

Quantitative Textual Analysis in Sociology I

The migrant stories analysis

What we can learn about migration from migrant stories?



Homework

Recent approaches to quantitative textual analysis in sociology

How to overcome the problem of coding/dictionaries?

- exploratory research questions rather than testing hypotheses
- inductively constructed classifications rather than pre-established ones
- observing relationship between textual and non-textual data

Sociology and Quantitative Textual Analysis

Sociology works with quantitative data from its beginning.

Data are primarily non-textual, such as socio-demographic data, opinions, attitudes, and behaviour.

If textual data enter into analysis, they are coded and codes further analyzed. For example, answers to open-ended questions in a survey questionnaire (What do you feel is the most important issue facing the world today?) are coded according to the problems mentioned (climate change, terrorism, inequality etc.)

Coding textual data: problems of validity and (inter-coder) reliability.

Data processing

Pre-processing: cleaning, tokenization, stemming, lemmatization, stop words removal, etc.

Feature extraction: TF-IDF, word embeddings, etc.

Modeling: Naive Bayes, Logistic Regression, SVM, etc.

Evaluation: accuracy, precision, recall, F1 score, etc.

Analysis

Our analysis will be exploratory.

What we know about narrators: gender, the original and current country of residence, and GDP per capita for countries. We then constructed following independent variables:

female - male migrant
 immigrant - homecomer
 intracontinental - intercontinental migrant
 higher - equal - lower GDP migrant

A dependent variable - the migrant's story

We will try to search if there is any relationship between independent and dependent variables. For example, what kind of impact, if any, gender of narrators has on their stories.

The analysis was carried out in R, a free software environment for statistical computing and graphics (<https://www.r-project.org/>).

Data

- 1018 short biographic narratives of migrants published on www.migrant.org.uk
- the stories have been adapted for publication by groups or organizations supporting the story and eventually selected by OMA, first UK organization providing help for migrants
- it is not a representative nor ordered sample of migrant experiences over the world
- it is a very heterogeneous sample of migrants' stories



Sociology and Quantitative Textual Analysis

Sociology works with quantitative data from its beginning.

Data are primarily non-textual, such as socio-demographic data, opinions, attitudes, and behaviour.

If textual data enter into analysis, they are coded and codes further analyzed. For example, answers to open-ended questions in a survey questionnaire (What do you feel is the most important issue facing the world today?) are coded according to the problems mentioned (climate change, terrorism, inequality etc.)

Coding textual data: problems of validity and (inter-coder) reliability.

The General Inquirer (1966)

- the first form of computer-aided content analysis
- a computer program capable to search for recurrent patterns within textual data
- based on universal and custom dictionaries
- problems with coding irregularities, generalization, input, expansion, continuity

- human coders are employed to solve the problem of validity
- involves a new problem of bias, learning and fatigue, not one of inter-coder reliability or agreement
- various measures of inter-coder agreement

Activity

The General Inquirer (1966)

- the first form of computer-aided content analysis
- a computer program capable to search for recurrent patterns within textual data
- based on universal and custom dictionaries
- problems with coding/dictionaries: lemmatization, tropes, translation, contextuality

Certainty: Language indicating resoluteness, inflexibility, completeness, and a tendency to speak ex cathedra

Activity: Language featuring movement, change, the implementation of ideas and the avoidance of inertia

Optimism: Language endorsing or highlighting the positive entailments of some person, group, concept, or event

Realism: Language describing tangible, immediate, recognizable matters that affect people's everyday lives

Commonality: Language highlighting the agreed-upon values of a group and rejecting idiosyncratic modes of engagement

Inter-coder reliability.

- human coders are employed to solve the problem of validity
- it creates a new problem of time, training and budget, but also of inter-coder reliability or agreement
- various measures of inter-coder agreement

Recent approaches to quantitative textual analysis in sociology

How to overcome the problem of coding/dictionaries?

- exploratory research questions rather than testing hypotheses
- inductively constructed classifications rather than pre-established ones
- observing relationship between textual and non-textual data



Paterson, L. L., & Gregory, I. N. (2019). *Representations of Poverty and Place: Using Geographical Text Analysis to Understand Discourse*. Palgrave Macmillan.

Combined textual and geographic analysis to understand the representation of poverty in UK.

Procedure:

1. Construction of the corpus: Guardian & Daily Mail articles containing the word "poverty" and names of localities for a certain period of time.

Example: Mike Barry, once a debt adviser with Citizens Advice, and now operations director of the town's credit union, is dismayed—both by **Blackpool's** worsening **poverty** and by the rise of the corporate moneylenders (Guardian: news, July 2013)

2. Mapping localities linked to poverty for both newspapers.

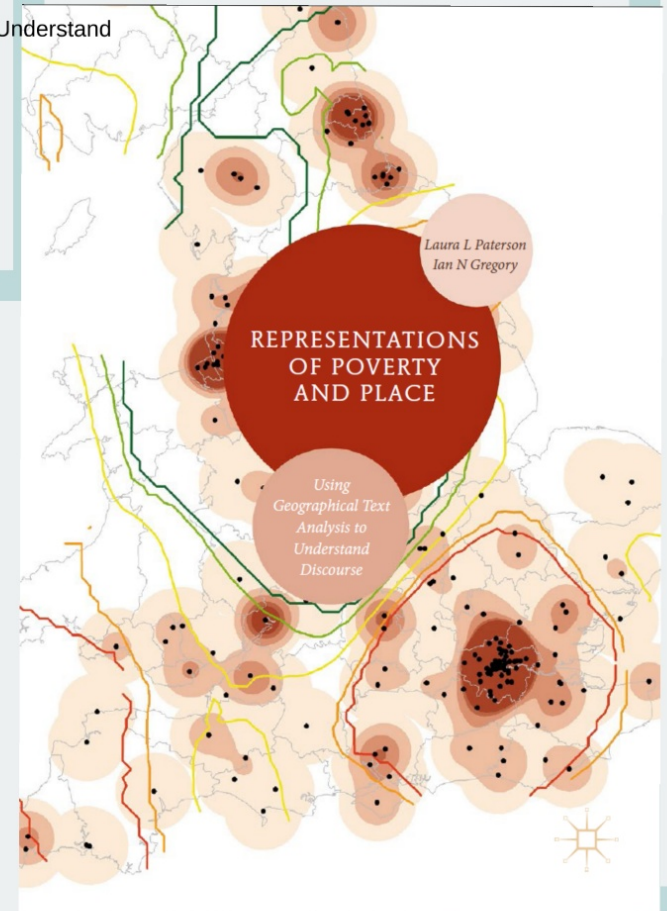
3. Identification of "poverty" collocations in different localities.



Fig. 4.1 Density smoothed maps of "<poverty*>" PNCs: a Guardian b Daily Mail

Table 4.3 PNC keywords in the *Daily Mail* comparing the co-text of "<poverty*>" in London with the rest of the UK

Place	Sig.	Keywords
London	<.01	two, one, rate, reached, country, protest, he, Victorian, years
London	<.05	slum, caf, says, reason, cereals, imported, walked, Keely, poorest, improved, great, compared, depression, Facebook, home, become, during, high, rose, class, byword, nostalgia, despite, situation, capital, broadly, miles, men, walk, show
Rest of UK	<.01	too, risks, lowest, people, left, care, we
Rest of UK	<.05	level, making, road, citing, wherever, small, fuel, herself, find, families, local, her, poorer, five, here, taking, miss, ignoring



The migrant stories analysis

What we can learn about migration from migrant stories?

The story of Jonathan (country of origin: France, country of residence: Vietnam)

My life back in France was not a very good experience for me due to many unfortunate incidents in my personal life. I did not feel happy and was not satisfied with it in a way. Leaving was just a dream for me, a dream to start over with a new life, a dream to have a second chance. It was a distant dream and I had never thought that I would be able to make it happen.

With a nice surprise, the opportunity came when a friend of mine in France asked me to join him in opening new business in Vietnam. Before France I had lived in Brazil for three months but I had no previous experience in Asia. Vietnam was not necessarily my first choice, but I like to think of it as my destiny. So I did not think twice before choosing to follow my instinct and taking part in this adventure.

The first chapter of my life in Vietnam was not easy at all. I could not speak English or Vietnamese...

The migrant stories analysis

What we can learn about migration from migrant stories?

The story of Jonathan (country of origin: France, country of residence: Vietnam)

My life back in France was not a very good experience for me due to many unfortunate incidents in my personal life. I did not feel happy and was not satisfied with it in a way. Leaving was just a dream for me, a dream to start over with a new life, a dream to have a second chance. It was a distant dream and I had never thought that I would be able to make it happen.

With a nice surprise, the opportunity came when a friend of mine in France asked me to join him in opening new business in Vietnam. Before France I had lived in Brazil for three months but I had no previous experience in Asia. Vietnam was not necessarily my first choice, but I like to think of it as my destiny. So I did not think twice before choosing to follow my instinct and taking part in this adventure.

The first chapter of my life in Vietnam was not easy at all. I could not speak English or Vietnamese...

- we can try to understand migrant's experience from individual stories = a hermeneutic (qualitative) approach

The migrant stories analysis

What we can learn about migration from migrant stories?

The story of Jonathan (country of origin: France, country of residence: Vietnam)

My life back in France was not a very good experience for me due to many unfortunate incidents in my personal life. I did not feel happy and was not satisfied with it in a way. Leaving was just a dream for me, a dream to start over with a new life, a dream to have a second chance. It was a distant dream and I had never thought that I would be able to make it happen.

With a nice surprise, the opportunity came when a friend of mine in France asked me to join him in opening new business in Vietnam. Before France I had lived in Brazil for three months but I had no previous experience in Asia. Vietnam was not necessarily my first choice, but I like to think of it as my destiny. So I did not think twice before choosing to follow my instinct and taking part in this adventure.

The first chapter of my life in Vietnam was not easy at all. I could not speak English or Vietnamese...

- we can try to understand migrant's experience from individual stories = a hermeneutic (qualitative) approach
- we can compare stories of different categories of migrants, observe similarities and differences in what they told and make inferences about migration (how is narrated and experienced) = quantitative textual approach

The migrant stories analysis

What we can learn about migration from migrant stories?

The story of Jonathan (country of origin: France, country of residence: Vietnam)

My life back in France was not a very good experience for me due to many unfortunate incidents in my personal life. I did not feel happy and was not satisfied with it in a way. Leaving was just a dream for me, a dream to start over with a new life, a dream to have a second chance. It was a distant dream and I had never thought that I would be able to make it happen.

With a nice surprise, the opportunity came when a friend of mine in France asked me to join him in opening new business in Vietnam. Before France I had lived in Brazil for three months but I had no previous experience in Asia. Vietnam was not necessarily my first choice, but I like to think of it as my destiny. So I did not think twice before choosing to follow my instinct and taking part in this adventure.

The first chapter of my life in Vietnam was not easy at all. I could not speak English or Vietnamese...

- we can compare stories of different categories of migrants, observe similarities and differences in what they told and make inferences about migration (how is narrated and experienced) = quantitative textual approach

Data

- 1018 short biographic narratives of migrants published on iamamigrant.org site
- the stories have been adapted for publication by people or organizations submitting the story and eventually selected by IOM, the UN organization providing help for migrants
- it is not a representative nor unbiased sample of migrant experiences over the world
- it is a very heterogeneous sample of migrants' stories

The screenshot displays the homepage of iamamigrant.org. At the top, there is a navigation bar with the site name and links for 'about', 'stories', 'blog', and 'videos'. Language options (DE, EN, ES, FR, IT, EL) are also present. The main content area is divided into several sections:

- latest story:** Features a large image of a man with a quote: "My fellow brothers need to know that this route is one of sacrifice. You don't know how you will die, but chances are it will happen." - @ideyemach, 3,608 km.
- participate:** A form for users to submit their own stories, including fields for Name, Occupation, Current country, and Country of Origin.
- stories map:** A world map showing the locations of various stories.
- videos:** A small video thumbnail with the title "My name is Fatima".
- featured stories:** A grid of story cards, each with a photo, a distance, and a quote. The cards are filtered by 'country of origin' and 'current country'.
 - Eliza, 7,211 km:** "I help scientists reach out to school children and share their research and love of science."
 - Marlon, 369 km:** "I can live and work wherever I want. Borders are unimportant to me."
 - Immae, 1,514 km:** "I expected to discover the perfect Eldorado, but reality quickly caught up with me." - @out
 - Shahin, 3,768 km:** "Refugees are like newborn babies. You have to start over completely - learning to live in a new culture and abide by new rules." - @iamamigrant
 - Yolande, 6,293 km:** "Migrants are here. It's a fact so there should be strong policies of inclusion." - @iamamigrant
 - Youssef, 19,320 km:** "My business is running well and my income is enough to support myself and my family." - @out
 - Quinn, 10,056 km:** "Integration is not being scared of what is around you. It's when you see what people do around you as normal." - @iamamigrant
 - Victoria, 2,129 km:** "I started a London-based Ad agency."

At the bottom, there is a 'view more' button, the 'TOGETHER' logo, and links to the website and the Sustainable Development Goals page. The footer contains copyright information for the IOM Media and Communications Division and a link to the Privacy Policy.

Data processing

- data were downloaded from <https://iamamigrant.org> using wget utility for retrieving files from the Internet
- html files (web pages) were converted into plain text files
- names, countries (origin & destination) and stories were extracted from the text files
- all above done by Jiří Mírovský

- paratextual elements were removed from stories, e.g. "Watch Natasha's full story on Immigrant Stories."; stop words were also removed
- gender of narrators was identified semiautomatically (gendernamefinder.com + manual search)
- GDP per capita for countries was downloaded from UN statistics
- countries were manually grouped into continents

Analysis

Our analysis will be exploratory.

What we know about narrators: gender, the original and current country of residence; and GDP per capita for countries. We then constructed following independent variables:

female - male migrant
immigrant - homecomer
intracontinental - intercontinental migrant
higher - equal - lower GDP migrant

A dependent variable - the migrant's story.

We will try to search if there is any relationship between independent and dependent variables. For example, what kind of impact, if any, gender of narrators has on their stories.

The analysis was carried out in R, a free software environment for statistical computing and graphics (<https://www.r-project.org/>).

Text analytical techniques used in the study

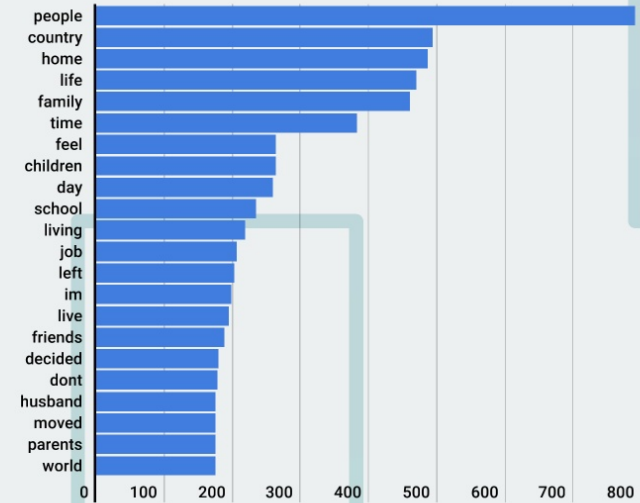
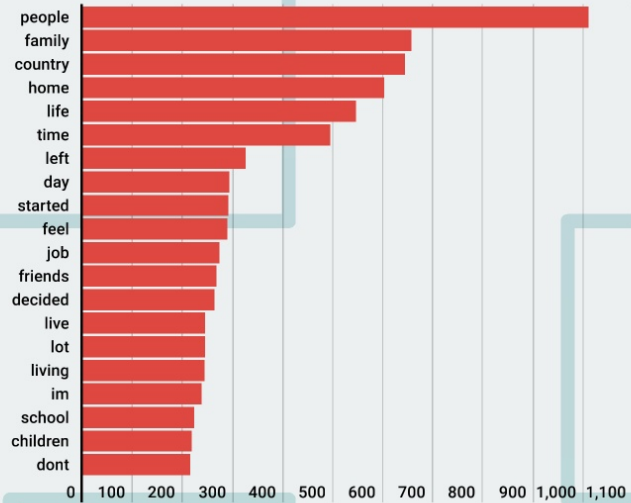
Word frequencies: calculation of most frequent semantic words; it gives us basic information about words through which the migrant experience is expressed

Bigrams: identification of frequent pairs (collocations) of semantic words occurring in the narratives; it informs us about typical entities, such as "primary school" or "speak English"

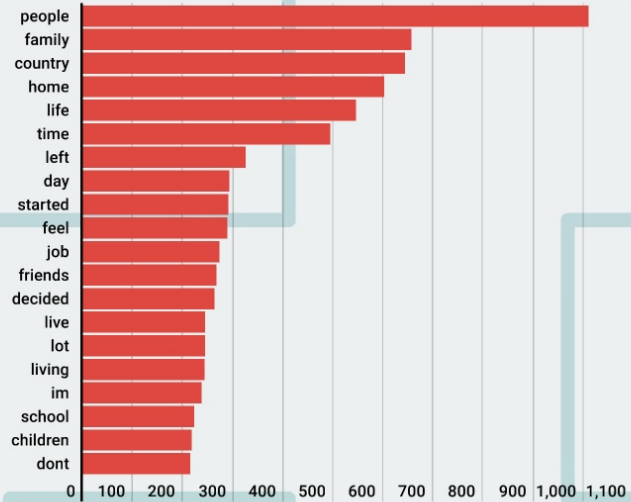
Topic modeling: identifying clusters or recurring patterns of co-occurring words (called "topics"); it provides information about potential themes or topics in the narratives

increase in analytical complexity

Findings I. - word frequencies



Findings I. - word frequencies

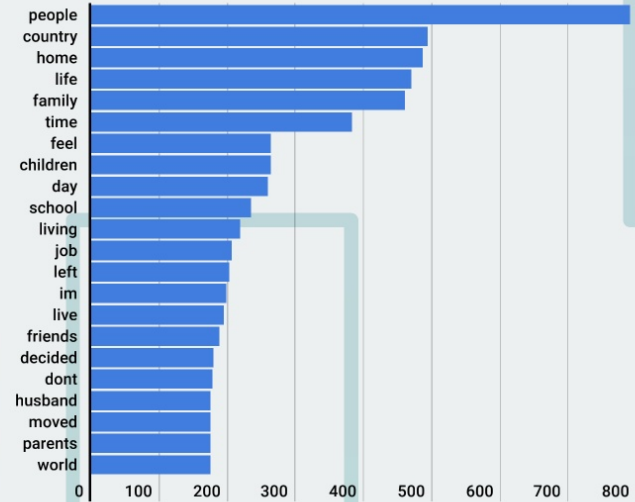


people

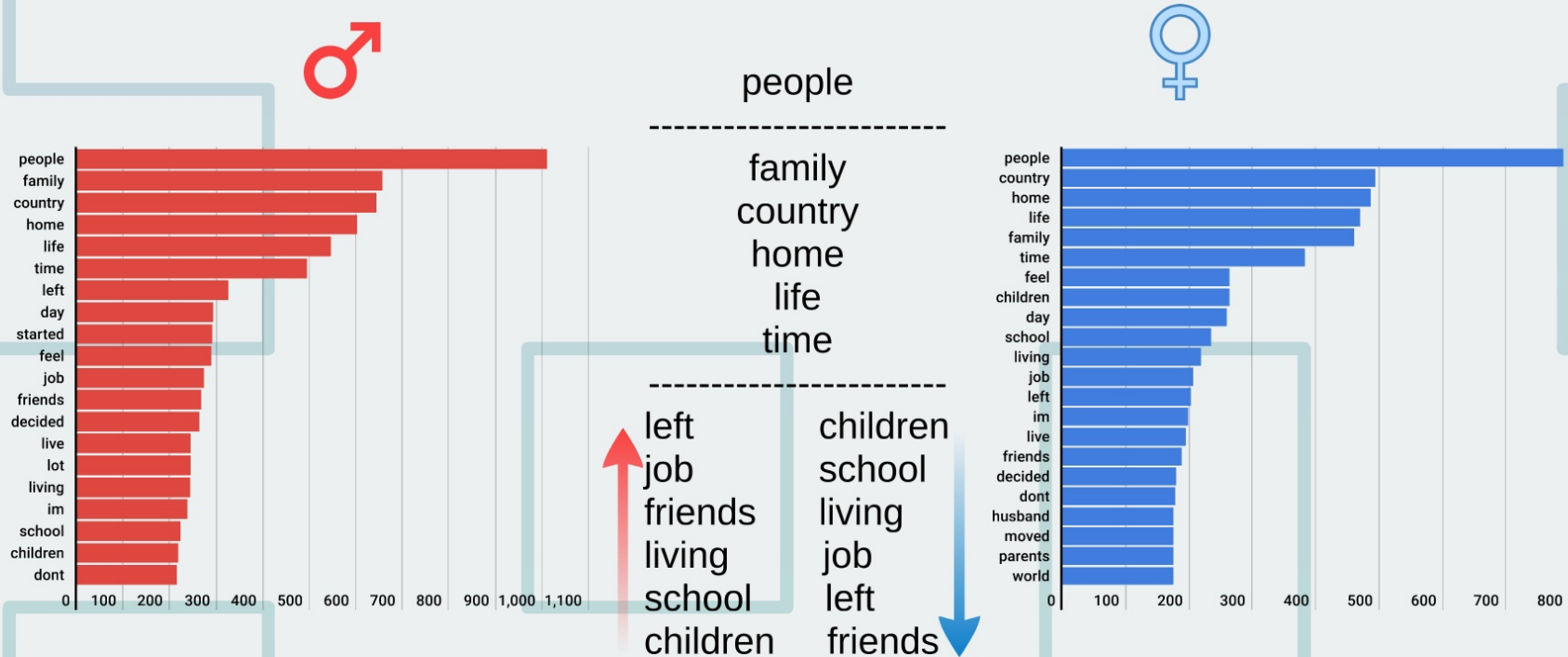
family
country
home
life
time

left
job
friends
living
school
children

children
school
living
job
left
friends



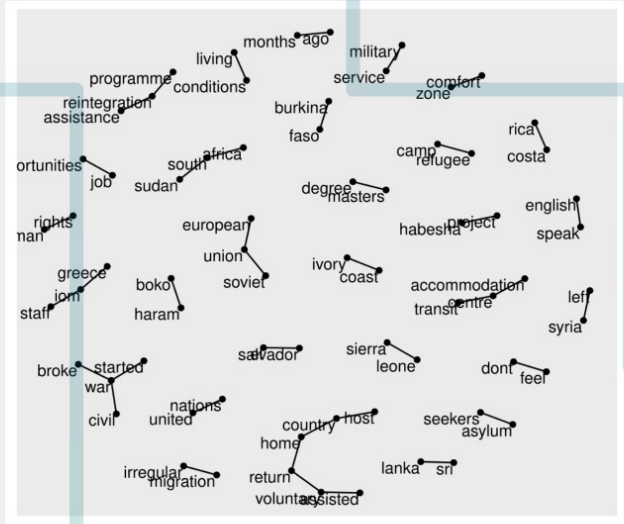
Findings I. - word frequencies



Results: There are no surprising differences between male and female migrants' stories on the level of words. It seems that migration is linked to families; migrants of both genders express it through similar words (-- share similar experiences?).

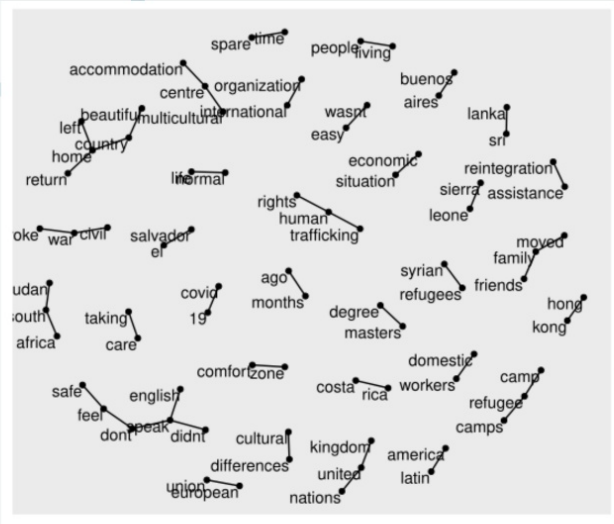
Findings II. - bigrams

men



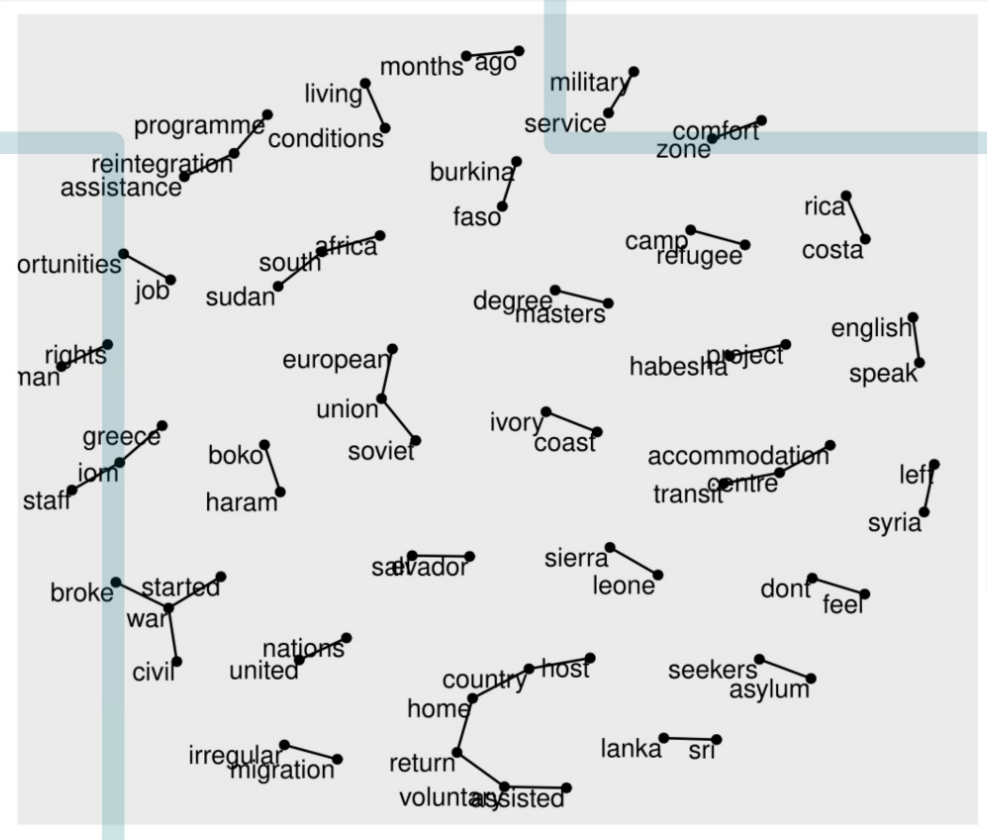
minimal occurrence = 9

women



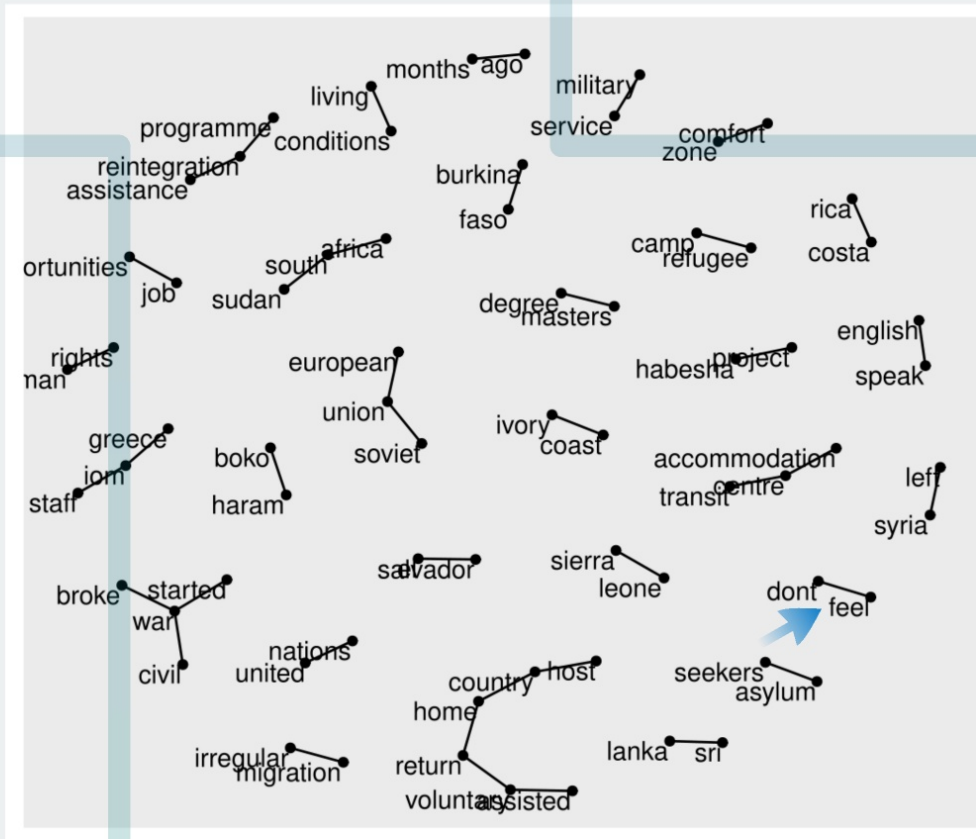
minimal occurrence = 6

men



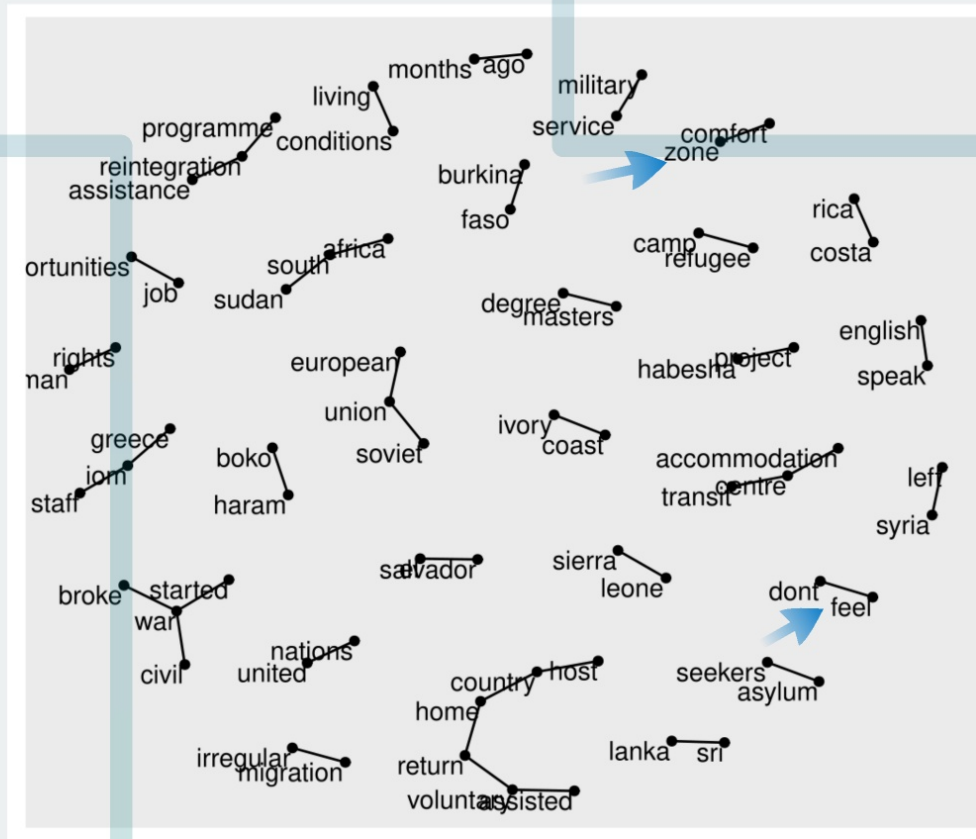
minimal occurrence = 9

men



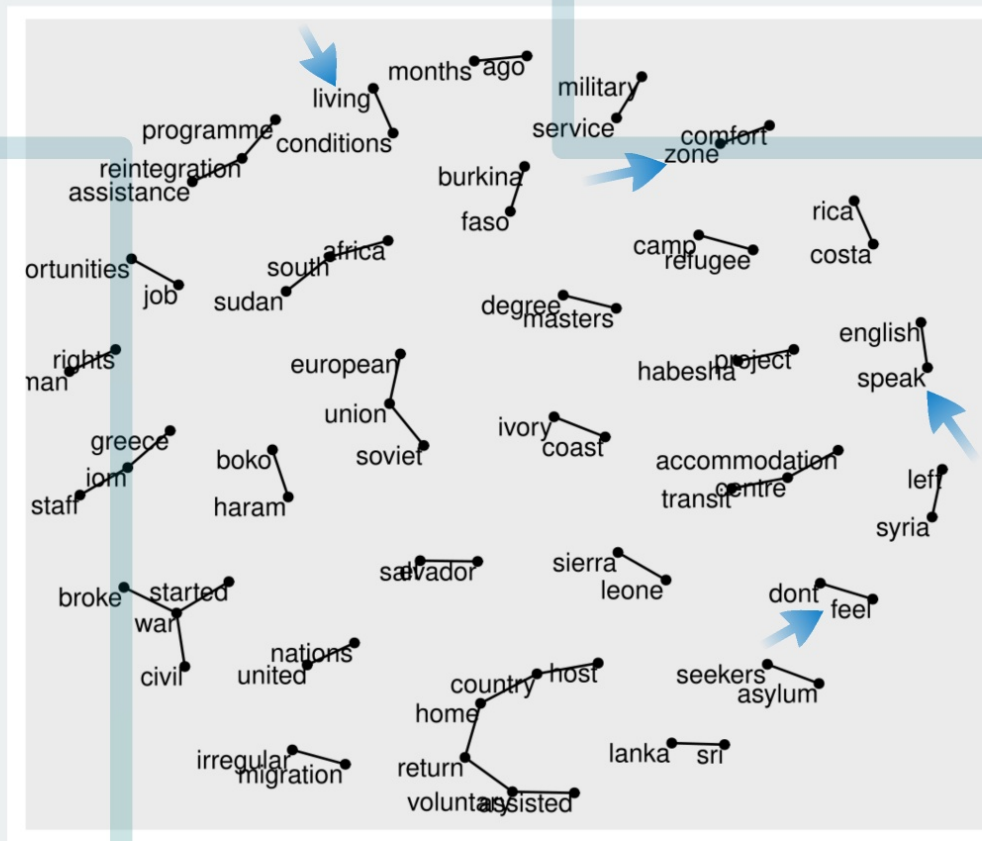
minimal occurrence = 9

men



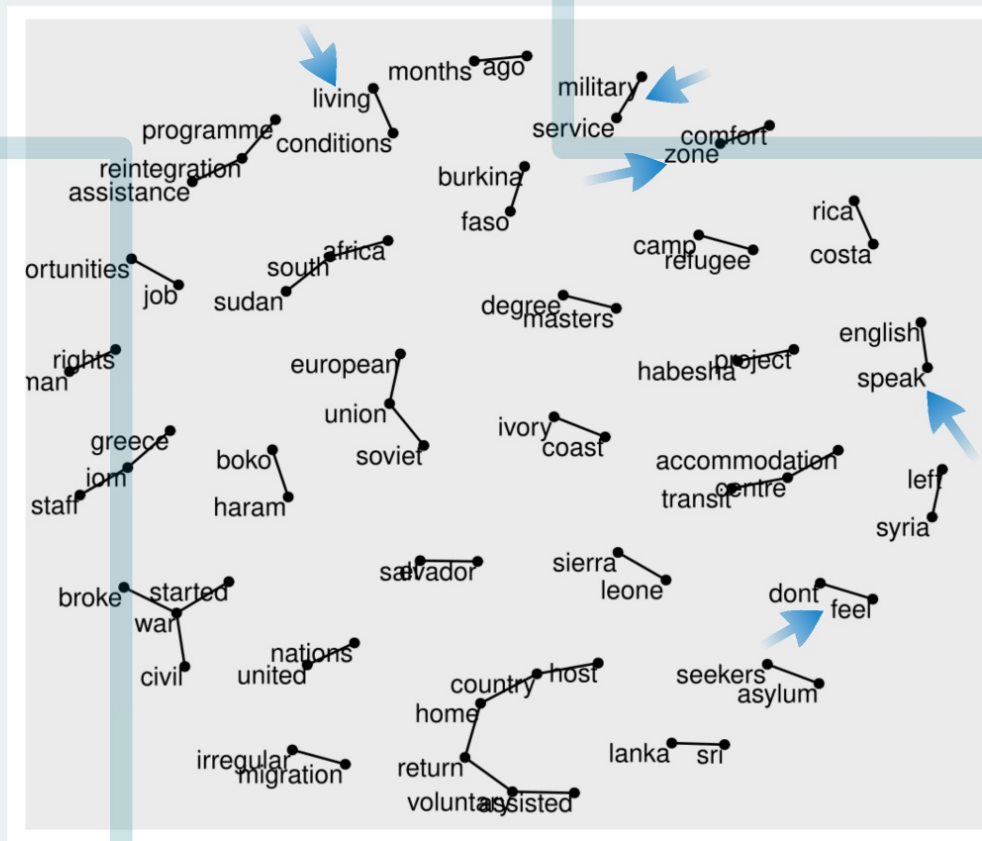
minimal occurrence = 9

men



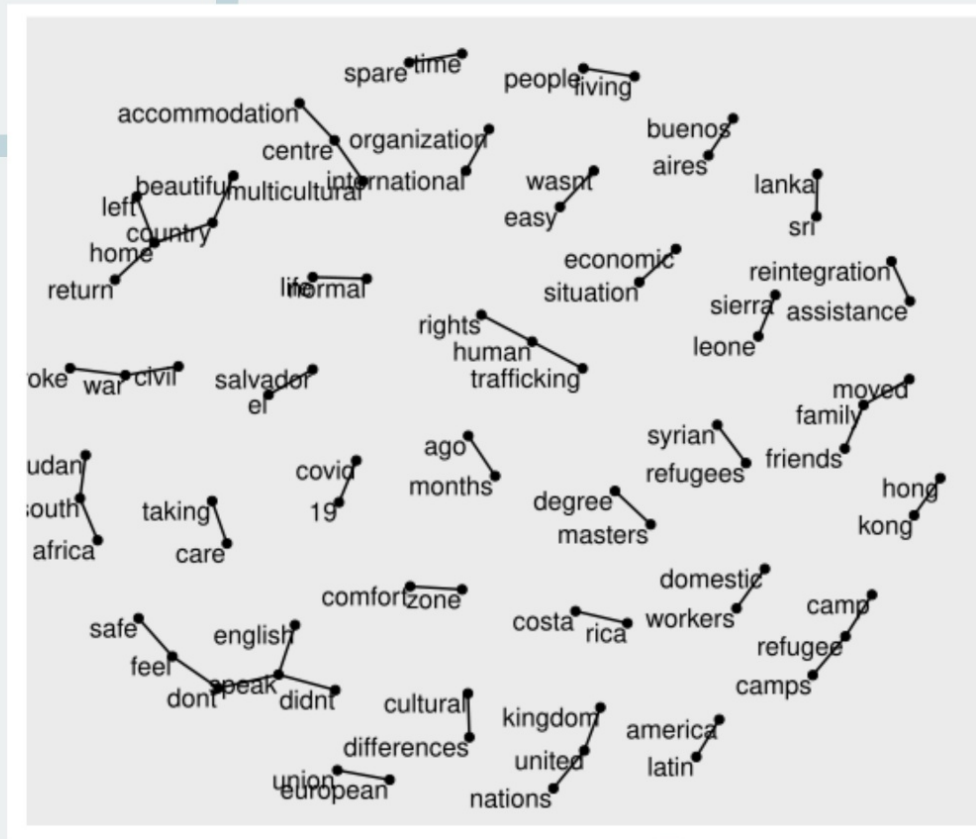
minimal occurrence = 9

men



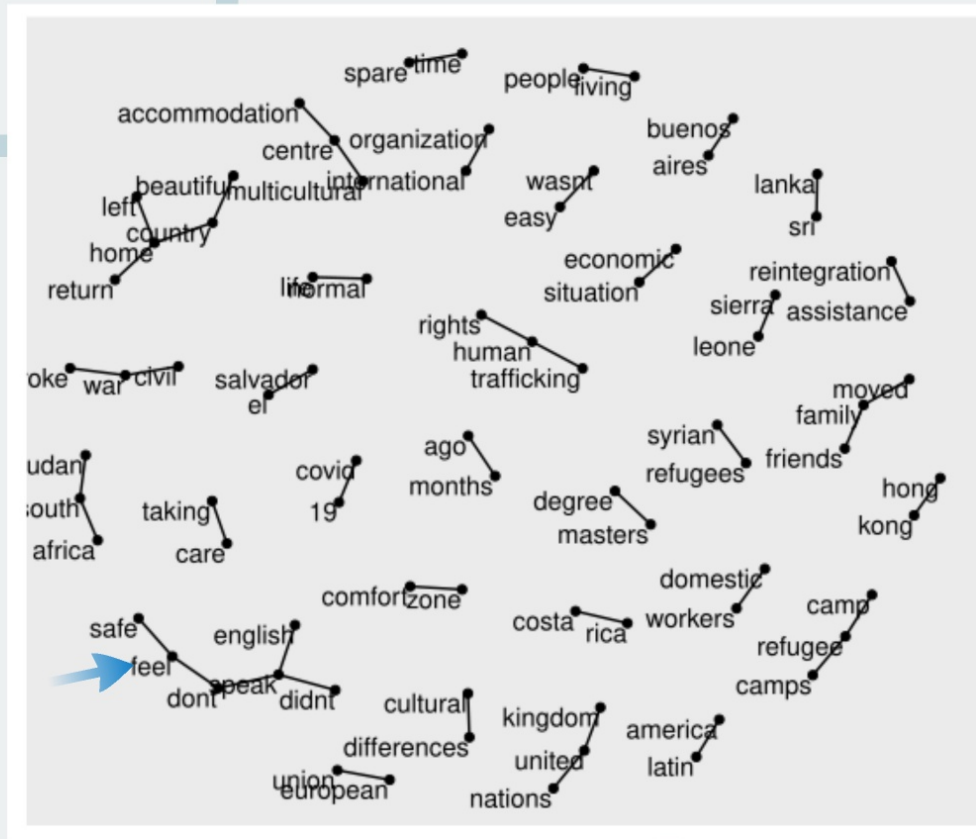
minimal occurrence = 9

women



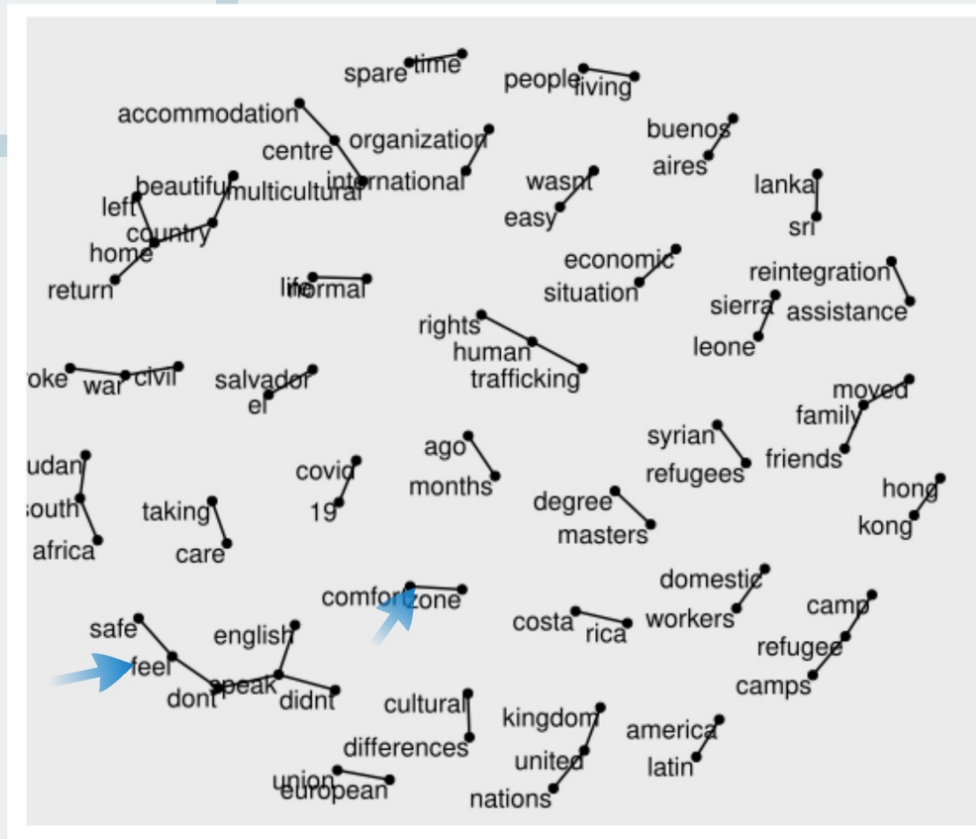
minimal occurrence = 6

women



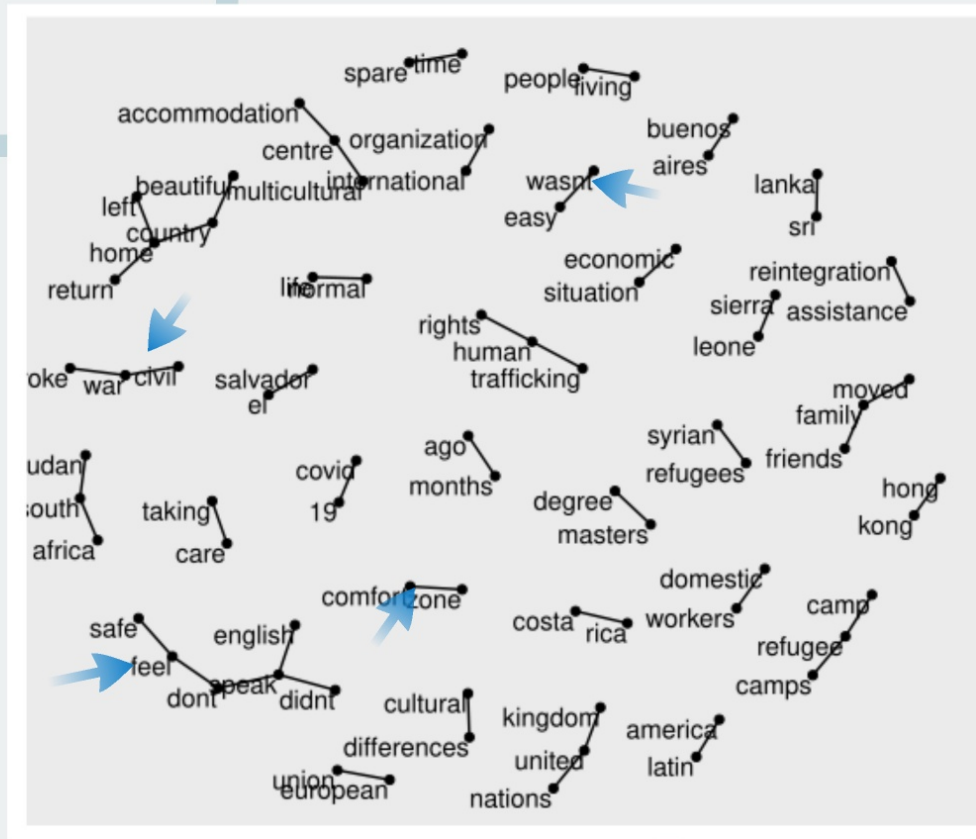
minimal occurrence = 6

women



minimal occurrence = 6

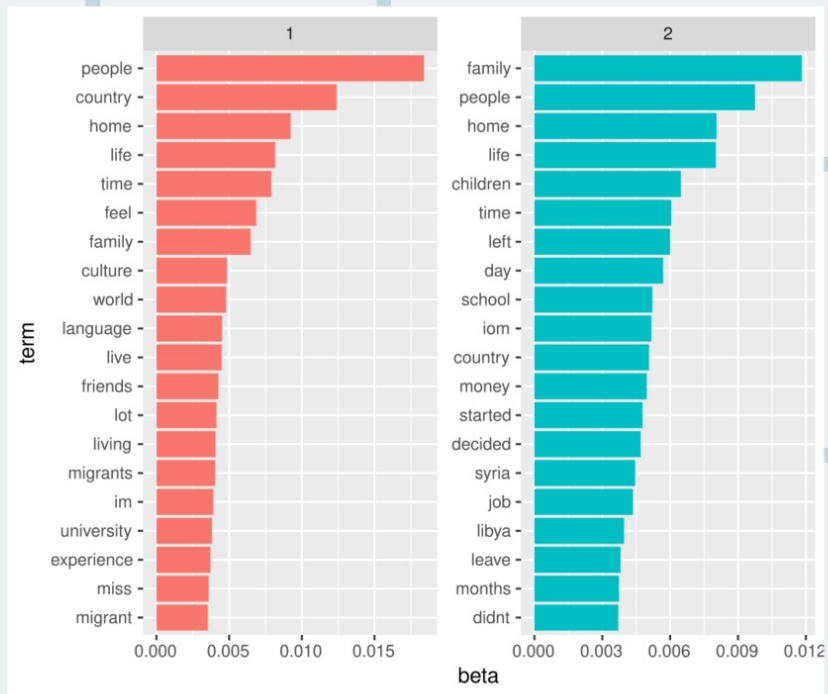
women



minimal occurrence = 6

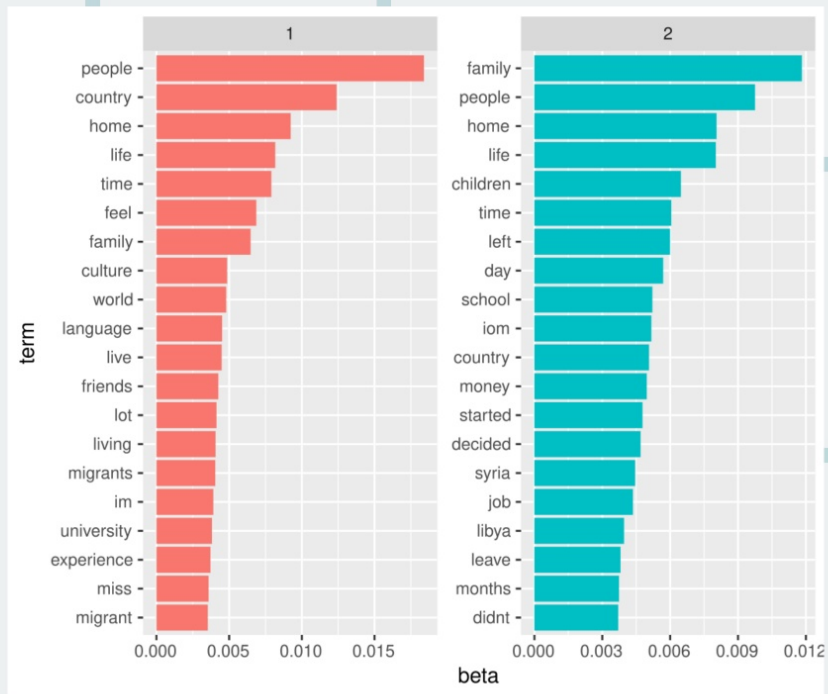
Findings III. - topic modeling

- computing macro patterns of words ("topics") in the narratives
- a single narrative can contain one or more topics
- probabilistic estimation = it can generate valid results but also artifacts



Findings III. - topic modeling

- computing macro patterns of words ("topics") in the narratives
- a single narrative can contain one or more topics
- probabilistic estimation = it can generate valid results but also artifacts



Topic 1

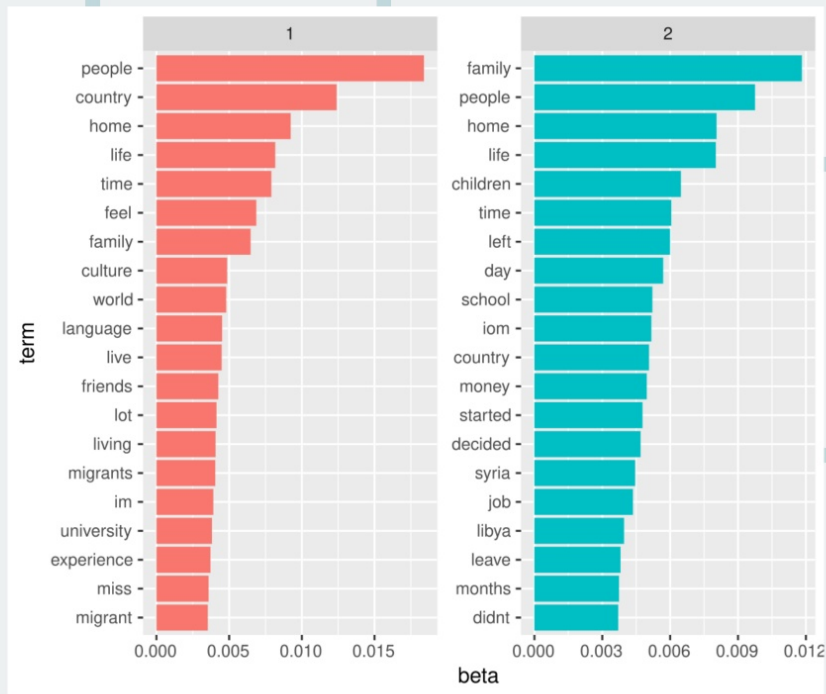
country, feel, culture, world, language, friends, migrants, uni, experience, miss

Topic 2

family, children, left, school, IOM, money, started, decided, Syria, Libya

Findings III. - topic modeling

- computing macro patterns of words ("topics") in the narratives
- a single narrative can contain one or more topics
- probabilistic estimation = it can generate valid results but also artifacts



Topic 1

country, feel, culture, world, language, friends, migrants, uni, experience, miss

Topic 2

family, children, left, school, IOM, money, started, decided, Syria, Libya

Interpretation

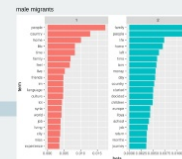
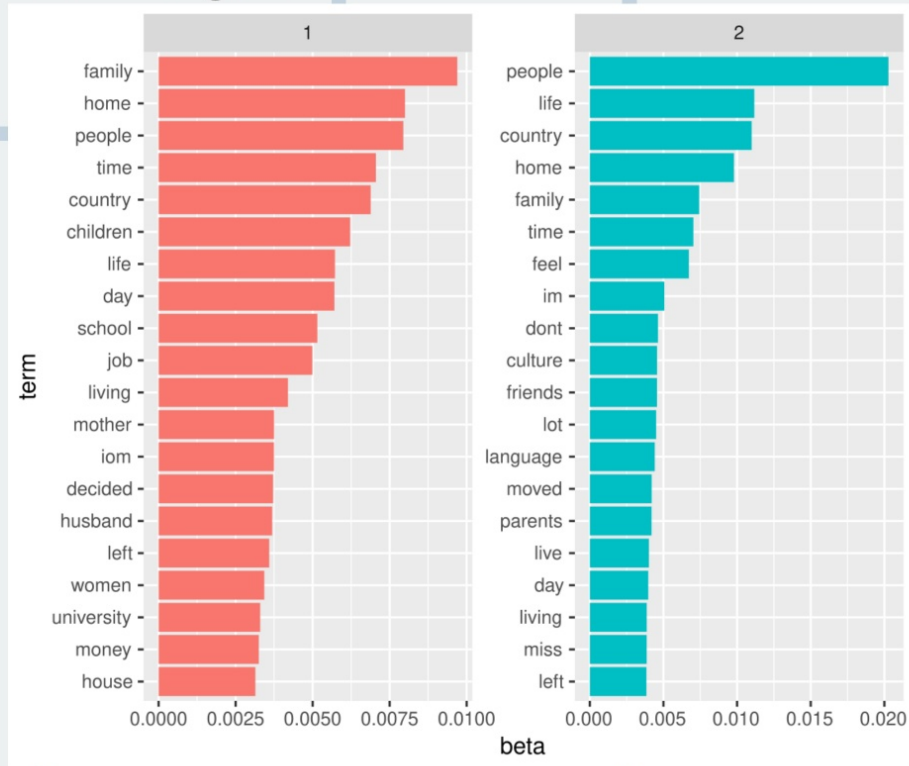
Topic 1 - general narratives of migrants about their experience of moving from one country to another, encountering different cultures, languages and feelings associated with this change

Topic 2 - a particular theme of war migrants (refugees) who fled from their homes and sought asylum elsewhere. They often come from Syria or transit through or end in Libya.

Findings III. - topic modeling (cont.)

- gendered decomposition

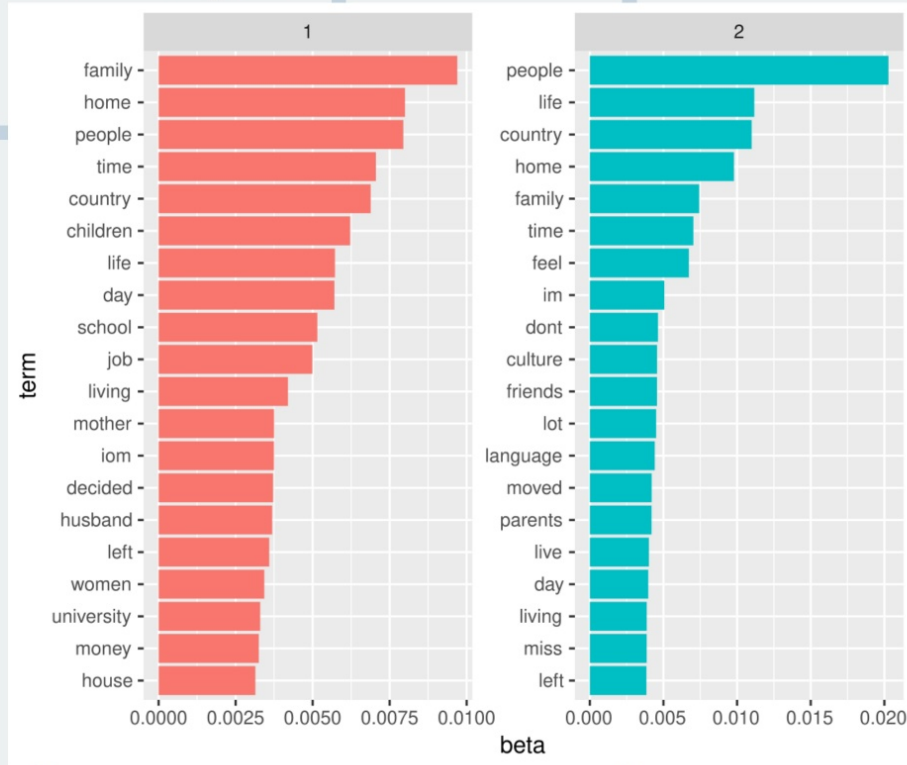
female migrants



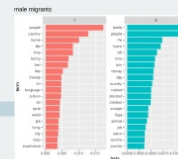
Findings III. - topic modeling (cont.)

- gendered decomposition

female migrants



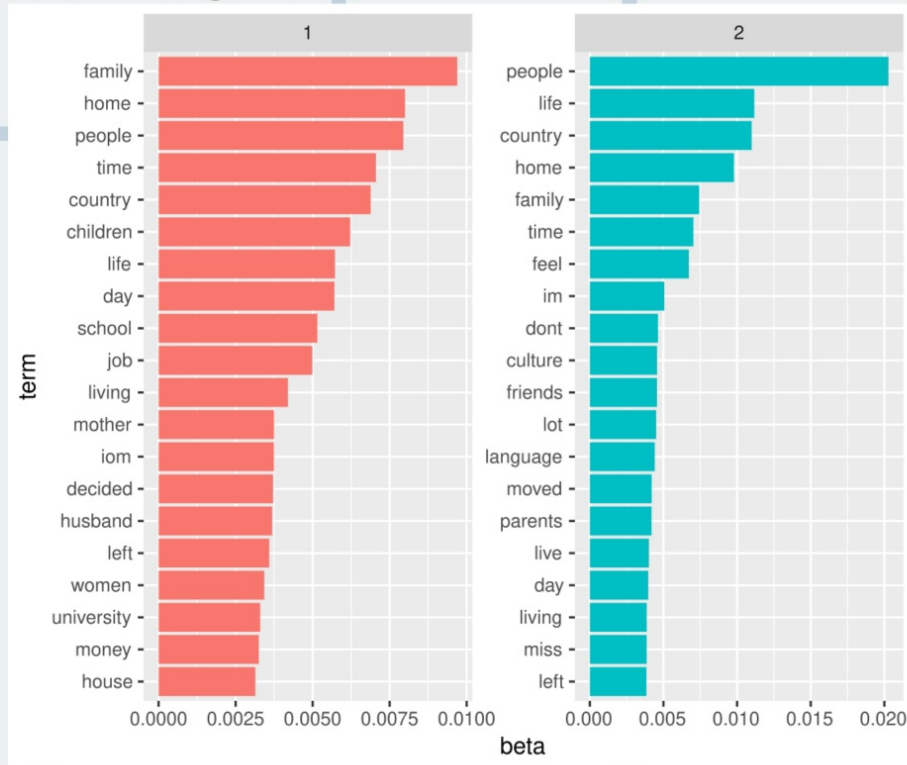
War migrant topic is dominant, general migrant experience topic is secondary.



Findings III. - topic modeling (cont.)

- gendered decomposition

female migrants



War migrant topic is dominant, general migrant experience topic is secondary.

Interpretation

A. There is more war than non-war female migrants in the sample as they are less often voluntary migrants.

B. Migration is always challenging for female migrants; therefore, they express their experience similarly to war migrants.



Homework

Carry out your own analysis on the same sample of narratives. Compare groups in any subset of narratives:

- immigrants vs. homecomers
- transcontinental vs. intracontinental
- immigrants to countries with (much) higher, equal or (much) lower GDP
- men vs. women.

Download data from google drive <https://tinyurl.com/26vpzrj6>

Use a web application Voyant (<https://voyant-tools.org/>) to explore similarities and differences between groups of stories. Play with the app and try different analytical tools.

Make a report with your findings and tentative interpretation.



Deadline: Sunday 10.4.2023 midnight.

Quantitative Textual Analysis in Sociology I

Sociology and Quantitative Textual Analysis

Sociology works with quantitative data from its beginning.

Data are primarily non-textual, such as socio-demographic data, opinions, attitudes, and behaviour.

If textual data enter into analysis, they are coded and codes further analyzed. For example, answers to open-ended questions in a survey questionnaire (What do you feel is the most important issue facing the world today?) are coded according to the problems mentioned (climate change, terrorism, inequality etc.)

Coding textual data: problems of validity and (inter-coder) reliability.

Analysis

Our analysis will be exploratory.

What we know about narrators: gender, the original and current country of residence, and GDP per capita for countries. We then constructed following independent variables:

female - male migrant
immigrant - homecome
intracontinental - intercontinental migrant
higher - equal - lower GDP migrant

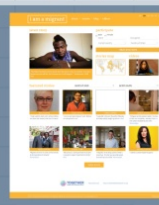
A dependent variable - The migrant's story

We will try to search if there is any relationship between independent and dependent variables. For example, what kind of impact, if any, gender of narrators has on their stories.

The analysis was carried out in R, a free software environment for statistical computing and graphics (<https://www.r-project.org/>).

Data

- 1018 short biographic narratives of migrants published on samigrant.org site
- The stories have been adapted for publication by groups or organizations submitting the story and eventually selected by IOM, first UN organization providing help for migrants
- It is not a representative but unbiased sample of migrant experiences over the world
- It is a very heterogeneous sample of migrants' stories



The migrant stories analysis

What we can learn about migration from migrant stories?

- we can compare stories of different categories of migrants, observe similarities and differences in what they said and make inferences about migration (how is narrated and experienced) = quantitative textual approach

Recent approaches to quantitative textual analysis in sociology

How to overcome the problem of coding/dictionaries?

- exploratory research questions rather than testing hypotheses
- inductively constructed classifications rather than pre-established ones
- observing relationship between textual and non-textual data

Data processing

1. Data cleaning and preprocessing
2. Data normalization
3. Data segmentation
4. Data encoding
5. Data storage and retrieval
6. Data visualization
7. Data analysis and interpretation
8. Data evaluation and validation
9. Data documentation and archiving
10. Data security and privacy protection

