1

2

3       **Unraveling the Associations Between Voice Pitch and Major Depressive Disorder: A**

4                                    **Multisite Genetic Study**

5

6       Yazheng Di[1,2.], B.Sc., Elior Rahmani[3.], Ph.D., Joel Mefford[4.], Ph.D., Jinhan Wang[5.], M.Sc., Vijay

7       Ravi[5.], M.Sc., Aditya Gorla[6.], B.Sc., Abeer Alwan[5.], Ph.D., Kenneth S. Kendler[7,8], M.D., Tingshao

8       Zhu[1,2.*], Ph.D., Jonathan Flint[9.]*, M.D.

9

10      1. CAS Key Laboratory of Behavioral Science, Institute of Psychology, Beijing 100101,
11         China.
12      2. Department of Psychology, University of Chinese Academy of Sciences, Beijing 100049,
13         China.
14      3. Department of Computational Medicine, University of California Los Angeles, Los
15         Angeles, CA, USA.
16      4. Department of Neurology, University of California Los Angeles, Los Angeles, CA, USA.
17      5. Department of Electrical and Computer Engineering, University of California Los Angeles,
18         Los Angeles, CA, USA.
19      6. Bioinformatics Interdepartmental Program, University of California Los Angeles, Los
20         Angeles, CA, USA
21      7. Virginia Institute for Psychiatric and Behavioral Genetics, Richmond, VA, USA
22      8. Department of Psychiatry, Virginia Commonwealth University School of Medicine,
23         Richmond, VA, USA
24      9. Department of Psychiatry and Biobehavioral Sciences, Brain Research Institute, University
25         of California Los Angeles, Los Angeles, CA, USA.
26

27      *Corresponding authors: Jonathan Flint & Tingshao Zhu

28      Email: JFlint@mednet.ucla.edu(J.F.), tszhu@psych.ac.cn(T.Z.)

29      *Disclosures:* No conflict of interest.

31

2

32

33 **Abstract**

34 Major depressive disorder (MDD) often goes undiagnosed due to the absence of clear biomarkers.

35 We sought to identify voice biomarkers for MDD and separate biomarkers indicative of MDD

36 predisposition from biomarkers reflecting current depressive symptoms. Using a two-stage meta-

37 analytic design to remove confounds, we tested the association between features representing vocal

38 pitch and MDD in a multisite case-control cohort study of Chinese women with recurrent

39 depression. Sixteen features were replicated in an independent cohort, with absolute association

40 coefficients (beta values) from the combined analysis ranging from 0.24 to 1.07, indicating

41 moderate to large effects. The statistical significance of these associations remained robust, with

42 P-values ranging from $7.2 \times 10^{-6}$ to $6.8 \times 10^{-58}$. Eleven features were significantly associated

43 with current depressive symptoms. Using genotype data, we found that this association was driven

44 in part by a genetic correlation with MDD. Significant voice features, reflecting a slower pitch

45 change and a lower pitch, achieved an AUC-ROC of 0.90 (sensitivity of 0.85 and specificity of

46 0.81) in MDD classification. Our results return vocal features to a more central position in clinical

47 and research work on MDD.

48

49

50

3

51

## Introduction

53  Changes in human pitch and tone of speech have been noted as an important sign of depression

54  for over a century(1,2). Although not contained in symptomatic criteria for major depressive

55  disorder (MDD) in DSM-III(3), DSM-IIIR(4), DSM-IV(5), or DSM-5(6), they are found in 26 out

56  of 28 detailed clinical descriptions of melancholia published from 1880-1900(1) and in 19 out of

57  21 of such descriptions of depression published in the 20th century(2). Given the current

58  challenges in diagnosing MDD(2,7,8), where a large proportion of cases (ranging from 50% to

59  90%) remain untreated(9–11), the transformation of voice phenomena into diagnostic biomarkers

60  could aid in both clinical and research arenas.

61  Clinical observations describe the speech patterns of depressed patients as slow, weak, low-

62  pitched, and monotonous(1,2,12,13). These phenomena are typically quantified by increased pause

63  time, lower volume, lower pitch, and reduced pitch variability(14–16). Many studies have sought

64  to develop features from pitch as a biomarker for depression(17–26), but none have, to date,

65  achieved sufficient accuracy and precision for clinical utility. The large number of both vocal

66  features and confounds(24,25) imposes a multi-testing burden that requires large sample sizes

67  which few studies have obtained(16). Also, a critical distinction between current mood and

68  susceptibility to MDD on effects on voice has never been addressed(27–29). Furthermore, MDD

69  is likely heterogeneous(30,31): studies not accounting for this may be underpowered (16).

70  Our study of the relationship between voice features and MDD was designed to be well-

71  powered, using thousands of subjects, to be more robust to heterogeneity by analyzing cases with

72  recurrent depression of one sex only, and able to separate susceptibility to MDD from the effect

73  of current mood on voice features by using genetic data. We used a large case-control study of

4

74   MDD where we could replicate findings in an independent sample, both from China. Since, of the

75   four groups of speech features (source, spectral, prosodic, and formant features(16)), the

76   association between depression and a key component of prosody, pitch, has repeatedly been

77   observed(15,16), our analysis was restricted to examining pitch-related features. Our results return

78   vocal features to a more central position in clinical and research work on MDD.

79

80   **Results**

81   *Subjects*

82   We used recordings conducted as part of the CONVERGE(32) (China, Oxford, and VCU

83   Experimental Research on Genetic Epidemiology) study (3,968 cases and 4,354 controls).

84   CONVERGE recruited only women with recurrent MDD from hospital settings and compared

85   them with matched controls with no history of MDD, thus reducing heterogeneity in both genetic

86   and vocal signals(33,34).   A summary of demographic data from cases and controls is provided

87   in **Table S1**; the relation of these to depression is reported in earlier publications (31,35–40). All

88   recordings, obtained during diagnostic interviews, were listened to, and segments that contained

89   only the patient's voice at an adequate quality for the analyses (see the Method section in

90   **Supplementary** for details) were identified. In this study, a "segment" refers to the longest

91   continuous portion of an audio recording that contains only the patient's voice, uninterrupted by

92   other speakers. Segments are not split at pauses within the patient's speech but are instead defined

93   by changes in the speaker, ensuring that each segment represents uninterrupted speech from the

94   patient alone. It can be as short as a single word or as long as a complex sentence or multiple

95   sentences. This resulted in 364,929 voice segments with a duration greater than two seconds from

96   7,654 subjects. The selection of subjects for each component of the study is shown in **Figure S1**,

5

97   which provides an overview of the design of the project (**Figure S1a**), and a PRISMA diagram to

98   indicate how many cases and controls were discarded at different stages and pathways of analyses

99   (**Figure S1b)**.

100

101   *Feature identification*

102   The perceptual attribute of pitch corresponds to the physical measurement of fundamental

103   frequency, or F0(15). Our interest in prosodic features of speech, particularly pitch (hereafter

104   referred to as the F0) and change in pitch (ΔF0),  led us to choose the INTERSPEECH 2016

105   Computational Paralinguistics Evaluation (COMPARE16) feature set (41,42).  The primary

106   application for the feature set has been depression detection(24,43–47) and it captures typical

107   temporal information and long-term information either statically, through the use of utterance level

108   statistics/functionals, or dynamically, through frame-based delta (ΔF0) coefficients, reflecting

109   differences between the adjacent vectors' feature coefficients(48,49). Its feature extraction process

110   is well documented(50), with standardized and well-referenced methodologies that facilitate

111   reproducibility and validation by other researchers in many languages (including Chinese)(51–54).

112        The COMPARE16 feature set contains 83 F0/ΔF0-based features, many of which are

113   highly correlated (**Figure S2**). Implementing a feature selection process to remove redundancy

114   (described in **Supplemental Methods**), we extracted 30 voice features (**Table S2**, their

115   distributions are in **Figure S3** and **Figure S4**), providing a comprehensive characterization of the

116   speaker's prosodic patterns, pitch, and intonation(41,42) We calculated statistics and functions

117   based on the time series of F0, and its differential values, namely pitch change speed (ΔF0). These

118   statistics and functions include mean values, quartiles, range, and regression coefficients, which

119   capture pitch trends and dynamics in speech.

6

120

121  *A two-stage meta-analysis identifies 20 features associated with MDD*

122  Our association analysis had to account for a number of potential confounds. Not everyone in the

123  study spoke the same language: 60% of the subjects spoke in standard Mandarin, whereas the rest

124  spoke either local languages or Mandarin with local accents. This might not matter if language

125  differences were randomized with respect to case status, but uneven case/control ratios between

126  hospitals could confound the analysis. Similarly, the quality of the recordings varied, potentially

127  confounding the association testing. While we made every attempt to ensure that the location of

128  interviews was comparable (in outpatient departments) and that the interviews were carried out in

129  the same way by clinically experienced interviewers that we had trained (Supplementary Methods),

130  these and other, unknown confounders might impact the voice features.

131       We dealt with these issues as follows. First, during the process of identifying the patients'

132  voice segments, we annotated background noise and language. The noise was categorized into five

133  levels, and a binary indicator tagged whether the subjects' speech was in standard Mandarin or not.

134  We included these features as covariates in our analyses.

135       Second, to account for differences between hospitals, we implemented a two-stage meta-

136  analysis, in which associations were first calculated at the hospital level, including demographic

137  features as covariates (**Table S1**), noise levels, and speech indicators, and subsequently pooling

138  results using a random-effects model. We selected 27 hospitals with at least 100 individuals,

139  yielding a total subject count of N = 5,681 (**Figure S1**). By analyzing cases and controls within

140  each hospital first and then combining the results in a meta-analysis, we alleviated the risk of site-

141  specific confounders. We identified 20 features significantly associated with MDD at a 5% FDR

142  threshold (**Table 1**; a description of the vocal features is given in **Table S2**). All features were

143  standardized  (to give a standard deviation of 1) before analysis so that beta coefficients can be

144  compared and interpreted. Eighteen features were significant under a family-wise error rate control

145  using Bonferroni (P-value < 0.0017), with 14 features showing absolute β coefficients greater than

146  0.3. The most significantly associated feature is a $\Delta F0$ measure (interquartile range; $\beta =$

147  $-1.07, \text{SE} = 0.07, P_{\text{FDR}} = 1.1 \times 10^{-49}$).

148  To test the sensitivity of our results to differences between cases and controls, we compared

149  analyses with and without the inclusion of 20 genetic principal components (PCs). The results of

150  this analysis are presented in **Table S3**, and a comparison of the betas with and without adjusting

151  for genetic PCs is presented in **Figure S5**. The correlation between betas in the two analyses was

152  $r = 0.99, P = 8.7 \times 10^{-28}$. These results demonstrate a high degree of consistency in the estimated

153  association effects, irrespective of the adjustment for genetic PCs. There is no overall decrease in

154  significance with the inclusion of the genetic covariates, implying that cases and controls are

155  overall adequately matched by location, which implicitly also means matching by accent.

156

157  *Replication in an independent sample*

158  We evaluated the 20 associated features in an independent sample. The replication sample was

159  collected six years after the discovery CONVERGE data, using the same selection criteria and

160  interview protocol (described in **Supplemental Methods**). The replication sample collected data

161  from three new hospitals, and one that was part of the CONVERGE sample. It used none of the

162  same interviewers and the participants did not overlap with the discovery sample. A description of

163  the sample is given in **Table S4**.  While the replication sample is smaller than the discovery sample

164  (1,084, **Figure S1**), power to detect twelve of the features was greater than 80% (**Table 1**). Again,

165  we listened to all recordings, annotated them for quality and accent, extracted prosodic features,

8

166    analyzed the association within each hospital, and combined results by meta-analysis. As **Table 1**

167    shows, 14 features exceeded a Bonferroni corrected threshold of 0.0025 (0.05/20) and 16 exceeded

168    an FDR 5% threshold. We argue that this strong replication finding in an independent sample,

169    excludes systematic bias in the way recordings were made, the way interviews were conducted,

170    and differences in accent among subjects and differences between hospitals.

171         Finally, we jointly analyzed the discovery and replication samples by meta-analysis and

172    show the results in **Table 1**. Eighteen features exceeded the Bonferroni corrected threshold and

173    all exceeded the 5% FDR threshold. The three most significantly associated features were $\Delta$F0

174    interquartile range ($\beta = -1.07, \mathrm{SE} = 0.07, \mathrm{P_{FDR}} = 6.8 \times 10^{-58}$), which measures the range

175    between the 25th and 75th percentile of pitch change speed ($\Delta$F0), $\Delta$F0 maximum ($\beta =$

176    $-0.97, \mathrm{SE} = 0.07, \mathrm{P_{FDR}} = 1.8 \times 10^{-48}$), which measures the highest value of pitch change speed,

177    and time with F0>90th percentile ($\beta = -0.80, \mathrm{SE} = 0.06, \mathrm{P_{FDR}} = 1.3 \times 10^{-44}$), which measures

178    the amount of time the pitch stays above the 90th percentile of its range. We found that there was

179    no significant heterogeneity in the association effects between Mandarin speakers and non-

180    Mandarin speakers (**Supplemental Results**).

181

182    *Differences between case and control interviews do not account for the associations*

183    We considered next one additional potential confound: the possible impact of the questions asked

184    at interview. The interview for cases is typically more than twice as long as for controls, as we ask

185    about past occurrences of depression and associated stressful life events. Could the emotion

186    associated with this questioning alter speech in such a way as bias our findings? We conducted a

187    sensitivity analysis on responses to neutral questions to check if the effects remain consistent

188    across contexts.

9

189    We selected two questions from the demographic section of the interview based on their

190    high response rates (**Table S5**) and neutral nature. These question (D2.A: "What is your date of

191    birth?" and D10: "How much do you weigh while wearing indoor clothing?") were chosen because

192    they are unlikely to trigger emotional differences between MDD cases and controls. We identified

193    533 subjects with voice response to question D2.A and 617 to question D10. The average segment

194    durations were 3.37 seconds (SD=3.05), and 8.47 seconds (SD=2.69), respectively. For each

195    question, we used the corresponding segments to extract the 16 pitch features that were associated

196    with MDD in our main analysis. Using the two-stage meta-analysis method again, we re-estimated

197    their associations. Due to the small sample sizes, power to detect effects was low so we used a

198    one-sided binomial sign test to test consistency in the direction of association effects between the

199    two analyses.

200    The estimated association effects in context-constrained analysis are reported in **Table S6**.

201    We found that for question D10, three out of 16 pitch features maintained significant associations

202    with MDD at FDR<0.05. Remarkably, 15 out of 16 features showed the same direction of

203    association effects, a fraction significantly higher than chance (Binomial P= 0.00026). For D2.A,

204    despite the average duration being only 3.37 seconds, three features achieved nominal significance

205    for associations (uncorrected P<0.05), and 12 out of 16 pitch features showed consistent directions

206    of association effects (Binomial P= 0.038). In total, 12 out of 16 voice features showed consistent

207    directions of association effects across all four analyses. We conclude that the findings from the

208    main analyses are not biased by the context of the interview.

209    As a summary for these analyses, **Figure S6** shows the effect sizes (beta coefficients) and

210    the 95% confidence intervals for the association between 16 voice F0/ΔF0 features and MDD in

211    the discovery (CONVERGE), replication and single segment analyses.

10

212

213 *Genetic Correlations Between Pitch Features and MDD*

214 Cases for the CONVERGE study were identified as those who have a history of recurrent MDD,

215 and though all were ascertained through hospitals, many were in remission. This raises the

216 important question as to what the association with voice features represents: does it reflect their

217 current low mood, compared to controls, or does it reflect their history of MDD? We addressed

218 this question in the following way.

219 To see if any voice features correlated with current mood, we used a standard assessment

220 of current mood for subjects, the depressive symptom checklist (SCL)(55). These data were only

221 available for the replication sample. The distributions of SCL scores for cases and controls are

222 presented in **Figure S7**. Of the 16 pitch features, 11 showed a significant association with current

223 depressive symptoms, after FDR correction (**Table 2**).

224 These results confirm that most of the features we found to be associated with MDD are

225 correlated with current mood (the relatively smaller sample for this analysis cannot exclude the

226 possibility that all features are thus associated). To examine whether the association reflected a

227 genetic effect common to both variability in vocal features and susceptibility to MDD we estimated

228 the SNP-based heritability for each of the 16 pitch features (this analysis was carried out with

229 CONVERGE data, the only group for which there are genetic data(32)). Results are presented in

230 **Table 3**. Four features were heritable at FDR<0.05. We repeated the heritability analyses adjusting

231 for more genetic PCs to determine whether population structure might contribute to the correlation

232 and found that the heritability remained significant even after adjusting for as many as 60 genetic

233 PCs (**Table S7**). While we cannot rule out the possibility that all vocal features are to some extent

234 heritable (our sample size is too small to confidently detect heritabilities of less than 10%), **Table**

11

235 **3** shows that SNP-based heritability varies significantly: the estimated confidence interval of the

236 heritability for two, $\Delta$F0_maxPos and $\Delta$F0_qregc1, lie outside those for $\Delta$F0_iqr1-3. The low

237 heritability of these features indicate the genetic effects are unlikely to be the only contributing

238 factor to the association with MDD. We did not find any genome-wide significant SNPs for these

239 heritable features (**Supplemental Results**), presumably owing to the limited sample size.

240 We estimated the genetic correlation with MDD for the four features with evidence of

241 heritability and found that three $\Delta$F0 features had significant genetic correlations (**Table 3**). They

242 were: 1) the interquartile range (IQR1-3), quantifying the variation of speed in pitch change; 2)

243 the kurtosis, signaling the extremity of speed in pitch change; and 3) the maximum, representing

244 the speed of the fastest pitch change. There was no detectable genetic correlation with MDD for

245 one heritable feature, $\Delta$F0_kurtosis. Again, while we cannot exclude the possibility of some

246 degree of genetic correlation for this and other features, our results indicate that genetic effects

247 alone cannot explain the association with MDD for all features. The vocal features index a

248 composite of heritable and non-heritable contributions to mood change.

249

250 *Associations Between Pitch Features and MDD Symptoms, Risk Factors, and Comorbidities.*

251 The deep set of phenotypes available in CONVERGE, which includes MDD symptoms,

252 environmental risk factors, comorbid disease, and suicidality, permit us to explore other

253 associations for the vocal features associated with MDD. For these exploratory analyses we

254 included all 30 voice features and tested association with 30 traits (detailed in **Table S8**). The

255 results of within-case two-stage meta-analysis are shown in **Table S8**. We categorized the traits

256 into six classes (MDD symptoms, MDD clinical features, suicidal features, co-morbid psychiatric

12

257    disease, neuroticism and stressful life events) as we were interested in determining the effects on

258    these categories

259         98 associations are significant at an uncorrected 5% significance threshold (where 45 are

260    expected by chance). Surprisingly, features assessing stressful life events showed the greatest

261    enrichment of low P-values. After applying a Bonferroni correction for the 900 tests

262    ($P<0.05/900=5.5 \times 10^{-5}$) five associations were significant, four for stressful life events and one

263    for the personality trait neuroticism. Two features replicated (corrected threshold P<0.05/5=0.01).

264    Both associations were between the total number of stressful life events and $\Delta F0$ features,

265    including the IQR1-3 of $\Delta F0$ (16 hospitals in CONVERGE, total N=2,064, $\beta = -0.21, SE =$

266    $0.03,$ uncorrected $P = 1.9 \times 10^{-11}$ ; four hospitals in the replication, total N=295, $\beta =$

267    $-0.20, SE = 0.07,$ uncorrected $P = 0.0078$) and maximum of $\Delta F0$ (16 hospitals in CONVERGE,

268    total N=2,064, $\beta = -0.19, SE = 0.03,$ uncorrected $P = 6.5 \times 10^{-10}$ ; four hospitals in the

269    replication, total N=295, $\beta = -0.20, SE = 0.08,$ uncorrected $P = 0.0074$).

270

271    *Classification Performance*

272    If the voice features are to have any clinical utility, they must not just be associated with MDD,

273    but they must predict it accurately. We took advantage of access to our two independently collected

274    samples (discovery and replication), using the discovery group for training data (n=7,654) and the

275    replication sample (n=1,189) to test the classification performance.

276         We compared the classification performance of a full model against a null model. The null

277    model was a logistic regression (LR) model trained on the covariates. We then trained a full model

278    using the covariates together with the identified voice features in discovery, based on the same LR

279    method. **Table 4** illustrates these comparisons. Integrating voice data significantly enhanced the

13

280   predictive accuracy of our models. Adding voice features to the LR model increased the AUC-

281   ROC from 0.70 to 0.83 and the accuracy from 0.63 to 0.76.

282   We then tested whether the classification results were robust to different machine learning

283   methods. We evaluated this on three established methods suitable for our dataset, support vector

284   machine (SVM), extreme gradient boosting (XGBoost), and multi-layer perceptron (MLP).

285   Results improved prediction, with XGBoost delivering an AUC of 0.90 (sensitivity of 0.85,  and

286   specificity of 0.81). **Figure 1** plots the ROC curves for a null model (using covariates only to

287   classify depression) and the four models using voice features. The precision-recall curve of these

288   models are in **Figure S8**.

289

290   **Discussion**

291   We set out to find voice pitch features associated with MDD. By using a large and homogeneous

292   case-control cohort, a two stage meta-analysis and an independent replication, we provide robust

293   evidence that certain pitch features distinguished MDD cases from matched controls. The

294   associated features were a slower change in pitch and a particularly uneven distribution of these

295   variations. Features measuring the variability and extremity of pitch change speed were heritable

296   and had genetic correlations with MDD, which we interpret to mean that at least some of the

297   association between variation in pitch and susceptibility to depression is genetic in origin.

298   Classification of those with and without depression, based on vocal features, was achieved with

299   an AUC of 0.90, highlighting the potential use of these features as biomarkers for MDD detection

300   and secondary prevention.

301   Establishing a robust, replicable association between voice features and MDD is difficult

302   because of the numerous confounds that could potentially introduce systematic differences

14

303    between cases and controls and thus corrupt our findings. We addressed this concern by using a

304    large, and as far as possible homogeneous sample of depression. Our study used only women with

305    recurrent MDD in a population where many comorbid disorders, such as smoking, alcohol, and

306    drug abuse, are rare or practically non-existent(32). By adopting a two stage meta-analysis to take

307    into account variation between hospitals, and using an independent replication sample, our results

308    are unlikely to be explained by differences between hospitals, location, accent, quality of the

309    recording or the interview questions.

310        Our findings support and extend previous studies which have indicated potential links

311    between pitch patterns and MDD, but were limited by smaller sample sizes or more heterogeneous

312    cohorts(15,16). First, our large sample size provided adequate power to test several pitch features

313    from a standardized features set, providing more fine-grained quantitative evidence for the

314    descriptions of the monotonous speech pattern in MDD than in previous studies. Previous studies

315    have found that depressed people speak more slowly with lower pitch and decreased

316    variability(16,25). Here, our study showed MDD was negatively associated with features

317    measuring how fast pitch changes (the maximum, the 3rd quartile, and the root quadratic mean of

318    $\Delta F0$, **Table 1**), indicating that the reduced rate of change in pitch is a characteristic of voice in

319    MDD patients. We also found that MDD patients spend less time in their upper vocal range (Time

320    with $F0>90^{th}$ percentile, **Table 1**), affirming the "low-pitched" pattern.

321        Second, our results indicate that MDD's pitch dynamics involve more than reduced

322    variability, showing a broader pitch range and more extreme values (range and kurtosis of F0,

323    **Table 1**). Our analysis also revealed an uneven distribution of the speed with which an MDD

324    patient's pitch changes, as shown by the negative association between MDD and the flatness of

325    $\Delta F0$ and the positive association with the kurtosis of $\Delta F0$ (**Table 1**). Overall, these various features

326    enrich our understanding of pitch dynamics, demonstrating a pattern of slower change in pitch, yet

327    with more frequent occurrences of extreme values and pitch change speed.

328         Third, our research examined the relationship between voice features and the effects of

329    current low mood, and the effects of susceptibility to MDD. Some vocal features might be more

330    reflective of a person's underlying propensity towards developing MDD, while others could be

331    more indicative of a current depressive state. We found that some vocal features were indeed

332    heritable, but still correlated with changes in current mood: individuals with an increased genetic

333    risk of MDD may have a smaller value of speed for the fastest pitch change, thus being unable to

334    speak as fast as those without depression. They may show a narrower IQR of pitch change speeds

335    and more frequently occurring extreme changes of pitch (higher kurtosis).  Shared genetic effects

336    exist between at least some $\Delta$F0 features and MDD, but while our low power to detect heritability

337    and genetic correlations raises the possibility that the other features may also be associated in this

338    way, our findings are consistent with vocal features' association with both current low mood and

339    susceptibility to MDD.

340         We also found that two heritable voice features were associated with the number of

341    stressful life events. The reason for these associations is unclear, but suggests the possibility that

342    stressful life events reveal a latent predisposition to depression(56,57), evidenced through a change

343    in vocal features.

344         It is interesting to consider physiological interpretations of our findings. Possibly, changes

345    in pitch could reflect tiredness, or the psychomotor retardation that characterizes an episode of

346    MDD. We do not have physiological assessments of our subjects that characterize these features

347    (all of our data are from interviews, and therefore reflect the subjects' perceptions) but it is worth

348    pointing out that at some level physiological  and psychological contributions will be confounded.

16

349  For example retardation of thought (psychological) can result in a slowness of expression (motor

350  effect), so that in many cases the distinction may not be relevant.

351      Could the features we identified provide clinically useful predictions? A key aim of our

352  research was to find vocal biomarkers that could take on this role. Applying XGBoost to the

353  independent test dataset we obtained an AUC-ROC of 0.90, and a level of accuracy indicating that

354  the features could be useful in identifying cases of MDD. We applied three models for

355  classification (SVM, XGBoost, and MLP) because relying on a single model would not provide

356  sufficient evidence for the robustness and generalizability of the voice features across different

357  machine learning approaches. SVM excels at handling high-dimensional data and constructing

358  optimal decision boundaries, but it assumes that the data is linearly separable in the kernel-

359  transformed feature space(58). If the relationship between the features and depression states is

360  highly non-linear, SVM may struggle to find an optimal solution. In contrast, XGBoost(59), an

361  ensemble of decision trees, can capture complex feature interactions and handle non-linear

362  relationships. However, it may not be as effective as deep learning models like MLP in capturing

363  hierarchical representations of the data. MLP, with its deep learning architecture, can learn

364  intricate patterns and hierarchical representations, but it is prone to overfitting, particularly when

365  the dataset is small or the network is overly complex(60). By applying the voice features and

366  geographical covariates to all three models, we provide a robust justification for the usefulness and

367  generalizability of the voice features in depression state classification.

368      Our results should be assessed with respect to several limitations. First, we only recruited

369  Han Chinese women with recurrent MDD. Our results may not extrapolate to men, those with

370  single episode MDD, or to non-Chinese speakers. Second, our analysis focused solely on pitch

371  features, and future studies should explore other feature types. Neural network-based approaches,

17

372  which can learn directly from raw audio signals, also hold promise for detecting depression but

373  face challenges in portability and interpretability, particularly when addressing complex

374  confounding factors(61). Third, although our context-constrained analysis demonstrates that the

375  signals we found are persistent across speech content, we cannot separate pitch differences due to

376  word choice from pitch differences due to emotional content without additional experiments

377  directly controlling the linguistic context.

378      While we don't know how far results will generalize outside the female Chinese cohort,

379  the findings reveal that vocal features can be used to identify MDD cases with high accuracy and

380  we expect that with improvements, such as the inclusion of additional voice features, even higher

381  predictive accuracy may be obtainable. Our hope is that these findings will further encourage

382  efforts to assess changes in the voice, long understood by experienced clinicians to be a valuable

383  sign, returning it to a more central position in clinical and research work on MDD.

384

385  **Methods**

386  *Participants*

387  We used data from the CONVERGE(32) study, in which women with recurrent MDD were

388  recruited from 58 provincial mental health centers and psychiatric departments of general medical

389  hospitals in 45 cities and 23 provinces of China. Participants were aged between 30-60, with  two

390  or more episodes of MDD that met the DSM-IV criteria(5), with the first episode occurring

391  between ages 14-50. Cases were excluded if they had pre-existing bipolar disorder, nonaffective

392  psychosis, smoking/nicotine dependence (alcohol and substance abuse were virtually absent in this

393  study, so it was not assessed), or mental retardation. Control subjects, screened to exclude a history

394  of MDD, were recruited from patients undergoing minor surgical procedures at general hospitals

18

395  and individuals attending local community centers. The replication study(43) used the same

396  inclusion/exclusion criteria as CONVERGE, and recruited samples from 20 different hospitals in

397  China, with a final sample size of 1,189 (**Figure S1**). This study was approved by Institutional

398  Review Boards at UCLA, Bio-x Center, Shanghai Jiao Tong University (M16033), and local

399  hospitals. All participants provided written informed consent.

400

401  *Data Collection*

402  All subjects went through a semi-structured interview using a computerized assessment system as

403  outlined previously(43) and described in **Supplementary Methods**. Recordings for cases were

404  obtained in outpatient clinics. Controls were recorded in outpatient clinics and in community

405  health centers. Recordings were not standardized and varied in quality and content. All participants

406  provided DNA samples for genetic analysis. Details of DNA sequencing and genotype imputation

407  have been previously reported(32) and described briefly in **Supplementary Methods**.

408

409  *Covariates*

410  The covariates were five demographic variables and two recording quality variables. The

411  demographic variables were age, education level, occupation, marital status, and social class. The

412  recording quality referred to noise level and accent. The noise level and accent label were

413  determined subjectively by the listeners during the process of identifying the patients' voice

414  segments. The noise was categorized into four levels: 1) No noise; 2) Slight noise but the subject's

415  speech was clear; 3) Noise present but the content of the subject's speech could be clearly heard;

416  4) High noise levels and unclear speech. Note that noise level 4 means that the speech, although

417  difficult to understand, can still be comprehended with extra effort. And samples were excluded

19

418     during quality controls stage if the speech is not able to comprehended at all. The accent was a

419     binary label that indicated whether the subjects' speech was in standard Mandarin or not.

420

421     *Voice Data Preprocessing*

422     8,322 subjects the subjects recruited in CONVERGE had interview recordings (**Figure S1**). To

423     obtain the subjects' utterances, a group of undergraduates listened to the recordings to identify any

424     voice segments from the subjects with a duration >2 seconds. Audio samples were excluded where

425     the noise was so prominent that the content of the subject's speech could not be understood, which

426     yielded 7,654 subjects with available segments. All segments from the same subject were

427     concatenated in the order in which they occur in the interview and down sampled to 8 kHz. Two

428     postgraduate psychological students listened to all the segments to ensure that no speech voice

429     other than the subjects was included in the segments and that no words were cut off mid-way.

430     The preprocessing procedure in the replication study was the same as in CONVERGE. Of the

431     initially recruited 1,301 participants (551 cases), 1,189 subjects (including 490 cases) had available

432     voice segments (**Figure S1**). All segments from the same subject were concatenated into one, to

433     extract voice features for replicating the association between voice and MDD and identifying the

434     voice features associated with SCL (see distribution of the audio segment length in **Figure S9**).

435     The speech segments in the replication study were further annotated to indicate the specific

436     question that prompted each spoken response (**Table S5**). This additional level of data analysis

437     was introduced to better understand the relationship between the interview content and the

438     participants' speech patterns. We additionally extracted the same voice features on the non-

439     concatenated segments corresponding to a single question for a sensitivity analysis, which we

440     referred to as, the single-item analysis (described below).

20

441

*Voice features*

443 We used the INTERSPEECH 2016 Computational Paralinguistics Evaluation(41,42).

444 Calculations were implemented in the openSMILE python package v2.4.2(62) and described in

445 **Supplementary Methods**. Given that many of the features were highly correlated (for example,

446 the arithmetic and root-quadratic mean of F0, as shown in **Figure S2**), we removed redundant

447 features (described in **Supplementary Methods**), resulting in a set of 15 F0-based features and 15

448 ΔF0-based features. We provide in Supplemental **Table S2** technical definitions of the 30 features

449 used, along with non-technical explanations of what each feature measures.

450

*Two-stage meta-analysis*

452 We used a two stage meta-analytic framework to take into account differences between hospitals.

453 In the first stage, for each hospital a linear regression model was fitted for each F0-related feature

454 as the dependent variable using MDD and covariates as the predictor variables. We applied rank-

455 based inverse normal transformation to the voice features. At stage 2, beta coefficients for MDD

456 and standard errors from stage 1 were pooled using random-effects meta-analysis(63), assuming

457 that the true effect sizes in different sites are not exactly the same but are drawn from a distribution

458 of effect sizes. P-values were FDR-adjusted(64). In the second stage, we repeated the analyses in

459 four hospitals with sample sizes ≥ 100 (N=1,084, **Figure S1**). We performed the same procedure

460 as in the two-stage meta-analysis above. We combined results from both discovery and replication

461 cohorts using random-effects meta-analysis(63).

462

21

463   *Heritability and Genetic Correlations*

464   Heritability and genetic correlations were estimated on the 7,654 subjects in CONVERGE (**Figure**

465   **S1**). The SNP-based heritability used a generalized REML (restricted maximum likelihood)

466   method implemented in LDAK(65). We applied rank-based inverse normal transformation to the

467   voice features and incorporated the above covariates and 20 genetic PCs. P-values were FDR-

468   adjusted.  For heritable voice features, we estimated their genetic correlation with MDD, adjusting

469   for these same covariates and 20 genetic PCs. The genetic correlation was calculated through a

470   bivariate GREML analysis implemented in GCTA(66,67). P-values were FDR-adjusted based on

471   the number of heritable voice features.

472

473   *Associations between pitch features and current mood*

474   To identify biomarkers for current mood, subjects in the replication cohort were given a 16-item,

475   self-administered questionnaire assessing the severity of depression-related symptoms on a five-

476   point distress scale over the past 30 days (subscales for depression symptom checklist, SCL)(55).

477   We used the same two-stage meta-analysis method to estimate the association between the 16

478   voice features and SCL scores. All four hospitals from the replication cohort with sample sizes ≥

479   100 were selected (N=1,084, **Figure S1**). At stage 1, for each hospital, a linear regression model

480   was fitted for each pitch feature as the dependent variable using SCL scores and the covariates as

481   the predictor variables. At stage 2, beta coefficients for SCL and standard errors from stage 1 were

482   pooled using random-effects meta-analyses(63). SCL scores were standardized using rank-based

483   inverse normal transformation. P-values were FDR-adjusted.

484

485   *Classification Model*

22

486    We employed a logistic regression model using seven covariates (age, education level, occupation,

487    marital status, social class, noise level, and accent) to establish a null model. Full models that

488    incorporated both these covariates and voice features were then developed.  We included all 20

489    voice features identified as associated with MDD during the discovery stage, including those not

490    replicated. We used samples available in the discovery group for training data (n=7,654). For the

491    test data set we used the replication sample (n=1,189). There was no re-estimation of weights in

492    the test sample.

493        We compared the results of the full model with the null model based on logistic regression.

494    Then we tested the classification performances using SVM, XGBoost, and MLP as the classifiers.

495    The performance of models was evaluated across several metrics: accuracy, sensitivity, specificity,

496    AUC-ROC, AUC-PR, and F1-score. To identify the best hyperparameters for each model, a grid

497    search with 5-folds cross-validation was employed on the training dataset. The above process was

498    implemented in python with package scikit-learn(68) v1.2.2 and xgboost(59) v2.0.3.

499

500    *Associations Between Pitch Features and MDD Symptoms, Risk Factors, and Comorbidities.*

501    We examined the relationship between the 30 voice F0/$\Delta$F0 features with 33 variables related to

502    MDD, including eight risk factor variables, 11 comorbidity variables, seven symptoms, three

503    variables about suicidality, age of onset, number of MDD episodes, neuroticism, and premenstrual

504    syndrome score (summarized in **Table S8**). We again employed the two-stage meta-analysis

505    procedure. At stage 1, for each hospital, a multivariate linear regression model was fitted for each

506    pitch feature as the dependent variable using one of the above variables and covariates as the

507    independent variables. At stage 2, we used the Q statistics to measure the heterogeneity of the

508    pooled beta coefficients and standard errors. If the heterogeneity test is significant (P<0.05), we

509  then used the random-effects model for meta-analyses, otherwise we used fixed-effects(64). We

510  applied a Bonferroni correction to obtain a 5% significance threshold. Due to the high endorsement

511  rates for certain variables within some hospitals, the hospitals included in the meta-analysis varied

512  depending on the variable being analyzed (for example, if all cases from one hospital did not have

513  suicidal attempts, this hospital would be excluded for the analysis of suicidal attempts at stage 1.).

514  We reported the number of hospitals and sample size for each association along with the meta-

515  results.

516  *Context-Constrained Analysis*

517  We counted the total number of available voice segments for each question in the replication cohort

518  and selected the two most frequently answered questions from the demographic section of the

519  interview: D2.A ("What is your date of birth?") and D10 ("How much do you weigh while wearing

520  indoor clothing?"). For each question, we used the corresponding segments to extract the 16 voice

521  features that were associated with MDD in our previous analysis. Finally, we re-assessed the

522  associations between these voice features and MDD through the two-stage meta-analysis method.

523  The limited voice duration and small sample size reduced power to detect a significant signal. We

524  applied the one-sided binomial sign test to determine whether the number of voice features

525  demonstrating consistent directions of association effects between the concatenated segments and

526  the context-constrained segments was greater than expected by chance (that is, a one-sided test of

527  whether this fraction is greater than 0.5).

528

529

530

24

## References

1.  Kendler KS. The genealogy of major depression: symptoms and signs of melancholia from 1880 to 1900. Mol Psychiatry. 2017 Nov;22(11):1539–53.

2.  Kendler KS. The Phenomenology of Major Depression and the Representativeness and Nature of DSM Criteria. AJP. 2016 Aug;173(8):771–80.

3.  American Psychiatric Association. Diagnostic and Statistical Manual of Mental Disorders, Third Edition. Washington, D.C: American Psychiatric Association; 1980.

4.  American Psychiatric Association. Diagnostic and Statistical Manual of Mental Disorders, Revised Third Edition. Washington, D.C: American Psychiatric Association; 1987.

5.  American Psychiatric Association. Diagnostic and statistical manual of mental disorders, Fourth Edition. Washington, D.C: American Psychiatric Association; 1994.

6.  American Psychiatric Association. Diagnostic and statistical manual of mental disorders (DSM-5®). Washington, D.C: American Psychiatric Association; 2013.

7.  Lux V, Kendler KS. Deconstructing major depression: a validation study of the DSM-IV symptomatic criteria. Psychological Medicine. 2010 Oct;40(10):1679–90.

8.  Hyman SE. Can neuroscience be integrated into the DSM-V? Nat Rev Neurosci. 2007 Sep;8(9):725–32.

9.  Kessler RC, Berglund P, Demler O, Jin R, Koretz D, Merikangas KR, et al. The Epidemiology of Major Depressive Disorder Results From the National Comorbidity Survey Replication (NCS-R). JAMA. 2003;289(23):3095–105.

10. Lu J, Xu X, Huang Y, Li T, Ma C, Xu G, et al. Prevalence of depressive disorders and treatment in China: a cross-sectional epidemiological study. The Lancet Psychiatry. 2021 Nov;8(11):981–90.

11. Thornicroft G, Chatterji S, Evans-Lacko S, Gruber M, Sampson N, Aguilar-Gaxiola S, et al. Undertreatment of people with major depressive disorder in 21 countries. The British Journal of Psychiatry. 2017 Feb;210(2):119–24.

12. Guislain J. Orales sur Les Phrénopathies, ou Traitê Thêorique Et Pratique Des Maladies Mentales: Cours Donné A La Clinique Des Êtablissements D'Aliénés A Gand. Vol. 1. Paris, & Bonn,: Gand; 1852.

13. Kraepelin E. Manic-depressive insanity and paranoia. Edinburgh: E. & S. Livingstone; 1921.

14. Sobin C. Psychomotor Symptoms of Depression. A m J Psychiatry. 1997;15.

15. Cummins N, Scherer S, Krajewski J, Schnieder S, Epps J, Quatieri TF. A review of depression and suicide risk assessment using speech analysis. Speech Communication. 2015;71:10–49.

16. Low DM, Bentley KH, Ghosh SS. Automated assessment of psychiatric disorders using speech: A systematic review. Laryngoscope Investigative Otolaryngology. 2020 Feb;5(1):96–116.

17. Nilsonne Å. Acoustic analysis of speech variables during depression and after improvement. Acta Psychiatrica Scandinavica. 1987;76(3):235–45.

18. Nilsonne Å. Speech characteristics as indicators of depressive illness. Acta Psychiatrica Scandinavica. 1988;77(3):253–63.

19. Mundt JC, Snyder PJ, Cannizzaro MS, Chappie K, Geralts DS. Voice acoustic measures of depression severity and treatment response collected via interactive voice response (IVR) technology. Journal of Neurolinguistics. 2007 Jan 1;20(1):50–64.

20. Mundt JC, Vogel AP, Feltner DE, Lenderking WR. Vocal Acoustic Biomarkers of Depression Severity and Treatment Response. Biological Psychiatry. 2012 Oct 1;72(7):580–7.

21. Kuny St, Stassen HH. Speaking behavior and voice sound characteristics in depressive patients during recovery. Journal of Psychiatric Research. 1993 Jul 1;27(3):289–307.

22. Cannizzaro M, Harel B, Reilly N, Chappell P, Snyder PJ. Voice acoustical measurement of the severity of major depression. Brain and Cognition. 2004 Oct 1;56(1):30–5.

23. Alpert M, Pouget ER, Silva RR. Reflections of depression in acoustic measures of the patient's speech. Journal of Affective Disorders. 2001 Sep;66(1):59–69.

24. Pan W, Flint J, Shenhav L, Liu T, Liu M, Hu B, et al. Re-examining the robustness of voice features in predicting depression: Compared with baseline of confounders. Li Z, editor. PLoS ONE. 2019 Jun 20;14(6):e0218172.

25. Wang J, Zhang L, Liu T, Pan W, Hu B, Zhu T. Acoustic differences between healthy and depressed people: a cross-situation study. BMC Psychiatry. 2019 Dec;19(1):300.

26. Schultebraucks K, Yadav V, Shalev AY, Bonanno GA, Galatzer-Levy IR. Deep learning-based classification of posttraumatic stress disorder and depression following trauma utilizing visual and auditory markers of arousal and mood. Psychol Med. 2022 Apr;52(5):957–67.

27. Di Y, Wang J, Liu X, Zhu T. Combining Polygenic Risk Score and Voice Features to Detect Major Depressive Disorders. Frontiers in Genetics. 2021;12:2451.

28. Flint J. The genetic basis of major depressive disorder. Mol Psychiatry [Internet]. 2023 Jan 26 [cited 2023 Jan 31]; Available from: https://www.nature.com/articles/s41380-023-01957-9

29. Hasler G, Drevets WC, Manji HK, Charney DS. Discovering Endophenotypes for Major Depression. Neuropsychopharmacol. 2004 Oct;29(10):1765–81.

30. Kendler KS, Aggen SH, Neale MC. Evidence for multiple genetic factors underlying DSM-IV criteria for major depression. JAMA psychiatry. 2013;70(6):599–607.

26

597    31. Peterson RE, Cai N, Dahl AW, Bigdeli TB, Edwards AC, Webb BT, et al. Molecular Genetic
598        Analysis Subdivided by Adversity Exposure Suggests Etiologic Heterogeneity in Major
599        Depression. AJP. 2018 Jun;175(6):545–54.

600    32. CONVERGE consortium. Sparse whole-genome sequencing identifies two loci for major
601        depressive disorder. Nature. 2015 Jul;523(7562):588–91.

602    33. Andrianopoulos MV, Darrow KN, Chen J. Multimodal Standardization of Voice Among Four
603        Multicultural Populations: Fundamental Frequency and Spectral Characteristics. Journal of
604        Voice. 2001 Jun 1;15(2):194–219.

605    34. Kendler KS, Gardner C, Neale M, Prescott C. Genetic risk factors for major depression in men
606        and women: similar or different heritabilities and same or partly distinct genes? Psychological
607        medicine. 2001;31(4):605.

608    35. Tao M, Li Y, Xie D, Wang Z, Qiu J, Wu W, et al. Examining the relationship between lifetime
609        stressful life events and the onset of major depression in Chinese women. Journal of Affective
610        Disorders. 2011 Dec;135(1–3):95–9.

611    36. Gao J, Li Y, Cai Y, Chen J, Shen Y, Ni S, et al. Perceived parenting and risk for major
612        depression in Chinese women. Psychol Med. 2012 May;42(5):921–30.

613    37. Gan Z, Li Y, Xie D, Shao C, Yang F, Shen Y, et al. The impact of educational status on the
614        clinical features of major depressive disorder among Chinese women. Journal of Affective
615        Disorders. 2012 Feb;136(3):988–92.

616    38. Yang F, Li Y, Xie D, Shao C, Ren J, Wu W, et al. Age at onset of major depressive disorder
617        in Han Chinese women: Relationship with clinical features and family history. Journal of
618        Affective Disorders. 2011 Dec;135(1–3):89–94.

619    39. Shi J, Zhang Y, Liu F, Li Y, Wang J, Flint J, et al. Associations of Educational Attainment,
620        Occupation, Social Class and Major Depressive Disorder among Han Chinese Women. PLOS
621        ONE. 2014 Jan 31;9(1):e86674.

622    40. Li Y, Shi S, Yang F, Gao J, Li Y, Tao M, et al. Patterns of co-morbidity with anxiety disorders
623        in Chinese women with recurrent major depression. Psychol Med. 2012 Jun;42(6):1239–48.

624    41. Schuller B, Steidl S, Batliner A, Hirschberg J, Burgoon JK, Baird A, et al. The
625        INTERSPEECH 2016 Computational Paralinguistics Challenge: Deception, Sincerity &
626        Native Language. In: Interspeech 2016 [Internet]. ISCA; 2016 [cited 2023 Apr 19]. p. 2001–
627        5.        Available        from:        https://www.isca-
628        speech.org/archive/interspeech_2016/schuller16_interspeech.html

629    42. Weninger F, Eyben F, Schuller BW, Mortillaro M, Scherer KR. On the Acoustics of Emotion
630        in Audio: What Speech, Music, and Sound have in Common. Front Psychol [Internet]. 2013
631        [cited        2021        Dec        20];4.        Available        from:
632        http://journal.frontiersin.org/article/10.3389/fpsyg.2013.00292/abstract

27

633    43. Di Y, Wang J, Li W, Zhu T. Using i-vectors from voice features to identify major depressive
634        disorder. Journal of Affective Disorders. 2021 Jun;288:161–6.

635    44. Afshan A, Guo J, Park SJ, Ravi V, Flint J, Alwan A. Effectiveness of Voice Quality Features
636        in Detecting Depression. Interspeech 2018 [Internet]. 2018 Sep [cited 2023 Apr 19]; Available
637        from:        https://par.nsf.gov/biblio/10098305-effectiveness-voice-quality-features-detecting-
638        depression

639    45. Alghowinem S, Goecke R, Epps J, Wagner M, Cohn J. Cross-Cultural Depression Recognition
640        from Vocal Biomarkers. In: Interspeech 2016 [Internet]. ISCA; 2016 [cited 2023 May 23]. p.
641        1943–7.            Available            from:            https://www.isca-
642        speech.org/archive/interspeech_2016/alghowinem16_interspeech.html

643    46. Quatieri TF, Malyska N. Vocal-source biomarkers for depression: a link to psychomotor
644        activity. In: Interspeech 2012 [Internet]. ISCA; 2012 [cited 2022 Jul 7]. p. 1059–62. Available
645        from: https://www.isca-speech.org/archive/interspeech_2012/quatieri12_interspeech.html

646    47. Syed ZS, Schroeter J, Sidorov K, Marshall D. Computational Paralinguistics: Automatic
647        Assessment of Emotions, Mood and Behavioural State from Acoustics of Speech. In:
648        Interspeech 2018 [Internet]. ISCA; 2018 [cited 2023 Nov 27]. p. 511–5. Available from:
649        https://www.isca-speech.org/archive/interspeech_2018/syed18_interspeech.html

650    48. Schuller B, Batliner A, Steidl S, Seppi D. Recognising realistic emotions and affect in speech:
651        State of the art and lessons learnt from the first challenge. Speech Communication. 2011 Nov
652        1;53(9):1062–87.

653    49. Schuller B, Steidl S, Batliner A, Burkhardt F, Devillers L, Müller C, et al. Paralinguistics in
654        speech and language—State-of-the-art and the challenge. Computer Speech & Language. 2013
655        Jan 1;27(1):4–39.

656    50. Eyben F. Real-time speech and music classification by large audio feature space extraction.
657        Springer; 2015.

658    51. Mao K, Wu Y, Chen J. A systematic review on automated clinical depression diagnosis. npj
659        Mental Health Res. 2023 Nov 20;2(1):1–17.

660    52. Xu S, Yang Z, Chakraborty D, Chua YHV, Tolomeo S, Winkler S, et al. Identifying psychiatric
661        manifestations in schizophrenia and depression from audio-visual behavioural indicators
662        through a machine-learning approach. Schizophr. 2022 Nov 7;8(1):1–13.

663    53. Ringeval F, Schuller B, Valstar M, Cummins Ni, Cowie R, Tavabi L, et al. AVEC 2019
664        Workshop and Challenge: State-of-Mind, Detecting Depression with AI, and Cross-Cultural
665        Affect Recognition. arXiv:190711510 [cs, stat] [Internet]. 2019 Jul 10 [cited 2021 Jan 21];
666        Available from: http://arxiv.org/abs/1907.11510

667    54. Hansen L, Rocca R, Simonsen A, Olsen L, Parola A, Bliksted V, et al. Speech- and text-based
668        classification of neuropsychiatric conditions in a multidiagnostic setting. Nat Mental Health.
669        2023 Dec;1(12):971–81.

28

670    55. Derogatis LR. SCL-90: an outpatient psychiatric rating scale-preliminary report.
671         Psychopharmacol Bull. 1973;9:13–28.

672    56. Kendler KS, Karkowski-Shuman L. Stressful life events and genetic liability to major
673         depression: genetic control of exposure to the environment? Psychological Medicine. 1997
674         May;27(3):539–47.

675    57. Kendler KS, Kessler RC, Walters EE, MacLean C, Neale MC, Heath AC, et al. Stressful Life
676         Events, Genetic Liability, and Onset of an Episode of Major Depression in Women. FOC.
677         2010 Jul;8(3):459–70.

678    58. Suthaharan S. Support Vector Machine. In: Suthaharan S, editor. Machine Learning Models
679         and Algorithms for Big Data Classification: Thinking with Examples for Effective Learning
680         [Internet]. Boston, MA: Springer US; 2016. p. 207–35. Available from:
681         https://doi.org/10.1007/978-1-4899-7641-3_9

682    59. Chen T, Guestrin C. XGBoost: A Scalable Tree Boosting System. In: Proceedings of the 22nd
683         ACM SIGKDD International Conference on Knowledge Discovery and Data Mining
684         [Internet]. San Francisco California USA: ACM; 2016 [cited 2024 May 15]. p. 785–94.
685         Available from: https://dl.acm.org/doi/10.1145/2939672.2939785

686    60. Glorot X, Bengio Y. Understanding the difficulty of training deep feedforward neural networks.
687         In: Proceedings of the Thirteenth International Conference on Artificial Intelligence and
688         Statistics [Internet]. JMLR Workshop and Conference Proceedings; 2010 [cited 2024 May 16].
689         p. 249–56. Available from: https://proceedings.mlr.press/v9/glorot10a.html

690    61. Wang J, Ravi V, Flint J, Alwan A. Speechformer-CTC: Sequential modeling of depression
691         detection with speech temporal classification. Speech Communication. 2024 Sep
692         1;163:103106.

693    62. Eyben F, Wöllmer M, Schuller B. Opensmile: the munich versatile and fast open-source audio
694         feature extractor. In: Proceedings of the 18th ACM international conference on Multimedia
695         [Internet]. Firenze Italy: ACM; 2010 [cited 2023 May 24]. p. 1459–62. Available from:
696         https://dl.acm.org/doi/10.1145/1873951.1874246

697    63. Viechtbauer W. Conducting Meta-Analyses in R with the metafor Package. Journal of
698         Statistical Software. 2010 Aug 5;36:1–48.

699    64. Benjamini Y, Hochberg Y. Controlling the False Discovery Rate: A Practical and Powerful
700         Approach to Multiple Testing. Journal of the Royal Statistical Society Series B
701         (Methodological). 1995;57(1):289–300.

702    65. Speed D, Cai N, Johnson MR, Nejentsev S, Balding DJ. Reevaluation of SNP heritability in
703         complex human traits. Nat Genet. 2017 Jul;49(7):986–92.

704    66. Lee SH, Yang J, Goddard ME, Visscher PM, Wray NR. Estimation of pleiotropy between
705         complex diseases using single-nucleotide polymorphism-derived genomic relationships and
706         restricted maximum likelihood. Bioinformatics. 2012 Oct 1;28(19):2540–2.

29

707    67. Yang J, Lee SH, Goddard ME, Visscher PM. GCTA: A Tool for Genome-wide Complex Trait
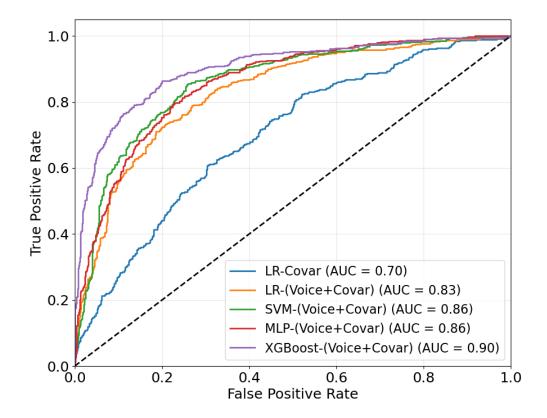708         Analysis. The American Journal of Human Genetics. 2011 Jan;88(1):76–82.

709    68. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-learn:
710         Machine Learning in Python. Journal of Machine Learning Research. 2011;12:2825–30.

711

712

30

713    **Figures**

714    **Figure 1. Receiver Operating Characteristic (ROC) Curve.**



715

716

717    The figure shows the ROC for four models for predicting depression from voice features, and a

718    null model, a logistic regression model trained on demographic covariates only (LR-Covar). The

719    full models are logistic regression (LR), support vector machine (SVM), multi-layer perceptron

720    (MLP), and extreme gradient boosting (XGBoost), trained on voice and covariates (Covar+Voice).

721    AUC: area under the curve.

722

# Tables

## Table 1

| Features | Statistical Functionals | CONVERGE | | | Replication | | | | Combined | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Beta | P | P_FDR | Beta | P | P_FDR | Power | Beta | P | P_FDR |
| ΔF0_iqr1-3 | Interquartile range (3rd -1st) | -1.07 | 3.5E-51 | 1.1E-49 | -1.09 | 1.46E-07 | 0 | 1 | -1.07 | 3.4E-59 | 6.7E-58 |
| ΔF0_percentile99.0 | Maximum (99th percentile) | -0.97 | 1.2E-43 | 1.8E-42 | -0.97 | 5.35E-06 | 0 | 1 | -0.97 | 1.8E-49 | 1.8E-48 |
| F0_upleveltime90 | Time with F0>90th percentile | -0.80 | 1.5E-41 | 1.5E-40 | -0.78 | 5.37E-05 | 0.0001 | 1 | -0.80 | 1.9E-45 | 1.3E-44 |
| ΔF0_kurtosis | Kurtosis | 0.87 | 1.5E-39 | 1.1E-38 | 0.83 | 2.69E-04 | 0.0005 | 1 | 0.86 | 6.4E-42 | 3.2E-41 |
| ΔF0_rqmean | Root quadratic mean | -0.69 | 6.3E-31 | 3.8E-30 | -0.76 | 2.99E-04 | 0.0005 | 0.99 | -0.70 | 6.3E-34 | 2.1E-33 |
| ΔF0_quartile3 | 3rd quartile | -0.78 | 5.4E-30 | 2.7E-29 | -0.88 | 1.16E-05 | 0 | 0.98 | -0.80 | 4.1E-35 | 1.6E-34 |
| ΔF0_minPos | Position of the minimum | -0.60 | 3.3E-25 | 1.4E-24 | -0.57 | 7.05E-07 | 0 | 0.96 | -0.59 | 7.2E-31 | 2.0E-30 |
| ΔF0_amean | Mean | 0.43 | 5.8E-20 | 2.0E-19 | 0.56 | 5.85E-08 | 0 | 0.88 | 0.45 | 3.0E-26 | 7.4E-26 |
| ΔF0_linregc1 | Slope of linear regression | -0.39 | 6.0E-20 | 2.0E-19 | -0.55 | 3.12E-03 | 0.0042 | 0.88 | -0.41 | 8.7E-21 | 1.5E-20 |
| F0_range | Range | 0.57 | 1.0E-19 | 3.1E-19 | 0.48 | 2.35E-03 | 0.0034 | 0.88 | 0.55 | 3.7E-22 | 8.2E-22 |
| ΔF0_maxPos | Position of the maximum | -0.55 | 1.8E-19 | 4.9E-19 | -0.61 | 1.36E-03 | 0.0023 | 0.87 | -0.56 | 5.4E-22 | 1.1E-21 |
| ΔF0_qregc1 | 1st quadratic regression coefficient | 0.38 | 1.0E-17 | 2.6E-17 | 0.47 | 1.04E-04 | 0.0002 | 0.83 | 0.40 | 9.2E-22 | 1.7E-21 |
| ΔF0_flatness | Flatness | -0.43 | 5.5E-14 | 1.3E-13 | -0.39 | 2.85E-05 | 0.0001 | 0.68 | -0.42 | 6.6E-18 | 1.0E-17 |
| F0_lpc2 | 2nd linear prediction coding coefficient | 0.33 | 1.3E-09 | 2.7E-09 | 0.41 | 2.22E-05 | 0.0001 | 0.44 | 0.34 | 6.0E-13 | 8.6E-13 |
| F0_lpc4 | 4th linear prediction coding coefficient | -0.28 | 4.9E-09 | 9.8E-09 | -0.25 | 1.07E-01 | 0.1259 | 0.4 | -0.27 | 2.7E-09 | 3.6E-09 |
| F0_kurtosis | Kurtosis | 0.23 | 1.1E-05 | 2.0E-05 | 0.30 | 1.81E-03 | 0.0028 | 0.19 | 0.24 | 3.6E-08 | 4.5E-08 |
| F0_lpc0 | 0th linear prediction coding coefficient | -0.29 | 9.0E-05 | 1.6E-04 | -0.39 | 1.97E-02 | 0.0246 | 0.14 | -0.30 | 6.1E-06 | 7.2E-06 |
| ΔF0_risetime | Time with which ΔF0 is rising | -0.15 | 4.1E-04 | 6.9E-04 | -0.25 | 1.39E-01 | 0.1539 | 0.1 | -0.16 | 8.9E-05 | 9.9E-05 |
| F0_ff0_maxSegLen | Maximum length of voiced segments with F0 > 0 | 0.19 | 6.3E-03 | 1.0E-02 | 0.12 | 2.69E-01 | 0.2837 | 0.05 | 0.19 | 3.2E-03 | 3.4E-03 |
| F0_rqmean | Root quadratic mean | 0.14 | 9.0E-03 | 1.3E-02 | 0.00 | 9.97E-01 | 0.9974 | 0.05 | 0.12 | 1.7E-02 | 1.7E-02 |

The table shows the names of the prosodic phenotypes, explained in Table S2. Beta values (Beta) P-values (P) and FDR corrected (FDR) are from the logistic regression analysis. Beta coefficients are derived from analyses of normalized voice features (with standard deviation of 1).Results from the CONVERGE study and an independently collected replication are shown. The column headed 'Power' shows the power of the replication study to detect the effect found in the discovery sample. The last three columns ('Combined') are results from a meta-analysis of CONVERGE and the Replication sample

**Table 2. Voice pitch features associated with SCL scores.**

| Feature | Beta | SE | P | P_FDR |
|---|---|---|---|---|
| ΔF0_maxPos | -0.222 | 0.04 | 3.4E-08 | 5.5E-07 |
| ΔF0_percentile99.0 * | -0.296 | 0.074 | 6.7E-05 | 5.3E-04 |
| ΔF0_rqmean | -0.238 | 0.067 | 3.6E-04 | 0.001 |
| ΔF0_quartile3 | -0.286 | 0.081 | 3.9E-04 | 0.001 |
| F0_upleveltime90 | -0.246 | 0.066 | 1.9E-04 | 0.001 |
| ΔF0_iqr1-3 * | -0.309 | 0.097 | 0.001 | 0.004 |
| ΔF0_kurtosis * | 0.225 | 0.081 | 0.005 | 0.01 |
| F0_lpc2 | 0.118 | 0.042 | 0.005 | 0.01 |
| ΔF0_amean | 0.133 | 0.053 | 0.01 | 0.02 |
| F0_kurtosis * | 0.097 | 0.041 | 0.02 | 0.03 |
| F0_range | 0.128 | 0.057 | 0.02 | 0.03 |
| ΔF0_flatness | -0.116 | 0.056 | 0.04 | 0.05 |
| ΔF0_qregc1 | 0.104 | 0.062 | 0.09 | 0.12 |
| ΔF0_linregc1 | -0.103 | 0.067 | 0.12 | 0.14 |
| ΔF0_minPos | -0.13 | 0.086 | 0.13 | 0.14 |
| F0_lpc0 | -0.047 | 0.104 | 0.65 | 0.65 |

The associations between pitch features and SCL scores were estimated in the replication sample using a two-stage meta-analysis. Asterisks indicate heritable features, from **Table 3**.

2

**Table 3. Heritable voice pitch features and their genetic correlation with MDD.**

| Feature | SNP heritability | | | | Genetic Correlation with MDD | | | |
|---|---|---|---|---|---|---|---|---|
| | **h2** | **95% CI** | **P** | **P_FDR** | **rg** | **95% CI** | **P** | **P_FDR** |
| **ΔF0_iqr1-3** | **0.171** | **(0.071, 0.272)** | **0.0004** | **0.006** | **-0.45** | **(-0.77, -0.13)** | **0.03** | **0.04** |
| **F0_kurtosis** | **0.134** | **(0.035, 0.234)** | **0.004** | **0.03** | -0.3 | (-1.14, 0.54) | 0.2 | 0.2 |
| **ΔF0_kurtosis** | **0.125** | **(0.025, 0.225)** | **0.007** | **0.03** | **0.55** | **(0.23, 0.88)** | **0.01** | **0.02** |
| **ΔF0_percentile99.0** | **0.121** | **(0.022, 0.221)** | **0.008** | **0.03** | **-0.7** | **(-1.28, -0.11)** | **4.2E-05** | **0.0002** |
| ΔF0_flatness | 0.105 | (0.007, 0.203) | 0.02 | 0.05 | | | | |
| F0_lpc2 | 0.093 | (-0.005, 0.191) | 0.03 | 0.08 | | | | |
| ΔF0_rqmean | 0.058 | (-0.040, 0.157) | 0.12 | 0.25 | | | | |
| F0_upleveltime90 | 0.044 | (-0.053, 0.141) | 0.18 | 0.32 | | | | |
| ΔF0_minPos | 0.03 | (-0.067, 0.128) | 0.27 | 0.40 | | | | |
| ΔF0_amean | 0.024 | (-0.073, 0.120) | 0.31 | 0.40 | | | | |
| ΔF0_quartile3 | 0.019 | (-0.078, 0.116) | 0.35 | 0.40 | | | | |
| F0_lpc0 | 0.019 | (-0.079, 0.116) | 0.35 | 0.40 | | | | |
| F0_range | 0.001 | (-0.096, 0.098) | 0.49 | 0.49 | | | | |
| ΔF0_linregc2 | -0.004 | (-0.101, 0.093) | 0.47 | 0.49 | | | | |
| ΔF0_qregc1 | -0.028 | (-0.124, 0.068) | 0.28 | 0.40 | | | | |
| ΔF0_maxPos | -0.065 | (-0.160, 0.030) | 0.09 | 0.21 | | | | |

3

**Table 4. Classification performance using the identified voice pitch features.**

LR, Logistic Regression. XGBoost, Extreme Gradient Boosting SVM, Support Vector Machine.

MLP, Multi-layer Perceptron

| Model | Method | Feature | AUC-ROC | AUC-PR | Sensitivity | Specificity | F1 Score | Accuracy |
|---|---|---|---|---|---|---|---|---|
| NULL | LR | Demographic Covariates | 0.70 | 0.62 | 0.75 | 0.53 | 0.63 | 0.62 |
| FULL | LR | Demographic Covariates + Voice | 0.83 | 0.77 | 0.81 | 0.70 | 0.73 | 0.75 |
| | SVM | Demographic Covariates + Voice | 0.86 | 0.80 | 0.88 | 0.67 | 0.75 | 0.76 |
| | MLP | Demographic Covariates + Voice | 0.85 | 0.80 | 0.86 | 0.69 | 0.75 | 0.76 |
| | XGBoost | Demographic Covariates + Voice | 0.90 | 0.88 | 0.85 | 0.81 | 0.80 | 0.82 |