

Feature Extraction Approaches for Biological Sequences: A Comparative Study of Mathematical Models

Robson Parmezan Bonidia^{a,b,*}, Lucas Dias Hiera Sampaio^a, Fabrício Martins Lopes^a, André Carlos Ponce de Leon Ferreira de Carvalho^b, Danilo Sipoli Sanches^a

^a*Department of Computer Science, Bioinformatics Graduate Program (PPGBIOINFO), Federal University of Technology - Paraná, UTFPR, Campus Cornélio Procópio, Brazil.*

^b*Institute of Mathematics and Computer Sciences, University of São Paulo - USP, São Carlos, 13566-590, Brazil*

Abstract

The number of available biological sequences has increased significantly in recent years due to various genomic sequencing projects, creating a huge volume of data. Consequently, new computational methods are needed to analyze and extract information from these sequences. Machine learning methods have shown broad applicability in computational biology and bioinformatics. The utilization of machine learning methods has helped to extract relevant information from various biological datasets. However, there are still several problems that motivate new algorithms and pipeline proposals, mainly involving feature extraction problems, in which extracting significant discriminatory information from a biological set is challenging. Considering this, our work proposes to study and analyze a feature extraction pipeline based on mathematical models (Numerical Mapping, Fourier, Entropy, and Complex Networks). As a case study, we analyze Long Non-Coding RNA sequences. Moreover, we divided this work into two studies, e.g., (I) we assessed our proposal with the most addressed problem in our review, e.g., lncRNA vs. mRNA; (II) we tested its generalization on different classification problems, e.g., circRNA vs. lncRNA. The experimental results demonstrated three main contributions: (1) An in-depth study of several mathematical models; (2) a new feature extraction pipeline and (3) its generalization and robustness

*Corresponding author

Email address: robserveridor@gmail.com (Robson Parmezan Bonidia)

for distinct biological sequence classification.

Keywords: Feature Extraction; Long Non-Coding RNAs; Biological Sequences; Numerical Mapping Techniques; Fourier; Complex Networks; Shannon; Tsallis.

1 1. Background

2 In recent years, due to advances in DNA sequencing, an increasing num-
3 ber of biological sequences have been generated by thousands of sequencing
4 projects [1], creating a huge volume of data [2]. During the last decade,
5 Machine Learning (ML) methods have shown broad applicability in compu-
6 tational biology and bioinformatics [3]. Consequently, the ability to process
7 and analyze biological data has advanced significantly [4]. Tools have been
8 applied in gene networks, protein structure prediction, genomics, proteomics,
9 protein-coding genes detection, disease diagnosis, and drug planning [5, 6].
10 Fundamentally, ML investigates how computers can learn (or improve their
11 performance) based on the data. Moreover, ML is a specialization of com-
12 puter science related to pattern recognition and artificial intelligence [7].

13 Based on this, several works have focused on investigating sequences of
14 DNA and RNA molecules. Applying ML methods in these sequences has
15 helped to extract important information from various datasets to explain
16 biological phenomena [3]. The development of efficient approaches benefits
17 the mathematical understanding of the structure of biological sequences [1],
18 e.g., Precision cancer diagnostics [8] and the Coronavirus epidemic [9, 10].
19 However, according to [3, 11], there are still several challenging biological
20 problems that motivated the emergence of proposals for new algorithms.
21 Fundamentally, biological sequence analysis with ML presents one major
22 problem: Feature Extraction [12].

23 Feature extraction seeks to generate a feature vector, optimally trans-
24 forming the input data [12]. This procedure is exceptionally relevant for the
25 success of the ML application. Another primary goal of feature extraction is
26 to extract important information from input data compactly, as well as re-
27 moving noise and redundancy to increase the accuracy of ML models [13, 12].
28 Furthermore, the feature extraction is an inevitable method, especially in the
29 stage of biological sequence preprocessing [14].

30 Necessarily, several methods in bioinformatics apply ML algorithms for
31 sequence classification, and as many algorithms can deal only with numerical

32 data, sequences need to be translated into sequences of numbers. Thereby,
33 modern applications extract relevant features from sequences based on several
34 biological properties, e.g., physicochemical, Open Reading Frames (ORF)-
35 based, usage frequency of adjoining nucleotide triplets, GC content, among
36 others. This approach is common in biological problems, but these implemen-
37 tations are often difficult to reuse or adapt to another specific problem, e.g.,
38 ORF features are an essential guideline for distinguishing Long non-coding
39 RNAs (lncRNA) from protein-coding genes [15], but not useful features for
40 classifying lncRNA classes [2]. Consequently, the feature extraction problem
41 arises, in which extracting a set of useful features that contain significant
42 discriminatory information becomes a fundamental step in the construction
43 of a predictive model [16].

44 Therefore, these problems make the process of biological sequence clas-
45 sification a challenging task, creating a growing need to develop new tech-
46 niques and methods to analyze sequences effectively and efficiently. Thereby,
47 this work studies the performance of different feature extraction methods
48 for biological sequence analysis, using mathematical models, e.g., numerical
49 mapping, Fourier transform, entropy, and graphs. As a case study, we will
50 use lncRNA sequences, which are fundamentally unable to produce proteins
51 [17] and have recently casted doubt on its functionality [18].

52 LncRNAs present several problem classes (e.g., lncRNA vs. mRNA [19,
53 20] and lncRNA vs. circRNA [21]), thus enabling us to create a scenario to
54 answer the questions raised in this work. Fundamentally, our main objective
55 is to propose generalist techniques, demonstrating their efficiency concerning
56 biological features. We consider biological approaches, those characteristics
57 that present a bias to the analyzed problem or some biological explanation,
58 e.g., ORF for lncRNA vs. mRNA [6, 15], as well as mathematical approaches
59 and information quantity measures such as entropy. Based on this context
60 and objectives, we assume the following hypothesis:

- 61 • **Hypothesis:** Feature extraction approaches based on mathematical
62 models are as efficient and generalist as biological approaches.

63 Considering this, our work contributes to the area of computer science
64 and bioinformatics. Specifically, it introduces new ideas and analysis for
65 the feature extraction problem in biological sequences. Thereby, we present
66 four new contributions: (1) A feature extraction pipeline using mathematical
67 models; (2) Analysis of 9 different mathematical models; (3) Analysis of 6

68 numerical mappings with Fourier, proposing statistical characteristics; (4)
69 The generalization and robustness of mathematical approaches for the feature
70 extraction in biological sequences.

71 **2. Related Works**

72 Essentially, as emphasized, we adopt lncRNA sequences as a case study, a
73 class of Non-Coding RNAs (ncRNAs). Fundamentally, ncRNAs are unable to
74 produce proteins. However, these ncRNAs contain unique information that
75 produces other functional RNA molecules [22, 17]. Moreover, they demon-
76 strate essential roles in cellular mechanisms, playing regulatory roles in a
77 wide variety of biological reactions and processes [22]. The ncRNAs can be
78 classified by length into two classes: Long Non-Coding RNA (lncRNA - 200
79 nucleotides (nt) or more) and short ncRNA (less than 200 nt) [23, 24]. The
80 lncRNAs are sequences with a length greater than 200 nucleotides [25], and
81 according to recent studies, play essential roles in several critical biological
82 processes [26, 27, 28], including transcriptional regulation [29], epigenetics
83 [30], cellular differentiation [31], and immune response [32]. Moreover, they
84 are correlated with some complex human diseases, such as cancer and neu-
85 rodegenerative diseases [6, 33, 34].

86 In plants, according to [6, 35], the lncRNAs act in gene silencing, flowering
87 time control, organogenesis in roots, photomorphogenesis in seedlings, stress
88 responses [36, 37], and reproduction [38]. Furthermore, lncRNAs are present
89 in large numbers in genome [39] and have similar sequence characteristics
90 with protein-coding genes, such as 5' cap, alternative splicing, two or more
91 exons [40], and polyA+ tails [41]. They are also observed in almost all living
92 beings, not only in animals and plants but also yeasts, prokaryotes, and even
93 viruses [42, 43].

94 According to [39], lncRNAs do not contain functional ORFs. However,
95 recent studies have found bifunctional RNAs [44], raising the possibility that
96 many protein-coding genes may also have non-coding functions. Further-
97 more, lncRNAs can be grouped into five broad categories. The classifi-
98 cation occurs conforming to the genomic location, that is, where they are
99 transcribed, concerning well-established markers, e.g., protein-coding genes.
100 Among the categories are [45, 40]: sense, antisense, bidirectional, intronic,
101 intergenic. The genomic context does not necessarily provide some informa-
102 tion about the lncRNAs function or evolutionary origin; nevertheless, it can
103 be used to organize these broad categories [46].

104 In this context, we have conducted an in-depth review of the lncRNAs
 105 classification methods, in which several approaches have been developed,
 106 such as: CPC [47], CPAT [48], CNCI [49], PLEK [50], lncRNA-MFDL [51],
 107 LncRNA-ID [52], lncRScan-SVM [53], LncRNApred [54], DeepLNC [55],
 108 PlantRNA_Sniffer [56], PLncPRO [57], RNAplonc [58], BASiNET [59], and
 109 LncFinder [20]. For better understanding, Figure 1 presents these works
 110 divided into Mathematical, Biological, and Hybrid approaches.

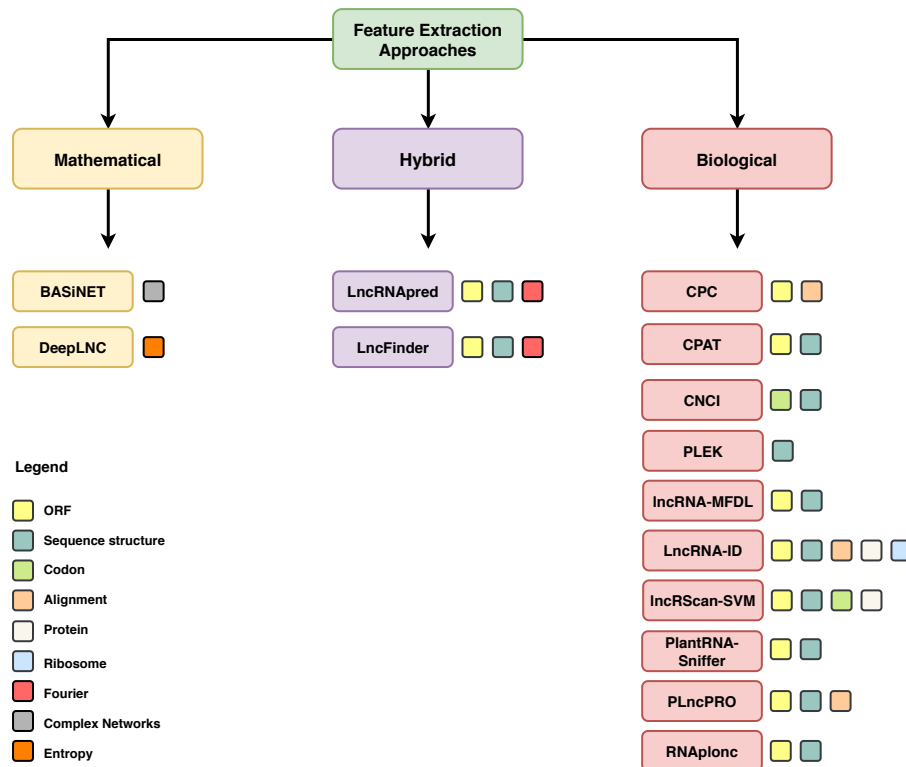


Figure 1: Feature extraction approaches in our case study divided into: Mathematical, Biological, and Hybrid.

111 The CPC uses the extent and quality of the ORF, and derivation of the
 112 BLASTX [60] search to measure the protein-coding potential of a transcript.
 113 In the classification, the authors applied the LIBSVM package to train a Sup-
 114 port Vector Machine (SVM) model, using the standard radial basis function
 115 kernel. CPAT classifies transcripts of coding and non-coding using the Logis-
 116 tic Regression (LR) classifier. This approach implements four features: ORF

117 coverage, ORF size, hexamer usage bias, and Fickett TESTCODE statis-
118 tic. CNCI was induced with SVM and applies profiling Adjoining Nucleotide
119 Triplets, and most-like CDS (MLCDS).

120 In contrast, PLEK (2014) is based on the k -mer scheme ($k = 1, \dots, 5$)
121 to predict lncRNA, also applying the SVM classifier. lncRNA-MFDL uses
122 Deep Learning (DL) and multiple features, among them: ORF, K-mer ($k =$
123 $1, 2, 3$), secondary structure (minimum free energy), and MLCDS. LncRNA-
124 ID predicts lncRNAs with Random Forest (RF) through ORF (length and
125 coverage), sequence structure (Kozak motif), ribosome interaction, alignment
126 (profile Hidden Markov Mode - profile HMM), and protein conservation.

127 lncRScan-SVM uses stop codon count, GC content, ORF (score, CDS
128 length and CDS percentage), transcript length, exon count, exon length, and
129 average PhastCons scores. LncRNApred classified lncRNAs with RF and
130 features based on ORF, signal to noise ratio, k -mer ($k = 1, 2, 3$), sequence
131 length, and GC content. DeepLNC uses only the k -mer scheme with entropy
132 and Deep Neural Network (DNN). PlantRNA_Sniffer was developed in 2017
133 to predict Long Intergenic Non-Coding RNAs (lincRNAs). The method ap-
134 plied SVM and extracted features from ORF (proportion and length) and
135 nucleotide patterns.

136 PLncPRO is based on machine learning and uses RF. The features se-
137 lected include ORF quality (score and coverage), number of hits, significance
138 score, total bit score, and frame entropy. RNAplnc classified sequences
139 with the REPTree algorithm, considering 16 features (ORF, GC content, K-
140 mer scheme ($k = 1, \dots, 6$), sequence length). BASiNET classifies sequences
141 based on the feature extraction from complex network measurements. Lastly,
142 LncFinder tests five classifiers (LR, SVM, RF, Extreme Learning Machine,
143 and Deep Learning), to apply the algorithm that obtains the highest ac-
144 curacy. The authors extract features from ORF, secondary structural, and
145 EIIP-based physicochemical properties.

146 In general, the aforementioned works apply supervised learning methods
147 using binary classification (two classes - lncRNAs and protein-coding genes
148 (mRNA)). There is a considerable amount of research on humans, followed
149 by animals and plants. Regarding feature extraction, we observed a full do-
150 main of ORF and sequence-structure descriptors. As seen in Figure 1, there
151 is a frequent use of biological features. On the other hand, some works have
152 explored mathematical approaches for feature extraction, such as Genomic
153 Signal Processing (GSP), DNA Numerical Representation (DNR) [54, 20],
154 and Complex Networks [59]. Nevertheless, the authors used these charac-

155 teristics in conjunction with other biological feature extraction techniques
156 or without testing other mathematical features. Practically no papers have
157 focused on several mathematical approaches. Based on this, the objective of
158 this section was to summarize the main methods of the literature and their
159 characteristic descriptors. Therefore, we will not use the works shown for
160 comparison, but the most applied features.

161 3. Materials and Methods

162 In this section, we describe the methodological approach used to achieve
163 the proposed objectives, as shown in Figure 2. Essentially, we divided our
164 study into five stages: (1) Data selection and preprocessing; (2) Feature
165 extraction; (3) Training; (4) Testing; (5) Performance analysis. Hence, each
166 stage of the study is described, as well as information about the adopted
167 process.

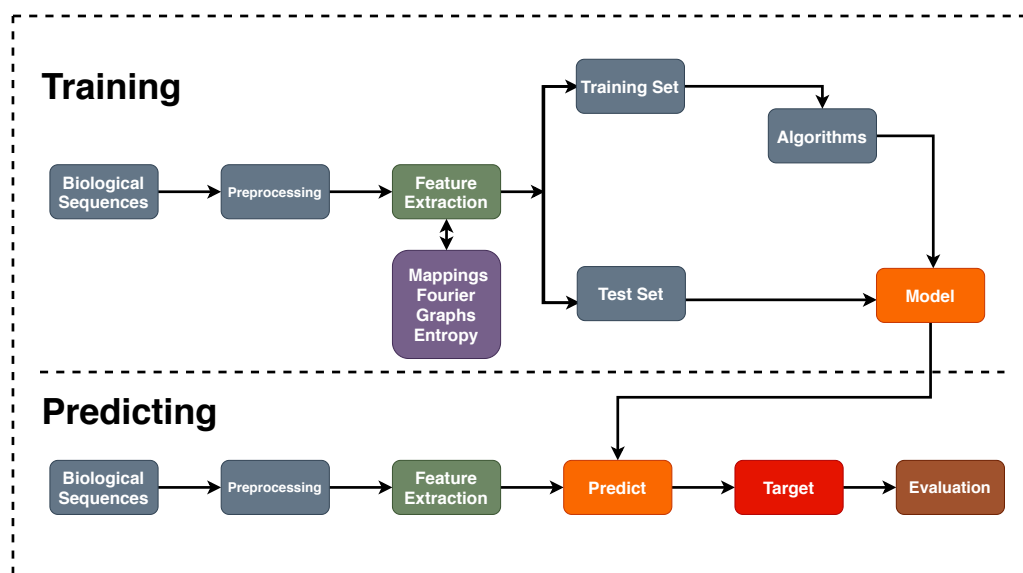


Figure 2: Proposed Pipeline. Essentially, (1) datasets are preprocessed; (2) Feature extraction techniques are applied to each dataset; (3) Machine learning algorithms are executed in the training set to induce predictive models; (4) Induced models are applied to the test set; Finally, (5) the models are evaluated.

168 This work was also divided into two case studies: (I) We assessed our
169 mathematical approaches with the most addressed problem in our review,

170 e.g., lncRNA vs. mRNA; (II) We tested its generalization on different clas-
171 sification problems.

172 3.1. Data Selection

173 As previously mentioned, we chose the lncRNAs classification problem,
174 because it is a new and relevant theme in the literature, in which, recently,
175 it has presented several works, mainly with ML, as explored in Section 2.
176 However, we will also adopt other datasets to assess the generalization of
177 mathematical features. As preprocessing, we used only sequences longer
178 than 200nt [50], and we also removed sequence redundancy. Moreover, the
179 sampling method was adopted in our dataset, since we are faced with the
180 *imbalanced data problem* [2]. Therefore, we applied random majority under-
181 sampling, which consists of removing samples from the majority class (to
182 adjust the class distribution) [61]. Finally, we divided this paper into two
183 case studies.

184 3.1.1. Case Study I

185 Sequences of five plant species were adopted to validate the proposed
186 approaches. The summary of the dataset can be seen in Table 1. According to
187 the literature approaches, this study also adopts two classes for the datasets:
188 the positive class, with lncRNAs, and the negative class, with protein-coding
189 genes (mRNAs).

Table 1: Adopted species to create the datasets.

Species	Sequences	Samples	Preprocessing	Selected
<i>A. trichopoda</i>	lncRNA	5698	4556	4556
	mRNA	26846	22326	4556
<i>A. thaliana</i>	lncRNA	2540	2540	2540
	mRNA	13973	13973	2540
<i>C. sinensis</i>	lncRNA	2562	2215	2215
	mRNA	46147	45846	2215
<i>C. sativus</i>	lncRNA	1929	1730	1730
	mRNA	30364	29829	1730
<i>R. communis</i>	lncRNA	4198	3487	3487
	mRNA	31221	29042	3487

190 The mRNA data of the *Arabidopsis thaliana* (obtained from CPC2 [19])
191 were built from the RefSeq database with protein sequences annotated by

192 Swiss-Prot [19], and lncRNA data from the **Ensembl** (*v87*) and **Ensembl**
193 **Plants** (*v32*) database. The mRNA transcript data of the *Amborella tri-*
194 *chopoda*, *Citrus sinensis*, *Cucumis sativus* and *Ricinus communis* were ex-
195 tracted from **Phytozome** (version 13) [62]. The lncRNAs data from these
196 species were extracted from **GreenC** (version 1.12) [63].

197 3.1.2. Case Study II

198 In this case study, we will apply the best mathematical models (con-
199 sidering accuracy) of case study I to different classification problems with
200 lncRNAs, in order to test their generalization. Thus, divided this part into
201 three problems:

202 • **Problem 1** (lncRNA vs. sncRNA): Dataset with only non-coding
203 sequences (lncRNA and Small non-coding RNAs (sncRNAs), also ob-
204 tained from [19])

205 – lncRNA: 1291 sequences — sncRNA: 1291 sequences

206 • **Problem 2** (lncRNA vs. Antisense): Dataset with lncRNAs and long
207 noncoding antisense transcripts (obtained from [64]).

208 – lncRNA: 57 sequences — Antisense: 57 sequences

209 • **Problem 3** (circRNA vs. lncRNA): Dataset with lncRNA and circu-
210 lar RNAs (cirRNAs) sequences (circRNA obtained from PlantcircBase
211 [65]. This problem was based on [66] and [21], in order to classify
212 circRNA from other lncRNAs.

213 – circRNA: 2540 sequences — lncRNA: 2540 sequences

214 It is important to emphasize that we used only sequences from *Arabidop-*
215 *sis thaliana* in this second case study because it is the model species in
216 plants. Moreover, plant sequences is the least addressed field by the studies,
217 consequently presenting more challenges.

218 3.2. Feature Extraction

219 In this section, 9 feature extraction approaches are shown: 6 numer-
220 ical mapping techniques with Fourier transform, Entropy, Complex Net-
221 works. It is necessary to emphasize that we denote a biological sequence
222 $\mathbf{s} = (s[0], s[1], \dots, s[N - 1])$ such that $\mathbf{s} \in \{A, C, G, T\}^N$ [2].

223 3.3. Fourier Transform and Numerical Mappings

224 To extract features based on a Fourier model, we applied the Discrete
225 Fourier Transform (DFT), widely used for digital image and signal processing
226 (here GSP), which can reveal hidden periodicities after transformation of
227 time domain data to frequency domain space [67]. According to Yin and
228 Yau [68], the DFT of a signal with length N , $\mathbf{x} \in \mathbb{R}^N$, at frequency k , can
229 be defined by Equation (1):

$$X[k] = \sum_{n=0}^{N-1} x[n] e^{-j\frac{2\pi}{N}kn}, \quad k = 0, 1, \dots, N-1. \quad (1)$$

230 This method is has been widely studied in bioinformatics, mainly for
231 analysis of periodicities and repetitive elements in DNA sequences [69] and
232 protein structures [70]. This approach is shown in Figure 3 and was based
233 on [2].

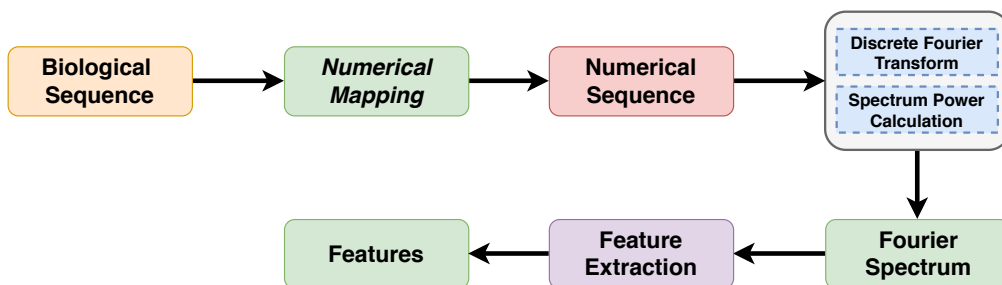


Figure 3: Fourier Transform and Numerical Mapping Pipeline. (1) Each sequence is mapped to a numerical sequence; (2) DFT is applied to the generated sequence; (3) The spectrum power is calculated; (4) The Feature Extraction is performed; Finally, (5) the features are generated.

234 To calculate DFT, we will use the Fast Fourier Transform (FFT), that
235 is a highly efficient procedure for computing the DFT of a time series [71].
236 However, to use GSP techniques, a numeric representation should be used
237 for the transformation or mapping of genomic data. In the literature, dis-
238 tinct DNR techniques have been developed [72]. According to Mendizabal-
239 Ruiz et al. [73], these representations can be divided into three categories:
240 single-value mapping, multidimensional sequence mapping, and cumulative
241 sequence mapping. Thereby, we study 6 numerical mapping techniques (or
242 representations), which will be presented below: Voss [74], Integer [73, 75],
243 Real [76], Z-curve [77], EIIP [78] and Complex Numbers [72, 79, 80].

244 *3.3.1. Voss Representation*

245 This representation can use single or multidimensional vectors. Funda-
 246 mentally, this approach transforms a sequence $\mathbf{s} \in \{A, C, G, T\}^N$ into a
 247 matrix $\mathbf{V} \in \{0, 1\}^{4 \times N}$ such that $\mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3, \mathbf{v}_4]^T$, where T is the trans-
 248 pose operator and each \mathbf{v}_i array is constructed according to the following
 249 relation:

$$v_i[n] = \begin{cases} 1, & s[n] = \alpha[i] \\ 0, & s[n] \neq \alpha[i] \end{cases}, \text{ where } \alpha = (A, C, G, T), \quad n = 0, 1, \dots, N - 1. \quad (2)$$

250 As a result, each row of matrix \mathbf{V} may be seen as an array that marks each
 251 base position such that the first row denotes the presence of base A , row two
 252 for base C , row three base G and the last row for base T . For example, let $\mathbf{s} =$
 253 $(G, A, G, A, G, T, G, A, C, C, A)$ be a sequence that needs to be represented
 254 using Voss representation, therefore, $\mathbf{v}_1 = (0, 1, 0, 1, 0, 0, 0, 1, 0, 0, 1)$, which
 255 represents the locations of bases A , $\mathbf{v}_2 = (0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0)$ for bases
 256 C , $\mathbf{v}_3 = (1, 0, 1, 0, 1, 0, 1, 0, 0, 0, 0)$ for the G bases, $\mathbf{v}_4 = (0, 0, 0, 0, 0, 1, 0,$
 257 $0, 0, 0, 0)$ for T bases. Then, using the DFT in the indicator sequences shown
 258 above, we obtain (see Equation 3):

$$V_i[k] = \sum_{n=0}^{N-1} v_i[n] e^{-j \frac{2\pi}{N} kn}, \quad \forall i \in [1, 4], \quad k = 0, 1, \dots, N - 1. \quad (3)$$

259 The power spectrum of a biological sequence can be obtained by Equation
 260 (4):

$$P_V[k] = \sum_{i=1}^4 |V_i[k]|^2, \quad k = 0, 1, \dots, N - 1. \quad (4)$$

261 *3.3.2. Integer Representation*

262 This representation is one-dimensional [75, 73]. This mapping can be
 263 obtained by substituting the four nucleotides (T, C, A, G) of a biological
 264 sequence for integers (0, 1, 2, 3), respectively, e.g., let $\mathbf{s} = (G, A, G, A, G,$
 265 $T, G, A, C, C, A)$, thus, $\mathbf{d} = (3, 2, 3, 2, 3, 0, 3, 2, 1, 1, 2)$, as exposed in
 266 Equation (5). The DFT and power spectrum are presented in Equation (6).

$$d[n] = \begin{cases} 3, & s[n] = G \\ 2, & s[n] = A \\ 1, & s[n] = C \\ 0, & s[n] = T \end{cases} \quad n = 0, 1, \dots, N - 1. \quad (5)$$

$$D[k] = \sum_{n=0}^{N-1} d[n]e^{-j\frac{2\pi}{N}kn}, \quad P_D[k] = |D[k]|^2, \quad k = 0, 1, \dots, N - 1. \quad (6)$$

267 3.3.3. Real Representation

268 In this representation, Chakravarthy et al. [76] use real mapping based on
 269 the complement property of the complex mapping of [69]. This mapping ap-
 270 plies negative decimal values for the purines (A, G), and positive decimal val-
 271 ues for the pyrimidines (C, T), e.g., let $\mathbf{s} = (G, A, G, A, G, T, G, A, C, C, A)$,
 272 thus, $\mathbf{r} = (-0.5, -1.5, -0.5, -1.5, -0.5, 1.5, -0.5, -1.5, 0.5, 0.5, -1.5)$, as Equation
 273 (7) and Equation (8).

$$r[n] = \begin{cases} -0.5, & s[n] = G \\ -1.5, & s[n] = A \\ 0.5, & s[n] = C \\ 1.5, & s[n] = T \end{cases} \quad n = 0, 1, \dots, N - 1. \quad (7)$$

274

$$R[k] = \sum_{n=0}^{N-1} r[n]e^{-j\frac{2\pi}{N}kn}, \quad P_R[k] = |R[k]|^2, \quad k = 0, 1, \dots, N - 1. \quad (8)$$

275 3.3.4. Z-curve Representation

276 The Z-curve scheme is a three-dimensional curve presented by [77], to
 277 encode DNA sequences with more biological semantics. Essentially, we can
 278 inspect a given sequence $s[n]$ of length N , taking into account the n -th el-
 279 ement of the sequence ($n = 1, 2, \dots, N$). Then, we denote the cumulative
 280 occurrence numbers A_n, C_n, G_n and T_n for each base A, C, G and T , as the
 281 number of times that a base occurred from $s[1]$ up until $s[n]$. Fundamentally,
 282 this method reduces the number of indicator sequences from four (Voss) to
 283 three (Z-curve) in a symmetrical way for all four components [81]. Therefore:

$$A_n + C_n + G_n + T_n = n \quad (9)$$

284 Where the Z-curve consists of a series of nodes P_1, P_2, \dots, P_N , whose
 285 coordinates $x[n]$, $y[n]$, and $z[n]$ ($n = 1, 2, \dots, N$) are uniquely determined by
 286 the Z-transform, shown in Equation (10):

$$P[n] = \begin{cases} x[n] = (A_n + G_n) - (C_n + T_n) \\ y[n] = (A_n + C_n) - (G_n + T_n), \\ z[n] = (A_n + T_n) - (C_n + G_n) \end{cases} \quad (10)$$

$$x[n], y[n], z[n] \in [-n, n], \quad n = 1, 2, \dots, N.$$

287 The coordinates $x[n]$, $y[n]$, and $z[n]$ represent three independent distri-
 288 butions that fully describe a sequence [72]. Therefore, we will have three dis-
 289 tributions with definite biological significance: (1) $x[n]$ = purine/pyrimidine,
 290 (2) $y[n]$ = amino/keto, (3) $z[n]$ = weak hydrogen bonds/strong hydro-
 291 gen bonds [77], e.g., let $\mathbf{s} = (\text{G}, \text{A}, \text{G}, \text{A}, \text{G}, \text{T}, \text{G}, \text{A}, \text{C}, \text{C}, \text{A})$, thus,
 292 $\mathbf{x} = (1, 2, 3, 4, 5, 4, 5, 6, 5, 4, 5)$; $\mathbf{y} = (-1, 0, -1, 0, -1, -2, -3, -2, -1, 0, 1)$;
 293 $\mathbf{z} = (-1, 0, -1, 0, -1, 0, -1, 0, -1, -2, -1)$. Essentially, the difference be-
 294 tween each dimension at the n -th position and the previous $(n - 1)$ position
 295 can be either 1 or -1 [77]. Therefore, we may define the following set of
 296 equations in order to update the values of each dimension array considering
 297 that $x[-1] = y[-1] = z[-1] = 0$:

$$x[n] = \begin{cases} x[n-1] + 1, & s[n] = A \text{ or } G \\ x[n-1] - 1, & s[n] = C \text{ or } T \end{cases} \quad (11)$$

$$y[n] = \begin{cases} y[n-1] + 1, & s[n] = A \text{ or } C \\ y[n-1] - 1, & s[n] = G \text{ or } T \end{cases} \quad n = 1, 2, \dots, N. \quad (12)$$

$$z[n] = \begin{cases} z[n-1] + 1, & s[n] = A \text{ or } T \\ z[n-1] - 1, & s[n] = G \text{ or } C \end{cases} \quad (13)$$

298 Finally, the DFT and power spectrum of the Z-Curve representation may
 299 be defined as [82]:

$$X[k] = \sum_{n=1}^N x[n] e^{-j \frac{2\pi}{N} kn}, \quad Y[k] = \sum_{n=1}^N y[n] e^{-j \frac{2\pi}{N} kn}, \quad Z[k] = \sum_{n=1}^N z[n] e^{-j \frac{2\pi}{N} kn}. \quad (14)$$

$$P_C[k] = |X[k]|^2 + |Y[k]|^2 + |Z[k]|^2, \quad k = 1, 2, \dots, N. \quad (15)$$

300 *3.3.5. EIIP Representation*

301 Nair and Sreenadhan [78] proposed EIIP values of nucleotides to repre-
 302 sent biological sequences and to locate exons. According to the authors, a
 303 numerical sequence representing the distribution of free electron energies can
 304 be called "EIIP indicator sequence", e.g., let $\mathbf{s} = (G, A, G, A, G, T, G, A,$
 305 $C, C, A)$, thus, $\mathbf{b} = (0.0806, 0.1260, 0.0806, 0.1260, 0.0806, 0.1335, 0.0806,$
 306 $0.1260, 0.1340, 0.1340, 0.1260)$, as shown in Equation (16). The DFT and
 307 power spectrum of this representation are presented in Equation (17).

$$b[n] = \begin{cases} 0.0806, & s[n] = G \\ 0.1260, & s[n] = A \\ 0.1340, & s[n] = C \\ 0.1335, & s[n] = T \end{cases}, \quad n = 0, 1, \dots, N - 1. \quad (16)$$

$$B[k] = \sum_{n=0}^{N-1} b[n]e^{-j\frac{2\pi}{N}kn}, \quad P_B[k] = |B[k]|^2, \quad k = 0, 1, \dots, N - 1. \quad (17)$$

308 *3.3.6. Complex Numbers Representation*

309 This numerical mapping has the advantage of better translating some of
 310 the nucleotides features into mathematical properties [80] and represents the
 311 complementary nature of AT and CG pairs [72]; e.g., let $\mathbf{s} = (G, A, G, A,$
 312 $G, T, G, A, C, C, A)$, thus, $\bar{\mathbf{r}} = (-1 - j, 1 + j, -1 - j, 1 + j, -1 - j, 1 - j,$
 313 $-1 - j, 1 + j, -1 + j, -1 + j, 1 + j)$, as shown in Equation (18). The DFT
 314 and power spectrum of this representation are presented in Equation (19).

$$\bar{r}[n] = \begin{cases} -1 - j, & s[n] = G \\ 1 + j, & s[n] = A \\ -1 + j, & s[n] = C \\ 1 - j, & s[n] = T \end{cases}, \quad n = 0, 1, \dots, N - 1. \quad (18)$$

$$\bar{R}[k] = \sum_{n=0}^{N-1} \bar{r}[n]e^{-j\frac{2\pi}{N}kn}, \quad P_{\bar{R}}[k] = |\bar{R}[k]|^2, \quad k = 0, 1, \dots, N - 1. \quad (19)$$

315 *3.3.7. Features*

316 The feature extraction is applied in each representation with Fourier
 317 transform, adopting Peak to Average Power Ratio (PAPR), mistakenly con-
 318 fused with the Signal to Noise Ratio (SNR), average power spectrum, me-
 319 dian, maximum, minimum, sample standard deviation, population stan-
 320 dard deviation, percentile (15/25/50/75), amplitude, variance, interquartile

321 range, semi-interquartile range, coefficient of variation, skewness, and kurto-
322 sis. Since according to [83] the RNA has a statistical phenomenon known as
323 period-3 behavior or 3-base periodicity, where the peak power will always be
324 at the sample $N/3$. Nevertheless, the PAPR is defined as [84]:

$$PAPR = \frac{\max_{0 \leq k \leq N-1} (P[k])}{\frac{1}{N} \sum_{k=0}^{N-1} P[k]} \quad (20)$$

325 3.4. Entropy

326 Information theory has been widely used in bioinformatics [85, 86]. Based
327 on this, we consider the study of [87], which applied an algorithmic and math-
328 ematical approach to DNA code analysis using entropy and phase plane.
329 Fundamentally, according to [86], entropy is a measure of the uncertainty
330 associated with a probabilistic experiment. To generate a probabilistic ex-
331 periment, we use a known method in bioinformatics, the k-mer (our pipeline
332 is shown in Figure 4).

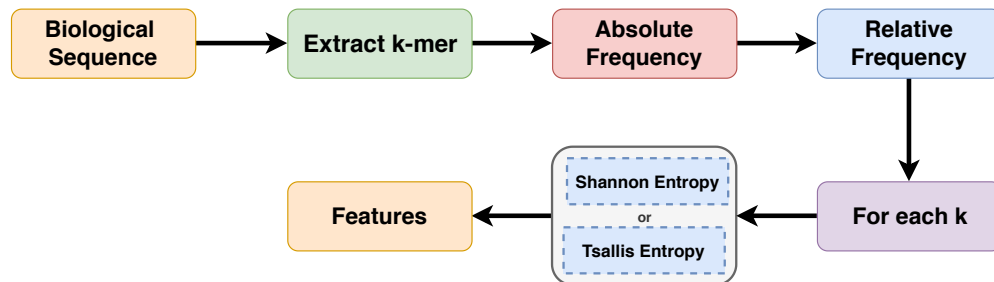


Figure 4: Entropy Pipeline. (1) Each sequence is mapped in k -mers; (2) The absolute frequency of each k is calculated; (3) Based on absolute frequency, the relative frequency is calculated; (4) The Tsallis or Shannon entropy is applied to each k ; Finally, (5) features are generated.

333 In this method, each sequence is mapped in the frequency of neighboring
334 bases k , generating statistical information. The k -mer is denoted in this work
335 by P_k , corresponding to Equation (21).

$$P_k(\mathbf{s}) = \frac{c_i^k}{N - k + 1} = \left(\frac{c_1^1}{N - 1 + 1}, \dots, \frac{c_4^1}{N - 1 + 1}, \frac{c_{4+1}^2}{N - 2 + 1}, \dots, \frac{c_i^k}{N - k + 1} \right) \quad k = 1, 2, \dots, 24. \quad (21)$$

336 We applied this equation to each sequence with frequencies of $k = 1, 2,$
 337 $\dots, 24$. Where, c_i^k is the number of substring occurrences with length k in a
 338 sequence (\mathbf{s}) with length N , in which the index $i \in \{1, 2, \dots, 4^1 + \dots + 4^k\}$
 339 represents the analyzed substring. For a better understanding, Figure 5
 340 demonstrated an example with $k = 6$ and $k = 9$.

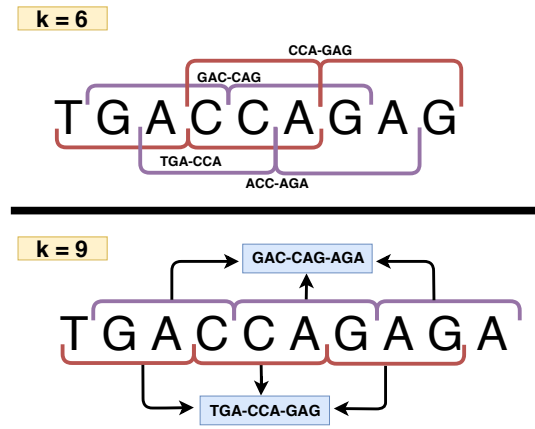


Figure 5: k -mer Workflow. Example with $k = 6$ and $k = 9$.

341 Basically, histograms with short bins are adopted, such as $[\{A\}, \{C\},$
 342 $\{G\}, \{T\}]$, that occur for $k = 1$, up to histograms with long sequence counting
 343 bins such as $[\{GGGGGGGGGGGG\}, \dots, \{AAAAAAAAAAAA\}]$, that
 344 result for $k = 12$. Where, after counting the absolute frequencies of each k ,
 345 we generate relative frequencies (see Equation (21)), and then apply Shannon
 346 and Tsallis entropy to generate the features.

347 3.4.1. Shannon and Tsallis Entropy

348 Fundamentally, we chose Shannon entropy, because it quantifies the amount
 349 of information in a variable [88], that is, we can reach a single value that

350 quantifies the information contained in different observation periods (e.g.,
351 our case: k -mer). However, according to [89], it is important to explore a
352 generalized form of the Shannon's entropy. Based on this, we have opted for
353 a generalized entropy proposed by Tsallis, applied by several works in the
354 literature [90, 91]. Thereby, for a discrete random variable F taking values in
355 $\{f[0], f[1], f[2], \dots, f[N-1]\}$ with probabilities $\{p[0], p[1], p[2], \dots, p[N-1]\}$,
356 represented as $P(F = f[n]) = p[n]$. The Shannon (Equation 22) and Tsallis
357 (Equation 23) entropy associated with this variable is given by the following
358 expressions:

$$H_S[k] = - \sum_{n=0}^{N-1} p_k[n] \log_2 p_k[n] \quad k = 1, 2, \dots, 24. \quad (22)$$

$$H_T[k] = \frac{1}{q-1} \left(1 - \sum_{n=0}^{N-1} p_k[n]^q \right) \quad k = 1, 2, \dots, 24. \quad (23)$$

359 Where k represents the analyzed k -mer, N the number of possible events
360 and $p[n]$ the probability that event n occurs.

361 3.5. Complex Networks

362 Complex networks are widely used in mathematical modeling and have
363 been an extremely active field in recent years [92], as well as becoming an
364 ideal research area for mathematicians, computer scientists, and biologists.
365 Based on this, we consider the study of [59], in which we propose a feature
366 extraction model based on complex networks, as shown in Figure 6.

367 Each sequence is mapped to the frequency of neighboring bases k ($k = 3$
368 - see Figure 5). This mapping is converted into an undirected graph repre-
369 sented by an adjacency matrix, in which we applied a threshold scheme for
370 feature extraction, thus generating our characteristic vector. Fundamentally,
371 we represent our structure by undirected weighted graphs. According to [92],
372 a graph $G = \{V, E\}$ is structured by a set V of vertices (or nodes) connected
373 by a set E of edges (or links). Each edge reflects a link between two vertices,
374 e.g., $e_p = (i, j)$ connection between the vertices i and j [92]. If there is an
375 edge connecting the vertices i and j , the elements a_{ij} are equal to 1, and
376 equal to 0 otherwise.

377 In our case, the graph is undirected, that is, the adjacency matrix A is
378 symmetric, e.g., elements $a_{ij} = a_{ji}$ for any i and j [92]. Furthermore, we
379 apply a threshold scheme presented by [59], in which we extract weight of

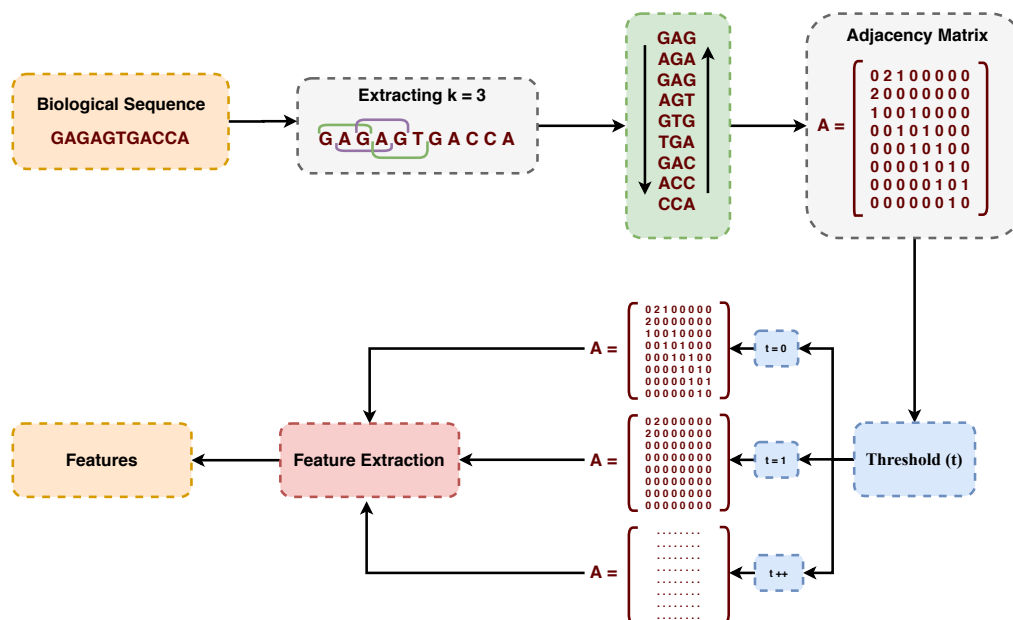


Figure 6: Complex Networks Pipeline. (1) Each sequence is mapped in the frequency of neighboring bases k ($k = 3$); (2) This mapping is converted to a undirected graph represented by an adjacency matrix; (3) Feature extraction is performed using a threshold scheme; Finally, (4) the features are generated.

380 the edges to capture adjacencies at different frequencies. Finally, as fea-
 381 tures, several network characterization measures were obtained, based on
 382 [59, 93], among them: Betweenness, assortativity, average degree, average
 383 path length, minimum degree, maximum degree, degree standard deviation,
 384 frequency of motifs (size 3 and 4), clustering coefficient.

385 3.6. Normalization, Training and Evaluation Metrics

386 Data normalization is a preprocessing technique often applied to a dataset.
 387 Essentially, features can have different dynamic ranges. This problem may
 388 have a stronger effect in the induction of a predictive model, mainly for
 389 distance-based ML algorithms. Consequently, the application of a normal-
 390 ization procedure makes the ranges similar, reducing this problem [94]. We
 391 used the min-max normalization, which reduces the data range to 0 and 1
 392 (or -1 to 1, if there are negative values) [2]. The general formula is given as
 393 (Equation (24)) [95]:

$$x'_{ij} = \frac{x_{ij} - \min(j)}{\max(j) - \min(j)}. \quad (24)$$

394 Where x is the original value and x'_{ij} is its normalized version. Further-
 395 more, $\min(j)$ and $\max(j)$ are, respectively, the smallest and largest values of
 396 a feature j [6, 95]. Next, we investigate three classification algorithms, such
 397 as Random Forest (RF) [96], AdaBoost [97] and CatBoost [98]. We chose
 398 these ML algorithms because they induce interpretable predictive models
 399 when humans can easily understand the internal decision-making process.
 400 Thus, domain experts can validate the knowledge used by the models for
 401 the classification of new sequences [6]. Finally, to induce our models, we
 402 used 70% of samples for *training* (with 10-fold cross-validation) and 30% for
 403 *testing*, as shown in Table 2.

Table 2: Number of sequences used for training and testing in each dataset.

Case Study	Dataset	Samples	Training	Testing
I	<i>A. trichopoda</i>	9112	6378	2734
	<i>A. thaliana</i>	5080	3556	1524
	<i>C. sinensis</i>	4430	3101	1329
	<i>C. sativus</i>	3460	2422	1038
	<i>R. communis</i>	6974	4881	2093
II	<i>lncRNA vs. snRNA</i>	2582	1807	775
	<i>lncRNA vs. Antisense</i>	114	79	35
	<i>circRNA vs. lncRNA</i>	5080	3556	1524

404 The methods were evaluated with four measures: Sensitivity (SE - Equa-
 405 tion 26), Specificity (SPC - Equation 27), Accuracy (ACC - Equation 25),
 406 and Cohen's kappa coefficient [99] (Equation 28).

$$ACC = \frac{TP + TN}{TN + FP + TP + FN} \quad (25) \quad SPC = \frac{TN}{TN + FP} \quad (27)$$

$$SE = \frac{TP}{TP + FN} \quad (26) \quad Kappa = \frac{p_o - p_e}{1 - p_e} \quad (28)$$

407 These measures use True Positive (TP), True Negative (TN), False Posi-
 408 tive (FP) and False Negative (FN) values, where: TP measures the correctly

409 predicted positive label; TN represents the correctly classified negative label;
 410 FP describes all those negative entities that are incorrectly classified as positive and;
 411 FN represents the positive label that are incorrectly classified as the negative label.
 412

413 4. Results

414 This section shows experimental results from 9 feature extraction approaches with mathematical models for biological sequences, divided into
 415 two parts: Case Study I and Case Study II.
 416

417 4.1. Case Study I

418 Initially, we induced models with the RF, AdaBoost, and CatBoost classifiers in the training set of three datasets (*A. trichopoda*, *A. thaliana*, and
 419 *R. communis*). Our initial goal is to choose the best classifier to follow in the testing phases. Thereby, to estimate the real accuracy, we applied 10-fold
 420 cross-validation. Thereby, to estimate the real accuracy, we applied 10-fold
 421 cross-validation, as shown in Table 3.
 422

Table 3: Accuracy for the training set (*A. trichopoda*, *A. thaliana*, and *R. communis*) using 10-fold cross-validation.

Dataset	Model	RF	AdaBoost	CatBoost
<i>A. trichopoda</i>	Z-curve	0.90 (\pm 0.03)	0.91 (\pm 0.02)	0.92 (\pm 0.02)
	Binary	0.92 (\pm 0.02)	0.94 (\pm 0.02)	0.94 (\pm 0.02)
	Real	0.91 (\pm 0.02)	0.93 (\pm 0.02)	0.94 (\pm 0.02)
	Integer	0.91 (\pm 0.02)	0.93 (\pm 0.02)	0.94 (\pm 0.02)
	EIIP	0.92 (\pm 0.02)	0.94 (\pm 0.02)	0.94 (\pm 0.02)
	Complex	0.92 (\pm 0.03)	0.94 (\pm 0.02)	0.94 (\pm 0.02)
	Graphs	0.92 (\pm 0.02)	0.94 (\pm 0.02)	0.94 (\pm 0.02)
	Shannon	0.92 (\pm 0.02)	0.94 (\pm 0.02)	0.94 (\pm 0.02)
	Tsallis	0.92 (\pm 0.02)	0.94 (\pm 0.02)	0.94 (\pm 0.02)
<i>A. thaliana</i>	Z-curve	0.95 (\pm 0.02)	0.93 (\pm 0.03)	0.94 (\pm 0.02)
	Binary	0.94 (\pm 0.02)	0.94 (\pm 0.02)	0.94 (\pm 0.02)
	Real	0.95 (\pm 0.02)	0.94 (\pm 0.02)	0.95 (\pm 0.02)
	Integer	0.94 (\pm 0.02)	0.94 (\pm 0.02)	0.94 (\pm 0.02)
	EIIP	0.95 (\pm 0.02)	0.94 (\pm 0.02)	0.95 (\pm 0.03)
	Complex	0.94 (\pm 0.02)	0.94 (\pm 0.02)	0.94 (\pm 0.01)
	Graphs	0.94 (\pm 0.02)	0.94 (\pm 0.02)	0.95 (\pm 0.02)

	Shannon	0.94 (\pm 0.02)	0.94 (\pm 0.02)	0.95 (\pm 0.02)
	Tsallis	0.94 (\pm 0.02)	0.94 (\pm 0.02)	0.94 (\pm 0.02)
<i>R. communis</i>	Z-curve	0.93 (\pm 0.02)	0.92 (\pm 0.02)	0.93 (\pm 0.02)
	Binary	0.95 (\pm 0.01)	0.95 (\pm 0.02)	0.95 (\pm 0.02)
	Real	0.95 (\pm 0.02)	0.94 (\pm 0.02)	0.94 (\pm 0.02)
	Integer	0.94 (\pm 0.01)	0.94 (\pm 0.01)	0.94 (\pm 0.02)
	EIIP	0.95 (\pm 0.02)	0.95 (\pm 0.02)	0.95 (\pm 0.01)
	Complex	0.95 (\pm 0.02)	0.95 (\pm 0.01)	0.95 (\pm 0.01)
	Graphs	0.95 (\pm 0.01)	0.95 (\pm 0.01)	0.95 (\pm 0.02)
	Shannon	0.95 (\pm 0.02)	0.95 (\pm 0.02)	0.95 (\pm 0.01)
	Tsallis	0.95 (\pm 0.01)	0.95 (\pm 0.01)	0.95 (\pm 0.01)

423 Assessing each classifier, we noted that the best performance was of the
 424 CatBoost with all mathematical models in *A. trichopoda*, followed by Ad-
 425 aBoost (6 best results) and RF (no better results). In *A. thaliana*, CatBoost
 426 kept the best performance (7 best results), followed by RF (6 best results)
 427 and AdaBoost (3 best results). In contrast, the RF classifier obtained the
 428 best results (6) in *R. communis*, followed by CatBoost (5 best results) and
 429 AdaBoost (3 best results). Based on this, we continued testing the models
 430 with the CatBoost classifier. Thus, in Table 4, we present the results of all
 431 mathematical models using 4 evaluation metrics.

Table 4: Performance analysis. This table compares the sensitivity, specificity, accuracy and kappa metrics for each model in the test sets using CatBoost classifier.

Dataset	Model	SE	SPC	ACC	Kappa
<i>A. trichopoda</i>	Z-curve	0.9744	0.8566	0.9155	0.8310
	Binary	0.9795	0.9005	0.9400	0.8800
	Real	0.9802	0.8837	0.9320	0.8639
	Integer	0.9773	0.8822	0.9298	0.8595
	EIIP	0.9781	0.8990	0.9386	0.8771
	Complex	0.9802	0.9012	0.9407	0.8815
	Graphs	0.9737	0.9020	0.9378	0.8756
	Shannon	0.9781	0.9020	0.9400	0.8800
	Tsallis	0.9795	0.9005	0.9400	0.8800
		Z-curve	0.9777	0.9383	0.9580
	Binary	0.9619	0.9449	0.9534	0.9068

<i>A. thaliana</i>	Real	0.9803	0.9409	0.9606	0.9213
	Integer	0.9698	0.9436	0.9567	0.9134
	EIIP	0.9646	0.9449	0.9547	0.9094
	Complex	0.9724	0.9409	0.9567	0.9134
	Graphs	0.9685	0.9423	0.9554	0.9108
	Shannon	0.9738	0.9462	0.9600	0.9200
	Tsallis	0.9764	0.9409	0.9587	0.9173
	<i>C. sinensis</i>	Z-curve	0.9021	0.8707	0.8864
Binary		0.8901	0.8707	0.8804	0.7607
Real		0.9142	0.8571	0.8856	0.7713
Integer		0.8825	0.8692	0.8758	0.7517
EIIP		0.8840	0.8526	0.8683	0.7367
Complex		0.9081	0.8496	0.8789	0.7577
Graphs		0.9006	0.8632	0.8819	0.7637
Shannon		0.9172	0.8586	0.8879	0.7758
Tsallis		0.9262	0.8541	0.8901	0.7803
<i>C. sativus</i>	Z-curve	0.8979	0.8478	0.8728	0.7457
	Binary	0.9056	0.8459	0.8757	0.7514
	Real	0.9268	0.8439	0.8854	0.7707
	Integer	0.9056	0.8536	0.8796	0.7592
	EIIP	0.8979	0.8459	0.8719	0.7437
	Complex	0.9326	0.8343	0.8834	0.7669
	Graphs	0.9075	0.8536	0.8805	0.7611
	Shannon	0.9326	0.8382	0.8854	0.7707
	Tsallis	0.9403	0.8401	0.8902	0.7803
	<i>R. communis</i>	Z-curve	0.9446	0.9140	0.9293
Binary		0.9417	0.9589	0.9503	0.9006
Real		0.9589	0.9408	0.9498	0.8997
Integer		0.9465	0.9456	0.9460	0.8920
EIIP		0.9455	0.9551	0.9503	0.9006
Complex		0.9398	0.9561	0.9479	0.8958
Graphs		0.9455	0.9542	0.9498	0.8997
Shannon		0.9388	0.9589	0.9489	0.8978
Tsallis		0.9417	0.9608	0.9513	0.9025

432 As can be seen, all models presented excellent results, with the worst per-
 433 formance (ACC) of 0.8901 (*C. sinensis*) and the best of 0.9606 (*A. thaliana*).
 434 That is, all models were robust in different datasets without a high loss of

435 performance. Assessing each metric individually, we realized that in SE,
436 the best performance was from Real representation (3 datasets), followed by
437 Tsallis (2 datasets) and Complex numbers (1 dataset). In SPC, the best
438 results were from Entropy (3 datasets), followed by Graphs (2 datasets).
439 In ACC, Tsallis presented the best performance (3 datasets), followed by
440 Real representation and Complex numbers (1 dataset). For each dataset, we
441 can see in *A. trichopoda* the best ACC was 0.9407 (Complex); *A. thaliana*
442 with 0.9606 (Real); *C. sinensis* with 0.8901 (Tsallis); *C. sativus* with 0.8902
443 (Tsallis); and *R. communis* with 0.9513 (Tsallis).

444 4.2. Case Study II

445 After evaluating all methods in 5 different datasets (lncRNA of different
446 species) and observing their results, we applied a second case study, where we
447 used only three mathematical models for generalization analysis, including
448 GSP (Fourier + complex numbers), entropy (Tsallis) and graphs (complex
449 networks). Here, our objective was to analyze how each model behaved in
450 different biological sequence classification problems. For this, we tested 3
451 new datasets established in Section 3.1.2, as can be seen in Figure 7.

452 Again, all showed robust results, in which, graph-based models are the
453 best in 2 of the 3 problems analyzed, followed by entropy and GSP. In the
454 first three datasets, our methods achieved excellent accuracy. Furthermore, if
455 we analyze at the last problem (circRNA vs. lncRNA), our approaches were
456 effective when compared to our references that reached an ACC of 0.7780 [66]
457 and 0.7890 [21] in their datasets against 0.8307 from our best model (graph
458 - using these comparisons as an (indirect) reference indicator).

459 4.3. Statistical Significance Tests

460 The statistical significance was assessed in both case studies (difference
461 in ACC), using Friedman's statistical test and the Conover post-hoc test.
462 Thereby, our null hypothesis ($H_0 = M(1) = M(2) = \dots = M(k)$), is tested
463 against the alternative hypothesis ($H_A =$ at least one model has statistical
464 significance ($\alpha = 0.05$, $p < \alpha$)). First, we apply the global test in the case
465 study I, in which the Friedman test indicates significance ($\chi^2(8) = 17.34$, p -
466 value = 0.0268), that is, we can reject H_0 , as $p < 0.05$. Thus, it is essential to
467 execute the post-hoc statistical test. Conover statistics values were obtained,
468 as well as p -values (see Table 5), using 95% of significance ($\alpha = 0.05$).

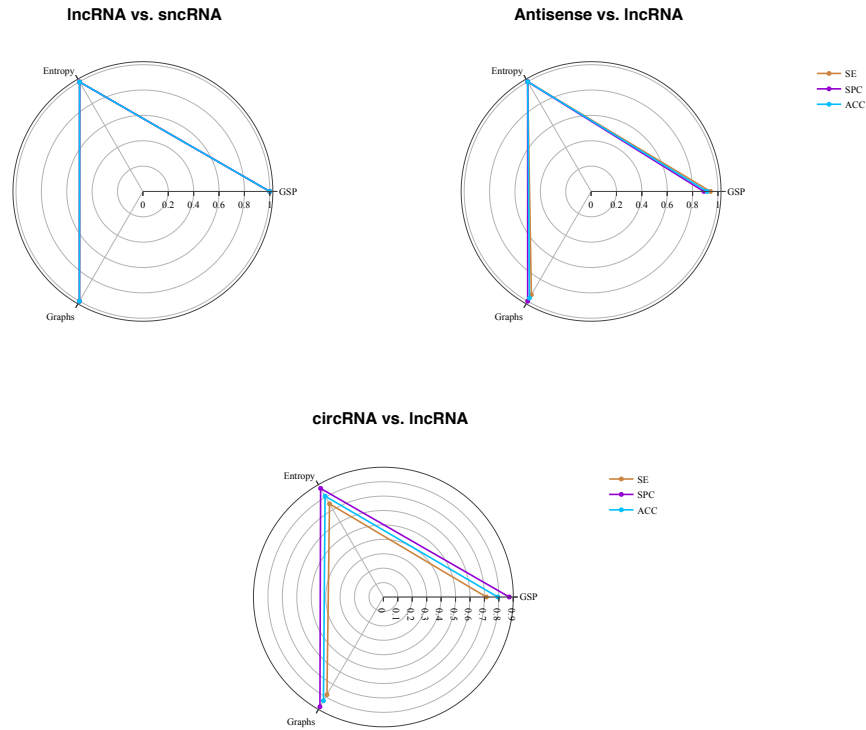


Figure 7: Performance analysis of three mathematical models, GSP (fourier + complex numbers), entropy (Tsallis) and graphs (complex networks), for different problems.

Table 5: Conover statistics values - The accepted alternative hypothesis is in bold (p -values for $\alpha = 0.05$).

	Z-curve	Binary	Real	Integer	EIIP	Complex	Graphs	Shannon
Binary	0.5580	-	-	-	-	-	-	-
Real	0.1416	0.3671	-	-	-	-	-	-
Integer	0.7896	0.3956	0.0852	-	-	-	-	-
EIIP	0.9574	0.5230	0.1284	0.8309	-	-	-	-
Complex	0.3671	0.7489	0.5580	0.2451	0.3399	-	-	-
Graphs	0.5580	1.0000	0.3671	0.3956	0.5230	0.7489	-	-
Shannon	0.0687	0.2057	0.7089	0.0390	0.0616	0.3399	0.2057	-
Tsallis	0.0146	0.0550	0.2898	0.0075	0.0128	0.1050	0.0550	0.4892

469 Concerning to the Conover post-hoc test, entropy-based models have
 470 highly significant differences for the Z-curve ($p < 0.0146$), Integer ($p < 0.0075$)

471 - Tsallis and $p < 0.0390$ - Shannon), and EIIP ($p < 0.0128$). Possibly,
472 these results indicate that entropy has a more significant performance when
473 compared to representations with Fourier. However, other mathematical
474 models in case study I do not differ significantly, indicating their efficiency
475 in all datasets. Now, evaluating case study II, we realized that the global
476 test with Friedman's statistical test is not significant, in which we obtained
477 $\chi^2(2) = 1.64$, p -value = 0.4412, indicating that the three studied feature ex-
478 traction techniques show a similar performance in the problems, once more
479 confirming the effectiveness and robustness of all mathematical models.

480 4.4. Computational Time

481 In addition, we also assessed the computational time cost of each tested
482 model. To do this, we ran three models, GSP (Fourier + complex numbers),
483 entropy (Tsallis) and graphs (complex networks), in 1291 random sequences,
484 as shown in Figure 8.

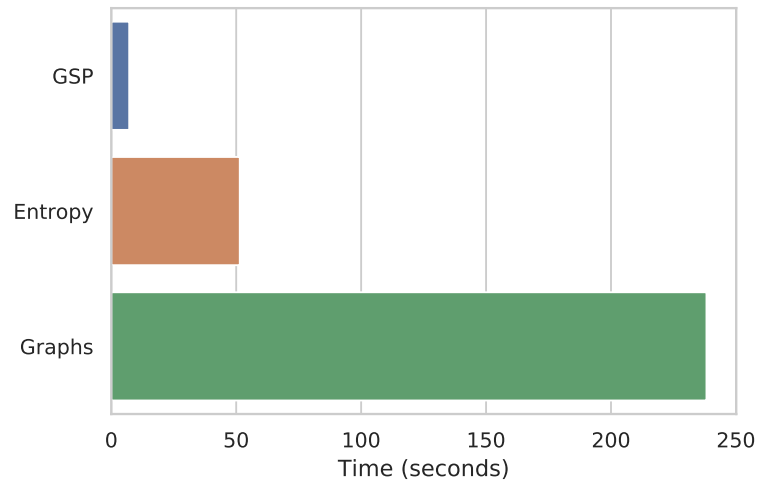


Figure 8: Execution Time.

485 We performed the experiments using Intel Core i3-9100F CPU (3.60GHz),
486 16GB memory, and running in Debian GNU/Linux 10. The lowest cost in
487 computational time is for models based on GSP (0m7.183s) and entropy
488 (0m51.427s), while graphs (3m58.208s) have a much higher cost. These re-
489 sults demonstrated that, although the models present a similar performance,
490 the computational time efficiency is significantly different.

491 5. Discussion

492 This section discusses our findings in terms of whether they support our
493 hypothesis (*feature extraction approaches based on mathematical models are*
494 *as efficient and generalist as biological approaches*). Overall, several exper-
495 imental tests were assumed in this research, in which all feature extraction
496 approaches based on mathematical models showed excellent results, as can
497 be seen in Table 4 and Figure 7. Regarding its performance in distinct clas-
498 sification problems, case study II, we used only three mathematical models
499 for generalization analysis, including GSP (Fourier + complex numbers), en-
500 tropy (Tsallis) and graphs (complex networks). In which, entropy and graph-
501 based models reported the best performance followed by GSP. Furthermore,
502 all models maintained robust results in different sequence classification prob-
503 lems.

504 Furthermore, to fully support our hypothesis, we also compare three
505 mathematical models shown in Figure 7 concerning a biological and hybrid
506 approach, in four datasets ((lncRNA vs. mRNA (case study I)); (lncRNA vs.
507 snRNA; lncRNA vs. Antisense; circRNA vs. lncRNA (case study II)). Thus,
508 we generate our biological model using some of the most applied features in
509 Figure 1. Thus, features used by the models are:

- 510 • **Biological:** The features were provided by [19]: Fickett TESTCODE
511 score, isoelectric point, open reading frame (ORF) length, and ORF
512 integrity.
- 513 • **Hybrid:** The features were generated by one of the most current ap-
514 proaches in the literature (lncFinder [20] - 2018). We classify this model
515 as a hybrid because it uses a combination of biological and mathemati-
516 cal features. Among the biological characteristics is Logarithm-distance
517 of hexamer on ORF, length and coverage of the longest ORF. Regard-
518 ing mathematical features, [20] uses an EIIP-based physicochemical
519 property with Fourier Transform (similar to our approach with GSP,
520 but using only EIIP mapping).

521 For a fair comparison, the new experiments follow the same methodology
522 (70% training, 30% test, and CatBoost classifier), as shown in Table 6.

523 As can be seen, the hybrid model (0.9915) reported the best performance
524 in the first dataset (lncRNA vs. mRNA), followed by the biological (0.9816)
525 and our mathematical model (Entropy - 0.9587), with only a difference of

Table 6: Performance analysis of three mathematical models against a biological and hybrid model for different sequence classification problems.

lncRNA vs. mRNA				lncRNA vs. snRNA			
Models	SE	SPC	ACC	Models	SE	SPC	ACC
GSP	0.9724	0.9409	0.9567	GSP	1.0000	1.0000	1.0000
Entropy	0.9764	0.9409	0.9587	Entropy	0.9974	0.9974	0.9974
Graphs	0.9685	0.9423	0.9554	Graphs	1.0000	1.0000	1.0000
Biological	0.9869	0.9764	0.9816	Biological	0.7855	0.8273	0.8065
Hybrid	0.9895	0.9934	0.9915	Hybrid	0.9509	0.9485	0.9497

lncRNA vs. Antisense				circRNA vs. lncRNA			
Models	SE	SPC	ACC	Models	SE	SPC	ACC
GSP	0.9412	0.8889	0.9143	GSP	0.7139	0.8727	0.7933
Entropy	1.0000	1.0000	1.0000	Entropy	0.7467	0.8701	0.8084
Graphs	0.9412	1.0000	0.9714	Graphs	0.7822	0.8793	0.8307
Biological	0.8889	0.9412	0.9143	Biological	0.6024	0.7612	0.6818
Hybrid	0.9412	0.7778	0.8571	Hybrid	0.7283	0.8819	0.8051

526 0.0328 and 0.0229, respectively. However, it is relevant to highlight that
527 the biological and hybrid models use the ORF descriptor, a highly employed
528 feature for discovering coding sequences and which, according to [15, 6] is an
529 essential guideline for distinguishing lncRNAs from mRNA. In other words,
530 this explains the great result, but, as mentioned at the beginning of this
531 manuscript, this type of feature with a biological insight is often difficult
532 to reuse or adapt to another specific problem. Thereby, our study has an
533 gain in terms of generalization, since this would not be possible only with
534 the ORF. If we analyze at the hybrid model, in this first dataset, the gain
535 was minimal compared to the biological (0.0099), which again confirms the
536 efficiency of the previously mentioned features. This is different from our
537 approaches, which showed an excellent result without using bias features for
538 the analyzed problem.

539 Consequently, this hypothesis is proven in the other three datasets, where
540 our mathematical models perform much better than the biological model,
541 mainly in the fourth dataset (circRNA vs. lncRNA), in which we obtained
542 a gain of 0.1489 in ACC. Regarding the hybrid model, it can be observed
543 that the mixture of biological and mathematical characteristics helped to

544 keep the model competitive in all datasets, indicating the effectiveness of
545 mathematical features. Even so, our models showed the best results in three
546 of the four proposed problems. Therefore, our pipeline is efficient in terms
547 of generalization to classify lncRNA from mRNA, as well as other biological
548 sequence classification problems. We also assessed the statistical significance
549 of the mathematical versus biological approach in the previously applied
550 tests, in which entropy ($p < 0.0480$) and graphs ($p < 0.0200$) indicated
551 significant results concerning the biological model. Lastly, considering all
552 these findings, we fully support the suggested hypothesis.

553 **6. Conclusion**

554 This work proposed to analyze feature extraction approaches for biolog-
555 ical sequence classification. Specifically, we concentrated our work on the
556 study of feature extraction techniques using mathematical models. We ana-
557 lyzed mathematical models to propose efficient and generalist techniques for
558 different problems. As a case study, we used lncRNA sequences. Moreover,
559 we divided this paper into two case studies. In our experiments, as a start-
560 ing point, 9 mathematical models for feature extraction were analyzed: 6
561 numerical mapping techniques with Fourier transform; Tsallis and Shannon
562 entropy; Graphs (complex networks). Thereby, several biological sequence
563 classification problems were adopted to validate the proposed approach.

564 As a result, all models presented excellent results, with performances
565 (ACC) between 0.8901-0.9606 in case study I. In case study II, once more,
566 all showed excellent results with models based on entropy and graphs show-
567 ing the best performance, followed by GSP. Furthermore, to validate our
568 study, we compared the performance of three mathematical models against
569 a biological and hybrid approach, in four different datasets. In which, our
570 models demonstrated suitable results, and was superior or competitive and
571 robust in terms of generalization. In our experiments, we verified that math-
572 ematical approaches perform as accurately as biological approaches and have
573 a better generalization capacity since they outperform biological features in
574 scenarios not designed for them. Finally, among the different mathematical
575 models tested in this work, the combination of k-mer and entropy, as well
576 as graph-based models performs better than GSP at the cost of a significant
577 increase in computational complexity.

578 **Declaration of Competing interests**

579 All authors declare that they have no conflict of interest.

580 **Financial support**

581 This project has been supported by a master scholarship from Federal
582 University of Technology - Paraná (UTFPR) (Grant: April/2018) and CAPES
583 (April/2019 and PROEX-11919694/D).

584 **Acknowledgements**

585 The authors would like to thank UTFPR-CP, ICMC-USP, and CAPES
586 for the financial support given to this research.

References

- [1] H. Lou, M. Schwartz, J. Bruck, F. Farnoud, Evolution of k-mer frequencies and entropy in duplication and substitution mutation systems, *IEEE Transactions on Information Theory* (2019).
- [2] R. P. Bonidia, L. D. H. Sampaio, F. M. Lopes, D. S. Sanches, Feature extraction of long non-coding rnas: A fourier and numerical mapping approach, in: I. Nyström, Y. Hernández Heredia, V. Milián Núñez (Eds.), *Progress in Pattern Recognition, Image Analysis, Computer Vision, and Applications*, Springer International Publishing, Cham, 2019, pp. 469–479.
- [3] R. Min, *Machine Learning Approaches to Biological Sequence and Phenotype Data Analysis*, University of Toronto, 2010.
- [4] M.-R. Cao, Z.-P. Han, J.-M. Liu, Y.-G. Li, Y.-B. Lv, J.-B. Zhou, J.-H. He, Bioinformatic analysis and prediction of the function and regulatory network of long non-coding rnas in hepatocellular carcinoma, *Oncology letters* 15 (5) (2018) 7783–7793.
- [5] W. J. d. S. Diniz, F. Canduri, *Bioinformatics: an overview and its applications*, *Genet Mol Res* 16 (1) (2017).

- [6] R. Parmezan Bonidia, A. C. Ponce de Leon Ferreira de Carvalho, A. Rossi Paschoal, D. Sipoli Sanches, Selecting the most relevant features for the identification of long non-coding rnas in plants, in: 2019 8th Brazilian Conference on Intelligent Systems (BRACIS), 2019, pp. 539–544. doi:10.1109/BRACIS.2019.00100.
- [7] V. I. Jurtz, A. R. Johansen, M. Nielsen, J. J. Almagro Armenteros, H. Nielsen, C. K. Sønderby, O. Winther, S. K. Sønderby, An introduction to deep learning on biological sequence data: examples and solutions, *Bioinformatics* 33 (22) (2017) 3685–3690. doi:10.1093/bioinformatics/btx531.
URL <https://doi.org/10.1093/bioinformatics/btx531>
- [8] M. E. Maros, D. Capper, D. T. Jones, V. Hovestadt, A. von Deimling, S. M. Pfister, A. Benner, M. Zucknick, M. Sill, Machine learning workflows to estimate class probabilities for precision cancer diagnostics on dna methylation microarray data, *Nature Protocols* (2020) 1–34.
- [9] J. Li, W. Liu, Puzzle of highly pathogenic human coronaviruses (2019-ncov), *Protein & Cell* (2020) 1–4.
- [10] D. Benvenuto, M. Giovanetti, A. Ciccozzi, S. Spoto, S. Angeletti, M. Ciccozzi, The 2019-new coronavirus epidemic: Evidence for virus evolution, *Journal of Medical Virology* 92 (4) (2020) 455–459. arXiv:<https://onlinelibrary.wiley.com/doi/pdf/10.1002/jmv.25688>, doi:10.1002/jmv.25688.
URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/jmv.25688>
- [11] C. Xu, S. A. Jackson, BioSeq-Analysis: a platform for DNA, RNA and protein sequence analysis based on machine learning approaches, *Genome Biology* 20 (2019) 1–4. doi:<https://doi.org/10.1186/s13059-019-1689-0>.
URL <https://doi.org/10.1186/s13059-019-1689-0>
- [12] D. Storcheus, A. Rostamizadeh, S. Kumar, A survey of modern questions and challenges in feature extraction, in: *Feature Extraction: Modern Questions and Challenges*, 2015, pp. 1–18.
- [13] I. Guyon, S. Gunn, M. Nikravesh, L. A. Zadeh, *Feature extraction: foundations and applications*, Vol. 207, Springer, 2008.

- [14] R. Saidi, S. Aridhi, E. M. Nguifo, M. Maddouri, Feature extraction in protein sequences classification: a new stability measure, in: *Proceedings of the ACM Conference on Bioinformatics, Computational Biology and Biomedicine*, ACM, 2012, pp. 683–689.
- [15] J. Baek, B. Lee, S. Kwon, S. Yoon, Incrnanet: Long non-coding rna identification using deep learning, *Bioinformatics* 1 (2018) 9.
- [16] R. Muhammod, S. Ahmed, D. Md Farid, S. Shatabda, A. Sharma, A. Dehzangi, PyFeat: a Python-based effective feature generation tool for DNA, RNA and protein sequences, *Bioinformatics* 35 (19) (2019) 3831–3833. arXiv:<http://oup.prod.sis.lan/bioinformatics/article-pdf/35/19/3831/30061688/btz165.pdf>, doi:10.1093/bioinformatics/btz165. URL <https://doi.org/10.1093/bioinformatics/btz165>
- [17] Q. Abbas, S. M. Raza, A. A. Biyabani, M. A. Jaffar, A review of computational methods for finding non-coding rna genes, *Genes* 7 (12) (2016) 113.
- [18] N. Amin, A. McGrath, Y.-P. P. Chen, Evaluation of deep learning in non-coding rna classification, *Nature Machine Intelligence* 1 (5) (2019) 246.
- [19] Y.-J. Kang, D.-C. Yang, L. Kong, M. Hou, Y.-Q. Meng, L. Wei, G. Gao, Cpc2: a fast and accurate coding potential calculator based on sequence intrinsic features, *Nucleic acids research* 45 (W1) (2017) W12–W16.
- [20] S. Han, Y. Liang, Q. Ma, Y. Xu, Y. Zhang, W. Du, C. Wang, Y. Li, Lncfinder: an integrated platform for long non-coding rna identification utilizing sequence intrinsic composition, structural information and physicochemical property, *Briefings in Bioinformatics* (2018).
- [21] L. Chen, Y.-H. Zhang, G. Huang, X. Pan, S. Wang, T. Huang, Y.-D. Cai, Discriminating cirrnas from other lncrnas using a hierarchical extreme learning machine (h-elm) algorithm with feature selection, *Molecular Genetics and Genomics* 293 (1) (2018) 137–149.
- [22] S. R. Eddy, Non-coding rna genes and the modern rna world, *Nature Reviews Genetics* 2 (12) (2001) 919.

- [23] P. Kapranov, J. Cheng, S. Dike, D. A. Nix, R. Dutttagupta, A. T. Willingham, P. F. Stadler, J. Hertel, J. Hackermüller, I. L. Hofacker, et al., Rna maps reveal new rna classes and a possible function for pervasive transcription, *Science* 316 (5830) (2007) 1484–1488.
- [24] Y. Zhang, Y. Tao, Q. Liao, Long noncoding rna: a crosslink in biological regulatory network, *Briefings in bioinformatics* (2017).
- [25] A. Li, Q. Zang, D. Sun, M. Wang, A text feature-based approach for literature mining of lncrna–protein interactions, *Neurocomputing* 206 (2016) 73–80.
- [26] Y. Wang, Y. Li, Q. Wang, Y. Lv, S. Wang, X. Chen, X. Yu, W. Jiang, X. Li, Computational identification of human long intergenic non-coding rnas using a ga–svm algorithm, *Gene* 533 (1) (2014) 94–99.
- [27] L. Wang, L. Kuang, S. Ye, M. F. B. Iqbal, T. Pei, et al., A novel method for lncrna-disease association prediction based on an lncrna-disease association network, *IEEE/ACM Transactions on Computational Biology and Bioinformatics* (2018).
- [28] W. Zhang, Q. Qu, Y. Zhang, W. Wang, The linear neighborhood propagation method for predicting long non-coding rna–protein interactions, *Neurocomputing* 273 (2018) 526–534.
- [29] Q.-Z. Zhou, B. Zhang, Q.-Y. Yu, Z. Zhang, Bmncrnadb: a comprehensive database of non-coding rnas in the silkworm, *bombyx mori*, *BMC bioinformatics* 17 (1) (2016) 370.
- [30] M. Q. Hassan, C. E. Tye, G. S. Stein, J. B. Lian, Non-coding rnas: Epigenetic regulators of bone development and homeostasis, *Bone* 81 (2015) 746–756.
- [31] C. Ciaudo, N. Servant, V. Cognat, A. Sarazin, E. Kieffer, S. Viville, V. Colot, E. Barillot, E. Heard, O. Voinnet, Highly dynamic and sex-specific expression of micrnas during early es cell differentiation, *PLoS genetics* 5 (8) (2009) e1000620.
- [32] X. Peng, L. Gralinski, C. D. Armour, M. T. Ferris, M. J. Thomas, S. Proll, B. G. Bradel-Tretheway, M. J. Korth, J. C. Castle, M. C. Biery, et al., Unique signatures of long noncoding rna expression in response to

- virus infection and altered innate immune signaling, *MBio* 1 (5) (2010) e00206–10.
- [33] C. Pastori, C. Wahlestedt, Involvement of long noncoding rnas in diseases affecting the central nervous system, *RNA biology* 9 (6) (2012) 860–870.
- [34] Q. Zhang, Y. Wei, Z. Yan, C. Wu, Z. Chang, Y. Zhu, K. Li, Y. Xu, The characteristic landscape of lncrnas classified by rbp–lncrna interactions across 10 cancers, *Molecular bioSystems* 13 (6) (2017) 1142–1151.
- [35] H.-L. V. Wang, J. A. Chekanova, Long noncoding rnas in plants, in: *Long Non Coding RNA Biology*, Springer, 2017, pp. 133–154.
- [36] C. Di, J. Yuan, Y. Wu, J. Li, H. Lin, L. Hu, T. Zhang, Y. Qi, M. B. Gerstein, Y. Guo, et al., Characterization of stress-responsive lncrnas in arabidopsis thaliana by integrating expression, epigenetic and structural features, *The Plant Journal* 80 (5) (2014) 848–861.
- [37] D. Wang, Z. Qu, L. Yang, Q. Zhang, Z.-H. Liu, T. Do, D. L. Adelson, Z.-Y. Wang, I. Searle, J.-K. Zhu, Transposable elements (te s) contribute to stress-related long intergenic noncoding rna s in plants, *The Plant Journal* 90 (1) (2017) 133–146.
- [38] Y.-C. Zhang, J.-Y. Liao, Z.-Y. Li, Y. Yu, J.-P. Zhang, Q.-F. Li, L.-H. Qu, W.-S. Shu, Y.-Q. Chen, Genome-wide screening and functional analysis identify a large number of long noncoding rnas involved in the sexual reproduction of rice, *Genome biology* 15 (12) (2014) 512.
- [39] Y. Fang, M. J. Fullwood, Roles, functions, and mechanisms of long non-coding rnas in cancer, *Genomics, proteomics & bioinformatics* 14 (1) (2016) 42–54.
- [40] T. Derrien, R. Johnson, G. Bussotti, A. Tanzer, S. Djebali, H. Tilgner, G. Guernec, D. Martin, A. Merkel, D. G. Knowles, et al., The gencode v7 catalog of human long noncoding rnas: analysis of their gene structure, evolution, and expression, *Genome research* 22 (9) (2012) 1775–1789.
- [41] J. Cheng, P. Kapranov, J. Drenkow, S. Dike, S. Brubaker, S. Patel, J. Long, D. Stern, H. Tammana, G. Helt, et al., Transcriptional maps of

- 10 human chromosomes at 5-nucleotide resolution, *Science* 308 (5725) (2005) 1149–1154.
- [42] L. Ma, V. B. Bajic, Z. Zhang, On the classification of long non-coding rnas, *RNA biology* 10 (6) (2013) 924–933.
- [43] R. Hu, X. Sun, Incrnatargets: a platform for lncrna target prediction based on nucleic acid thermodynamics, *Journal of bioinformatics and computational biology* 14 (04) (2016) 1650016.
- [44] S. Chooniedass-Kothari, E. Emberley, M. Hamedani, S. Troup, X. Wang, A. Czosnek, F. Hube, M. Mutawe, P. Watson, E. Leygue, The steroid receptor rna activator is the first functional rna encoding a protein, *FEBS letters* 566 (1-3) (2004) 43–47.
- [45] Y. He, X.-M. Meng, C. Huang, B.-M. Wu, L. Zhang, X.-W. Lv, J. Li, Long noncoding rnas: Novel insights into hepatocellular carcinoma, *Cancer letters* 344 (1) (2014) 20–27.
- [46] J. T. Kung, D. Colognori, J. T. Lee, Long noncoding rnas: past, present, and future, *Genetics* 193 (3) (2013) 651–669.
- [47] L. Kong, Y. Zhang, Z.-Q. Ye, X.-Q. Liu, S.-Q. Zhao, L. Wei, G. Gao, Cpc: assess the protein-coding potential of transcripts using sequence features and support vector machine, *Nucleic acids research* 35 (suppl.2) (2007) W345–W349.
- [48] L. Wang, H. J. Park, S. Dasari, S. Wang, J.-P. Kocher, W. Li, Cpat: Coding-potential assessment tool using an alignment-free logistic regression model, *Nucleic acids research* 41 (6) (2013) e74–e74.
- [49] L. Sun, H. Luo, D. Bu, G. Zhao, K. Yu, C. Zhang, Y. Liu, R. Chen, Y. Zhao, Utilizing sequence intrinsic composition to classify protein-coding and long non-coding transcripts, *Nucleic acids research* 41 (17) (2013) e166–e166.
- [50] A. Li, J. Zhang, Z. Zhou, Plek: a tool for predicting long non-coding rnas and messenger rnas based on an improved k-mer scheme, *BMC bioinformatics* 15 (1) (2014) 311.

- [51] X.-N. Fan, S.-W. Zhang, *lncrna-mfdl*: identification of human long non-coding rnas by fusing multiple features and using deep learning, *Molecular BioSystems* 11 (3) (2015) 892–897.
- [52] R. Achawanantakun, J. Chen, Y. Sun, Y. Zhang, *Lncrna-id*: Long non-coding rna identification using balanced random forests, *Bioinformatics* 31 (24) (2015) 3897–3905.
- [53] L. Sun, H. Liu, L. Zhang, J. Meng, *lncscan-svm*: a tool for predicting long non-coding rnas using support vector machine, *PloS one* 10 (10) (2015) e0139654.
- [54] C. Pian, G. Zhang, Z. Chen, Y. Chen, J. Zhang, T. Yang, L. Zhang, *Lncrnaped*: Classification of long non-coding rnas and protein-coding transcripts by the ensemble algorithm with a new hybrid feature, *PloS one* 11 (5) (2016) e0154567.
- [55] R. Tripathi, S. Patel, V. Kumari, P. Chakraborty, P. K. Varadwaj, *Deeplnc*, a long non-coding rna prediction tool using deep neural network, *Network Modeling Analysis in Health Informatics and Bioinformatics* 5 (1) (2016) 21.
- [56] L. M. Vieira, C. Grativol, F. Thiebaut, T. G. Carvalho, P. R. Hardoim, A. Hemerly, S. Lifschitz, P. C. G. Ferreira, M. E. M. Walter, *Plantrna_sniffer*: a svm-based workflow to predict long intergenic non-coding rnas in plants, *Non-coding RNA* 3 (1) (2017) 11.
- [57] U. Singh, N. Khemka, M. S. Rajkumar, R. Garg, M. Jain, *Plncpro* for prediction of long non-coding rnas (*lncrnas*) in plants and its application for discovery of abiotic stress-responsive *lncrnas* in rice and chickpea, *Nucleic acids research* 45 (22) (2017) e183–e183.
- [58] T. d. C. Negri, W. A. L. Alves, P. H. Bugatti, P. T. M. Saito, D. S. Domingues, A. R. Paschoal, Pattern recognition analysis on long non-coding rnas: a tool for prediction in plants, *Briefings in bioinformatics* (2018).
- [59] E. A. Ito, I. Katahira, F. F. d. R. Vicente, L. F. P. Pereira, F. M. Lopes, *Basinet*—biological sequences network: a case study on coding and non-coding rnas identification, *Nucleic acids research* (2018).

- [60] S. F. Altschul, T. L. Madden, A. A. Schäffer, J. Zhang, Z. Zhang, W. Miller, D. J. Lipman, Gapped blast and psi-blast: a new generation of protein database search programs, *Nucleic acids research* 25 (17) (1997) 3389–3402.
- [61] A. C. Liu, The effect of oversampling and undersampling on classifying imbalanced text datasets, The University of Texas at Austin (2004).
- [62] D. M. Goodstein, S. Shu, R. Howson, R. Neupane, R. D. Hayes, J. Fazo, T. Mitros, W. Dirks, U. Hellsten, N. Putnam, et al., Phytozome: a comparative platform for green plant genomics, *Nucleic acids research* 40 (D1) (2011) D1178–D1186.
- [63] A. Paytuví Gallart, A. Hermoso Pulido, I. Anzar Martínez de Lagrán, W. Sanseverino, R. Aiese Cigliano, Greenc: a wiki-based database of plant lncrnas, *Nucleic acids research* 44 (D1) (2015) D1161–D1166.
- [64] D. Chen, C. Yuan, J. Zhang, Z. Zhang, L. Bai, Y. Meng, L.-L. Chen, M. Chen, PlantNATsDB: a comprehensive database of plant natural antisense transcripts, *Nucleic Acids Research* 40 (D1) (2011) D1187–D1193. arXiv:<https://academic.oup.com/nar/article-pdf/40/D1/D1187/9481672/gkr823.pdf>, doi:10.1093/nar/gkr823. URL <https://doi.org/10.1093/nar/gkr823>
- [65] Q. Chu, X. Zhang, X. Zhu, C. Liu, L. Mao, C. Ye, Q.-H. Zhu, L. Fan, Plantcircbase: a database for plant circular rnas, *Molecular plant* 10 (8) (2017) 1126–1128.
- [66] X. Pan, K. Xiong, Predcircrna: computational classification of circular rna from other long non-coding rna using hybrid features, *Molecular Biosystems* 11 (8) (2015) 2219–2226.
- [67] C. Yin, Y. Chen, S. S.-T. Yau, A measure of dna sequence similarity by fourier transform with applications on hierarchical clustering, *Journal of theoretical biology* 359 (2014) 18–28.
- [68] C. Yin, S. S.-T. Yau, A fourier characteristic of coding sequences: origins and a non-fourier approximation, *Journal of computational biology* 12 (9) (2005) 1153–1165.

- [69] D. Anastassiou, Genomic signal processing, *IEEE signal processing magazine* 18 (4) (2001) 8–20.
- [70] L. Marsella, F. Sirocco, A. Trovato, F. Seno, S. C. Tosatto, Repetita: detection and discrimination of the periodicity of protein solenoid repeats by discrete fourier transform, *Bioinformatics* 25 (12) (2009) i289–i295.
- [71] W. T. Cochran, J. W. Cooley, D. L. Favin, H. D. Helms, R. A. Kaenel, W. W. Lang, G. C. Maling, D. E. Nelson, C. M. Rader, P. D. Welch, What is the fast fourier transform?, *Proceedings of the IEEE* 55 (10) (1967) 1664–1674.
- [72] M. Abo-Zahhad, S. M. Ahmed, S. A. Abd-Elrahman, Genomic analysis and classification of exon and intron sequences using dna numerical mapping techniques, *International Journal of Information Technology and Computer Science* 4 (8) (2012) 22–36.
- [73] G. Mendizabal-Ruiz, I. Román-Godínez, S. Torres-Ramos, R. A. Salido-Ruiz, J. A. Morales, On dna numerical representations for genomic similarity computation, *PloS one* 12 (3) (2017) e0173288.
- [74] R. F. Voss, Evolution of long-range fractal correlations and $1/f$ noise in dna base sequences, *Physical review letters* 68 (25) (1992) 3805.
- [75] P. D. Cristea, Conversion of nucleotides sequences into genomic signals, *Journal of cellular and molecular medicine* 6 (2) (2002) 279–303.
- [76] N. Chakravarthy, A. Spanias, L. D. Iasemidis, K. Tsakalis, Autoregressive modeling and feature analysis of dna sequences, *EURASIP Journal on Applied Signal Processing* 2004 (2004) 13–28.
- [77] R. Zhang, C.-T. Zhang, Z curves, an intuitive tool for visualizing and analyzing the dna sequences, *Journal of Biomolecular Structure and Dynamics* 11 (4) (1994) 767–782.
- [78] A. S. Nair, S. P. Sreenadhan, A coding measure scheme employing electron-ion interaction pseudopotential (eiip), *Bioinformation* 1 (6) (2006) 197.
- [79] D. Anastassiou, Genomic signal processing, *IEEE Signal Processing Magazine* 18 (4) (2001) 8–20. doi:10.1109/79.939833.

- [80] N. Yu, Z. Li, Z. Yu, Survey on encoding schemes for genomic data representation and feature learning—from signal processing to machine learning, *Big Data Mining and Analytics* 1 (3) (2018) 191–210.
- [81] J. Shao, X. Yan, S. Shao, Snr of dna sequences mapped by general affine transformations of the indicator sequences, *Journal of mathematical biology* 67 (2) (2013) 433–451.
- [82] C.-T. Zhang, A symmetrical theory of dna sequences and its applications, *Journal of theoretical biology* 187 (3) (1997) 297–306.
- [83] C. Yin, S. S.-T. Yau, Prediction of protein coding regions by the 3-base periodicity analysis of a dna sequence, *Journal of theoretical biology* 247 (4) (2007) 687–694.
- [84] H. Nikookar, Peak-to-average power ratio, in: *Wavelet Radio: Adaptive and Reconfigurable Wireless Systems Based on Wavelets*, Cambridge University Press, 2013, pp. 93–111. doi:10.1017/CBO9781139084697.006.
- [85] I. Pritišanac, R. M. Vernon, A. M. Moses, J. D. Forman Kay, Entropy and information within intrinsically disordered protein regions, *Entropy* 21 (7) (2019) 662.
- [86] S. Vinga, Information theory applications for biological sequence analysis, *Briefings in bioinformatics* 15 (3) (2013) 376–389.
- [87] J. T. Machado, A. C. Costa, M. D. Quelhas, Shannon, rényie and tsallis entropy analysis of dna using phase plane, *Nonlinear Analysis: Real World Applications* 12 (6) (2011) 3135–3144.
- [88] A. Lesne, Shannon entropy: a rigorous notion at the crossroads between probability, information theory, dynamical systems and statistical physics, *Mathematical Structures in Computer Science* 24 (3) (2014).
- [89] M. P. De Albuquerque, I. A. Esquef, A. G. Mello, Image thresholding using tsallis entropy, *Pattern Recognition Letters* 25 (9) (2004) 1059–1065.
- [90] F. M. Lopes, E. A. de Oliveira, R. M. Cesar, Inference of gene regulatory networks from time series by tsallis entropy, *BMC systems biology* 5 (1) (2011) 61.

- [91] A. Ramírez-Reyes, A. R. Hernández-Montoya, G. Herrera-Corral, I. Domínguez-Jiménez, Determining the entropic index q of tsallis entropy in images through redundancy, *Entropy* 18 (8) (2016) 299.
- [92] L. d. F. Costa, F. A. Rodrigues, A. S. Cristino, Complex networks: the key to systems biology, *Genetics and Molecular Biology* 31 (3) (2008) 591–601.
- [93] X. F. Wang, Complex networks: topology, dynamics and synchronization, *International journal of bifurcation and chaos* 12 (05) (2002) 885–916.
- [94] B. K. Singh, K. Verma, A. Thoke, Investigations on impact of feature normalization techniques on classifier’s performance in breast tumor classification, *International Journal of Computer Applications* 116 (19) (2015).
- [95] M. C. de Souto, D. S. de Araujo, I. G. Costa, R. G. Soares, T. B. Luderemir, A. Schliep, Comparative study on normalization procedures for cluster analysis of gene expression datasets, in: *Neural Networks, 2008. IJCNN 2008.(IEEE World Congress on Computational Intelligence). IEEE International Joint Conference on, IEEE, 2008*, pp. 2792–2798.
- [96] L. Breiman, Random forests, *Machine learning* 45 (1) (2001) 5–32.
- [97] T. Hastie, S. Rosset, J. Zhu, H. Zou, Multi-class adaboost, *Statistics and its Interface* 2 (3) (2009) 349–360.
- [98] A. V. Dorogush, V. Ershov, A. Gulin, Catboost: gradient boosting with categorical features support, *arXiv preprint arXiv:1810.11363* (2018).
- [99] J. Cohen, A coefficient of agreement for nominal scales, *Educational and psychological measurement* 20 (1) (1960) 37–46.