

# Detect feature edges for diagnosis of bacterial vaginosis

Jie Li<sup>1</sup>, yaotang Li<sup>Corresp. 1</sup>

<sup>1</sup> School of mathematics and statistics, Yunnan University, kunming, yunnan, China

Corresponding Author: yaotang Li  
Email address: liyaotang@ynu.edu.cn

One of the most common diseases among women of reproductive age is bacterial vaginosis (BV). However, the etiology of BV is still unknown. In this paper, we use a network to model a temporal sample of the vaginal microbiome and study the relationship between the edges of the network and BV. We use the machine learning algorithms decision tree and ReliefF to select the network feature edges that are related to BV and then validate those features using logistic regression and support vector machine. We discover that a few features are required to achieve high BV classification accuracy; logistic regression and support vector machine performs nearly identically under the same feature edges; decision tree feature edges outperform ReliefF feature edges in classification performance, and the feature edges selected by those two algorithms are very different. The feature edges may serve as indicators for personalized diagnosis of BV and may aid in the clarification of a more mechanistic interpretation of its etiology.

## 1 Detect feature edges for diagnosis of bacterial vaginosis

2 Jie Li<sup>1</sup>, YaoTang Li<sup>1</sup>

3 <sup>1</sup> School of mathematics and statistics, Yunnan University, Kunming, Yunan, China

4 Corresponding Author: YaoTang Li

5 YaoTang Li<sup>1</sup>

6 No. 2, Cuihu North Road, Wuhua District, Kunming, Yunan, 650091, China

7 Email address: liyaotang@ynu.edu.cn

8

### 9 Abstract

10 One of the most common diseases among women of reproductive age is bacterial vaginosis  
11 (BV). However, the etiology of BV is still unknown. In this paper, we use a network to model a  
12 temporal sample of the vaginal microbiome and study the relationship between the edges of the  
13 network and BV. We use the machine learning algorithms decision tree and ReliefF to select the  
14 network feature edges that are related to BV and then validate those features using logistic  
15 regression and support vector machine. We discover that a few features are required to achieve  
16 high BV classification accuracy; logistic regression and support vector machine performs nearly  
17 identically under the same feature edges; decision tree feature edges outperform ReliefF feature  
18 edges in classification performance, and the feature edges selected by those two algorithms are  
19 very different. The feature edges may serve as indicators for personalized diagnosis of BV and  
20 may aid in the clarification of a more mechanistic interpretation of its etiology.

21 **Keywords:** Bacterial vaginosis; machine learning; network; feature edges;

22

### 23 Introduction

24 Bacterial vaginosis (BV) has been identified to be an independent risk factor for women's health  
25 (Koumans et al. 2001), including preterm delivery and low infant birth weight, the development of  
26 pelvic inflammatory disease increased susceptibility to HIV infection, and other chronic health  
27 issues (Hay et al., 1994; Ness et al., 2005; Sha et al., 2005a; Atashili et al., 2008; van de Wijgert et al.,  
28 2008; Ma et al. 2012). BV is frequently characterized by changes in the vaginal microbiome;  
29 however, the causes of these changes are unknown (Redelinguys et al.,2020). Historically, BV  
30 has been diagnosed using the Nugent score and/or Amsel's clinical criteria (Nugent et al., 1991;  
31 Amsel et al., 1983). The Nugent score is based on the presence or absence of lactobacilli on the  
32 Gram stain and generates a score ranging from 1 to 10. A score of seven or greater indicates a  
33 positive BV diagnosis. Three of the following four Amsel's criteria yield a positive diagnosis: 1)  
34 the presence of a fishy-like odor, 2) the presence of a white discharge, 3) a vaginal pH of >4.5,  
35 and 4) a minimum of 20% "clue cells" detection. The "gold standard" for BV diagnosis is  
36 Amsel's criteria and the Nugent scoring system. These methods have the disadvantages of being  
37 difficult to standardize and subject to interobserver variability because the assessment of the  
38 diagnostic criteria is dependent on the observer's skill and experience (Klebanof et al., 2004;  
39 Modak et al., 2011).

40 The recent advancement of molecular and high-throughput sequencing technologies allows for  
41 the detection of a large number of unculturable microorganisms from clinical samples (Adzitey et  
42 al., 2013). As a result, high-throughput biomolecular data are used to track the history of BV or to  
43 identify the pathogens of BV (Srinivasan et al. 2010; Ravel et al. 2011, 2013, White et al. 2011, Gajer  
44 et al. 2012, Hickey et al. 2012, Ma et al. 2012, Romero et al. 2014; Doyle et al. 2018). Ravel et al.  
45 (2013), for example, report on the temporal dynamics of 25 vaginal communities over 10 weeks  
46 using daily samples collected from women who were diagnosed with symptomatic BV,  
47 asymptomatic BV, and healthy. Srinivasan et al. (2010) conducted deep sequencing of the 16S  
48 rRNA gene in an attempt to investigate the variety and composition of vaginal bacteria in BV-  
49 positive women.

50

51 In the meantime, machine learning techniques have been used in this field. Baker et al. (2014)  
52 used genetic programming, random forests, and logistic regression machine learning methods on  
53 two BV datasets in the hopes of discovering BV-related microbial relationships. Later, Beck and  
54 Foster (2015) used random forests and logistic regression, in conjunction with ReliefF, to  
55 diagnose BV. *Aerococcus*, *Atopobium*, *Dialister*, *Eggerthella*, and *Gardnerella* were identified as  
56 the most important bacteria associated with BV in their findings. Pérez-Gómez et al (2020) used  
57 a decision tree and the ReliefF algorithm as feature selectors, as well as the support vector  
58 machine and the logistic regression algorithm as classifiers to identify bacteria associated with  
59 BV. Loquet et al. (2021) designed classification and regression trees for BV to diagnosis in pregnant  
60 women. These works fall into the category of discovering BV-related feature bacteria (or OTU,  
61 Operational taxonomic unit).

62

63 Existing research indicates that BV is a systemic abnormality caused by multiple bacteria and  
64 that interactions between bacteria also play a role in the onset of BV (Srinivasan et al. 2010; White  
65 et al. 2011; Ravel et al., 2011, 2013; Gajer et al. 2012; Romero et al. 2014; Doyle et al. 2018). As a  
66 result, studying bacterial interactions is required to gain insight into BV signaling pathways. The  
67 bacterium can be defined as a network node, and interactions between bacteria can be defined as  
68 network edges. The challenge now is to identify network edges (interactions between two  
69 bacteria) that can characterize the state of the vaginal microbiome. Efforts to find reliable feature  
70 edges rely on information about bacterial interactions, so temporal sample datasets are required.  
71 The dataset reported by Ravel et al. (2013) provides ideal material to investigate this topic. In  
72 this paper, we model each temporal dataset of the vaginal community from Ravel et al. to a  
73 network, and then we build 25 networks. We apply supervised machine learning methods to 25  
74 networks to find feature edges that are related to BV. We hope that these feature edges will aid in  
75 the diagnosis of BV and promote research into the pathogenesis of BV.

76

## 77 **Materials & Methods**

### 78 Vaginal Microbiome Dataset

79 The dataset was originally reported by Ravel et al. (2013). Ravel et al. (2013) sequenced vaginal  
80 communities collected daily for ten weeks from 25 women diagnosed with symptomatic BV  
81 (SBV:  $n = 15$  women), asymptomatic BV (ABV:  $n = 6$ ), or healthy (HEA:  $n = 4$ ). In total, Ravel  
82 et al. (2013) sequenced 1,657 samples (median = 67 per woman) and obtained 420 8,757,681  
83 high-quality sequenced reads of the V1–V3 hypervariable region of 16S-rRNA genes, with a  
84 median of 5,093 reads per sample. The dataset is freely accessible to the public (Ravel et al.,  
85 2013).

### 86 Feature Selection Algorithms

87 Feature selection aims to find the optimal subset of features. Feature selection can be used to  
88 eliminate irrelevant or redundant features, reduce the number of features, filter out features  
89 related to class information, and improve model accuracy. The general process of feature  
90 selection:

91 1. Generate subsets: search for feature subsets and provide feature subsets for the evaluation  
92 function;

93 2. Evaluation function: evaluate the quality of the feature subset;

94 3. Stopping criteria: related to the evaluation function, generally a threshold, the search can be  
95 stopped after the evaluation function reaches a certain standard;

96 4. Verification process: verify the validity of the selected feature subset on the verification data  
97 set.

98 Decision tree (Bramer 2007) and ReliefF (Robnik-Šikonja et al., 2003) are used in this work,  
99 they are belonging to the surprised feature selection method. These methods are implemented  
100 function by function in the Python modules skfeature (Li et al, 2018) and sklearn.

101

### 102 Classification Algorithms

103 A classification algorithm has two phases: learning and classification. The classification model is  
104 trained on the given dataset and its label information during the learning phase; during the  
105 classification phase, the classification model assigns the label to the new dataset. The  
106 classification model in this paper uses logistic regression (Han et al., 2011) and support vector  
107 machine (SVM, Wang et al., 2018), both of which are classic binary classification models that are  
108 widely used in a variety of fields.

109 Given a training dataset of feature space  $T = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$ , where  $x_i \in R^n$ ,

110  $y_i \in \{+1, -1\}$ ,  $i = 1, 2, \dots, N$ ,  $x_i$  is  $i$ th feature vector,  $y_i$  is the class label. For a given dataset  $T$  and

111 hyperplane  $w \cdot x + b = 0$ , the distance between the sample point and hyperplane can be defined as

112  $\gamma_i = y_i \left( \frac{w}{\|w\|} \cdot x_i + \frac{b}{\|w\|} \right)$ . The minimum value of the geometric interval of the hyperplane with

113 respect to all sample points is  $\gamma = \min_i \gamma_i$ . According to the above definition, SVM can be

114 represented as

$$115 \quad \max_{w,b} \gamma \quad s.t. \quad y_i \left( \frac{w}{\|w\|} \cdot x_i + \frac{b}{\|w\|} \right) \geq \gamma, i = 1, 2, \dots, N.$$

116

117 In logistic regression, for a given dataset  $T$ , the aim of the algorithm is still to find the decision  
118 boundary  $w \cdot x + b = 0$  (named hyperplane in SVM). Based on the likelihood theory in statistics,  
119 the optimization model of the model is

$$120 \quad L(w) = \prod [p(x_i)]^{y_i} [1 - p(x_i)]^{1-y_i}, \text{ where } p(x) = \frac{1}{1 + e^{-(w^T x + b)}}.$$

121 Leave-One-Out Validation

122 The dataset will divide into the training set and validation set. The training set is used to train the  
123 model, while the validation set is used to assess the model's generalizability. If the size of the  
124 dataset  $D$  is  $N$ , then use  $N - 1$  pieces of data for training, and use the remaining piece of data as  
125 validation. A total of times  $N$  are calculated for each group taken from  $D$  as the verification set  
126 until all samples have been verified as the set, and finally the verification error is averaged. This  
127 method is called leave-one-out cross-validation (Torgo 2010).

128 Performance Measures

129 Assume there are only two categories (positive and negative, usually the class of interest is the  
130 positive class, and the other classes are the negative class). The confusion matrix is as follows  
131 (table 1):

132 Accuracy: the ratio of correctly classified samples to total samples, is calculated as  $\frac{TP + TN}{P + N}$ .

133 Precision: the ratio of the number of true positive cases to the number of positive cases judged as  
134 positive, is calculated as  $\frac{TP}{TP + FP}$ .

135 Recall: the ratio of the number of positive cases correctly determined to the total number of  
136 positive cases, is calculated as  $\frac{TP}{P}$ .

## 137 Experimental Studies

138 We construct the networks and perform machine learning algorithms to find BV related feature  
139 edges.

140 1. The OTU in each table is taken as the network node for the OTU table of 25 vaginal  
141 microbiomes; the Spearman's rank correlation coefficients are calculated as the weight of edge  
142 between the OTUs, then 25 networks  $A_i, i = 1, 2, \dots, 25$  are obtained. Each network assigned labels  
143 SBV, ABV, and H according to the diagnosis of the corresponding women.

144 2. Divided the 25 networks into four groups according to the research intention (BV =  
145 ABV+SBV vs. H; SBV vs. H; ABV vs. H; SBV vs. ABV)

146 3. To find feature edges in each group, use a feature selection and classification algorithm. The  
147 specific procedure is as follows: the significance of each edge is scored using a feature selection  
148 algorithm under leave-one-out cross-validation, and the scores are recorded in each run. After  
149 leave-one-out cross-validation, the importance scores of each edge are averaged. Edges are  
150 sorted in descending order by mean importance score. Again, using cross-validation, according  
151 to the mean value of the edge's importance score, a select subset of edges as feature edges to  
152 train the classification model on the training set and classify on the prediction set. The indicators  
153 (accuracy, precision, and recall) are used to evaluate classification performance after cross  
154 validation. The process is depicted in the diagram below.

155

## 156 Results

157

158 We performed the results of the four groupings, as shown in Table 2.

159 From the calculation results, we get the following conclusions.

160 1. Machine learning can distinguish different vaginal microbiome states (BV, ABV, SBV, H)  
161 based on bacterial interaction. It captures the difference between BV, SBV, ABV, and H better  
162 than that between SBV and ABV is weak.

163 2. Selecting the top 5 feature edges of importance can achieve the best accuracy for the feature  
164 selection and classification model. In some cases, the increase of the number of feature edges  
165 will reduce the performance of the classification algorithm.

166 3. The feature edges selected by decision tree outperform those selected by ReliefF in terms of  
167 classification algorithm logistic regression and SVM performance; however, there is almost no  
168 difference between classification algorithm logistic regression and SVM on the same feature  
169 edges.

170 4. The two feature selection algorithms have great differences in the importance of ranking of  
171 edges. Using the top 5 edge set as an example, the feature edges chosen by the two algorithms  
172 have almost no intersection.

173

## 174 Discussion

175 The feature edges that we discovered can distinguish the state of the vaginal microbiome (BV vs.  
176 H; SBV vs. H; ABV vs. H); however, the ability to distinguish between SBV and ABV is  
177 limited. In conclusion, our results show that there are differences in the expression of feature  
178 edges (interaction between the bacteria) under different vaginal environmental conditions. As a  
179 result, these feature edges may be useful in the diagnosis of BV. The feature edges chosen by  
180 different feature selection algorithms are inconsistent, a problem that has also been observed in  
181 previous studies (Baker et al., 2014; Beck and Foster, 2015). This adds to the complexity of the  
182 interpretability of feature edges. Similarly, Ma et al., (2021a) found 15 different types of network  
183 markers (motif, interactions among three species) that were present only in the BV microbiome  
184 and absent in the healthy microbiome, and which were validated on other BV datasets. We take  
185 the result of the decision tree algorithm to compare with the result of Ma et al. (2021a). We  
186 found that there was no overlap between them. It implies that the identification of BV associate  
187 feature edges may not be unique and that finding universal feature edges is difficult and  
188 complex, necessitating the mining of more sample data.

189 Further insights can be shed on the ecological mechanisms of BV by distinguishing key bacteria,  
190 or by the identifications of the critical pathway of interactions. However, BV is still poorly  
191 understood. However, the BV “single causative agent” theory is no longer widely accepted.  
192 Alternatively, BV is thought to be polymicrobial in nature. There are evidences that interspecies  
193 interactions characterize the vaginal microbiota with BV. *Gardnerella* spp. may provide a  
194 favorable environment for the growth of other BV-associated bacteria during the onset of BV,  
195 according to Pybus and Onderdonk (1997). Srinivasan and Fredricks (2008) proposed that BV  
196 occurs when BV-associated bacteria enter the vagina and displace lactobacilli. Furthermore, BV-  
197 associated bacteria (*Bacteroides* spp., *Enterococcus faecalis*, Vaginal G., *Mobiluncus* spp., and  
198 *Peptococcus* spp.) can inhibit *Lactobacillus* growth. And in a healthy vaginal environment,  
199 *Lactobacillus* species produce hydrogen peroxide (H<sub>2</sub>O<sub>2</sub>) to inhibit the overgrowth of anaerobic  
200 bacteria. The reduction of *Lactobacillus* spp. was therefore considered to indicate vaginal  
201 dysbiosis. Those arguments imply, logically, that interactions between certain bacteria are related  
202 to BV. Feature edges (interaction between bacteria) have the potential to reveal the dysbiosis  
203 pathway and signaling associated with BV. However, several risk factors have been identified in  
204 the pathogenesis of BV, such as age, socio-economic status, antibiotic usage, sexual behavior,  
205 and ethnicity (Brumley, 2012; Singh et al., 2015; Ranjit et al., 2018). As a result, while the road to  
206 discovering the full face of BV remains long, our research provides important candidate  
207 materials (feature edges) and tools to further our understanding of BV risk and etiology.

## 208 Conclusion

209 The feature edges discovered by the machine learning algorithm can accurately distinguish BV  
210 and the health status of the vaginal microbiome. These features can also help reveal the  
211 pathogenesis of BV. However, different machine learning algorithms find different feature

212 edges, which increases the complexity of feature interpretation. Furthermore, the data set used in  
213 the study is insufficient, and the sample size is unbalanced. Because only the Spearman  
214 correlation coefficient is used when building the sample network, more work is required. In the  
215 future, we will also try to use different correlation measures to build a sample network, collect  
216 more data, and consider sample balance for research, in the hopes of obtaining more reliable  
217 results and promoting BV diagnosis and pathogenesis.

218

## 219 **Acknowledgements**

220 The authors thank Ravel et al. for their data support to this study.

## 221 **References**

- 222 Koumans, EH, Kendrick, JS. 2001. Preventing adverse sequelae of bacterial vaginosis: a public health  
223 program and research agenda. *Sexually Transmitted Diseases*, 28(5):292–297.
- 224 Hay, P. E., Lamont, R. F., Taylor-Robinson, D., Morgan, D. J., Ison, C., and Pearson, J. 1994. Abnormal  
225 bacterial colonisation of the genital tract and subsequent preterm delivery and late miscarriage. *BMJ* 308,  
226 295–298.
- 227 Ness, R. B., Kip, K. E., Hillier, S. L., Soper, D. E., Stamm, C. A., Sweet, R. L., et al. 2005. A cluster  
228 analysis of bacterial vaginosis-associated microflora and pelvic inflammatory disease. *Am. J. Epidemiol.*  
229 162, 585–590.
- 230 Ma, B, LJ Forney, J. Ravel. 2012. The Vaginal Microbiome: Rethinking Health and Disease. *Annual*  
231 *Review of Microbiology*. 66:371-389.
- 232 Sha, B. E., Zariffard, M. R., Wang, Q. J., Chen, H. Y., Bremer, J., Cohen, M. H., et al. 2005a. Female  
233 genital-tract HIV load correlates inversely with *Lactobacillus* species but positively with bacterial  
234 vaginosis and *Mycoplasma hominis*. *J. Infect. Dis.* 191, 25–32.
- 235 Atashili, J., Poole, C., Ndumbe, P. M., Adimora, A. A., and Smith, J. S. 2008. Bacterial vaginosis and  
236 HIV acquisition: a meta-analysis of published studies. *AIDS* 22, 1493–1501.
- 237 van de Wijgert, J. H. H. M., Morrison, C. S., Cornelisse, P. G., Munjoma, M., Moncada, J., Awio, P., et  
238 al. 2008. Bacterial vaginosis and vaginal yeast, but not vaginal cleansing, increase HIV-1 acquisition in  
239 African women. *J. Acquir. Immune Defic. Syndr.* 48, 203–210.
- 240 Redelinghuys MJ, Geldenhuys J, Jung H, Kock MM. Bacterial Vaginosis: Current Diagnostic Avenues  
241 and Future Opportunities. *Front Cell Infect Microbiol.* 2020 Aug 11; 10:354.
- 242 Nugent, R. P., Krohn, M. A., and Hillier, S. L. 1991. Reliability of diagnosing bacterial vaginosis is  
243 improved by a standardized method of gram stain interpretation. *J. Clin. Microbiol.* 29, 297–301.
- 244 Amsel, R., Totten, P. A., Spiegel, C. A., Chen, K. C. S., Eschenbach, D., and Holmes, K. K. 1983.  
245 Nonspecific vaginitis: diagnostic criteria and microbial and epidemiologic associations. *Am. J. Med.* 74,  
246 14–22.
- 247 Klebanoff, M. A., Schwebke, J. R., Zhang, J., Nansel, T. R., Yu, K. F., and Andrews, W. W. 2004.  
248 Vulvovaginal symptoms in women with bacterial vaginosis. *Obstet. Gynecol.* 104, 267–272.
- 249 Modak, T., Arora, P., Agnes, C., Ray, R., Goswami, S., Ghosh, P., et al. 2011. Diagnosis of bacterial  
250 vaginosis in cases of abnormal vaginal discharge: comparison of clinical and microbiological criteria. *J.*  
251 *Infection Dev. Countries* 5, 353–360.
- 252 Adzitey, F., Huda, N., and Ali, G. R. R. 2013. Molecular techniques for detecting and typing of bacteria,  
253 advantages and application to foodborne pathogens isolated from ducks. *Biotechnology* 3, 97–107.

- 254 Srinivasan S, Liu C, Mitchell CM, Fiedler TL, Thomas KK, Agnew KJ, et al. 2010. Temporal variability  
255 of human vaginal bacteria and relationship with bacterial vaginosis. 2010, PloS One 2010, 5:e10197.
- 256 Ravel J., Gajer, P., Abdo, Z., Schneider, G. M., Koenig, S. S. K., McCulle, S. L., Karlebach, S., et al.  
257 2011. Vaginal microbiome of reproductive-age women. Proceedings of the US National Academy of  
258 Sciences, 108 Suppl. 1, 4680–4687.
- 259 Ravel, J. Jacques Ravel, Rebecca M Brotman, Pawel Gajer, Bing Ma, Melissa Nandy, Douglas W  
260 Fadrosch, Joyce Sakamoto, Sara SK Koenig, Li Fu, Xia Zhou, Roxana J Hickey, Jane R Schwebke, Larry J  
261 Forney. 2013. Daily temporal dynamics of vaginal microbiota before, during and after episodes of  
262 bacterial vaginosis. Microbiome 2013, 1:29.
- 263 White, BA, D. J. Creedon, K. E. Nelson, B. A. Wilson. 2011. The vaginal microbiome in health and  
264 disease. Trends in Endocrinology and Metabolism, 2011, 22(1): 389-393.
- 265 Gajer, P., R. M. Brotman, G. Bail, J. Sakamoto, U. M. E. Schütte, X. Zhong, S. S. Koeng, L. Fu, Z. Ma,  
266 X. Zhou, Z. Abdo, L. J. Forney and J. Ravel. 2012. Temporal Dynamics of the Human Vaginal  
267 Microbiota. Science Translational Medicine, 4(132):132ra52.
- 268 Hickey RJ, Zhou X, Pierson JD, Ravel J, Forney LJ. 2012. Understanding vaginal microbiome  
269 complexity from an ecological perspective. Transl. Res, 2012, 160:267-282.
- 270 Romero, R, SS Hassan, P. Gajer, AL Tarca, DW Fadrosch, L. Nikita, M. Galuppi, RF Lamont, P.  
271 Chaemsaitong, J. Miranda, T. Chaiworapongsa and J. Ravel. 2014. The composition and stability of the  
272 vaginal microbiota of normal pregnant women is different from that of non-pregnant women.  
273 Microbiome, 2014, 2:4.
- 274 Doyle R, Gondwe A, Fan Y, et al. 2018. A Lactobacillus-Deficient Vaginal Microbiota Dominates  
275 Postpartum Women in Rural Malawi. Applied and Environmental Microbiology, 84(6): e02150-17.
- 276 Baker, Y.S.; Beck, D.; Agrawal, R.; Dozier, G.; Foster, J.A. Detecting Bacterial Vaginosis using machine  
277 learning. In Proceedings of the 2014 ACM Southeast Regional Conference, Kennesaw, GA, USA, 28–29  
278 March 2014.
- 279 Beck, D.; Foster, J.A. Machine learning classifiers provide insight into the relationship between microbial  
280 communities and bacterial vaginosis. BioData Min. 2015, 8, 23.
- 281 Pérez-Gómez, J.F., Canul-Reich, J., Hernández-Torruco, J., & Hernández-Ocaña, B. 2020. Predictor  
282 Selection for Bacterial Vaginosis Diagnosis Using Decision Tree and Relief Algorithms. Applied  
283 Sciences.
- 284 Loquet A, Le Guern R, Grandjean T, Duployez C, Bauduin M, Kipnis E, Brabant G, Subtil D, Desein R.  
285 Classification and Regression Trees for Bacterial Vaginosis Diagnosis in Pregnant Women Based on  
286 High-Throughput Quantitative PCR. J Mol Diagn. 2021 Feb;23(2):234-241.
- 287 Bramer, M. Principles of Data Mining; Springer: London, UK, 2007; Volume 180.
- 288 Robnik-Šikonja, M.; Kononenko, I. Theoretical and empirical analysis of ReliefF and RReliefF. Mach.  
289 Learn. 2003, 53, 23–69.
- 290 Jundong Li, Kewei Cheng, Suhang Wang, Fred Morstatter, Robert P. Trevino, Jiliang Tang, and Huan  
291 Liu. 2017. Feature Selection: A Data Perspective. ACM Comput. Surv. 50, 6, Article 94 (November  
292 2018), 45 pages.
- 293 Han, J.; Pei, J.; Kamber, M. Data Mining: Concepts and Techniques, 3rd. ed.; University of Illinois at  
294 Urbana-Champaign: Champaign, IL, USA; Simon Fraser University: Burnaby, BC, Canada; Elsevier:  
295 Amsterdam, The Netherlands, 2011; ISBN 978-0-12-381479-1.
- 296 Wang, H.; Zheng, B.; Yoon, S.W.; Ko, H.S. A support vector machine-based ensemble algorithm for  
297 breast cancer diagnosis. EJOR 2018, 267, 687–699.

298 Torgo, L. *Data Mining with R: Learning with Case Studies*; eBook; Chapman and Hall/CRC: New York,  
299 NY, USA, 2010; ISBN 9780429292859.

300 Ma, Z. S., and Ellison, A. M. (2021a). In silico trio-biomarkers for bacterial vaginosis revealed by species  
301 dominance network analysis. *Comput. Struct. Biotechnol. J.* 19, 2979–2989.

302 Brumley, J. (2012). *Testing a model of bacterial vaginosis among black women.* (PhD Dissertation).  
303 University of South Florida. Available online at: <https://scholarcommons.usf.edu/etd/3995> (accessed 16  
304 December 2019).

305 Singh, H. O., Singh, A., Dhole, T. N., and Sumitra, N. (2015). Factor associated to bacterial vaginosis in  
306 non-pregnant women of North Indian population. *J. Biotechnol. Biomater.* 05:195.

307 Ranjit, E., Raghubanshi, B. R., Maskey, S., and Parajuli, P. (2018). Prevalence of bacterial vaginosis and  
308 its association with risk factors among nonpregnant women: a hospital-based study. *Int. J. Microbiol.*  
309 2018:8349601.

**Table 1** (on next page)

The indicators to evaluate the classification performance.

\* The AB mode in the table: the first indicates whether the predicted result was correct or incorrect, and the second indicates the predicted category. For example, TP means True Positive, that is, the correct prediction is a positive class; FN means, False Negative, that is, the wrong prediction is a negative class.

1

---

Actual category	Prediction category		
	Positive	Negative	Summarize
Positive	$TP$	$FN$	$P$ (Actually positive)
Negative	$FP$	$TN$	$N$ (Actually negative)

---

2

**Table 2** (on next page)

Tables corresponding to calculation results.

\* The results of feature selection only list the top 5 feature edges of importance.

1

Groups	Feature selection	Performance of classifiers
BV vs. H	Table 3	Table 4
SBV vs. H	Table 5	Table 6
ABV vs. H	Table 7	Table 8
SBV vs. ABV	Table 9	Table 10

2

**Table 3** (on next page)

Top 5 feature edges of importance for BV vs. H group selection by Decision tree and ReliefF feature selection algorithms.

1

Decision tree		Relieff	
Feature edges	Import value	Feature edges	Import value
Streptococcus anginosus Veillonellaceae	0.38	Lactobacillus iners Lactobacillus crispatus	1429.8
Actinomycetales Prevotella buccalis	0.228	Atopobium vaginae Megasp3aera	1340.52
Peptonip3ilus Stap3ylococcus	0.04	Atopobium vaginae Stap3ylococcus aureus	1311.08
Streptococcus anginosus Prevotella buccalis	0.04	Lactobacillus iners Atopobium vaginae	1273.84
Atopobium vaginae Parvimonas micra	0.03	Clostridiales Prevotella	1251.24

2

**Table 4**(on next page)

The performance of two classification algorithms on different quantities of features for BV vs. H.

1

Feature Selection Classifiers	Decision tree/ReliefF					
	Logistic regression			SVM		
Feature number	Acc	Pre	Recall	Acc	Pre	Recall
5	0.92/0.84	0.95/0.84	0.95/1	0.92/0.84	0.95/0.84	0.95/1
10	1.00/0.84	1.00/0.84	1.00/1	1.00/0.84	1.00/0.84	1.00/1
15	1.00/0.84	1.00/0.84	1.00/1	1.00/0.84	1.00/0.84	1.00/1
20	1.00/0.84	1.00/0.84	1.00/1	1.00/0.84	1.00/0.84	1.00/1
25	0.96/0.84	0.95/0.84	1.00/1	0.96/0.84	0.95/0.84	1.00/1
30	0.96/0.84	0.95/0.84	1.00/1	0.96/0.84	0.95/0.84	1.00/1
50	0.96/0.84	0.95/0.84	1.00/1	0.96/0.84	0.95/0.84	1.00/1
80	0.84/0.84	0.84/0.84	1.00/1	0.84/0.84	0.84/0.84	1.00/1
100	0.84/0.84	0.84/0.84	1.00/1	0.84/0.84	0.84/0.84	1.00/1
200	0.84/0.84	0.84/0.84	1.00/1	0.84/0.84	0.84/0.84	1.00/1
Mean	0.93/0.84	0.93/0.84	1.00/1	0.93/0.84	0.93/0.84	1.00/1

2

**Table 5** (on next page)

Top 5 feature edges of importance for SBV vs. H group selection by Decision tree / ReliefF feature selection algorithms.

1

Decision tree		Relieff	
Feature edges	Import value	Feature edges	Import value
Lactobacillus iners Bifidobacteriaceae	0.21	Lactobacillus jensenii Streptococcus salivarius	1347.42
Actinomycetales Prevotella buccalis	0.17	Stap3ylococcus Eggert3ella	1315.68
Lactobacillus iners Eggert3ella	0.13	Lactobacillus iners Atopobium vaginae	1314.79
Streptococcus anginosus Veillonellaceae	0.09	Lactobacillus iners BVAB2	1314.26
Megasp3aera sp. type 2 Streptococcus anginosus	0.05	Lactobacillus iners Lactobacillus jensenii	1299.16

2

**Table 6** (on next page)

Top 5 feature edges of importance for ABV vs. H group selection by Decision tree/ReliefF feature selection algorithms.

1

Feature selection Classifiers	Decision tree/ReliefF					
	Logistic regression			SVM		
Feature number	Acc	Pre	Recall	Acc	Pre	Recall
5	0.95/0.95	0.94/0.94	1/1	0.95/0.95	0.94/0.94	1/1
10	0.95/0.89	0.94/0.88	1/1	0.95/0.89	0.94/0.88	1/1
15	0.95/0.89	0.94/0.88	1/1	0.95/0.89	0.94/0.88	1/1
20	0.95/0.89	0.94/0.88	1/1	0.95/0.89	0.94/0.88	1/1
25	0.95/0.89	0.94/0.88	1/1	0.95/0.89	0.94/0.88	1/1
30	0.95/0.84	0.94/0.83	1/1	0.95/0.84	0.94/0.83	1/1
50	0.95/0.79	0.94/0.79	1/1	0.95/0.79	0.94/0.79	1/1
80	0.79/0.79	0.79/0.79	1/1	0.79/0.79	0.79/0.79	1/1
100	0.79/0.79	0.79/0.79	1/1	0.79/0.79	0.79/0.79	1/1
200	0.79/0.79	0.79/0.79	1/1	0.79/0.79	0.79/0.79	1/1
Mean	0.90/0.85	0.90/0.85	1/1	0.90/0.85	0.90/0.85	1/1

2

**Table 7** (on next page)

Top 5 feature edges of importance for ABV vs. H group selection by Decision tree / ReliefF feature selection algorithms.

1

Decision tree		Relieff	
Feature edges	Import value	Feature edges	Import value
Sneat3ia sanguinegens Megasp3aera	0.2	Megasp3aera sp. type 1 Megasp3aera	1431.7
Lactobacillus iners Peptostreptococcus	0.1	BVAB2 Peptonip3ilus asacc3arolyticus	1384
Bacteria Lactobacillus vaginalis	0.1	Lactobacillus crispatus Gammaproteobacteria	1377.4
Dialister sp. type 2 Sneat3ia sanguinegens	0.1	Megasp3aera sp. type 2 Streptococcus	1374.3
BVAB2 Clostridiales	0.1	Atopobium vaginae Prevotella bivia	1352.7

2

**Table 8** (on next page)

The performance of two classification algorithms on different quantities of features for ABV vs. H.

1

Feature Selection Classifiers	Decision tree/RelieFF					
	Logistic regression			SVM		
Feaure number	Acc	Pre	Recall	Acc	Pre	Recall
5	1/0.3	1/0.43	1/0.5	1/0.3	1/0.43	1/0.5
10	1/0.3	1/0.43	1/0.5	1/0.3	1/0.43	1/0.5
15	1/0.5	1/0.56	1/0.83	1/0.5	1/0.56	1/0.83
20	1/0.5	1/0.56	1/0.83	1/0.5	1/0.56	1/0.83
25	1/0.5	1/0.56	1/0.83	1/0.5	1/0.56	1/0.83
30	1/0.6	1/0.6	1/1	1/0.6	1/0.6	1/1
50	1/0.6	1/0.6	1/1	1/0.6	1/0.6	1/1
80	1/0.6	1/0.6	1/1	1/0.6	1/0.6	1/1
100	0.7/0.6	0.67/0.6	1/1	0.7/0.6	0.67/0.6	1/1
200	0.6/0.6	0.6/0.6	1/1	0.6/0.6	0.6/0.6	1/1
Mean	0.93/0.51	0.93/0.55	1/0.85	0.93/0.51	0.93/0.55	1/0.85

2

**Table 9** (on next page)

Top 5 feature edges of importance for SBV vs. ABV group selection by Decision tree/ReliefF feature selection algorithms.

1

Decision tree		ReliefF	
Features	Import value	Features	Import value
Megasp3aera sp. type 2 Stap3ylococcus	0.29	BVAB1 Firmicutes	1246.48
Megasp3aera sp. type 2 Enterococcus faecalis	0.22	Prevotella genogroup 1 Prevotella buccalis	1240.57
Megasp3aera sp. type 2 Actinomycetales	0.19	Gemella Sneat3ia sanguinegens	1240.10
Stap3ylococcus aureus Megasp3aera	0.05	Actinobacteria .class. Clostridiales Family XI. Incertae Sedis	1224.48
Prevotella buccalis Bifidobacterium	0.05	Eggert3ella Prevotella genogroup 3	1219.38

2

**Table 10**(on next page)

The performance of two classification algorithms on different quantities of features for SBV vs. ABV

1

Feature Selection Classifiers	Decision tree/RelieFF					
	Logistic regression			SVM		
Feature number	Acc	Pre	Recall	Acc	Pre	Recall
5	0.86/0.67	0.93/0.70	0.87/0.93	0.86/0.67	0.93/0.70	0.87/0.93
10	0.67/0.67	0.7/0.70	0.93/0.93	0.67/0.67	0.7/0.70	0.93/0.93
15	0.62/0.71	0.68/0.71	0.87/1	0.62/0.71	0.68/0.71	0.87/1
20	0.57/0.71	0.67/0.71	0.8/1	0.57/0.71	0.67/0.71	0.8/1
25	0.57/0.71	0.67/0.71	0.8/1	0.57/0.71	0.67/0.71	0.8/1
30	0.62/0.71	0.68/0.71	0.87/1	0.62/0.71	0.68/0.71	0.87/1
50	0.67/0.71	0.7/0.71	0.93/1	0.67/0.71	0.7/0.71	0.93/1
80	0.67/0.71	0.7/0.71	0.93/1	0.67/0.71	0.7/0.71	0.93/1
100	0.71/0.71	0.71/0.71	1/1	0.71/0.71	0.71/0.71	1/1
200	0.71/0.71	0.71/0.71	1/1	0.71/0.71	0.71/0.71	1/1
Mean	0.67/0.70	0.72/0.71	0.9//0.99	0.67/0.70	0.72/0.71	0.9//0.99

2

# Figure 1

Conceptual diagram of experimental process

