# Transfer Learning for Stance Analysis in COVID-19 Tweets

Erik Tjong Kim Sang, Marijn Schraagen, Shihan Wang, Mehdi Dastani
Netherlands eScience Center / Utrecht University

CLIN 31 - 9 July 2021

netherlands eScience center

# RIVM: Applying behavioural science to COVID-19

National Institute for Public Health
and the Environment
Ministry of Health, Welfare and Sport

**RIVM** Committed to *health and sustainability*

🏠  Topics    About RIVM    Publications    International    Contact

Search

Home › COVID 19 › Research on COVID 19 in the Netherlands › Applying behavioural science to COVID 19

## Applying behavioural science to COVID-19

Study on behavioural measures and well-being
during the COVID-19 pandemic

Preventive behaviour plays an important role in working together to gain control of the
coronavirus SARS-CoV-2. Healthy behaviour is also vital in staying healthy during the
coronavirus pandemic. Research on behaviour and health provides insights on how to help
people keep following behavioural rules – with a focus on their own health and the people
around them.

## Study on behavioural measures and well-being

The measures taken by the government in the fight against the coronavirus have a major impact on the daily lives of
everyone in the Netherlands. The government would like to know whether people can follow these rules, and what
they think of them.

## Table of contents

## Verandering in het houden aan de (basis) gedragsregels
### Meting 1 t/m 13



- — hoesten/niezen in elleboog (1)
- — geen handen schudden (2)
- — mondkapje in openbaar vervoer (2)
- — mondkapje in publieke binnenruimtes (2)
- — niet op drukke plek geweest of elke keer omgekeer...
- — testen bij klachten (2)
- — thuisblijven bij klachten (2)
- — handen wassen als het nodig is (1)
- — voldoende afstand houden van anderen (1)
- — thuiswerken (3)
- — maximaal aantal bezoekers thuis (2)

**Timeline**
R1 17 - 24 April 2020
R2 7 - 12 May 2020
R3 27 May - 1 June 2020
R4 17 - 21 June 2020
R5 8 - 12 July 2020
R6 19 - 13 August 2020
R7 30 September – 4 October 2020
R8 11 - 15 November 2020
R9 30 December – 3 January 2021
R10 10 - 14 February 2021
R11 24 - 28 March 2021
R12 5 – 9 May 2021
R13 16 - 20 June 2021

**Translation**
- sneezing in elbow (1)
- not shaking hands (2)
- facemask in public transport (2)
- Facemask in public spaces (2)
- Avoid busy areas
- Test with symptoms (2)
- Stay at home with symptoms (2)
- Wash hands when necessary (1)
- Social distancing (1)
- Work at home (3)
- Maximum of home vistors (2)

(1) % of times; (2) % of participants; (3) % of work hours

https://www.rivm.nl/gedragsonderzoek/maatregelen-welbevinden/naleven-gedragsregels

## Goal:

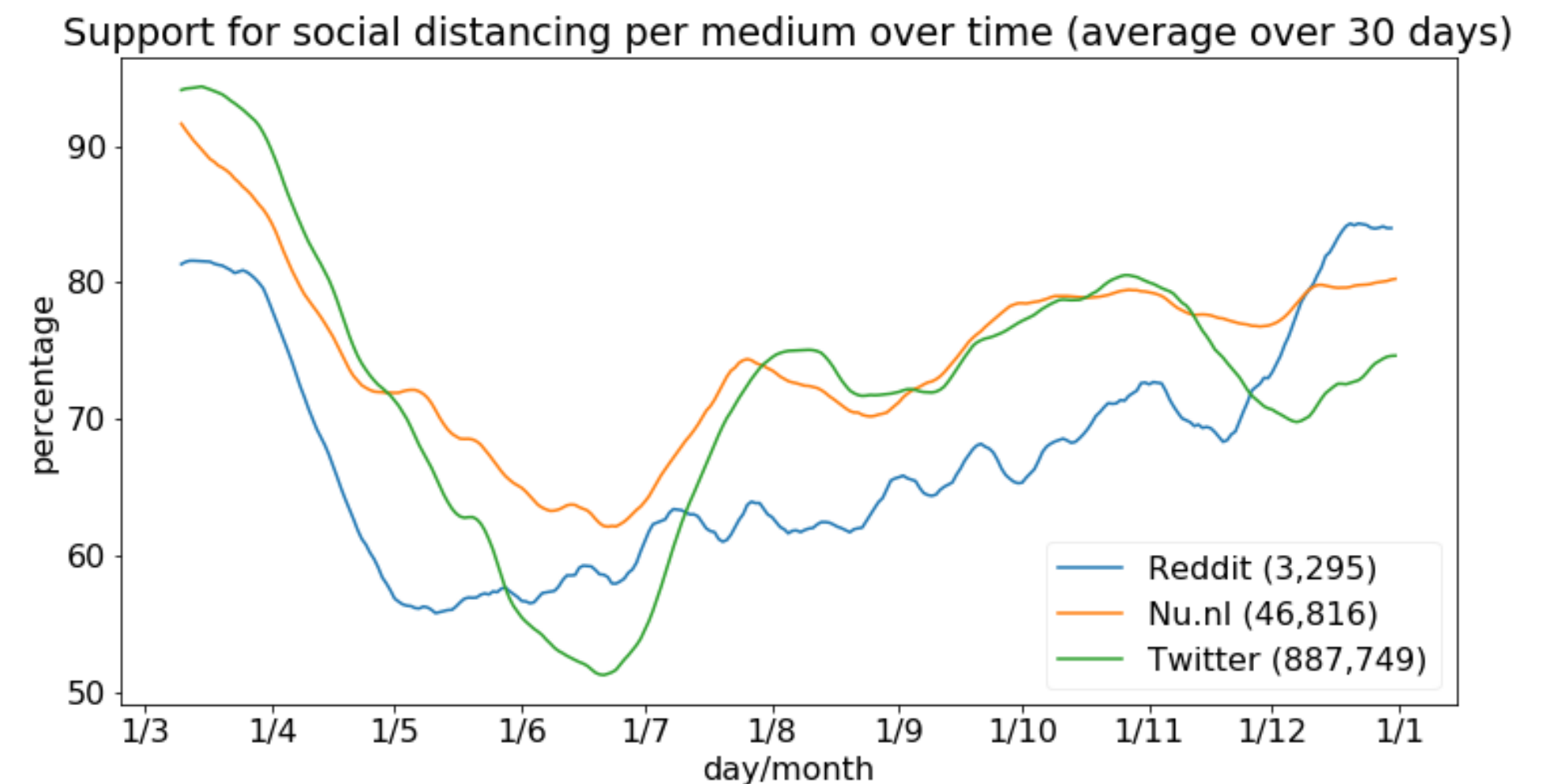- Derive stances on Dutch pandemic measures from social media data

## Method

1. Select relevant data with keyword search

2. Estimate stances in data with machine learning

3. Visualize results

## Challenge

- We need annotated data on many topics



Support for social distancing per medium over time (average over 30 days)

Reddit (3,295)
Nu.nl (46,816)
Twitter (887,749)

More information: https://github.com/puregome/notebooks

We tested different approaches of transfer learning / domain adaptation:

TGTONLY: build models from (limited) in-domain data only (baseline)

SRCONLY: build models from out-of-domain data only

ALL: build models from both in-domain and out-of-domain data

RLVONLY: like ALL but only use relevant out-of-domain data

FEATAUG: like RLVONLY but use different features for in-domain and out-of-domain data (Daumé III, 2007)
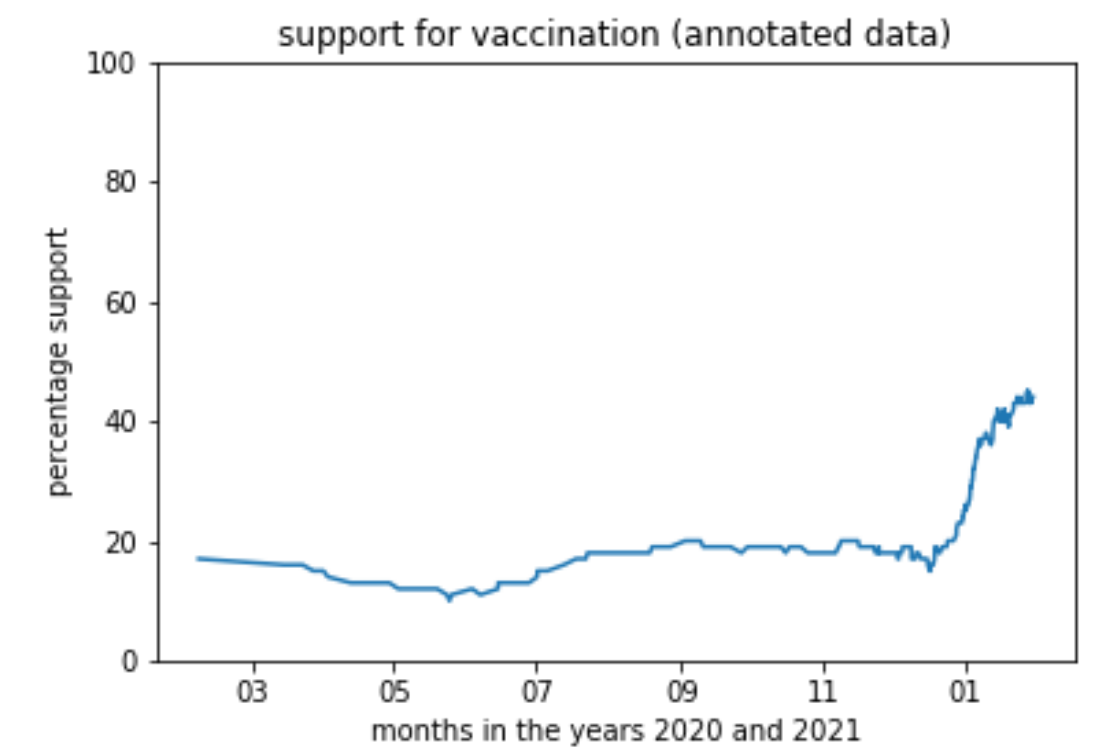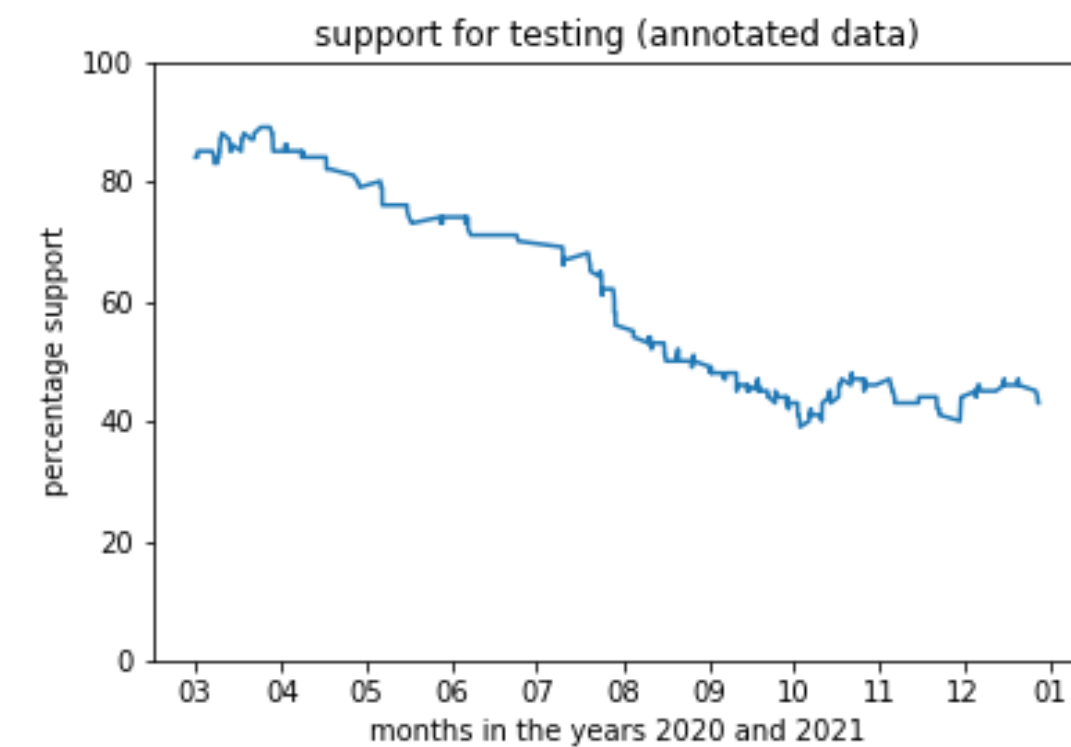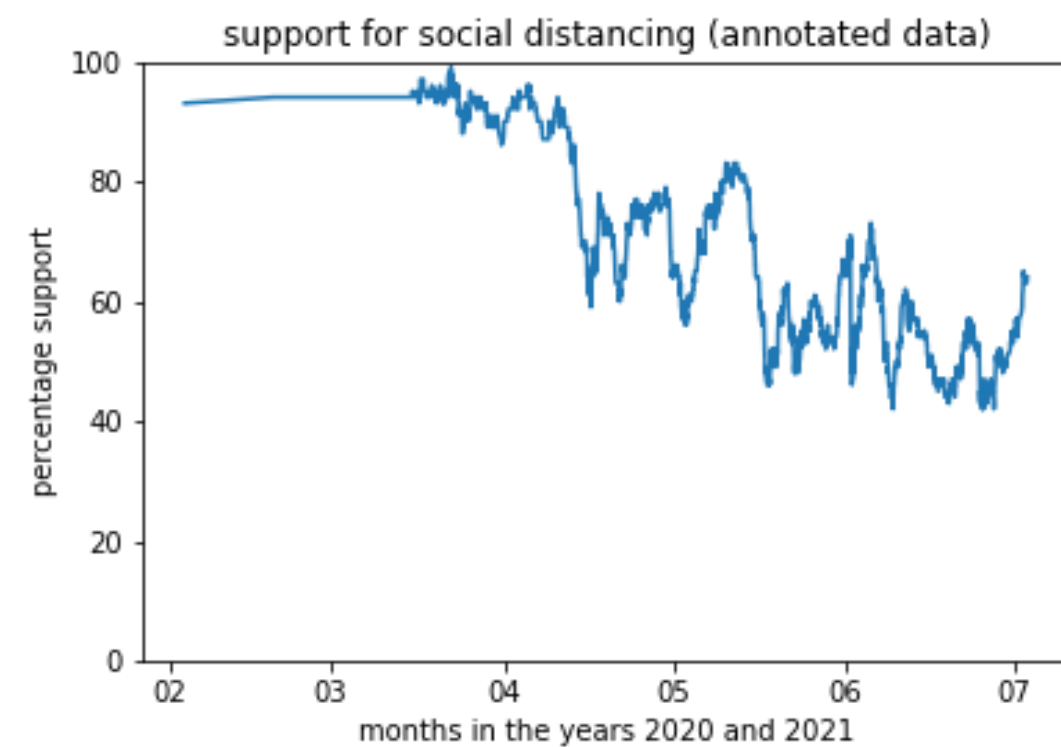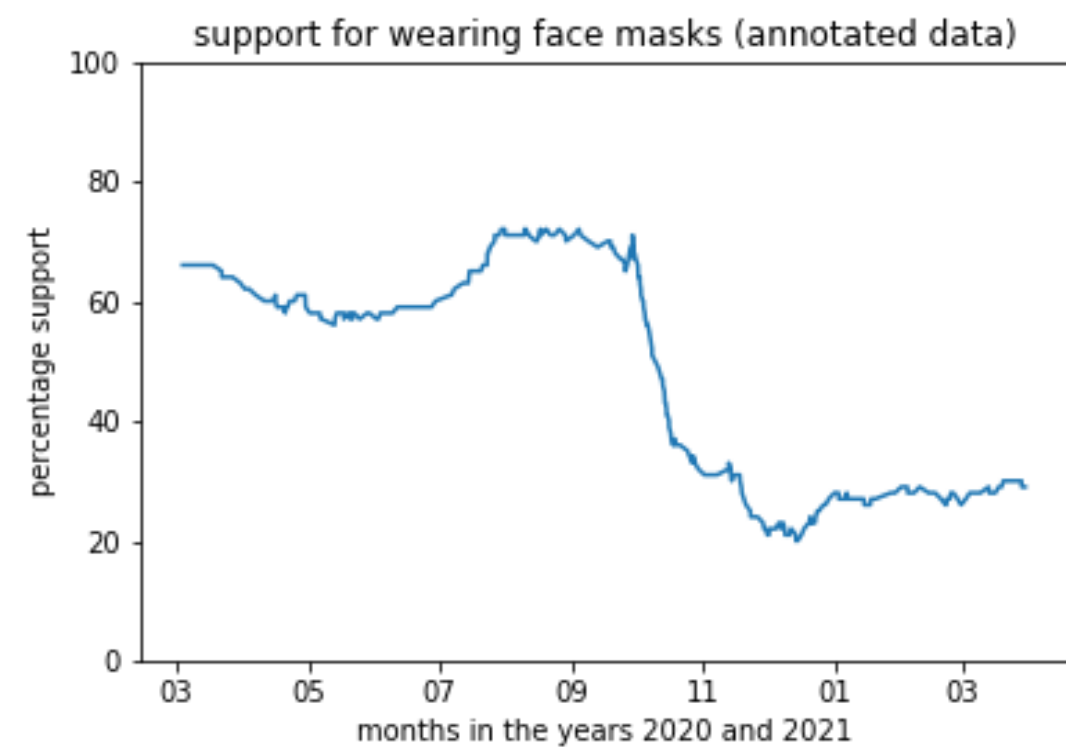
## Data, topics and keyword filtering

| Topic | Number of tweets | Time frame |
|---|---|---|
| Dutch tweets | 346,672,687 | February 2020 – June 2021 |
| Pandemic | 41,175,162 | February 2020 – June 2021 |
| Face masks | 1,847,245 | February 2020 – June 2021 |
| Social distancing | 1,441,262 | February 2020 – June 2021 |
| Testing | 3,831,931 | February 2020 – June 2021 |
| Vaccination | 7,643,870 | February 2020 – June 2021 |

| Topic | Keyword filter |
|---|---|
| Face masks | mondkapje |
| Social distancing | 1[.,]5[ -]*m | afstand.*hou | hou.*afstand | anderhalve[ -]*meter |
| Testing | \btest | getest | sneltest | pcr |
| Vaccination | vaccin | ingeënt | ingeent | inent | prik | spuit | bijwerking | --> | (emoji) | pfizer | moderna | astrazeneca | astra | zeneca | novavax | biontec |

In the keyword filters, "|" stands for OR and "\b" represents a word boundary

**Data annotation**

| Topic | Tweets | Time frame | Supports | Rejects | Irrelevant |
|---|---|---|---|---|---|
| Face masks | 1,011 | March 2020 - March 2021 | 21.2% | 23.6% | 55.2% |
| Social distancing | 5,977 | February  2020 – July 2020 | 56.2% | 20.0% | 23.8% |
| Testing | 1,181 | March 2020 – December 2020 | 28.2% | 17.8% | 53.4% |
| Vaccination | 1,007 | January 2020 - January 2021 | 9.6% | 24.3% | 66.1% |

For text classification we used the system FastText: a linear classifier which represents texts as an average of word vectors:

- Training started from a language model developed from 230 million Dutch tweets from the year 2020

- A limited grid search was done for finding the best parameters for the social distancing data

- Non-default parameter values used: vector length: 300, learning rate: 0.3, number of epochs: 200

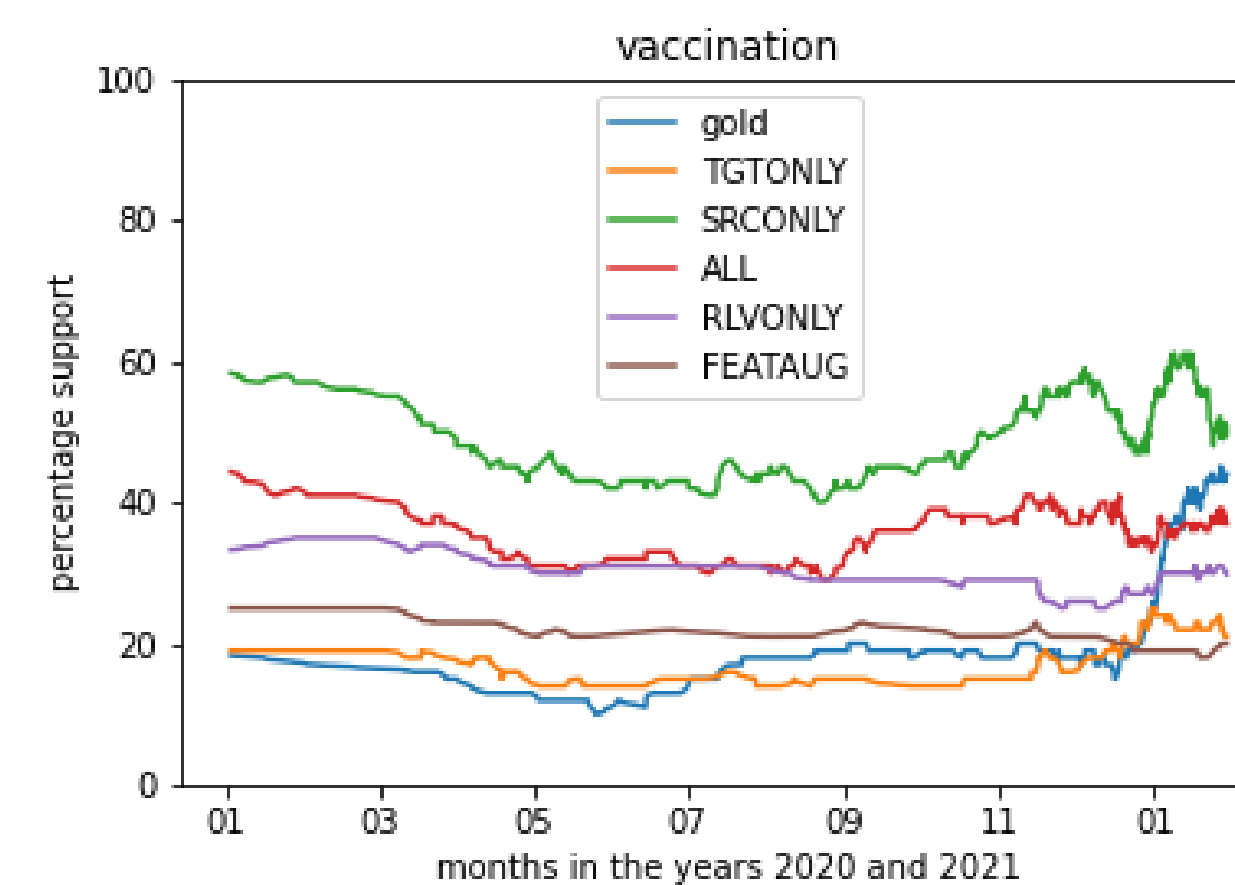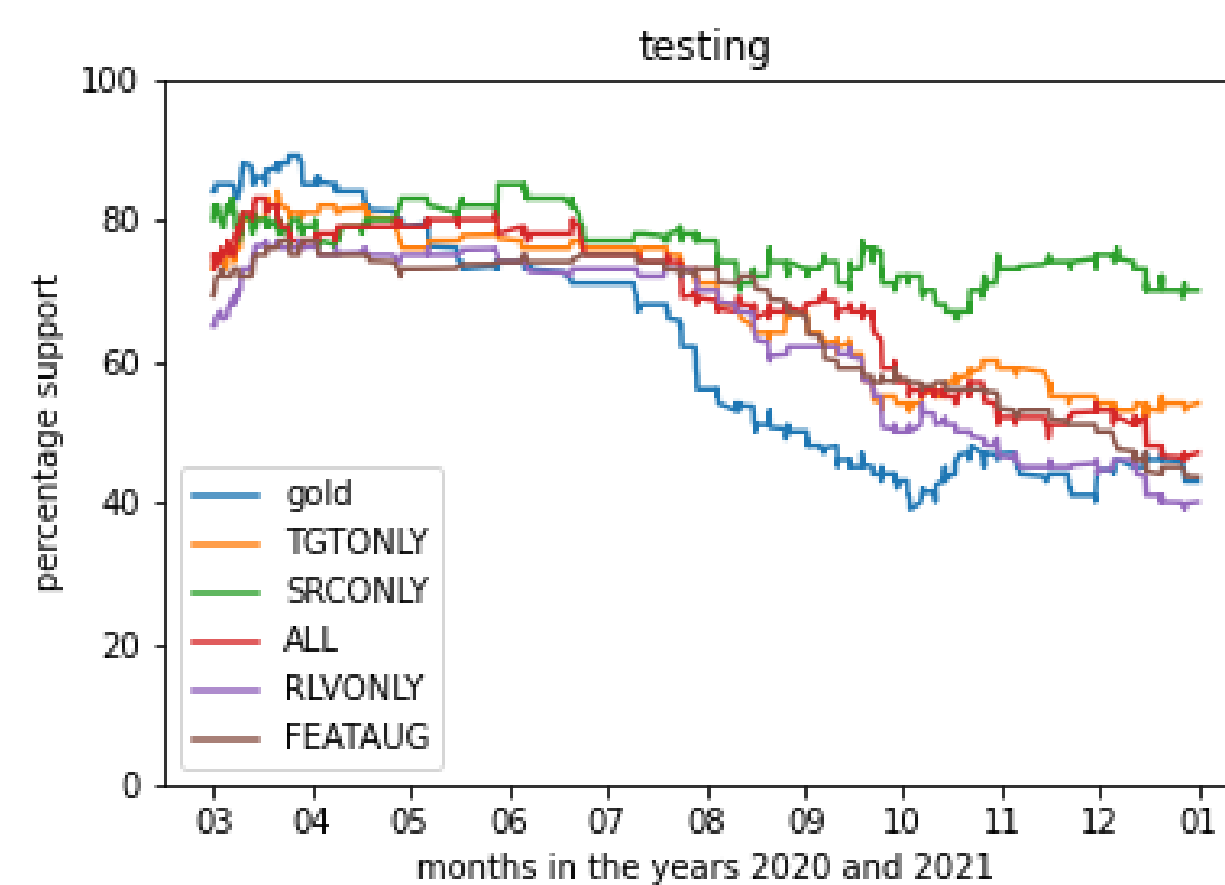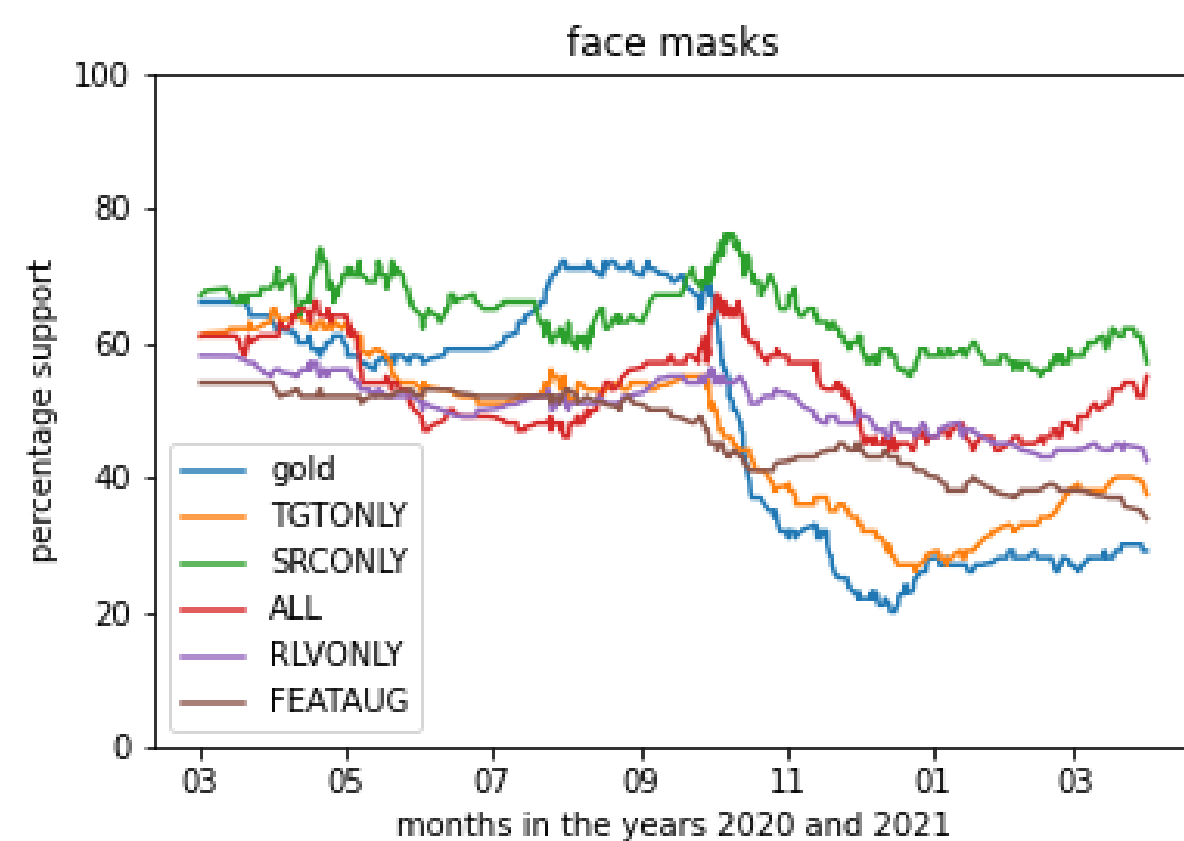**Transfer learning results for social distancing data as out-of-domain data**

| Accuracy | TGTONLY | SRCONLY | ALL | RLVONLY | FEATARG |
|---|---|---|---|---|---|
| Face masks | **0.592** | 0.414 | 0.543 | 0.569 | 0.583 |
| Testing | 0.561 | 0.452 | 0.541 | **0.572** | **0.572** |
| Vaccination | 0.608 | 0.488 | 0.572 | 0.585 | **0.616** |
| Social distancing | 0.656 | | | | |

- SRCONLY and ALL perform poorly, because of inclusion of the relevance task

- RLVONLY does better: (1) selecting relevant tweets with in-domain data and then (2) determining stance with all tweets

- FEATARG performs even better: in step (2) use separate features for in-domain and out-of-domain data

- Absolute accuracy gains are small

# Graph-based assessment of transfer learning (annotated data graphs used as gold data)

| Pearson r | TGTONLY | SRCONLY | ALL | RLVONLY | FEATARG |
|---|---|---|---|---|---|
| Face masks | **0.877** | 0.533 | 0.435 | 0.743 | 0.806 |
| Testing | **0.949** | 0.819 | 0.882 | 0.856 | 0.848 |
| Vaccination | 0.361 | 0.586 | **0.592** | 0.528 | -0.824 |

| Absolute difference | TGTONLY | SRCONLY | ALL | RLVONLY | FEATARG |
|---|---|---|---|---|---|
| Face masks | **0.080** | 0.177 | 0.143 | 0.133 | 0.125 |
| Testing | 0.093 | 0.169 | 0.097 | **0.087** | 0.102 |
| Vaccination | **0.078** | 0.260 | 0.149 | 0.107 | 0.102 |

**Concluding remarks and future work**

- We evaluated four transfer learning methods for predicting COVID-19 measure stances from social media data

- None of the methods consistently outperformed the baseline: using only in-domain data

- The problem with reproducing the vaccination graph suggests that it is important to have annotated in-domain data for all available time periods

- Future work: repeat the experiments with BERT as machine learner