**Research Popularity Prediction**

# Table of Contents
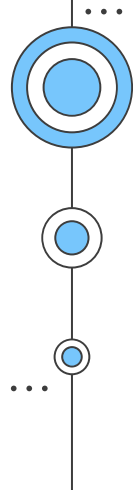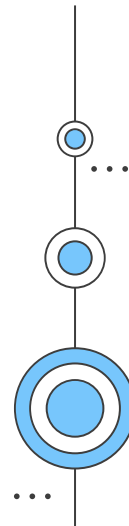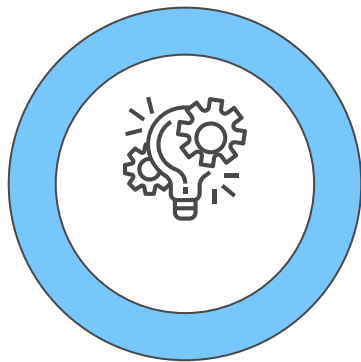
# 01
## Introduction

# Problem Description

We're going to analyse what factors contribute to social science research articles popularity and predict the popularity of articles.

...

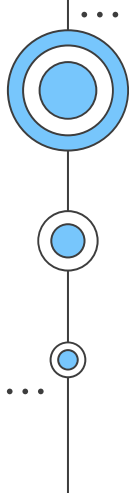# 02
## Data Description

- We scraped the data from researchgate.
- All articles are from social science departments.

# Dataset

**(1899,12)**

| column | Description | Type |
|--------|-------------|------|
| **title** | article's title | string |
| **author** | author name | string |
| **abstract** | Abstract full-text | string |
| **category** | article, literature review, conference paper..etc | string |
| **date_published** | date the article was published | date |
| **date_added** | date the article was uploaded to researchgate | date |
| **figures** | 1 if there're figures & 0 if not | int |
| **full_text?** | availability of full text ( using download or Request full-text as keywords) | string |
| **citation** | number of times that paper was cited | int |
| **reads** | number of views (target variable) | int |

· · ·

# Data Collection Limitations

Have To use organization email Address

Limited Quota Per Day

# 03

## Data preprocessing

# Preprocessing flow

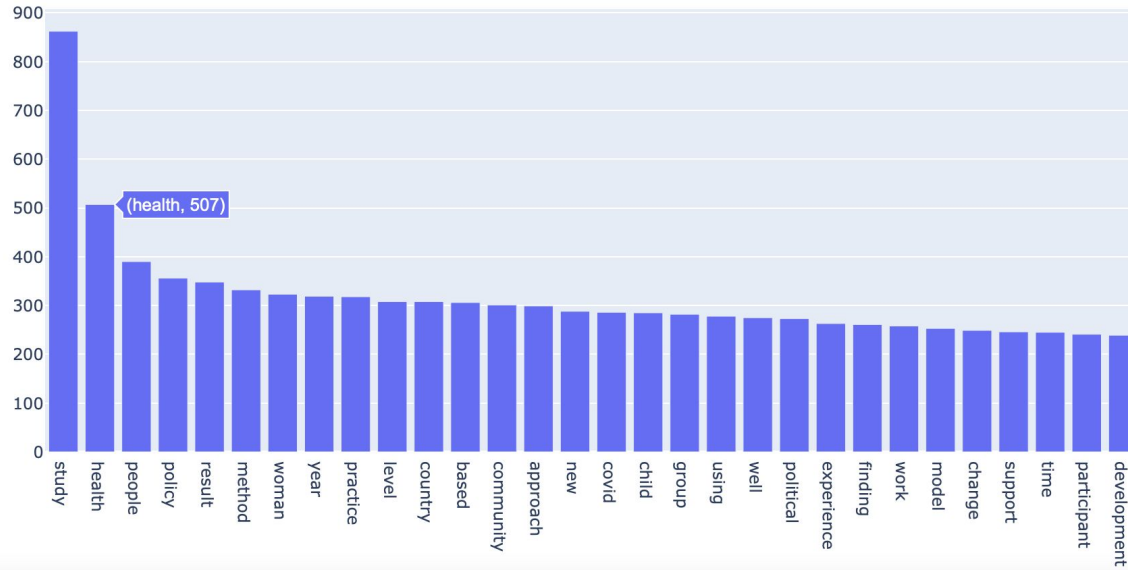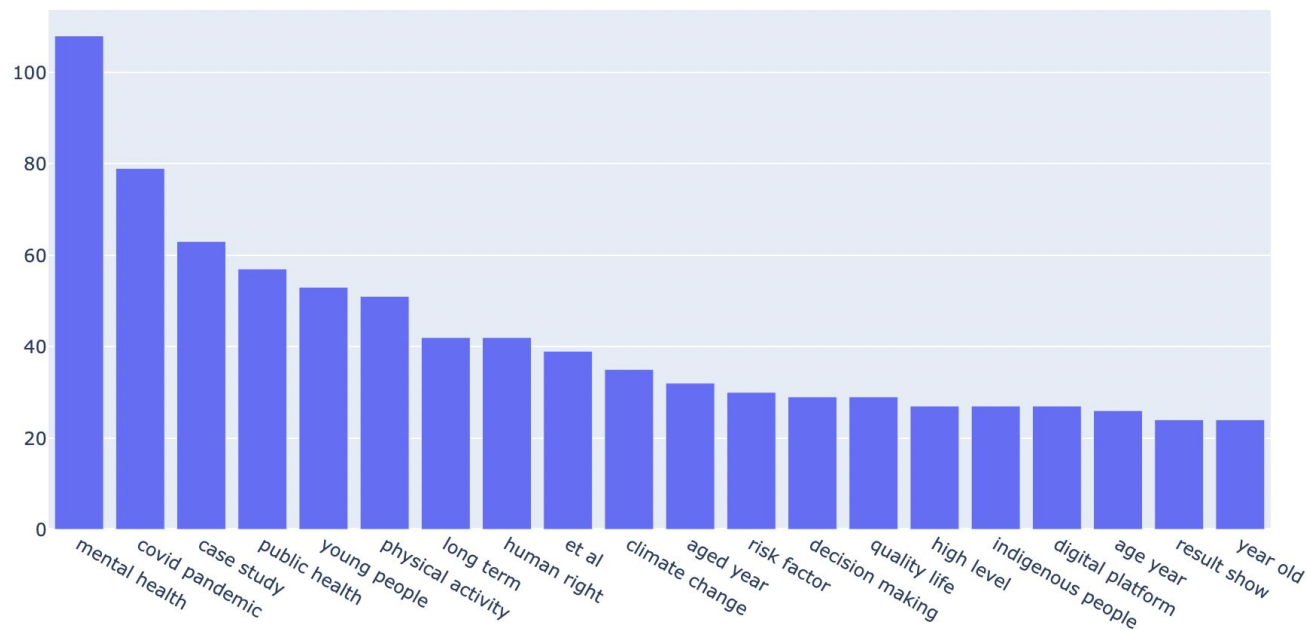Removed non-english articles.

**(1899,12)**

**(1185,12)**

Text Preprocessing

Removed nulls from abstract column

Pos and Tokenization

Removed Stopwords

Applied Lemmatization WordNetLemmatizer

Vectorization: Tfidf / CV

1 2 3 4 5 6 7

# 04

## Topic Modeling

NMF, LDA, CorEx

Top 30 unigrams in the abstract after removing stop words and lemmatization



(health, 507)

Top 20 bigrams in the abstract after removing stop words and lemmatization

# NMF Topic Modeling

TFIDF Vectorization

**16>12>10>9**

Number of topics

# NMF Topic Modeling

**Topic 0**
political, policy, state, new, economic, power, right

**Topic 1**
sexual, drug, violence, gender, minority, health, sex, sexuality, identity, lgbtq

**Topic 2**
covid, health, pandemic, mental, lockdown, country, financial, psychologica

**Topic 3**
woman, party, gender, labour, policy, migrant, feminist, men

**Topic 4**
police, crime, justice, knife, policing, criminal, violence, armed,

**Topic 5**
student, school, education, university, learning, digital, online, teaching

**Topic 6**
patient, age, study, factor, disease, ci, risk, cancer

**Topic 7**
child, parent, family, food, childrens, abuse, parental

**Topic 8**
care, people, service, community, disability, health

# CorEX Topic Modeling

**Count Vectorization**

**12**
Number of topics

Used words from NMF as anchor words, with correlation strength of **7**

# CorEX Topic Modeling



Topic #0
political right
migrant
state
surveillance
politics power
regime
biopolitical

Topic #1
online covid
pandemic protective cov
lockdown
coronavirus
syndrome virus

Topic #2
education
school undergraduate
learning student
teaching academic
university
educational

Topic #3
access people care
young
community
professional training
service support

# CorEX Topic Modeling

### Topic #4
emotion mental sd recruited psychological depression distress stress burnout

### Topic #5
armed violence crime knife police policing rape officer criminal

### Topic #6
childrens child family birth sexually childhood abuse parent friend

### Topic #7
sexual queer sex removal festival survivor sexuality lgbtq marriage

# CorEX Topic Modeling



**Topic #8**
country  poverty
economy  market
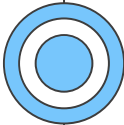economic  income
financial
security
investment

**Topic #9**
environment  city
ecological  climate
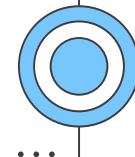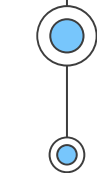environmental
critique  field  issue
sustainability

**Topic #10**
men  woman
percent  gender
experience
gendered  identify
bisexual  feminist

**Topic #11**
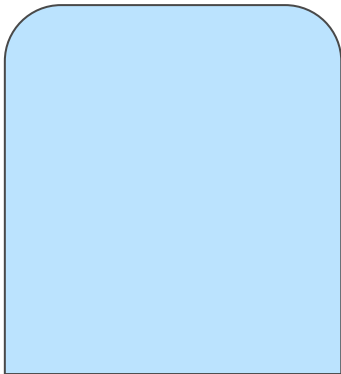age  health
drug  patient
pain
disease
cancer  chronic
physical

# NMF Topics in Bigrams

# CorEx Topics in Bigrams
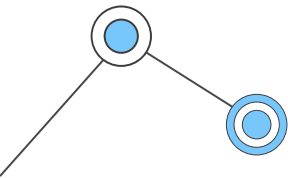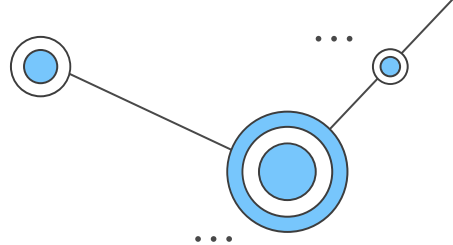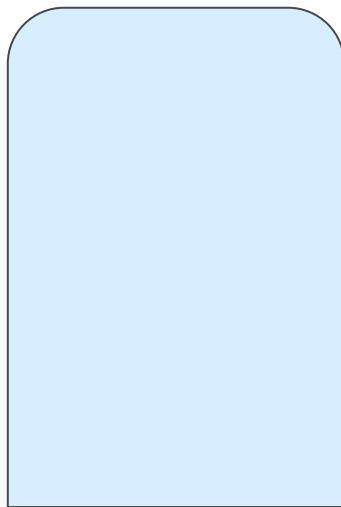
# Chosen Model

**Corex**

Doesn't give a probability

**NMF**

Well separated topics, probability

# Topics

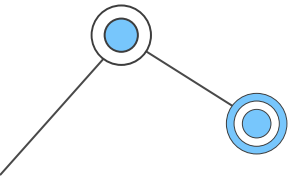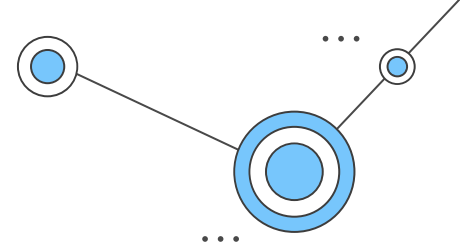Politics And Economics

Sexology

COVID

Women

Criminology

Education

Health & Medicine

Family & Child

Community Service

Generating Subtopics

# Subtopics example: COVID

Topic 0 **( Covid & Mental Health )**
mental, health, covid, psychological, distress, gi, symptom, impact, disorder, participant

Topic 1 **( Covid & economy)**
financial, well, inclusion, study, economy, country, understand, effect, system, pandemic

Topic 2 **( lockdown )**
behavior, lockdown, change, time, diary, risk, activity, location, uk, population

Topic 3 **( Policies )**
health, policy, covid, response, pandemic, medium, crisis, state, news, public

# Subtopics example: Sexology

Topic 0 **( Sexual Health )**
sexual, health, people, adult, young, minority, study, lockdown, older, identity

Topic 1 **(?)**
coumarins, coumarin, compound, various, widely, performed, perfume, crystalline, explicates, scent

Topic 2 **( Sex & Drugs )**
drug, lgbtq, gender, consumption, practice, sex, particular, effect, consumer, policy

Topic 3 **( Sexual Violence )**
violence, festival, sexual, music, argue, feminist, harm, framework, experience, victim

# 05

## Regression

# Predicting the number of Reads

# Feature Engineering

**Adding dominant topic column**

**Title & abstract sentiment polarity**

**Reformatting addition date to days**

**Encoding Journal column**

**(1185,43)**

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |

**(1185,12)**

**Adding probability per topic columns**

**Title & Abstract length in words / sentences**

**Number of authors**

**Mapping Full_text column to 0 and 1**

Correlation Between Reads and other Features
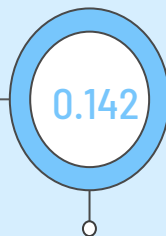
Baseline Model

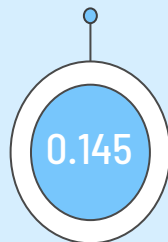Fill null values with zero

score of the train

score of the test

0.145

0.142

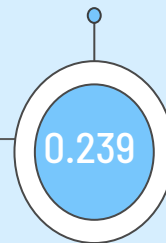0.233

score of the validations

# Data prep & Experimentation
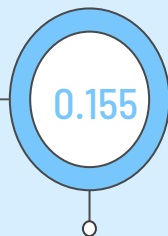
## Fill null values with mode/mean
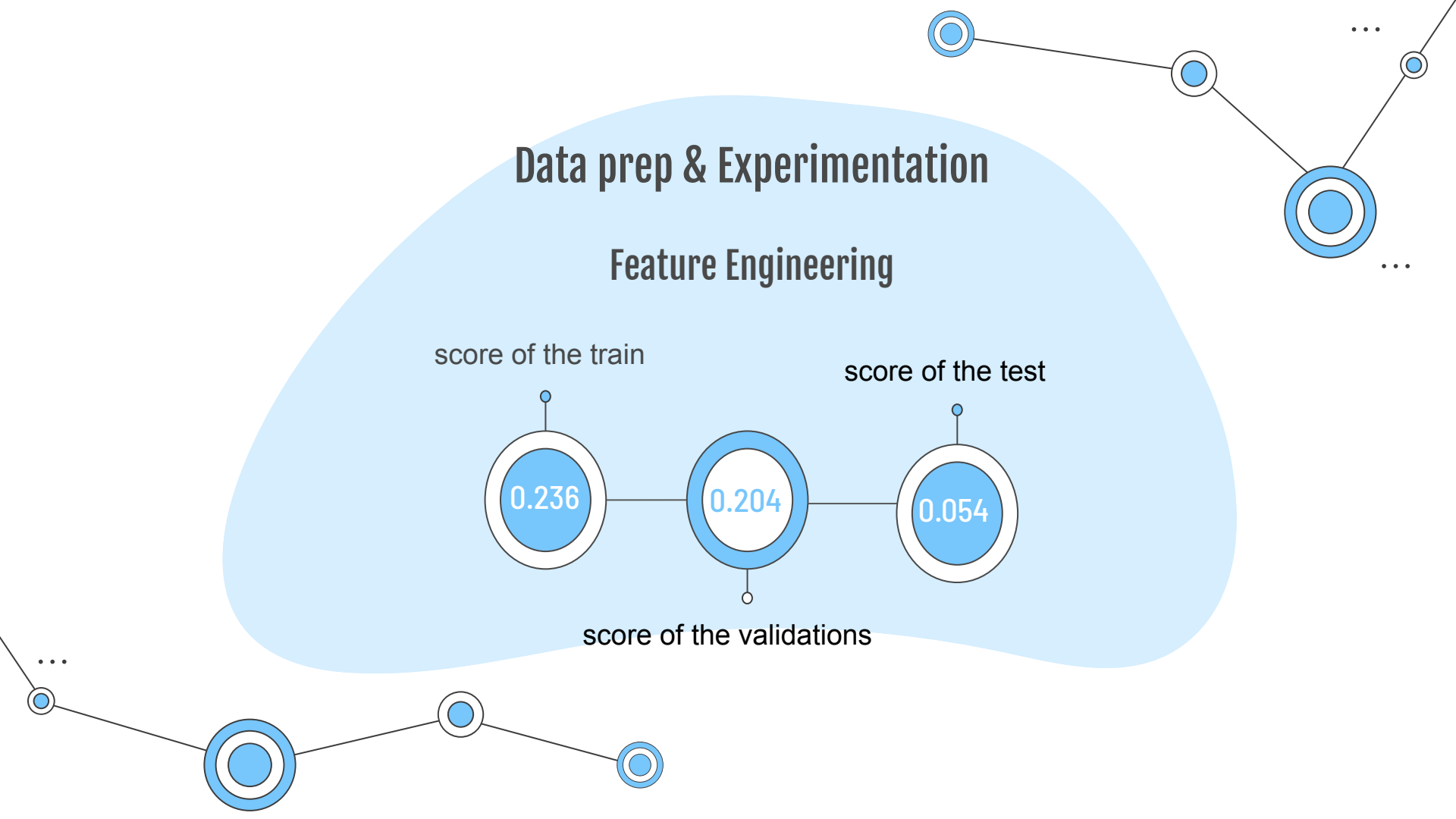
score of the train

score of the test
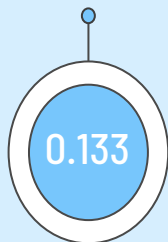
0.145

0.155

0.239

score of the validations

Data prep & Experimentation

Ridge

score of the train

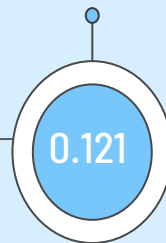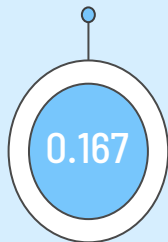score of the test

0.133
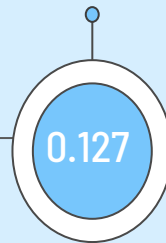
0.172

0.121

score of the validations

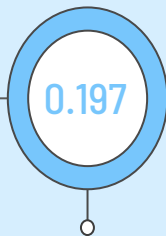Data prep & Experimentation

Lasso

score of the train

score of the test

0.167    0.197    0.127

score of the validations

Based on the unsatisfactory results in the linear regression, we decided to for the problem as a classification problem.
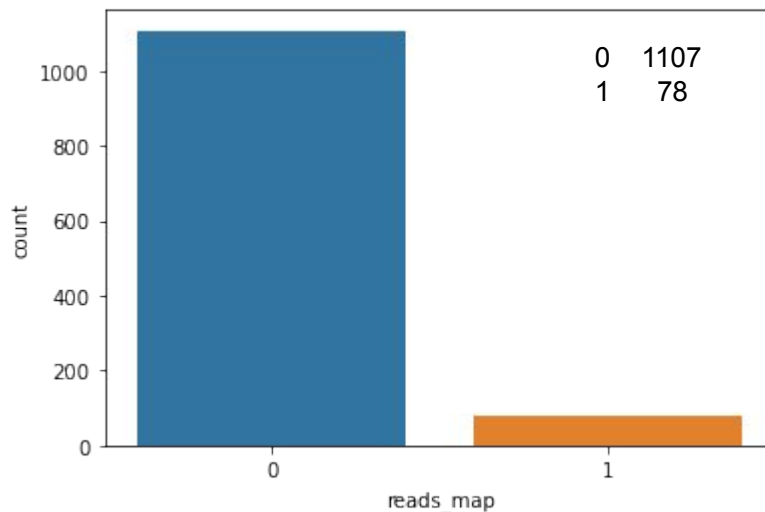
# 06

## Classification

# Data Distribution After Mapping

**Average #reads = 57**

**Popular if Reads >= 200**

# Discretization!

grouping continuous data into a
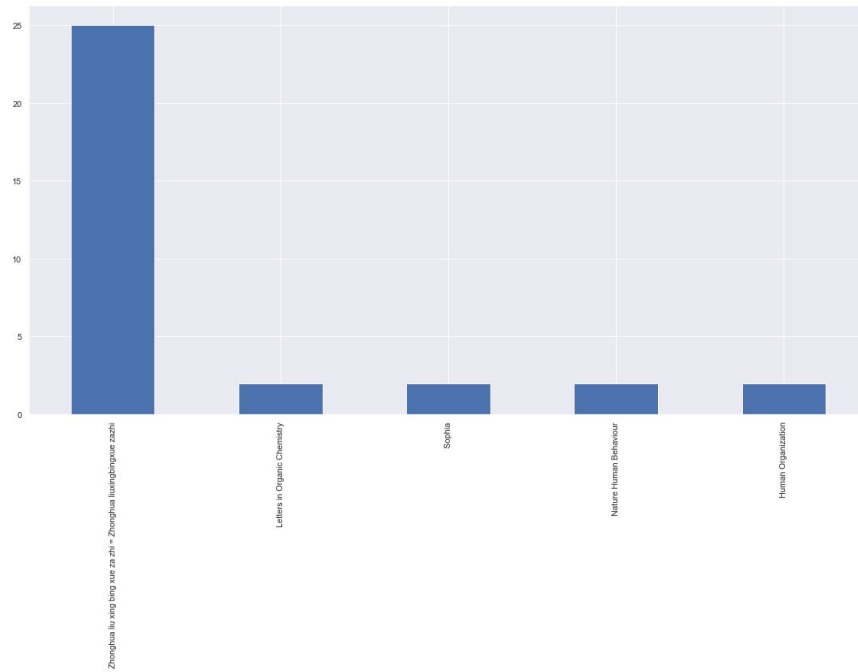number of intervals or bins
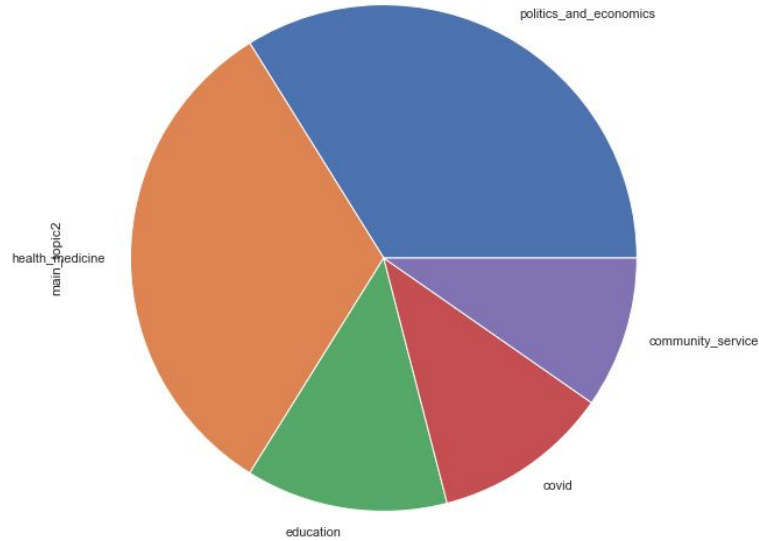
# Popular Articles numerical description
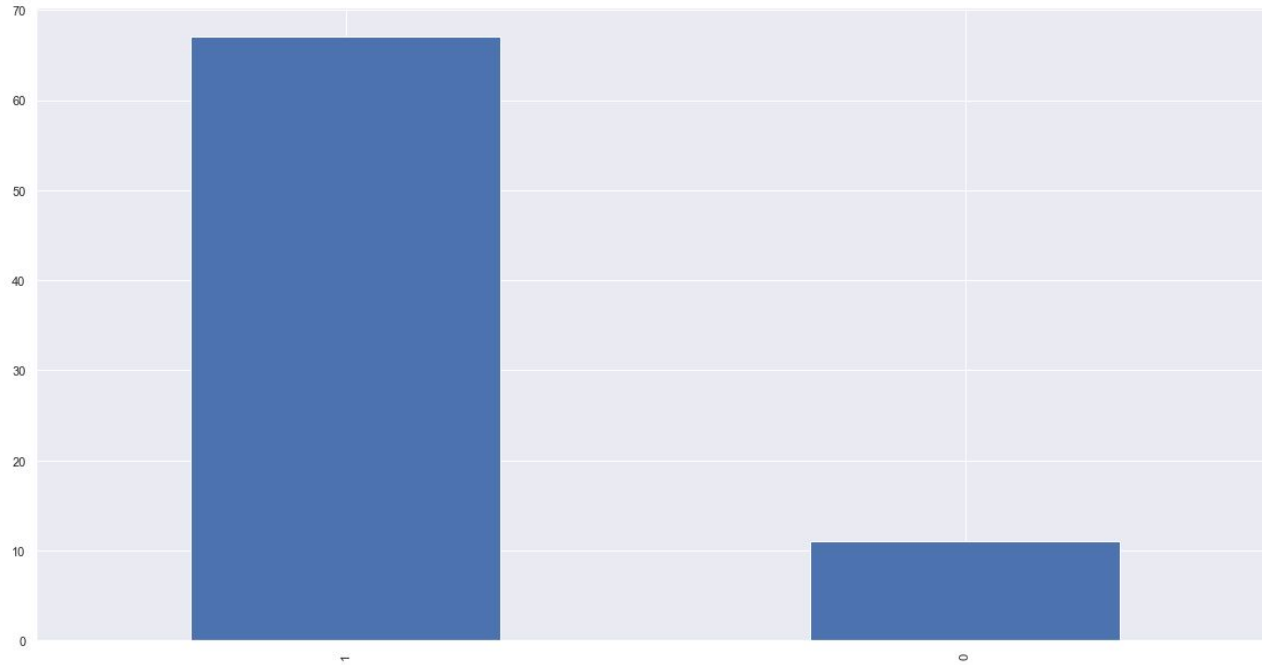
Count : 78

Mean Reads : 398

Max Reads : 936

# EDA : journal

# EDA : Popular Topics

# EDA : Full Text availability

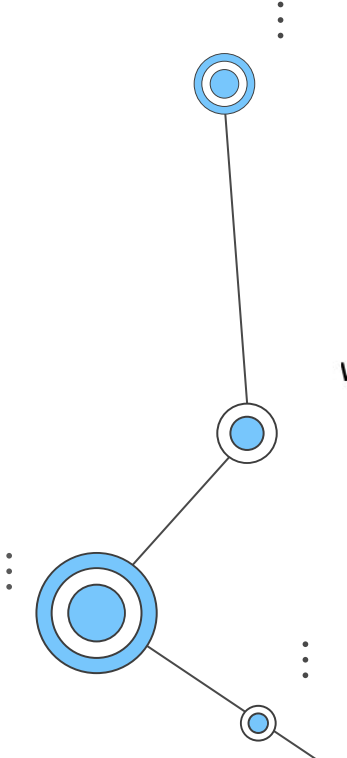# Experiments

# Uninclude Features

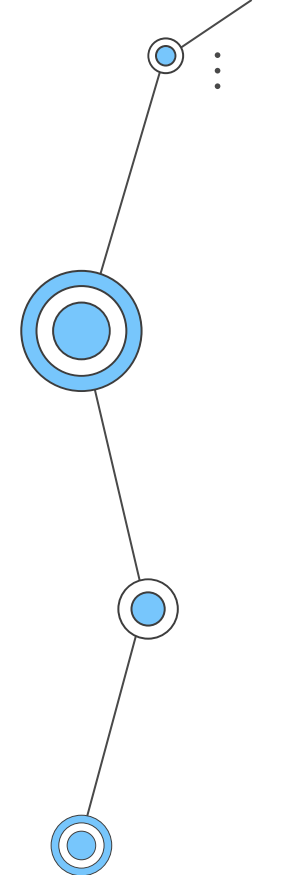**Research Interest**

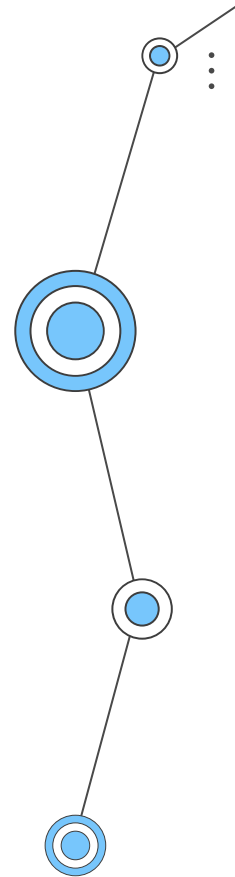Calculated from other features
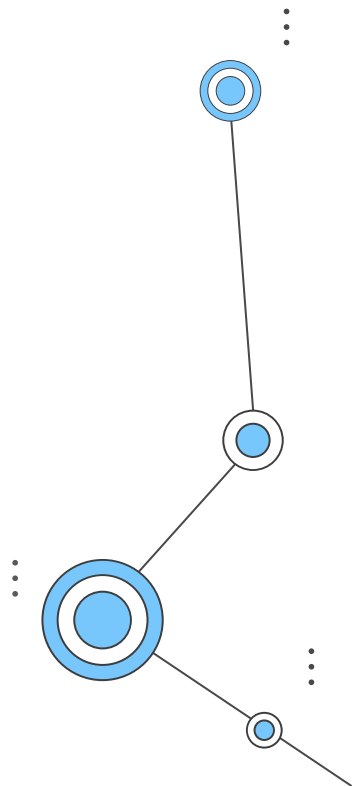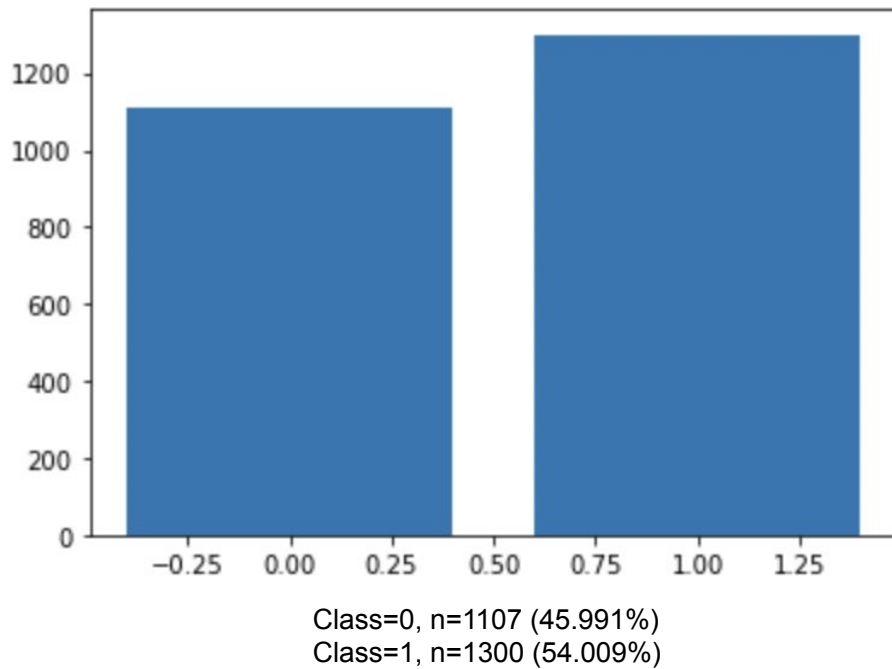
**Citation**

Popular articles are usually cited

# Result Before Balancing Techniques

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0            | 0.96      | 1.00   | 0.98     | 182     |
| 1            | 0.00      | 0.00   | 0.00     | 8       |
| accuracy     |           |        | 0.96     | 190     |
| macro avg    | 0.48      | 0.50   | 0.49     | 190     |
| weighted avg | 0.92      | 0.96   | 0.94     | 190     |

# balanced data : BorderLine SMOTE



Class=0, n=1107 (45.991%)
Class=1, n=1300 (54.009%)

# Results

Train Accuracy

0.89

0.87

0.90

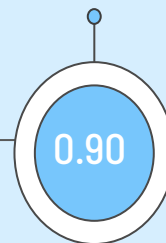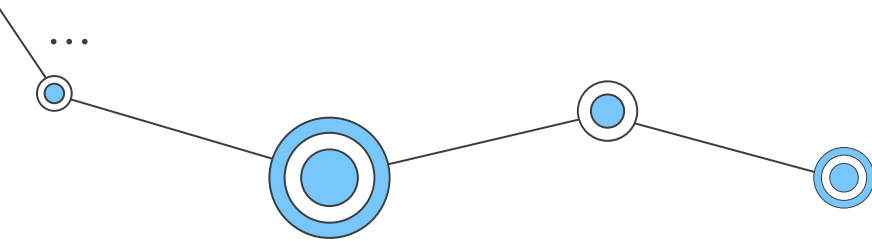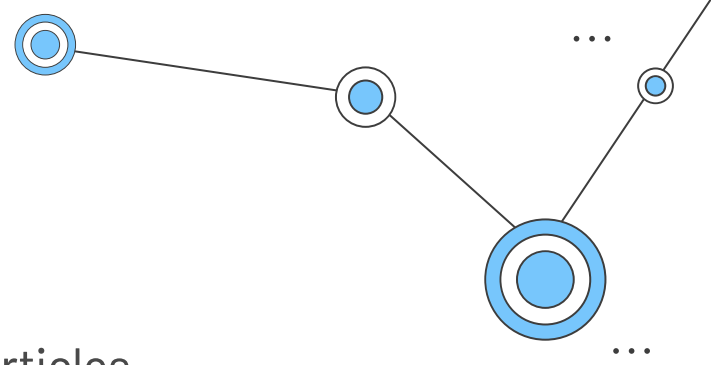Test Accuracy

Validation Accuracy

# Conclusion

We found 9 main topics in the articles, predicting the reads with the available data isn't possible, more data needs to be collected, but it's possible to classify articles according to popularity.

# Thanks!

Do you have any questions?