# The Burden of Being a Bridge: Analysing Subjective Well-Being of Twitter Users during the COVID-19 Pandemic⋆

Ninghan Chen[1], Xihui Chen[2], Zhiqiang Zhong[1], Jun Pang[1,2] ✉

[1] Faculty of Science, Technology and Medicine,
University of Luxembourg, Esch-sur-Alzette, Luxembourg
[2] Interdisciplinary Centre for Security, Reliability and Trust,
University of Luxembourg, Esch-sur-Alzette, Luxembourg
{ninghan.chen, xihui.chen, zhiqiang.zhong, jun.pang}@uni.lu

**Abstract.** The outbreak of the COVID-19 pandemic triggers *infodemic* over online social media, which significantly impacts public health around the world, both physically and psychologically. In this paper, we study the impact of the pandemic on the mental health of influential social media users, whose sharing behaviours significantly promote the diffusion of COVID-19 related information. Specifically, we focus on subjective well-being (SWB), and analyse whether SWB changes have a relationship with their *bridging performance* in information diffusion, which measures the speed and wideness gain of information transmission due to their sharing. We accurately capture users' bridging performance by proposing a new measurement. Benefiting from deep-learning natural language processing models, we quantify social media users' SWB from their textual posts. With the data collected from Twitter for almost two years, we reveal the greater mental suffering of influential users during the COVID-19 pandemic. Through comprehensive hierarchical multiple regression analysis, we are the first to discover the strong relationship between social users' SWB and their bridging performance.

**Keywords:** Subjective well-being · COVID-19 · Information diffusion

## 1 Introduction

Since its outbreak, COVID-19 has become an unprecedented global health crisis and incited a worldwide *infodemic*. The term "infodemic" outlines the perils of misinformation during disease outbreaks mainly on social media [7, 15]. Apart from accelerating virus transmission by distracting social reactions, the infodemic increases cases of psychological diseases such as anxiety, phobia and depression during the pandemic [10]. As a result, the infodemic impairs the UN's

sustainable development goals (SDGs), especially SDG3 which aims to promote mental health and well-being.

To combat infodemic, both governments and healthcare bodies have launched a series of social media campaigns to diffuse trustworthy information. To amplify the speed and wideness of information spread, users with a large number of followers are invited to help share messages [33, 1]. Healthcare professionals and social activists also voluntarily and actively participate in relaying information they deem as useful with their social media accounts. All these people actually play a bridging role on social media delivering information to the public, although their *bridging performance* differs. We use bridging performance as an analogy to estimate how efficient and wide information can spread across social media due to the sharing of a user.

Subjective well-being (SWB), one important indicator of SDG3, evaluates individuals' cognitive (e.g., life satisfaction) and affective (i.e., positive and negative) perceptions of their lives [19]. Since the onset of the COVID-19 pandemic, the decrease of SWB has been unanimously recognised across the world. With studies for various sub-populations [17, 12], many factors have been discovered correlating to SWB changes such as professions, immigration status and gender. In this paper, we concentrate on influential social media users who play the bridging role in diffusing COVID-19 information, and study the impact of the pandemic on their SWB. We further examine whether their active participation in diffusing COVID-19 information is a predictor of the SWB changes. To the best of our knowledge, we are the first to study the mental health of this specific group of people during the pandemic.

We identify two main challenges to overcome before conducting our analysis. First, there are no measurements available that can accurately quantify users' real bridging performance in diffusing COVID-19 related information. The measurements, widely used in crisis communications and online marketing, rely on social connections, and have been found insufficient in capturing users' actual bridging performance, especially in such a global health crisis [27]. For instance, although some healthcare professionals are not super tweeters with thousands of followers, their professional endorsement significantly promotes the popularity of the posts they retweeted [27]. The second challenge is the access to the SWB levels of a large number of social media users whose bridging performance is simultaneously available.

In this paper, we take advantage of the information outbreak on social media incurred by the COVID-19 pandemic and the advances of artificial intelligence to address the two challenges. For the first challenge, we propose a new bridging performance measurement based on *information cascades* [29] which abstract both information spread processes and social connections. To address the second challenge, we leverage the success of deep learning in Natural Language Processing (NLP) and estimate users' SWB by referring to the sentiments expressed in their textual posts. In spite of the inherent biases, the power of social media posts has been shown in recent studies [19] for robust extraction of well-being with supervised data-driven methods. In this paper, instead of manually con-

structed features, we use the state-of-the-art transformer-based text embedding to automatically learn the representative features of textual posts.

*Our contributions.* We collect data from Twitter generated from *the Greater Region of Luxembourg* (GR). GR is a cross-border region centred around Luxembourg and composed of adjacent regions of Belgium, Germany and France. One important reason to select this region is its intense inter-connections of international residents from various cultures, which is unique as a global financial centre. Moreover, they well represent the first batch of countries administering COVID vaccines. Our collection spans from October 2019 to the end of 2021 for over 2 years, including 3 months before the outbreak of the COVID-19 pandemic. Our contributions are summarised as follows:

− We propose a new measurement to capture the actual bridging performance of individual users in diffusing COVID-19 related information. Compared to existing social connection-based measurements, it is directly derived from information diffusion history. Through manual analysis of the collected dataset, our measurement allows for identifying the accounts of influential health professionals and volunteers that are missed before in addition to super tweeters.
− Through deep learning-based text embedding methods, we implement a classification model which can accurately extract the sentiments expressed in social media messages. With the sentiments of posts, we quantitatively estimate individual users' SWB, and confirm the greater suffering of influential users in their SWB during the pandemic.
− Through the hierarchical multiple regression model, we reveal that users' SWB has a strong negative relationship with their bridging performance in COVID-19 information diffusion, but weak relationship with their social connections.

Our research provides policy makers with an effective method to identify influential users in the fight against infodemic. Moreover, we contribute to the realisation of SDG3 by highlighting the necessity to pay special attention to the mental well-being of people who actively participate in transmitting information in health crises like COVID-19.

## 2   Related Work

### 2.1   Measuring bridging performance

A considerable amount of literature has been published quantifying users' bridging performance based on social connections to identify amplifiers in social media. We can divide the measurements into two types. The first type of measurements implicitly assume that influential users are likely to hold certain topology properties on social networks such as large degrees, strong betweenness centrality or community centrality [14]. The second type of measurements assume that influential users tend to be more likely reachable from other users through random walks. PageRank [25] and its variant TwitterRank [30] among the representative benchmarks of this type of measurements. PageRank is calculated only with

network structures while TwitterRank additionally takes into account topic similarities between users. All the two categories of measurements have been widely applied in practice, from public health crisis communication [23] to online marketing [21]. However, recent studies pointed out that they may not truly capture users' actual bridging performance in information diffusion during a specific public healthy crisis [23]. Although new measurements are proposed by extending existing ones with fusion indicators, their poor efficiency prevents them from being applied to real-world large-scale networks like Twitter and Facebook [18].

### 2.2   Subjective well-being extraction

Subjective well-being is used to measure how people subjectively rate their lives both in the present and in the near future [9]. Many methods have been used to assess subjective well-being, from traditional self-reporting methods [8] to the recent ones exploiting social media [32].Studies have cross-validated SWB extracted from social media data with the Gallup-Sharecare Well-Being Index survey,[1] a classic reference used to investigate public SWB, and found that SWB extracted from social media is a reliable indicator of SWB [19]. Twitter-based studies usually calculate SWB as the overall scores of positive or negative emotions (i.e., sentiment or valence) [19]. Sentiment analysis has developed from the original lexicon-based approaches [3] to the data-driven ones which ensure better performance [19]. We adopt the recent advances of the latter approaches, and make use of the pre-trained XLM-RoBERTa [24], a variant of RoBERTa [22], to automatically learn the linguistic representation of textual posts. As a deep learning model, RoBERTa and its variants have been shown to overwhelm traditional machine learning models in capturing the linguistic patterns of multilingual texts [2].

## 3   The GR-ego Twitter Dataset

In this section, we describe how we build our Twitter dataset, referred to as *GR-ego*. In addition to its large number of active users, we have another two considerations to select Twitter as our data source. First, the geographical addresses of posters are attached with tweets and thus can be used to locate users. Second, tweet status indicates whether a tweet is retweeted. If a tweet is retweeted, the corresponding original tweet ID is provided. Together with the time stamps, we can track the diffusion process of an original tweet. Our GR-ego dataset consists of two components: (i) the social network of GR users recording their following relations; (ii) the tweets posted or retweeted by GR users during the pandemic. We follow three sequential steps to collect our GR-ego dataset. Table 1 summarises its main statistics.
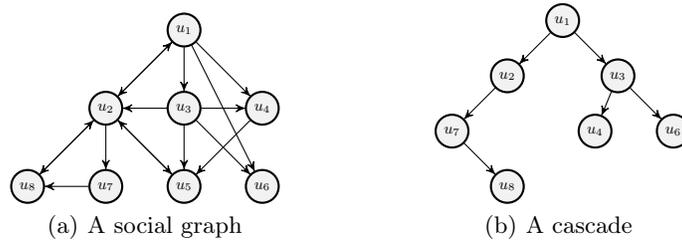
---

[1] `https://www.gallup.com/175196/gallup-healthways-index-methodology.aspx`

**Table 1.** Statistics of the GR-ego dataset.

| Social network | #node | 5,808,938 |
|---|---|---|
| | #edge | 12,511,698 |
| | Average degree | 2.15 |
| Timeline tweets | #user | 14,756 |
| | #tweet before COVID | 5,661,949 |
| | #tweet during COVID | 18,523,099 |
| | #tweet per user before COVID | 388.44 |
| | #tweet per user during COVID | 1255.29 |

*Step 1. Meta data collection.* Our purpose of this step is to collect seed users in GR who actively participate in COVID-19 discussions. Instead of directly searching by COVID-19 related keywords, we make use of a publicly available dataset of COVID-19 related tweets for the purpose of efficiency [4]. Restricted by Twitter's privacy policies, this dataset only consists of tweet IDs. We extract the tweet IDs posted between October 22nd, 2019, which is about three months before the start of the COVID-19 pandemic, and December 31st, 2021. Then with these IDS, we download the corresponding tweet content. On Twitter, geographical information, i.e., the locations of tweet posters and original users if tweets are re-tweeted, is either maintained by Twitter users, or provided directly by their positioning devices. We stick to the device-input positions, and only use user-maintained ones when such positions are unavailable. Due to the ambiguity of user-maintained positions, we leverage the geocoding APIs, Geopy and ArcGis Geocoding to regularise them into machine-parsable locations. With regularised locations, we filter the crawled tweets and only retain those from GR. In total, we obtain 128,310 tweets from 8,872 GR users.

*Step 2. Social network construction.* In this step, we search GR users from the seed users and construct the GR-ego social network. We adopt an iterative approach to gradually enrich the social network. For each seed user, we obtain his/her followers and only retain those who have a mutual following relation with the seed user, because such users are more likely to reside in GR [6]. We then extract new users' locations from their profile data and regularise them. Only users from GR are added to the social network as new nodes. New edges are added if there exist users in the network with following relationships with the newly added users. After the first round, we continue going through the newly added users by adding their mutually followed friends that do not exist in the current social network. This process continues until no new users can be added. Our collection takes 5 iterations before termination. In the end, we take the largest weakly connected component as the *GR-ego social network*.

*Step 3. COVID-19 related timeline tweets crawling.* In this step, we collect tweets originally posted or re-tweeted by the users in our dataset. These tweets will be used to extract users' SWB. Thus, the collected tweets are *not limited* to those relevant to the COVID-19 pandemic. Due to the constraints of Twitter, it is not

(a) A social graph          (b) A cascade

**Fig. 1.** Example of a cascade.

tractable to download all the users' past tweets. We select a sufficiently large number of representative users who actively participated in retweeting COVID-19 related messages, and then crawl their history tweets. In detail, we choose 14,756 users who (re)tweeted at least three COVID-19 related messages. With the newly released Twitter API which allows for downloading 500 tweets of any given month for each user, we collect $37,281,824$ tweets spanning between October 22nd, 2019 and December 31st, 2021. This period also contains the last three months before the pandemic is officially claimed. We release the IDs of our collected tweets via Github.[2]

## 4    Data Processing

### 4.1    Cascade computation

A cascade records the process of the diffusion of a message. It stores all activated users and the time when they are activated. In our dataset, a user is activated in diffusing a message when he/she retweets the message. In this paper, we adopt the widely accepted *cascade tree* to represent the cascade of a message [29, 6, 5].

The first user who posted the message is regarded as the root of the cascade tree. Users who retweeted the message, but received no further retweeting comprise the leaf nodes. Note that a tweet with the quotation to another tweet is also considered as a retweet of the quoted message. An edge from $u$ to $u'$ is added to the cascade if $u'$ follows $u$ and $u'$ re-tweeted the message after $u$, indicating $u$ activated $u'$. If many of the users who $u'$ follows ever retweeted the message, meaning $u'$ may be activated by any of them, we select the one who lastly retweeted as the parent node of $u'$. Figure 1(b) shows a cascade of the social network in Figure 1(a). In this example, user $u_4$ can be activated by the messages retweeted by either $u_1$ or $u_3$. Since $u_3$ retweeted after $u_1$, we add the edge from $u_3$ to $u_4$ indicating that the retweeting of $u_3$ activated $u_4$.

We denote the root node of a cascade $C$ by $r(C)$. We call a path that connects the root and a leaf node a *cascade path*, which is actually a sequence of nodes ordered by their activation time. For instance, $(u_1, u_3, u_4)$ is a cascade path in our example indicating that the diffusion of a message started from $u_1$ and

---

[2] https://github.com/NinghanC/SWB4Twitter

reached $u_4$ in the end through $u_3$. In this paper, we represent a cascade tree as a set of cascade paths. For instance, the cascade in Figure 1(b) is represented by the following set $\{(u_1, u_2, u_7, u_8), (u_1, u_3, u_4), (u_1, u_3, u_6)\}$.
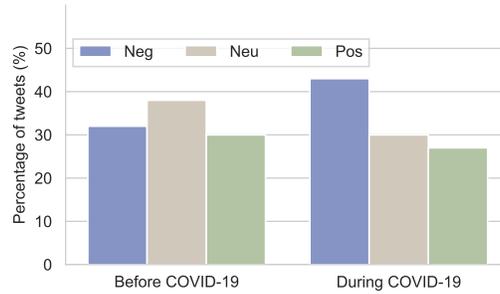
For our study, we follow the method in [20] to construct tweet cascades. Recall that when a tweet's status is '*Retweeted*', the ID number of the original tweet is also recorded. We first create a set of original tweets with all the ones labelled in our meta data as '*Original*'. Second, for each original tweet, we collect the IDs of users that have retweeted the message. At last, we generate the cascade for every original tweet based on the following relations in our GR-ego social network and their retweeting time stamps. We eliminate cascades with only two users where messages are just retweeted once. In total, 614,926 cascades are built and the average size of these cascades is 7.13.

## 4.2  Sentiment analysis

Previous works [34] leverage user-provided mood (e.g., angry, excited) or status to extract users' sentiment (i.e., positive or negative) and use them to approximately estimate affective subjective well-being. However, such information is not available on Twitter. We refer to the sentiments expressed in textual posts to extract users' SWB. In this paper, we treat sentiment extraction as a tri-polarity sentiment analysis for short texts, and classify a tweet as *negative*, *neutral* or *positive*. In order to deal with the multilingualism of our dataset, we benefit from the advantages of deep learning in sentiment analysis [2], and build an end-to-end deep learning model to conduct the classification. Our model is composed of three components. The first component uses a pre-trained multilingual language model, i.e., XLM-RoBERTa [24], to calculate the representation of tweets. The representations are then sent to the second component, a fully-connected ReLU layer with dropout. The last component is a linear layer added on the top of the second component's outputs with sigmoid as the activation function. We use cross-entropy as the loss function and optimise it with the Adam optimiser.

*Model training and testing.* We train our model on the *SemEval-2017 Task 4A* dataset [26], which has been used for sentiment analysis on COVID-19 related messages [11]. The dataset contains 49,686 messages which are annotated with one of the three labels, i.e., positive, negative and neutral. We shuffle the dataset and take the first 80% for training and the rest 20% for testing. We assign other training parameters following the common principles in existing works. We run 10 epochs with the maximum string length as 128 and dropout ratio as 0.5. When tested with macro-average F1 score and accuracy metrics, we achieve an accuracy of 70.09% and macro-average F1 score of 71.31%.

Despite its effectiveness on classifying *SemEval-2017 Task 4A* data, in order to check whether such performance will persist on our GR-ego dataset, we construct a new testing dataset. This dataset consists of 500 messages, 100 for each of the top 5 most popular languages. We hire two annotators to manually label the selected tweets and the annotated labels are consistent between them

**Fig. 2.** Sentiment distribution of users' timeline tweets.

with Cohen's Kappa coefficient $k = 0.93$. When applied on this new manually annotated dataset, our trained model achieves a similar accuracy of about 87%.

*Analysing our GR-ego dataset.* Before applying our sentiment classification model on our GR-ego dataset, we clean tweet contents by removing all URLs, and mentioned usernames. Figure 2 summarises the statistics obtained from user timeline tweets before and during the pandemic. The numbers of users' timeline tweets are consistent with previous studies. For instance, users tend to become more negative after the outbreak of the COVID-19 pandemic [12, 17].

## 5    Bridging Performance of Users in Information Diffusion

We devote this section to addressing the first challenge regarding identifying users that play the bridging role in transmitting COVID-19 related information.

### 5.1    Measuring user bridging performance

We evaluate users' overall performance in the diffusion of observed COVID-19 related tweets. As a user can participate in diffusing a number of tweets, we first focus on her/his importance in the diffusion of one single tweet and then combine her/his importance in all tweets into one single measurement. We consider a user *more important* in diffusing a tweet when his/her retweeting behaviour activates more users, or leads to a given number of activated users with fewer subsequent retweets. In other words, a more important user promotes wider acceptance of the information or accelerates its propagation. Given a cascade path $S = (u_1, u_2, \ldots, u_n)$, we use $S^*(u_i)$ $(1 \leq i < n)$ to denote the subsequence composed of the nodes after $u_i$ (including $u_i$), i.e., $(u_i, u_{i+1}, \ldots, u_n)$. For any $u$ that does not exist in $S$, we have $S^*(u) = \varepsilon$ where $\varepsilon$ represents an empty sequence and its length $|\varepsilon| = 0$.

**Definition 1 (Cascade bridging value)** *Given a cascade tree $C$ and a user $u$ ($u \neq r(C)$), the cascade bridging value of $u$ in $C$ is calculated as:*

$$\alpha_C(u) = \left( \sum_{S \in C} \frac{\mid S^*(u) \mid}{\mid S \mid} \right) / |C|.$$

Note that our purpose is to evaluate the importance of users as transmitters of messages. Therefore, the concept of cascade bridging value is not applicable to the root user, i.e., the message originator.

*Example 1.* In Figure 1(b), $u_3$ participated in two cascade paths, i.e., $S_1 = (u_1, u_3, u_4)$ and $S_2 = (u_1, u_3, u_6)$. Thus, $S_1^* = (u_3, u_4)$ and $S_2^* = (u_3, u_6)$. We then have $\alpha_C(u_3) = \frac{2/3 + 2/3}{3} \approx 0.44$.

In Definition 1, we do not simply use the proportion of users activated by a user in a cascade to evaluate her/his bridging performance. This is because it only captures the number of activated users and ignores the speed of the diffusion. Taking $u_2$ in Figure 1(b) as an example, according to our definition, $\alpha_C(u_2) = 0.25$ which is smaller than $\alpha_C(u_3)$. This is due to the fact that $u_2$ activated two users through two retweets while $u_3$ only used one. However, if we only consider the proportion of activated users, the values of these two users will be the same.

With a user's bridging value calculated in each cascade, we define *user bridging magnitude* to evaluate her/his overall importance in the diffusion of a given set of observed messages. Intuitively, we first add up the bridging values of a user in all his/her participated cascades and then normalise the sum by the maximum number of cascades participated by a user. This method captures not only the bridging value of a user in each participated cascade, but also the number of cascades she/he participated in. This indicates that, a user who is more active in sharing COVID-19 related information is considered more important in information diffusion.

**Definition 2 (User bridging magnitude (UBM))** *Let $\mathcal{C}$ be a set of cascades on a social network and $\mathcal{U}$ be the set of users that participate in at least one cascade in $\mathcal{C}$. A user $u$'s user bridging magnitude (UBM) is calculated as:*

$$\omega_{\mathcal{C}}(u) = \frac{\sum_{C \in \mathcal{C}} \alpha_C(u)}{\max_{u' \in \mathcal{U}} |\{C \in \mathcal{C} | \alpha_C(u') > 0\}|}.$$
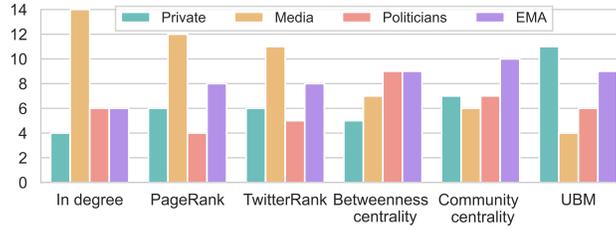
With this measurement, we can compare the UBM values of any two given users, and learn which one plays a more important role in information diffusion.

## 5.2   Validation of UBM

**Experimental results.** We compare the effectiveness of our UBM to five widely used topology-based measurements in the literature, i.e., in-degree, PageRank [25], TwitterRank [30], betweenness centrality [14] and community centrality [14]. We

**Table 2.** Comparison of bridging performance with benchmarks.

| | in-degree | PageRank | TwitterRank | Betweenness centrality | Community centrality | UBM |
|---|---|---|---|---|---|---|
| Avg. #activated user/minute | 0.042 | 0.057 | 0.064 | 0.043 | 0.056 | 0.104 |
| Avg. #activated users | 13.99 | 16.84 | 17.68 | 15.54 | 17.00 | 23.81 |
| %impacted user | 32.17 | 52.54 | 57.44 | 43.44 | 56.54 | 71.66 |



**Fig. 3.** Profile distribution of the top 30 accounts with highest bridging performance

randomly split the set of cascades into two subsets. The first subset accounts for 80% of the cascades and is used to calculate the bridging performance of all users. Then we select the top 20% users with the highest bridging performance in every adopted measurement and use the other subset to compare their actual influences in information diffusion. We adopt three measurements to quantitatively assess the effectiveness of UBM and the benchmarks. We use the *average number of activated users per minute* to evaluate the efficiency of the information diffusion. The more users activated in a minute, the faster information can be spread when it is shared by the influential users. The *average number of activated users* counts the users who received the information after the retweeting behaviour of an identified influential user. It is meant to evaluate the expected wideness of the spread once an influential user retweets a message. The *percentage of impacted users* gives the proportion of users that have ever received a message due to the sharing behaviours of identified influential users. This measurement is to compare the overall accumulated influence of all the selected influential users. We show the results of UBM and other benchmark measurements in Table 2. We can observe that it takes less time on average for the influential users identified according to UBM to activate an additional user, with 0.104 users activated a minute due to their retweets. With 23.81 users activated, UBM allows for finding the users whose retweeting action can reach more than 35% users than those identified by the benchmarks. In the end, the top 20% influential users identified by UBM spread their shared information to 71% users in our dataset, which overwhelms that of the best benchmark by about 15%. From the above analysis in terms of the three measurements, we can see that our UBM can successfully identify influential users whose sharing on social media manages to promote both the wideness and the speed of the diffusion of COVID-19 related information.

*Manual analysis.* In order to understand the profiles of the calculated influential users by the measurements, we select the top 30 users with the highest bridging performance of each measurement. We identify four types of user profiles: *private*, *media*, *politicians* and *emergency management agencies* (EMA). Figure 3 shows the distributions of their profiles. We can observe that the distributions vary due to the different semantics of social connections captured by the measurements. For instance, due to the large numbers of followers, Twitter accounts managed by traditional media are favoured by in-degree. This obviously underestimates the importance of accounts such as those of EMAs in publishing pandemic updates. With reachability and importance in connecting users and communities considered, more accounts of politicians and EMAs stand out. The proportion of private accounts also starts to increase. When UBM is applied, the percentage of private accounts becomes dominant. A closer check discovers that 10 out of the 11 private accounts belong to health professionals and celebrities. This is consistent with the literature [16] which highlights the importance of health professionals and individuals in broadcasting useful messages about preventive measures and healthcare suggestions in the pandemic.

## 6    Impact of COVID-19 on the SWB of Influential Users

### 6.1    Measuring SWB

We extend the definition proposed in [34] to measure the level of subjective well-being of users based on the sentiment expressed in their past tweets. Specifically, we extend it from bi-polarity labels, i.e., negative and positive affection, to tri-polarity with neutral sentiment by multiplying a scaling factor to simulate the trustworthiness of the bi-polarity SWB.

**Definition 3 (Social media Subjective well-being value (SWB))** *We use $N_p(u)$, $N_{neg}(u)$ and $N_{neu}(u)$ to denote the number of positive, negative and neutral posts of a user $u$, respectively. The subjective well-being value of $u$, denoted by $swb(u)$, is calculated as:*

$$\frac{N_p(u) - N_{neg}(u)}{N_p(u) + N_{neg}(u)} \cdot \left( \frac{N_p(u) + N_{neg}(u)}{N_p(u) + N_{neg}(u) + N_{neu}(u)} \right)^{\frac{1}{2}} .$$

If all messages are neutral, then $swb(u)$ is 0.

*Discussion.* Note that i) consistent with [34], we focus on affective SWB (i.e., positive and negative) in this paper, while ignoring its cognitive dimension; ii) users' SWB is evaluated based on their original messages: originally posted tweets and quotations; iii) for tweets with quotations to other messages, only the texts are considered without the quoted messages. As retweets may not explicitly include users' subjective opinions, we exclude them from the SWB calculation.
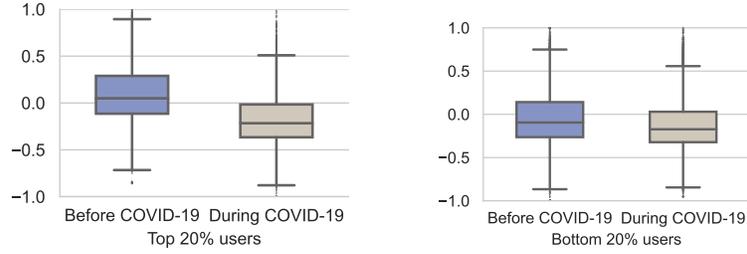
**Fig. 4.** SWB changes after the outbreak of the pandemic.

## 6.2 Analysing SWB changes of influential users

With the proposed SWB measurement, we study how users' subjective well-being changes due to the outbreak of the COVID-19 pandemic. We calculate the UBM values of the users in our collected dataset and order them descendingly. Then we select the top 20% users as well as the bottom 20% users and compare the two groups' responses to the pandemic. For each group, we calculate users' SWBs according to their posts before the pandemic and after the pandemic to capture the changes. Note that we only consider the users with more than 5 posts in each time period. In Figure 4, we show the SWB distributions of the two user groups. On average, the users with high UBM have positive SWB of 0.11 before the pandemic while the users with low UBM are negative. *The SWB of both user groups decreases after the pandemic but the SWB of the top 20% users drops more significantly.* Specifically, their SWB falls by 0.33, which is two times as much as that of the bottom 20% users. The lowest value of the top 20% users' SWB slightly decreases after the pandemic, while the lowest value of the bottom 20% of users does not change significantly. Note that the minimum values here do not include outliers that lie outside the box whiskers. This indicates that the top 20% users become even more negative than the bottom 20% users, in terms of mean and minimum values. To sum up, the pandemic causes more negative mental impacts on the social media users who play a more important bridging role in transmitting COVID-19 related information.

## 6.3 Relation between SWB and bridging performance

We conduct the first attempt to study if a user's bridging performance has a relationship with the SWB changes of the users actively participating the diffusion of COVID-19 related information. In addition to UBM and the five benchmark measurements used in Section 5.2, we consider two additional variables: *out-degree* and *activity*. Out-degree is used to check whether the number of accounts a user follows correlates with SWB changes. The activity variable evaluates how active a user is engaged in the online discourse and is quantified by the number of messages he/she posted. In order to isolate the impacts of these variables, we adopt the method of *hierarchical multiple regression* [28]. The intuitive idea is

to check whether the variables of interest can explain the SWB variance after accounting for some variables.

To check the validity of applying hierarchical multiple regression, we conduct first-line tests to ensure a sufficiently large sample size and independence between variables. We identify the variables corresponding to community centrality and TwitterRank fail to satisfy the multi-collinearity requirement. We thus ignore them in our analysis. The ratio of the number of variables to the sample size is 1:2108, which is well below the requirement of 1:15 [28]. This indicates the sample size is adequate. We iteratively input the variables into the model with three stages. The results are shown in Table 3. In the first stage, we input the variables related to network structures, i.e., in-degree, out-degree, Pagerank and Betweenness centrality. The combination of the variables can explain 4.30% of the SWB variance ($F = 4.379, p < 0.05$). Note that an F-value of greater than 4 indicates the linear equation can explain the relation between SWB and the variables. This demonstrates that there exists a positive relationship between the topology-based variables and SWB, but this relationship is rather weak. A closer check on the $t$-values show that out-degree is irrelevant to SWB and the rest three variables are weakly related. In the second stage, we add the variable of activity to the model. After controlling all the variables of the first stage, we observe that user activity does not significantly contribute to the model with $t$-value of 0.716. This suggests that user activity is not a predictor of SWB. In the third stage, we introduce UBM to the model. The addition of UBM, with the variables in the previous two stages controlled, reduces the R value from -0.219 to -0.579. UBM contributes significantly to the overall model with $F = 147.82$ ($p < 0.001$) and increases the predicted SWB variance by 23.8%. Together with the $t$-value of -11.469 ($p < 0.001$), we can see there exists a strong negative relation between UBM and SWB, and UBM is a strong predictor for SWB.

*Discussion.* The results illustrate that UBM is strongly related to SWB, while in-degree, Pagerank and betweenness centrality are weakly related. This difference further shows that UBM can more accurately capture users' behaviour changes after the outbreak of the pandemic while topology features remain similar to those before the pandemic. This may be explained by the recent studies [17] that once considered as a change in life after the pandemic outbreak, this extra bridging responsibility in diffusing COVID-19 related messages is likely to associate with lower life satisfaction.

## 7   Conclusion and Limitation

In this paper, we concentrated on the social media users whose sharing behaviours significantly promote the popularity of COVID-19 related messages. By proposing a new measurement for bridging performance, we identified these influential users. With our collected Twitter data of an international region, we successfully show the influential users suffer from more decrease in their subjective well-being compared to those with smaller bridging performance. We then

**Table 3.** Hierarchical multiple regression model examining variance in SWB explained by independent variables, $*p < 0.05; **p < 0.001$

| Variable | $B$ | $SEB$ | $b$ | $t$ | $R$ | $R^2$ | $\Delta R^2$ |
|---|---|---|---|---|---|---|---|
| **Stage 1** | | | | | -0.207 | 0.043 | 0.043 |
| In-degree | 0.234 | 0.103 | 0.160 | 2.272* | | | |
| Out-degree | 0.861 | 0.680 | 0.054 | 1.267 | | | |
| Pagerank | 3.081 | 0.148 | 0.180 | 2.082* | | | |
| Betweenness centrality | -3.287 | 0.728 | -1.453 | -4.515** | | | |
| **Stage 2** | | | | | -0.312 | 0.097 | 0.054 |
| In-degree | 0.228 | 0.102 | 0.158 | 2.239* | | | |
| Out-degree | 0.075 | 0.080 | 0.050 | 0.945 | | | |
| Pagerank | 0.307 | 0.150 | 0.180 | 2.049* | | | |
| Betweenness centrality | -3.268 | 0.723 | -1.390 | -4.520** | | | |
| Activity | 0.861 | 0.123 | 0.037 | 0.716 | | | |
| **Stage 3** | | | | | -0.579 | 0.335 | **0.238** |
| In-degree | 0.158 | 0.123 | 0.107 | 1.125* | | | |
| Out-degree | 0.516 | 0.45 | 0.050 | 1.147 | | | |
| Pagerank | 0.191 | 0.143 | 0.168 | 1.338* | | | |
| Betweenness centrality | -1.105 | 0.541 | -0.509 | -2.066** | | | |
| Activity | 0.067 | 0.133 | 0.053 | 0.508 | | | |
| UBM | -2.254 | -0.196 | -1.797 | **-11.469** | | | |

conducted the first research to reveal the strong relationship between a user's bridging performance in COVID-19 information diffusion and his/her SWB. Our research provides a cautious reference to public health bodies that some users can be mobilised to help spread health information, but special attention should be paid to their psychological health.

This paper has a few limitations that deserve further discussion. First, we only focused on the affective dimension of subjective well-being while noticing its multi-dimensional nature. This allows us to follow previous SWB studies to convert the calculation of SWB to sentiment analysis, but does not comprehensively evaluate users' cognitive well-being, such as life satisfaction. In our following research, we will attempt to leverage more advanced AI models to investigate cognitive aspects such as *happy* and *angry*. Second, extracting SWB from users' online disclosure inevitably incurs bias compared to social surveys although it supports analysis of an unprecedented large number of users. Third, socio-demographic information of users is not taken into account in this paper. It is known that SWB varies among different socio-demographic groups, and such variation may have an impact on the results of the hierarchical multiple regression [19]. Currently deep learning based models exist for socio-demographic inference. In our future work, we will use the models to extract users' socio-demographic information such as age, gender, income and political orientation to ascertain whether the regression results will change due to the variations of socio-demographic information. Last, we notice that the region we targeted at may introduce additional bias in our results. As a continuous work, we will ex-

tend our study to a region of multiple European countries and cross-validate our findings with other published results in social science.

*Ethical considerations.* This work is based completely on public data and does not contain private information of individuals. Our dataset is built in accordance with the FAIR data principles [31] and Twitter Developer Agreement and Policy and related policies. Meanwhile, there have been a significant amount of studies on measuring users' subjective well-being through social media data. It has become a consensus that following the terms of service of social media networks is adequate to respect users' privacy in research [13]. To conclude, we have no ethical violation in the collection and interpretation of data in our study.

## References

1. Banerjee, D., Meena, K.S.: COVID-19 as an "Infodemic" in public health: Critical role of the social media. Frontiers in Public Health **9**, 231–238 (2021)
2. Barbieri, F., Camacho-Collados, J., Espinosa Anke, L., Neves, L.: TweetEval: Unified benchmark and comparative evaluation for tweet classification. In: Proc. 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP). pp. 1644–1650. Association for Computational Linguistics (2020)
3. Bradley, M.M., Lang, P.J.: Affective norms for English words (ANEW): Instruction manual and affective ratings. Tech. rep., the Centre for Research in Psychophysiology, University of Florida (1999)
4. Chen, E., Lerman, K., Ferrara, E.: Tracking social media discourse about the COVID-19 pandemic: development of a public coronavirus Twitter data set. JMIR Public Health and Surveillance **6**(2), e19273 (2020)
5. Chen, N., Chen, X., Zhong, Z., Pang, J.: From #jobsearch to #mask: Improving COVID-19 cascade prediction with spillover effects. In: Proc. 2021 International Conference on Advances in Social Networks Analysis and Mining(ASONAM). pp. 455–462. ACM (2021)
6. Chen, N., Chen, X., Zhong, Z., Pang, J.: Exploring spillover effects for COVID-19 cascade prediction. Entropy **24**(2) (2022)
7. Cinelli, M., Quattrociocchi, W., Galeazzi, A., Valensise, C.M., Brugnoli, E., Schmidt, A.L., Zola, P., Zollo, F., Scala, A.: The COVID-19 social media infodemic. Scientific reports **10**(1), 1–10 (2020)
8. Diener, E., Emmons, R.A., Larsen, R.J., Griffin, S.: The satisfaction with life scale. Journal of Personality Assessment **49**(1), 71–75 (1985)
9. Diener, E., Oishi, S., Lucas, R.E.: Personality, culture, and subjective well-being: Emotional and cognitive evaluations of life. Annual Review of Psychology **54**(1), 403–425 (2003)
10. Dubey, S., Biswas, P., Ghosh, R., Chatterjee, S., Dubey, M.J., Chatterjee, S., Lahiri, D., Lavie, C.J.: Psychosocial impact of COVID-19. Diabetes & Metabolic Syndrome: Clinical Research & Reviews **14**(5), 779–788 (2020)
11. Duong, V., Luo, J., Pham, P., Yang, T., Wang, Y.: The ivory tower lost: How college students respond differently than the general public to the COVID-19 pandemic. In: Proc. 2020 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM). pp. 126–130. IEEE (2020)

12. Engel de Abreu, P.M., Neumann, S., Wealer, C., Abreu, N., Coutinho Macedo, E., Kirsch, C.: Subjective well-being of adolescents in Luxembourg, Germany, and Brazil during the COVID-19 pandemic. Journal of Adolescent Health **69**(2), 211–218 (2021)
13. Fernando, S., López, J.A.D., Şerban, O., Gómez-Romero, J., Molina-Solana, M., Guo, Y.: Towards a large-scale Twitter observatory for political events. Future Generation Computer Systems **110**, 976–983 (2020)
14. Freeman, L.C.: Centrality in social networks conceptual clarification. Social Networks **1**(3), 215–239 (1978)
15. Guarino, S., Pierri, F., Giovanni, M.D., Celestini, A.: Information disorders during the COVID-19 infodemic: The case of Italian Facebook. Online Social Networks Media **22**, 100124 (2021)
16. Hernandez, R.G., Hagen, L., Walker, K., O'Leary, H., Lengacher, C.: The COVID-19 vaccine social media infodemic: Healthcare providers' missed dose in addressing misinformation and vaccine hesitancy. Human Vaccines & Immunotherapeutics **17**(9), 2962–2964 (2021)
17. Hu, Z., Lin, X., Kaminga, A.C., Xu, H.: Impact of the COVID-19 epidemic on lifestyle behaviors and their association with subjective well-being among the general population in mainland China: Cross-sectional study. Journal of Medical Internet Research **22**(8), e21176 (2020)
18. Huang, S., Lv, T., Zhang, X., Yang, Y., Zheng, W., Wen, C.: Identifying node role in social network based on multiple indicators. PLoS One **9**(8), e103733 (2014)
19. Jaidka, K., Giorgi, S., Schwartz, H.A., Kern, M.L., Ungar, L.H., Eichstaedt, J.C.: Estimating geographic subjective well-being from Twitter: A comparison of dictionary and data-driven language methods. Proceedings of the National Academy of Sciences **117**(19), 10165–10171 (2020)
20. Kupavskii, A., Ostroumova, L., Umnov, A., Usachev, S., Serdyukov, P., Gusev, G., Kustarev, A.: Prediction of retweet cascade size over time. In: Proc. 2012 International Conference on Information and Knowledge Management (CIKM). pp. 2335–2338. ACM (2012)
21. Li, Y.M., Lai, C.Y., Chen, C.W.: Discovering influencers for marketing in the blogosphere. Information Sciences **181**(23), 5143–5157 (2011)
22. Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., Stoyanov, V.: RoBERTa: A robustly optimized bert pretraining approach (2019)
23. Mirbabaie, M., Bunker, D., Stieglitz, S., Marx, J., Ehnis, C.: Social media in times of crisis: Learning from Hurricane Harvey for the coronavirus disease 2019 pandemic response. Journal of Information Technology **35**(3), 195–213 (2020)
24. Ou, X., Li, H.: Ynu_oxz @ haspeede 2 and AMI : XLM-RoBERTa with ordered neurons LSTM for classification task at EVALITA 2020. In: Proc. 2020 Evaluation Campaign of Natural Language Processing and Speech Tools for Italian (EVALITA). vol. 2765 (2020)
25. Page, L., Brin, S., Motwani, R., Winograd, T.: The PageRank citation ranking: Bringing order to the web. Tech. rep., Stanford InfoLab (1999)
26. Rosenthal, S., Farra, N., Nakov, P.: Semeval-2017 task 4: Sentiment analysis in Twitter. In: Proc. 2017 International Workshop on Semantic Evaluation (SemEval). pp. 502–518 (2017)
27. Struweg, I.: A Twitter social network analysis: The South African health insurance bill case. Responsible Design, Implementation and Use of Information and Communication Technology **12067**,  120 (2020)

28. Tabachnick, B.G., Fidell, L.S., Ullman, J.B.: Using Multivariate Statistics. Pearson Education (2007)
29. Wang, Y., Shen, H., Liu, S., Gao, J., Cheng, X.: Cascade dynamics modeling with attention-based recurrent neural network. In: Proc. 2017 International Joint Conference on Artificial Intelligence (IJCAI). pp. 2985–2991. IJCAI (2017)
30. Weng, J., Lim, E.P., Jiang, J., He, Q.: Twitterrank: finding topic-sensitive influential twitters. In: Proc. 2010 ACM International Conference on Web Search and Data Mining (WSDM). pp. 261–270 (2010)
31. Wilkinson, M.D., Dumontier, M., Aalbersberg, I.J., Appleton, G., Axton, M., Baak, A., Blomberg, N., Boiten, J.W., da Silva Santos, L.B., Bourne, P.E., et al.: The fair guiding principles for scientific data management and stewardship. Scientific Data **3**(1), 1–9 (2016)
32. Yang, C., Srinivasan, P.: Life satisfaction and the pursuit of happiness on Twitter. PLoS One **11**(3), e0150881 (2016)
33. Zarocostas, J.: How to fight an infodemic. The Lancet **395**(10225), 676 (2020)
34. Zhou, X., Jin, S., Zafarani, R.: Sentiment paradoxes in social networks: Why your friends are more positive than you? In: Proc. 2020 International Conference on Web and Social Media (ICWSM). pp. 798–807. AAAI Press (2020)