

Detection and Mitigation of Political Bias in Natural Language Processing: A Literature Review

Anonymous ACL submission

Abstract

With the increasing importance of Natural Language Processing (NLP) tools, their implications on the propagation of societal biases become more and more relevant. In this context, the analysis of political bias in manually written and automatically generated text is a relatively understudied field. Political bias refers to the preference or prejudice towards one political ideology over another. To increase the discourse in this subject area, we analyze contemporary studies on detecting and mitigating political bias in this literature review. We further discuss the benefits and potential drawbacks of the considered methods and look at the ethical considerations involved with political bias in NLP, before we give suggestions for future studies.

1 Introduction

With the rising integration of Natural Language Processing (NLP) models in everyday applications, their effects on the propagation of societal biases become increasingly relevant. In this context, not only political bias detection but also mitigation approaches are needed to expose and alleviate such biases.

The definition of political bias across studies is inconsistent. According to [Chen et al. \(2020\)](#), politically unbiased means to "report on an event without taking a political position, characterization, or terminology" while [Gangula et al. \(2019\)](#) define it as not "selectively publishing articles to specifically choosing to highlight some events, parties and leaders". In the review at hand, we base our political bias definition on that of [Chen et al. \(2020\)](#). Political objectivity in our context hence means that an event is being reported without taking a political stance and without adapting an ideology-specific terminology, i.e. to be politically unbiased, as well as fair with regards to the

report of original facts rather than opinionated statements ([Chen et al., 2020](#); [Ad Fontes Media, 2021](#)). For example, a phrase like 'death tax' would be considered politically biased towards the conservative party in the United States where the term is used to describe a tax that is imposed on property that gets transferred to another person after the owner's death. On the other hand, liberals, who are in favor of this concept, call it 'estate tax' ([Graetz, 2016](#)).

Looking at bias in a machine-learning context, previous studies found that it can be exhibited in multiple components in NLP systems such as pre-trained word embeddings or training data ([Zhao et al., 2018](#); [Bolukbasi et al., 2016](#); [Caliskan et al., 2017](#)). In this paper, we are specifically interested in which detection and mitigation approaches exist to deal with such biased sub-components and to prevent the amplification of political bias through text.

The topic is relevant because objective reporting is necessary for an unbiased societal discourse on potentially controversial topics that in turn can shape the political agenda as well as corresponding initiatives on a national and international level ([Dardis et al., 2008](#)). While NLP models can be used to identify political bias in human-generated text, they can also be the source of said bias in generative language models. Especially given the large amount of data online, automatic detection and mitigation methods are necessary since the manual identification of political bias becomes increasingly infeasible.

In the paper at hand, we elaborate on the occurrence of political bias, corresponding ethical considerations as well as the necessity for future research in the field in Section 2. Due to the different nature of the detection and mitigation

041
042
043
044
045
046
047
048
049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081

task in NLP, we decided to split our corresponding review into two separate sections: In Section 3, we analyze approaches to detect political bias before we examine the current state of political bias mitigation in NLP in Section 4. An overview of the considered studies is given in Table 1. We conclude our work with an overview of future research directions in Section 5.

We make two **contributions** to the current research:

- To the best of our knowledge, we put together the first review of political bias in NLP which builds a basis for future discussions in the field.
- We critically discuss current detection as well as debiasing methods to identify optimization potential and future research directions.

2 Background

Recently, ethical implications of bias in NLP have been the focus of research efforts (Sheng et al., 2020; Bordia and Bowman, 2019). Politically biased texts can be both human- and machine-made. In both cases, the consumers of such texts can be influenced in their decisions and perceptions of the world and hence need to be aware of potential ideology-specific tendencies exhibited by such texts. In addition, consumers need to be able to access unbiased, fair reports about their topic of interest.

Any form of bias can be categorized into two different categories: allocation or representation bias (Crawford, 2017). The former means that certain groups are being preferred in the allocation of resources. In an NLP context, this occurs when models perform better on the majority data. On the other hand, representation bias occurs when considered subgroups (e.g., a specific political ideology) are associated with specific concepts in parameters or embeddings. Both the application and the representation bias can deepen political misconceptions and hence have implications for the national political agenda and respective initiatives. Hence, the increasing use of NLP models poses the risk of propagating and amplifying damaging stereotypes in society.

The most prominent example of an area where political bias can occur is the media. At the article level, published texts might implicitly convey the author’s or the news outlet’s political ideology, i.e., exhibit a right or left bias. In an extreme form of said ideology-specific tendencies, articles can be classified as propaganda (Rashkin et al., 2017; Da San Martino et al., 2019). Any form of media could be biased so that people are not aware of it. For example, word choices and the selective or misrepresentative reporting of events can influence the reader’s perception. A relevant ethical consideration in this context is whether, and if so, in what way, politically biased reporting should be exposed, a) to allow media organizations to stay credible and b) to give people the control over which content they consume and which texts influence their opinions. In this context, the political bias of a medium is further essential to detect so-called fake news (Horne et al., 2018), or to fact-check a claim (Nguyen et al., 2018) and hence to ensure that the reader is either informed about the reliability of the respective source or that the bias is mitigated in the first place.

Political bias is also relevant in the virtual space: In online communities such as social networking sites, complex profiling of users that include psychological characteristics, demographics and meta-data has occurred. Such profiles were subsequently used to micro-target users with politically biased content to gain some form of political advantage (Lazer et al., 2018; Vosoughi et al., 2018). Another aspect of the online realm impacted by political bias is hate speech detection. Hate speech has become an increasing issue in online communities, and detection methods for this phenomenon were developed. However, said models can be impacted by undesired political bias in the training data, which can negatively impact the performance of hate speech classifiers (Wich et al., 2020). This unfavorable effect, in turn, can lead to issues regarding the freedom of speech or to hampering the social discourse if articles are falsely identified as politically-biased hate speech. On the other hand, false identification as non-hate-speech could negatively impact the attacked people, so both misclassifications need to be addressed.

Finally, large-scale language models have re-

Study	Purpose	Method	Data Source(s)
Iyyer et al. (2014)	Detection	RNN	Convote & subset of the Ideological Books Corpus
Chen et al. (2017)	Detection	Opinion-aware Knowledge Graph	Convote & subset of the Ideological Books Corpus & a collection of political Tweets
Jiang et al. (2019)	Detection	CNN & Batch Normalization	automatically labeled articles from Kiesel et al. (2019) & a collection of manually labeled articles
Chen et al. (2020)	Detection	RNN & Reverse Feature Analysis	allsides.com & adfontesmedia.com & a collection of manually labeled articles
Baly et al. (2020)	Detection	Multi-Task Ordinal Regression	allsides.com
Liu et al. (2021)	Mitigation	Reinforced Calibration	Media Cloud API & survey data from the Pew Research Center

Table 1: Overview of the studies considered along with their purpose, the respective employed method as well as the data source(s) used.

cently been the focus of research efforts to advance human-like text generation (Zhang et al., 2020; Peng et al., 2020). Other applications of such models are machine translations (Zhu et al., 2020). Given that these language models have been trained on sizeable unsupervised text corpora – for example, GPT-2 (Radford et al., 2019) was trained on 8 million web pages –, they can potentially inherit the (political) bias that was present in the training data and propagate it in the subsequently generated text. This propagation can lead to the amplification of political bias through such models in society, and hence a potentially unethical influence on public opinion.

Our paper is the first review of methods to detect and mitigate political bias in NLP. We provide a basis for future discussions and suggest research directions to advance the current state of the field.

3 Political Bias Detection

As seen in the last section, political bias is a phenomenon that can affect people in their opinion formation. Due to the increasing volume of distributed text online, people are more exposed to politically biased work. In addition, the rate of information dissemination in the online realm is much faster, and manual detection of political bias

is not feasible in most cases. For this reason, methods to support an automatic bias detection process have been explored in NLP. They are the focus of this chapter. However, to date, there is no standardized data set for politically biased text, and hence the considered papers in our review all rely on individually constructed corpora, limiting the direct comparability of these studies.

Recursive Neural Networks

Iyyer et al. (2014) created a balanced data set by subsampling to account for label imbalances. They used a filtered subset of texts from two sources: the U.S. Congressional floor debate (Convote) data set (Thomas et al., 2006) and a manually labeled subset of the Ideological Books Corpus (IBC) (Gross et al., 2013). The former includes transcripts of debates from the U.S. Congress in 2005 labeled with the speaker’s parties (Democrat, Republican, or Independent), which was taken as a proxy for the text’s political bias. The modified IBC data set included texts written by authors with well-known political leanings. Iyyer et al. (2014) subsequently hired crowdsourcers to obtain annotations on a 3-point scale (left, neutral, right) for these texts on the sentence and phrase level. While they only included sentences on which at

233 least two labelers agreed, this approach introduces
234 an uncertainty since the crowd workers were not
235 specialized in political science and hence might
236 have either propagated their views in the labeling
237 process or misjudged the presence or direction of
238 political bias in the data. Post labeling, the authors
239 trained and tested their recursive neural network
240 (RNN) architectures on the two different data sets
241 and found better results for the one with shorter
242 sentences and more training data. They suggested
243 that that is likely the case because a) more training
244 data brings significant improvements for RNN and
245 b) information is lost at every propagation step, i.e.,
246 the meaning of shorter sentences is captured easier
247 than that of longer ones. Their best performing
248 RNN reached an accuracy of 70.2 %.

250 [Chen et al. \(2020\)](#), constructed a binary-labeled
251 corpus at the document-level from [allsides.com](#) and
252 [adfontesmedia.com](#), two platforms that provide
253 assessments of news articles' topics, political bias,
254 and unfairness. Compared to the work by [Iyyer
255 et al. \(2014\)](#), the authors hence only evaluated
256 whether political bias is present, but not whether
257 the text shows a left or right tendency. On the
258 one hand, this makes the assessment of bias more
259 reliable since it is easier to determine whether
260 something is politically biased than the task of
261 additionally identifying the direction of the bias.
262 At the same time, however, this causes information
263 about the political tendency of the text to be
264 omitted. The authors subsequently approached the
265 exposure of imbalanced news coverage with an
266 RNN. Their choice was motivated by the fact that
267 such networks can capture syntactic and semantic
268 composition when provided with textual input that
269 keeps the word order. This capturing is possible
270 by considering the hierarchical nature of language:
271 Each word in a sentence is represented as a vector,
272 and so are higher-level linguistic constructs like
273 phrases and sentences with the exact vector dimen-
274 sions as the words they are built on. That way,
275 the underlying vector representations are trained
276 to retain the meaning of a sentence ([Iyyer et al.,
277 2014](#)). For example, if a vector represents a liberal
278 linguistic construct, i.e., a phrase or a sentence, it
279 should significantly differ from the corresponding
280 vector representation of a right-wing sentence.
281 This property of RNN is especially relevant when
282 identifying more advanced social constructs like
283 political bias, which are only identifiable at higher

284 levels of sentence structures rather than at the word
285 level.

286
287 To avoid the learning of media-outlet-specific
288 features, [Chen et al. \(2020\)](#) removed portal-
289 identifying information from the text in their
290 study. The authors achieved an accuracy of 75.42
291 % for political non-objectivity detection. When
292 considering individual results, it can be noted
293 that the prediction of objective articles in this
294 research tended to be more accurate than the
295 prediction of non-objective articles, presumably
296 due to the uneven distribution of biased and
297 non-biased articles in the training data. Compared
298 to the previous study by [Iyyer et al. \(2014\)](#), [Chen
299 et al. \(2020\)](#) did not create a balanced data set to
300 account for label imbalances, which might be an
301 explanation for this outcome.

302
303 An advantage of RNNs, in general, is that seman-
304 tic information of close-by words, but also of con-
305 structs that are further apart, are detected. However,
306 this mechanic only works with sufficient training
307 data, as was suggested by the finding of [Iyyer et al.
308 \(2014\)](#) that better results were obtained with the
309 more extensive data set. For example, the construct
310 'should not be used as an instrument to achieve char-
311 itable or social ends' got misinterpreted by their
312 network as non-biased instead of being liberally
313 oriented because formulations with 'should not'
314 did not appear often enough in the training data
315 for the RNN to pick up on it. Another issue that
316 needs to be taken into account is that the semantic
317 information that the network captures depends on
318 the text's overall context. Sarcasm and idioms will
319 most likely not be correctly detected by the RNN
320 architectures in the two studies considered.

321 **Opinion-aware Knowledge Graph**

322 A different approach was taken by [Chen et al.
323 \(2017\)](#), who created an opinion-aware knowledge
324 graph. Specifically, they used a background knowl-
325 edge graph ([Bizer et al., 2009](#)) containing enti-
326 ties and semantic relations and infused it with
327 ideology-specific training data to estimate opinions
328 expressed towards entities in the graph as senti-
329 ment distributions over two ideological categories
330 (conservative vs. liberal). In the next step, the
331 opinion distributions were propagated based on the
332 semantic relations between the entities in the graph.
333 The final opinion-aware knowledge graph was then

used to detect the political ideology of the test data set by matching the test entities with the entities in the constructed graph and inferring the respective political orientation from these entities. The authors' graph was built on three different data sets: In addition to the Convote data set (Thomas et al., 2006) and the IBC (Iyyer et al., 2014) that were also used by Iyyer et al. (2014) (see Section 3 – Recursive Neural Networks), they added Tweets annotated for political bias. For that, they took a list by Bakshy et al. (2015) that contained media outlets with their respective ideological leanings. They subsequently found the corresponding Twitter accounts of these organizations and labeled Tweets from these accounts with the ideology of the respective source.

Such a knowledge graph has the advantage that factual and subjective information can be used for a joint inference based on texts and knowledge bases to detect the political bias of a sentence or document. This is supported by the results that Chen et al. (2017) achieved: The accuracies for their best-performing RNN and support vector machines (SVM) on the data were 70% and 76% respectively, while their knowledge graph achieved an accuracy of 81%.

ELMo Sentence Representation Convolutional Neural Network

Compared to the previously discussed approaches, Jiang et al. (2019) introduced an Embeddings from Language Model (ELMo) sentence representation convolutional neural network (CNN) to identify left- or right-wing hyperpartisanship. They first calculated sentence-level embeddings as the mean of ELMo word embeddings to represent documents as sequences of these sentence-level embeddings, which were subsequently used in a CNN to predict the political orientation. As part of the *SemEval-2019 Task 4: Hyperpartisan News Detection* competition (Kiesel et al., 2019), the authors were given two different data sets: Firstly, one that encompassed 750k articles that were classified by the political bias of the respective news source they were collected from. To obtain the source's bias, the organizers of the competition cross-checked two public media bias lists from BuzzFeed and Media Bias Fact Check (Kiesel et al., 2019). In addition, the second data set they provided included 645 manually labeled articles. For these articles, the bias was rated on a 5-point Likert scale by three

annotators (Vincent and Mestre, 2018).

Notably in this study is that Jiang et al. (2019) found that their best-performing model was only trained on the manually labeled articles while including the by-publisher data set worsened the accuracy on the test set. The authors achieved an accuracy of 82 % on the held-out test data in the competition (only using the manually labeled set) vs. an accuracy of 64 % for the model trained on the articles classified by publisher. This result indicates that the bias of a news outlet isn't necessarily propagated to the articles by the respective publishers – which amplifies the need for a data set labeled on the article level.

LIWC and Reverse Feature Analysis

Based on their previously described RNN model, Chen et al. (2020) further performed a reverse feature analysis to investigate how political bias is revealed on the word- and sentence-level as well as in the overall article structure. In each iteration, they removed text parts and re-calculated the bias associated with the text. The estimated bias of the removed segment was derived by subtracting the estimated bias of the new text from the estimated bias of the old text. This approach can be viewed as an attention-based model that outputs weights indicating feature importance (Bahdanau et al., 2014).

On a word level, Chen et al. (2020) correlated the most biased sentences with Linguistic Inquiry and Word Count (LIWC) categories (Pennebaker et al., 2001). They found that especially the classes 'negative emotions', 'focus present' and 'percept' were negatively correlated with political objectivity. This result means that authors of politically biased articles tended to use opinionated and feelings-related words such as 'angry' and 'disappoint'. Chen et al. (2020) also found that unfair articles, i.e., those that only report selected facts, tended to include more words from the category 'focus present', for example, 'admit' and 'determine'.

The investigation of higher-level linguistic structures, on the other hand, yielded that the bias strength in the first and second quarters of articles tended to be comparable for objective and non-objective articles. This outcome can be explained by the fact that most articles start with

434 a high-level summary, followed by background
435 information (Pöttker and Starck, 2003; Chen et al.,
436 2020), which shows a tendency to be written in a
437 neutral tone. The biased nature of articles typically
438 shows in later parts of the text, especially in the
439 last quarter. Chen et al. (2020) further found that
440 it was easiest to detect unfair articles, i.e., those
441 in which selected facts were reported in favor of
442 a party. According to the authors, this could be
443 the case because the word usage in such articles
444 tended to be more emotional, both with regards to
445 positive and negative feelings. This finding made
446 it easier to recognize the underlying bias.

447
448 An advantage of the unsupervised reverse feature
449 analysis proposed by the authors is that the unit to
450 be analyzed does not have to be defined before
451 the training of the model, as it is the case with
452 most other attention models (Chen et al., 2020). In
453 addition, the knowledge of how biased articles are
454 structured can make the detection of political bias
455 in these texts more efficient since future models
456 could be trained to focus on the last one or two
457 quarters of a text, for example. Looking at the
458 proposed word-level analysis, however, it can be
459 observed that most LIWC categories did not exhibit
460 a strong correlation. In general, words captured by
461 LIWC are limited since it is a human-made lexicon
462 that might not capture the most revealing terms for
463 political bias. Furthermore, the categories are only
464 used at the word level. However, more complex
465 constructs like political bias tend to show at higher-
466 level text granularity (Chen et al., 2018; Iyyer et al.,
467 2014), e.g., in phrases or sentences, and hence a
468 standalone LIWC analysis is not as meaningful in
469 the detection of political bias.

470 **Multi-Task Ordinal Regression**

471 Baly et al. (2020) proposed a multi-task ordinal
472 regression to detect the bias of entire news outlets
473 in combination with a trustworthiness estimation.
474 Specifically, the authors modeled the left-right bias
475 on a 7-point scale (extreme-left, left, center-left,
476 center, center-right, right, extreme-right) and
477 factuality, which has been used as a substitute
478 for trustworthiness by the authors, on a 3-point
479 scale (low, mixed, high). Their approach was
480 motivated by the observation that center media
481 tends to be more impartial than hyperpartisanship
482 media, which tends to be more emotional, i.e., less
483 factual reporting. Compared to previous studies

484 that looked at the detection of trustworthiness and
485 political bias independently, Baly et al. (2020)
486 reported significant performance improvements
487 for the joint model. They collected multiple
488 articles from the target medium and derived
489 part-of-speech tags, linguistic features as well as
490 word embeddings. The authors used a model to
491 approximate the learning of the joint probability
492 density function between political bias and
493 factuality. They found a joint model in which
494 political bias is considered on 3- and 5-point scales
495 as auxiliary tasks yielded the best performance
496 at a mean absolute error of 1.475. An accuracy
497 score to compare this approach to the ones pre-
498 sented in the previous subsections was not reported.

499
500 One limitation of the study by Baly et al. (2020)
501 is the fact that they evaluated the political bias of
502 entire news outlets. While their study is based on
503 sample articles from each outlet, their final results
504 refer to the outlet itself. However, the evaluated
505 bias of the outlet does not necessarily reflect the
506 bias of future articles. Furthermore, the 7-point
507 Likert scale used for classifying political bias goes
508 beyond the universal left-right classification and
509 can exhibit more regional idiosyncrasies (Tavits
510 and Letki, 2009), which could reduce the validity
511 of the results.

512 **4 Political Bias Mitigation**

513 The mitigation of political bias is a new field with
514 little published research up to date. The most
515 promising study to decrease political bias was pub-
516 lished by Liu et al. (2021), who proposed an ap-
517 proach called 'Reinforced Calibration' for automati-
518 cally generated text. This method is also the only
519 work on political bias mitigation we are aware of
520 to date.

521 **Reinforced Calibration**

522 When generating text based on language models,
523 text prompts like 'I think about marijuana because'
524 are used. Liu et al. (2021) found that attributes
525 such as gender, location, or topic have a significant
526 influence on the political bias of the subsequently
527 generated text. For example, for the sample
528 sentence above, a GPT-2 language model generates
529 the liberally-biased supplement 'I believe it should
530 be legal and not regulated.'. A noteworthy aspect
531 that the authors found was that even conservative
532 prompts were completed with liberal output by

GPT-2.

To approach the task of mitigating such political bias in generated text, Liu et al. (2021) kept the main GPT-2 architecture but added a debiasing stage, with which the original text generation was re-calibrated to produce unbiased output based on multiple steps of "reinforced optimization" (Liu et al., 2021). They defined a state at step t as all previously generated tokens and an action as the next output token. The policy in this context was the softmax output of the last hidden state, as this could be taken as the probability to choose a specific token, i.e., an action in this reinforcement learning setting, according to Dathathri et al. (2019). The authors further used a debias reward to guide the reinforced optimization. In this context, they employed two different rewards: a word-embedding-guided and a classifier-guided debias reward.

A word-embedding-guided debias reward was used in previous studies to force what are considered neutral words to be equally apart from topic-sensitive words in the embedding space, e.g., gender (Zhao et al., 2018; Park et al., 2018; Bolukbasi et al., 2016). Liu et al. (2021) used this approach to pick the next unbiased token at each time step. However, an issue with this approach is that political bias tends to occur at higher granularity levels (see Section 3 for more details) instead of at the word level. Furthermore, this approach is dependent on the quality of previously defined political bias words, which can have a significant impact on the final results (Zhou et al., 2019; Liu et al., 2021).

The classifier-guided debias that the authors additionally employed helped alleviate these issues. It was based on two different auxiliary tasks: Firstly, a political bias classifier was used to evaluate whether the text at hand was objective or not. Secondly, a constraint was introduced in the form of the Kullback-Leibler divergence between the original and the newly debiased policy to regulate the shift away from the vanilla softmax output, which might cause limited semantic coherence. Both components were balanced in the process of the reinforced optimization.

Liu et al. (2021) found that with regards to the

considered attributes in the prompts, Reinforced Calibration was able to reduce the political bias in the generated text while maintaining readability. Comparing their debiasing results, they further found that the word-embedding debias reward led to worse performance than the classifier-guided debias reward.

An advantage of this approach is that the underlying language model does not have to be accessed or retrained; instead, an additional debiasing layer can be added, significantly reducing the necessary computing power and time. While the authors used the GPT-2 architecture in the paper as a base, the idea is also easily expandable to other language models through the addition of the debias stage. However, a drawback is that in this study, the focus was only on binary outputs, i.e., left or right ideologies, and an extension to more fine-grained political bias distinctions is non-trivial. Another aspect to consider regarding the mitigation of political bias is that through the additional layer of the Reinforced Calibration approach, additional noise is introduced, which might cause the overall performance of the respective NLP model to decline.

5 Future Research Directions

This paper highlighted ethical implications of political bias in text and summarized contemporary studies that focus on the detection and mitigation of political bias in NLP. We further analyzed the advantages as well as drawbacks of the individual methods.

So far, the limited existing approaches have not been evaluated in a unified framework. This paper addressed this gap to allow for a more exhaustive discourse of the topic at hand. We also found that while multiple authors have addressed political bias detection, the mitigation of such bias remains understudied. In this final section, we present a non-exhaustive overview of how to address the most severe shortcomings in the research, which we identified in our review to foster research in the area of political bias in NLP.

Standardized Definitions, Benchmarks, and Data

Due to the relative recentness of the subject, standardized definitions, evaluation metrics, and benchmarks are missing to measure political bias in

633 text. While we recognize that different applications
634 might require different standards, this area should
635 be addressed in future research. Especially the use
636 of different data sets and labeling instructions
637 for politically biased text limits the comparability
638 of contemporary studies. This issue is aggravated
639 because political bias is evaluated on different lev-
640 els: Some authors consider political bias on a news
641 outlet, some on an article, and some on a sentence
642 level. This divergence ties in with the lack of a stan-
643 dardized gold-standard political bias data set at the
644 sentence level, limiting the progression of research
645 in the field and should therefore be addressed.

646 **Non-binary Political Bias**

647 In all reviewed studies, the political spectrum con-
648 sidered was limited. Most studies focused on a
649 binary left-right classification of political partisan-
650 ship. More nuanced political ideologies were be-
651 ing disregarded. Future work could follow two
652 directions regarding this issue: In supervised ap-
653 proaches, more nuanced political ideologies could
654 be taken into account. On the other hand, unsuper-
655 vised approaches could help discover the variety of
656 political ideologies present and prevent limitations
657 through pre-defined political affiliations.

658 **Application of Bias Mitigation Techniques from** 659 **Other Bias Types**

660 Methods from other NLP bias analyses could be
661 considered to mitigate political bias in NLP tasks.
662 For example, data augmentation methods could
663 be used to decrease political bias in generated
664 text. This approach could be successful if dis-
665 proportionate class distributions in the data
666 cause political bias in NLP applications. Data
667 augmentation was previously implemented for
668 gender, and race bias (Zmigrod et al., 2019; Yucer
669 et al., 2020). In the case of gender bias, (Zhang
670 et al., 2020) augmented the training data such
671 that the gender in sentences was swapped and the
672 algorithm was trained on the combination of the
673 old and the augmented data. Kusner et al. (2017),
674 on the other hand, used an approach in which
675 data samples were treated equally in actual and
676 counterfactual demographic groups, which could
677 be extended to political partisanship, too.

679 Another approach to consider would be embed-
680 ding manipulations. Garg et al. (2018) found that
681 societal biases are reflected in word embeddings,
682 which is likely valid for political bias as well.

683 With regards to gender bias, this was studied, for
684 example, by Bolukbasi et al. (2016). The authors
685 ensured that gender-neutral word embeddings
686 were orthogonal to a gender direction defined by
687 gender-bias words selected through a classifier.
688 Zhang et al. (2020) built on this method and tried
689 to force neutral words to have an equal distance
690 to pre-defined groups of sensitive words to obtain
691 a gender-neutral embedding space. In addition,
692 Zhou et al. (2019) retrained language models with
693 a fairness loss to ensure unbiased text generation.
694

695 These approaches rely on the retraining of the un-
696 derlying language model, which is often not avail-
697 able (such as in the case of GPT-3) or computationally
698 costly. Nevertheless, a comparison between
699 complete language model retraining approaches
700 and Reinforced Calibration that focused on adding
701 a debiasing layer should be conducted to evaluate
702 the performance in both settings and assess which
703 one is more effective in mitigating political bias.

704 **Mitigating and Detecting Political Bias in** 705 **Languages Other Than English**

706 The considered studies only focused on English
707 text. In future work, existing techniques could
708 also be applied to political bias in other languages.
709 However, this is non-trivial for two reasons: Firstly,
710 especially in countries other than the U.S., the party
711 landscape is often more diverse, and the differen-
712 tiation between political camps is more nuanced,
713 which might be harder to be picked up by NLP
714 models. Secondly, most politically-oriented cor-
715 pora are English, and hence there would be a need
716 to create complementary training data. With re-
717 gards to both detection and mitigation approaches,
718 an extensive training set is salient and needs to be
719 created before considering the transfer of existing
720 approaches to other languages.

721 **6 Conclusion**

722 Political bias detection and mitigation in NLP is an
723 emerging field. Due to the increased usage of NLP
724 and its potential to propagate societal biases, it is
725 vital to address such problems early to unify efforts
726 within the research community. To the best of our
727 knowledge, this is the first review paper to address
728 the state of the research in this area. We further
729 suggested research opportunities to advance the
730 detection and mitigation of political bias in NLP
731 methods.

732
733
734
735
736

737
738
739
740

741
742
743
744

745
746
747
748

749
750
751
752
753

754
755
756
757
758
759

760
761
762

763
764
765
766

767
768
769
770

771
772
773
774

775
776
777
778
779

780
781

782
783
784

References

Ad Fontes Media. 2021. How to Rate Veracity. Website. <https://www.adfontesmedia.com/wp-content/uploads/2020/12/2-How-to-Rate-Veracity-Guide-and-Video-Script.pdf>. Accessed: 2021-07-08.

Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*.

Eytan Bakshy, Solomon Messing, and Lada A Adamic. 2015. Exposure to ideologically diverse news and opinion on facebook. *Science*, 348(6239):1130–1132.

Ramy Baly, Giovanni Da San Martino, James Glass, and Preslav Nakov. 2020. We can detect your bias: Predicting the political ideology of news articles. *arXiv preprint arXiv:2010.05338*.

Christian Bizer, Jens Lehmann, Georgi Kobilarov, Sören Auer, Christian Becker, Richard Cyganiak, and Sebastian Hellmann. 2009. Dbpedia-a crystallization point for the web of data. *Journal of web semantics*, 7(3):154–165.

Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. 2016. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. *Advances in neural information processing systems*, 29:4349–4357.

Shikha Bordia and Samuel R Bowman. 2019. Identifying and reducing gender bias in word-level language models. *arXiv preprint arXiv:1904.03035*.

Aylin Caliskan, Joanna J Bryson, and Arvind Narayanan. 2017. Semantics derived automatically from language corpora contain human-like biases. *Science*, 356(6334):183–186.

Wei Chen, Xiao Zhang, Tengjiao Wang, Bishan Yang, and Yi Li. 2017. Opinion-aware knowledge graph for political ideology detection. In *IJCAI*, pages 3647–3653.

Wei-Fan Chen, Khalid Al-Khatib, Henning Wachsmuth, and Benno Stein. 2020. Analyzing political bias and unfairness in news articles at different levels of granularity. *arXiv preprint arXiv:2010.10652*.

Wei-Fan Chen, Henning Wachsmuth, Khalid Al Khatib, and Benno Stein. 2018. Learning to flip the bias of news headlines. In *Proceedings of the 11th International conference on natural language generation*, pages 79–88.

Kate Crawford. 2017. The trouble with bias. Keynote at Neural Information Processing Systems.

Giovanni Da San Martino, Seunghak Yu, Alberto Barrón-Cedeno, Rostislav Petrov, and Preslav Nakov. 2019. Fine-grained analysis of propaganda in news

article. In *Proceedings of the 2019 conference on empirical methods in natural language processing and the 9th international joint conference on natural language processing (EMNLP-IJCNLP)*, pages 5636–5646. 785
786
787
788
789

Frank E Dardis, Frank R Baumgartner, Amber E Boyd-stun, Suzanna De Boef, and Fuyuan Shen. 2008. Media framing of capital punishment and its impact on individuals’ cognitive responses. *Mass Communication & Society*, 11(2):115–140. 790
791
792
793
794

Sumanth Dathathri, Andrea Madotto, Janice Lan, Jane Hung, Eric Frank, Piero Molino, Jason Yosinski, and Rosanne Liu. 2019. Plug and play language models: A simple approach to controlled text generation. *arXiv preprint arXiv:1912.02164*. 795
796
797
798
799

Rama Rohit Reddy Gangula, Suma Reddy Duggenpudi, and Radhika Mamidi. 2019. Detecting political bias in news articles using headline attention. In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 77–84. 800
801
802
803
804
805

Nikhil Garg, Londa Schiebinger, Dan Jurafsky, and James Zou. 2018. Word embeddings quantify 100 years of gender and ethnic stereotypes. *Proceedings of the National Academy of Sciences*, 115(16):E3635–E3644. 806
807
808
809
810

Michael J Graetz. 2016. Death tax politics. *BCL Rev.*, 57:801. 811
812

Justin H Gross, Brice Acree, Yanchuan Sim, and Noah A Smith. 2013. Testing the etch-a-sketch hypothesis: a computational analysis of mitt romney’s ideological makeover during the 2012 primary vs. general elections. In *APSA 2013 Annual Meeting Paper, American Political Science Association 2013 Annual Meeting*. 813
814
815
816
817
818
819

Benjamin D Horne, William Dron, Sara Khedr, and Sibel Adali. 2018. Assessing the news landscape: A multi-module toolkit for evaluating the credibility of news. In *Companion Proceedings of the The Web Conference 2018*, pages 235–238. 820
821
822
823
824

Mohit Iyyer, Peter Enns, Jordan Boyd-Graber, and Philip Resnik. 2014. Political ideology detection using recursive neural networks. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1113–1122. 825
826
827
828
829
830

Ye Jiang, Johann Petrak, Xingyi Song, Kalina Bontcheva, and Diana Maynard. 2019. Team berthavon suttner at semeval-2019 task 4: Hyperpartisan news detection using elmo sentence representation convolutional network. In *Proceedings of the 13th International Workshop on Semantic Evaluation*, pages 840–844. 831
832
833
834
835
836
837

Johannes Kiesel, Maria Mestre, Rishabh Shukla, Emmanuel Vincent, Payam Adineh, David Corney, 838
839

