

# Bias In Machine Learning Applications In Psychiatry

Luca Anna Kosina\*

University of Copenhagen

## 1 Abstract

The integration of machine learning (ML) into psychiatric practice offers benefits in diagnosis and treatment prediction, but it also creates challenges related to bias. Unlike physical healthcare, psychiatric diagnoses rely on subjective assessments and patient-reported symptoms, introducing complexities in applying ML models to diagnosis as well as treatment.

This paper examines the potential biases in two specific ML applications in psychiatry: suicide prediction and autism assessment. Drawing on existing psychological research on bias in these domains, it highlights biases in psychiatric diagnoses, particularly concerning ethnicity and gender. The first case study explores a suicide prediction model using electronic health record data which might adopt and enhance racial biases in risk scores. The second case study focuses on an observation-based classifier for autism, emphasizing the gender bias inherent in diagnostic tools like the Autism Diagnostic Observation Schedule (ADOS). The results of the analysis support how important it is to consider existing biases in psychiatric concepts when developing ML applications.

## 2 Introduction

The utilization of machine learning (ML) models in healthcare and medical settings has given us new opportunities for diagnostic precision and treatment prediction. Especially for physical healthcare, the objectivity of methods, such as imaging and biomarkers, enables the evaluation of model performance against human doctors. However, this is different in psychiatric healthcare, where the diagnosis and treatment is not as straightforward in to comparison physical health contexts.

For example, imaging, biomarkers etc. show very concrete illness markers and how they develop

through certain treatment. In psychiatric practice, however, finding diagnosis and treatment is based on psychological factors and depends on human evaluation. Not only is the diagnosis therefore influenced by the subjectivity of the assessment of the doctor, but also on the content of the patient's report of symptoms. Additionally, cultural norms and historical factors are relevant influences in definition and classification of different disorders. Since diagnoses hinge on subjective client descriptions and feedback, the potential for distortions in diagnostic assessments across different demographic groups is also heightened.

This brings benefits and risks of applying ML methods in the field of psychiatry: on one hand using machine learning in neuro-imaging could enhance the objectivity of diagnosis processes since ML models could find more specific differences in brain functions of mentally ill patients than the human eye. On the other hand, especially in methods based on questionnaire and behavioral data, bias that already exists in the field of psychiatry might be reflected and amplified when training ML models with them.

The stakes in psychiatric ML applications are particularly high seeing as errors in diagnosis or treatment predictions can have severe and lasting effects on patients. While ML therefore offers an opportunity to improve clinical processes and introduce objectivity, there is a risk of maintaining existing biases into these models. Therefore, many researchers have already noted that reflection on bias in ML models in psychiatry is relevant as well as necessary and introduced the theoretical basis to controlling biases. (see e.g. (4))

Previous studies, such as that conducted by Obermeyer et. al. (1), have also already shed light on the biased outcomes of AI algorithms in physical healthcare, revealing differences in health risk scores that disadvantage Black patients compared to their White counterparts.

\*e-mail: tpn338@alumni.ku.dk

Looking at bias in applications of ML models in psychiatry bias is relevant in the following respects. Firstly, we can find historical and socio-cultural bias in psychiatric data. Psychological research has been crafted and evaluated primarily for privileged white and western groups of people. This has led to a deficiency in socio-cultural sensitivity of diagnosis criteria, inclusiveness, and adaptability in our treatments. (2) New technological solutions are, however, often especially created to treat more people and such who do not usually access mental health services. (3) Therefore, it is especially important to check for bias in AI applications and how models can reinforce and sustain past societal disparities - since otherwise these solutions continue to discriminate or even amplify biased access and treatment to certain people and groups. Building ML models for psychiatry could therefore create a feedback loop in which biased models lead to more biased data.

Additionally, unbalanced samples or choosing features based on only one patient group can introduce representational bias when training the models. They are then discriminating against certain groups which, for example, need different features representations to assess their diagnosis, risk scores or medication success.

There are different ways in which machine learning has already been applied in psychiatry which are sensitive to bias. These different areas include diagnosis and classification, prediction of the course of illness, treatment through interventions and monitoring. (3)

In this paper I want to look at two specific examples of these types of application, specifically in the field of diagnosis and prediction.

I will focus on concepts which provide sufficient knowledge from psychological research how biased our classification of diagnosis and assessment potentially are. This will be used to reflect how this knowledge is relevant for the ML models, how they should be developed or assessed and adapted.

For example, with more research coming up on autism in females and shifting towards a less gender biased diagnosis, utilizing machine learning models in this field could be a chance to consider and include more variability. This is why one of the studies to be analyzed is on autism assessment with ML.

The other study will focus on suicide prediction with a ML model and I will look at how racial bias in psychiatric diagnosis could influence the bias in

this prediction of suicidal behavior.

By investigating and critiquing these cases, I seek to contribute to a nuanced understanding of the challenges and opportunities presented by ML applications in psychiatric settings and to inform strategies for mitigating bias in future developments.

### 3 Related Work

Different works have analyzed the relevant principles for avoiding bias in machine learning methods.

Timmons et. al. (3) have analyzed areas in which AI can find implementation in the mental health fields as well as therefore can be influenced by the bias of AI. These areas include diagnosis, intervention, engagement, maintenance monitoring and prediction.

I will exemplarily look at works in two of these areas in which machine learning models were utilized in the fields of diagnosis and prediction and analyze different possible bias in these studies.

Tay et. al. (4) provide different types of bias which can occur in this field of research: There can be bias in elements like ground truth, platform-based construct, behavioral expression, and feature computing. Similarly, algorithm-training bias may emerge during the development of algorithms when there is a lack of equivalence in the connection between extracted features and the ground truth. They state that one big limitation in the research is the absence of clear methodological guidelines. They also differentiate measurement bias from socio-cognitive bias, which involves errors in human cognitive judgments or attributions. Measurement bias is characterized by a distinct relationship between the latent score (i.e., psychological construct score) and the predicted observed score.

Since many studies do not analyze which socio-cognitive bias is relevant for ML models or their measurement bias for different subgroups I will focus on how certain studies should analyze subgroups to find potential differences between prediction and what has been found in psychological research as differences and bias between different groups. I will therefore focus on the influence of socio-cognitive bias on applying ML for e.g. diagnosis and induces measurement bias.

There is only limited work in which researchers analyze concrete ML model results for their bias. For example, Mosteiro et.al. (5) have developed a machine learning model predicting future benzo-

182 diazepam administrations, revealing unexpected  
183 gender-related bias and analyzed potential bias-  
184 mitigation strategies. In many existing projects using  
185 ML in psychiatry researchers do not, however,  
186 analyze datasets regarding their potential for bias  
187 and how the model performs for different groups.  
188 This is why in the following I will critique how certain  
189 studies should be analyzed regarding potential  
190 sources of bias.

## 191 **4 Machine Learning for predicting** 192 **suicidal behavior, racial bias and bias** 193 **against minorities**

### 194 **4.1 Method and results**

195 In their study Barak-Corren et. al. (6) tested the  
196 application of machine learning to predict suicidal  
197 behavior for patients based on various medical and  
198 demographic backgrounds. They used EHR data  
199 from medical centers in Boston and defined suicidal  
200 behavior through the ICD-9 diagnostic codes.

201 Naive Bayesian classifier models were utilized  
202 to assess individual suicide risks based on in-  
203 dependent input variables such as demographic  
204 characteristics, diagnostic codes, laboratory re-  
205 sults (normal/low/high), and prescribed medica-  
206 tions (true/false values). They had a control group  
207 of patients without suicidal behavior for which they  
208 analyzed data from the entire observed time peri-  
209 ods - in contrast, for the case group, the group  
210 of patients with suicidal behavior, they analyzed  
211 data from before the first suicidal event. Also,  
212 the separated the female from the male patient  
213 data and build two separate models for men and  
214 women which therefore minimizes gender bias in  
215 this model.

216 For the model training they assigned each in-  
217 dependent input variable (e.g., diagnoses, medica-  
218 tions, etc.) a risk score based on the ratio of its  
219 occurrence among patients with suicide attempts  
220 compared to the control group patients. Each pa-  
221 tient had a risk score at each time point and with  
222 the model partial risk scores were calculated for  
223 each variable. To interpret the patient's score, they  
224 utilized the predefined thresholds selected during  
225 the training phase to achieve 90 percent and 95  
226 percent specificities for each of which the sensitivity  
227 and timeliness of prediction were assessed.

228 Relevant results were that the model detected  
229 44 percent of suicides among men and 46 percent  
230 of suicides among women at a 90 percent speci-  
231 ficity level and was therefore more successful than

232 models only using bivariate combinations of the ac-  
233 cepted risk factors of depression, substance abuse  
234 and other mental disorders. Furthermore, using the  
235 model for even more specific subcategories such  
236 as specific age groups (e.g. women aged 45–65)  
237 performed better with 54 percent sensitivity.

238 Important risk factors in the model are person-  
239 ality disorders as well as bipolar disorder which  
240 were 7–10 times more common among case sub-  
241 jects compared to the control group. Also, opioid  
242 abuse had a 16 times higher likelihood among the  
243 case group.

### 244 **4.2 Bias in psychological research**

245 To assess which influences could bring bias in the  
246 developed model we have to look at how the vari-  
247 ables used in the model are known to be biased  
248 among different groups. Especially diagnosis of  
249 disorders, but also the received type of medication  
250 differ between ethnicities. Additionally, the way  
251 people from different groups access professional  
252 psychiatric help differs.

253 Firstly, regarding diagnosis of personality disor-  
254 ders there is evidence that e.g. people with African  
255 Caribbean ethnicity are less likely to be diagnosed  
256 properly with a personality disorder with differing  
257 rates in different personality disorders. Especially  
258 high is the bias for example for anti-social person-  
259 ality disorders. (7)

260 In patients with bipolar disorder, racial bias leads  
261 to misdiagnosis of schizophrenia instead of ma-  
262 nia with psychotic features in African American  
263 patients. (8) Additionally, treatment of African  
264 Americans is less likely to include anti-psychotic  
265 medication, but more likely to be treated with mood  
266 stabilizers such as lithium. (8)

267 Also, when having co-morbid disorders to a sub-  
268 stance abuse disorder, White patients are more  
269 likely to be diagnosed and treated for these co-  
270 morbid disorders. Garb et. al. (7) write, for ex-  
271 ample: "Compared to Whites, Blacks with these  
272 co-morbid disorders were significantly less likely  
273 to receive services for mood or anxiety, equally  
274 likely to receive services for alcohol use disorders,  
275 and more likely to receive services for drug use  
276 disorders."

277 Anderson et. al. (9) investigates differences be-  
278 tween races in responses of adolescences in suicide-  
279 related psychiatric assessments and noted that pre-  
280 vious research has not explored the possibility of  
281 non-response, even though evidence suggests that  
282 individuals with higher levels of suicidality may be

less willing to disclose such information. (9) While there were less significant differences for questions about suicide plans, there was a statistically significant relationship between race and the response rate to questions about suicide attempts. African American students stood out in terms of non-response. Nearly 20 percent of African American adolescents omitted the question about suicide attempts, while only 7.6 percent of Caucasian youth omitted this item. (9)

Additionally, Rockett et. al. (13) investigated potential for misclassification of suicide for different ethnicities and found that Blacks and Hispanics had higher potential suicide misclassification relative to Whites.

Bommersbach et. al. (?) similarly found that patients with Black and Hispanic ethnicity were less likely than White individuals to report suicidal ideation, however more likely to report a suicide attempt. They also found that among people with suicidal ideation, for Black adults the likelihood of also having a suicidal attempt is higher than for White people. This indicates that that stigma is higher or help-seeking behavior lower in the group of Black people leading to them reaching a more severe health state when receiving help.

Lastly, there we can find an interaction between different demographic factors such as age and race regarding suicide: Garlow et. al. (12) showed the median age of patients who attempted suicide was more than a decade lower for black patients with 34 years compared to white patients with 44 years.

### 4.3 Analysis of Bias in ML model

While suicide rates among White people in the US remain higher than the ones among Black as well as Asian Pacific and Hispanic people, increases of suicide rate from 2018 to 2019 were higher for the latter while number of suicides reduced for White people. (11) This highlights the need to have especially sensitive early prediction of suicide among populations of non-White people.

However, we can find different way in which the model by Barak-Corren et. al. might neglect sensitivity towards differences between different race groups.

Firstly, when analyzing the dataset used to train the model in regards to racial equality we notice unbalanced sampling since the dataset has unequal ratios of various ethnicities with the majority of patients being white. (6)

If however, as we have seen, risk factors and

medical records might differ between different ethnicities, this can lead to inaccurate or inefficient associations between other risk factors and ethnicity and mix different risk profiles.

In general, cause of a lack of number of patients with other ethnicities than White to compare suicide risk factors between groups might be the mentioned misclassification of suicide in some ethnicities. Also, we know from Anderson et. al. (9) that there might be differences in disclosure of suicide attempts of Black patients - therefore the number of unrecorded cases of unsuccessfully suicide attempts might also be lower for this population group. In predicting suicide the risk of Black patients might therefore be underestimated.

However, in this study we can find a different sampling problem. The features the model uses in suicide prediction include demographic data as well as medical data about diagnoses, laboratory results and medications. The demographic data also includes ethnicity of the patient and one might suspect it to therefore influence the suicide risk assessment through suicide rates being higher for White people. However, the model results in higher risk scores for African Americans with 1.63 for men and 1.29 for women as well as Hispanic men with 1.43 and Hispanic women with 1.69 compared to Caucasian patients with 0.99 for both men and women. This leads us to discover another imbalance in the dataset used to train the model: While for the case and control group the ratio of White patients to other ethnicities is similar, the percentage of Hispanic and African American patients compared to other ethnicities in the case group is higher than in the control group.

However, this does not represent the ratios of the general population: While Hispanic and African American people make up 7.4 and 6.4 percent of the control group, the make up 12.3 and 8.3 percent of the case group. This can explain why suicide risk is higher for these ethnicities in the model than it would be based on population wide suicide risks. For the US population overall, suicide rate for Hispanic and Black individuals are 7.3 and 7.4 per 100.000 compared to 17.6 in 100.00 for White individuals. (11)

Reason for the imbalance in the dataset from the hospitals might be that - as we could conclude from Bommersbach et. al. (10) - Black and Hispanic people seek help when their health state is more severe, when they had a suicidal attempt and not "only" suicidal ideation. There might therefore be

386 a higher percentage of case subjects than control  
387 subjects from these ethnicities.

388 While it might therefore seem like the model  
389 makes up for bias in suicide assessment and clas-  
390 sification since risk scores are lower for White pa-  
391 tients, the cause of this might be the sample choice  
392 and imbalance. We can ask how risk scores would  
393 be different when trained on different datasets with  
394 different ratios and how the model can be applied  
395 to groups with different ratios of ethnicities.

396 Not only could the risk scores for certain ethnic-  
397 ities be influenced by the ethnicity's percentage in  
398 the sample groups but also through interaction with  
399 other features. On one hand, seemingly objective  
400 medical data can, as analyzed, include bias which  
401 is perpetuated in the model. On the other, the risk  
402 scores of certain factors might be different for dif-  
403 ferent ethnicities and we should pay attention how  
404 these differences should be accounted for in the  
405 model.

406 For example, the model calculated especially  
407 high risk scores for personality disorders, bipolar  
408 disorder and opioid abuse. (6) We also know that  
409 diagnosis of personality disorders, bipolar and co-  
410 morbid disorders to substance abuse are affected by  
411 racial bias in psychiatry which could lead to false  
412 increases as well as decreases of the risk score in  
413 the ML model. Therefore, if, for example, the  
414 group of African American people is not correctly  
415 diagnosed with bipolar, but instead with schizophre-  
416 nia, their risk scores might be different for different  
417 disorders - similarly for African Caribbean patients  
418 not being correctly diagnosed with personality dis-  
419 orders. Also, if their substance abuse disorder is  
420 not correctly co-diagnosed with other disorders as  
421 well as not treated properly, their risk scores for  
422 substance abuse might be different. Since medi-  
423 cation is also a feature in the model, bias through  
424 wrong or missing diagnosis and treatment has a  
425 second influence on the risk prediction for different  
426 ethnicities.

427 All of this can lead to higher lower risk scores  
428 for certain ethnicities because of different diagno-  
429 sis patterns and co-morbidities or less diagnoses  
430 overall. It could also lead to higher risk scores  
431 through incorrect diagnosis of disorders or differ-  
432 ent medication.

433 It should therefore be taken into account that  
434 there are interactions of different risk factors, not  
435 only between, for example, suicide risk and race  
436 itself, but also through racial bias in other features.  
437 As an additional example we have also noted age

differences in suicide attempts for different ethnic-  
ities. For that reason having more specific sub-  
groups for age and race might result in more accu-  
rate prediction. The fact that the model already per-  
formed better when choosing specific age groups  
for women and men (6), is a sign that even more  
specification might be a valuable practice.

Just like there were two different models trained  
for female vs. male individuals, it might also make  
sense to train for different ethnicities or - first of  
all - test their the predictions for groups of differ-  
ent ethnicities to analyze limitations of applying  
this model to different demographics. Then differ-  
ent specifications and combinations of age groups,  
ethnicities etc. can be tested to see which lead to  
higher performance.

## 5 Machine Learning for autism and gender bias

### 5.1 Method and results

In their study Duda et. al. (14) develop a  
machine-learning classifier for differentiating cases  
of autism spectrum disorder (ASD) from non-  
spectrum. They use an observation-based classifier  
(OBC) which holds eight behaviors, namely fre-  
quency of vocalization directed to others, eye con-  
tact, social smile, shared enjoyment in interaction,  
showing, initiation of joint attention, functional  
play with objects, imagination/creativity. It calcu-  
lates a score through an Alternating Decision Tree  
algorithm to assess a patient's class (spectrum or  
non-spectrum) and also gives back the confidence  
level of the classification. The scores were com-  
pared to patient's ADOS (Autism Diagnostic Ob-  
servation Schedule) scores as a common autism  
diagnostic instrument to calculate sensitivity and  
specificity and the correlation of classification out-  
comes of the OBC and ADOS (using Spearman  
rank correlation).

Comparing to two version of ADOS, ADOS-G  
and ADOS-2, the OBC hat a sensitivity of >97  
percent for both and a specificity of 77 percent  
for ADOS-G and 84 percent for ADOS-2. (14)  
The OBC results also had significant correlation  
to the ADOS-2 comparison score, so the authors  
concluded that the OBC reflects severity of the  
autism phenotype and that the OBC is a possibility  
for rapid classification to shorten waiting times for  
autism diagnosis. (14)

## 5.2 Bias in psychological research

Autism Spectrum Disorder presents - as the name suggests - as a wide spectrum of symptoms and severity levels. However, a lot of varieties and differences between groups have been a topic of research only in recent years: ASD has been considered a dominantly "male" disorder with 4:1 percent of diagnosed people being male compared to female. (15) In research on autism there has however been a sampling bias through, for example majority of the studied patients being male and most studies ignoring variables such as gender/sex in their analysis. (15)

Now, however, there is now a shift in research towards a female phenotype of ASD and more female children receive an autism diagnosis. Causes for diagnostic bias might be that male children on the autism spectrum show more obvious or easily identifiable repetitive patterns in behavior, interests, activities (RRBIs) than girls since the latter rather hyper-fixate on animals or dolls. (15) Also, females with ASD have better socio-emotional reciprocity (nonverbal communication, appropriate facial communication, offering comfort) and have a tendency to masking their autistic symptoms while still perceiving them more than men. (15)

The discovery of these differences in the female phenotype result in a shift of symptom structure and classification and why it motivates research on how biased the diagnosis through ADOS is. (16) Definitions in the DSM and ICD and clinical tools such as the ADOS have been based on Kanner and Asperger's descriptions and therefore their sensitivity has been fit to the male phenotype. (15)

Adamou et. al. (16) have tested the gender bias of the ADOS classification of patients into autistic vs. non-spectrum. The ADOS contains five modules each with a protocol of activities or social processes in which items are given scores from zero indicating "no abnormality of type specified" to three indicating "moderate to severe abnormality".

In their sample of people diagnosed with ASD based on the DSM-5, they found that the mean ADOS score for men with ASD diagnosis was 11.5, whereas women had a mean score of 6 - which is below the diagnostic threshold for ASD. For the people without ASD the mean for men was 5.7 and for women 4. They conclude that relying excessively on ADOS scores therefore may introduce gender bias against females and that the ADOS sen-

sitivity is lower for autism in females. (16) Instead the thresholds for ADOS score indicating a ASD diagnosis should probably be adjusted for females - or the tested items and their rating adjusted to the female phenotype.

## 5.3 Analysis of Bias in ML model

Firstly, bias could be introduced in the OBC through unbalanced sampling since - similarly to other studies on autism - the female-male ratio of the studied cohort was unbalanced with 76.8 percent male patients, 18.5 percent female and 4.7 percent with unknown gender. It is, however, with the background on how autism research shifts towards differentiating symptoms of female and male children with ASD, relevant to test the OBC model on a more balanced dataset to see if this influences the performance of the model.

Furthermore, using the ADOS as a tool of comparison to assess the model might turn out to not be an optimal or appropriate measure due to the bias we have seen that the ADOS has.

Since the mean ADOS score for autistic as well as non-autistic females seems to have a high likelihood to be below the diagnostic cut-off, if the model returns similar results as ADOS, it is therefore also likely to wrongly categorize females as not on the autism spectrum. The authors don't provide a general analysis how differently the model performs in classifying female compared to male children. One could compare how many false positives are male compared to female and how the correlation to ADOS scores differs between female and male participants. Also, analyzing false negatives for different gender groups would be relevant.

Additionally, for debiasing this application of ML it would be necessary to assess how girls are classified differently in each behavioral category. For example, there would be indication that with better non-verbal communication there scores for smile and showing might be better.

It also indicates that we should ask if it is advisable to create separate models for female and male patients due to the weights of different factors such as socio-emotional reciprocity being different.

## 6 Conclusion

We have seen that applying machine learning models without taking into consideration the existing socio-cognitive bias in psychological diagnosis and treatment processes could amplify these biases in

586 models and that there is still too little consideration  
587 of this risk in the field.

588 In the ML model by Barak-Corren et. al. gender  
589 bias was considered by training separate models for  
590 female and male, however, we argued that there are  
591 racial bias influences. We found overrepresentation  
592 of certain ethnicities in the case group, while also  
593 assessing potential for wrong risk scores through  
594 missing or incorrect diagnosis and medication.

595 In Duda's et. al. model to diagnose autism gen-  
596 der bias must be taken into consideration and I  
597 stressed the necessity of adjusting diagnostic tools  
598 and ML models to account for the evolving under-  
599 standing of ASD in females.

600 In both studies a more nuanced analysis of per-  
601 formance results for different subgroups is conse-  
602 quently advisable to improve model performance.  
603 Also, I suggested ways in which features for sub-  
604 groups could differ and why therefore training dif-  
605 ferent models might be useful.

606 We can conclude that future work on ML ap-  
607 plications in psychiatry should be more aware of  
608 bias analysis, especially taking into account the  
609 existing research in psychology on how bias psy-  
610 chiatric concepts might be. Studies should take into  
611 account how datasets might have to split into sub-  
612 groups either for training or to analyze afterwards  
613 if results of the model differ between these groups.  
614 Here it is important - as previous papers have noted  
615 - to compare the ratio of a model's predictions by  
616 group and by the intersection of groups. For each  
617 group ratios should be the same or maintained from  
618 the ratios in reality and also performance should be  
619 equal to avoid social inequities. (3)

620 In the future, the spread of a more uniform stan-  
621 dard and widely used bias detection or mitigation  
622 strategy would be advisable to avoid progress in  
623 the field benefiting only certain groups and instead  
624 leveraging the beneficial effect of machine learning  
625 for objectivity and inclusiveness in psychiatry.

## 626 References

- 627 [1] Ziad Obermeyer et al. 2019. "Dissecting  
628 racial bias in an algorithm used to manage  
629 the health of populations." *Science* 366, 447-  
630 453. DOI:10.1126/science.aax2342
- 631 [2] Roberts, Steven O., Carmelle Bareket-Shavit, For-  
632 rest A. Dollins, Peter D. Goldie, and Eliza-  
633 beth Mortenson. 2020. "Racial Inequality in Psy-  
634 chological Research: Trends of the Past and  
635 Recommendations for the Future." *Perspectives*

on *Psychological Science* 15 (6): 1295–1309.  
https://doi.org/10.1177/1745691620927709.

- [3] Timmons, Adela C., Jacqueline B. Duong, Na-  
talia Simo Fiallo, Theodore Lee, Huong Phuc  
Quynh Vo, Matthew W. Ahle, Jonathan S. Comer,  
LaPrincess C. Brewer, Stacy L. Frazier, and  
Theodora Chaspari. 2023. "A Call to Action on  
Assessing and Mitigating Bias in Artificial Intel-  
ligence Applications for Mental Health." *Perspec-  
tives on Psychological Science* 18 (5): 1062–96.  
https://doi.org/10.1177/17456916221134490.
- [4] Tay, Louis, Sang Eun Woo, Louis Hickman, Brandon  
M. Booth, and Sidney D'Mello. 2022. "A Conceptual  
Framework for Investigating and Mitigating Machine-  
Learning Measurement Bias (MLMB) in Psychologi-  
cal Assessment." *Advances in Methods and Practices  
in Psychological Science* 5 (1): 25152459211061337.  
https://doi.org/10.1177/25152459211061337.
- [5] Mosteiro, Pablo, Jesse Kuiper, Judith Masthoff,  
Floortje Scheepers, and Marco Spruit. 2022.  
"Bias Discovery in Machine Learning Models  
for Mental Health." *Information* 13 (May): 237.  
https://doi.org/10.3390/info13050237.
- [6] Barak-Corren, Yuval, Victor M. Castro, Solomon  
Javitt, Alison G. Hoffnagle, Yael Dai, Roy H.  
Perlis, Matthew K. Nock, Jordan W. Smoller, and  
Ben Y. Reis. 2017. "Predicting Suicidal Behav-  
ior From Longitudinal Electronic Health Records."  
*American Journal of Psychiatry* 174 (2): 154–62.  
https://doi.org/10.1176/appi.ajp.2016.16010077.
- [7] Garb, Howard N. 2021. "Race Bias and Gender Bias  
in the Diagnosis of Psychological Disorders." *Clini-  
cal Psychology Review* 90 (December): 102087.  
https://doi.org/10.1016/j.cpr.2021.102087.
- [8] Akinhanmi, Margaret O, Joanna M Biernacka,  
Stephen M Strakowski, Susan L McElroy, Joyce  
E Balls Berry, Kathleen R Merikangas, Shervin  
Assari, et al. 2018. "Racial Disparities in Bipo-  
lar Disorder Treatment and Research: A Call  
to Action." *Bipolar Disorders* 20 (6): 506–14.  
https://doi.org/10.1111/bdi.12638.
- [9] Anderson, Laura M., Lynda S. Lowry, and  
Karl L. Wuensch. 2015. "Racial Differences  
in Adolescents' Answering Questions About  
Suicide." *Death Studies* 39 (10): 600–604.  
https://doi.org/10.1080/07481187.2015.1047058.
- [10] Bommersbach, Tanner J., Robert A. Rosenheck,  
and Taeho Greg Rhee. 2023. "Racial and Eth-  
nic Differences in Suicidal Behavior and Mental  
Health Service Use among US Adults, 2009–2020."  
*Psychological Medicine* 53 (12): 5592–5602.  
https://doi.org/10.1017/S003329172200280X.
- [11] Ramchand, Rajeev, Joshua A. Gordon, and  
Jane L. Pearson. 2021. "Trends in Suicide Rates  
by Race and Ethnicity in the United States."  
*JAMA Network Open* 4 (5): e2111563–e2111563.  
https://doi.org/10.1001/jamanetworkopen.2021.11563.

- 693 [12] Garlow, Steven J., David Purselle, and Michael  
694 Heninger. 2005. "Ethnic Differences in Pat-  
695 terns of Suicide Across the Life Cycle." *Ameri-*  
696 *can Journal of Psychiatry* 162 (2): 319–23.  
697 <https://doi.org/10.1176/appi.ajp.162.2.319>.
- 698 [13] Rockett, Ian RH, Shuhui Wang, Steven Stack,  
699 Diego De Leo, James L. Frost, Alan M. Ducatman,  
700 Rheeda L. Walker, and Nestor D. Kapusta. 2010.  
701 "Race/Ethnicity and Potential Suicide Misclassifica-  
702 tion: Window on a Minority Suicide Paradox?" *BMC*  
703 *Psychiatry* 10 (1): 35. [https://doi.org/10.1186/1471-](https://doi.org/10.1186/1471-244X-10-35)  
704 [244X-10-35](https://doi.org/10.1186/1471-244X-10-35).
- 705 [14] Duda, M, J A Kosmicki, and D P Wall.  
706 2014. "Testing the Accuracy of an Observation-  
707 Based Classifier for Rapid Detection of Autism  
708 Risk." *Translational Psychiatry* 4 (8): e424–e424.  
709 <https://doi.org/10.1038/tp.2014.65>.
- 710 [15] Young, H., M.-J. Oreve, and M. Speranza.  
711 2018. "Clinical Characteristics and Problems  
712 Diagnosing Autism Spectrum Disorder in  
713 Girls." *Archives de Pédiatrie* 25 (6): 399–403.  
714 <https://doi.org/10.1016/j.arcped.2018.06.008>.
- 715 [16] Adamou, Marios, Maria Johnson, and Bronwen  
716 Alty. 2018. "Autism Diagnostic Observation Sched-  
717 ular (ADOS) Scores in Males and Females Diagnosed  
718 with Autism: A Naturalistic Study." *Advances in*  
719 *Autism* 4 (2): 49–55. [https://doi.org/10.1108/AIA-](https://doi.org/10.1108/AIA-01-2018-0003)  
720 [01-2018-0003](https://doi.org/10.1108/AIA-01-2018-0003).