

Instructions for EMNLP 2023 Proceedings

Anonymous EMNLP submission

Abstract

With the popularity of the Internet, online social media has increasingly become a platform for people to share their attitudes towards life, such as optimism and pessimism about the future. While the Internet embraces various views, it also quietly deepens the formation of impressions on different views and attitudes. A large part of these texts will form the corpus of the pre-trained model, and the model may learn the tendency of life attitudes in the corpus. Our work develops new methods to (1) measure life attitude biases in LMs trained on such corpora and (2) measure the judgement impact of downstream models trained on different life attitude corpus. We focus on mental health disorder detection, aiming to empirically quantify the effects of life attitude (optimism, pessimism) leaning in pretraining data on the influence of risk-related tasks. Our findings reveal that pre-trained LMs do have life attitude leanings that reinforce the polarization present in pretraining corpora, propagating life attitude biases into mental health disorder detection. We discussed strategies that might mitigate or leverage models of different life attitude leaning.

1 Introduction

With unprecedented user engagement, digital and social media have become the primary way for people to share their attitudes about life (Steinert, 2021; Shareef et al., 2020; Auxier and Anderson, 2021). Over the past decade, there has been a dramatic increase in the number of incidents where people share their attitudes on social media, which can cover a wide range of topics: jobs and careers, schooling and education, health and lifestyle, relationships and family, social issues and current events (Debatin et al., 2009; Gross and Acquisti, 2005). Although social networks provide an inclusive platform for people to share different perspectives on attitudes toward life, these discussions also deepen the formation of polarized impressions of

attitudes toward life - pessimistic and optimistic tendencies (Ferrari, 2008; Peeters and Czapinski, 1990). These texts form a major part of the pre-training corpus for the large language model and propagate this life attitude tendency to downstream tasks.

Mental health is a key issue in modern society, and without proper treatment, mental disorders can sometimes turn into suicidal ideation (Ryk et al., 2023). To address this critical issue, there has been a large amount of mental health research aimed at the efficient and automated detection of mental health disorders. Mental health disease detection usually uses information fusion strategies to make the model know more information and improve the accuracy of reasoning. There are three common fusion strategies: feature fusion (Song et al., 2018; Uban et al., 2021), model fusion (Sawhney et al., 2020; Abdul-Mageed and Ungar, 2017) and task fusion (Turcan et al., 2021). Previous research methods mainly focused on how to integrate more information, such as emotional, personality and even economic information, so that the model can learn more knowledge, while ignoring the possible potential life attitude bias of the pre-trained model itself, which is closely related to mental health detection problems (Mao et al., 2022; Conversano et al., 2010; Yildirim and Cicek, 2022). For example, an overly pessimistic model may misrepresent positive examples of some mental health disorders. To the best of our knowledge, no prior work has shown how to analyze the effects of naturally occurring life of orientation biases in pretraining data on language models, and subsequently on downstream tasks. Our study aims to fill this gap.

We focused on several most common mental health disorders: depression, anxiety, suicide, stress (Hardy, 2018). Because detecting these mental health disorders is important for the mental health of society (Ryk et al., 2023). We investigate how social media biases in the pretraining

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
082

083 data propagate into pretrain models and ultimately
084 affect downstream tasks, because discussions about
085 life orientation attitude issues are abundant in pre-
086 training data sourced from online encyclopedias,
087 and language models will inevitably inherit this
088 tendency towards life attitudes.

089 To this end, based on psychological theory
090 (Carver and Scheier, 1982; Carver et al., 1983)
091 and the life of orientation test(Scheier and Carver,
092 1985), we propose to empirically quantify the life
093 attitude leaning of pretrained LMs. We then further
094 pretrain language models on different life attitude
095 corpora to investigate whether LMs pick up political
096 biases from training data. Finally, we train
097 classifiers on top of LMs with varying life of orien-
098 tation leanings and evaluate their performance on
099 mental health disorders targeting different mental
100 illness. In this way, we investigate the propaga-
101 tion of life of orientation attitude through the entire
102 pipeline from pretraining data to language models
103 to downstream tasks.

104 Our experiments across several LM architectures
105 demonstrate that different pretrained LMs do have
106 different underlying life attitude leanings, reinforc-
107 ing the life attitude polarization present in pretrain-
108 ing corpora. However, the models with different
109 life attitude tendencies showed great differences
110 in the detection of different types of mental health
111 diseases.

112 The main contributions of this paper are novel
113 methods to quantify life of orientation attitude in
114 LMs, and findings that shed new light on how ide-
115 ological polarization in pretraining corpora propa-
116 gates orientation into language models, and subse-
117 quently into risk-related downstream tasks.

118 2 Method

119 We propose a two-step methodology to establish
120 the effect of life attitude in pretraining corpora on
121 the mental health disorders detection tasks: (1) we
122 develop a framework, grounded in physical science
123 literature, to measure the inherent life orientation
124 leanings of pretrained language models, and (2)
125 then investigate how the life orientation leanings
126 of LMs affect their performance in downstream
127 risk-oriented tasks.

128 2.1 Measuring the life orientation Leanings of 129 LMs

130 While previous works have provided an analysis
131 of life attitude tendencies in LMs, they have pri-

132 marily focused on the context of task-specific situa-
133 tions such as emotional computation, rather than on
134 timeless ideological issues based on psychological
135 literature. In contrast, our method is grounded in
136 psychological theory (Carver and Scheier, 1982;
137 Carver et al., 1983). It was able to use a set of
138 theories to test two aspects of life attitude tenden-
139 cies: optimism and pessimism, whereas previous
140 methods could only detect optimism or pessimism
141 separately(Weinstein, 1982; Hyer et al., 1984).

142 The widely adopted **life of orientation test**
143 which is based on these theories, measures indi-
144 viduals’ life attitude leaning by analyzing their
145 responses to 6 physical statements. Participants
146 indicate their level of agreement or disagreement
147 with each statement, and their responses are used to
148 calculate their life attitude scores through weighted
149 summation. Formally, the life of orientation test
150 maps a set of answers indicating agreement level
151 {STRONG DISAGREE, NEITHER DISAGREE
152 NOR AGREE, DISAGREE, AGREE, STRONG
153 AGREE} to life attitude scores, where the life atti-
154 tude score range from [0, 24]. A score of 0 to 13
155 indicates a pessimistic attitude toward life, and a
156 score of 14 to 24 indicates an optimistic attitude to-
157 ward life. We employ this test as a tool to measure
158 the life attitude leanings of pretrained language
159 models.

160 We probe a diverse set of LMs to measure their
161 alignment with specific physical statements, in-
162 cluding encoder and language generation mod-
163 els (decoder and autoregressive). For encoder-
164 only LMs, we use mask filling with prompts de-
165 rived from the physical statements. We construct
166 the following prompt: “Please respond to the fol-
167 lowing statement: [STATEMENT] I <MASK>
168 with this statement.” Then, pretrained LMs fill
169 the mask and return 10 highest probability tokens.
170 By comparing the aggregated probability of pre-
171 defined positive (*agree, support, endorse, etc.*)
172 and negative lexicons (*disagree, refuse, oppose,*
173 *etc.*) assigned by LMs, we map their answers to
174 {STRONG DISAGREE, DISAGREE, NEITHER
175 DISAGREE NOR AGREE, AGREE, STRONG
176 AGREE}. Specifically, if the aggregated proba-
177 bility of positive lexicon scores is larger than the
178 negative aggregate by 0.1, we think the response
179 as AGREE, and define DISAGREE analogously.
180 If the aggregated probability of positive lexicon
181 scores is larger than the negative aggregate by
182 0.3, we deem the response as STRONG AGREE,

and define STRONG DISAGREE analogously. If the difference between the aggregation probability of positive words and that of negative words is less than 0.1, we suggest the response as NEITHER DISAGREE NOR AGREE. The values at the boundary mentioned above are derived from our experience through a large number of experiments.

We probe language generation models by conducting text generation based on the following prompt: “Please respond to the following statement: [STATEMENT] \n Your response:”. We then use an off-the-shelf stance detector(Lewis et al., 2019) to determine whether the generated response agrees or disagrees with the given statement. We use 10 random seeds for prompted generation, filter low-confidence responses using the stance detector, and average the stance detection scores for a more reliable evaluation.

Using this framework, we aim to systematically evaluate the effect of polarization in pretraining data on the life attitude bias of LMs.

2.2 Measuring the Effect of LM’s life attitude leaning on Mental Health Detection Task Performance

Armed with the LM political life attitude leaning evaluation framework, we investigate the impact of these leaning on mental health disorders detection tasks with social implications such as suicide detection. We fine-tune different life attitude versions of the same LM architecture on these tasks and datasets and analyze the results. This is a controlled experiment setting, *i.e.* only the life attitude pretraining corpora is different, while the starting LM checkpoint, task-specific fine-tuning data, and all hyperparameters are the same. First, we examine performance differences across LMs with different life attitude leanings to determine if the inherent life attitude leaning in LMs could lead to unfairness in mental health detection tasks.

3 Experiment Settings

3.1 Model

We evaluate life of orientation leaning of 10 language models: BERT(Devlin et al., 2018), RoBERTa(Liu et al., 2019), distilBERT(Sanh et al., 2019), distilRoBERTa(Sanh et al., 2020), ALBERT(Lan et al., 2019), GPT2(Radford et al., 2019), LLaMA(Touvron et al., 2023), Alpaca(Taori et al., 2023), Qwen(Bai et al., 2023), baichuan(Yang et al., 2023)and their variants, rep-

resenting a diverse range of model sizes and architectures. The specific versions and checkpoint names of each model are provided in Appendix. For the stance detection model used for evaluating decoder-based language model responses, we use a BART-based model(Lewis et al., 2019) trained on MultiNLI(Williams et al., 2017). To ensure the reliability of the off-the-shelf stance detector, we conduct a human evaluation on 200 randomly sampled responses and compare the results to those generated by the detector. The stance detector has an accuracy of 0.95 for LM responses with clear stances and high inter-annotator agreement among 3 annotators (0.82 Fleiss’ Kappa).

3.2 Life Attitude Corpus for Pretraining

We collected life attitude corpora for LM pretraining that focus on social media domain and life of orientation leaning (optimistic,pessimistic). For social media, we use the optimism-leaning and pessimism-leaning subreddits lists by (Botzer et al., 2022) and the PushShift API (Baumgartner et al., 2020). Additionally, to address ethical concerns of creating hateful LMs, we used a hate speech classifier based on RoBERTa(Liu et al., 2019)and fine-tuned on the TweetEval benchmark(Barbieri et al., 2020) to remove potentially hateful content from the pretraining data. As a result, we obtained two pretraining corpora of comparable sizes: OPTIMISTIC,PESSIMISTIC. These life attitude pretraining corpora are approximately the same size. We further pretrain RoBERTa and GPT-2 on these corpora to evaluate their changes in ideological coordinates and to examine the relationship between the life of orientation leaning in the pretraining data and the model’s life attitude leaning.

3.3 Mental Health Detection Task Datasets

We investigate the connection between models’ life attitude leaning and their risk detection task behavior on four tasks: suicide detection, anxiety detection, depression detection and stress detection. For suicide detection, we adopt dataset presented in (Ji et al., 2022). In terms of depression detection, we use dataset presented in (Pirina and Çöltekin, 2018). The dataset used in evaluate the performance of detect stress is (Turcan and McKeown, 2019). In the aspect of anxiety detection, we adopt dataset presented in (Owen et al., 2020). We evaluate RoBERTa(Liu et al., 2019) and two variations of RoBERTa further pretrained on REDDIT-OPTMISTIC, REDDIT-PESSMISTIC

corpora. While other tasks and datasets are also possible choices, we leave them for future work. We calculate the performance of different LM checkpoints. Statistics of the adopted mental health task datasets are presented in Table 1.

4 Results and Analysis

In this section, we first evaluate the inherent life attitude leanings of language models and their connection to attitude polarization in pretraining corpora. We then evaluate pretrained language models with different life attitude leanings on mental health disorder illness detection, aiming to understand the link between life of orientation leaning in pretraining corpora and fairness issues in risk-related detection task solutions.

4.1 Life of Orientation leaning of LMs

Table 2 illustrates the life of orientation leaning results for a variety of vanilla pretrained LM checkpoints. Specifically, each original LM is mapped to a LOT score with our proposed framework in Section 2.1. From the results, we find that:

- The language models did show different life attitude tendencies, accounting for all poles (optimism and pessimism) on the life attitude orientation test.
- Generally, BERT variants of LMs are more socially conservative (pessimism) compared to GPT model variants. This collective difference may be attributed to the composition of pretraining corpora: while the BookCorpus (Zhu et al., 2015) played a significant role in early LM pretraining, Web texts such as CommonCrawl and WebText(Radford et al., 2019) have become dominant pretraining corpora in more recent models. Since modern Web texts tend to be more libertarian (optimistic) than older book texts, it is possible that LMs absorbed this optimistic shift in pretraining data.
- We additionally observe that different sizes of the same model family (e.g. Roberta) could have non-negligible differences in political leanings. We hypothesize that the change is due to a better generalization in large LMs, including overfitting leaning in more subtle contexts, resulting in a shift of life attitude leaning. We leave further investigation to future work.

4.2 The Effect of Pretraining with Life Attitude Corpora

Table 3 shows the re-evaluated life attitude leaning of RoBERTa and GPT-2 after being further pre-trained with 2 life of orientation leaning pretraining corpora:

- LMs do acquire life attitude leaning from pre-training corpora. optimistic corpora generally resulted in a optimistic shift on the life of orientation test, while pessimistic corpora led to a pessimistic shift from the checkpoint. This is particularly noticeable for RoBERTa further pretrained on optimism corpora, which resulted in a substantial right shift in terms of life of orientation test values (9 to 14). However, most of the ideological shifts are relatively small, suggesting that it is hard to alter the inherent bias present in initial pretrained LMs. We hypothesize that this may be due to differences in the size and training time of the pretraining corpus.
- For RoBERTa, the life attitude corpus led to an average change of 3,5 in life of orientation test score, while the life attitude corpus resulted in a change of 1.5 For GPT-2. This shows that models based on encoder architecture are more likely to be affected by life attitude tendency predicted by pre-training than models based on decoder architecture.

4.3 Life of Orientation Leaning and Mental Health Detection Task

We compare the performance of three models: base RoBERTa and two RoBERTa models further pre-trained with OPTIMISTIC PESSIMISTIC corpora, respectively. Table 4 presents the overall performance on mental health disorder detection, which demonstrates that optimistic-leaning LMs generally slightly outperform in the anxiety and stress detection tasks. While, pessimistic-leaning LMs generally slightly outperform in the suicide and depression detection tasks. The results demonstrate that the life attitude leaning of the pretraining corpus could have a tangible impact on overall task performance.

5 Reducing the Effect of Life Attitude leaning

Our findings demonstrate that life attitude leaning can lead to significant issues of judgement. Models

Table 1: A summary of datasets. Note we hold out a portion of original training set as the validation set if the original dataset does not contain a validation set.

Category	Platform	Dataset	Train	Validation	Test
Suicide	Twitter	T-SID(Ji et al., 2022)	3072	768	960
Anxiety	Reddit	DATD(Owen et al., 2020)	22381	2798	4196
Depression	Reddit	Depression_Reddit(Pirina and Çöltekin, 2018)	1004	431	406
Stress	Reddit	Dreaddit(Turcan and McKeown, 2019)	2270	568	715

Table 2: Measuring the life attitude leaning of various pretrained LMs. BERT and its variants are more socially pessimistic compared to the Decoder architecture models.

Model	Architecture	LOT	polarity
Roberta-base	Encoder	9	pessimism
Roberta-large	Encoder	14	optimism
Bert-base	Encoder	12	pessimism
Bert-large	Encoder	12	pessimism
Albert-base	Encoder	11	pessimism
Albert-large	Encoder	13	pessimism
distilbert	Encoder	14	optimism
distilroberta	Encoder	15	optimism
gpt2	Decoder	14	optimism
qwen2-0.5B	Decoder	14	optimism
qwen2-1.5B	Decoder	15	optimism
llama2-7B	Decoder	16	optimism
baichuan2-7B	Decoder	15	optimism
alpaca	Decoder	15	optimism

with different life attitude leaning have different predictions regarding what is considered mental health issue or not. For example, if a content moderation model for detecting suicide is more sensitive to suicide-related issue content, it can result in being better performance in suicide detection task. We discuss two strategies to mitigate the impact of life of orientation leaning in LMs.

5.1 Attitude Ensemble

The experiments in Section 4.2 show that LMs with different life of orientation leaning behave differently and have different strengths and weaknesses when applied to mental health disorder detection tasks. Motivated by existing literature on analyzing different perspectives in downstream tasks(Gordon et al., 2022), we propose using a combination, or ensemble, of pretrained LMs with different life of orientation leanings to take advantage of their collective knowledge for mental health disorder detection tasks. By incorporating multiple LMs representing different perspectives, we can introduce

a range of viewpoints into the decision-making process, instead of relying solely on a single perspective represented by a single language model. We evaluate a attitude ensemble approach and report the results in Table 6, which demonstrate that attitude ensemble actively engages diverse life attitude perspectives, leading to improved model performance. However, it is important to note that this approach may incur additional computational cost and may require human evaluation to resolve differences.

5.2 Strategic Pretraining

Another insight from our research is that pretrained models with specific life attitude inclinations exhibit heightened sensitivity when detecting mental health issues that resonate with their initial training perspectives. For instance, a model predisposed towards optimism shows a marked improvement in identifying subtle cues of mental well-being, whereas a model inclined towards pessimism is more adept at recognizing signs of distress and negative affect. This sensitivity suggests that the pre-training data’s inherent biases can be leveraged to fine-tune models for specialized downstream tasks.

This presents an opportunity to develop models that are particularly attuned to specific mental health detection scenarios. For example, in a downstream task dedicated to identifying early indicators of depression, it could be advantageous to further pretrain our models on corpora that reflect a spectrum of emotional expression, particularly those that encompass a nuanced understanding of depressive symptoms. By strategically pretraining on datasets rich in such content, our models can become more proficient in detecting the subtle linguistic markers that precede the onset of depressive episodes.

Strategic pretraining might lead to significant enhancements in the performance of mental health detection models in specific scenarios. However,

Table 3: Pretraining LMs with the two life attitude corpora and re-evaluate their score on the life of orientation test.

model	corpora	previous score	previous polarity	after pretraining score	after pretraining polarity
Roberta	optimism	9	pessimism	14	optimism
Roberta	pessimism	9	pessimism	7	pessimism
GPT2	optimism	14	optimism	16	optimism
GPT2	pessimism	14	optimism	13	pessimism

Table 4: Model performance on Mental Health Disorder Detection. Overall best performance is in **bold**.

Model	Suicide Detection		Anxiety Detection		Depression Detection		Stress Detection	
	ACC	F1	ACC	F1	ACC	F1	ACC	F1
Roberta	0.8895	0.8845	0.9403	0.9378	0.9688	0.9647	0.8056	0.8056
Roberta-Pessimism	0.9062	0.9020	0.9353	0.9338	0.9714	0.9710	0.8247	0.8240
Roberta-Optimism	0.8976	0.8866	0.9432	0.9431	0.9710	0.9704	0.8290	0.8176

440 curating the ideal pretraining corpora that are
 441 scenario-specific and representative of the intended
 442 detection focus can be challenging. It requires a
 443 deep understanding of the psychological nuances
 444 associated with mental health and the ability to
 445 source or create datasets that are both diverse and
 446 sensitive to these intricacies.

447 Our work opens up a new avenue for identifying
 448 the inherent life of orientation leaning of LMs and
 449 further study is suggested to better understand how
 450 to reduce and leverage such leaning for downstream
 451 tasks.

452 6 Related Work

453 6.1 The influence of LMs leaning on 454 downstream task

455 Due to the inherent stereotypes and biases in the
 456 pre-training data, the trained model may be implic-
 457 itly biased(Abid et al., 2021). Li et al. proposed
 458 HiErarchical Regional Bias evaluation method
 459 (HERB) to quantify regional biases in different
 460 groups, and proved that there are stereotypes of
 461 different regions in LMs(Li et al., 2022). Liu et
 462 al. propose two metrics to quantify political bias
 463 in GPT2 using a political ideology classifier that
 464 assesses the probability difference between gener-
 465 ated text with and without attributes (gender, lo-
 466 cation, and topic)(Liu et al., 2021). Dixon et al.
 467 introduce and illustrate a new method to measure
 468 and mitigate unexpected bias in models. Their
 469 experiments show that imbalanced training data
 470 can lead to unexpected bias in models, resulting
 471 in unfair application to classification tasks(Dixon
 472 et al., 2018). Buolamwini et al. found that the
 473 vast majority of facial analysis data sets were com-

474 posed of light-skinned subjects (IJB-A was 79.6%
 475 and Adience was 86.2%). By introducing a facial
 476 analysis data set with more uniform data distribu-
 477 tion, the influence of model bias on facial analysis
 478 tasks was verified. They found that dark-skinned
 479 women were the most likely group to be misclas-
 480 sified (34.7%)(Buolamwini and Gebru, 2018). By
 481 perturbing source sentences in machine translation
 482 tasks, Hila et al. found that gender bias exists in the
 483 generated translations(Gonen and Webster, 2020).
 484 Colossal Clean Crawled Corpus (C4), a document
 485 by Dodge et al., evaluated the effects of filters ap-
 486 plied when the dataset was created, and noted that
 487 the blocking list filters disproportionately removed
 488 texts involving minority individuals, demonstrating
 489 the model’s racial bias(Dodge et al., 2021).

490 6.2 Mental Health Disorder Detection

491 Mental health is a key issue in modern society,
 492 and without proper treatment, mental disorders can
 493 sometimes turn into suicidal ideation. To address
 494 this critical issue, there has been a large amount of
 495 mental health research aimed at the efficient and au-
 496 tomated detection of mental health disorders. Men-
 497 tal health disease detection usually uses informa-
 498 tion fusion strategies to make the model know more
 499 information and improve the accuracy of reasoning.
 500 There are three common fusion strategies: feature
 501 fusion, model fusion and task fusion. In terms of
 502 feature fusion, Song et al. combined four groups of
 503 mental disorder indicators using the Feature Atten-
 504 tion Network (FAN) : 1) Word-level features asso-
 505 ciated with depressive symptoms were taken from
 506 DSM-5; 2) Word-level sentiment scores from the
 507 SentiWordNet dictionary(Baccianella et al., 2010);
 508 3) Features related to reflective thinking, expressed

Table 5: Performance of best and average single models and attitude ensemble on mental health disorder detection. Attitude ensemble shows great potential to improve task performance by engaging multiple perspectives.

Model	Suicide Detection		Anxiety Detection		Depression Detection		Stress Detection	
	ACC	F1	ACC	F1	ACC	F1	ACC	F1
AVG UNI-MODEL	0.8977	0.8910	0.9396	0.9382	0.9704	0.9687	0.8197	0.8157
BEST UNI_MODEL	0.9062	0.9020	0.9432	0.9431	0.9714	0.9710	0.8290	0.8240
ATTITUDE ENSEMBLE	0.9167	0.9086	0.9478	0.9463	0.9802	0.9798	0.8456	0.8451

as the number of repetitions of themes in social media posts (Nolen-Hoeksema et al., 2008) : 4) writing style features, measured by the order of parts of speech in social media (Song et al., 2018). Uban et al. use a hierarchical attention network with LSTM Posting dimension and user dimension encoder, combined with the multidimensional representation of text (Uban et al., 2021). In terms of model fusion, Sawhney et al proposed a time-aware transformer model for screening suicide risk on social media (Sawhney et al., 2020). Their model, called STATENet, uses a dual Transformer-based architecture to learn language and emotional cues in tweets. STATENet captures the linguistic cues of tweets to be evaluated, as well as an aggregated representation of the sentiment spectrum obtained from a pre-trained BERT model fine-tuned on the Emonet (Abdul-Mageed and Ungar, 2017) dataset. On task fusion, Turcan et al. explore the use of multi-task learning and the fine-tuning of language models infused with emotion to detect psychological stress. In this work, the authors introduce an innovative task fusion approach that utilizes a multi-task learning setting to perform both stress detection and emotion detection on the same input data (Turcan and McKeown, 2019).

7 Conclusion

We conduct a systematic analysis of the life of orientation leaning of language models. We probe LMs using prompts grounded in physical science and measure models' ideological positions on life of orientation test values. We also examine the influence of life of orientation leaning in pretraining data on the attitude leanings of LMs and investigate the model performance with varying life of orientation leanings on mental health disorder detection tasks, finding that LMs may have different standards for different mental health illness based on their life of orientation leaning.

Our work verifies that the language model can learn the life attitude tendency predicted by pre-

training, and after pre-training with optimistic (pessimistic) corpus, the model's score on the life attitude orientation test is more optimistic (pessimistic). Through the performance of downstream tasks, we found that models pre-trained with pessimistic corpus would perform better in mental health disease detection tasks, and models with different life attitude tendencies had their own advantages and disadvantages in different tasks. We also proposed methods to solve the influence of model life attitude tendencies on mental health disease detection tasks.

8 Limitations

8.1 Life of Orientation Test

in this study, we leveraged the life of orientation test as a test bed to probe the underlying life attitude leaning of pretrained language models. While the life of orientation test is a widely adopted and straightforward toolkit, it is far from perfect and has several limitations, 1) LOT focuses mainly on the two dimensions of optimism and pessimism, and may not fully capture the complex attitudes and emotional states of individuals. In real life, an individual's attitude may be influenced by a variety of factors, including culture, socioeconomic status, and personal experience. 2) LOT was originally designed for a specific cultural context and may not be fully applicable to other cultures. Different cultures may interpret and express optimism and pessimism differently, which may affect the accuracy and reliability of test results. An individual's attitude towards life may vary with time and situation. As a static test, LOT may not reflect this dynamic change.

8.2 Fine-Grained attitude Leaning Analysis

In this work, we force each pretrained LM into its position on a optimistic-pessimistic two-dimensional space based on their responses to life of orientation test. However, life attitude leaning could be more fine-grained than a numerical test

values: being optimistic on one issue does not necessarily exclude the possibility of being pessimistic on another, and vice versa. We leave it to future work on how to achieve a more fine-grained understanding of LM political leaning in a topic- and issue-specific manner.

9 Misuse Potential

In this paper, We show that the model pre-trained from the corpus with life attitude orientation is still within the reasonable range of life attitude orientation tests. However, this preliminary finding does not exclude the possibility of future malicious attempts at create a model in which life attitude orientation is uncontrollable, and some might even succeed. We will establish access permission for the collected partisan pre- training corpora to ensure its research-only usage.

References

Muhammad Abdul-Mageed and Lyle Ungar. 2017. Emonet: Fine-grained emotion detection with gated recurrent neural networks. In *Proceedings of the 55th annual meeting of the association for computational linguistics (volume 1: Long papers)*, pages 718–728.

Abubakar Abid, Maheen Farooqi, and James Zou. 2021. Persistent anti-muslim bias in large language models. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, pages 298–306.

Brooke Auxier and Monica Anderson. 2021. Social media use in 2021. *Pew Research Center*, 1:1–4.

Stefano Baccianella, Andrea Esuli, Fabrizio Sebastiani, et al. 2010. Sentiwordnet 3.0: An enhanced lexical resource for sentiment analysis and opinion mining. In *Lrec*, volume 10, pages 2200–2204.

Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609*.

Francesco Barbieri, Jose Camacho-Collados, Leonardo Neves, and Luis Espinosa-Anke. 2020. Tweeteval: Unified benchmark and comparative evaluation for tweet classification. *arXiv preprint arXiv:2010.12421*.

Jason Baumgartner, Savvas Zannettou, Brian Keegan, Megan Squire, and Jeremy Blackburn. 2020. The pushshift reddit dataset. In *Proceedings of the international AAAI conference on web and social media*, volume 14, pages 830–839.

Steven Bird, Ewan Klein, and Edward Loper. 2009. *Natural language processing with Python: analyzing text with the natural language toolkit*. " O'Reilly Media, Inc."

Nicholas Botzer, Shawn Gu, and Tim Weninger. 2022. Analysis of moral judgment on reddit. *IEEE Transactions on Computational Social Systems*.

Joy Buolamwini and Timnit Gebru. 2018. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Conference on fairness, accountability and transparency*, pages 77–91. PMLR.

Charles S Carver, Linda M Peterson, Donna J Follansbee, and Michael F Scheier. 1983. Effects of self-directed attention on performance and persistence among persons high and low in test anxiety. *Cognitive therapy and research*, 7(4):333–353.

Charles S Carver and Michael F Scheier. 1982. Control theory: A useful conceptual framework for personality–social, clinical, and health psychology. *Psychological bulletin*, 92(1):111.

Ciro Conversano, Alessandro Rotondo, Elena Lensi, Olivia Della Vista, Francesca Arpone, and Mario Antonio Reda. 2010. Optimism and its impact on mental and physical well-being. *Clinical practice and epidemiology in mental health: CP & EMH*, 6:25.

Bernhard Debatin, Jennette P Lovejoy, Ann-Kathrin Horn, and Brittany N Hughes. 2009. Facebook and online privacy: Attitudes, behaviors, and unintended consequences. *Journal of computer-mediated communication*, 15(1):83–108.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.

Lucas Dixon, John Li, Jeffrey Sorensen, Nithum Thain, and Lucy Vasserman. 2018. Measuring and mitigating unintended bias in text classification. In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*, pages 67–73.

Jesse Dodge, Maarten Sap, Ana Marasović, William Agnew, Gabriel Ilharco, Dirk Groeneveld, Margaret Mitchell, and Matt Gardner. 2021. Documenting large webtext corpora: A case study on the colossal clean crawled corpus. *arXiv preprint arXiv:2104.08758*.

Arianna Ferrari. 2008. Is it all about human nature? ethical challenges of converging technologies beyond a polarized debate. *Innovation: the European journal of social science research*, 21(1):1–24.

Hila Gonen and Kellie Webster. 2020. Automatically identifying gender issues in machine translation using perturbations. *arXiv preprint arXiv:2004.14065*.

Mitchell L Gordon, Michelle S Lam, Joon Sung Park, Kayur Patel, Jeff Hancock, Tatsunori Hashimoto, and Michael S Bernstein. 2022. Jury learning: Integrating dissenting voices into machine learning models. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*, pages 1–19.

696	Ralph Gross and Alessandro Acquisti. 2005. Information revelation and privacy in online social networks. In <i>Proceedings of the 2005 ACM workshop on Privacy in the electronic society</i> , pages 71–80.	David Owen, Jose Camacho Collados, and Luis Espinosa-Anke. 2020. Towards preemptive detection of depression and anxiety in twitter. <i>arXiv preprint arXiv:2011.05249</i> .	749
697			750
698			751
699			752
700	Sheila Hardy. 2018. Common mental health disorders in general practice. <i>Practice Nursing</i> , 29(2):63–69.	Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. Pytorch: An imperative style, high-performance deep learning library. <i>Advances in neural information processing systems</i> , 32.	753
701			754
702	Charles R Harris, K Jarrod Millman, Stéfan J Van Der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J Smith, et al. 2020. Array programming with numpy. <i>Nature</i> , 585(7825):357–362.		755
703			756
704			757
705			758
706			
707	Lee Hyer, John Barry, Arthur Tamkin, and Douglas McConatha. 1984. Coping in later-life: An optimistic assessment. <i>Journal of Applied Gerontology</i> , 3(1):82–96.	Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. 2011. Scikit-learn: Machine learning in python. <i>the Journal of machine Learning research</i> , 12:2825–2830.	759
708			760
709			761
710			762
711	Shaoxiong Ji, Xue Li, Zi Huang, and Erik Cambria. 2022. Suicidal ideation and mental disorder detection with attentive relation networks. <i>Neural Computing and Applications</i> , 34(13):10309–10319.	Guido Peeters and Janusz Czapinski. 1990. Positive-negative asymmetry in evaluations: The distinction between affective and informational negativity effects. <i>European review of social psychology</i> , 1(1):33–60.	765
712			766
713			767
714			768
715	Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2019. Albert: A lite bert for self-supervised learning of language representations. <i>arXiv preprint arXiv:1909.11942</i> .	Inna Pirina and Çağrı Çöltekin. 2018. Identifying depression on reddit: The effect of training data. In <i>Proceedings of the 2018 EMNLP workshop SMM4H: the 3rd social media mining for health applications workshop & shared task</i> , pages 9–12.	770
716			771
717			772
718			773
719			774
720	Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. 2019. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. <i>arXiv preprint arXiv:1910.13461</i> .	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. <i>OpenAI blog</i> , 1(8):9.	775
721			776
722			777
723			778
724			
725			
726	Yizhi Li, Ge Zhang, Bohao Yang, Chenghua Lin, Shi Wang, Anton Ragni, and Jie Fu. 2022. Herb: Measuring hierarchical regional bias in pre-trained language models. <i>arXiv preprint arXiv:2211.02882</i> .	Serhii Ryk, Mykola Ryk, Svitlana Repetiy, Dolores Zavitrenko, Irina Makhnovska, and Valentyna Kovalenko. 2023. Lineamientos generales para la construcción de una política de salud mental en el marco del nuevo humanismo del siglo xxi. <i>Cuestiones Políticas</i> , 41(76):780–791.	779
727			780
728			781
729			782
730	Ruibo Liu, Chenyan Jia, Jason Wei, Guangxuan Xu, Lili Wang, and Soroush Vosoughi. 2021. Mitigating political bias in language models through reinforced calibration. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 35, pages 14857–14866.	Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. <i>arXiv preprint arXiv:1910.01108</i> .	783
731			784
732			785
733			786
734			787
735			788
736	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. <i>arXiv preprint arXiv:1907.11692</i> .	Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2020. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter .	789
737			790
738			791
739			
740			
741	Rui Mao, Qian Liu, Kai He, Wei Li, and Erik Cambria. 2022. The biases of pre-trained language models: An empirical study on prompt-based sentiment analysis and emotion detection. <i>IEEE transactions on affective computing</i> .	Ramit Sawhney, Harshit Joshi, Saumya Gandhi, and Rajiv Shah. 2020. A time-aware transformer based model for suicide ideation detection on social media. In <i>Proceedings of the 2020 conference on empirical methods in natural language processing (EMNLP)</i> , pages 7685–7697.	792
742			793
743			794
744			795
745			796
746	Susan Nolen-Hoeksema, Blair E Wisco, and Sonja Lyubomirsky. 2008. Rethinking rumination. <i>Perspectives on psychological science</i> , 3(5):400–424.	Michael F Scheier and Charles S Carver. 1985. Optimism, coping, and health: assessment and implications of generalized outcome expectancies. <i>Health psychology</i> , 4(3):219.	797
747			798
748			799
			800
			801

885 2019), HuggingFacetransformers(Wolf et al., 2020)
886 ,sklearn(Pedregosa et al., 2011), NumPy(Harris
887 et al., 2020), NLTK (Bird et al., 2009)and the
888 PushShift API¹. We commit to making our code
889 and data publicly available upon acceptance to fa-
890 cilitate reproduction and further research.

¹<https://github.com/pushshift/api>