

Nfeature: A platform for computing features of nucleotide sequences

Megha Mathur[#], Sumeet Patiyal[#], Anjali Dhall[#], Shipra Jain[#], Ritu Tomer, Akanksha Arora,
Gajendra P. S. Raghava*

Department of Computational Biology, Indraprastha Institute of Information Technology,
Okhla Phase 3, New Delhi-110020, India.

Emails of authors

Megha Mathur (MM): megha19104@iiitd.ac.in

Sumeet Patiyal (SP): sumeetp@iiitd.ac.in

Anjali Dhall (AD): anjalid@iiitd.ac.in

Shipra Jain (SJ): shipra@iiitd.ac.in

Ritu Tomer (RT): ritut@iiitd.ac.in

Akanksha Arora (AA): akankshaar@iiitd.ac.in

Gajendra P. S. Raghava (GPSR): raghava@iiitd.ac.in

#Contributed Equally

***Corresponding author**

Prof. G.P.S. Raghava,

Head of Department, Department of Computational Biology, Indraprastha Institute of
Information Technology, Okhla Phase 3, New Delhi-110020, India.

E-mail address: raghava@iiitd.ac.in

Phone No: +91-11-26907444

Abstract

In the past few decades, public repositories on nucleotides have increased with exponential rates. This pose a major challenge to researchers to predict the structure and function of nucleotide sequences. In order to annotate function of nucleotide sequences it is important to compute features/attributes for predicting function of these sequences using machine learning techniques. In last two decades, several software/platforms have been developed to elicit a wide range of features for nucleotide sequences. In order to complement the existing methods, here we present a platform named Nfeature developed for computing wide range of features of DNA and RNA sequences. It comprises of three major modules namely Composition, Correlation, and Binary profiles. Composition module allow to compute different type of compositions that includes mono-/di-tri-nucleotide composition, reverse complement composition, pseudo composition. Correlation module allow to compute various type of correlations that includes auto-correlation, cross-correlation, pseudo-correlation. Similarly, binary profile is developed for computing binary profile based on nucleotides, di-nucleotides, di-/tri-nucleotide properties. Nfeature also allow to compute entropy of sequences, repeats in sequences and distribution of nucleotides in sequences. In addition to compute feature in whole sequence, it also allows to compute features from part of sequence like split-composition, N-terminal, C-terminal. In a nutshell, Nfeature amalgamates existing features as well as number of novel features like nucleotide repeat index, distance distribution, entropy, binary profile, and properties. This tool computes a total of 29217 and 14385 features for DNA and RNA sequence, respectively. In order to provide, a highly efficient and user-friendly tool, we have developed a standalone package and web-based platform (<https://webs.iitd.edu.in/raghava/nfeature>).

Introduction

In recent years, the amount of available biological sequence data has increased exponentially due to the significant increase in genome sequencing projects [1-8]. With an increasing data set, it is important to capture relevant information to solve biological questions [9-10]. In literature, several studies reported the application of machine learning techniques in annotation of biological molecules (DNA/RNA) [11-17]. It is important to compute features of nucleotide sequences, in order to implement machine learning techniques for developing models for annotating biological molecules [18]. In the past, various computational tools have been developed (such as PseKNC [19], RepDNA [20], RepRNA [21], BioTriangle [22], PyBioMed

[23], BioSeq-Analysis2.0 [24]) for computing various sequence-based features. Despite several methods developed in the past, some important features have not been integrated into those platforms.

In order to supplement previous efforts, we have made a systematic attempt to developed a webserver platform “Nfeature” that integrates most of the features discovered in the past along with the incorporation of new features. In this study, we have introduced new features Nucleotide Repeat Index (NRI), Distance distribution of Nucleotides (DDN) and Entropy at sequence level as well as nucleotide level as a novel feature for DNA/RNA sequences. We have also incorporated Binary profile-based features for a given nucleotide sequence. These features are essential for motif predictions, factors/enhancers binding sites, etc. Using these modules, user can easily calculate the binary fingerprints of each nucleotide in a given sequence. Nfeature is a comprehensive platform to fetch all relevant information from a given nucleotide sequence in the form of vectors, which can be directly used for developing prediction models. A user-friendly web server and standalone package have been developed to facilitate users in computing features of nucleotide sequences (Figure 1).

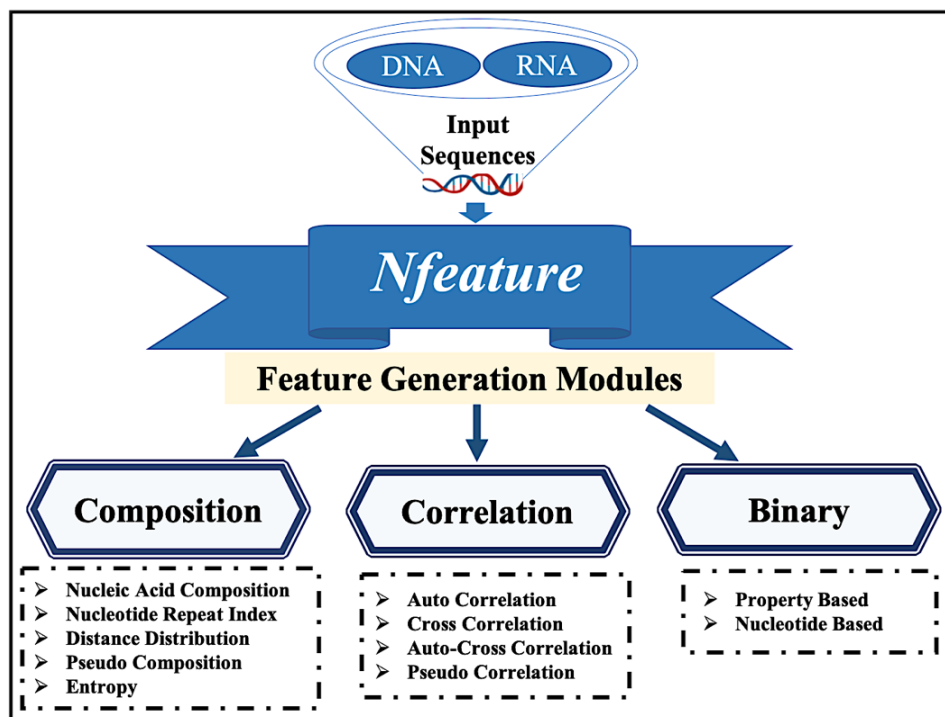


Figure 1: Schematic representation of major modules of Nfeature

Nfeature Overview

Composition/distance distribution-based features

This module aims to calculate nucleotide composition-based features in nucleotide sequences. It enable users to compute nucleic acid composition, distance distribution of nucleotides (DDN), nucleotide repeat index (NRI), pseudo composition and entropy of a sequence. We have incorporated most of the features used in previous studies. In addition, we integrate new features like entropy where we compute entropy at sequence and at nucleotide level. Several past studies have shown that nucleotide repeats have important biological functions. For example, repeated DNA residues are essential for the expression of unique coding sequence which further form nucleoprotein complexes. Whereas, in some cases these repeated sequences causes biological disorders like “GGGGCC” repeat sequences in the *C9orf72* gene causes neurodegenerative disease [25-32]. Additionally, Pfeature, a computational tool of feature extraction also incorporates residue repeat and distance distribution feature generation methods for the protein sequences [33]. Hitherto, no method capture such information of nucleotides from nucleotide sequences. In this study, first time we used NRI to calculate the repeating nucleotide information of DNA/RNA sequence. It measures the number of continuous runs of a nucleotide in a biological sequence. To calculate the distance distribution of nucleotide residues, a module is used for computing distance distribution information for a nucleotide sequence. Residue repeat and distance distribution of nucleotide sequence can be calculated by using the given formula:

$$NRI_i = \frac{\sum_{j=1}^N (W_j)^2}{\sum_{j=1}^N W_j} \quad (1)$$

$$DDN_i = \frac{(W_{NT})^2 + \sum_{j=1}^N (W_j)^2 + (W_{CT})^2}{(L - F_i) + 1} \quad (2)$$

Here, NRI_i and DDN_i stands for nucleotide repeat index and distance distribution of nucleotide type i , where N stands for the maximum number of occurrences, W_j stands for the number of repeats in occurrence “ j ” for nucleotide type “ i ”, W_{NT} - nucleotide distance from N-terminal; W_j - Inter-distance between nucleotides “ i ”; W_{CT} - Nucleotide distance from C-terminal, L - Total length of nucleic acid sequence; F_i – Frequency of nucleotide type “ i ”.

Shannon entropy plays a significant role in the field of information theory. Recently several studies have shown the importance of entropy in DNA sequences for example, to investigate exons and introns in the DNA sequences [34], to identify DNA sequence diversity between different alleles within one individual [35]. To the best of our knowledge there is no method

incorporating entropy-based features. So, to get the Shannon entropy information of DNA and RNA sequences, we first time introduce an entropy-based module which computes entropy of DNA/RNA sequences at residue as well as sequence level. Sequence and residue level entropy is calculated using following equations 3, and 4.

$$H(X) = - p_i \log_2 p_i \quad (3)$$

$$H(X) = - \sum_{i=1}^4 p_i \log_2 p_i \quad (4)$$

Here “i” is the nucleotide in the sequence and X is any nucleotide sequence, and p_i is the probability of a given nucleotide in the sequence.

Correlation based features

In this module, correlation-based features of DNA and RNA sequences are calculated. Correlation is defined as a relation between properties/features i.e. if a feature variable is related to its own then it is defined as auto-correlation and if there exists some correlation between two features/variables then it is known as cross-correlation. Correlation based features basically convert the different length DNA and RNA sequences into fixed length vectors, so that machine learning techniques can be applied to the extracted features. These descriptors identify features based on nucleotide properties along the sequence.

The correlation module of Nfeature calculates the diverse range of features for DNA and RNA sequences, including autocorrelation, cross-correlation, auto-cross correlation, and pseudo correlation. The autocorrelation module further subdivided into seven categories based on the properties and correlation types, such as dinucleotide based autocorrelation (DAC), trinucleotide based autocorrelation (TAC), dinucleotide based Moran autocorrelation (DMAC), trinucleotide based Moran autocorrelation (TMAC), dinucleotide based Geary autocorrelation (DGAC), trinucleotide based Geary autocorrelation (TGAC), and normalized Moreau-Broto autocorrelation (NMBAC) based on 38 properties. Likewise, the cross-correlation and auto cross-correlation module is further divided into two levels, dinucleotide-based cross-correlation (DCC) and trinucleotide-based cross-correlation (TCC), and dinucleotide-based auto-cross correlation (DACC) and trinucleotide-based auto cross-correlation (TACC). Finally, the pseudo correlation module is further divided into four sub-modules, parallel correlation pseudo dinucleotide composition (PC_PDNC), parallel correlation pseudo trinucleotide composition (PC_PTNC), serial correlation pseudo-dinucleotide composition (SC_PDNC), and serial correlation pseudo-trinucleotide

composition (SC_PTNC). For DNA sequences, all the correlation-based features are generated, whereas, for RNA sequences, features incorporating only dinucleotide-based properties are calculated by Nfeature. Features based on dinucleotide properties include 38 properties, and trinucleotide properties-based features incorporate 12 properties. Our correlation-based module can calculate a total of 3,526 and 2,959 descriptors of DNA and RNA sequence (with default parameters).

Binary profile-based features

This module covers binary profile-based features of a given DNA/ RNA sequence. Binary profile-based features are important to motif predictions, factors/enhancers binding sites, etc. Using these modules, users can compute binary equivalents of each nucleotide in a given sequence. Generally, overlapping windows are used to create all possible patterns of fixed length for a given sequence. Then each pattern is converted to a binary profile (1 is used for presence and 0 is used to depict absence) as a numerical representation of each profile. Composition and correlation-based features are adequate to explain the function of the nucleic acids as a whole but fails to capture the residue level information. In order to predict the function of nucleic acids at the residue level, such as protein or biomolecule interacting nucleic acids, it is essential to represent the nucleotides in such a way or feature that apprehends the order and position of the nucleotides. Binary profiles or one-hot encoding meet the requirements mentioned above, where each or set of nucleotides designated by the vectors consist of ones and zeroes. In Nfeature, binary profiles are broadly categorized into two categories, such as property-based and nucleotide-based binary profiles. In DNA, property-based binary profiles are further divided into dinucleotide and tri-nucleotide property based binary profiles. In contrast, for RNA, only dinucleotide properties based binary profile is calculated. Whereas mono -, di- and tri-nucleotide based binary profiles are calculated for DNA and RNA under nucleotide-based binary profiles.

Nucleotide-Based Features

In this module, three sub-modules were involved for DNA and RNA, such as binary profile for mono-nucleotide (BPM), dinucleotide (BPD), and trinucleotide (BPT). In BPM, each nucleotide is represented by the vector of length four; for instance, A is defined as (1,0,0,0); G is designated as (0,1,0,0); C is indicated as (0,0,1,0), and T in case of DNA and U in case of

RNA, is referred as (0,0,0,1). Similarly, for BPD, each dinucleotide is represented by the vector of size 16 with one element as 1, and the rest are zeros, such as dinucleotide AA is indicated as (1,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0). Likewise, in BPT, each trinucleotide is represented by the vector of size 64. This module can compute 840 descriptors of a single DNA/RNA sequence with length of at least 10 residues.

Property-Based Features

This module is one of the novel modules of Nfeature, representing the binary profile of the sequence based on dinucleotide physicochemical properties (BP_DP) of DNA and RNA and also trinucleotide physicochemical properties (BP_TP) of DNA. 38 dinucleotide and 12 trinucleotide properties are involved in Nfeature. Each sequence in BP_DP is expressed by the vector of length $16*(N-1)*p$, where N is the length of the sequence, and p is the number of selected features between 1 to 38. Likewise, the sequence in BP_TP is signified by the vector of length $64*(N-2)*p$, where N is the length of the sequence, and p is the number of selected features between 1 to 12. The property-based module of Nfeature can compute highest number of features i.e., 11,616 of a single DNA sequence with length of at least 10 residues. In case of RNA based features, it can compute 6798 features of a single sequence with a length of 10 residue sequence.

Web Implementation and Standalone Package

In order to facilitate scientific community, we have developed a user-friendly webserver named as “Nfeature”. Webserver was integrated using Apache software on Linux/Ubuntu operating system. All the web pages have been developed with the help of HTML, CSS3 and PHP5. It is compatible with number of devices such as smart phone, laptop, Desktop, iPad. The submit page of server permit the user to submit nucleotide sequences (DNA/RNA) in FASTA format. The result page of web server allows the user to download the output in csv format. Figure 1 represents the description of all the modules of Nfeature tool (Figure 2). Additionally, the standalone package incorporates a readme file, description manual and separate codes for both DNA and RNA modules in respective directories.

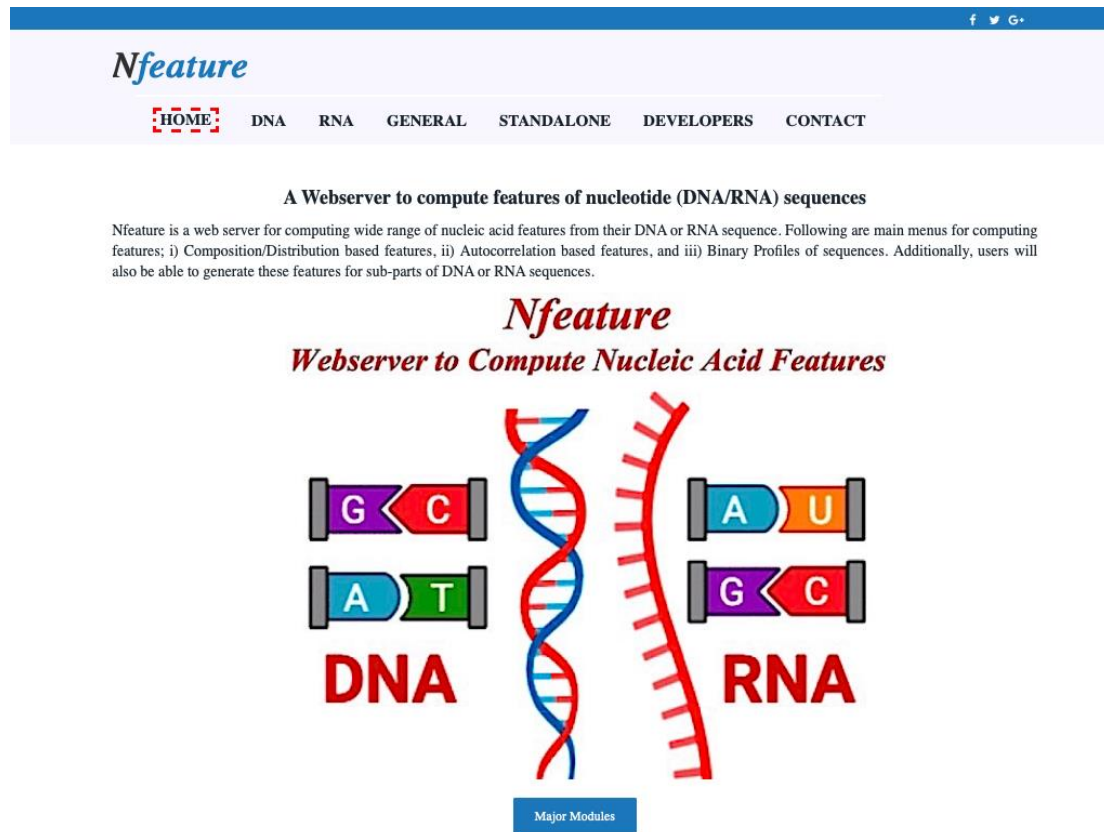


Figure 2: Web-interface of Nfeature to compute features of DNA/RNA sequences

Comparison with other tools

In Table 1, we have compared Nfeature with the existing softwares/web servers based on the platform compatibility, package development and current running status. We showed that most of the software's and packages are either not working or are not platform compatible with most of the frequently used operating systems. Nfeature and BioSeq-Analysis 2.0 are found to be available as a working web server and standalone package. Both of the tools were found to be compatible with widely used platforms such as Windows, Linux, and Mac OS.

Table 1: Comparison between the existing methods and Nfeature

Tools	Web Server	Standalone/Package	Platform Compatibility
Nfeature	✓	✓	Windows, Linux, and Mac OS
PseKNC	✗	✗	Not Working

repDNA	✗	✓	Windows, Linux, and Mac OS
repRNA	✓	✗	Windows, Mac OS
BioTriangle	✓	✗	Windows, Linux, and Mac OS
BioSeq-Analysis2.0	✓	✓	Windows, Linux, and Mac OS
PyBioMed	✗	✓	Windows, and Linux

Software's/web servers developed by various groups to compute various features based on input DNA/ RNA sequences, have their own limitations. Each tool/ web server aims to provide a unique set of features to serve user requirements. Nfeature is developed with an idea to add novel features to the existing features developed by other tools and also to provide conceivable feature generation techniques in a user-friendly fashion on a single platform. Nfeature integrates all the existing tools to calculate features for DNA and RNA sequences as shown in Table 1 and has added novel features.

Table 1: Comparison of features integrated in different platform/software. These descriptors are computed at nucleotide level, can be used to compute overall function/structure of a DNA/RNA sequence.

Features Integrated in Nfeature			Feature Calculation Software (DNA) [#]	Feature Calculation Software (RNA) [#]
	Type of Descriptors	Number [^]		
Nucleic acid Composition	Mono-nucleotide	4	All	Most
	Di-nucleotide	16	All	Most
	Tri-nucleotide	64	All	Most
Reverse complement K-Mer composition	Mono-nucleotide	2	Most	Few
	Di-nucleotide	10	Most	Few
	Tri-nucleotide	32	Most	Few
Pseudo Composition	Pseudo Di-nucleotide	16+1m	All	Few
	Pseudo Tri-nucleotide	64+1m	All	Few
Auto Correlation	Di-Nucleotide	Pd*lag	All	Most
	Tri-Nucleotide	Pt*lag	Most	Few
	Dinucleotide Moran	Pd*lag	Few	Few
	Trinucleotide Moran	Pt*lag	Few	Few

	Dinucleotide Geary	$Pd*lag$	Few	Few
	Trinucleotide Geary	$Pt*lag$	Few	Few
	Normalized Moreau-Broto	$Pd*lag$	Few	Few
	Dinucleotide Cross Correlation	$Pd(Pd-1)*lag$	Most	Few
	Trinucleotide Cross Correlation	$Pt(Pt-1)*lag$	Most	Few
	Auto Dinucleotide -Cross Correlation	$(Pd(Pd-1)*lag)+Pd$	Most	Few
	Auto Trinucleotide - Cross Correlation	$(Pt(Pt-1)*lag)+Pt$	Most	Few
	Parallel Correlation Pseudo Dinucleotide	$16+lm$	Most	Most
	Parallel Correlation Pseudo Trinucleotide	$64+lm$	Most	Few
	Serial Correlation Pseudo Dinucleotide	$(L-1)*Pd*16$	Most	Most
	Serial Correlation Pseudo Trinucleotide	$(L-1)*Pt*64$	Most	Few
Binary Profile	Mono-nucleotide	$4*L$	Few	Few
Number of descriptors for whole nucleotide			DNA=16020	RNA=8740
Total descriptors (Whole nucleotide+ N-Term + C-Term + RN-Term + Split etc.)				128160

#All=Available in all software's, Most=Present in atleast 4 software's, Few= Present in 3 or less than 3 software's;

^Pd=38, Pt=12, lm=1, lag=1, L=10, Split=2,

This platform provides new feature generation methods such as Nucleotide Repeat Index, Distance Distribution, Sequence level Entropy, Nucleotide level Entropy and Binary profiles of inputs DNA/ RNA sequences as represented in Table 2.

Table 2: Novel features for DNA and RNA sequences incorporated in Nfeature Tool

Features	Type of Descriptors	Number [^]	Category
Nucleotide Repeat Index	Mono-nucleotide	4	DNA/RNA
Distance Distribution	Mono-nucleotide	4	DNA/RNA
Entropy	Nucleotide-level	1	DNA/RNA
Entropy	Sequence-level	4	DNA/RNA
Binary Profile Properties	Di-nucleotide properties	$(L-1)*Pd*16$	DNA/RNA
Binary Profile Properties	Tri-nucleotide properties	$(L-1)*Pt*64$	DNA
Binary Profile	Di-Nucleotide	$16*L$	DNA/RNA
Binary Profile	Tri-Nucleotide	$64*L$	DNA
Number of descriptors for whole nucleotide			DNA = 13197 RNA = 5645
Total descriptors (Whole nucleotide + N-Term + C-Term + RN-Term + Split etc.)			105576

^Pd=38, Pt=12, lm=1, lag=1, L=10, Split=2

Discussion and Conclusion

Due to advancement in technology in next generation sequencing, databases like NCBI [36], GenBank [4], EMBL [2], INSDC - DDBJ [1] are growing with exponential rate. In order to address numerous unsolved biological questions, there is an urgent need to develop computer-aided tools to annotate new sequences in above databases. In order to annotate any sequence, most important step is computation of numerical vector that represent characteristics of a sequence. In simple term computation of features or descriptors of a sequence is an important and essential step for computing function or structure of a sequence. In the past various package and web-based platform has been developed to compute wide range of features of proteins and nucleotide sequences. For example, pseudo K-tuple nucleotide composition (PseKNC), method used to generate composition and few correlation-based features from DNA/RNA sequences [19]. RepDNA and RepRNA have been developed for calculating various features for DNA and RNA sequences respectively [20,21]

BioTriangle is a web server that generate features for chemicals, proteins and nucleotide sequences and their interactions [22]. They have reported to calculate 14 type features from DNA/RNA sequences. PyBioMed also allows to generate features for chemicals, proteins nucleotides [23]. This package generates compositions, autocorrelation and pseudo nucleic acid composition-based feature vectors. BioSeq-Analysis [37], which is recently updated to a new version, known as BioSeq-Analysis 2.0 [24]. This web server also incorporates prediction models. None of the above software allow to compute entropy-based features, Nucleotide repeat Index and Distance distribution of DNA/RNA sequences, unlike Nfeature. On the other hand, Nfeature integrates all existing features available for DNA/RNA sequences (Table 1) and adds eight novel features for better understanding of the sequence insights as represented in Table 2.

In summary, a number of tools have been developed to compute various features based on DNA/RNA sequences. These tools provide a unique set of features to cater different user requirements. Aim of developing Nfeature was to complement existing tools and to provide all possible feature generation techniques at a single platform in a user-friendly mode. Also, in order to overcome the limitations of the existing tools we have integrated all the existing tools features for DNA and RNA sequences with few new features in Nfeature platform. Overall, it

is a comprehensive, easy-to-use web server/standalone package that allow users to calculate various features.

Acknowledgements

Authors are thankful to the Department of Computational Biology, IIIT-Delhi for infrastructure, Department of Biotechnology (DBT), Department of Science and Technology (DST-INSPIRE) and Council of Scientific and Industrial Research (CSIR), Govt. of India for financial support and fellowships.

Author contribution

MM and SP wrote all the scripts. AD, SP, RT, AA, and SJ developed the web interface. RT and AA prepared the manual. SJ and AD prepared the first draft of manuscript. MM, SP AD, SJ, RT, AA, and GPSR prepared the final version of manuscript. GPSR conceived the idea and coordinated the entire project.

Funding

The current work has not received any specific grant from any funding agencies.

Conflict of Interest

The authors declare no competing financial and non-financial interests.

References

- [1] G. Cochrane, I. Karsch-Mizrachi, T. Takagi, C. International Nucleotide Sequence Database, The International Nucleotide Sequence Database Collaboration, *Nucleic Acids Res*, 44 (2016) D48-50.
- [2] G. Stoesser, W. Baker, A. van den Broek, E. Camon, M. Garcia-Pastor, C. Kanz, T. Kulikova, R. Leinonen, Q. Lin, V. Lombard, R. Lopez, N. Redaschi, P. Stoehr, M.A. Tuli, K. Tzouvara, R. Vaughan, The EMBL Nucleotide Sequence Database, *Nucleic Acids Res*, 30 (2002) 21-26.

- [3] C.S. Pareek, R. Smoczynski, A. Tretyn, Sequencing technologies and genome sequencing, *J Appl Genet*, 52 (2011) 413-435.
- [4] D.A. Benson, M. Cavanaugh, K. Clark, I. Karsch-Mizrachi, D.J. Lipman, J. Ostell, E.W. Sayers, GenBank, *Nucleic Acids Res*, 41 (2013) D36-42.
- [5] R. Leinonen, R. Akhtar, E. Birney, L. Bower, A. Cerdeno-Tarraga, Y. Cheng, I. Cleland, N. Faruque, N. Goodgame, R. Gibson, G. Hoad, M. Jang, N. Pakseresht, S. Plaister, R. Radhakrishnan, K. Reddy, S. Sobhany, P. Ten Hoopen, R. Vaughan, V. Zalunin, G. Cochrane, The European Nucleotide Archive, *Nucleic Acids Res*, 39 (2011) D28-31.
- [6] L. Hood, L. Rowen, The Human Genome Project: big science transforms biology and medicine, *Genome Med*, 5 (2013) 79.
- [7] F.S. Collins, L. Fink, The Human Genome Project, *Alcohol Health Res World*, 19 (1995) 190-195.
- [8] M.C. Schatz, Biological data sciences in genome research, *Genome Res*, 25 (2015) 1417-1422.
- [9] G.W. Brodland, How computational models can help unlock biological systems, *Semin Cell Dev Biol*, 47-48 (2015) 62-73.
- [10] S.K.S. Sabyasachi Dash, Mohit Sharma & Sandeep Kaushik, Big data in healthcare: management, analysis and future prospects, *Journal of Big Data*, (2019).
- [11] K. Huang, C. Xiao, L.M. Glass, C.W. Critchlow, G. Gibson, J. Sun, Machine learning applications for therapeutic tasks with genomics data, *Patterns (N Y)*, 2 (2021) 100328.
- [12] A. Yang, W. Zhang, J. Wang, K. Yang, Y. Han, L. Zhang, Review on the Application of Machine Learning Algorithms in the Sequence Data Mining of DNA, *Front Bioeng Biotechnol*, 8 (2020) 1032.
- [13] M. Mahmud, M.S. Kaiser, T.M. McGinnity, A. Hussain, Deep Learning in Mining Biological Data, *Cognit Comput*, (2021) 1-33.
- [14] C. Xu, S.A. Jackson, Machine learning and complex biological data, *Genome Biol*, 20 (2019) 76.
- [15] M.R.G.M. Jonathan Schmidt, Silvana Botti, Miguel A. L. Marques Recent advances and applications of machine learning in solid-state materials science, *npj computational materials*, (2019).
- [16] V.I. Jurtz, A.R. Johansen, M. Nielsen, J.J. Almagro Armenteros, H. Nielsen, C.K. Sonderby, O. Winther, S.K. Sonderby, An introduction to deep learning on biological sequence data: examples and solutions, *Bioinformatics*, 33 (2017) 3685-3690.
- [17] U.S.a.Z. Usman, Chapter 4 Biological Sequence Analysis.

- [18] I.Y. Abdurakhmonov, *Bioinformatics: Basics, Development, and Future*, (2016).
- [19] W. Chen, T.Y. Lei, D.C. Jin, H. Lin, K.C. Chou, PseKNC: a flexible web server for generating pseudo K-tuple nucleotide composition, *Anal Biochem*, 456 (2014) 53-60.
- [20] B. Liu, F. Liu, L. Fang, X. Wang, K.C. Chou, repDNA: a Python package to generate various modes of feature vectors for DNA sequences by incorporating user-defined physicochemical properties and sequence-order effects, *Bioinformatics*, 31 (2015) 1307-1309.
- [21] B. Liu, F. Liu, L. Fang, X. Wang, K.C. Chou, repRNA: a web server for generating various feature vectors of RNA sequences, *Mol Genet Genomics*, 291 (2016) 473-481.
- [22] J. Dong, Z.J. Yao, M. Wen, M.F. Zhu, N.N. Wang, H.Y. Miao, A.P. Lu, W.B. Zeng, D.S. Cao, BioTriangle: a web-accessible platform for generating various molecular representations for chemicals, proteins, DNAs/RNAs and their interactions, *J Cheminform*, 8 (2016) 34.
- [23] J. Dong, Z.J. Yao, L. Zhang, F. Luo, Q. Lin, A.P. Lu, A.F. Chen, D.S. Cao, PyBioMed: a python library for various molecular representations of chemicals, proteins and DNAs and their interactions, *J Cheminform*, 10 (2018) 16.
- [24] B. Liu, X. Gao, H. Zhang, BioSeq-Analysis2.0: an updated platform for analyzing DNA, RNA and protein sequences at sequence level and residue level based on machine learning approaches, *Nucleic Acids Res*, 47 (2019) e127.
- [25] A. Chanou, S. Hamperl, Single-Molecule Techniques to Study Chromatin, *Front Cell Dev Biol*, 9 (2021) 699771.
- [26] M. Rajewska, K. Wegrzyn, I. Konieczny, AT-rich region and repeated sequences - the essential elements of replication origins of bacterial replicons, *FEMS Microbiol Rev*, 36 (2012) 408-434.
- [27] D.E. Handy, R. Castro, J. Loscalzo, Epigenetic modifications: basic mechanisms and role in cardiovascular disease, *Circulation*, 123 (2011) 2145-2156.
- [28] I. Malik, C.P. Kelley, E.T. Wang, P.K. Todd, Molecular mechanisms underlying nucleotide repeat expansion disorders, *Nat Rev Mol Cell Biol*, 22 (2021) 589-607.
- [29] A. De Bustos, A. Cuadrado, N. Jouve, Sequencing of long stretches of repetitive DNA, *Sci Rep*, 6 (2016) 36665.
- [30] J.J. Miret, L. Pessoa-Brandao, R.S. Lahue, Instability of CAG and CTG trinucleotide repeats in *Saccharomyces cerevisiae*, *Mol Cell Biol*, 17 (1997) 3382-3387.
- [31] J.A. Shapiro, R. von Sternberg, Why repetitive DNA is essential to genome function, *Biol Rev Camb Philos Soc*, 80 (2005) 227-250.

- [32] A.A. Abugable, J.L.M. Morris, N.M. Palminha, R. Zaksauskaite, S. Ray, S.F. El-Khamisy, DNA repair and neurological disease: From molecular understanding to the development of diagnostics and model organisms, *DNA Repair (Amst)*, 81 (2019) 102669.
- [33] S.P. Akshara Pande, Anjali Lathwal, Chakit Arora, Dilraj Kaur, Anjali Dhall, Gaurav Mishra, Harpreet Kaur, Neelam Sharma, Shipra Jain, Salman Sadullah Usmani, Piyush Agrawal, Rajesh Kumar, Vinod Kumar, Gajendra P.S. Raghava, Computing wide range of protein/peptide features from their sequence and structure, (2019).
- [34] J. Li, L. Zhang, H. Li, Y. Ping, Q. Xu, R. Wang, R. Tan, Z. Wang, B. Liu, Y. Wang, Integrated entropy-based approach for analyzing exons and introns in DNA sequences, *BMC Bioinformatics*, 20 (2019) 283.
- [35] W.B. Sherwin, *Entropy and Information Approaches to Genetic Diversity and its Expression: Genomic Geography, Entropy in Genetics and Computational Biology*, (2010).
- [36] N.R. Coordinators, Database resources of the National Center for Biotechnology Information, *Nucleic Acids Res*, 44 (2016) D7-19.
- [37] B. Liu, BioSeq-Analysis: a platform for DNA, RNA and protein sequence analysis based on machine learning approaches, *Brief Bioinform*, 20 (2019) 1280-1294.