

CLASSIFICATION OF BIPOLAR DISORDER, MAJOR DEPRESSIVE DISORDER, AND HEALTHY STATE USING VOICE

MASAKAZU HIGUCHI^{1*}, SHINICHI TOKUNO¹, MITSUTERU NAKAMURA¹, SHUJI SHINOHARA²,
SHUNJI MITSUYOSHI², YASUHIRO OMIYA³, NAOKI HAGIWARA³, TAKESHI TAKANO³, HIROYUKI TODA⁴,
TAKU SAITO⁴, HIROO TERASHI⁵, HIROSHI MITOMA⁶

¹Graduate School of Medicine, The University of Tokyo, Tokyo, Japan. ²Department of Bioengineering, Graduate School of Engineering, The University of Tokyo, Tokyo, Japan. ³PST Inc., Kanagawa, Japan. ⁴Department of Psychiatry, National Defense Medical College, Saitama, Japan. ⁵Department of Neurology, Tokyo Medical University, Tokyo, Japan. ⁶Medical Education Promotion Center, Tokyo Medical University, Tokyo, Japan. Email: higuchi@m.u-tokyo.ac.jp

Received: 21 December 2017, Revised and Accepted: 5 February 2018

ABSTRACT

Objective: In this study, we propose a voice index to identify healthy individuals, patients with bipolar disorder, and patients with major depressive disorder using polytomous logistic regression analysis.

Methods: Voice features were extracted from voices of healthy individuals and patients with mental disease. Polytomous logistic regression analysis was performed for some voice features.

Results: With the prediction model obtained using the analysis, we identified subject groups and were able to classify subjects into three groups with 90.79% accuracy.

Conclusion: These results show that the proposed index may be used as a new evaluation index to identify depression.

Keywords: Voice, Bipolar disorder, Major depressive disorder, Polytomous logistic regression analysis.

© 2018 The Authors. Published by Innovare Academic Sciences Pvt Ltd. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>) DOI: <http://dx.doi.org/10.22159/ajpcr.2018.v11s3.30042>

INTRODUCTION

Due to the stress in modern society, mental health care has become an important issue. In recent years, mental disorders resulting from stress were shown to cause major social loss by declining labor productivity [1]. Therefore, it is important to create a system where emotional problems can be discovered in the early stages and to develop a technology that can easily assess depression and stress.

The current screening methods for mental disorders use biomarkers such as saliva [2], blood [3], electrocardiogram [4], and electroencephalogram [5], but these are invasive and high cost because special equipment and medicines are needed. Some non-invasive methods include self-administered psychological tests such as the General Health Questionnaire [6] and the Beck Depression Inventory [7]. Self-administered psychological tests are relatively easy, but reporting bias cannot be eliminated. Reporting bias is defined as selective under- or overestimation of certain information influenced by the responder's consciousness/unconsciousness [8]. Bipolar disorder and major depressive disorder are two types of depression and divided based on their respective symptoms. Bipolar disorder is a mental disease with alternating manic and depressive states, with a difficult differential diagnosis from unipolar depressive state during a depressive episode [9]. It is particularly difficult to diagnose using a single self-administered psychological test, and the differences between these two diseases can be difficult to identify.

On the contrary, it is empirically known that changes in mood are expressed in a person's voice, and in a previous study, the authors developed a method to estimate mental health states, such as depression and stress, using a person's voice [10-12]. Voice analysis is non-invasive, does not require a specialized device, and can be

performed easily and remotely. Furthermore, it may solve the reporting bias associated with self-administered psychological tests and other various problems encountered while detecting mental diseases. Thus, the stress evaluation method using voice analysis has recently been garnering attention.

The authors had conducted a study on voice index that could detect bipolar disorder and major depressive disorder [13]. In this previous study, we identified healthy individuals and patients with mental disease based on a single classic voice index and showed that, with another single classic voice index, bipolar disorder and major depressive disorder could be identified. However, the detection precision was not sufficient.

Therefore, in this study, we targeted groups of healthy individuals, patients with bipolar disorder, and patients with the major depressive disorder and proposed a new voice evaluation index to identify these three groups with improved precision using polytomous logistic regression analysis [14,15].

METHODS

Subjects

Our subjects were outpatients at the National Defense Medical College Hospital who were undergoing treatment for bipolar disorder and major depressive disorder. They were diagnosed by psychiatrists using the Mini-International Neuropsychiatric Interview [16]. Subjects who had no issues with their daily living were enrolled as healthy volunteers.

There were eight patients with bipolar disorder and 14 patients with major depressive disorder. There were nine and 23 healthy people at the National Defense Medical College Hospital and Tokyo Medical

University Hospital, respectively, with a total of 32 healthy individuals. Depending on their situations, patients received multiple examinations, and their voices were recorded at each examination. In total, there were 14 data of voices with bipolar disorder; 30 data of voices with major depressive disorder, and 32 healthy individuals.

Table 1 shows the details of the data, where values outside the brackets represent the number of patients, and values inside the brackets represent the number of voice data.

The mean age of the healthy group was 50.48±13.45 years (age could not be confirmed for three patients). The mean age of the patients with bipolar disorder was 46.50±13.06 years. The mean age of the patients with major depressive disorder was 43.71±11.57 years.

As other tests, all patients recorded scores of Hamilton Depression Rating Scale [17] and Young Mania Rating Scale [18] at each examination.

Voice recording

Voice recording was performed in examination rooms at each hospital, and patients were asked to read a fixed sentence consisting of 17 phrases. For healthy individuals, voice recording was conducted for the nine individuals at the National Defense Medical College Hospital and for the 23 individuals at the Tokyo Medical University Hospital (six individuals overlapped) in the same environments as that of the patient voice recordings. Voices were recorded with a ME52W microphone (OLYMPUS, Tokyo, Japan) attached on the chest, about 100 mm below the subject’s mouth. The recording device used was the Portable Recorder R-26 (Roland, Shizuoka, Japan). The sampling rate of recording was 96 kHz, and the data resolution was 24 bit.

Voice analysis

After this, healthy individuals are referred to as “HE,” patients with bipolar disorder as “BP,” and patients with major depressive disorder as “MD.”

To extract features from the voices, we used the free software openSMILE (v. 2.3) [19].

The openSMILE uses a script that automatically extracts a set of various features from the voice, and in this study, we extracted feature sets used for emotion recognition (The large openSMILE emotion feature set) from each voice. In this manner, we extracted 6.552 voice features from each voice. From these features, we extracted features that suited the classifier of healthy and patient groups. The procedure was as follows:

1. Eliminating the difference in recording environment at the National Defense Medical College Hospital and the Tokyo Medical University Hospital. We separated healthy individuals into two groups: Those recorded at the National Defense Medical College Hospital (HE_N) and those recorded at Tokyo Medical University Hospital (HE_T), and for each extracted feature, we calculated the effect size between the two groups (HE_N and HE_T). Effect size is an index that evaluates the difference in mean values between two groups, and this difference is defined based on the equations given below, as a value standardized with standard deviation:

$$ES = \frac{|\mu_A - \mu_B|}{\sigma} \text{ where } \sigma = \sqrt{\frac{(n_A - 1)\sigma_A^2 + (n_B - 1)\sigma_B^2}{n_A + n_B - 2}} \tag{1}$$

In these equations, μ_A and μ_B , σ_A and σ_B , and n_A and n_B represent mean values, standard deviations, and the number of samples for Groups A and B, respectively. If the effect size is small, the difference is lower

between the two groups. In this study, we selected features for which the effect size was lower than 0.5.

Furthermore, to eliminate data dependence on effect size, we tested the difference in mean values for each feature between two of the groups. For this, we used a *t*-test if feature distribution of HE_N and HE_T satisfied normalcy and the Mann-Whitney *U*-test if it did not. We also selected features in which the *p* value obtained in the test was larger than 0.05.

Using these processes, we selected 1.553 features from 6.552.

2. We selected features that were effective for identifying HE, BP, and MD. For each one of the 1.553 features, we calculated the effect size between two of the groups - HE-BP, HE-MD, and BP-MD - using the equations from the first step of this procedure. In addition, we performed a multiple comparison test to compare the three groups (HE, BP, and MD). We used the Steel-Dwass test, which does not restrict the distribution shape of the three groups [20].

In this study, we selected features in which the effect size between

all combinations of two groups among the three groups (HE, BP, and MD) was >0.5, and *p* values of paired comparison tests were all smaller than 0.1. Results led to a selection of nine features from 1.553 features.

Polytomous logistic regression analysis

A polytomous logistic regression analysis is a multivariate analysis that classifies data into three or more groups based on predicted values. It is an expansion of a normal logistic regression analysis that classifies data into two groups. Model equations for three groups (A, B, and C) are shown below:

$$\log\left(\frac{P_A}{P_B}\right) = \alpha_A + \beta_{1A}x_1 + \beta_{2A}x_2 + \dots + \beta_{nA}x_n, \tag{2}$$

$$\log\left(\frac{P_B}{P_C}\right) = \alpha_B + \beta_{1B}x_1 + \beta_{2B}x_2 + \dots + \beta_{nB}x_n, \tag{3}$$

In these equations, P_A , P_B , and P_C indicate occurrence probabilities for each group, (x_1, x_2, \dots, x_n) indicates one data set, and $\alpha_A, \beta_{1A}, \dots, \beta_{nA}, \alpha_B, \beta_{1B}, \dots, \beta_{nB}$ indicate model coefficients. To perform polytomous logistic regression analysis, a standard group must be set first; (2) and (3) are model equations that use the Group C as the standard group. Groups other than the standard group are target groups. Logit of the model equation expresses the likelihood of a target group relative to that of the standard group; if the logit is smaller, the standard group is more likely, and if the logit is larger, the target group is more likely. If the coefficient of the model equation is not statistically 0, the variable for such coefficient has an impact on the logit. The occurrence probability for each group based on (2) and (3) can be calculated using the following equations:

$$P_A = \frac{\exp\{\log(P_A / P_C)\}}{1 + \exp\{\log(P_A / P_C)\} + \exp\{\log(P_B / P_C)\}}, \tag{4}$$

$$P_B = \frac{\exp\{\log(P_B / P_C)\}}{1 + \exp\{\log(P_A / P_C)\} + \exp\{\log(P_B / P_C)\}}, \tag{5}$$

$$P_C = 1 - P_A - P_B \tag{6}$$

For each data set, P_A , P_B , and P_C are estimated from model equations, and the data are classified as the most probable group.

In this study, we used categorical information of each group (HE, BP, and MD) as dependent variables and the nine selected voice features as

Table 1: Details of voice data

State	Male	Female	Total
Healthy	19 (19)	13 (13)	32 (32)
Bipolar disorder	2 (3)	6 (11)	8 (14)
Major depressive disorder	12 (27)	2 (3)	14 (30)

independent variables. We performed the analysis using the HE group as the standard group. For statistical analysis, we used the statistical analysis free software R version 3.4.2 [21].

RESULTS

We recorded the effectiveness of the nine voice features that were selected through voice analysis as follows:

1. mfcc_sma1_quartile2. The median of the three-point moving average in mel-frequency sub-band (second).
2. mfcc_sma2_linregc2. The intercept of a regression line of the three-point moving average in mel-frequency sub-band (third).
3. mfcc_sma2_qregc3. The secondary regression curve coefficient (constant term) of the three-point moving average in mel-frequency sub-band (third).
4. mfcc_sma2_quartile1. The first quartile of the three-point moving average in mel-frequency sub-band (third).
5. mfcc_sma2_amean. The arithmetic mean of the three-point moving average in mel-frequency sub-band (third).
6. mfcc_sma3_linregerrA. The primary error of the regression line and the original data of the three-point moving average in mel-frequency sub-band (fourth).
7. mfcc_sma10_quartile1. The first quartile of three-point moving average in mel-frequency sub-band (11th).
8. F0env_sma_qregerrQ. The square error of the secondary regression curve and the original data of the three-point moving average for the F0 envelope (smoothed by the attenuating exponential function).
9. pcm_fftMag_spectralCentroid_sma_linregerrA. The primary error of the regression line and the original data of the three-point moving average for gravity center frequency of spectral power.

Table 2 shows the results of the polytomous logistic regression analysis, which also used stepwise regression where the categorical information of subject groups (HE, BP, and MD) was used as dependent variables, the abovementioned nine voice features were dependent variables, and the HE group was the standard group.

Coefficients in the table represent coefficients to the independent variables from the prediction model. The features ultimately selected by the stepwise regression were the following seven:

- mfcc_sma1_quartile2,
- mfcc_sma2_linregc,
- mfcc_sma2_amean,
- mfcc_sma3_linregerrA,

Table 2: The polytomous logistic regression analysis results

Prediction model for the BP group versus the HE group	Coefficients	SE	p value
Intercept	-0.85	1.12	0.45
mfcc_sma1_quartile2	-1.96	0.88	0.025*
mfcc_sma2_linregc2	4.80	4.35	0.27
mfcc_sma2_amean	-4.43	4.44	0.32
mfcc_sma3_linregerrA	-1.31	0.90	0.14
mfcc_sma10_quartile1	-1.50	0.72	0.037*
F0env_sma_qregerrQ	3.26	1.30	0.012*
pcm_fftMag_spectralCentroid_sma_linregerrA	-1.81	1.13	0.11
Prediction model for the MD group versus the HE group	Coefficients	SE	p value
Intercept	1.03	0.79	0.19
mfcc_sma1_quartile2	-1.67	0.69	0.015*
mfcc_sma2_linregc2	-3.69	3.27	0.26
mfcc_sma2_amean	5.11	3.39	0.13
mfcc_sma3_linregerrA	-1.32	0.73	0.072
mfcc_sma10_quartile1	0.14	0.61	0.82
F0env_sma_qregerrQ	2.60	1.26	0.040*
pcm_fftMag_spectralCentroid_sma_linregerrA	-2.04	0.85	0.016*

- mfcc_sma10_quartile1,
- F0env_sma_qregerrQ, and
- pcm_fftMag_spectralCentroid_sma_linregerrA.

The features mfcc_sma2_qregc3 and mfcc_sma2_quartile1 were excluded because they were determined to be features that did not contribute to the dependent variables. Therefore, the coefficient for these two features is 0, and the prediction equation is expressed as follows:

$$\begin{aligned} \chi_1 &= mfcc_sma1_quartile2, \\ \chi_2 &= mfcc_sma2_linregc2, \\ \chi_3 &= mfcc_sma2_amean, \\ \chi_4 &= mfcc_sma3_linregerrA, \\ \chi_5 &= mfcc_sma10_quartile1, \\ \chi_6 &= F0env_sma_qregerrQ, \\ \chi_7 &= pcm_fftMag_spectralCentroid_sma_linregerrA, \end{aligned} \tag{7}$$

$$\log\left(\frac{P_{BP}}{P_{HE}}\right) = -0.85 - 1.96x_1 + 4.80x_2 - 4.43x_3 - 1.31x_4 - 1.50x_5 + 3.26x_6 - 1.81x_7, \tag{8}$$

$$\log\left(\frac{P_{MD}}{P_{HE}}\right) = 1.03 - 1.67x_1 - 3.69x_2 + 5.11x_3 - 1.32x_4 + 0.14x_5 + 2.60x_6 - 2.04x_7. \tag{9}$$

In these equations, P_{HE} , P_{BP} , and P_{MD} express the probability of the data being classified as the HE group, BP group, and MD group, respectively.

To determine if the prediction model was significant, we performed a likelihood ratio test with a model using only the intercept. The results had a p value of $p=1.22e^{-14}<0.01$, which confirmed the significance of the prediction model.

To eliminate sex dependence in features useful to the model, we tested the difference between two groups based on the sex of healthy individuals. If the two groups satisfied normalcy, we performed a *t*-test, and if they did not, we performed a Mann-Whitney *U*-test. Results are shown in Table 3.

We calculated the probabilities of being classified into each of the three groups for each data used in the analysis using the prediction model equations (P_{HE} , P_{BP} , and P_{MD}), and by classifying the data in the most probable group, we organized the data. Results are shown in Table 4.

Table 3: Results of sex differences in the healthy group

Features	p value of the HE group sex difference test
mfcc_sma1_quartile2	0.14
mfcc_sma2_linregc2	0.66
mfcc_sma2_amean	0.54
mfcc_sma3_linregerrA	0.45
mfcc_sma10_quartile1	0.74
F0env_sma_qregerrQ	0.68
pcm_fftMag_spectralCentroid_sma_linregerrA	0.43

Table 4: Data identification results

Actual group/predicted group	HE	BP	MD
HE	29	1	2
BP	1	12	1
MD	1	1	28

The HE group was identified with 90.63% accuracy. The BP group was identified with 85.71% accuracy. The MD group was identified with 93.33% accuracy. Overall accuracy was 90.79%.

Fig. 1 shows probability distribution of each subject being identified in each group.

Fig. 1 shows that the probability of subjects in the HE group being classified as the HE group was much higher than the probability of being classified into either of the other two groups. Although the probability of subjects in the BP group being classified as the BP group has a wider range, it was still higher than the probability of being classified into either of the other two groups. The probability of subjects in the MD group being classified as the MD group was much higher than the probability of being classified into either of the other two groups.

DISCUSSION

Within the present analytical data, there were data recorded at two different locations for the healthy group. The analysis using only voice data from the National Defense Medical College Hospital led to a bias in the number of male and female subjects in the healthy group, and as the number of samples was small, it was difficult to obtain a model with sufficient precision. Therefore, to eliminate the sex bias and acquire a sufficient number of samples, we added voice data from the Tokyo Medical University Hospital. This addition required that we eliminate the impact of different recording locations on the model, but in the first step of feature selection, such impact was mostly eliminated.

The predicted model equation had a significant coefficient that was 43% of the whole (excluding the intercept), which is a relatively low percentage. However, the overall model had statistical significance, and the data fit was good. In the present analysis, the feature selection condition was somewhat relaxed, which might have been the reason for the coefficient not being significant. In the future, we will apply the prediction model equation obtained from the analysis to another data set and verify the impact of the coefficient on classification precision.

For the seven features useful for the model, the difference between sexes could not be detected for healthy individuals in the present number of samples. In other words, the way in which patients are classified into two groups based on the seven features is not influenced by differences in sexes.

Using the prediction model equation based on the BP group data, one datum was classified as belonging to the major depressive disorder group. This datum was recorded during the second examination of the patient, and this patient's Hamilton Depression Scale score at that time was higher than it was at the first examination. The same patient's Young Mania Rating Scale score was lower than it was at the first examination. Thus, as the patient's symptoms of major depressive

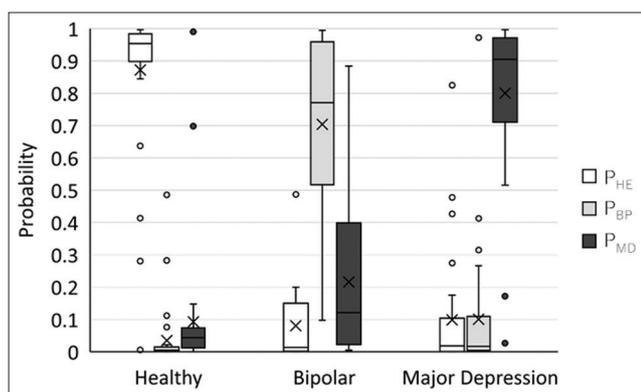


Fig. 1: Probability distribution of each subject being identified in each group

disorder became more prominent than those of bipolar disorder, the patient was classified as belonging to the major depressive disorder group. There was another datum in the MD group that was classified as healthy using the prediction model equation. This datum was recorded during the second examination of the patient, and it is possible that, after outpatient treatment, the patient's depressive symptoms improved, and the patient could be classified as healthy. It is our future challenge to verify how patients are classified when symptoms change with time.

In this study, we did not discuss the impact of the voice features used in the prediction model on the model itself. It is another future challenge for us to examine the features that can most effectively identify the diseases and the characteristics of patients' voices that each feature captures.

CONCLUSION

In this study, we targeted groups of healthy individuals, patients with bipolar disorder, and patients with major depressive disorder and, based on their voices, extracted a large-scale voice feature set using voice feature extraction software. We selected voice features useful as the classifier. Based on the selected features, we performed the polytomous logistic regression analysis and proposed a voice index that identifies healthy individuals, patients with bipolar disorder, and patients with major depressive disorder. Using the prediction model obtained from the analysis, subjects could be classified into the three groups with 90.79% accuracy. The above results indicated that the proposed index could be useful as a new evaluation index to identify depression.

CONFLICTS OF INTEREST

All authors have none to declare.

REFERENCES

- Okumura Y, Higuchi T. Cost of depression among adults in Japan. *Prim Care Companion CNS Disord* 2011;13:PCC.10m01082.
- Izawa S, Sugaya N, Shirotaki K, Yamada KC, Ogawa N, Ouchi Y, et al. Salivary dehydroepiandrosterone secretion in response to acute psychosocial stress and its correlations with biological and psychological changes. *Biol Psychol* 2008;79:294-8.
- Sekiyama A. Interleukin-18 is Involved in Alteration of Hypothalamic-pituitary-adrenal axis Activity by Stress. *Society of Biological Psychiatry Annual Meeting (San Diego)*; 2007.
- Garcia RG, Valenza G, Tomaz CA, Barbieri R. Instantaneous bispectral analysis of heartbeat dynamics for the assessment of major depression. *Computing Cardiol (Nice)* 2015;42:781-784.
- Acharya UR, Sudarshan VK, Adeli H, Santhosh J, Koh JE, Puthankatti SD, et al. A novel depression diagnosis index using non-linear features in EEG signals. *Eur Neurol* 2015;74:79-83.
- Goldberg DP. *Manual of the General Health Questionnaire*. Windsor: NFER Publishing; 1978.
- Beck AT, Ward CH, Mendelson M, Mock J, Erbaugh J. An inventory for measuring depression. *Arch General Psychiatry* 1961;4:561-71.
- Delgado-Rodriguez M, Llorca J. Bias. *J Epidemiol Commun Health* 2004;58:635-41.
- American Psychiatric Association. *Diagnostic and Statistical Manual of Mental Disorders*. 5th ed. Arlington, VA, US: American Psychiatric Association; 2013.
- Higuchi M, Shinohara S, Nakamura M, Omiya Y, Hagiwara N, Mitsuyoshi S, et al. Study on Depression Evaluation Indicator in the Elderly using Sensibility Technology. *Proceedings of the 3rd International Conference on Information and Communication Technologies for Ageing Well and e-Health (Porto)*; 2017. p. 70-7.
- Shinohara S, Mitsuyoshi S, Nakamura M, Omiya Y, Tsumatori G, Tokuno S. Validity of a Voice-based Evaluation Method for Effectiveness of Behavioural Therapy. *Parvasive Computing Paradigms for Mental Health*, 5th International Conference, MindCare 2015, Milan, Italy, September 24-25, 2015, Revised Selected Papers. Switzerland: Springer International Publishing; 2016. p. 604, 43-51.
- Tokuno S, Mitsuyoshi S, Suzuki G, Tsumatori G. Stress Evaluation by Voice: A Novel Stress Evaluation Technology. *The 9th International*

- Conference on Early Psychosis (Tokyo); 2014. p. 17-9.
13. Nakamura M, Omiya Y, Shinohara S, Mitsuyoshi S, Higuchi M, Hagiwara N, *et al.* Feasibility Study of Classifying Major Depressive Disorder and Bipolar Disorders using Voice Features. WPA XVII World Congress of Psychiatry (Berlin); 2017.
 14. Agresti A. Categorical Data Analysis. 3rd ed. New Jersey: Wiley-Interscience; 2012.
 15. Venables WN, Ripley BD. Modern Applied Statistics with S. 4th ed. New York: Springer; 2002.
 16. Sheehan DV, Lecrubier Y, Sheehan KH, Amorim P, Janavs J, Weiller E, *et al.* The Mini-international neuropsychiatric interview (M.I.N.I): The development and validation of a structured diagnostic psychiatric interview for DSM-IV and ICD-10. J Clin Psychiatry 1998;59 Suppl 20:22-33.
 17. Hamilton M. A rating scale for depression. J Neurol Neurosurg Psychiatry 1960;23:56-62.
 18. Young RC, Biggs JT, Ziegler VE, Meyer DA. A rating scale for mania: Reliability, validity and sensitivity. Br J Psychiatry 1978;133:429-35.
 19. Eyben F, Wöllmer M, Schuller B. openSMILE–The Munich Versatile and Fast Open-Source Audio Feature Extractor, Proceedings of the 18th ACM international conference on Multimedia (Firenze); 2010;1459-1462.
 20. Steel RG. A rank sum test for comparing all pairs of treatments. Technometrics 1960;2:197-207.
 21. R Core Team. R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing. Vienna: R Core Team; 2017. Available from: <https://www.R-project.org/>. [Last accessed on 2018 Sep 18].