

Russia invaded Ukraine on 24 February 2022, marking a steep escalation of the Russo-Ukrainian War, which began in 2014 following the Ukrainian Revolution of Dignity. The invasion has caused Europe's largest refugee crisis since World War II,^{[17][18]} with more than 5.5 million Ukrainians leaving the country^[19] and a quarter of the population displaced.^{[20][21]}

At the start of the war in 2014 Russia annexed Crimea, and Russian-backed separatists seized part of the south-eastern Donbas region of Ukraine, sparking a regional war there.^{[202][203]} In 2021, Russia began a large military build-up along its border with Ukraine, amassing up to 190,000 troops along their equipment. In a televised address shortly before the invasion, Russian president Vladimir Putin expressed *irrefragable views*,^[204] questioned Ukraine's right to statehood,^{[202][205]} and falsely^[206] accused Ukraine of being organized by neo-Nazis who persecute the ethnic Russian minority.^[207] Putin also said the North Atlantic Treaty Organization (NATO) constitutes a threat to Russia's national security by expanding eastward since the end of the second world war.^[208] He also accused the United States of provoking Russian opposition and of attacking Ukraine from every angle in the strategic alliance.^[209] The United States and others accused Russia of planning to attack or invade Ukraine, which Russian officials repeatedly denied as late as 23 February 2022.^[210]

On 21 February 2022, **Russia** controlled the **Donetsk People's Republic** and the **Luhansk People's Republic**, two self-proclaimed states in Donbas organised by pro-Russian separatists.^[36] The following day, the **Federation Council of Russia** authorised the use of military force to invade, and Russian troops overtly entered both territories.^[36] The invasion began on the morning of 24 February,^[37] when Putin announced a "special military operation" to "demilitarise and denazify" Ukraine.^[38] Minutes later, missiles and airstrikes hit across Ukraine, including the capital **Kyiv**, shortly followed by a large ground invasion from multiple directions.^{[40a][41]} In response, Ukrainian President **Volodymyr Zelenskyy** issued **martial law** and **general mobilisation** of all male Ukrainian citizens for between the ages of 18 and 60, who were banned from leaving the country.^{[42a][43]}

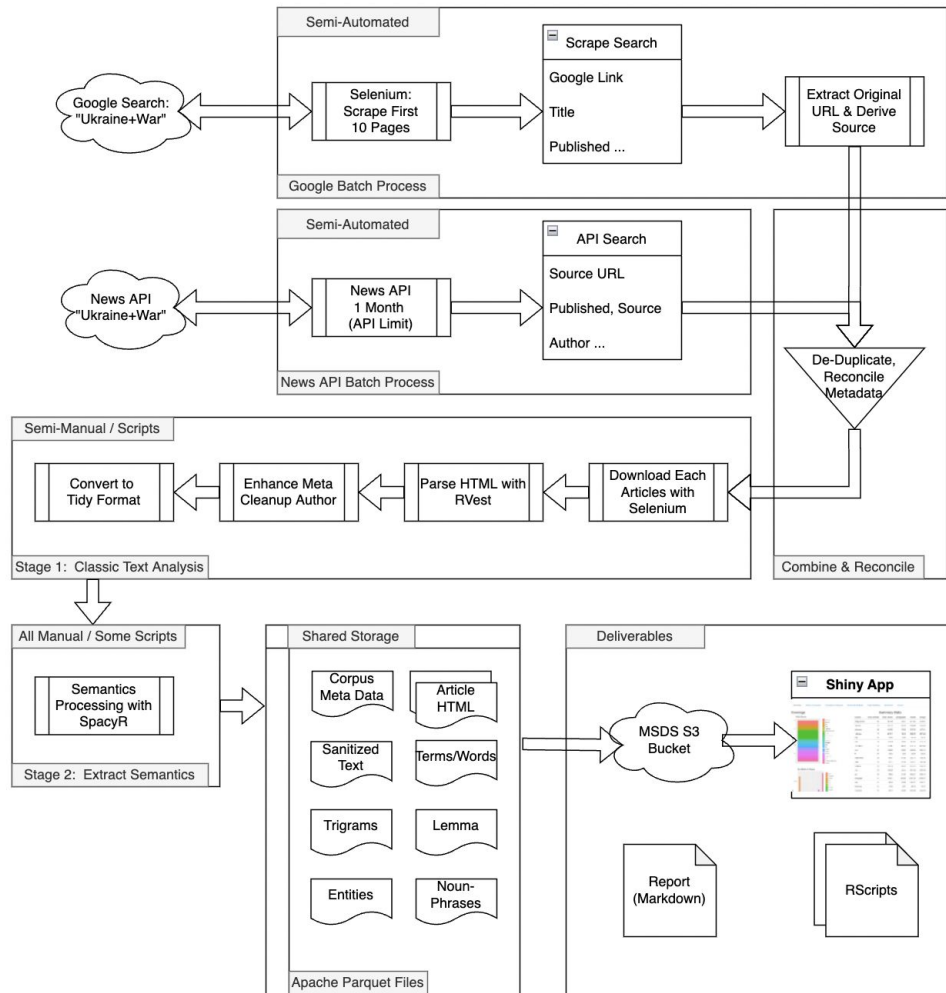
At the start of the invasion on 24 February, the northern front was launched out of Belarus and targeting Kyiv with the northeastern front launched at the city of Kharkiv; the southeastern front was conducted as two separate spearhead fronts including a southern front (originating in Crimea) and a separate probeable southeastern front [launched at the cities of Luhansk and Donetsk]^{[34][45]}. On 1 April, the Russian military had announced that all its troops and divisions deployed in southeastern Ukraine would be united under General Aleksandr Dvornikov, who was placed in charge of combined military operations, including the redeployed probeable fronts originally assigned to the northern front and the north-eastern front which were subsequently withdrawn and reassigned to the second phase to the south.^[46] By 17 April, progress on the southeastern front appeared to be impeded by residual troops continuing to hold-out in abandoned locations in Mariupol, Sloviansk, Ilovaisk, and Kostiantynivka. The Ukrainian defense force referred to as an attack group from the east front extending from Kharkiv to Donetsk and Luhansk, with simultaneous missile attacks again directed at Kyiv in the north and Lviv in western Ukraine^[47].

The invasion has been widely condemned internationally as an act of aggression [49][50]. The United Nations General



General Flow

- Identify news providers: Google search and [News API](#) were most straightforward (discarded Guardian, NY Times API due to complexity - value tradeoff)
- Run Google searches via Selenium in batches
 - Deal with Google and other sources' throttling
- Call News API iteratively due to limits on free account (1 month of free data) to accumulate data
- Collect or parse all available metadata (dates published, authors, title, url, etc.) over 2 months
- Store metadata in a separate data frame (provider)
- Download files to avoid issues with broken links, re-running process,etc.
- Parse HTML and cleanup auth, etc.
- Sanitize (stop words, etc.) pre-calculate tokens, trigrams and calculate stats for performance
- Convert to Tidy Format and extract semantics (mostly manual)
- Shuffle files to S3 so Shiny App is detached from backend data collection and processing
- Lots of try/catch, error handling for poorly formatted data, Selenium erroring out, throttling, etc.
- Robustness at scale is really hard and time consuming and did not meet the goal!

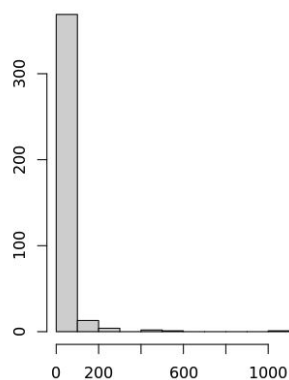


Summary Statistics

- **10,390** Articles
 - 6,996 from Google News including 1,009 attribute to Google but actually made up of random feeds, eg. youtube videos, unknown sources
 - 3,395 from News API
- Timeframe: Jan 31st to April 16th : paused data collection and cleaning
- 679 unique sources but highly skewed with a relatively small # of sources contributing most content (long tail possibly an artifact of data collection)
- News API: 533 secondary sources, attributed to other publications such as AP (262), Reuters(75), VOA (Voice of America) News, ABC, France, New York Times, Bloomberg
- Surprised that NewsAPI, an expensive service, starting at \$450 month, has really unclean metadata (especially authors)
- Hundreds of articles were attributed to AP, other news organizations which makes sense but having the author would have enhanced the analysis

Summary Stats

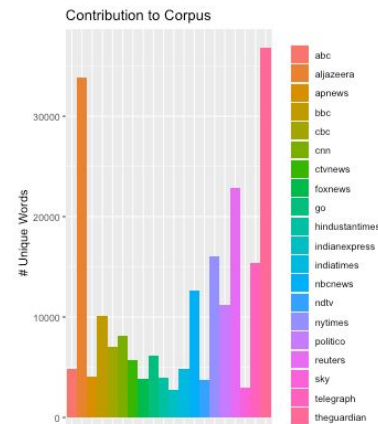
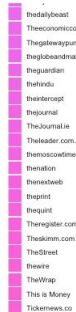
source	total_articles	total_words	paragraphs	words	unique
aljazeera	514	350722	51.80	682.34	372.78
theguardian	500	311412	27.18	622.82	356.92
reuters	408	194237	34.49	476.07	331.19
ndtv	280	71234	20.61	254.41	179.12
wsj	242	23738	10.88	98.09	70.47
cnr	199	95359	38.17	479.19	297.91
npr	196	60100	20.24	306.63	195.72
bbc	171	89608	55.36	524.02	356.10
Yahoo Entertainment	170	161822	57.69	951.89	623.66
Daily Mail	158	125653	73.33	795.27	446.21
indiatimes	151	105693	70.01	699.95	297.17
nypost	143	76473	30.59	534.78	299.19
hindustantimes	126	90834	50.41	720.90	427.82
nbcnews	122	93631	55.02	767.47	467.01
google	111	99474	67.65	896.16	475.89
Freerepublic.com	109	46459	37.14	426.23	296.54
apnews	103	52117	29.03	505.99	337.26
Independent	103	47618	55.07	462.31	276.09
abc	101	89355	89.34	884.70	543.14
newsweek	94	168648	122.80	1794.13	767.54



AP	58
Feed / Feeder	57
Facebook (BBC)	53
AP Wire	25
Post Editorial Board	20
ABC Newshour	19
Socialists & Democrats in EU Parliament	18
Alex Millson	14
Staff	13
Angela Dennis	12
Caroline mimbs nyce	11
Press operations	11
SA Transcripts	11
Eustance Huang	10
LW	9
Star Tribune Staff	9
AFP	8
Bob Brigham	8
Deutsche Welle	8

-

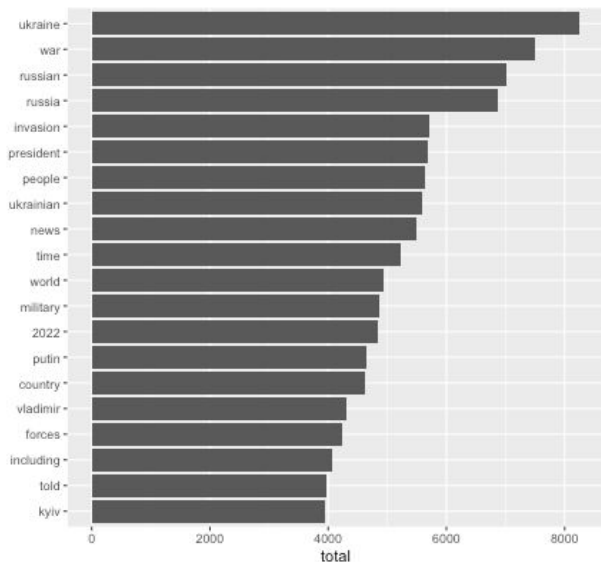
CBC News	marketWatch
CBS News	MarketWatch
Chical.com	middleeasteye
cnn.com	military
Copa.org	Minneapolis Star Tribune
rgm	Mirror Online
charlotteobserver	moneycontrol
Chicago Tribune	newspapers
chicagotributes.net	National Post
Christian.com	Melby Fuel Australia
CleanTechnica	enr.com
CNA	MSNBC
cnn	National Observer
CNBC	National Post
cnn.com	NationalSecurity.com
CNET	nature
com	NBC News
CNN	ribochicago
Comcast	Not asic
Commonwealth.ca	npr
cp24	NYT
CP24 Toronto's Breaking News	NYT News





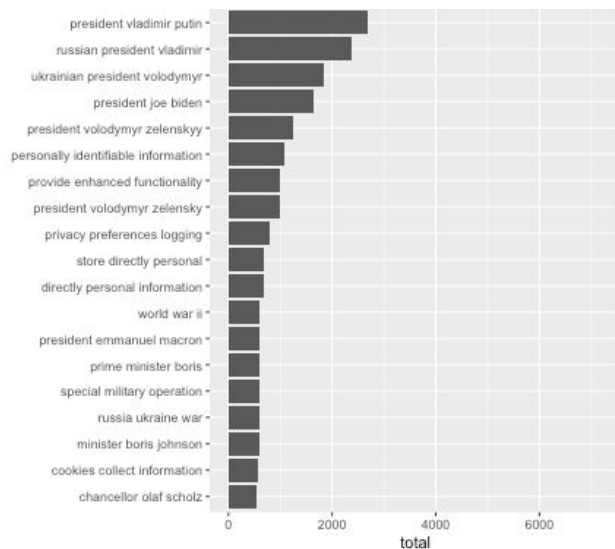
Terms & Concepts

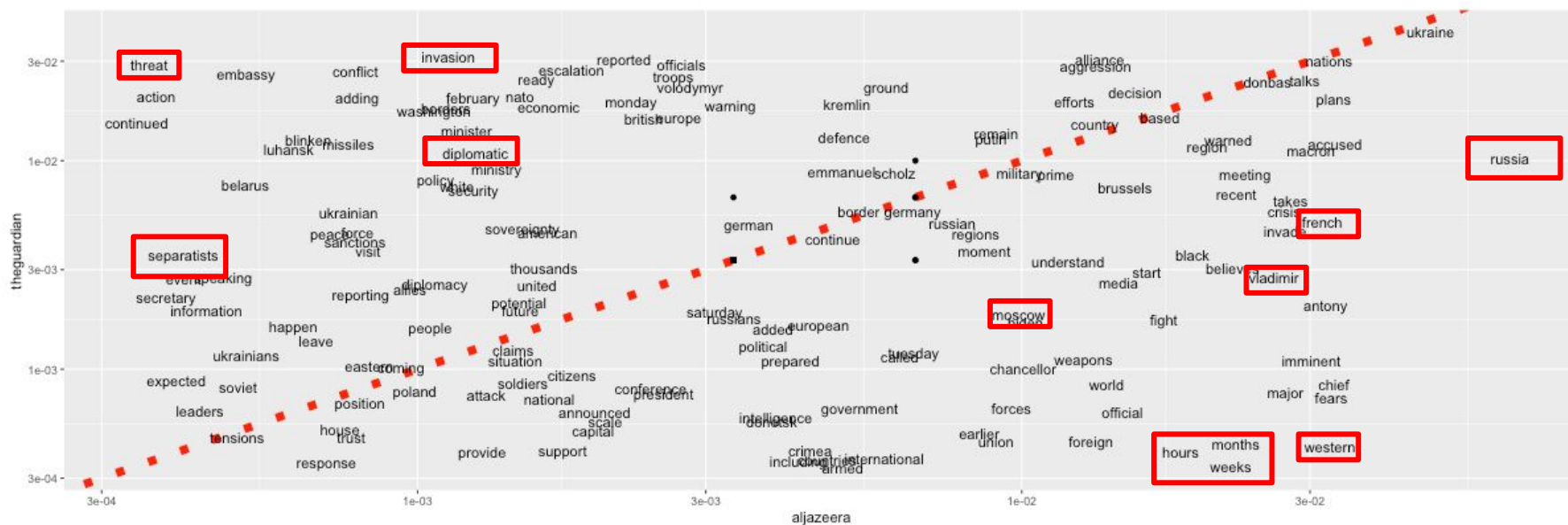
Top Words



- Overall most articles have a high % of unique words which means the proportion of top words is a low % of the total set of words in the corpus
- Sparsity of words is an issues when computing frequency, etc.
- Terms and trigrams are as expected, mostly bubbling up the key players (Putin, Boris Johnson, Biden) and concept
- “World War II”, which seem to contribute to fear sentiments pops-up relatively high on the list
- Some data oddities such as “personally identifiable information” and “cookies collect information” due to “boilerplate” web content. not filtered in parsing (improvement)

Top Trigrams







Comparative Analysis

- From the word comparison and unique terms, some indication that the content of the **The Guardian** differs from **Al Jazeera**
- Both are highly credible and well regarded as sources though The Guardian has a reputation as being somewhat left-leaning and biased toward Britain, and its role and leadership in the conflict
- Some evidences of emphasis on high level diplomacy and threat posed by Ukraine conflict to European, Britain's collective security and action
- Al Jazeera seems to have a more balanced view of the conflict that emphasis broader context; i.e. West's ambiguous relationship with Russia, potentially culpability in the runup to the conflict
- Al Jazeera seems more empirical, objective and focused on the dates and facts but also at times seems more tolerant of Russia
- Unrelated commercial terms (ads) and browser/cookie (boilerplate notice) - small glitch with "belarus"
- Vladimir / Putin is conspicuously prevalent on Al Jazeera's along with French diplomatic efforts and Germany's role (especially early on hesitation to break with Russia and send weapons)

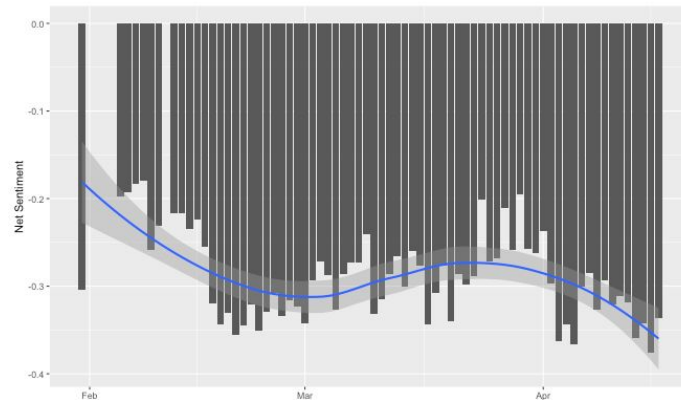
word	word
guardian	jazeera
ability	administration
abruptly	advertisers
access	advertising
adviser	aggregated
analysis	alliance
appeared	anonymized
artillery	assault
attempt	assistance
believed	belarus
benefit	black
biggest	block
billionaire	borders
boris	browser
britain	build
casualties	business
chancellor	ceasefire
changing	central
claimed	change
claims	china
clock	citizens
collapse	collection
commercial	concerns
communities	confirmed
	consent

Guardian Only

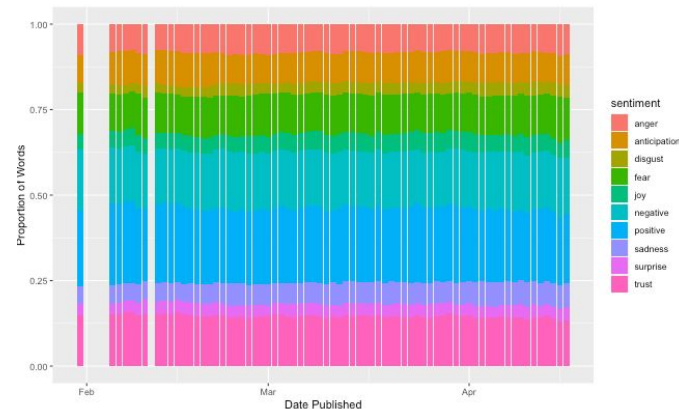
Al Jazeera Only

Sentiment Analysis

- Sentiment analysis using **bing** and **nrc** to get both a positive/negative and some categorization
- As expected sentiment started out strongly negative and declines precipitously over time as the reality of the war settles in as well as realization that it's not likely to end
- Anticipation was highest at the beginning of the conflict both in terms of if invasion was happening and later on the (misplaced) hope for a peaceful (short term upswing)
 - Less pronounced when looking at 2 ½ months vs 2 week sample
- Expected anger and disgust to increase relative to other emotions as evidence of war crimes surfaced, possibly the sentiment is diluted by sheer number of sources?
- Sadness seems to be increasing but not as much as expected
- Some evidence in the data of fear increasing marginally as terms such as WW3 and nuclear conflict started occurring more frequently
- Sentiment analysis might have been more insightful if run on top 20 news sources and with pruning (separate positive and negative)
- Can we measure fatigues (apathy) somehow?



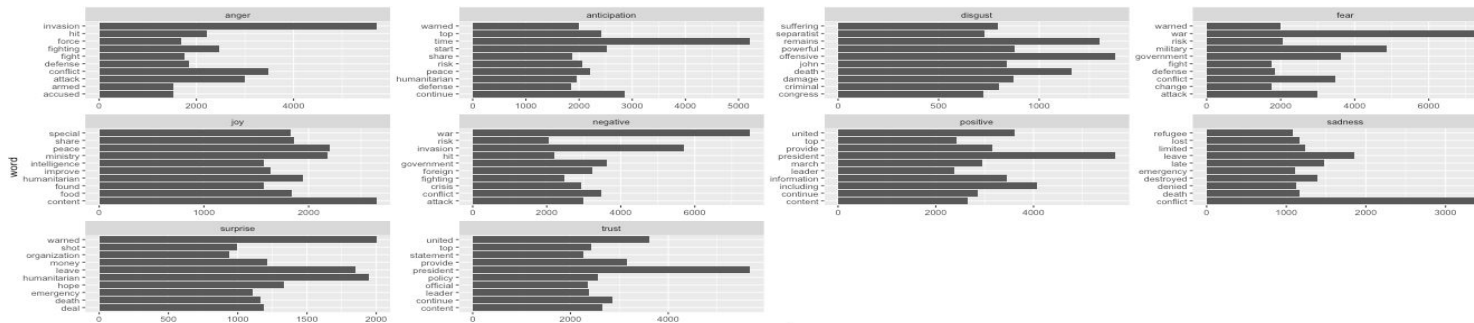
Categories of Sentiments



Sentiment Analysis - Emotions & Terms

- Given the content, the associated terms aren't too surprising: i.e. invasion, force, fighting associates with anger/fear
- Likewise for refugee, death, emergency with sadness
- Also as expected, temporal (start, continue, time) associate with anticipation: humanitarian (corridors) and peace (conference) do bubble up which tracks with trends
- Some associations are obviously irrelevant ("intelligence" with joy?) - limits of Sentiment Analysis?
- News coverage emphasis the "unexpected" (surprise) of the breadth and depth of humanitarian assistance - is this unexpected or news accentuating the positives and negatives (extremes)?
- Hard to gauge trust: is there significant trust in policy, leaders or even president?
- Sentiment analysis is somewhat limited expect when it bubbles up an unexpected result

Top Terms by Sentiment





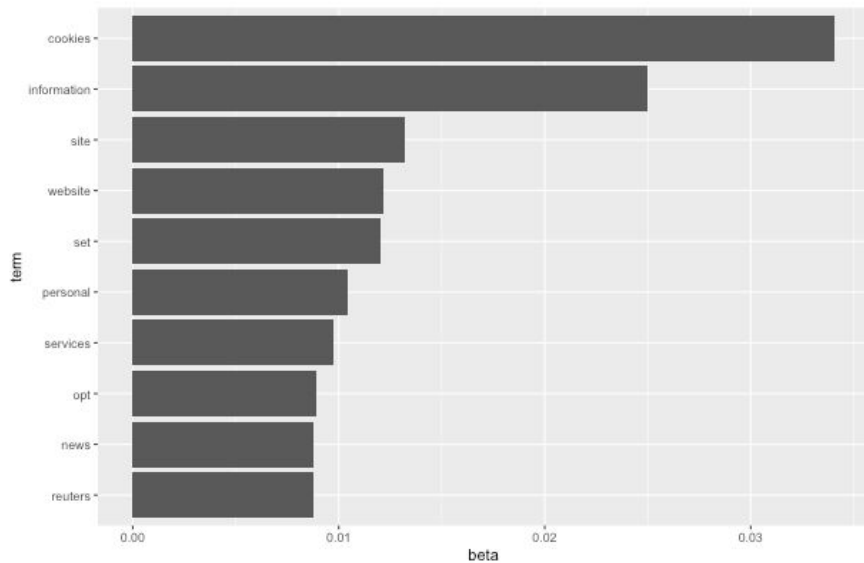
Topic Modeling & Linear Discriminant Analysis

- Ran LDA with default parameters and K=8 and Fisher's method: derived experimentally (trial and error)
 - Needed to convert tidy text with terms sparse matrix
 - `dfMatrix <- dfNewsTidyTerms %>% cast_dtm(document, term, count)`
- Goal was to see if LDA could extract the dimensions of the conflict (political, military, humanitarian, financial (natural resources), etc.
 - Ideally combine this with sources to tag and differentiate sources according to dimension
- Need to explore [coherence score](#) to determine the right number of K and it would help to narrow down the sources to the most trustworthy, highest regarded (see earlier comments)
- Potentially use topic modeling to remove irrelevant search terms, trigrams (advertisement, browser/cookie, etc.)
- Potentially categorize and subdivide the Corpus into subsets and analyze each (drill-down) - automation?
- Some duplication of common words (Ukraine, Russia) is difficult to escape as the words have different context (interesting problem)

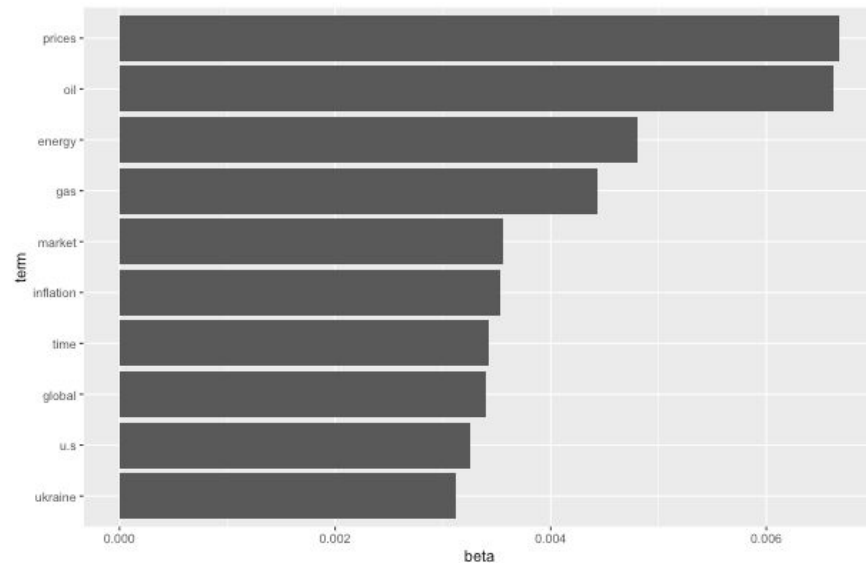


Topic Modeling Results - Part 1

Topic 1



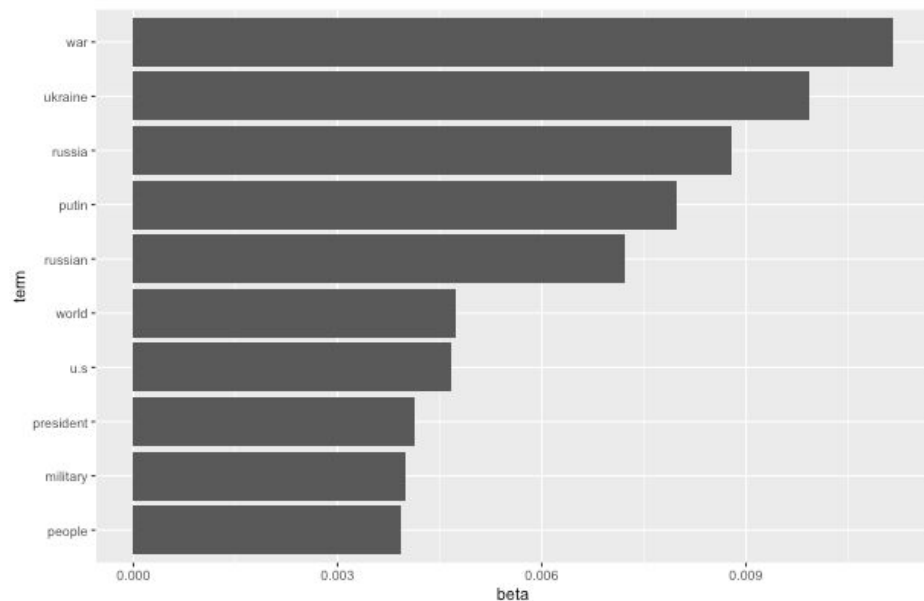
Topic 2



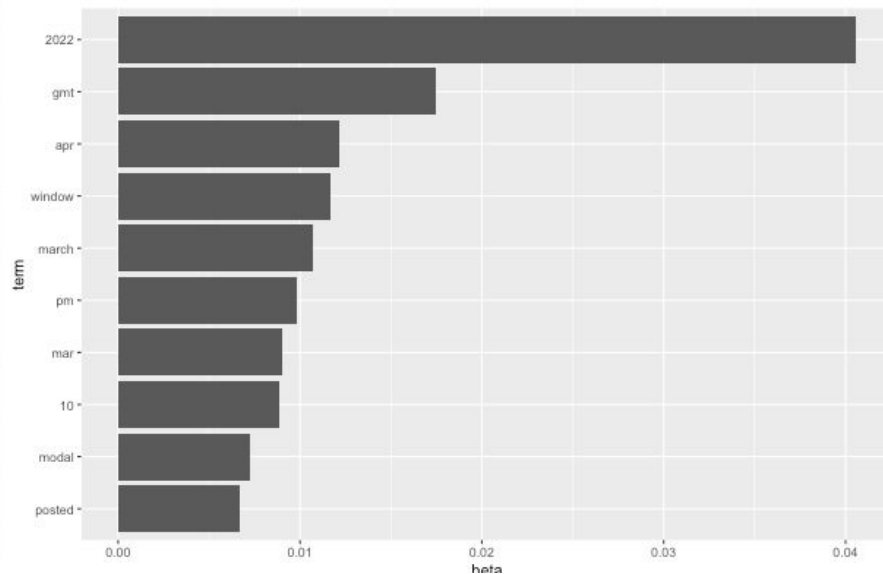


Topic Modeling Results - Part 2

Topic 3



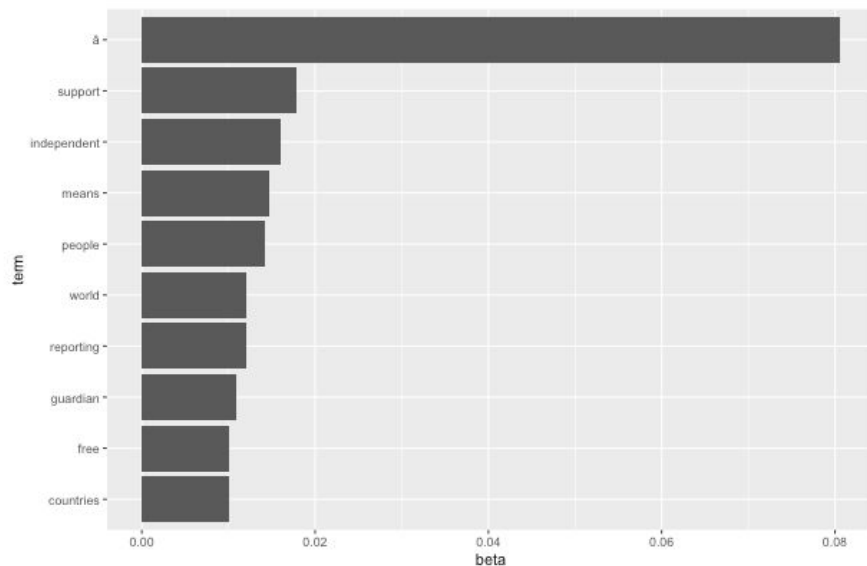
Topic 4



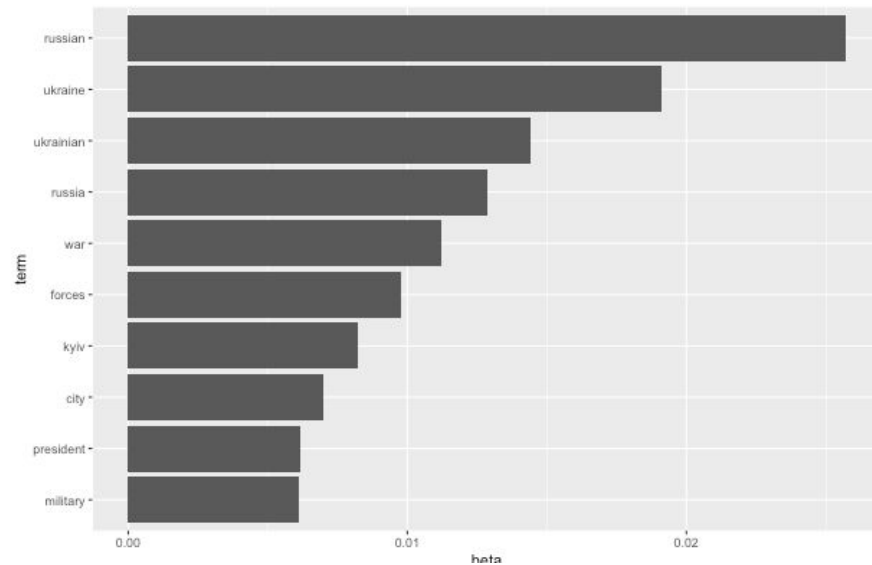


Topic Modeling Results - Part 3

Topic 5

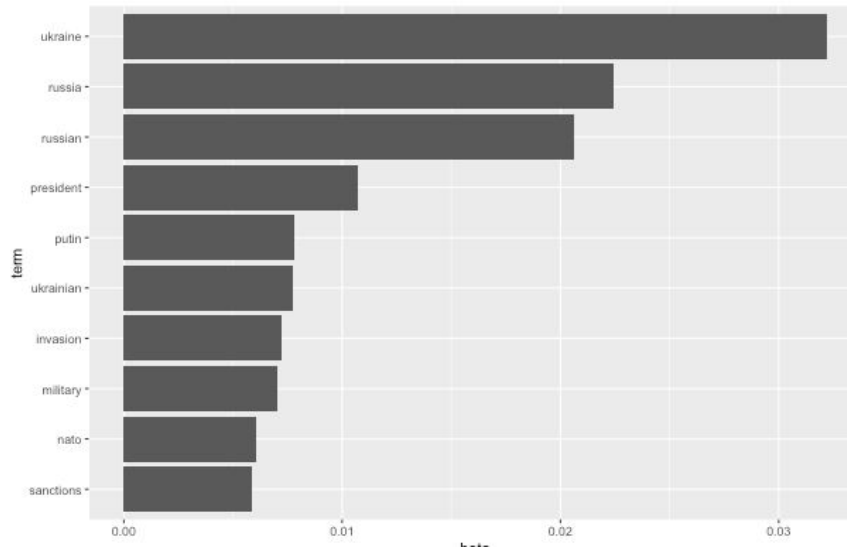


Topic 6

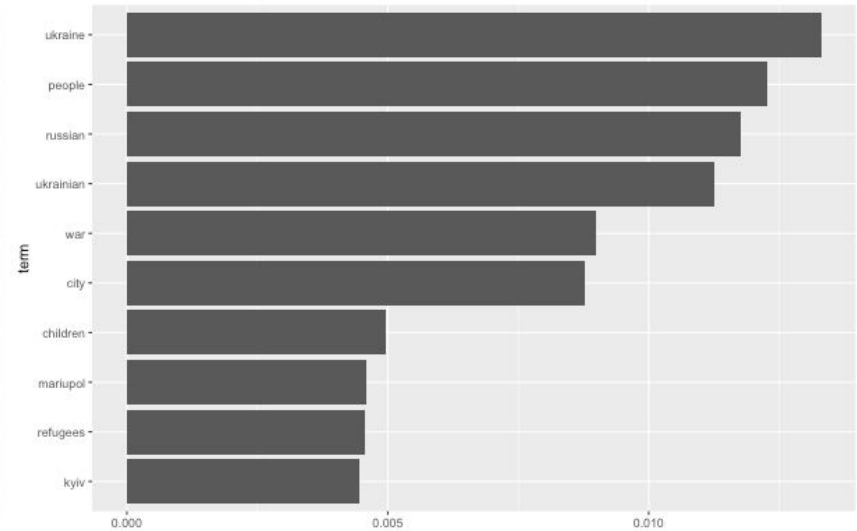


Topic Modeling Results - Part 4

Topic 7



Topic 8



Shiny App, Sources, Misc.

- Build a simple visualization application with some filtering controls (source, timeline)
- Due to issues with Google search throttling & blocking requests and NewsAPI rate limiting, processed offline and uploaded files to S3
- App loads files from S3 and builds visualizations slowly: from <10 seconds for “small” sample set of 1K articles and up to 3-5 minutes for 10K articles and 8 topic LDA (M1 Mac)
- Sources exposed via Shiny UI, can access original articles, if still available; spot checked some at random and link back to stats: another improvement?

Summary Terms & Concepts Comparative Analysis Sentiment Analysis Topic Modeling Semantics Source

Show 25 entries Search: Ukraine crisis: Why is

	id	published	source	authors	link
1	c6f5fb44-3cd5-4065-93b1-064ed8aff0d7	2022-01-31	aljazeera		Ukraine crisis: Why is Germany out of step with the US, Europe?
457	cea88d62-106f-423f-a230-42d3afe7aabe	2022-02-13	republicworld		Timeline of Russia-Ukraine crisis: Why did conflict start? What ...
573	418154b6-7861-4e45-8b95-473213e98ee4	2022-02-14	republicworld		Timeline of Russia-Ukraine crisis: Why did conflict start? What ...
1363	b3950d2f-fbfe-4b01-a12e-62a4b6d9833e	2022-02-22	aljazeera		Ukraine crisis: Why is Macron taking on the role of mediator?
1378	7d78dc04-6714-4534-8e90-61d09db64e0	2022-02-22	google		Russia - Ukraine Crisis: As West race to prevent war, why ...
1442	0d50e0c6-1549-4138-b23b-af40d3722efe	2022-02-23	aljazeera		Ukraine crisis: Why is Macron taking on the role of mediator?

Showing 1 to 6 of 6 entries (filtered from 10,390 total entries)

Previous 1 Next

News | Russia-Ukraine war

Ukraine crisis: Why is Germany out of step with the US, Europe?

Germany has refused to send weapons to Ukraine and adopted a relatively softer line on Russia.





Next Steps & Improvements

- SpacyR: semantics will be included in the final report and deliverable, some performance issues running on local machine with > 100 articles; hard to compare with “classic” text mining
 - Possible 1st application of SparklyR and Databricks
- Project report: summarize results and learnings in Markdown file

Future

- Automate data collection, especially scraping, by wrapping R script in Docker container and deploying to AWS Fargate (should help with URL throttling, rate limiting)
 - Look at proxy and other services as an alternative
- Automate the processing and allow for user selectable search argument
- Parse “byline” from Google search articles to extract more authors
- Topic modeling:
 - Explore coherence and algorithms for identifying K automatically
 - Potentially as a post-processing filter - feedback into pipeline
 - Subcategories of topics (multiple corpus) - multi-resolution analysis
- Analyze sentiment across sources and/or filter on trustworthy sources
- Derive a “trust” and “uniqueness” scores for each source to determine if the information is reliable and/or unique (adds to the bodies of data)
- Shiny App performance improvements: Shiny caching of reactive and non-reactive data