

iSuc-ChiDT: A Computational Method for Identifying Succinylation Sites using Statistical Difference Table Encoding and the Chi-Square Decision Table Classifier

Ying Zeng

Hunan Institute of Engineering <https://orcid.org/0000-0002-9358-8047>

Yuan Chen

Hunan Agricultural University

Zheming Yuan (✉ zhmyuan@sina.com)

Hunan Agricultural University

Research

Keywords: succinylation site, chi-square statistical difference table, ChiDT, imbalanced dataset, feature selection

Posted Date: June 10th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-590597/v1>

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

[Read Full License](#)

Version of Record: A version of this preprint was published at BioData Mining on February 10th, 2022. See the published version at <https://doi.org/10.1186/s13040-022-00290-1>.

1 **iSuc-ChiDT: a computational method for identifying succinylation**
2 **sites using statistical difference table encoding and the chi-square**
3 **decision table classifier**

4 Ying Zeng¹, Yuan Chen², Zheming Yuan^{2*}

5 1. College of Computer and Communication, Hunan Institute of Engineering,
6 Xiangtan 411104, China;

7 2. Hunan Engineering & Technology Research Center for Agricultural Big Data
8 Analysis & Decision-making, Hunan Agricultural University, Changsha 410128,
9 China;

10 ***Corresponding author:** Zheming Yuan (zhmyuan@sina.com)

11

12

13 **Abstract**

14 **Background:** Lysine succinylation is a type of protein post-translational modification which is
15 widely involved in cell differentiation, cell metabolism and other important physiological
16 activities. To study the molecular mechanism of succinylation in depth, succinylation sites need to
17 be accurately identified, and because experimental approaches are costly and time-consuming,
18 there is a great demand for reliable computational methods. Feature extraction is a key step in
19 building succinylation site prediction models, and the development of effective new features
20 improves predictive accuracy. Because the number of false succinylation sites far exceeds that of
21 true sites, traditional classifiers perform poorly, and designing a classifier to effectively handle
22 highly imbalanced datasets has always been a challenge.

23 **Results:** We propose a new computational method, iSuc-ChiDT, to identify succinylation sites in
24 proteins. In iSuc-ChiDT, chi-square statistical difference table encoding is developed to extract
25 positional features, and has the highest predictive accuracy and fewest features compared to binary
26 encoding and physicochemical property encoding. The chi-square decision table (ChiDT)
27 classifier is designed to implement imbalanced pattern classification. With a training set of
28 4748:50,551(true: false sites), independent tests showed that ChiDT significantly outperformed
29 traditional classifiers (including random forest, artificial neural network and relaxed variable
30 kernel density estimator) in predictive accuracy and only taking 17s. Using an independent testing
31 set of experimentally identified succinylation sites, iSuc-ChiDT achieved sensitivity of 70.47%,
32 specificity of 66.27%, Matthews correlation coefficient of 0.205, and a global accuracy index Q^9
33 of 0.683, showing a significant improvement in sensitivity and overall accuracy compared to
34 PSuccE, Success, SuccinSite and other existing succinylation site predictors.

35 **Conclusions:** iSuc-ChiDT shows great promise in predicting succinylation sites and is expected
36 to facilitate further experimental investigation of protein succinylation.

37

38 **Keywords**

39 succinylation site, chi-square statistical difference table, ChiDT, imbalanced dataset,
40 feature selection

41

42 **Background**

43 Protein post-translational modifications (PTMs) regulate cellular physiology and significantly
44 increase protein diversity and complexity. Lysine succinylation is an evolutionarily conserved
45 PTM present in both prokaryotic and eukaryotic cells where a succinyl group is covalently bonded

46 to specific lysine residues by enzymatic or non-enzymatic processes [1, 2]. Succinylation can
47 promote remarkable changes in protein structure and function, and may play a role in many
48 diseases, such as tuberculosis [3], allergic dermatitis [4], and inflammation [5]. Therefore,
49 elucidating the molecular mechanism of succinylation will provide valuable information for both
50 biomedical research and drug development.

51 Accurate identification of succinylation sites is critical to succinylation research, and because
52 experimental methods are costly and time-consuming, and have been unable to keep up with the
53 exponential growth of the number of sequenced proteins, efficient *in silico* methods are in great
54 demand. To date, many predictors for identifying succinylation sites have been developed, such as
55 SucPred [6], SuccinSite [7], pSuc-Lys [8], PSuccE [9], and so on, but with their limited overall
56 accuracy and poor sensitivity, numerous true succinylation sites remain undetected. Actually, what
57 interested us more is the information on true succinylation sites. Therefore, it is necessary to
58 further improve predictive accuracy, especially sensitivity. Two key components, feature
59 extraction and classifier construction, can greatly affect the accuracy of a computational method.

60 Commonly used features include positional features [7, 9-11], sequence composition [7-11],
61 evolutionary information [12-14], and protein secondary structure [13-15]. Positional information
62 of amino acids is basic but important to a protein sequence. While binary encoding [7, 9] is the
63 most intuitive method to extract positional features, the feature matrix is very sparse. The binary
64 encodings are the same for the same residue at different positions, and so it cannot reflect
65 positional differences. Physicochemical property encoding [7, 9, 11] is another position-based
66 amino acid encoding scheme that is frequently used. The AAindex [16] database records 531
67 physicochemical properties of 20 standard amino acids. Since it is not known in advance which

68 physicochemical properties are related to classification, physicochemical property encoding means
69 each position needs to be represented by 531 physicochemical properties, resulting in many
70 irrelevant and redundant features.

71 Traditional classifiers including support vector machine (SVM) [6, 9-11, 13], random forest
72 (RF) [7, 8] and decision tree [12, 15] have been applied in succinylation site prediction. The
73 number of false succinylation sites (non-succinylated lysine residues) far exceeds that of true sites,
74 for example, the dataset from Hasan *et al.* [7] contains 5004/53524 true/false succinylation sites (a
75 ratio of positive to negative samples of about 1:10). Training any traditional classifier with such
76 highly imbalanced datasets could strongly bias classification results [17], and the large number of
77 training samples would make the training time of some classifiers (*e.g.* SVM) unbearable. To
78 address this, some methods (*e.g.* SucPred, SuccinSite) balanced the class distribution by
79 under-sampling the negative samples, but this might lead to the loss of some potential
80 classification information due to the mass discarding of negative samples; some methods (*e.g.*
81 pSuc-Lys, PSuccE) designed classifier ensemble algorithms, however, they were still integrated
82 results of several individual classifiers trained with a balanced subset where positive samples were
83 repeatedly used.

84 With a highly imbalanced dataset, we developed an efficient approach called iSuc-ChiDT for
85 predicting succinylation sites. We used chi-square statistical difference table encoding to extract
86 positional features, then incorporated amino acid composition (AAC) and undirected pair-coupled
87 amino acid composition (undirected-PCAAC) features. After feature selection with the
88 Chi-MIC-share [18] algorithm, the ChiDT classifier was designed to resolve the imbalanced
89 classification problem. The flow chart of our method is shown in Fig. 1.

90 **Methods**

91 **Datasets**

92 From Uni-ProtKB/Swiss-Prot [19] database and NCBI protein sequence database [20], Ning
93 *et al.* [9] obtained 2322 succinylated proteins with 5009 experimentally verified lysine
94 succinylation sites, by applying a 30% homology-reducing screening procedure with CD-HIT [21],
95 and then randomly singled out 124 succinylated proteins to build an independent testing set, and
96 used the remaining 2198 succinylated proteins as a training set. We used the same training and
97 independent testing dataset as in Ning *et al.*, which were freely available via the web link [22].
98 Our training set, namely Tr_data, contains 4748/50,551 true/false succinylation sites; and our
99 testing set, namely Te_data, contains 254/2977 true/false succinylation sites.

100 Each true/false succinylation site was represented by a sequence fragment with an initial
101 length of 51 amino acid residues, where the candidate site (lysine residue) was at the central
102 position 0, and the upstream positions were successively labeled as -1, -2, ..., -25, and the
103 downstream positions labeled 1, 2, ..., 25. If the number of up- or downstream residues of the
104 candidate site was less than 25, amino acids were created through mirror extension to make up the
105 difference [8]. For example, the original sequence of the succinylated protein “SP-P0ABS8” is
106 “MLKNLAKLDQTEM₀DKVNVDLAAAGVAFKE...”. The first lysine (*K*) is the candidate site
107 and therefore the sequence fragment generated by mirror extension is
108 “KFAVGAAALDVNVKDMETQDLKAMLKNLAKLDQTEM₀DKVNVDLAAAGVAFK”. All
109 sequence samples contain only the 20 standard amino acids.

110 **Compressing the 2×20 contingency table of each position with chi-square tests**

111 The maximal information coefficient (MIC) [23] is a novel measure proposed to capture

112 dependences between paired variables. The MIC score ranges from 0 to 1, and only approaches 0
113 if two variables are statistically independent. To calculate the MIC score of the paired variables x
114 and y , the ApproxMaxMI [23] algorithm sets the $n_x \times n_y < B(n)$, where $B(n) = n^{0.6}$ is the maximal
115 grid size restriction, and n is the sample size, and n_x, n_y are the number of partition bins on x and y ,
116 respectively. The MIC score for two independent variables calculated by ApproxMaxMI depends
117 on the ratio between $B(n)$ and n [24], and it is close to 0 only when n approaches infinity. For two
118 independent variables under finite samples (especially for small sample sizes), ApproxMaxMI
119 leads to a large deviation between the calculated MIC score and 0, meaning that the MIC will
120 capture false associations. To address this drawback, Chen *et al.* [25] proposed an improved
121 algorithm, ChiMIC [25], which uses chi-square test to determinate optimal bin size for the
122 calculating of MIC score. For two independent variables with 100 sample points, ApproxMaxMI
123 tends to fall into the maximal grid size ($100^{0.6} \approx 16$), and the corresponding grid partition is a 2×8
124 grid, and the MIC score is 0.24. With ChiMIC, the MIC score is only 0.06, and the corresponding
125 grid partition is a 2×2 or 2×3 grid. This shows that the grid partition searched by ChiMIC is more
126 reasonable and that compressing a 2×8 grid into a 2×2 or 2×3 grid is wise.

127 Similarly, for each position in succinylation site-containing sequences, we can construct a
128 2×20 contingency table by respectively counting the occurrence frequencies of the 20 standard
129 amino acids in the positive and negative samples. For instance, Fig. 2 gives the 2×20 table of
130 position -10 in Tr_data. What we need to investigate is whether the 2×20 table (2×20 grid) is
131 reasonable, and could it be compressed into a 2×10 , or even a 2×2 table? A similar attempt was
132 made in donor splice site prediction. For each position in donor site-containing sequences, a 2×4
133 contingency table can be built by counting the frequencies of 4 bases in the positive and negative

134 samples. Following on from ChiMIC, Zeng *et al.* [26] compressed the 2×4 table of each position
135 into a $2 \times l$ ($2 \leq l \leq 4$) table using chi-square tests, and developed a high-performance approach to
136 predict donor splice sites based on this compression strategy.

137 Encouraged by the successful application of the compression strategy on nucleotide
138 sequences, we applied it to protein sequences. For the 2×20 contingency table for each position in
139 succinylation site-containing sequences, the compression procedure is described below.

140 Step 1: Set the initial value of r (r is an integer) to 20.

141 Step 2: The $2 \times r$ contingency table is compressed by merging two columns corresponding to
142 two different residues, and some $2 \times (r-1)$ contingency tables are obtained, then select a $2 \times (r-1)$
143 contingency table with the maximum chi-square value, denoted as $max_{2 \times (r-1)}$.

144 Step 3: A local 2×2 contingency table is constructed based on the merged residues in $max_{2 \times (r-1)}$
145 and perform a chi-square test. If the p -value is lower than a given threshold, $max_{2 \times (r-1)}$ is
146 unreasonable and will be backtracked to the $2 \times r$ contingency table and the compression procedure
147 is terminated. If the p -value is greater than a given threshold, $max_{2 \times (r-1)}$ is reasonable, and a further
148 compression of $max_{2 \times (r-1)}$ is attempted following these steps: 1) set $r=r-1$; 2) if $r \geq 3$, repeat Step
149 2~3; otherwise, terminate compression.

150 Taking position -10 in Tr_data as an example, its 2×20 contingency table was finally
151 compressed into a 2×3 table (Fig. 2). The 20 original status values of position -10 were therefore
152 turned into 3 status values, *i.e.*, “ARDGTV”, “NQEHILKM” and “CFPSWY”, where, “ARDGTV”
153 indicated that A, R, D, G, T, V at position -10 were regarded as the same status value, and the
154 others were similar.

155 **Key positions selection and window size determination**

156 For each position in the sequences with 51 residues, a $2 \times r$ ($2 \leq r \leq 20$) contingency table can be
157 obtained after compression based on the training set. A chi-square test was then performed on the
158 $2 \times r$ contingency table and the corresponding chi-square value was calculated. Higher chi-square
159 values indicate that the corresponding positions are more important for discriminating positives
160 from negatives. Fig. 3 shows the chi-square values of 50 positions (-25~+25, excluding position 0)
161 in Tr_data, and the chi-square tests of all the positions are significant. We calculate the average of
162 the chi-square values of all the positions, denoted as χ_{ave}^2 , then set χ_{ave}^2 as the threshold to select
163 key positions. The chi-square values of positions -8, -4~-1, 1, 2, 5, 7 are above $\chi_{ave}^2 = 92.797$ (see
164 the red line in Fig. 3), therefore these 9 positions are regarded as the key positions. Furthermore,
165 the contiguous 16 residues (positions -8~+7) are determined as the window size.

166 Positional feature extraction

167 For the 9 key positions in each sequence sample, we extracted 9 positional features, denoted
168 as $P_{-8}, P_{-4}, P_{-3}, P_{-2}, P_{-1}, P_1, P_2, P_5$ and P_7 respectively, where, P_{-8} represents the positional feature
169 of position -8, P_{-4} represents the positional feature of position -4, and so forth. The detailed
170 process is described as follows.

171 In the training set, the occurrence frequencies of the 20 standard amino acids were counted at
172 the i^{th} ($i=1, 2, \dots, 9$) position in the positive and negative samples, and then a 2×20 contingency
173 table was built (Table 1).

174 Table 1 Frequency distribution of amino acids at the i^{th} position

Sample	Amino acid residue						Total
	1	2	...	j	...	20	
Positive	$f_{i,1}^+$	$f_{i,2}^+$...	$f_{i,j}^+$...	$f_{i,20}^+$	f_i^+
Negative	$f_{i,1}^-$	$f_{i,2}^-$...	$f_{i,j}^-$...	$f_{i,20}^-$	f_i^-
Total	$f_{i,1}$	$f_{i,2}$...	$f_{i,j}$...	$f_{i,20}$	N

175 In Table 1, $f_{i,j}^+$ represents the frequency of the j^{th} ($j=1,2,\dots,20$) residue at the i^{th} position in
 176 the positive samples, $f_{i,j}^-$ represents the corresponding frequencies in the negative samples, f_i^+
 177 and f_i^- represent the total number of positive and negative samples, and N represents the total
 178 number of samples. The chi-square value corresponding to the i^{th} position is calculated by:

$$179 \quad \chi^2 = \frac{N^2}{f_i^+ \times f_i^-} \left[\sum_{j=1}^{20} \frac{f_{i,j}^{+2}}{f_{i,j}} - \frac{f_i^{+2}}{N} \right] \quad (1)$$

180 If a new training sample is added, and the j^{th} residue appears at the i^{th} position, first assume
 181 this new training sample is positive, replace $f_{i,j}^+$ with $f_{i,j}^+ + 1$, and calculate a chi-square value
 182 $\chi_{i,j}^{2+}$ using formula (1); then assume this new training sample is negative, replace $f_{i,j}^-$ with
 183 $f_{i,j}^- + 1$, and calculate a chi-square value $\chi_{i,j}^{2-}$ using formula (1). The score for the chi-square
 184 statistical difference table with the j^{th} residue at the i^{th} position is defined as:

$$185 \quad \Delta\chi_{i,j}^2 = \chi_{i,j}^{2+} - \chi_{i,j}^{2-} \quad (2)$$

186 Next, build a 20×9 chi-square statistical difference table (Table 2). Table 2 gives the scores of
 187 the various amino acid residues at each position. If the j^{th} residue appears at the i^{th} position, the i^{th}
 188 positional feature will be assigned a value of $\Delta\chi_{i,j}^2$. Table S1 (Additional file 1) shows the 20×9
 189 chi-square statistical difference table constructed based on 9 key positions in Tr_data.

190

Table 2 20×9 chi-square statistical difference table

Amino acid residue	Position				
	1	...	i	...	9
1	$\Delta\chi_{1,1}^2$...	$\Delta\chi_{i,1}^2$...	$\Delta\chi_{9,1}^2$
...
j	$\Delta\chi_{1,j}^2$...	$\Delta\chi_{i,j}^2$...	$\Delta\chi_{9,j}^2$
...
20	$\Delta\chi_{1,20}^2$...	$\Delta\chi_{i,20}^2$...	$\Delta\chi_{9,20}^2$

191 **Compositional feature extraction**

192 For each sequence sample with a window size of 16 residues, 230 compositional features are
193 extracted, including 20 AAC features and 210 undirected-PCAAC features.

194 The AAC features are defined as the occurrence frequencies of the 20 standard amino acids
195 in the sequence, respectively denoted as f_A, f_R, \dots, f_V , where, f_A represents the frequency of alanine
196 (A), f_R represents the frequency of arginine (R), and so forth.

197 The individual amino acid components are independent of each other, so the AAC features
198 cannot reflect any correlation between amino acids. The pair-coupled amino acid composition [27]
199 (PCAAC) features are composed of the occurrence frequencies of pairwise coupling between two
200 adjacent residues, which can reflect both sequence components and the most preliminary
201 association effect. To reduce feature dimension and solve the sparse problem of feature matrix, we
202 assume that the pairwise coupling has no direction. For example, A-R coupling is treated the same
203 as R-A coupling, and the corresponding pair-coupled component will be expressed by either f_{AR} or
204 f_{RA} , where $f_{AR}(f_{RA})$ is the sum of AR pair occurrence frequency and RA pair occurrence frequency
205 found in a sequence.

206 **Feature selection based on Chi-MIC-share**

207 Not all input features are equally important. Some may not be relevant to prediction, or there
208 may be redundancies, so feature selection is a necessary step in constructing a reliable model.
209 Minimum redundancy maximum relevance (mRMR) [28] is a popular feature selection method.
210 However, relevance measure and redundancy measure in mRMR are not comparable, mRMR only
211 gives the order of feature introduction and it is time-consuming to perform cross-validation in
212 training sets to get the optimal feature subset. To address this, Li *et al.* [18] used ChiMIC as the
213 unified measure of relevance and redundancy, and designed a redundancy sharing (rather than

214 redundancy removing) strategy to propose a novel feature selection method, Chi-MIC-share. We
 215 used the Chi-MIC-share algorithm in this study for feature selection.

216 Given an original feature set $\Omega = \{X_1, X_2, \dots, X_i, \dots, X_n\}$, $|\Omega|$ is the number of elements in Ω ,
 217 and $|\Omega|=n$. If the introduced feature set is represented by S , the complement of S is represented as
 218 $\Omega_S = \Omega - S$. Denoting the response variable as Y , the Chi-MIC-share algorithm is described as
 219 follows.

220 For an introduced feature X_i in S , the score after redundancy sharing is calculated by:

$$221 \quad \text{Chi-MIC-share}(X_i) = \sum_{X_j \in S} \frac{\text{Chi-MIC}(X_i; Y)}{\text{Chi-MIC}(X_i; X_j)} \quad (3)$$

222 The total score of all features in S after redundancy sharing is:

$$223 \quad \text{Chi-MIC-share}(S) = \sum_{X_i \in S} \frac{\text{Chi-MIC}(X_i; Y)}{\sum_{X_j \in S} \text{Chi-MIC}(X_i; X_j)} \quad (4)$$

224 If the next introduced feature is X_{next} , set $E = S + \{X_{next}\}$, then $|E|=|S|+1$. The criterion for
 225 introducing the next optimal feature is:

$$226 \quad \max_{X_{next} \in \Omega_S} [\text{Chi-MIC-share}(E)] = \sum_{X_i \in E} \frac{\text{Chi-MIC}(X_i; Y)}{\sum_{X_j \in E} \text{Chi-MIC}(X_i; X_j)} \quad (5)$$

227 If a new introduced feature no longer makes the total Chi-MIC-share score increase, this
 228 feature will be discarded and feature selection will be automatically terminated. Thus, the criterion
 229 for terminating feature introduction is:

$$230 \quad \text{Chi-MIC-share}(E) \leq \text{Chi-MIC-share}(S) \quad (6)$$

231 Furthermore, feature introduction can be forced to terminate according to the following
 232 criterion:

$$233 \quad \frac{\text{Chi-MIC-share}(E) - \text{Chi-MIC-share}(S)}{\text{Chi-MIC-share}(S)} \leq 0.01 \quad (7)$$

234 **Classifier construction**

235 To efficiently realize the imbalanced pattern classification, a new classifier, ChiDT, was
236 designed as follows.

237 **(1) Compress the $2 \times m$ contingency table of each retained feature**

238 For each feature retained by the Chi-MIC-share feature selection, its $2 \times m$ contingency table
239 (m is the number of original status values of the feature) was compressed according to the
240 previously described procedure to obtain a $2 \times r$ contingency table (r is the number of new status
241 values of the feature, $2 \leq r \leq m$). During the compression process, since the status values of each
242 retained feature are continuous, only adjacent status values could be merged together.

243 **(2) Introduce the retained features one by one**

244 Supposing the proportion of the k^{th} class samples in sample set D is p_k ($k=1, 2$), the
245 information entropy of D is defined as:

246
$$H(D) = -\sum_{k=1}^2 p_k \log_2 p_k \quad (8)$$

247 Given a Chi-MIC-share retained feature X_i , supposing it has r new status values as $\{s_1, s_2, \dots,$
248 $s_j, \dots, s_r\}$ after compressing, then the information gain that X_i brings for D can be calculated by:

249
$$\text{Gain}(D, X_i) = H(D) - \sum_{j=1}^r \frac{|D^j|}{|D|} H(D^j) \quad (9)$$

250 where D^j represents the samples in D whose X_i takes the status value as s_j ($1 \leq j \leq r$), while $H(D^j)$ is
251 the information entropy of D^j .

252 From the features whose information gains are above the average, pick out the one with the
253 highest gain ratio to be the first introduced feature. Here, the gain ratio of X_i is defined as:

254
$$\text{GainRatio}(D, X_i) = \frac{\text{Gain}(D, X_i)}{IV(X_i)} \quad (10)$$

255 where

256
$$IV(X_i) = -\sum_{j=1}^r \frac{|D^j|}{|D|} \log_2 \frac{|D^j|}{|D|} \quad (11)$$

257 and $IV(X_i)$ is the intrinsic value of X_i .

258 Next, the remaining features are introduced one by one with the following steps.

259 Step 1: Under the conditions in which the introduced features have existed, the $2 \times r$
 260 contingency table of each remaining feature is further compressed. If the r columns of the $2 \times r$
 261 contingency table are compressed into one column, the remaining feature cannot be introduced; if
 262 the r columns are not compressed into one column, the remaining feature will be considered as a
 263 candidate feature to be introduced.

264 Step 2: Calculate the information gain of every candidate feature. From the candidate features
 265 whose information gains are above the average, the one with the highest gain ratio is selected to be
 266 the next introduced feature.

267 Step 3: Repeat Step 1~2 until no further features can be introduced.

268 After this, the introduced features with their status values generate various rules. Taking
 269 Tr_data as an example, 10 Chi-MIC-share retained features were finally introduced and 137 rules
 270 were generated (see Additional file 1: Table S2).

271 **(3) Construct a balanced decision table for making decisions**

272 We counted the number of positive and negative training samples conforming to each rule
 273 then constructed a 2×137 imbalanced decision table (Table 3).

274

Table 3 Imbalanced decision table				
Sample	Rule*			Total
	$(P_{-1}=-2.028) \wedge (-0.907 \leq P_2 \leq 0.501) \wedge (-0.715 \leq P_{-8} \leq 0.066)$...	$(P_{-1}=1.839) \wedge (P_1=2.060)$	
Positive	23	...	32	4748
Negative	2907	...	83	50551

275 *For instance, " $(P_{-1}=-2.028) \wedge (-0.907 \leq P_2 \leq 0.501) \wedge (-0.715 \leq P_{-8} \leq 0.066)$ " represents P_{-1} taking a value of -2.028 and
 276 P_2 ranging from -0.907 to 0.501 and P_{-8} ranging from -0.715 to 0.066, where, " \wedge " denotes the logical conjunction.

277 The number of negative samples far exceeds the positives. To resolve the imbalanced
 278 classification problem, based on cost-sensitive learning [29], we adjust the decision weight of
 279 negative samples in each column of the imbalanced decision table, by multiplying the number of
 280 negative samples in each column of Table 3 by θ , where θ is defined as the ratio of the total
 281 number of positive and negative training samples, here, $\theta=4748/50551$. Then, a 2×137 balanced
 282 decision table is obtained (Table 4).

283

Sample	Rule*			Total
	$(P_{-1}=-2.028) \wedge (-0.907 \leq P_2 \leq 0.501) \wedge (-0.715 \leq P_{-8} \leq 0.066)$...	$(P_{-1}=1.839) \wedge (P_1=2.060)$	
Positive	23	...	32	4748
Negative	273.04	...	7.80	4748

284 We can then use the balanced decision table for making decisions. Suppose that a testing
 285 sample meets the rule " $(P_{-1}=-2.028) \wedge (-0.907 \leq P_2 \leq 0.501) \wedge (-0.715 \leq P_{-8} \leq 0.066)$ ". First, we assume
 286 that it is positive and replace 23 with 23+1, then calculate the corresponding chi-square value χ_{i+}^2 .
 287 We then assume that it is negative and replace 273.04 with 273.04+1, then calculate the
 288 corresponding chi-square value χ_{i-}^2 . If $\chi_{i+}^2 > \chi_{i-}^2$, the testing sample is predicted to be positive,
 289 if not, it is negative.

290 Performance evaluation

291 Sensitivity (SN), specificity (SP) and Matthews correlation coefficient (MCC) as the
 292 common indexes for evaluating binary classification are defined as follows:

$$293 \quad SN = \frac{TP}{TP + FN} \quad (12)$$

$$294 \quad SP = \frac{TN}{TN + FP} \quad (13)$$

$$295 \quad MCC = \frac{TP \times TN - FN \times FP}{\sqrt{(TP + FN) \times (TP + FP) \times (TN + FP) \times (TN + FN)}} \quad (14)$$

296 Here, TP, FP, TN, and FN denote the numbers of true positives, false positives, true negatives, and

297 false negatives. MCC is a balanced statistical index that considers SN and SP, but it is sensitive to
 298 class distribution in a testing set. As shown in Table 5, when a prediction model has an SN of 93%
 299 and SP of 95%, as the imbalance degree of the testing set grows, the MCC value declines. This
 300 shows that a low MCC value does not always indicate poor prediction performance as it may be
 301 caused by a highly imbalanced testing set.

302 Table 5 Various evaluation indexes on different ratios of positives to negatives

Positives/Negatives*	SN (%)	SP (%)	MCC	Q^9
100/100	93.00	95.00	0.880	0.939
100/1000	93.00	95.00	0.752	0.939
100/10000	93.00	95.00	0.371	0.939

303 * positive testing sample size/negative testing sample size

304 The content-balancing accuracy index Q^9 [30] is independent of the class distribution of the
 305 dataset and has been widely used to evaluate performance of many prediction programs including
 306 gene-finding, splice site prediction and protein secondary structure prediction [31-33]. As Table 5
 307 shown, the value of Q^9 remains unchanged across different ratios of positives to negatives. In this
 308 study, we introduced Q^9 as the measure of global accuracy to evaluate the prediction performance
 309 of models in case of an imbalanced testing set. Q^9 is defined as:

$$310 \quad Q^9 = (1 + q^9) / 2 \quad (15)$$

311 where

$$312 \quad q^9 = \begin{cases} (TN-FP)/(TN+FP), & \text{if } TP+FN=0 \\ (TP-FN)/(TP+FN), & \text{if } TN+FP=0 \\ 1 - \sqrt{2} \sqrt{[FN/(TP+FN)]^2 + [FP/(TN+FP)]^2}, & \text{if } TP+FN \neq 0 \text{ and } TN+FP \neq 0 \end{cases}$$

313 The value of Q^9 ranges from 0 to 1, and the larger the Q^9 value, the better the prediction
 314 performance.

315 Results and discussion

316 Features retained by Chi-MIC-share

317 Based on Tr_data, the Chi-MIC-share feature selection was performed on 239 original input
 318 features (9 positional features and 230 compositional features). As shown in Fig. 4, when the 37th
 319 feature was introduced, the Chi-MIC-share score peaked (0.12544), after which is began to decline
 320 and feature selection was automatically terminated. To improve computational efficiency, forced
 321 termination criteria were adopted and 10 features were retained (see the red line in Fig. 4). Table 6
 322 describes the retained features in detail. It can be seen that positional features contribute 80% of
 323 all the retained features, indicating that positional features have an important contribution to
 324 succinylation site prediction.

325 Table 6 Features retained by Chi-MIC-share

<i>No.</i>	Retained features	Type	<i>No.</i>	Retained features	Type
1	P_{-1}	position	6	P_7	position
2	P_1	position	7	f_R	AAC
3	P_{-2}	position	8	f_K	AAC
4	P_5	position	9	P_2	position
5	P_{-8}	position	10	P_{-3}	position

326 *No.* denotes the order of feature introduction

327 **Comparison of different classifiers**

328 Based on the same input features (10 features retained by Chi-MIC-share), the ChiDT
 329 classifier was compared to traditional classifiers including random forest (RF), artificial neural
 330 network (ANN) and relaxed variable kernel density estimator (RVKDE) [34]. We choose RVKDE
 331 for the comparison because it delivers the same level of accuracy as SVM when the number of
 332 training samples exceeds 10,000, with a significantly lower average time complexity of $O(n \log n)$
 333 [34]. RF and ANN classifiers were built with Weka 3.8.1 and the neural network toolbox in
 334 Matlab R2015a, respectively, and all parameters took the default values. The independent tests
 335 based on Tr_data and Te_data were employed for comparisons (Table 7).

336

337

Table 7 Independent test accuracy based on different classifiers

Classifier	SN (%)	SP (%)	MCC	Q^9
RF	2.75	99.83	0.115	0.312
ANN	0	99.90	-0.009	0.293
RVKDE	9.84	97.25	0.106	0.362
ChiDT	70.47	66.27	0.205	0.683

338 The results show that: 1) ChiDT achieved a significantly higher predictive accuracy and
339 effectively realized imbalanced pattern classification. When the training set was imbalanced (4748
340 positives/50,551 negatives), the prediction results of RF, ANN and RVKDE were biased to
341 negative samples, resulting in poor sensitivities (SNs<10%). With ANN, while specificity was up
342 to 99.9%, sensitivity was equal to 0. This meant that all positive samples were predicted to be
343 negative and thus the global accuracy of ANN was the lowest ($Q^9=0.293$). ChiDT built a balanced
344 decision table through weighted correction strategy to perform imbalanced classification and
345 obtained the highest accuracy ($Q^9=0.683$). 2) ChiDT has a satisfactory calculation speed and can
346 be applied to large samples. All simulations were run on an Intel Core i5-3320M 2.6 GHz/8 GB
347 RAM system, and the elapsed time of ChiDT and RVKDE were 17 seconds and 18 minutes,
348 respectively. ChiDT's high speed is achieved because there is no need for parameter optimization.

349 **Comparison of different position-based encoding schemes**

350 Based on the 9 key positions in Tr_data (here using 4748 positive samples and 4748 negative
351 samples), 5-fold cross validation was applied to evaluate binary encoding, physicochemical
352 property encoding (including 531 physicochemical properties [7] and 10 physicochemical
353 properties [9] for encoding) and chi-square statistical difference table encoding, respectively. The
354 results showed that chi-square statistical difference table encoding achieved the highest predictive
355 accuracy and the fewest features (Table 8).

356

357

Table 8 5-fold cross accuracy based on different encoding schemes

Encoding scheme	Feature dimension	SN (%)	SP (%)	MCC	Q ⁹
Binary	180	63.20	62.41	0.258	0.623
Physicochemical properties(531)	4779	58.86	60.39	0.188	0.593
Physicochemical properties(10)	90	59.77	62.59	0.225	0.607
Chi-square statistical difference table	9	65.91	62.91	0.289	0.641

358

Binary encoding means that each position is represented by 20 0/1-features and the corresponding feature matrix is therefore very sparse. When using binary encoding scheme, the encodings of the same residue at different positions are the same, which does not reflect positional difference, and for different residues at the same position, it does not reflect the degree of difference between residues. For example, the amino acid polarity indexes of residue S, T and R at the same position are 1.67, 1.66 and 52, respectively, indicating that the polarity difference between S and T is small, and between T and R is large, but the hamming distances of both S-T and S-R are equal to 2 when using binary encoding. As for physicochemical property encoding, when 531 amino acid indices in AAindex were all considered for sequence characterization, the number of features reached $531 \times 9 = 4779$ (Table 8), and a lot of irrelevant and redundant features would be seen. Ning *et al* [9] ranked 531 physicochemical properties according to their abilities to distinguish between true and false succinylation sites, then chose the top 10 physicochemical properties for sequence encoding, so that the feature dimension was greatly reduced. However, as shown in Table 8, no matter whether 531 or 10 physicochemical properties are used, the predictive accuracy is always lower than that of chi-square statistical difference table encoding.

373

Chi-square statistical difference table encoding reflects the difference of the same residue at different positions, as well as the degree of difference between different residues at the same position, thus, it could differentiate between the highly similar positive and negative samples. Another benefit of chi-square statistical difference table encoding is that it has a low feature

376

377 dimension, low redundancy, and a non-sparse feature matrix.

378 **Comparison of different window sizes**

379 Based on Tr_data and Te_data, independent tests were performed to compare the prediction
380 performance of the determined window size (-8~+7) with longer (e.g. -25~+25, -15~+15) and
381 shorter window sizes (e.g. -5~+5). The results (Table 9) show that the proposed model with a
382 window size of 16 residues (-8~+7) can achieve higher independent test accuracy compared to
383 other window sizes. This indicates that an overly long window size may introduce some irrelevant
384 information, while too short a window may lead to insufficient information collection, both of
385 which reduce predictive accuracy. This confirms that our window size determination is reliable.

386

Table 9 Independent test accuracy based on different window sizes

Window size	SN (%)	SP (%)	MCC	Q ⁹
51(-25~+25)	68.50	61.10	0.162	0.646
31(-15~+15)	64.57	66.48	0.174	0.655
16(-8~+7)	70.47	66.27	0.205	0.683
11(-5~+5)	62.20	65.03	0.152	0.636

387 **Necessity of Chi-MIC-share feature selection**

388 Based on Tr_data and Te_data, the independent test results with or without Chi-MIC-share
389 feature selection are shown in Table 10. They show that feature selection based on Chi-MIC-share
390 can: 1) improve predictive accuracy, with the Q⁹ value improving from 0.663 to 0.683, and 2)
391 reduce feature dimension and save computational time. After feature selection, the number of
392 original input features was reduced from 239 to 10, and the elapse time of ChiDT was reduced by
393 95%. Therefore, it is necessary and beneficial to perform a Chi-MIC-share feature selection.

394

Table 10 Independent test accuracy with and without Chi-MIC-share

Feature selection	Feature dimension	SN (%)	SP (%)	MCC	Q ⁹	Time(mm:ss)
No feature selection	239	70.08	62.95	0.182	0.663	06:14
Chi-MIC-share	10	70.47	66.27	0.205	0.683	00:17

395 **Comparison with existing methods**

396 To further evaluate the performance of our method (iSuc-ChiDT), we compared it with
 397 existing succinylation site predictors, SucPred, iSuc-PseAAC [10], SuccFind [11], SuccinSite,
 398 iSuc-PseOpt [35], pSuc-Lys, Success [13] and PSuccE, using the same independent testing set
 399 (Te_data). The results show that iSuc-ChiDT had a superior overall accuracy ($Q^9=0.683$) and
 400 sensitivity (70.47% vs. 12.20%~37.50%) (Table 11).

401

Table 11 Independent test accuracy for different methods

Method	SN (%)	SP (%)	MCC	Q^9
SucPred	27.20	67.30	-0.030	0.436
iSuc-PseAAC	12.20	88.70	0.013	0.374
SuccFind	25.20	79.20	0.029	0.451
SuccinSite	37.10	88.20	0.199	0.548
iSuc-PseOpt	30.30	75.80	0.038	0.478
pSuc-Lys	22.40	82.60	0.036	0.436
Success	14.20	86.80	0.007	0.386
PSuccE	37.50	88.60	0.204	0.551
iSuc-ChiDT	70.47	66.27	0.205	0.683

402 Positional information of amino acids is valuable for succinylation site prediction. Most
 403 compared methods used binary encoding or physicochemical property encoding to extract
 404 positional features. iSuc-ChiDT used chi-square statistical difference table encoding and our
 405 experiments showed that it was superior to these two encoding schemes (see Table 8). Moreover,
 406 iSuc-ChiDT combined positional features and compositional features to characterize samples.
 407 Employing the independent tests with Tr_data and Te_data, the MCC values of 9 positional
 408 features, 230 compositional features and 239 combinational features-based models were 0.153,
 409 0.099 and 0.182, respectively, confirming that feature fusion improved predictive accuracy. The
 410 ChiDT classifier outperformed traditional classifiers when dealing with imbalanced datasets (see
 411 Table 7), further supporting the observation that iSuc-ChiDT could achieve better prediction
 412 performance.

413 **Conclusion**

414 In this study, a computational method called iSuc-ChiDT is proposed to identify protein
415 succinylation sites, which incorporates chi-square statistical difference table encoding and the
416 ChiDT classifier. Chi-square statistical difference table encoding can differentiate between highly
417 similar positive and negative sequences, and its advantages include a reduced feature dimension
418 and non-sparse feature matrix, and the ChiDT classifier efficiently resolves the imbalanced dataset
419 problem, both of which contribute to the accurate prediction of succinylation sites. iSuc-ChiDT
420 greatly improved sensitivity and overall accuracy compared to previous predictors, and it will
421 serve as an useful complementary tool for detecting potential succinylation sites in proteins. In
422 future studies, we aim to explore more valuable features (*e.g.* evolutionary information, structural
423 information) for characterizing succinylation sites, in pursuit of better prediction performance.

424

425 **Additional file**

426 **Additional file 1: Table S1** (.xlsx 14KB). This table shows the 20×9 chi-square statistical
427 difference table constructed based on 9 key positions in Tr_data.

428 **Additional file 2: Table S2** (.xlsx 26KB). This table shows 137 rules generated based on Tr_data,
429 and lists the number of positive and negative training samples conforming to each rule.

430

431 **Abbreviations**

432 ChiDT: chi-square decision table; PTM: post-translational modification; SVM: support vector
433 machine; RF: random forest; ANN: artificial neural network; RVKDE: relaxed variable kernel
434 density estimator; AAC: amino acid composition; PCAAC: pair-coupled amino acid composition;
435 undirected-PCAAC: undirected pair-coupled amino acid composition; SN: sensitivity; SP:
436 specificity; TP: true positive; FP: false positive; TN: true negative; FN: false negative; MCC:
437 Matthews correlation coefficient; MIC: maximal information coefficient; mRMR: minimum
438 redundancy maximum relevance.

439

440 **Declarations**

441 **Ethics approval and consent to participate**

442 Not applicable.

443

444 **Consent for publication**

445 Not applicable.

446

447 **Availability of data and material**

448 All data generated or analyzed during this study are included in this published article and its
449 supplementary information files.

450

451 **Competing interests**

452 The authors declare that they have no competing interests.

453

454 **Funding**

455 This research was supported by the Doctoral Science Foundation of Hunan Institute of
456 Engineering and the Youth Key Research Project of Hunan Institute of Engineering (No. XJ2002).

457

458 **Authors' contributions**

459 YZ, ZMY conceived and designed the experiments. YZ performed the experiments and drafted the
460 manuscript. ZMY revised the manuscript. YC contributed software coding. All authors read and
461 approved the final manuscript.

462

463 **Acknowledgments**

464 We would like to thank all the anonymous reviewers for their constructive advices.

465

466 **Authors' information**

467 ¹College of Computer and Communication, Hunan Institute of Engineering, Xiangtan, Hunan,
468 China. ²Hunan Engineering & Technology Research Center for Agricultural Big Data Analysis &
469 Decision-making, Hunan Agricultural University, Changsha, Hunan, China.

470

471 **References**

- 472 1. Zhang ZH, Tan MJ, Xie ZY, Dai LZ, Chen Y, Zhao TM. Identification of lysine
473 succinylation as a new post-translational modification. *Nature Chemical Biology*.
474 2011(1);7:58-63.
- 475 2. Papanicolaou KN, O'Rourke B, Foster DB. Metabolism leaves its mark on the powerhouse:
476 recent progress in post-translational modifications of lysine in mitochondria. *Frontiers in*
477 *Physiology*. 2014;5:301.
- 478 3. Xu XY, Liu T, Yang J, Chen LH, Liu B, Wei CD, *et al*. The first succinylome profile of
479 *Trichophyton rubrum* reveals lysine succinylation on proteins involved in various key
480 cellular processes. *BMC Genomics*. 2017;18:577.
- 481 4. Shershakova N, Bashkatova E, Babakhin A, Andreev S, Nikonova A, Shilovsky L, *et al*.
482 Allergen-specific immunotherapy with monomeric allergoid in a mouse model of atopic

- 483 dermatitis. *PloS One*. 2015;10(8):e0135070.
- 484 5. Tannahill GM, Curtis AM, Adamik J, Pálsson-McDermott EM, McGettrick AF, Goel G, et al.
485 Succinate is an inflammatory signal that induces IL-1 β through HIF-1 α . *Nature*.
486 2013;496:238-42.
- 487 6. Zhao XW, Ning Q, Chai HT, Ma ZQ. Accurate in silico identification of protein
488 succinylation sites using an iterative semi-supervised learning technique. *Journal of*
489 *Theoretical Biology*. 2015;374:60-5.
- 490 7. Hasan MM, Yang SP, Zhou Y, Mollah MN. SuccinSite: a computational tool for the
491 prediction of protein succinylation sites by exploiting the amino acid patterns and properties.
492 *Molecular BioSystems*. 2016;12:786-95.
- 493 8. Jia JH, Liu Z, Xiao X, Liu BX, Chou KC. pSuc-Lys: Predict lysine succinylation sites in
494 proteins with PseAAC and ensemble random forest approach. *Journal of Theoretical Biology*.
495 2016;394:223-30.
- 496 9. Ning Q, Zhao XS, Bao LL, Ma ZQ, Zhao XW. Detecting succinylation sites from protein
497 sequences using ensemble support vector machine. *BMC Bioinformatics*. 2018;
498 19(1):237-46.
- 499 10. Xu Y, Ding YX, Ding J, Lei YH, Wu LY, Deng NY. iSuc-PseAAC: predicting lysine
500 succinylation in proteins by incorporating peptide position-specific propensity. *Scientific*
501 *Reports*. 2015;5(1):10184.
- 502 11. Xu HD, Shi SP, Wen PP, Qiu JD. SuccFind: a novel succinylation sites online prediction tool
503 via enhanced characteristic strategy. *Bioinformatics*. 2015;31(23): 3748-50.
- 504 12. Dehzangi A, López Y, Lal SP, Taherzadeh G, Michaelson J, Sattar A, *et al.* PSSM-Suc:
505 Accurately predicting succinylation using position specific scoring matrix into bigram for
506 feature extraction. *Journal of Theoretical Biology*. 2017;425:97-102.
- 507 13. López Y, Sharma A, Dehzangi A, Lal SP, Taherzadeh G, Sattar A, *et al.* Success:
508 evolutionary and structural properties of amino acids prove effective for succinylation site
509 prediction. *BMC Genomics*. 2018;19 Suppl 1:923-27.
- 510 14. Dehzangi A, López Y, Lal SP, Taherzadeh G, Sattar A, Tsunoda T, *et al.* Improving
511 succinylation prediction accuracy by incorporating the secondary structure via helix, strand
512 and coil, and evolutionary information from profile bigrams. *PLoS ONE*.
513 2018;13(2):e0191900.
- 514 15. López Y, Dehzangi A, Lal SP, Taherzadeh G, Michaelson J, Sattar A, *et al.* SucStruct:
515 Prediction of succinylated lysine residues by using structural properties of amino acids.
516 *Analytical Biochemistry*. 2017;527:24-32.
- 517 16. Kawashima S, Ogata H, Kanehisa M. AAindex: amino acid index database. *Nucleic Acids*
518 *Research*. 1999;27(1): 368-69.
- 519 17. Weiss GM, Provost F. The effect of class distribution on classifier learning: An empirical
520 study. Technical Report ML-TR-44. Department of Computer Science, Rutgers University,
521 2001.
- 522 18. Li YT, Dai ZJ, Cao D, Luo F, Chen Y, Yuan ZM. Chi-MIC-share: a new feature selection
523 algorithm for quantitative structure-activity relationship models. *RSC Advances*.
524 2020;10:19852-60.
- 525 19. UniProt Consortium. Ongoing and future developments at the Universal Protein Resource.
526 *Nucleic Acids Research*. 2011;39(Database issue):214-9.

- 527 20. NCBI protein sequence database. <https://www.ncbi.nlm.nih.gov/protein/>. Accessed 21 May
528 2021.
- 529 21. Li WZ, Godzik A. Cd-hit: a fast program for clustering and comparing large sets of protein or
530 nucleotide sequences. *Bioinformatics*. 2006;22(13):1658-9.
- 531 22. PSuccE. <https://github.com/ningq669/PSuccE>. Accessed 17 April 2021.
- 532 23. Reshef DN, Reshef YA, Finucane HK, Grossman SR, McVean G, Turnbaugh PJ, *et al.*
533 Detecting novel associations in large data sets. *Science*. 2011;334:1518-24.
- 534 24. Reshef DN, Reshef YA, Finucane HK, Grossman SR, McVean G, Turnbaugh PJ, *et al.*
535 Supporting Online Material for Detecting Novel Associations in Large Data Sets. *Science*.
536 2011; 334:1518-24.
- 537 25. Chen Y, Zeng Y, Luo F, Yuan ZM. A new algorithm to optimize maximal information
538 coefficient. *PLoS One*. 2016;11(6):e0157567.
- 539 26. Zeng Y, Yuan HJ, Yuan ZM, Chen Y. A high-performance approach for predicting donor
540 splice sites based on short window size and imbalanced large samples. *Biology Direct*.
541 2019;14:6.
- 542 27. Chou KC. Using pair-coupled amino acid composition to predict protein secondary structure
543 content. *Journal of Protein Chemistry*. 1999;18(4):473-80.
- 544 28. Peng H, Long F, Ding C. Feature selection based on mutual information criteria of
545 max-dependency, max-relevance, and min-redundancy. *IEEE Transactions on Pattern
546 Analysis & Machine Intelligence*. 2005;27(8), 1226-38.
- 547 29. Sun YM, Kamel MS, Wong AKC, Wang Y. Cost-sensitive boosting for classification of
548 imbalanced data. *Pattern Recognition*, 2007;40(12):3358-78.
- 549 30. Zhang CT, Zhang R. Evaluation of gene-finding algorithms by a content-balancing accuracy
550 index. *Journal of Biomolecular Structure and Dynamics*. 2002;19(6):1045-52.
- 551 31. Zhang CT, Ren Z. Q9, a content-balancing accuracy index to evaluate algorithms of protein
552 secondary structure prediction. *International Journal of Biochemistry and Cell Biology*.
553 2003;35:1256-62.
- 554 32. Zhang QW, Peng QK, Zhang Q, Yan YH, Li KK, Li J. Splice sites prediction of human
555 genome using length-variable Markov model and feature selection. *Expert Systems with
556 Applications*. 2010;37:2771-82.
- 557 33. Wei D, Zhang HL, Wei YJ, Jiang QS. A Novel Splice Site Prediction Method using Support
558 Vector Machine. *Journal of Computational Information Systems*. 2013; 20:8053-60.
- 559 34. Oyang YJ, Hwang SC, Ou YY, Chen CY, Chen ZW. Data classification with radial basis
560 function networks based on a novel kernel density estimation algorithm. *IEEE Transactions
561 on Neural Networks*. 2005;16(1):225-36.
- 562 35. Jia JH, Liu Z, Xiao X, Liu BX, Chou KC. iSuc-PseOpt: Identifying lysine succinylation sites
563 in proteins by incorporating sequence-coupling effects into pseudo components and
564 optimizing imbalanced training dataset. *Analytical Biochemistry*. 2016;497:48-56.

565

566 **Fig. 1 Flow chart of iSuc-ChiDT**

567 **Fig. 2 Illustration of compression procedure (position -10 in Tr_data)**

568 **Fig. 3 Chi-square values for different positions in Tr_data**

569 **Fig. 4 Chi-MIC-share scores after introduction of each feature.** The red line represents the
570 forced termination of feature introduction.

Figures

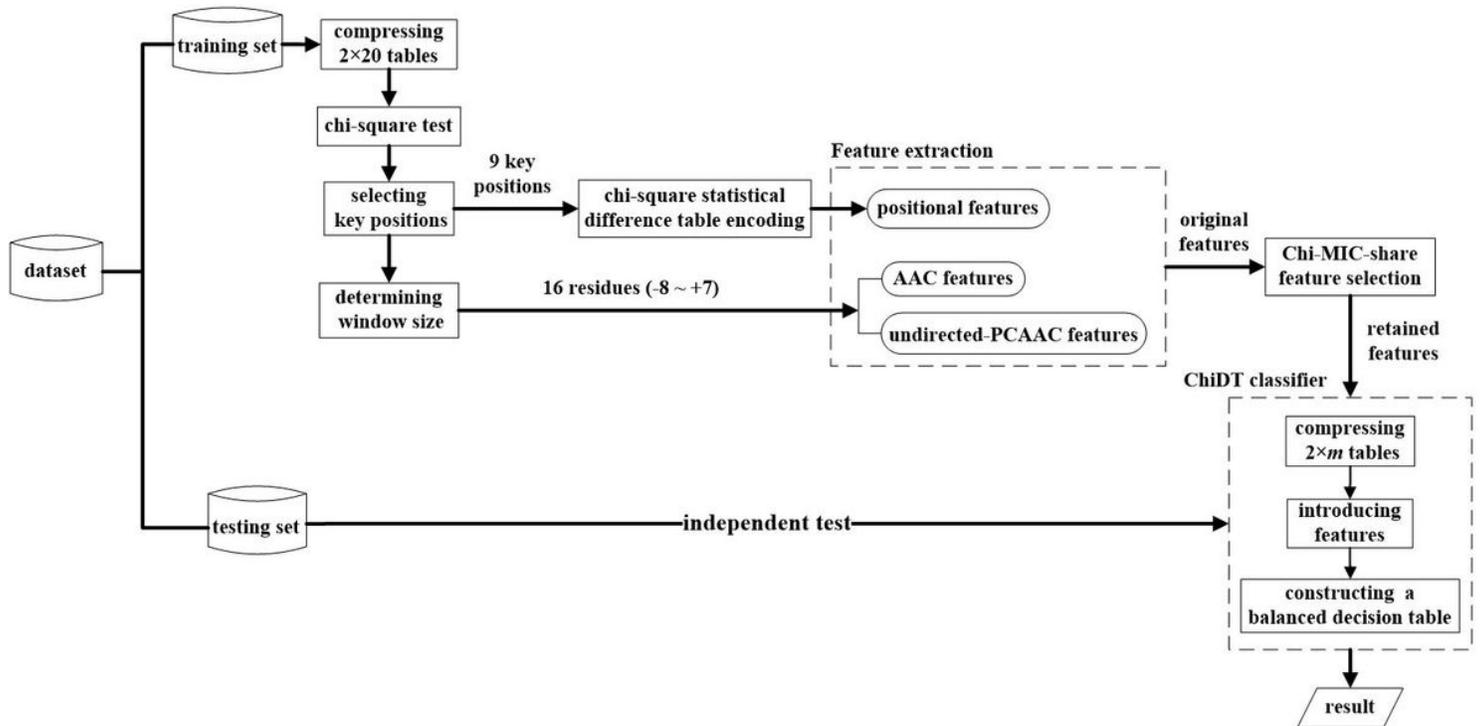


Figure 1

Flow chart of iSuc-ChiDT

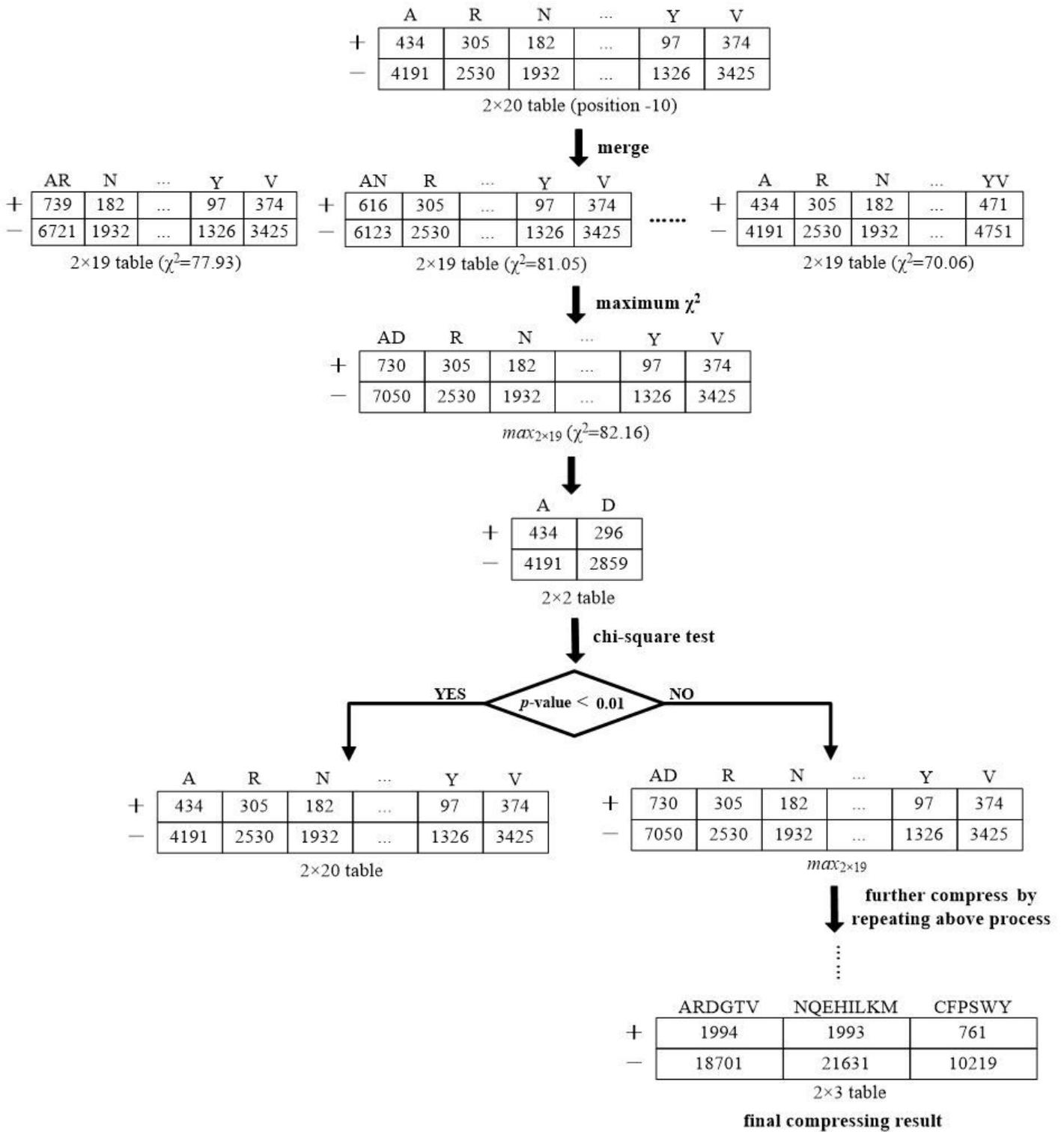


Figure 2

Illustration of compression procedure (position -10 in Tr_data)

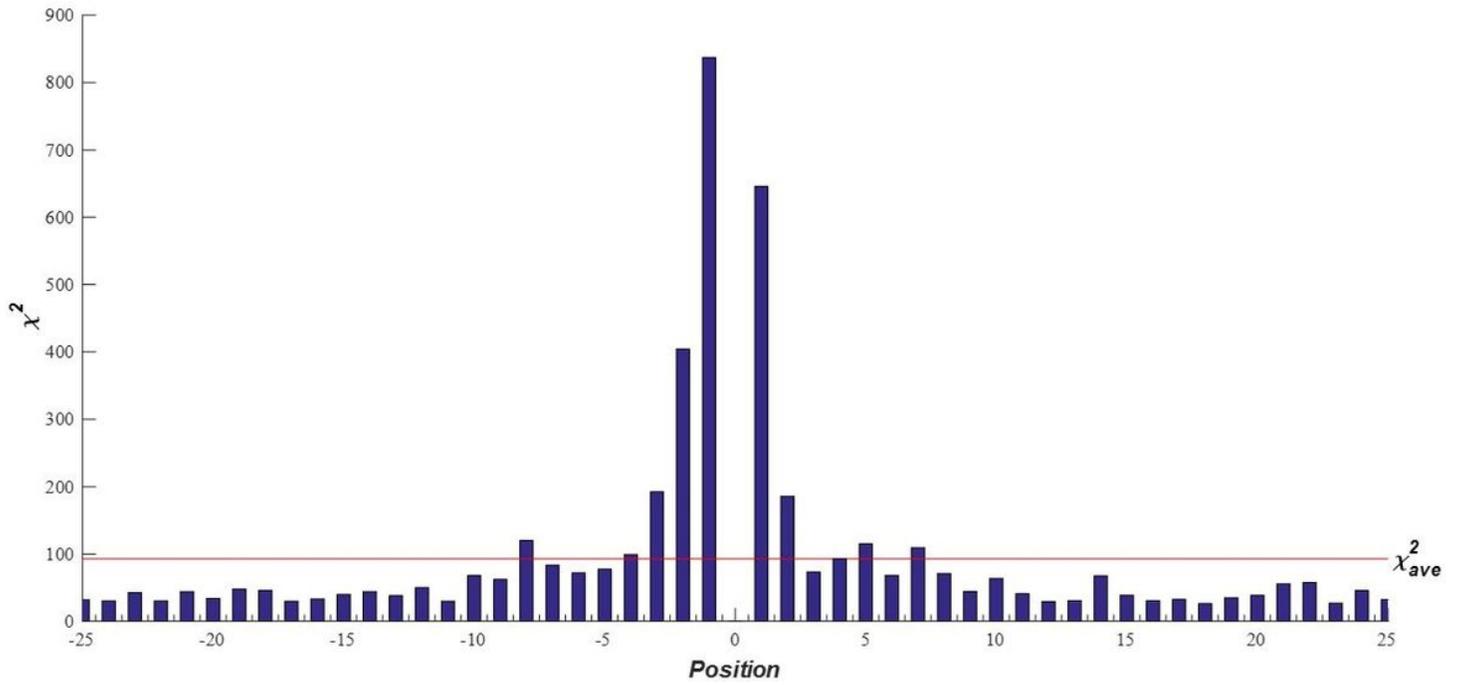


Figure 3

Chi-square values for different positions in Tr_data

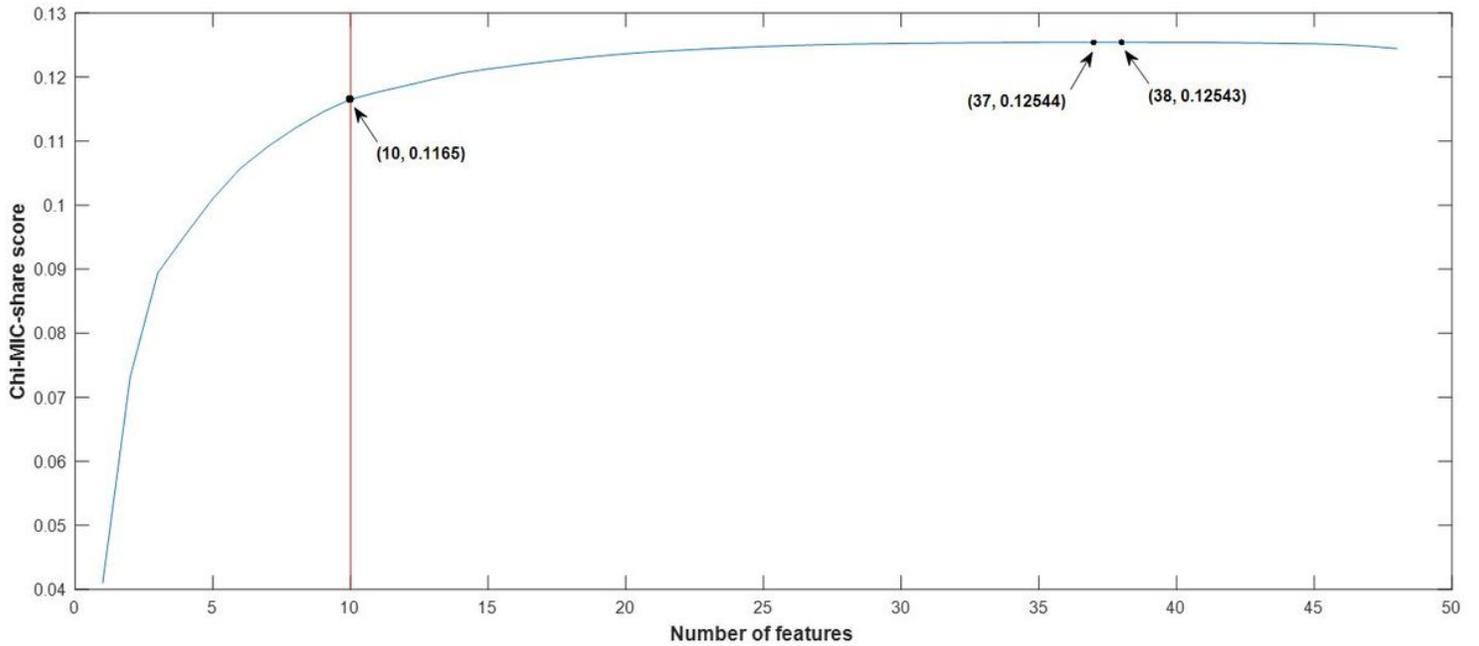


Figure 4

Chi-MIC-share scores after introduction of each feature. The red line represents the forced termination of feature introduction.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [TableS1.xlsx](#)
- [TableS2.xlsx](#)