

A mental health assessment method based on emotion level derived from voice

Shuji Shinohara (✉ shinohara@bioeng.t.u-tokyo.ac.jp)

The University of Tokyo <https://orcid.org/0000-0001-8442-836X>

Mitsuteru Nakamura

The University of Texas

Yasuhiro Omiya

PST Inc

Naoki Hagiwara

AGI Inc.

Shunji Mitsuyoshi

The University of Tokyo

Hiroyuki Toda

National Defense Medical College

Taku Saito

National Defense Medical College

Masaaki Tanichi

National Defense Medical College

Aihide Yoshino

National Defense Medical College

Shinich Tokuno

The University Tokyo

Research article

Keywords: mental health assessment, vitality, mental activity, voice emotion analysis, non-invasiveness

Posted Date: December 7th, 2019

DOI: <https://doi.org/10.21203/rs.2.18354/v1>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 **Title**

2 A mental health assessment method based on emotion level derived from voice

3 **Authors**

4 Shuji Shinohara^a, Mitsuteru Nakamura^b, Yasuhiro Omiya^c, Naoki Hagiwara^d, Shunji

5 Mitsuyoshi^a, Hiroyuki Toda^e, Taku Saito^e, Masaaki Tanichi^e, Aihide Yoshino^e, and

6 Shinich Tokuno^f

7 **Affiliations**

8 ^aDepartment of Bioengineering, Graduate School of Engineering, The University of

9 Tokyo, 7-3-1 Hongo, Bunkyo-ku 113-8656, Japan. Emails: [shinohara@bioeng.t.u-](mailto:shinohara@bioeng.t.u-
tokyo.ac.jp)

10 [tokyo.ac.jp](mailto:shinohara@bioeng.t.u-tokyo.ac.jp) (SS); mitsuyoshi@bioeng.t.u-tokyo.ac.jp (SM)

11 ^bHealth Science Center at San Antonio, The University of Texas, San Antonio, TX,

12 USA. Emails: m-nakamura@m.u-tokyo.ac.jp (MN)

13 ^cPST Inc., Yamashita-cho 2, Naka-ku, Yokohama, Japan. Email: [omiya@medical-](mailto:omiya@medical-
pst.com)

14 [pst.com](mailto:omiya@medical-pst.com) (YO)

15 ^dAGI Inc., 1-9-1-8F, Higashi-shimbashi, Minato-ku, Tokyo, Japan. Email:

16 hagiwara@agi-web.co.jp (NH)

17 °Department of Psychiatry, National Defense Medical College, Namiki 3-2,
18 Tokorozawa, Japan. Emails: toda1973@gmail.com (HT); t.saito3025@gmail.com,
19 (TS); mtanichi@gmail.com (MT); aihide@ndmc.ac.jp (AY)

20 †Graduate School of Medicine, The University of Tokyo, 7-3-1 Hongo, Bunkyo-ku,
21 Tokyo, Japan. Emails: tokuno@m.u-tokyo.ac.jp (ST)

22

23 **Corresponding author**

24 Shuji Shinohara

25 Tel: +81-3-5841-0439

26 Fax: +81-3-5841-7798

27 E-mail: shinohara@bioeng.t.u-tokyo.ac.jp

28 **Abstract**

29 **Background:** In many developed countries, mental health disorders have become a
30 problem, and the economic loss due to treatment costs and interference with work is
31 immeasurable. Therefore, a simple technique must be developed to determine
32 individuals' depressive state and stress levels. Voice analysis using smartphones is not
33 only noninvasive, it does not require a dedicated device; thus, it can be performed
34 conveniently and remotely. Consequently, we developed a method to assess individuals'
35 mental health levels using emotional components contained in the human voice.

36 **Methods:** We proposed two indices of mental health: a short-term index (vitality) and
37 mental activity calculated from long-term trends in vitality. We used the voices of
38 healthy individuals (men: $n = 10$, $M_{\text{age}} = 42.7 \pm 6.0$ years; women: $n = 4$, $M_{\text{age}} = 35.0 \pm$
39 14.4 years) and patients with major depression (men: $n = 19$, $M_{\text{age}} = 43.7 \pm 11.0$ years;
40 women: $n = 11$, $M_{\text{age}} = 53.9 \pm 8.2$ years). For patients, simultaneously with voice
41 collection, specialists assessed current depression severity using the Hamilton Rating
42 Scale for Depression (HAM-D).

43 **Results:** A significant negative correlation existed between the vitality extracted from
44 voice and HAM-D score ($r = -0.33$, $p < .05$). We could discriminate the voice data of
45 healthy individuals and patients with depression (judged as moderate or severe by the

46 specialists) with high accuracy using vitality ($p = .0085$, the area under the curve (AUC)
47 of the receiver operating characteristic curve = 0.87). However, there was no significant
48 difference between the vitality of the healthy individuals and the patients judged to be
49 the “no depression group with almost no depressive symptoms,” even if they were
50 outpatients with depression ($p > .1$, AUC = 0.64).

51 **Conclusions:** We developed a method to estimate stress through emotion instead of
52 analyzing stress directly from voice data. By daily monitoring of vitality using
53 smartphones, we can encourage hospital visits for people before they become depressed
54 or during the early stages of depression. This may lead to reduced economic loss due to
55 treatment costs and interference with work.

56 **Trial registration:** Not applicable.

57 **Keywords:** mental health assessment, vitality, mental activity, voice emotion analysis,
58 non-invasiveness

59

60 **Background**

61 In many developed countries, mental health disorders have become a problem
62 [1], and the economic loss due to treatment costs and interference with work is

63 immeasurable [2]. Therefore, a simple technique to determine individuals' depressive
64 state and stress level is desired.

65 Self-administered psychological tests, such as the General Health
66 Questionnaire (GHQ) [3,4] and Beck Depression Inventory (BDI) [5,6], can be used as
67 screening methods for patients with mental health disorders. In addition, a stress-check
68 method that uses biomarkers in saliva [7] and blood has been proposed [8]. Although
69 self-administered psychological tests are useful for early detection and as diagnostic
70 aids, there is a problem with reporting bias—in which specific information such as
71 smoking history and medical history are selectively suppressed or expressed by
72 participants [9]. Stress-check methods that use biomarkers also have problems such as
73 the cost of the test and the burden on the participants during specimen collection; i.e.,
74 they are not convenient.

75 On the other hand, with the recent widespread use of smartphones, pathological
76 analysis using voice data has become popular [10-12]. Voice analysis using
77 smartphones is not only noninvasive, it does not require a dedicated device; thus, it can
78 be performed conveniently and remotely.

79 The relationship between mental illness and voice has been observed in
80 previous studies; e.g., studies regarding the speaking rate of patients with depression

81 [13-15], studies on switching pause and percent pause of patients with depression
82 [15,16], etc. There are also studies in which Lyapunov exponents and Kolmogorov
83 entropy for the voice of patients with depression were measured using chaos analysis
84 [17]. A study that used frequency analysis showed that the shimmer and jitter in vowels
85 as voiced by patients with depression were higher than those of healthy people, and the
86 first and second formant frequency were low [18]. Zhou and colleagues proposed a new
87 feature derived from the Teager energy operator for the classification of voices under
88 stress [19]. In another study, a method was proposed to assess mental health from
89 envelope information for pitch and speech waveforms [20].

90 On the other hand, stress is known to have an impact on emotions [21], and a
91 method is being developed to estimate stress through emotion instead of analyzing
92 stress directly from voice data [22-24]. Mitsuyoshi and colleagues [22] proposed an
93 algorithm that estimates the expression of emotion from emotion components of the
94 voice—the vocal affect display. In addition, they experimentally analyzed the
95 relationship between this index and stress and estimated individuals' stress level from
96 their voice. In this study, sensibility technology (ST) that analyzes emotion in speakers'
97 voices was used [25-27]. The present study proposes a method to assess the mental

98 health of a speaker from the emotional components in his or her voice using ST with a
99 focus on the relationship between mental health and emotions.

100 **Methods**

101 *Acquisition of Voices*

102 In this study, we collected voice data in two categories—healthy individuals
103 and outpatients with depression. All participants provided written consent. Voice
104 acquisition of the patient group was performed intermittently from August 2013 to
105 October 2014 with outpatients at Kitahara rehabilitation hospital in Japan. Voices were
106 recorded during patients’ conversations with physicians during examination. All data
107 were then confirmed audibly; overlaps with other speakers and background noises were
108 removed manually.

109 Voices of healthy people were acquired from February to mid-May 2015.
110 During the acquisition period, participants worked normally at their jobs without
111 visiting medical facilities for a mental illness. Voice acquisition was continuously
112 performed once every several days; at each time, 14 types of fixed phrases were read
113 aloud twice. Voices were recorded in a quiet environment with little background noise.

114 Voices were recorded by a gun microphone (AT9944: audio-technica, Tokyo,
115 Japan) placed approximately 100 cm from participants, or by a pin microphone

116 (ME52W: OLYMPUS, Tokyo, Japan) attached to the chest at approximately 15 cm
 117 from participants' mouth. The recording device was MS-PST1 (NORITSU KOKI,
 118 Wakayama, Japan; not commercially available).

119 Table 1 shows participants' information per group. It should be noted that the
 120 number of participants and the number of data differed because data may have been
 121 collected multiple times from the same participant on different days. The average
 122 number of data collected per healthy person were 24.4 ± 33.3 for men and 6.3 ± 6.1 for
 123 women. For patients with depression, they were 6.0 ± 2.9 for men and 6.8 ± 3.2 for
 124 women. These collected data were used to create algorithms to calculate vitality and
 125 mental activity.

126 **Table 1. Experimental participant information for algorithm preparation**

Group	Sex	Number of participants	Mean age	Number of data
Healthy	Male	9	42.9 ± 5.6	220
	Female	4	33.3 ± 15.4	25
Major depression	Male	4	54.0 ± 12.0	24
	Female	5	49.4 ± 15.4	34

127 Regarding the above-described recorded voice, a healthy person's voice is a
 128 fixed-phrase utterance. On the other hand, a patient's voice is a free speech in the form
 129 of dialogue with a doctor, and the type of speech differed between a healthy person and
 130 a patient. Further, the recording location differed. To unify both speech types and

131 recording environments, data for algorithm verification were collected at the National
132 Defense Medical College Hospital in Japan with participants' consent. Participants were
133 informed that the anonymity and confidentiality of their data were guaranteed, and that
134 they were free to withdraw at any time. Participants were not rewarded for their
135 participation.

136 First, from December 2015 to June 2016, fixed-phrase reading voices were
137 collected from outpatients with major depression. Table 2 shows 17 types of Japanese
138 fixed phrases that were used for recording. At the time of voice collection, specialists
139 evaluated patients' depression severity using the Hamilton Rating Scale for Depression
140 (HAM-D) [28]. The HAM-D is not a self-assessment-type psychological test; rather,
141 experts such as doctors evaluate the characteristic items of depression symptoms. The
142 purpose of the HAM-D is for a professional to objectively quantify an individual's
143 depressive state. On the other hand, for voices of healthy individuals, in mid-December
144 2016, the same fixed-phrase reading voices as the patients were recorded in the same
145 examination room as the patients. However, for healthy people, severity assessment
146 using the HAM-D was not conducted.

147 **Table 2. Seventeen phrases used for recording**

No.	Phrase in Japanese	Purpose (meaning)
------------	---------------------------	--------------------------

1	I-ro-ha-ni-ho-he-to	Non-emotional (no means like “a-b-c”)
2	Honjitsu ha seiten nari	Non-emotional (It is fine today)
3	Tsurezurenaru mama ni	Non-emotional (Having nothing to do)
4	Wagahai ha neko dearu	Non-emotional (I am a cat)
5	Mukashi mukashi aru tokoro ni	Non-emotional (Once upon a time, there lived)
6	a-i-u-e-o	Check pronunciation of vowel sounds (no means like “a-b-c”)
7	Ga-gi-gu-ge-go	Check sonant pronunciation (no means like “a-b-c”)
8	Ra-ri-ru-re-ro	Check liquid sound pronunciation (no means like “a-b-c”)
9	Pa-pi-pu-pe-po	Check p-sound pronunciation (no means like “a-b-c”)
10	Omoeba tooku he kita monda	Non-emotional (When thinking, I have come to the far place)
11	Garapagosu shotou	Check pronunciation (Galápagos Islands)
12	Tsukarete guttari shiteimasu.	Emotional (I am tired/dead tired)
13	Totemo genki desu	Emotional (I am very cheerful)
14	Kinou ha yoku nemuremashita	Emotional (I was able to sleep well yesterday)
15	Shokuyoku ga arimasu	Emotional (I have an appetite)
16	Okorippoi desu	Emotional (I am irritable)
17	Kokoroga odayaka desu	Emotional (My heart is calm)

148 These voices were recorded by a pin microphone ME52W (OLYMPUS,
149 Tokyo, Japan) attached to the chest about 15 cm from participants’ mouth. The
150 recording device used was Portable Recorder R-26 (Roland, Shizuoka, Japan). Table 3
151 shows participants’ information for algorithm verification. The number of healthy
152 individuals for verification and the number of their voice data were the same because

153 they were collected only once from each healthy participant. Regarding patients, some
 154 participants performed multiple data acquisitions. Seven, three, and one performed data
 155 acquisition twice, three, and four times, respectively. Data were acquired only once
 156 from the remaining 19 people. The recording format of the voices was linear PCM, the
 157 sampling frequency was 11025 Hz, and the number of quantization bits was 16 bits.

158 **Table 3. Experimental participant information for algorithm verification**

Group	Sex	Number of participants	Mean age	Number of data
Healthy	Male	10	42.7 ± 6.0	10
	Female	4	35.0 ± 14.4	4
Major depression	Male	19	43.7 ± 11.0	34
	Female	11	53.9 ± 8.2	12

159 ***Voice Emotion Analysis System***

160 We used software ST Ver. 3.0 (AGI Inc., Tokyo, Japan) [25-27] to extract
 161 emotions from participants' voice. The categories of emotional elements detected by ST
 162 software are: "anger," "joy," "sorrow," "calmness" and "excitement." The strength of
 163 each emotion is represented as an integer value from 0 to 10. A value of 0 means that
 164 the input speech does not contain the emotion at all. A value of 10 means that the input
 165 speech contains the emotion most strongly. The unit of speech emotion analysis by ST
 166 software is "utterance." This is a part of continuous voice divided by breath. When a
 167 silent state changes to a speech state, it is considered that an utterance has started. When

168 the speech state continues for a certain period and changes to silence, it is considered
169 that the utterance has ended. Whether the silent state or the speech state is determined
170 from the volume using the threshold. The threshold was adjusted manually for each
171 recording, as the volume of the audio is affected by the participant and the condition of
172 the recording.

173 *Algorithm*

174 *Vitality and Mental Activity*

175 We proposed two scales—vitality and mental activity—as indices for the
176 degree of mental health obtained through voice analysis. Generally, “vitality” can be
177 defined in diverse ways; however, here, vitality refers to a scale that measures low for
178 patients with illnesses such as depression and high for healthy people. The main
179 difference between vitality and mental activity is the duration of the measurement.

180 Vitality is calculated from the emotional components of voice (calm, anger, joy, sorrow,
181 and excitement) based on short-term voice data such as a single phone conversation or a
182 hospital visit.

183 On the other hand, mental activity is calculated based on vitality data
184 accumulated over a certain period. Vitality changes based on the conditions at the time
185 of measurement in the same manner that blood pressure changes between post workout

186 and at rest. As accurate identification of high blood pressure is possible through long-
187 term monitoring, in this study, we aimed to accurately assess mental health by
188 introducing mental activity.

189 *Vivacity and Relaxation*

190 To calculate vitality, we introduced two new indices: “relaxation” and
191 “vivacity.” To define these indices, we used four out of five indices output by ST:
192 calmness, joy, sorrow, and excitement.

193 The fifth edition of the Diagnostic and Statistical Manual of Mental Disorders
194 describes the characteristics of a major depressive episode as a continuing depressive
195 state with loss of interests and happiness and feeling sorrow and emptiness [29]. In
196 contrast, if there is a component of joy more relative to sorrow in emotion, it is
197 considered a good mental state. Consequently, vivacity for an utterance was defined as
198 follows:

$$Vivacity = \frac{Joy}{Joy + Sorrow} \quad (1)$$

199 Stress and tension are major factors in mental health disorders. On the other hand, the
200 relaxed state is mentally positive; thus, relaxation for an utterance was defined as
201 follows:

$$Relaxation = \frac{Calm}{Calm + Excitement} \quad (2)$$

202 In other words, relaxation increases with the increasing calmness component of
203 emotion and decreasing excitement. Each emotional value output by ST, as well as the
204 excitement, are expressed with integers in the range of 0–10. Therefore, vivacity and
205 relaxation become real numbers in the range of 0.0–1.0. Vivacity and relaxation as
206 defined above were calculated for each utterance. Below, we define vivacity and
207 relaxation for an acquired voice as the mean value for each utterance contained in the
208 acquired voice.

209 *Vitality Calculation Algorithm*

210 Vitality was calculated as the weighted mean of vivacity and relaxation defined
211 in the previous section. Fig. 1 shows a scatter plot of relaxation and vivacity as
212 calculated from data for the algorithm preparation.

213

214 Fig. 1. Scatter plot of relaxation and vivacity. \times and \square show the data of the healthy and
215 patient groups for algorithm preparation. The straight line represents $0.60x + 0.40y =$
216 0.52 .

217

218 The symbol \times in the figure represents the data of the healthy group, while \square represents
219 the data of the patient group. Data are plotted for each voice acquisition. There are 245

220 data for the 13 people in the healthy group, and 58 for the 9 people in the patient group.
221 We added a straight line that separates the healthy group from the patient group ($0.60X$
222 $+ 0.40Y = 0.52$). Based on this line, the vitality for each acquired voice was defined as
223 follows:

$$Vitality = 0.60 \times Vivacity + 0.40 \times Relaxation \quad (3)$$

224 *Mental Activity Calculation Algorithm*

225 Vitality was calculated from short-term voice data such as a single examination
226 or consultation. Therefore, depending on participants' current mood, even healthy
227 people might score low in vitality, while patients may score high. To compensate for
228 such a weakness, mental activity was calculated from long-term trends in vitality.
229 Specifically, to express long-term trend in vitality, we calculated the mean of
230 accumulated vitality ($\overline{Vitality}$).

231 Furthermore, when vitality has little fluctuation and is stagnant at low values, it
232 is determined to have low mental activity. To actualize such a determination, we
233 introduce a new index: standard deviation ($VitalitySD$) that expresses variations in
234 vitalities for utterances contained in acquired voice. Then, the mean of vitality standard
235 deviation of the accumulated acquired voice ($\overline{VitalitySD}$) was calculated.

236

237 Fig. 2. Scatter plot of mean vitality and the mean (standard deviation) of vitality for
238 each participant. \times and \square indicate data for the healthy group and the patient group,
239 respectively. The straight line represents $0.75x + 0.25y = 0.426$.

240

241 Fig. 2 is a scatter plot of the mean vitality and mean standard deviation of
242 vitality for each participant calculated from the data for the algorithm preparation. The
243 number of data plotted were 13 people for the healthy group and 9 people for the patient
244 group ($N = 22$). When calculating the mean, we used all acquired voice data. In the
245 figure, we added a straight line that separates the healthy group and the patient group
246 ($0.75X + 0.25Y = 0.426$). Based on this line and using the mean and standard deviation
247 of vitality, we define mental activity as follows:

$$MindActivity = 0.75 \times \overline{Vitality} + 0.25 \times \overline{VitalitySD} \quad (4)$$

248

249 ***Method of Analysis***

250 According to the definition of Zimmerman and colleagues [30], the data of the
251 patient group were divided into 3 groups by HAM-D score: no depression (≤ 7), mild
252 (8–16), and moderate or severe (≥ 17).

253 The vitality of the four groups (i.e., these three and the healthy group) were compared
 254 among each other. P-values from Tukey-Kramer tests, the area under the curve (AUC)
 255 against Receiver Operating Characteristic (ROC), sensitivity, and specificity were used
 256 to evaluate the classification accuracy of Vitality.

257

258 **Results**

259 ***HAM-D score***

260 The mean values of HAM-D score in each group are shown in Table 4. In
 261 addition, the mild group and the moderate or severe group will be collectively referred
 262 to as the depression group (HAM-D score ≥ 8). The mean HAM-D score for the
 263 depression group was 16.1 ± 7.4 (n = 22). The number of participants in each group was
 264 11 men and 8 women in the no depression group and 5 men and 3 women in the mild
 265 group. All three participants in the moderate or severe group were men.

266 **Table 4. Average value of HAM-D score for each group**

Group	Number of participants	Number of data	HAM-D score
No depression (HAM-D ≤ 7)	19	24	3.1 ± 2.3
Mild (HAM-D = 8–16)	8	13	11.5 ± 3.2
Moderate or severe	3	9	22.8 ± 6.6

(HAM-D \geq 17)

267 ***Performance evaluation of vitality***

268 We evaluated the performance of vitality using the data for algorithm
269 verification shown in Table 3. Fig. 3 shows the relationship between HAM-D score and
270 vitality for 46 data obtained from the patient group. There was a significant negative
271 correlation between the two ($r = -0.33$, $n = 46$, $p < .05$).

272

273 Fig. 3. Relationship between HAM-D score and vitality in the data of patient group for
274 algorithm verification.

275

276 Fig. 4. Comparison of vitality for each group. (a) represents data distribution of a
277 healthy group, the no depression group, and the depression group. (b) shows the data
278 distribution when the depression group is divided into the mild group and the moderate
279 or severe group.

280

281 Fig 4 shows the distribution of vitality scores of the healthy group, the no
282 depression group, the mild group, the moderate or severe group, and the depression

283 group. The mean vitality in each group was 0.60 ± 0.10 ($n = 14$), 0.55 ± 0.10 ($n = 24$),
284 0.51 ± 0.13 ($n = 13$), 0.47 ± 0.07 ($n = 9$), and 0.49 ± 0.11 ($n = 22$), respectively.

285 The Tukey-Kramer test revealed significant differences between the healthy
286 group and the depression group, and between the healthy group and the moderate or
287 severe group ($P_s = .0085$ and $.020$, respectively). We used statistical analysis software R
288 [31].

289 Next, to evaluate the discrimination performance of vitality, the area under the
290 curve (AUC) of the receiver operating characteristic (ROC) curve, the sensitivity, and
291 the specificity were used. Figure 5 shows the ROC curves when using vitality to
292 identify whether the data for verification is for the healthy group or for each patient
293 group. Here, the horizontal axis represents 1-specificity (false positive rate), and the
294 vertical axis represents sensitivity (positive rate).

295

296 Fig. 5. Receiver operating characteristic curves when using vitality to identify
297 groups

298

299 Table 5 shows the performance when the data of the healthy group and each
300 group were distinguished using vitality. The AUC was 0.87, and the sensitivity and

301 specificity were 0.78 and 0.86, respectively regarding the discrimination performance
 302 between the healthy group and the moderate or severe group. On the other hand, both
 303 AUC were less than 0.7 regarding discrimination performance between the healthy
 304 group and the no depression group or mild group.

305 **Table 5. Discrimination ability of vitality**

Group	AUC	Sensitivity	Specificity
Health–Depression	0.76	0.55	0.93
Health–Moderate or severe	0.87	0.78	0.86
Health–Mild	0.69	0.46	0.93
Health–No depression	0.64	0.80	0.50

306 AUC: area under the curve of the receiver operating characteristic curve.

307 **Discussion**

308 In this study, we developed a method to measure mental health using emotional
 309 components contained in voice. Two indicators were proposed: vitality based on short-
 310 term voice data and mental activity calculated from long-term voice data. As shown in
 311 Fig 3, there was a significant negative correlation between vitality and HAM-D score
 312 (i.e., depression severity assessed by a physician). In addition, as shown in Fig. 4, the
 313 group with a higher severity of depression tended to have a lower mean vitality.

314 There was a significant difference between the healthy group and the
 315 depression group, and between the healthy group and the moderate or severe group in

316 vitality. On the other hand, there was no significant difference between the healthy
317 group and the no depression group with almost no depressive symptoms, even if they
318 were outpatients with depression. This suggests the possibility of measuring treatment
319 effects by vitality (i.e., voice). Moreover, as shown in Fig. 5 and Table 5, the voice data
320 of the healthy group and the voice data of the moderate or severe group could be
321 identified with high accuracy using vitality. This suggests the possibility of screening
322 for severe depression in individuals by using voice.

323 In our other study, we verified vitality with Romanian and Russian native
324 speakers [32]. In this verification, BDI tests were conducted simultaneously with voice
325 recordings. There was a significant difference between the mean vitality of the
326 depression high-risk group (BDI scores ≥ 17) and the mean vitality of the depression
327 low-risk group (BDI scores < 17 ; $p < .05$). Specifically, the scores for question 9—
328 concerning suicidal ideation—took a value that ranged 0–3. There was a significant
329 difference between the mean vitality of the suicide low-risk group (0 or 1 points) and
330 the mean vitality of the suicide high-risk group (2 or 3 points; $p < .01$). In the future, we
331 will examine the vitality of native speakers of other languages, such as English.

332 As a limitation of this research, only the fixed phrase read-out speech was used
333 for verification. To apply vitality to free speech such as a call, further verification is

334 required. Furthermore, in the verification data, the number of voices collected for each
335 participant, the sex ratio, and the age were not unified between groups. These
336 differences may be reflected in the features of voice. For example, all participants in the
337 moderate or severe group were men, and the number of participants was as small as
338 three. In the future, it is necessary to acquire a lot of voices of female patients,
339 especially those with severe depression, and to evaluate the performance level of
340 vitality.

341 Further, mental activity was not validated because continuous data could not be
342 collected sufficiently for the same participants in both the healthy group and the patient
343 group. However, comparing Figs. 1 and 2 showing data for algorithm preparation, there
344 is a possibility that mental activity can more accurately identify the data as compared to
345 vitality, which will be addressed in the future.

346 Vitality and the mental activity can be measured by only voice, and their
347 advantages are that they are non-invasive and less expensive as compared to self-
348 administered tests such as the GHQ-30 and BDI and stress-check methods using saliva
349 and blood. Moreover, it is also possible to record day-to-day state changes easily by
350 implementing them on a smartphone or the like.

351 We developed a smart phone application that implemented the algorithm for
352 vitality and mental activity—the Mind Monitoring System (MIMOSYS). We are
353 currently conducting world-wide demonstration experiments using the MIMOSYS [33].
354 In the future, we plan to verify the effectiveness of vitality and mental activity with such
355 large-scale data.

356 **Conclusions**

357 In this study, we developed a method to measure mental health from voice. The
358 algorithm to estimate stress through emotion instead of analyzing stress directly from
359 voice data is novel. The MIMOSYS implemented the algorithm for vitality and mental
360 activity, which is a cost-effective and convenient measurement device. If the correlation
361 between HAM-D score and vitality can be further enhanced, it may be used to aid
362 doctors' diagnoses in the future. By daily monitoring of vitality and mental activity
363 using the MIMOSYS, we can encourage hospital visits for people before they become
364 depressed or during the early stages of depression. This may lead to reduced economic
365 loss due to treatment costs and interference with work.

366

367 **List of abbreviations**

368 MIMOSYS: Mind Monitoring System

369 GHQ: General Health Questionnaire
370 BDI: Beck Depression Inventory
371 HAM-D: Hamilton Rating Scale for Depression
372 ST: sensibility technology
373 ROC: receiver operating characteristic
374 AUC: area under the curve

375

376 **Declarations**

377 *Ethics approval and consent to participate*

378 Ethical approval was obtained from the National Defense Medical Collage Ethics
379 Committee (no. 2248) and the Kitahara Rehabilitation Hospital Ethics Committee (no.
380 3).

381 *Consent for publication*

382 Not applicable.

383 *Availability of data and material*

384 The datasets used and/or analyzed during the current study are available from the
385 corresponding author on reasonable request.

386 *Competing interests*

387 The authors declare that they have no competing interests.

388 ***Funding***

389 This research was supported by the Center of Innovation Program from the Japan
390 Science and Technology Agency and by JSPS KAKENHI [grant numbers JP16K01408
391 and JP15H03002]. The funders had no role in study design, data collection and analysis,
392 decision to publish, or preparation of the manuscript.

393 ***Authors' contributions***

394 S. T. was responsible for the design of the clinical study. A. Y., H. T., T. S., and M. T.
395 were responsible for the execution of the clinical study including patient recruitment
396 and retention and data collection. S. S. conceived the algorithm, analyzed data, and
397 wrote the manuscript. M. N., Y. O., N. H., S. M., and S. T. contributed to the
398 interpretation of study findings. All authors participated in the editing and revision of
399 the final version of the manuscript.

400 ***Acknowledgements***

401 We thank Dr. Shinsuke Kondo for assistance with data collection and all participants for
402 participating. We also thank Editage [<http://www.editage.com>] for English-language
403 editing.

404

405 **References**

- 406 [1] World Health Organization. The global burden of disease: 2004 update. WHO Press,
407 Geneva, Switzerland; 2008. pp. 46-9.
- 408 [2] Kessler RC, Akiskal HS, Ames M, Birnbaum H, Greenberg P, Hirschfield, RM, et al.
409 Prevalence and effects of mood disorders on work performance in a nationally
410 representative sample of U.S. Workers. Am J Psychiatry. American Psychiatric
411 Association Publishing; 2006;163(9):1561-8.
- 412 [3] Goldberg DP, Blackwell B. Psychiatric illness in general practice: a detailed study
413 using a new method of case identification. Br Med J. 1970;2(5707):439-43.
- 414 [4] Goldberg D. Manual of the general health questionnaire. NFER Nelson; 1978.
- 415 [5] Beck AT. A systematic investigation of depression. Compr Psychiatry.
416 1961;2(3):163-70.
- 417 [6] Beck AT, Ward CH, Mendelson M, Mock J, Erbaugh J. An inventory for measuring
418 depression. Arch Gen Psychiatry. 1961;4(6):561-71.
- 419 [7] Takai N, Yamaguchi M, Aragaki T, Eto K, Uchihashi K, Nishikawa Y. Effect of
420 psychological stress on the salivary cortisol and amylase levels in healthy young
421 adults. Arch Oral Biol. 2004;49(12):963-8.

- 422 [8] Suzuki G, Tokuno S, Nibuya M, Ishida T, Yamamoto T, Mukai Y, et al. Decreased
423 plasma brain-derived neurotrophic factor and vascular endothelial growth factor
424 concentrations during military training. *PLoS One*. 2014;9(2):e89455.
- 425 [9] Miquel P. *A dictionary of epidemiology*. 6th ed. Oxford University Press; 2014.
- 426 [10] Arora S, Venkataraman V, Zhan A, Donohue S, Biglan KM, Dorsey ER, et al.
427 Detecting and monitoring the symptoms of Parkinson's disease using smartphones: a
428 pilot study. *Parkinsonism Relat Disord*. 2015;21(6):650-3.
- 429 [11] Rachuri KK, Musolesi M, Mascolo C, Rentfrow PJ, Longworth C, Aucinas A.
430 EmotionSense: a mobile phones based adaptive platform for experimental social
431 psychology research. In *Proceedings of the 12th ACM international conference on*
432 *Ubiquitous computing 2010 Sep 26 (pp. 281-290)*. ACM.
- 433 [12] Lu H, Frauendorfer D, Rabbi M, Mast MS, Chittaranjan GT, Campbell AT, et
434 al. Stresssense: Detecting stress in unconstrained acoustic environments using
435 smartphones. In *Proceedings of the 2012 ACM Conference on Ubiquitous Computing*
436 *2012 Sep 5 (pp. 351-360)*. ACM.
- 437 [13] Cannizzaro M, Harel B, Reilly N, Chappell P, Snyder PJ. Voice acoustical
438 measurement of the severity of major depression. *Brain Cognition*. 2004;56(1):30-5.

- 439 [14] Moore E, Clements M, Peifer J, Weisser L. Analysis of prosodic variation in
440 speech for clinical depression. In Proceedings of the 25th Annual International
441 Conference of the IEEE Engineering in Medicine and Biology Society (IEEE Cat.
442 No. 03CH37439) 2003 Sep 17 (Vol. 3, pp. 2925-2928). IEEE.
- 443 [15] Mundt JC, Snyder PJ, Cannizzaro MS, Chappie K, Geralts DS. Voice acoustic
444 measures of depression severity and treatment response collected via interactive voice
445 response (IVR) technology. *J Neurolinguistics*. 2007;20(1):50-64.
- 446 [16] Yang Y, Fairbairn C, Cohn JF. Detecting depression severity from vocal
447 prosody. *IEEE Transactions on Affective Computing*. 2012;4(2):142-50.
- 448 [17] Shimizu T, Furuse N, Yamazaki T, Ueta Y, Sato T, Nagata S. Chaos of
449 vowel/a/in Japanese patients with depression: a preliminary study. *J Occupat Health*.
450 2005;47(3):267-9.
- 451 [18] Vicsi K, Sztahó D, Kiss G. Examination of the sensitivity of acoustic-phonetic
452 parameters of speech to depression. In 2012 IEEE 3rd International Conference on
453 Cognitive Info Communications (CogInfoCom) 2012 Dec 2 (pp. 511-515). IEEE.
- 454 [19] Zhou G, Hansen JH, Kaiser JF. Nonlinear feature based classification of speech
455 under stress. *IEEE Trans Speech Audio Process*. 2001;9(3):201-16.

- 456 [20] Shinohara S, Mitsuyoshi S, Nakamura M, Omiya Y, Tsumatori G, Tokuno S.
457 Validity of a voice-based evaluation method for effectiveness of behavioural therapy.
458 In International Symposium on Pervasive Computing Paradigms for Mental Health
459 2015 Sep 24 (pp. 43-51). Springer, Cham.
- 460 [21] Lazarus RS. From psychological stress to the emotions: a history of changing
461 outlooks. *Ann Rev Psychol.* 1993;44(1):1-22.
- 462 [22] Mitsuyoshi S, Nakamura M, Omiya Y, Shinohara S, Hagiwara N, Tokuno S.
463 Mental status assessment of disaster relief personnel by vocal affect display based on
464 voice emotion recognition. *Disaster Mil Med.* Springer Nature; 2017;3(1).
- 465 [23] Tokuno S, Mitsuyoshi S, Suzuki G, Tsumatori G. Stress evaluation using voice
466 emotion recognition technology: a novel stress evaluation technology for disaster
467 responders. In *Proceedings of the XVI World Congress of Psychiatry 2014 Sep (Vol.*
468 *2, p. 301).*
- 469 [24] Tokuno S, Tsumatori G, Shono S, Takei E, Yamamoto T, Suzuki G, et al. Usage
470 of emotion recognition in military health care. 2011 Defense Science Research
471 Conference and Expo. IEEE; 2011 Aug.
- 472 [25] Mitsuyoshi S, Ren F, Tanaka Y, Kuroiwa, S. Non-verbal voice emotion analysis
473 system. *Int J ICIC.* 2006;2(4).

- 474 [26] Mitsuyoshi S, Shibasaki K, Tanaka Y, Kato M, Murata T, Minami T, et al.
475 Emotion voice analysis system connected to the human brain. 2007 International
476 Conference on Natural Language Processing and Knowledge Engineering. IEEE;
477 2007 Aug.
- 478 [27] Mitsuyoshi S, inventor; Advanced Generation Interface Inc, assignee. Emotion
479 recognizing method, sensibility creating method, device, and software. United States
480 patent US 7,340,393. 2008 Mar 4.
- 481 [28] Hamilton MA. Development of a rating scale for primary depressive illness. Br
482 J Soc Clin Psychol. 1967;6(4):278-96.
- 483 [29] America Psychiatric Association. Diagnostic and Statistical Manual of Mental
484 Disorders. 5th ed. Washington: D.C.; 2013.
- 485 [30] Zimmerman M, Martinez JH, Young D, Chelminski I, Dalrymple K. Severity
486 classification on the Hamilton depression rating scale. J Affect Disord.
487 2013;150(2):384-8.
- 488 [31] <https://cran.r-project.org/>
- 489 [32] Uruguchi T, Shinohara S, Denis N A, Țaicu M, Săvoiu G, Omiya Y, et al.
490 Evaluation of Mind Monitoring System (MIMOSYS) by subjects with Romanian and

491 Russian as their native language. 40th International Conference of the IEEE
492 Engineering in Medicine and Biology Society; 2018 July.

493 [33] Shinohara S, Omiya Y, Hagiwara N, Nakamura M, Higuchi M, Kirita T, et al.
494 Case studies of utilization of the mind monitoring system (MIMOSYS) using voice
495 and its future prospects. ESMSJ. 2017;7(1):7-12.

Figures

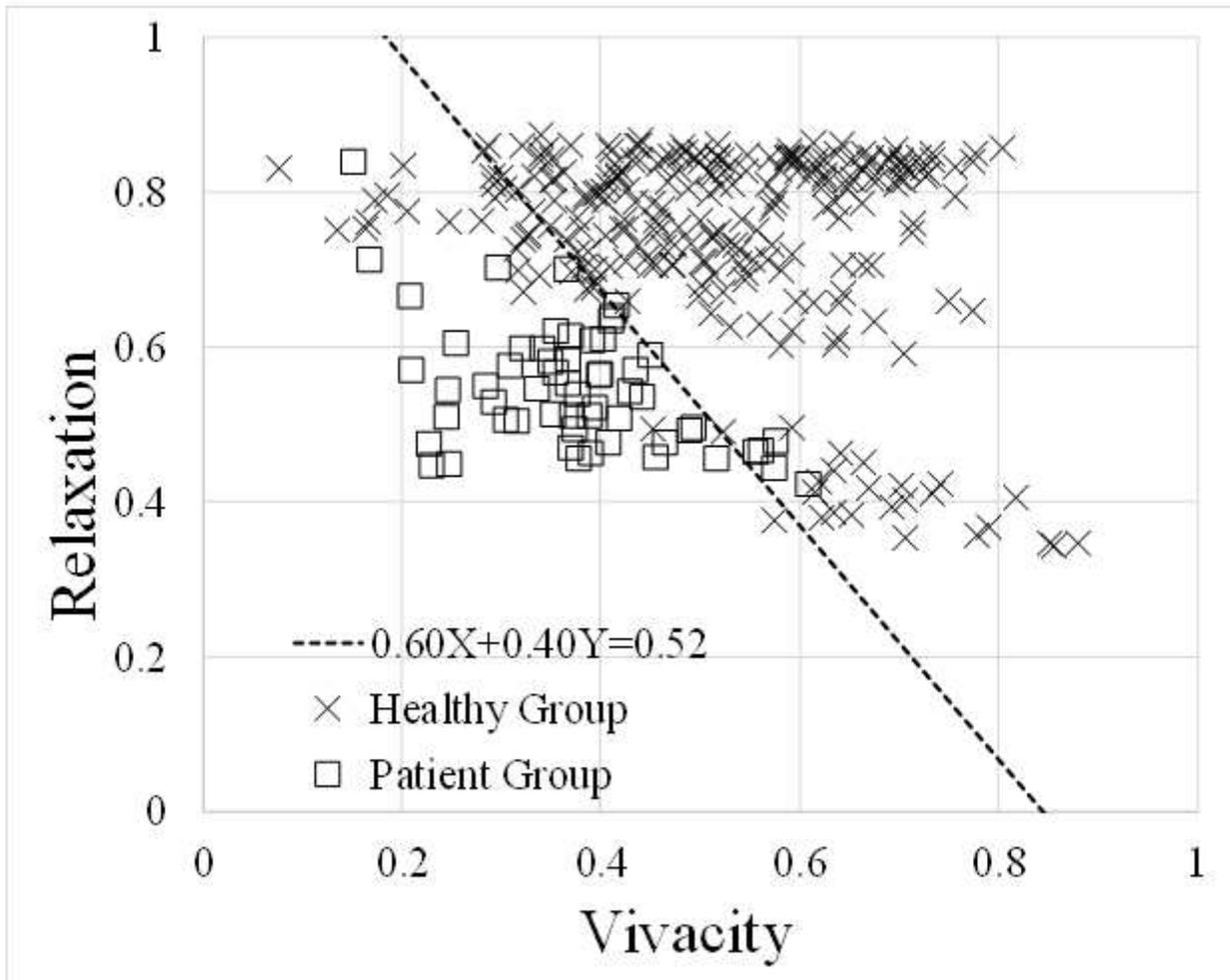


Figure 1

Scatter plot of relaxation and vivacity. \times and \square show the data of the healthy and patient groups for algorithm preparation. The straight line represents $0.60x + 0.40y = 0.52$.

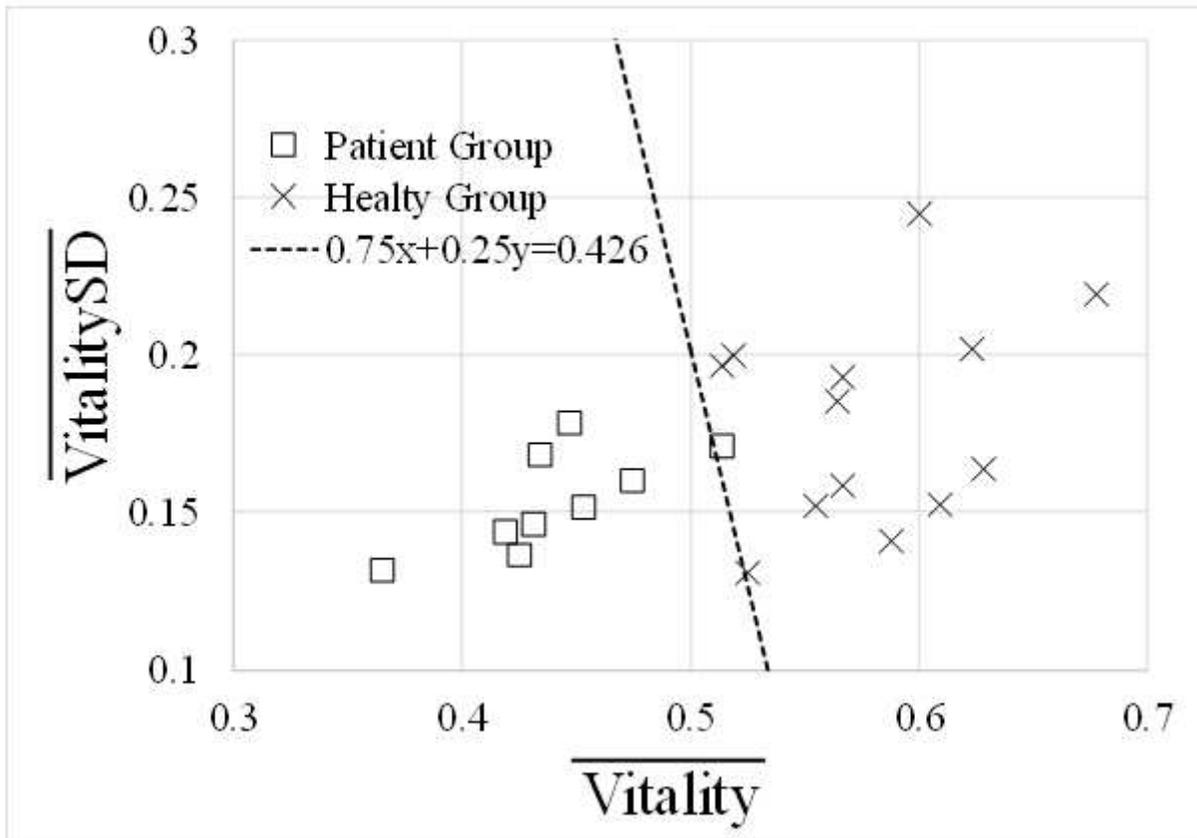


Figure 2

Scatter plot of mean vitality and the mean of standard deviation of vitality for each subject. × and □ indicate data for the healthy group and the patient group, respectively. The straight line represents $0.75x + 0.25y = 0.426$.

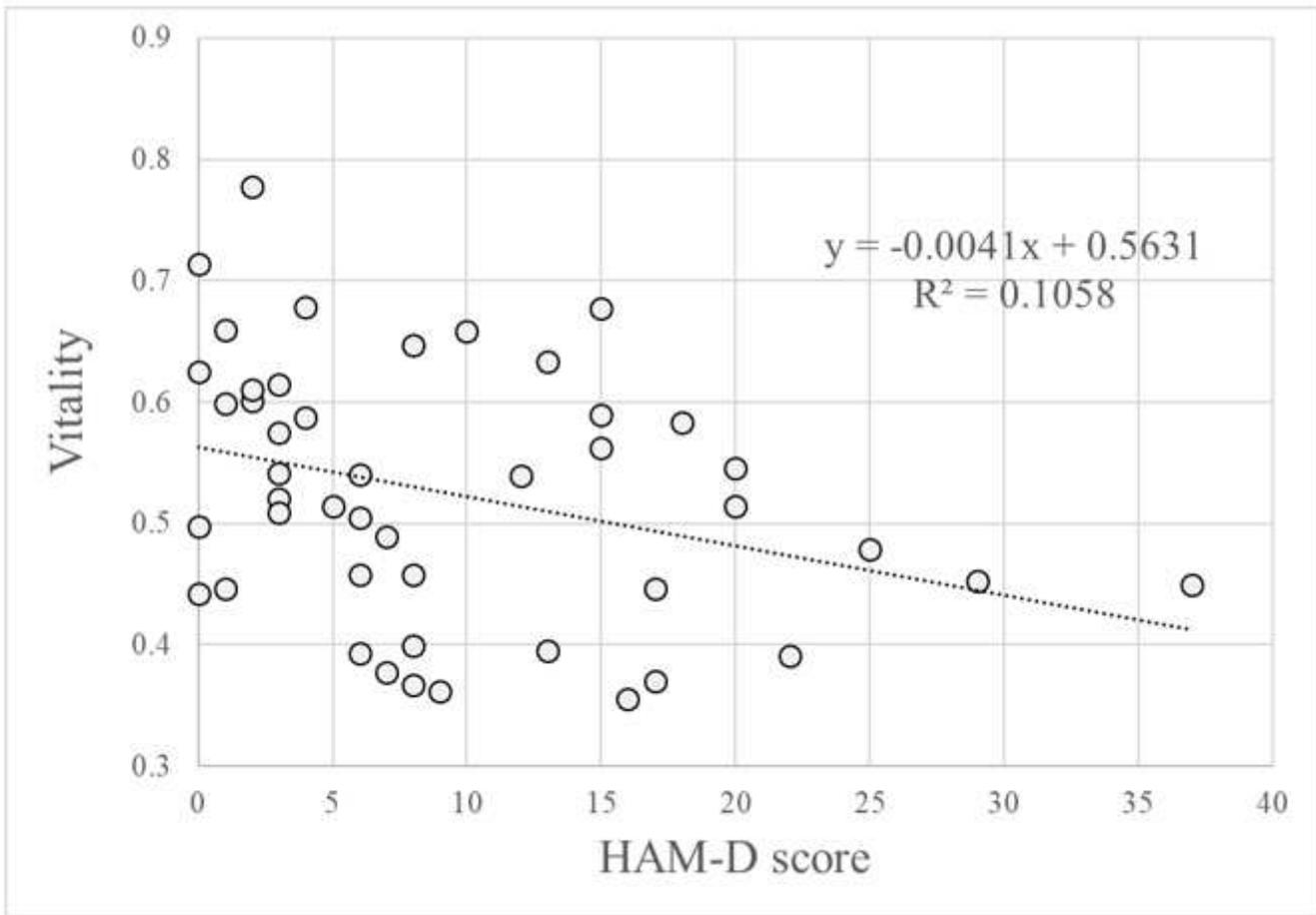
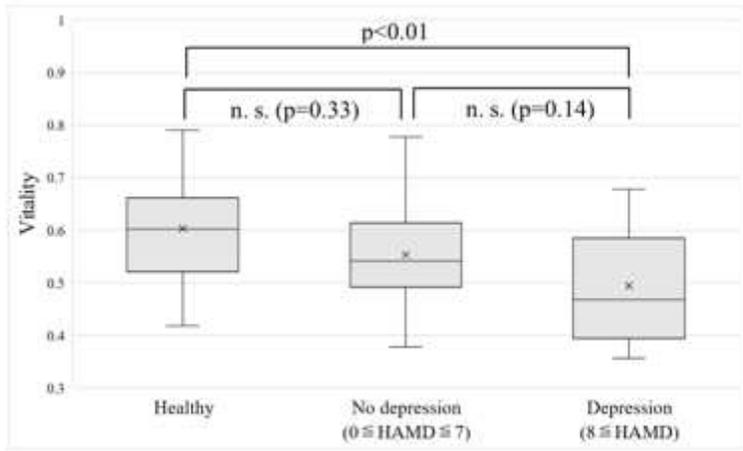
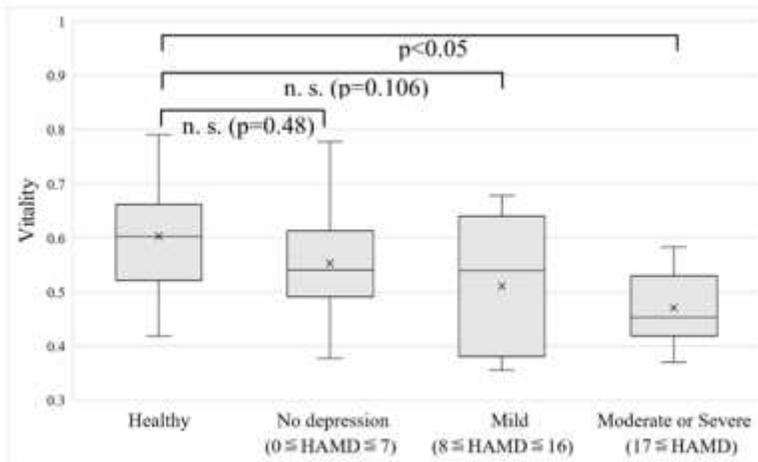


Figure 3

Relationship between HAM-D score and vitality in the data of patient group for algorithm verification.



(a)



(b)

Figure 4

Comparison of the vitality for each group. (a) represents data distribution of a healthy group, a no depression group, and a depression group. (b) shows the data distribution when the depression group is divided into the mild group and the moderate or severe group.

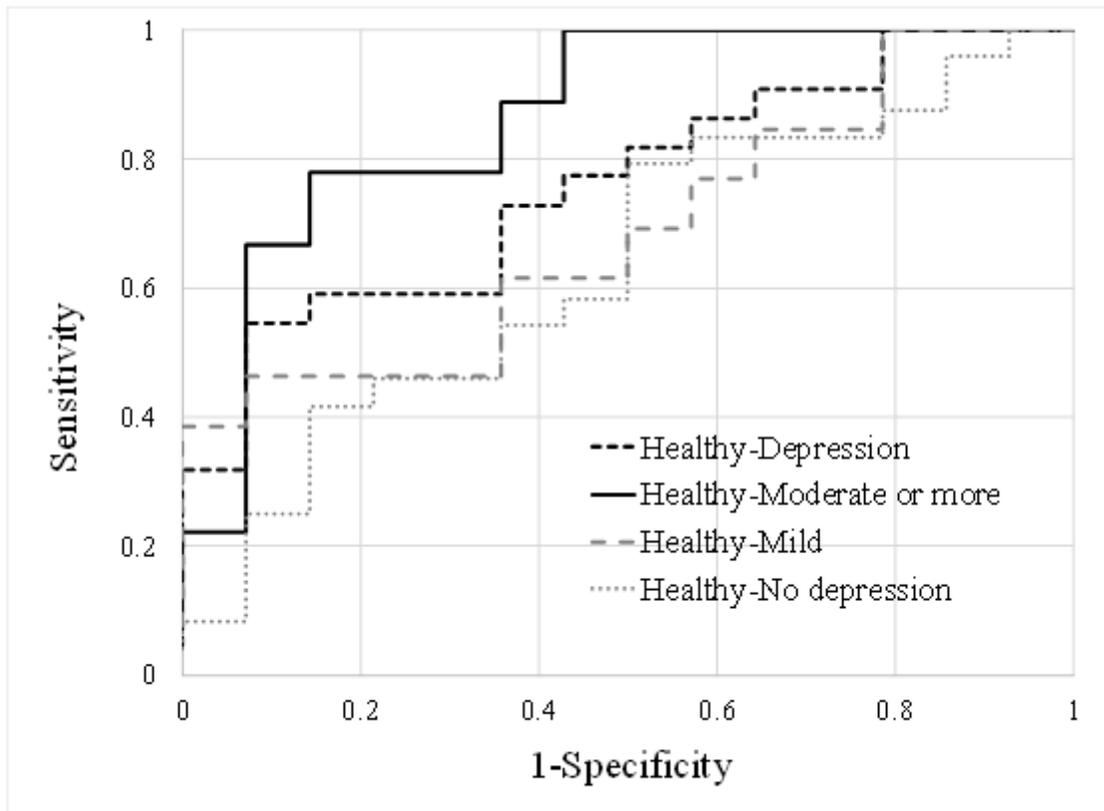


Figure 5

ROC curves when using the vitality to identify whether the data for verification is for the healthy group or for each patient group