

RESEARCH ARTICLE

## Public opinion evaluation on social media platforms: a case study of High Speed 2 (HS2) rail infrastructure project

Ruiqiu Yao<sup>1\*</sup> and Andrew Gillen<sup>2</sup>

### How to cite

Yao R, Gillen A. Public opinion evaluation on social media platforms: a case study of High Speed 2 (HS2) rail infrastructure project. *UCL Open: Environment*. 2023;(5):10. Available from: <https://doi.org/10.14324/111.444/ucloe.000063>

Submission date: 16 June 2022; Acceptance date: 30 June 2023; Publication date: 8 September 2023

### Peer review

*UCL Open: Environment* is an open scholarship publication, this article has been peer reviewed through the journal's standard open peer-review process. All previous versions of this article and open peer-review reports can be found online in the *UCL Open: Environment* Preprint server at [ucl.scienceopen.com](http://ucl.scienceopen.com)

### Copyright and open access

©2023 The Authors. Creative Commons Attribution Licence (CC BY) 4.0 International licence <https://creativecommons.org/licenses/by/4.0/>

### Open access

This is an open access article distributed under the terms of the Creative Commons Attribution Licence (CC BY) 4.0 <https://creativecommons.org/licenses/by/4.0/>, which permits unrestricted use, distribution and reproduction in any medium, provided the original author and source are credited.



### \*Corresponding author

E-mail: [ruiqiu.yao.19@ucl.ac.uk](mailto:ruiqiu.yao.19@ucl.ac.uk)

<sup>1</sup>Civil, Environmental and Geomatic Engineering, University College London, London, UK

<sup>2</sup>Department of Civil and Environmental Engineering, Northeastern University, Boston, MA, USA

## Abstract

Public opinion evaluation is becoming increasingly significant in infrastructure project assessment. The inefficiencies of conventional evaluation approaches can be improved with social media analysis. Posts about infrastructure projects on social media provide a large amount of data for assessing public opinion. This study proposed a hybrid model which combines pre-trained RoBERTa and gated recurrent units for sentiment analysis. We selected the United Kingdom railway project, High Speed 2 (HS2), as the case study. The sentiment analysis showed the proposed hybrid model has good performance in classifying social media sentiment. Furthermore, the study applies latent Dirichlet allocation topic modelling to identify key themes within the tweet corpus, providing deeper insights into the prominent topics surrounding the HS2 project. The findings from this case study serve as the basis for a comprehensive public opinion evaluation framework driven by social media data. This framework offers policymakers a valuable tool to effectively assess and analyse public sentiment.

**Keywords:** public opinion evaluation, civil infrastructure projects, machine learning, sentiment analysis, topic modelling

## Introduction

Infrastructure systems lay the foundation of the economy for a nation by providing primary transportation links, dependable energy systems and water management systems to the public.

In the United Kingdom (UK), the National Infrastructure Strategy 2020 reveals the determination of the UK government to deliver new infrastructure and upgrade existing infrastructure across the country to boost growth and productivity and achieve a net-zero objective by 2050 [1]. Although infrastructure projects positively affect the national economy, they can negatively impact the environment and society. For instance, they may disrupt the natural habitat of wildlife by filling up wetlands. As a result, the wildlife may have to migrate to other regions, causing problems to the ecology of certain regions [2].

Environmental impact assessments (EIA) are a critical part of the planning and delivery of large infrastructure projects. In EIA research, public participation schemes are becoming increasingly popular. O'Faircheallaigh [3] emphasised the importance of public participation in EIA decision-making processes. Social media platforms are gaining increasing ubiquity and are emerging methods for the public to participate in decision-making processes and raise environmental concerns. Thus, the research objective of this study is to evaluate the feasibility of using social media data to perform public participation analysis.

## Conventional approaches to public opinion evaluations

Public hearings and public opinion polling are the two most adopted public consultation approaches. Checkoway [4] stated some drawbacks of public hearings. For instance, the technical terms are hard to understand for the public, and participants often do not represent the actual population. As for polling, Heberlein [5] revealed that conducting polling can usually take a month or even years. As civil infrastructure projects typically have tight project timelines, there is a need for a more efficient public opinion evaluation method.

Moreover, Ding [6] argued that the data collection process is costly for conventional opinion polling. A typical 1000-participant telephone interview will cost tens of thousands of US dollars to carry out [7]. Besides conducting surveys, costs associated with data input and data analysis should also be considered [6].

Public hearings and polling are not ideal for obtaining public opinions for infrastructure projects. They can be costly, invasive and time-consuming. Therefore, researchers have drawn attention to developing an alternative method for obtaining and assessing public opinion. A new opportunity in acquiring and evaluating public opinion has emerged with the growing popularity of various social media platforms [8]. User-generated content on social media platforms provides a huge amount of data for text mining. This text data is an alternative resource for opinion evaluation toward civil infrastructure projects.

## Related work on public opinion evaluation with social media analysis

Kaplan and Haenlein [9] defined social media platforms as Internet-based applications adopting Web 2.0 (participative Web). Due to the number of active users on Facebook and Twitter, the massive amount of user-generated content provides valuable opportunities for researchers to study various social topics [10]. Moreover, with machine learning and natural language processing, researchers can perform advanced and automated algorithms on social media posts, such as sentiment analysis and topic modelling. Sentiment analysis can categorise the textual data in social media into different emotional orientations, providing an indicator of public opinion. Recent research on infrastructure project evaluation with social media analysis revealed the feasibility of using social media analysis as an alternative public opinion evaluation method.

Aldahawi [11] investigated social networking and public opinion on controversial oil companies by sentiment analysis of Twitter data. Kim and Kim [12] adopted lexicon-based sentiment analysis for public opinion sensing and trend analysis on nuclear power in Korea. Lexicon-based sentiment analysis with domain-specified dictionaries and topic modelling has also been used on public opinion data for California High-Speed Rail and the Three Gorges Project [6,8]. Lexicon-based sentiment analysis calculates the sentiment of documentation from the polarity of words [13]. In lexicon-based sentiment analysis, it is assumed that words have inherent sentiment polarity independent of their context. A user must establish dictionaries containing words with sentiment polarity to build a lexicon-based classifier. After building up the classifier, the polarity of a document is calculated in three phases: establishing word-polarity value pairs, replacing words in the document with polarity values and calculating the sentiment polarity for the document. Ding [6]

tailor-made a dictionary by removing unrelated words from a positive word list. Jiang et al. [8] built a dictionary for hydro projects by integrating the Nation Taiwan Sentiment Dictionary [14], Hownet (a Chinese/English bilingual lexicon database) [15] and a hydro project-related word list. Recent research showed the practicality of implementing the lexicon-based sentiment analysis for public opinion evaluation on civil projects. The recent developments in deep learning show a promising future for public opinion evaluation.

## Recent development of natural language processing

In 2014, Bahdanau et al. [16] introduced a novel neural network architecture named attention mechanisms. Attentional mechanisms are designed to mimic cognitive perception, which computes the attention weight on input sequences so that some parts of the input data obtain more attention than the rest. In 2017, Vaswani et al. [17] published their ground-breaking research paper 'Attention is all you need', where they proposed an influential neural network named transformer. The transformer architecture leverages self-attention and multi-head attention to enable parallel computation. Using multiple attention heads and a self-attention mechanism, the transformer architecture can obtain different aspects of input data through learning different functions. As a result, transformer architecture can handle increased model and data size. Kaplan et al. [18] demonstrated that transformer models have remarkable scaling behaviour as model performance increases with training size and model parameters. Hence, natural language processing can benefit from large-language models, such as generative pre-trained transformer (GPT) [19,20] and Bidirectional Encoder Representations from Transformers (BERT) [21].

## Research question and main contributions

The recent developments in deep learning research motivated this study to assess how state-of-art machine learning algorithms can help public opinion evaluation on infrastructure projects. The main contributions of this study include:

- (1) This study proposed a hybrid transformer-recurrent neural network model for sentiment analysis, which combines the pre-trained Robustly optimised BERT approach (RoBERTa) [22] and bidirectional gated recurrent neural networks [23].
- (2) This study employed tweets data of High Speed 2 (HS2) as a case study, utilising it to compare the performance of the proposed RoBERTa–bidirectional gated recurrent unit (BiGRU) with baseline classifiers. Moreover, this study applied topic modelling with latent Dirichlet allocation (LDA) on tweet corpus.
- (3) Based on the insights from the case study results, the study proposes a public opinion evaluation framework that leverages social media data with RoBERTa–BiGRU and topic modelling. This framework provides a valuable tool for policymakers to evaluate public opinion effectively.

The rest of this article is organised as follows: the machine learning models section provides a detailed exposition of the machine learning algorithms used in this study. The case study with the High Speed 2 project section delves into the specific details and findings. This is followed by the limitations of this research and suggests potential avenues for future research. Finally, the conclusion summarises the main findings and contributions.

## Machine learning models

This section provides a comprehensive overview of implementing machine learning algorithms for public opinion evaluation. The formulation of the multinomial naïve Bayes (MNB) classifier is presented. The proposed RoBERTa–BiGRU model is then introduced, highlighting its essential components and architecture. Finally, the topic modelling technique using LDA is discussed.

## Sentiment analysis with an MNB classifier

The naïve Bayes classifier is a family of probabilistic classification models based on the Bayes theorem [24]. The term 'naïve' means the naïve assumption of independence among each pair of features (attributes) and class variable values [25]. More specifically, the 'naïve' assumption means

that classifiers process the text data independently as a bag-of-words, ignoring the relationships among words, such as sequences, and only considering the word frequency in the document. The mathematical formula of the Bayes theorem Eq. (1) states that given  $n$  feature vectors  $x_1, \dots, x_n$  and class variable  $y$ , the probability distribution of  $y$  is:

$$P(y|x_1, \dots, x_n) = \frac{P(y)P(x_1, \dots, x_n|y)}{P(x_1, \dots, x_n)} \quad (1)$$

Because the probability distribution of feature vectors  $P(x_1, \dots, x_n)$  is given by the model input, the following classification rule Eq. (2) and Eq. (3) can be obtained [26]:

$$P(y|x_1, \dots, x_n) \propto P(y) \prod_{i=1}^n P(x_i|y) \quad (2)$$

$$\hat{y} = \underset{y}{\operatorname{argmax}} P(y) \prod_{i=1}^n P(x_i|y), \quad (3)$$

where  $P(y)$  is the frequency distribution of  $y$  in the training dataset and  $P(x_i|y)$  is determined by the naïve Bayes classifier assumptions. For example, the Gaussian naïve Bayes classifier assumes  $P(x_i|y)$  follows a Gaussian distribution.

In the case of the MNB classifier, the multinomial distribution is parameterised by  $(\theta_{y_1}, \dots, \theta_{y_n})$  vectors for each  $y$  with  $n$  features.  $\theta_{y_i}$  indicates the probability distribution of  $x_i$  under class  $y$  in the training set. In other words,  $\theta_{y_i} = P(x_i|y)$ . Then, smoothed maximum likelihood estimation [27] can be used to estimate  $\theta_{y_i}$ :

$$\hat{\theta}_{y_i} = \frac{N_{y_i} + \alpha}{N_y + \alpha n}, \quad (4)$$

where  $N_{y_i}$  is the number of occurrences of feature  $i$  for sentiment class  $y$ ;  $N_y$  is the number of occurrences of all features for  $y$ ; and  $\alpha$  is the smoothing prior, which is a hyperparameter to be tuned.

## Sentiment analysis with RoBERTa–BiGRU

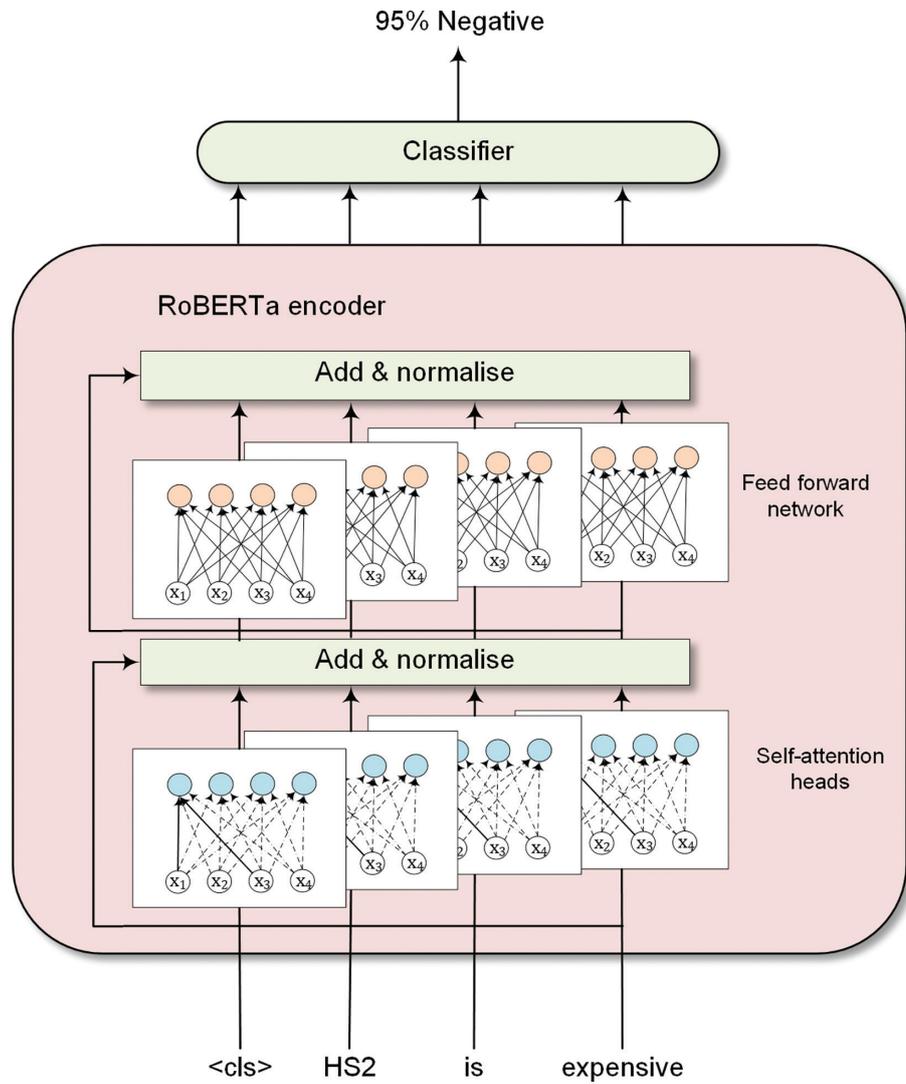
As mentioned in the recent development of natural language processing, transformer architectures have remarkable scaling ability to handle large training data sizes and model parameters. As a result, researchers have proposed fine-tuning a pre-trained large-scale transformer model for specific downstream natural language processing tasks. This approach is referred to as transfer learning and leverages knowledge learned from the large-scale database to other downstream tasks [28]. The Bidirectional Encoder Representations from Transformers (BERT) [21] is a large language model that has state-of-the-art for natural language processing performance. The BERT model encodes text data in a bidirectionally way such that BERT can process text tokens in both left-to-right and right-to-left directions. This study used a variant of the BERT model, named the Robustly optimised BERT approach (RoBERTa) [22], because RoBERTa is pre-trained on a much larger scale of text data than BERT.

Details of fine-tuning the RoBERTa model for sentiment analysis are shown in Fig. 1. RoBERTa used similar transformer architecture as BERT. The input token sequence is passed to multiple self-attention heads, followed by a layer normalisation [29]. The normalised data is subsequently sent to feed-forward networks and a second layer normalisation. Figure 1 shows the transformer architecture of a single encoder layer. The RoBERTa model contains multiple encoders based on model preference. A RoBERTa encoder's hidden states can then be fed into a classifier for classification tasks. Noticeably, the '<cls>' token indicates the global representation of input text [28].

The classifier can be different neural network architectures, such as feedforward neural networks (FNN) or recurrent neural networks (RNN). The long short-term memory (LSTM) architecture is a prevalent choice as the classifier [30]. The LSTM introduced internal states and gates in addition to RNN to process information in sequenced data [31]. The GRU architecture, proposed by Cho [23] in 2014, is a streamlined adaptation of LSTM architecture which retains internal states and gating mechanisms. This study adopted the GRU architecture as a classifier from RoBERTa outputs because gated recurrent unit (GRU) has a faster computation speed than LSTM with comparable performance [32].

Figure 1

Fine-tuning RoBERTa for sentiment analysis.



The GRU model consists of two internal gates: a reset gate and an update gate. The reset gate determines the extent to which information from the previous state is retained, while the update gate controls the proportion of the new state that replicates the old state. The mathematical formulate of the reset gate and update gate are:

$$R_t = \sigma(W_{ir} X_t + b_{ir} + W_{hr} H_{t-1} + b_{hr}) \tag{5}$$

$$Z_t = \sigma(W_{iz} X_t + b_{iz} + W_{hz} H_{t-1} + b_{hz}) \tag{6}$$

$$\sigma(x) = \frac{1}{1 + \exp(-x)}, \tag{7}$$

where  $X_t \in \mathbb{R}^{n \times d}$  is a minibatch input of a memory cell ( $n$  is the number of sample and  $d$  is the dimension of features);  $H_{t-1} \in \mathbb{R}^{n \times h}$  is the hidden state of the previous step ( $h$  is the number of hidden units of a GRU memory cell);  $W_{ir}, W_{hr} \in \mathbb{R}^{d \times h}$  and  $W_{iz}, W_{hz} \in \mathbb{R}^{n \times h}$  are model weights; and  $b_{ir}, b_{hr}, b_{iz},$  and  $b_{hz}$  are model bias parameters. The reset gate  $R_t \in \mathbb{R}^{n \times h}$  and update gate  $Z_t \in \mathbb{R}^{n \times h}$  are computed based on Eq. (5) and Eq. (6). In other words, two gates are fully connected layers with sigmoid activation function Eq. (7).

The reset gate is designed to yield a candidate hidden state  $N_t \in \mathbb{R}^{n \times h}$  with Eq. (8) and tanh activation function Eq. (9). The influences of previous information  $H_{t-1}$  in Eq. (8) is reduced by the Hadamard product of  $R_t$  and  $H_{t-1}$ . The candidate hidden state  $N_t$  is then passed to Eq. (10) to calculate the new hidden state  $H_t$ , in which the update gate  $Z_t$  controls the degree to which  $H_t$  resembles  $N_t$ .

$$N_t = \tanh(W_{in} X_t + b_{in} + R_t \odot (W_{hn} H_{t-1} + b_{hn})) \tag{8}$$

$$\tanh(x) = \frac{1 - \exp(-2x)}{1 + \exp(2x)} \tag{9}$$

$$\mathbf{H}_t = (1 - \mathbf{Z}_t) \odot \mathbf{N}_t + \mathbf{Z}_t \odot \mathbf{H}_{t-1}, \tag{10}$$

where  $\mathbf{W}_{in} \in \mathbb{R}^{d \times h}$  and  $\mathbf{W}_{hn} \in \mathbb{R}^{h \times h}$  are model weights;  $b_{in}$  and  $b_{hn}$  are bias parameters; and  $\odot$  is the Hadamard product, which is also referred to as the element-wise product.

Similar to the bidirectional setting of BERT, a two-layer GRU is also able to process the text data bidirectionally with a forward layer and a backward layer, as shown in Fig. 2. The hidden state of the forward layer and backward layer is denoted as  $\overrightarrow{\mathbf{H}}_t \in \mathbb{R}^{n \times h}$  and  $\overleftarrow{\mathbf{H}}_t \in \mathbb{R}^{n \times h}$ . The forward layer hidden states  $\overrightarrow{\mathbf{H}}_t$  is then multiplied with a dropout rate  $\delta$ , which is a Bernoulli random variable with  $\delta$  probability of being 0. The output of a GRU is a concatenate of  $\overrightarrow{\mathbf{H}}_{t,\delta}$  and  $\overleftarrow{\mathbf{H}}_t$  with dimension  $n \times 2h$ .

The RoBERTa model can be fine-tuned by optimising the loss function of the above-mentioned bidirectional GRU and connecting the output of a bidirectional GRU with a fully connected layer.

The loss function to be optimised in GRU is a cross entropy function [33]. Moreover, the fully connected layer uses the softmax activation function Eq. (11):

$$s(x_i) = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}, \tag{11}$$

where  $n$  is the number of sentiment classes. The fully connected layer converts the hidden states of the bidirectional GRU to the probability of each sentiment class.

Figure 3 demonstrates the complete structure of the RoBERTa-BiGRU model. Firstly, tweets are tokenised with the RoBERTa tokeniser. Then, the tokens are passed to 12 encoders with

Figure 2

Bidirectional GRU model.

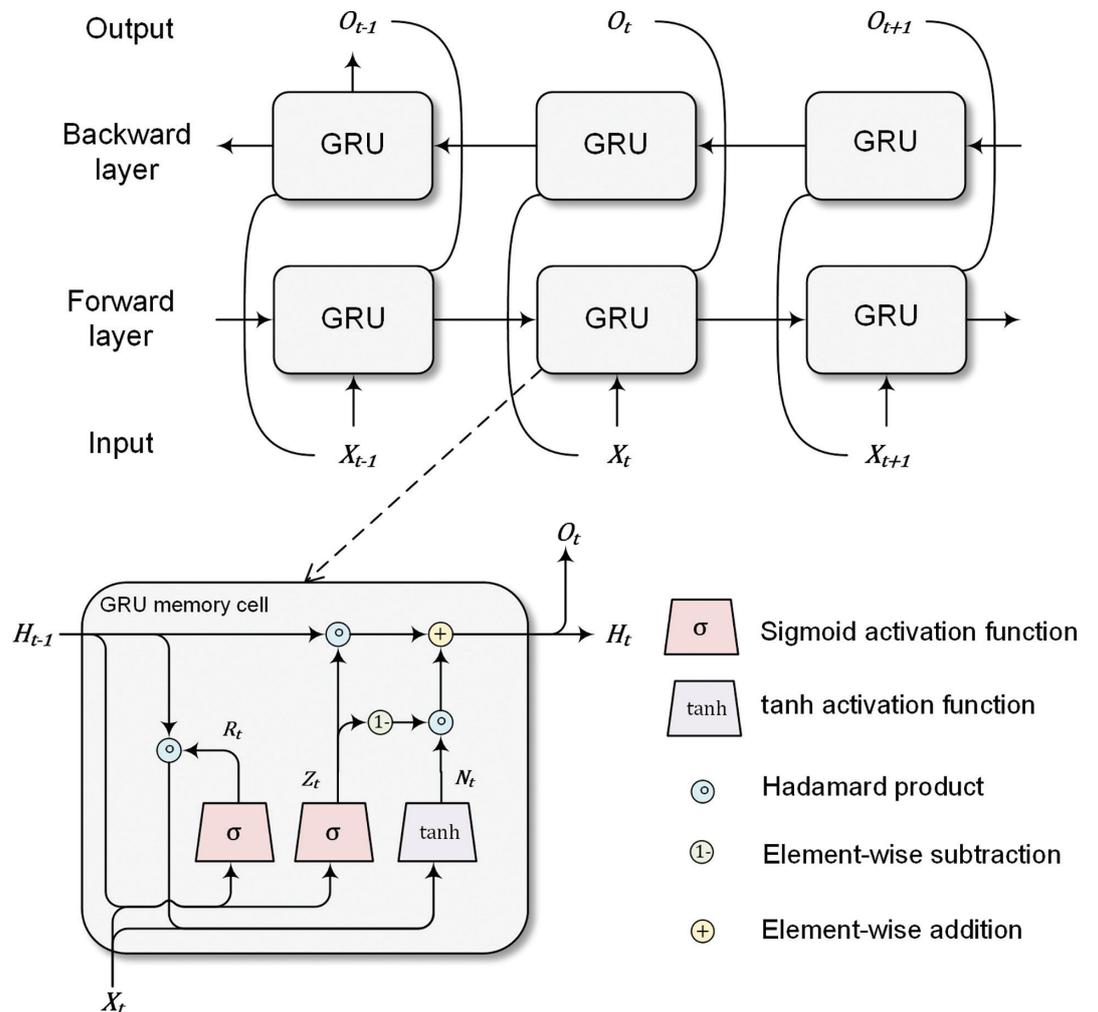
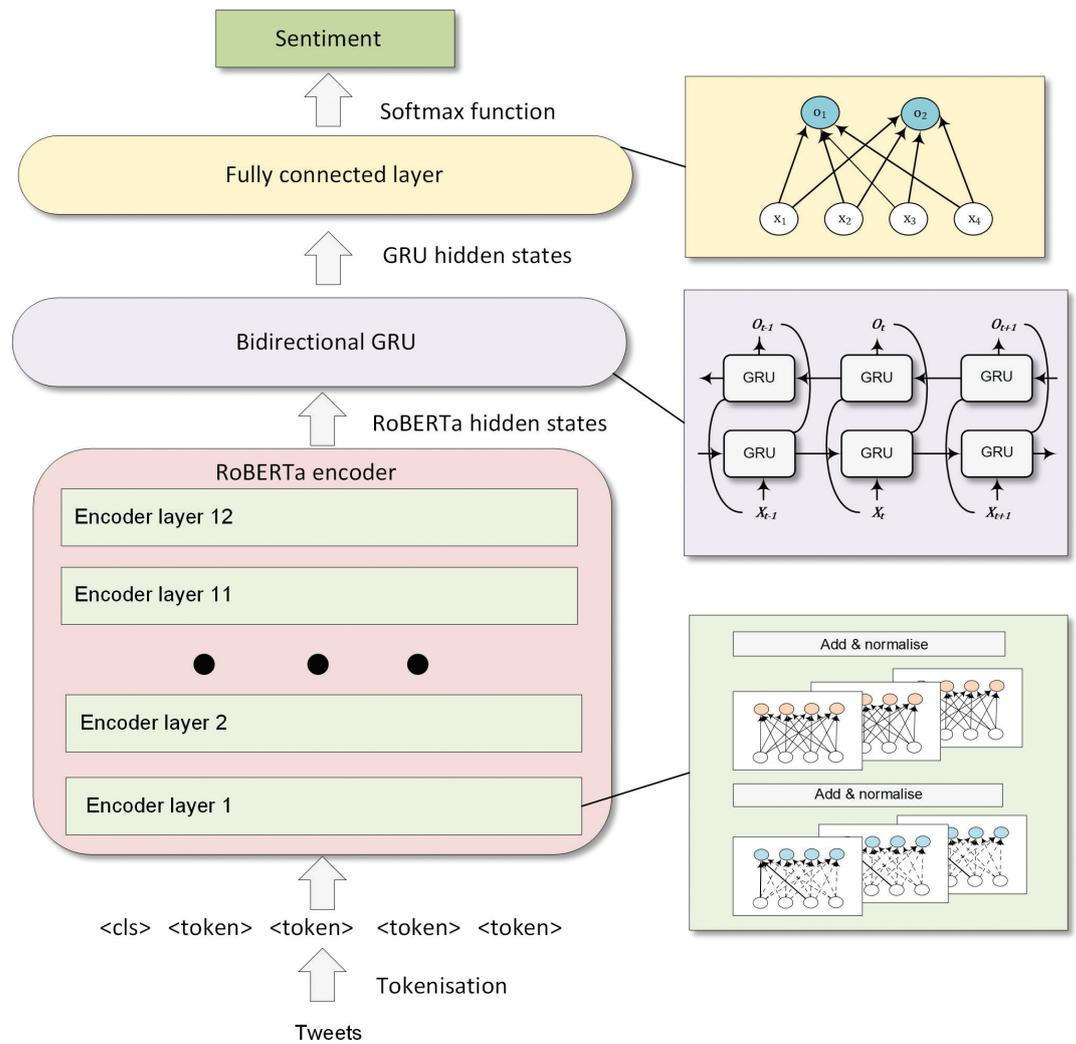


Figure 3

Structure of RoBERTa-BiGRU for sentiment analysis.



multiple self-attention heads to obtain 768 tweets’ hidden representations. The tweets’ hidden representations can then be allocated to sentiment classes through a bidirectional GRU and fully connected layer.

### Topic modelling with LDA

Deerwester et al. [34] proposed a latent semantic indexing method for topic modelling, applying singular value decomposition (SVD) to derive the latent semantic structure model from the matrix of terms from documents. SVD is a linear algebra technique to decompose an arbitrary matrix to its singular values and singular vectors [35]. Blei et al. [36] introduced LDA, which is a general probabilistic model of a discrete dataset (text corpus).

LDA is a Bayesian model, which models a document as a finite combination of topics. Each topic is modelled as a combination of topic probabilities. For example, an article that talks about the structural design of a building complex may have various topics, including ‘structural layout’ and ‘material’. The topic ‘structural layout’ may have high-frequency words related to structural design, such as ‘beam’, ‘column’, ‘slab’ and ‘resistance’. Also, the ‘material’ topic may have the words ‘concrete’, ‘steel’, ‘grade’ and ‘yield’. In short, a document has different topics with a probabilistic distribution, and each topic has different words with a probabilistic distribution. Human supervision is not required in LDA topic modelling, as LDA only needs a number of topics to perform an analysis.

Topic modelling with LDA has a wide range of applications in research. Xiao et al. [37] used LDA variant topic modelling to uncover the probabilistic relationship between adverse drug reaction topics. Xiao et al. found that the LDA variant topic modelling has higher accuracy than alternative methods. Jiang et al. [8] showed the feasibility of LDA topic modelling to extract topics about

the Three Gorges Project on a Chinese social media platform. Apart from focusing on extracting terms from the textual corpus, topic modelling is another trend-finding tool, as it will reveal the relationship between topics. Chuang et al. [38] proposed a method to visualise topics with circles in a two-dimensional plane, whose centre is determined by the calculated distance between topics. The distance is calculated using Jensen–Shannon divergence, and principal components analysis determines the size of the circle [39].

### Case study with the HS2 project

This section provides implementation details of sentiment classifiers and topic modelling methods for the HS2 case study. First, the background of the HS2 project is presented, offering insights into the rail infrastructure project. This is followed by an explanation on data collection and processing, detailing the methods employed to gather social media data related to HS2. Then the evaluation metrics used to assess the performance of sentiment classifiers are presented, enabling a thorough examination of sentiment classification models. The following two sections show the results of sentiment analysis and topic modelling, respectively. Finally, a framework for evaluating public opinion based on social media data is introduced.

### Background on the HS2 project

The transportation demand for the UK railway network has steadily grown over the past decades. According to the Department for Transport [40], rail demand has doubled since 1994–1995, with a rising rate of 3% every year. Therefore, the HS2 programme is proposed to construct a new high-speed and high-capacity railway, aiming to boost the economy in the UK, improve connectivity by shortening journey time, provide sufficient capacity to meet future railway network demand and reduce carbon emission by reducing long-distance driving. Figure 4 shows that HS2 will connect London, Leeds, Birmingham and Manchester, joining the existing railway infrastructure to allow passengers to travel to Glasgow, Newcastle and Liverpool [41].

Figure 4

HS2 infrastructure map [41].



## Data preparation

The collection of HS2-related tweets was carried out using Twitter application programming interfaces (API). Specifically, tweets that containing the hashtags '#HS2' and '#HighSpeed2' were collected. However, the number of collectable tweets is constrained by the limitations imposed by the Twitter API, which restricts the collection to under 10,000 tweets. Thus, the total number of tweets collected was 8623 tweets. The tweets were sampled over a 5-year period from 2017 to 2020. The tweet distribution across the years is: 2017 (1544 tweets), 2018 (1130 tweets), 2019 (2909 tweets) and 2020 (3040 tweets). Noticeably, the tweets collected were in an extended mode, allowing the retrieval of the complete text, surpassing the 140-character limit.

Data preprocessing involves cleaning and preparing data to increase the accuracy and performance of text-mining tasks, such as sentiment analysis and topic modelling. Tweet text data tend to contain uninformative text, such as URL links, Twitter usernames and email addresses. For MNB and lexicon-based classifiers, the stop words need to be removed. To be more specific, stop words are words that do not have sentiment orientation, such as 'me', 'you', 'is', 'our', 'him' and 'her'. As each word in text data is treated as a dimension, keeping stop words and uninformative text will complicate the text mining by making text mining a high dimension problem [42]. Other text preprocessing techniques for MNB and lexicon-based classifiers include text lowercasing and text stemming. Noticeably, the transformer architectures do not require removing stop words, lowercasing and text stemming, as transformers are able to handle the implied information in stop words.

Upon conducting a manual inspection of collected tweets, the number of tweets expressing positive sentiment was significantly lower than those with negative or neutral sentiment. The sentiment classification task is set to binary to address the imbalance issue. The task was designed to classify tweets as either having negative sentiment or non-negative sentiment (including neutral and positive sentiments). A set of 1400 tweets was carefully annotated to train classifiers in this case study. Within this annotated dataset, 700 tweets were labelled as negative sentiment, while the remaining 700 tweets were labelled as non-negative sentiment. To access the annotated training tweets, a GitHub link is provided in the [Open data and materials availability statement](#), facilitating transparency and reproducibility of this study. The annotated tweets were split into 70% training dataset (980 tweets) and 30% validation dataset (420 tweets).

## Sentiment analysis results

Three sentiment classifiers were used in this case study: (1) VADER [43], a rule-based lexicon sentiment classifier. (2) An MNB classifier, which is built following details in the background on the HS2 project. (3) A RoBERTa-BiGRU model that is developed based on the architecture presented in the data preparation section. The model details of each classifier are shown in Table 1. The hyperparameters in MNB and RoBERTa-BiGRU, such as smoothing priors  $\alpha$ , batch size, hidden units and dropout rate, were tuned by a grid search. The RoBERTa-BiGRU model was trained on a Tesla T4 GPU on Google Colab with a total training time of 2421.23 s for 100 epochs.

The performances of three classifiers were evaluated with accuracy and receiver operating characteristic (ROC) curve. Accuracy, as shown in Eq. (12), measures the accuracy of the classifier with all correctly identified cases overall. A ROC curve plots the true positive rate, as shown Eq. (13), along the y axis and the false positive rate, as shown in Eq. (14), along the x axis. An ROC curve shows the graphical interpretation of gain (true positive rate) and loss (false positive rate) [44]. The area under the curve (AUC) score calculates the total area under the ROC curve. The AUC

**Table 1. Model details of each classifier**

Name	Model parameters
VADER	Rules specified in [43]
MNB	Smoothing priors: $\alpha = 0.1$
RoBERTa-BiGRU	Batch size: 16 Hidden units: 256 Dropout rate: 0.5 Optimiser: AdamW Learning rate: $2 \times e^{-6}$ Epoch: 100

score quantitatively evaluates the performance of a classifier, which represents the possibility that a random positive datapoint ranks higher than a random negative datapoint [45].

$$\text{accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \tag{12}$$

$$\text{true positive rate (recall)} = \frac{TP}{TP + FN} \tag{13}$$

$$\text{false positive rate} = \frac{FP}{FP + TN}, \tag{14}$$

where  $TP$  = true positive,  $TN$  = true negative,  $FP$  = false positive and  $FN$  = false negative.

Table 2 demonstrates the accuracy of each sentiment classifier. The lexicon-based VADER has the lowest accuracy (70.24%) among the three classifiers. MNB and RoBERTa-BiGRU show better accuracy performance than VADER, whereas MNB and RoBERTa-BiGRU have increased accuracies of 12.38% and 19.28%, respectively. MNB and RoBERTa-BiGRU are then compared with respect to the AUC scores. MNB has an AUC score of 0.9023, while RoBERTa-BiGRU has a slightly lower AUC score of 0.8904. Both MNB and RoBERTa-BiGRU have an AUC score of around 0.9, which indicates that both models have a high level of ability to classify tweet sentiment. Noticeably, Fig. 5(b) has a much steeper curve. The steeper curve means that RoBERTa-BiGRU can achieve higher recall with a low FP rate, which is desirable behaviour in sentiment analysis. As a result, RoBERTa-BiGRU has the best performance in terms of both accuracy and the ROC curve. Thus, RoBERTa-BiGRU was used for sentiment analysis with all collected tweets.

Figure 6 shows the sentiment distribution of HS2-related tweets from 2017 to 2020. Notably, there was a substantial increase in the number of tweets in 2019, indicating a heightened presence of the HS2 project in social media discussions during and after that year. Moreover, it is worth mentioning that the majority of tweets collected across all time periods exhibited a negative sentiment. Specifically, negative tweets accounted for 57.77% in 2017, 53.32% in 2018, 60.64% in 2019 and 65.19% in 2020.

The substantial proportion of negative tweets in all periods indicates a prevailing negative sentiment among the public regarding HS2, highlighting the importance for policymakers and decision-makers to take this sentiment into consideration. However, it is essential to approach these findings with caution. While the high percentage of negative tweets may raise concerns, it is crucial to note that this alone does not necessarily imply a public relationship emergency for HS2. It is worth acknowledging that certain Twitter users might repeatedly express their negative sentiment towards HS2 [46], potentially influencing the overall sentiment distribution. Given the sentiment analysis

**Table 2. Model accuracy performance**

Name	Accuracy
VADER	70.24%
MNB	82.62%
RoBERTa-BiGRU	89.52%

**Figure 5**

- (a) ROC curve for MNB classifier.
- (b) ROC curve for RoBERTa-BiGRU.

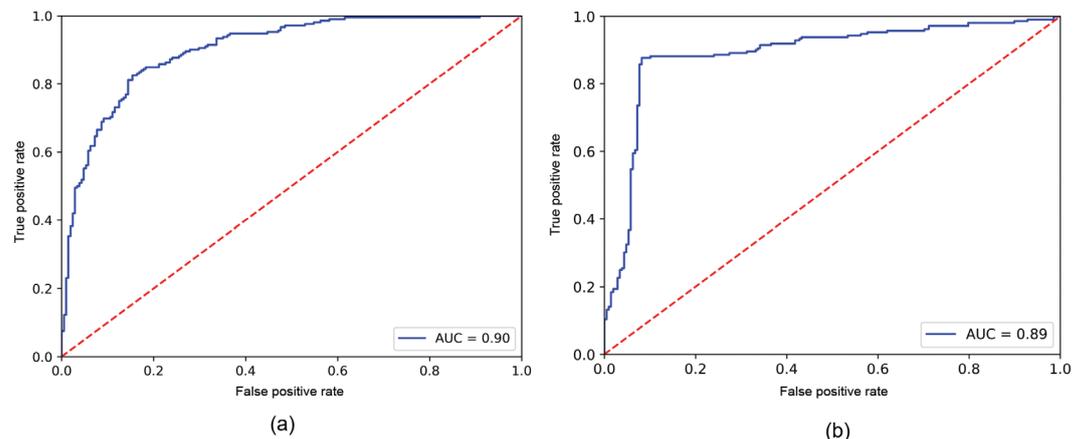
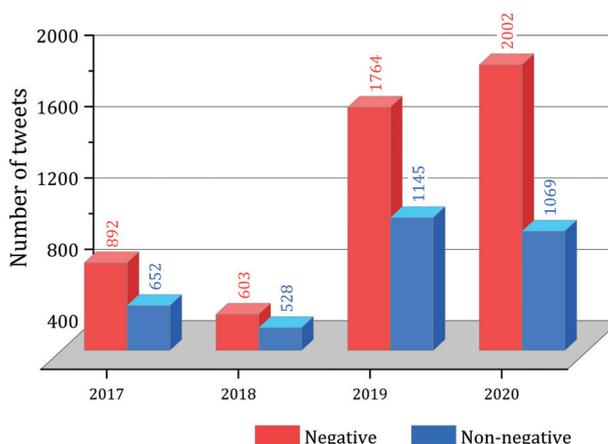


Figure 6

Sentiment analysis results for 2017–2020.



results, it is important to uncover the key topics within the tweets discussions, necessitating the application of topic modelling.

### Topic modelling results

The tweets dataset is then classified using the RoBERTa–BiGRU model into two collections: negative corpus and non-negative corpus. Each collection was performed individually with topic modelling and visualisation. Topic modelling with LDA was performed with genism, a collection of Python scripts developed by Rehurek and Sojka [47]. We used pyLDAvis, a Python library, for visualising topics such that we could determine the most suitable number of topics. Several models were constructed with a number of topics ranging from 3 to 20. We selected five as the number of topics through manual inspection of term distribution and topic relevance.

#### Negative tweets corpus

Table 3 shows major topics in the negative corpus. Topic 1 is the largest topic and accounts for 35.3% of the negative corpus. Topic 1 contains words such as ‘need’, ‘money’, ‘nhs’, ‘badly’ and ‘billion’. These words express the negative sentiment on HS2 budget spending. These tweets criticise the over-spending of HS2 and argue that the money should be invested in the National Health Service (NHS) rather than HS2. Topic 2 and Topic 4 have a similar focus. Topic 2 has words such as ‘government’, ‘protester’ and ‘social’, and Topic 4 include words such as ‘stophs2’, ‘petition’, ‘media’ and ‘political’. Both Topic 2 and Topic 4 discuss the campaign to stop HS2 project by a petition. Topic 3 and Topic 5 show some relevance. Topic 3 contains ‘stop’, ‘please’, ‘trees’, ‘contractors’, ‘changed’ and ‘essential’, which raises environmental concerns about construction work on woodlands. Topic 5 also discusses the environmental issues with the words ‘construction’, ‘damage’ and ‘destroy’.

#### Non-negative tweets corpus

Table 4 shows topics in a non-negative corpus. Topic 1 includes words such as ‘new’, ‘railway’, ‘good’, ‘midlands’ and ‘important’, where tweets express positive sentiment on HS2 by mentioning

Table 3. Topics in negative corpus

Topic number	Terms	Topic percentage
1	Borisjohnson, hs2, work, time, need, money, say, nhs, use, uk, course, amp, transport, nt, cancel, even local, badly, billion, ancient, public, needed, boris, way, think, country, rishisunak, trains, know	35.3%
2	Rail, government, going, still, protesters, like, news, case, go, social, could, economic, train, people, home, London, times, business, ltd, working, travel, back, road, north, sense, says, dont	24.2%
3	Stop, post, mps, please, another, anti, away, seems, trees, make, already, without, contractors, may, changed, control, steeple, long, big, bill, sign, essential, protest, claydon, likely, means, yet, billions, station, caught	13.9%
4	Sopths2, workers, petition, sites, via, take, destruction, ever, change, media, track, year, ukparliament, least, investment, everyone, account, despite, find, continue, political, wants, white, along, british, longer, evidence, called, massive, elephant	13.6%
5	Report, scrap, construction, costs, last, end, law, latest, true, tax, first, damage, full, job, trident, nesting, figures, wonder, share, read, unnecessary, questions, destroy, failed, coming, vital	13.1%

Table 4. Topics in non-negative corpus

Topic number	Terms	Topic percentage
1	Work, new, project, one, railway, station, first, time, may, ever, people, plans, common, good, midlands, find, watch, still, well, way, may, could, largest, part, back, important, day	35.4%
2	Construction, hs2ltd, rail, post, projects, train, business, build, track, road, read, network, phase, industry, latest, leaders, think, green, big, please, works, air, know, local, year, along	24.4%
3	High, speed, need, old, north, planning, would, capacity, built, engineering, course, Manchester, building, another, plan, recent, airport, must, benefit, needs, evidence, better, needed, chief, funding	15.9%
4	Government, news, trains, us, would, home, two, heathrow, cost, start, railways, service, suppliers, roads, update, every, keep, seems, question, longer, join, money	13.3%
5	Stations, use, lake, community, following, scheme, economic, really, opportunities, spending, committee, supply, benefits, due, chain, role, early, daily, fund, freight, article, essential, airports	11.1%

the positive effect on the Midlands. A similar result can be found in Topic 3, which includes words such as ‘planning’, ‘Manchester’, ‘airport’, ‘benefit’, ‘better’. Topic 3 highlights that the transportation infrastructure in Manchester could benefit from the HS2 project. Topic 2 discusses the business case of HS2 with words such as ‘project’, ‘business’, ‘build’, ‘network’ and ‘industry’. Topics 4 and 5 both discuss potential improvements on the accessibility to the airport with words ‘heathrow’, ‘airports’, ‘opportunities’. Overall, the LDA topic modelling showed good execution on obtaining key topics from the tweet corpus.

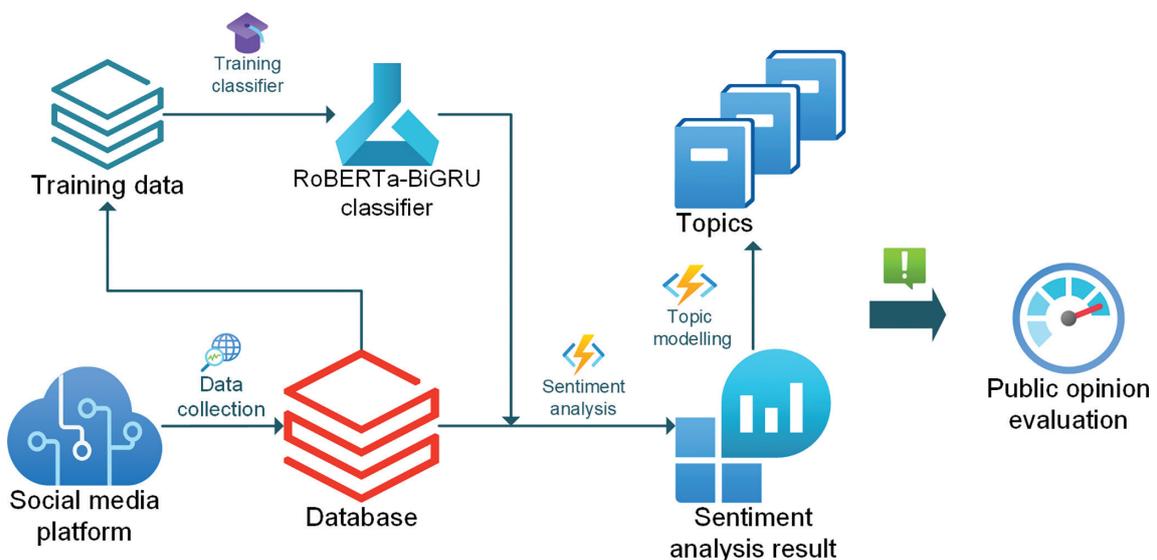
**Proposed public opinion evaluation framework using social media data**

The case study results showed that the RoBERTa-BiGRU and LDA topic modelling have a good performance in evaluating public opinion on HS2 with tweet data. Hence, a public opinion evaluation framework using social media data is proposed to facilitate the decision-making of policymakers.

Figure 7 presents the comprehensive public opinion evaluation framework that utilises social media data. The process begins by collecting social media data, such as tweets, and storing them in a database. Subsequently, the social media data is processed through sentiment annotation, which involves labelling the data to create training sets. These training sets are then utilised for training a sentiment classifier called RoBERTa-BiGRU. Once the RoBERTa-BiGRU sentiment classifier is trained, it is employed to categorise social media tweets into predefined sentiment labels. Additionally, leveraging LDA topic modelling, the framework extracts key topics from the social media data. Policymakers can subsequently utilise the sentiment analysis results and key topics to evaluate public opinion regarding infrastructure projects.

Figure 7

Public opinion evaluation framework.



## Limitation and future research direction

### Human factors in annotating tweets sentiments

Researchers usually assign multiple annotators (3 to 5) to tag the sentiment orientation to minimise the influence of human annotators [48]. However, in our study, all the training data was tagged by one annotator. As a result, the human factor may have affected the accuracy of the sentiment classifier. The future application of fine-tuning sentiment classifiers could benefit from multiple annotators.

Another impact of human factors could be different sentiment interpretations. For example, the following tweet may be tagged with different sentiment orientations. '#HS2 is a £100bn scheme to have slightly shorter journey times from Manchester and Birmingham to London, thereby solving Britain's biggest ever problem.' One annotator can argue that there are positive sentiment signs (shorter journey time and solving problems). In contrast, another annotator could also argue that the tweet used a sarcastic tone to express a negative sentiment towards over budget issue of HS2.

### Topic modelling challenges

Text documents are combinations of probabilistic distributions of topics, and each topic is a probabilistic distribution of words. However, tweets are short microblogs with character limitations (280 characters), which usually contain one topic. Therefore, LDA may have problems in calculating the probabilistic distribution of topics in tweets data. The performance of tweets topic modelling could be improved with the neural topic models, leveraging deep generative models [49]. Future research on public opinion evaluation with social media data could use Bayesian networks. In particular, gamma-belief networks showed promising results in yielding structure topics [50].

## Conclusion

This study utilised tweets data from the HS2 project as a case study. The tweets data were used to compare the performance of the proposed RoBERTa-BiGRU model with MNB and VADER. RoBERTa-BiGRU showed the best performance in terms of accuracy and ROC curves. Additionally, the study employs LDA to uncover key topics within the tweet corpus. This analysis enhances understanding of the prominent themes surrounding the HS2 project. The insights derived from the HS2 case study results lay the foundation for a public opinion evaluation framework. This framework, driven by social media data, is an invaluable tool for policymakers to evaluate public sentiment effectively. Overall, this study contributes to the field of public opinion evaluation by introducing a hybrid model, presenting a comprehensive case study analysis, and proposing a practical framework for public opinion evaluation.

## Funding

Not applicable to this article.

## Authorship contribution

This research was initially conducted as part of the requirements for the MSc in Civil Engineering at University College London. Mr Ruiqiu Yao was supervised by Dr Andrew Gillen for his MSc dissertation. The general topic and use of social media data were proposed by Dr Gillen, and they met regularly to discuss the research process. Mr Yao conducted the literature review as well as the data collection and analysis, identifying relevant sources of data and analytical tools. Mr Yao drafted the manuscript and Dr Gillen provided feedback on drafts.

## Open data and materials availability statement

The datasets generated during and/or analysed during the current study are available in the repository: <https://github.com/RY7415/OpinionAnalysisSocialMedia>. This includes the collected data (anonymised) and the Python source code.

## Declarations and conflicts of interest

### Research ethics statement

The authors conducted the research reported in this article in accordance with UCL Research Ethics standards.

## Consent for publication statement

The authors declare that research participants' informed consent to publication of findings – including photos, videos and any personal or identifiable information – was secured prior to publication.

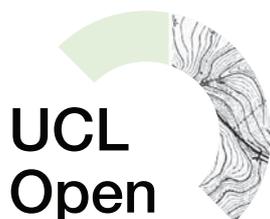
## Conflicts of interest statement

The authors declare no conflicts of interest with this work.

## References

- [1] HM Treasury. National infrastructure strategy. 2020. [Accessed 31 May 2023]. Available from: <https://www.gov.uk/government/publications/national-infrastructure-strategy>.
- [2] Hayes DJ. Addressing the environmental impacts of large infrastructure projects: making 'mitigation' matter. *Environ Law Rep.* 2014;44:10016. <https://heinonline.org/HOL/Page?handle=hein.journals/elrna44&id=18&collection=journals&index=#>. [Accessed 31 May 2023].
- [3] O'Faircheallaigh C. Public participation and environmental impact assessment: purposes, implications, and lessons for public policy making. *Environ Impact Assess Rev.* 2010;30(1):19–27. <https://doi.org/10.1016/j.eiar.2009.05.001>.
- [4] Checkoway B. The politics of public hearings. *J Appl Behav Sci.* 1981;17(4):566–82. <https://doi.org/10.1177/002188638101700411>.
- [5] Heberlein T. Some observations on alternative mechanisms for public involvement: the hearing, public opinion poll, the workshop and the quasi-experiment. *Nat Resour J.* 1976;16(1):197–212.
- [6] Ding Q. Using social media to evaluate public acceptance of infrastructure projects. Thesis. University of Maryland; 2018. <https://doi.org/10.13016/M27M0437D>.
- [7] O'Connor B, Balasubramanian R, Routledge BR, Smith NA. From tweets to polls: linking text sentiment to public opinion time series. In: Proceedings of the fourth international AAAI conference on weblogs and social media; 23–26 May 2010, Washington, USA. 2010.
- [8] Jiang H, Qiang M, Lin P. Assessment of online public opinions on large infrastructure projects: a case study of the Three Gorges Project in China. *Environ Impact Assess Rev.* 2016;61:38–51. <https://doi.org/10.1016/j.eiar.2016.06.004>.
- [9] Kaplan AM, Haenlein M. Users of the world, unite! The challenges and opportunities of social media. *Bus Horiz.* 2010;53(1):59–68. <https://doi.org/10.1016/j.bushor.2009.09.003>.
- [10] Park SB, Ok CM, Chae BK. Using Twitter data for cruise tourism marketing and research. *J Travel Tour Mark.* 2016;33(6):885–98. <https://doi.org/10.1080/10548408.2015.1071688>.
- [11] Aldahawi HA. Mining and analysing social network in the oil business: twitter sentiment analysis and prediction approaches. Cardiff University; 2015. [Accessed 31 May 2023]. <https://orca.cardiff.ac.uk/id/eprint/85006/1/2015aldahawihphd.pdf.pdf>.
- [12] Kim DS, Kim JW. Public opinion sensing and trend analysis on social media: a study on nuclear power on Twitter 1. *Int J Multimedia Ubiquitous Eng.* 2014;9(11):373–84. <https://doi.org/10.14257/ijmue.2014.9.11.36>.
- [13] Taboada M, Brooke J, Tofiloski M, Voll K, Stede M. Lexicon-based methods for sentiment analysis. *Comput Linguist.* 2011;37(2):267–307. [https://doi.org/10.1162/COLI\\_a\\_00049](https://doi.org/10.1162/COLI_a_00049).
- [14] Ku LW, Liang YT, Chen HH. Opinion extraction, summarization and tracking in news and Blog Corpora. 2006. [Accessed 31 May 2023]. <https://aaai.org/papers/0020-opinion-extraction-summarization-and-tracking-in-news-and-blog-corpora/>.
- [15] Dong Z, Dong Q. HowNet – a hybrid language and knowledge resource. In: NLP-KE 2003 – 2003 international conference on natural language processing and knowledge engineering, proceedings. Institute of Electrical and Electronics Engineers Inc; 26–29 October 2003, Beijing, China; 2003. p. 820–24. <https://doi.org/10.1109/NLPKE.2003.1276017>.
- [16] Bahdanau D, Cho K, Bengio Y. Neural machine translation by jointly learning to align and translate. [online]. September 2014. [Accessed 31 May 2023]. Available from: <http://arxiv.org/abs/1409.0473>.
- [17] Vaswani A, Shazeer N, Parmar N, Uszkoreit J, Jones L, Gomez AN, et al. Attention is all you need. [online]. June 2017. [Accessed 31 May 2023]. Available from: <http://arxiv.org/abs/1706.03762>.
- [18] Kaplan J, McCandlish S, Henighan T, Brown TB, Chess B, Child R, et al. Scaling laws for neural language models. [online]. January 2020. [Accessed 31 May 2023]. Available from: <http://arxiv.org/abs/2001.08361>.
- [19] Brown TB, Mann B, Ryder N, Subbiah M, Kaplan J, Dhariwal P, et al. Language models are few-shot learners. [online]. May 2020. [Accessed 31 May 2023]. Available from: <http://arxiv.org/abs/2005.14165>.
- [20] OpenAI. GPT-4 technical report. [online]. March 2023. [Accessed 31 May 2023]. Available from: <http://arxiv.org/abs/2303.08774>.
- [21] Devlin J, Chang MW, Lee K, Toutanova K. BERT: pre-training of deep bidirectional transformers for language understanding. [online]. October 2018. [Accessed 31 May 2023]. Available from: <http://arxiv.org/abs/1810.04805>.
- [22] Liu Y, Ott M, Goyal N, Du J, Joshi M, Chen D, et al. RoBERTa: a robustly optimized BERT pretraining approach. [online]. July 2019. [Accessed 31 May 2023]. Available from: <http://arxiv.org/abs/1907.11692>.
- [23] Cho K, van Merriënboer B, Gulcehre C, Bahdanau D, Bougares F, Schwenk H, et al. Learning phrase representations using RNN encoder–decoder for statistical machine translation. [online]. June 2014. [Accessed 31 May 2023]. Available from: <http://arxiv.org/abs/1406.1078>.
- [24] Bayes T. LII. An essay towards solving a problem in the doctrine of chances. By the late Rev. Mr. Bayes, F. R. S. communicated by Mr. Price, in a letter to John Canton, A. M. F. R. S. *Philos Trans R Soc Lond.* 1763;53:370–418. <https://doi.org/10.1098/rstl.1763.0053>.
- [25] Jiang L, Cai Z, Zhang H, Wang D. Naive Bayes text classifiers: a locally weighted learning approach. *J Exp Theor Artif Intell.* 2013;25(2):273–86. <https://doi.org/10.1080/0952813X.2012.721010>.
- [26] Zhang H. The optimality of naive Bayes. In: Proceedings of the seventeenth international Florida Artificial Intelligence Research Society conference. [online]. May

- 2004, Florida, USA: AAAI Press; 2004. p. 562–567. [Accessed 31 May 2023].
- [27] Manning CD, Raghavan P, Schütze H. Introduction to information retrieval. Cambridge: Cambridge University Press; 2008. <https://doi.org/10.1017/CBO9780511809071>.
- [28] Zhang A, Lipton ZC, Li M, Smola AJ. Dive into deep learning. Cambridge: Cambridge University Press; 2021.
- [29] Ba JL, Kiros JR, Hinton GE. Layer normalization. [online]. July 2016. [Accessed 31 May 2023]. Available from: <http://arxiv.org/abs/1607.06450>.
- [30] Tan KL, Lee CP, Anbananthen KSM, Lim KM. RoBERTa-LSTM: a hybrid model for sentiment analysis with transformer and recurrent neural network. IEEE Access. 2022;10:21517–25. <https://doi.org/10.1109/ACCESS.2022.3152828>.
- [31] Hochreiter S, Schmidhuber J. Long short-term memory. Neural Comput. 1997;9(8):1735–80. <https://doi.org/10.1162/neco.1997.9.8.1735>.
- [32] Chung J, Gulcehre C, Cho K, Bengio Y. Empirical evaluation of gated recurrent neural networks on sequence modeling. In: NIPS 2014 Workshop on Deep Learning, 13 December 2014, Montreal, Canada. 2014.
- [33] Gneiting T, Raftery AE. Strictly proper scoring rules, prediction, and estimation. J Am Stat Assoc. 2007;102(477):359–78. <https://doi.org/10.1198/016214506000001437>.
- [34] Deerwester S, Dumais ST, Furnas GW, Landauer TK, Harshman R. Indexing by latent semantic analysis. J Am Soc Inf Sci. 1990;41(6):391–407. [https://doi.org/10.1002/\(SICI\)1097-4571\(199009\)41:6<391::AID-AS11>3.0.CO;2-9](https://doi.org/10.1002/(SICI)1097-4571(199009)41:6<391::AID-AS11>3.0.CO;2-9).
- [35] Kanatani K. Linear algebra for pattern processing projection, singular value decomposition, and pseudoinverse. San Rafael, CA: Morgan & Claypool; 2021.
- [36] Blei DM, Ng AY, Jordan MI. Latent dirichlet allocation. J Mach Learn Res. 2003;3:993–1022.
- [37] Xiao C, Zhang P, Art Chaovaitwongse W, Hu J, Wang F. Adverse drug reaction prediction with symbolic latent Dirichlet allocation. In: Proceedings of the 31st AAAI conference on artificial intelligence, 4–9 February 2017. San Francisco, CA. <https://doi.org/10.1609/aaai.v31i1.10717>.
- [38] Chuang J, Ramage D, Manning CD, Heer J. Interpretation and trust: designing model-driven visualizations for text analysis. In: Conference on human factors in computing systems – proceedings. 5–10 May 2012, Texas, USA. 2012. p. 443–52. <https://doi.org/10.1145/2207676.2207738>.
- [39] Sievert C, Shirley K. LDAvis: a method for visualizing and interpreting topics. In: Proceedings of the workshop on interactive language learning, visualization, and interfaces. Association for Computational Linguistics (ACL); June 2014, Maryland, USA. p. 63–70. <https://doi.org/10.3115/v1/w14-3110>.
- [40] Department for Transport. Rail factsheet 2019. 2019. [Accessed 31 May 2023]. <https://www.gov.uk/government/statistics/rail-factsheet-2019>.
- [41] HS2 Ltd. High-speed network map. 2023. [Accessed 31 May 2023]. Available from: <https://www.hs2.org.uk/the-route/high-speed-network-map/>.
- [42] Haddi E, Liu X, Shi Y. The role of text pre-processing in sentiment analysis. In: Procedia computer science. Elsevier B.V.; January 2013. p. 26–32. <https://doi.org/10.1016/j.procs.2013.05.005>.
- [43] Hutto CJ, Gilbert E. VADER: a parsimonious rule-based model for sentiment analysis of social media text. In: Proceedings of the eighth international AAAI conference on weblogs and social media. 1–4 June 2014, Michigan, USA. [online]. 2014. p. 216–25. [Accessed 31 May 2023].
- [44] Davis J, Goadrich M. The relationship between precision-recall and ROC curves. In: Proceedings of the 23rd international conference on machine learning; 25–29 June 2006, Pennsylvania, USA. 2006.
- [45] Fawcett T. An introduction to ROC analysis. Pattern Recognit Lett. 2006;27(8):861–74. <https://doi.org/10.1016/j.patrec.2005.10.010>.
- [46] Rozema JG, Bond AJ. Framing effectiveness in impact assessment: discourse accommodation in controversial infrastructure development. Environ Impact Assess Rev. 2015;50:66–73. <https://doi.org/10.1016/j.eiar.2014.08.001>.
- [47] Rehurek R, Sojka P. Software framework for topic modelling with large corpora. In: Proceedings of the LREC 2010 workshop on new challenges for NLP frameworks. May 2010, Malta. 2010. p. 46–50.
- [48] Callison-Burch C. Fast, cheap, and creative: evaluating translation quality using Amazon’s Mechanical Turk. In: EMNLP 2009 – Proceedings of the 2009 conference on empirical methods in natural language processing: a meeting of SIGDAT, a Special Interest Group of ACL, Held in Conjunction with ACL-IJCNLP 2009, August 2009, Singapore. p. 286–95.
- [49] Zhao H, Phung D, Huynh V, Jin Y, Du L, Buntine W. Topic modelling meets deep neural networks: a survey. In: Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence; 19–27 August 2021, Montreal, Canada. 2021.
- [50] Zhang H, Chen B, Guo D, Zhou M. WHAI: Weibull hybrid autoencoding inference for deep topic modeling. In: Proceedings of 6th International Conference on Learning Representations. 3 May 2018, Vancouver, Canada. 2018.



UCLPRESS

### Extra information

*UCL Open: Environment* is an open scholarship publication, all previous versions and open peer-review reports can be found online in the *UCL Open: Environment* Preprint server at [ucl.scienceopen.com](http://ucl.scienceopen.com)