

1 **Prediction of antibiotic resistant strains of bacteria from their**
2 **beta-lactamases protein**

3 **Lubna Maryam^{†1}, Anjali Dhall^{†1}, Sumeet Patiyal^{†1}, Salman Sadullah Usmani¹, Neelam**
4 **Sharma¹, and Gajendra Pal Singh Raghava*¹**

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

7
8 **Emails of Authors:**

9 LM: lucymary20@gmail.com

10 AD: anjalid@iiitd.ac.in

11 SP: sumeetp@iiitd.ac.in

12 SSU: salman007usmani@gmail.com

13 NS: neelams@iiitd.ac.in

14 GPSR: raghava@iiitd.ac.in

15

16 [†]**These authors have contributed equally to this work.**

17 ***Corresponding author**

18 Prof. G.P.S. Raghava,

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

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

22 Phone No: +91-11-26907444

23

24

25

26

27 **Abstract**

28 Number of beta-lactamase variants have ability to deactivate ceftazidime antibiotic, which is
29 the most commonly used antibiotic for treating infection by Gram-negative bacteria. In this
30 study an attempt has been made to develop a method that can predict ceftazidime resistant
31 strains of bacteria from amino acid sequence of beta-lactamases. We obtained beta-
32 lactamases proteins from the β -lactamase database, corresponding to 87 ceftazidime-sensitive
33 and 112 ceftazidime-resistant bacterial strains. All models developed in this study were
34 trained, tested, and evaluated on a dataset of 199 beta-lactamases proteins. We generate 9149
35 features for beta-lactamases using Pfeature and select relevant features using different
36 algorithms in scikit-learn package. A wide range of machine learning techniques (like KNN,
37 DT, RF, GNB, LR, SVC, XGB) has been used to develop prediction models. Our random
38 forest-based model achieved maximum performance with AUROC of 0.80 on training dataset
39 and 0.79 on the validation dataset. The study also revealed that ceftazidime-resistant beta-
40 lactamases have amino acids with non-polar side chains in abundance. In contrast,
41 ceftazidime-sensitive beta-lactamases have amino acids with polar side chains and charged
42 entities in abundance. Finally, we developed a webservice “ABCRpred”, for the scientific
43 community working in the era of antibiotic resistance to predict the antibiotic
44 resistance/susceptibility of beta-lactamase protein sequences. The server is freely available at
45 (<http://webs.iiitd.edu.in/raghava/abcrpred/>).

46 **Keywords:** Antibiotic-resistance strains, Beta-lactamases, Ceftazidime antibiotic, Prediction
47 method, Machine learning techniques

48

49 **Key Points**

- 50 • Ceftazidime is commonly used to treat infection caused by Gram-negative bacteria.
- 51 • Beta-lactamase is responsible for lysing ceftazidime, make it resistant to bacteria.
- 52 • Comparison of resistant and sensitive variants of beta-lactamase.
- 53 • Classification of sensitive and resistant strain of bacteria based on beta-lactamase.
- 54 • Prediction models have been developed using different machine learning techniques.

55

56

57

58

59

60

Author's Biography

61

1. Lubna Maryam is currently working as a Post-Doctoral Fellow in the Department of Computational Biology, Indraprastha Institute of Information Technology, New Delhi, India.

62

63

64

2. Anjali Dhall is currently working as a Ph.D. scholar in Bioinformatics from Department of Computational Biology, Indraprastha Institute of Information Technology, New Delhi, India.

65

66

67

3. Sumeet Patiyl is currently working as a Ph.D. scholar in Bioinformatics from Department of Computational Biology, Indraprastha Institute of Information Technology, New Delhi, India.

68

69

70

4. Salman Sadullah Usmani has completed his Ph.D. in Bioinformatics from CSIR-IMTECH, Chandigarh, India and is now working as Research Associate-I in the Department of Computational Biology, Indraprastha Institute of Information Technology, New Delhi, India.

71

72

73

5. Neelam Sharma is currently working as a Ph.D. scholar in Bioinformatics from Department of Computational Biology, Indraprastha Institute of Information Technology, New Delhi, India.

74

75

76

6. G.P.S. Raghava is currently working as Professor and Head of Department of Computational Biology, Indraprastha Institute of Information Technology, New Delhi, India.

77

78

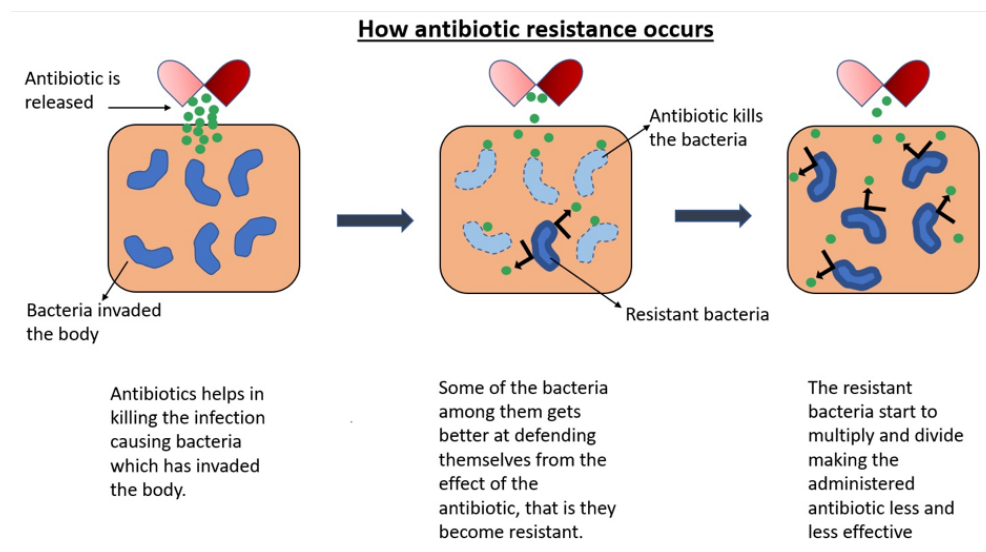
79

80

81

82 Introduction

83 Antimicrobial resistance (AMR) is the ability of bacteria to resist the effect of antibiotics that
84 are administered during infection (Figure 1). In 2020, WHO declared AMR as one of the
85 world's top 10 public health threats. There is a number of reasons for drug resistance that
86 include overuse of antibiotics and the emerging strain of bacteria. Moreover, the alarming
87 spread of multi-drug resistant bacteria (MDR) continues to lurk our capability to treat
88 common infections (e.g., sepsis, sexually transmitted infections, urinary tract infections,
89 diarrhea). There are a number of mechanisms adopted by the bacteria to evade killing by
90 antimicrobial molecules that include, production of antibiotic lysing enzymes (e.g., Beta-
91 lactamases), lowering the permeability of cell membrane, and modification of the antibiotics
92 binding site [1]. Beta-lactam antibiotics are the most prescribed antibiotics to fight broad
93 spectrum infections, i.e., 65% of the total antibiotics in the market [2]. These antibiotics have
94 four membered beta-lactam rings in their molecular structure, which are destroyed by beta-
95 lactamases [3]. Thus, beta-lactamases are responsible for multi-drug resistance against beta-
96 lactam antibiotics [4]. The number of beta-lactamases is continuously growing, around 7166
97 beta-lactamases have been already identified [5]. There are only a few variants of beta-
98 lactamases on which beta-lactam antibiotics are working (sensitive). Resistant beta-lactamase
99 genes that are spread via diffusion of mobile genetic elements, spread of epidemic plasmids,
100 dispersion of specific clones and horizontal gene transfer [6]. Therefore, there is an urgent
101 need to develop prediction models that can discriminate antibiotics sensitive and resistant
102 variants of beta-lactamases.



103

104

Figure 1: Pictorial representation of how antibiotic resistance occurs.

105

106 The standard technique to test the resistance and susceptibility of a strain to a particular
107 antibiotic is the disc diffusion test. This method is reliable and reproducible, but it is time
108 consuming and labour-intensive. Thus, there is a need to develop computational methods that
109 can predict antibiotic resistant strains of bacteria. Due to advancements in the next-generation
110 sequencing (NGS) technologies, it is in routine to sequence a gene or whole genome of
111 bacteria or metagenome. Several repositories have been already developed to maintain the
112 information regarding the genes, mutations, genomes and metagenomes of bacterial strains
113 [7]. This information has been used to develop methods for predicting drug resistant strains
114 of bacteria. These methods are mainly based on identification of antibiotic resistant gene,
115 mutation, whole genome and metagenome [8–11]. Almost all the existing tools are generic in
116 nature, where these methods predict whether a bacterial strain is resistant to all antibiotics. In
117 other words, these methods predict multi-antibiotics resistant bacteria. In the past, a large
118 number of antibiotics have been discovered to kill bacteria by a different mechanism. It is
119 possible that bacterial strain is only resistant to a particular antibiotic or class of antibiotics
120 but sensitive to other class of antibiotics. Thus, it is important to develop a method that can
121 predict antibiotic-specific sensitive or resistant bacterial strains; similar to personalized
122 medicine [12–14]. This is very important to manage treatment of bacterial infection using a
123 particular antibiotic which is sensitive to bacteria responsible for a given infection. In simple
124 worlds there is need to treat a bacterial infection using strain-specific antibiotics which is
125 similar to personalize drugs. Best of our knowledge, there is no computational tool that can
126 predict whether a bacterial strain is sensitive or resistant to a antibiotics. In this study, we first
127 time made an attempt to develop method for antibiotic ceftazidime that belongs to beta-
128 lactam group. We selected ceftazidime because it is routinely used for treatment of wide
129 range of bacterial infections like meningitis, sepsis, joint infection, urinary tract infection. In
130 addition, ceftazidime has been tested clinically on a number of bacterial strains where MIC
131 have been determined. In order to identify sensitive and resistant strain from MIC values,
132 European Committee on Antimicrobial Susceptibility Testing (2020) proposed that the
133 Enterobacterales are susceptible to ceftazidime when its concentration is less than or equal to
134 1 mg/ml and resistant when concentration is greater than 4 mg/ml. It is well known fact that
135 beta-lactamases are responsible for lysing ceftazidime or causing resistance. Thus, we have
136 designed a model for predicting beta-lactamase variants that make ceftazidime sensitive or
137 resistant to a bacterial strain. We used state of the arts techniques mainly based on machine
138 learning techniques to develop prediction models [15]. This will help in predicting

139 ceftazidime resistance/susceptibility towards beta-lactamase carrying bacterial species that
140 could emerge in the near future. The platform will provide vista to find out the beta-
141 lactamases strains which are sensitive to ceftazidime antibiotic.

142 **Methods and Material**

143 **Dataset Collection**

144 The main dataset was collected from the β -lactamases database [16]. It incorporates 2383
145 Minimum Inhibitory Concentration (MIC) values (in the presence and absence of beta-
146 lactamase genes) of 980 beta-lactamases (with their protein sequences) with different
147 antibiotics and the fold change in MIC values [16]. The database comprises experimentally
148 validated 21 different types of antibiotics corresponding to class-A, B, C, D beta-lactamase
149 proteins. In this study, we have considered the ceftazidime antibiotic dataset with different β -
150 lactamase protein sequences. Our final dataset included 199 β -lactamase protein sequences.
151 Further, we set a cutoff on MIC values of ceftazidime with β -lactamase proteins. The proteins
152 having (MIC value ≤ 4) were considered as antibiotic susceptible/sensitive proteins, and
153 proteins having (MIC value > 4) were taken as antibiotic resistant ones [17,18]. Finally, we
154 got 87 antibiotic-sensitive and 112 antibiotic-resistant unique proteins, referred to as positive
155 and negative dataset, respectively. Moreover, we have also collected 22 ceftazidime
156 resistance beta-lactamase protein sequences from Resistance Gene Identifier (RGI) database
157 for external validation [19].

158 **Generation of Features**

159 To generate a wide range of features from protein sequences, we have used Pfeature [20]. In
160 this study, we have used the standalone package of the Pfeature tool to compute thousands of
161 protein/peptide features. This tool also calculates the structural and functional properties of
162 protein sequences. We have generated 9149 composition-based features/descriptors using the
163 composition-based feature module of the Pfeature package. It incorporates 15 different type
164 of descriptors such as Amino acid composition (AAC), Dipeptide composition (DPC),
165 Tripeptide composition (TPC), Atomic and bond composition (ABC), Residue repeat
166 Information (RRI), Distance distribution of residue (DDOR), Shannon-entropy of protein
167 (SE), Shannon entropy of all amino acids (SER), Shannon entropy of physicochemical
168 property (SEP), Conjoint triad calculation of the descriptors (CTD), Composition-enhanced
169 transition distribution (CeTD), Pseudo amino acid composition (PAAC), Amphiphilic pseudo

170 amino acid composition (APAAC), Quasi-sequence order (QSO) and Sequence order
171 coupling number (SOCN).

172

173 **Pre-processing and Feature selection**

174 The biggest challenge is to find out the most important features/descriptors which can
175 classify the two classes more accurately. The standardization or scaling of the dataset is the
176 most common requirement for the machine learning techniques. In the current study, to
177 standardize the dataset, we used MinMaxScaler using the sklearn pre-processing package.
178 This scaling function converts the given values into a minimum and maximum range. After
179 the pre-processing step, we identified the best set of features from a huge dimension vector.
180 For determining the best features several dimension reduction methods are currently
181 available. We have used standard feature selection methods in which firstly we removed all
182 low variance features using the variance threshold method of the scikit-learn package [21]. It
183 removes all zero-variance features, so we were left with 275 features. Then, we applied the
184 SVC-L1 feature selection method for the selection of important set of features [21]. This
185 method is based on the support vector classifier (SVC) with linear kernel, penalized with L1
186 regularization. SVC-L1 method was performed on earlier deduced 275 features which
187 provided 33 features. Further we ranked the features based on their performance, using
188 feature selector tool. We developed our final machine learning models on selected 10, 20, and
189 33 features.

190 **Machine Learning**

191 In the present study, we have implemented several machine learning techniques to classify
192 ceftazidime antibiotic-resistant and sensitive/non-resistant proteins. We incorporated K-
193 nearest neighbors (KNNs), Decision tree (DT), Random Forest (RF), Gaussian Naive Bayes
194 (GNB), Logistic Regression (LR) and Support Vector Classifier (SVC) and XGBoost (XGB)
195 classification methods in the study (ref). These techniques are based on different algorithm
196 such as, KNN is a simple and supervised machine learning algorithm. It assumes the
197 similarity between the new data and the available data and put the new data into the category
198 that is most similar to the available categories [22]. DT is a tree-structured classifier based on
199 non-parametric machine learning models, which uses a decision tree as a model to go from
200 observations about a data to conclusions about new data. RF classification method uses
201 ensemble-based techniques which uses several decision trees for the training and prediction
202 of the outcome [23], GNB (Gaussian Naïve Bayes) are a group of supervised classification

203 algorithms based on Bayes theorem which uses probabilistic approach for the classification.
204 LR is a statistical model that measures the relationship between the categorial dependent
205 variable and one or more independent variable by guessing the likelihoods using a logistic
206 function [24]. SVC get the best fit of the data provided, the features can then be fed to see
207 what the predicted class is. XGB uses an iterative approach for the classification. It is a
208 decision tree-based ensemble machine learning technique that uses an approach where new
209 models are created that predict the errors of prior models and then added together to make the
210 final prediction. All these techniques were executed using python-library scikit-learn [21].

211 **Evaluation Techniques**

212 In order to evaluate the classification models, we have used five-fold cross-validation (CV)
213 and external validation method. For the training, testing, and evaluation, the dataset was
214 divided into 80:20 ratio. We have used the standard criterion for the evaluation, in which
215 80% of the data was used for training and 20% was used for external validation [25]. In 5-
216 fold CV, 80% of the data was divided into five equal portions/folds, one-fold was used for
217 testing, and four folds was used for the training purpose. A similar process was repeated five
218 times, in which each portion/fold was utilized for internal training and testing. Further, we
219 checked the performance of machine learning models on external dataset. In this study we
220 have used well established evaluation parameters [26]. It incorporates threshold-dependent
221 and independent parameters. We measured threshold-dependent parameters like sensitivity
222 (Sens), Specificity (Spec), Accuracy (Acc) and Matthews correlation coefficient (MCC) with
223 the help of following equations. The standard threshold-independent parameter is Area Under
224 the Receiver Operating Characteristic (AUROC) curve [27–29] which was computed to
225 estimate the performance of different modes.

226

$$227 \text{ Sensitivity} = \frac{TP}{TP+FN} \times 100 \dots\dots\dots (1)$$

$$228 \text{ Specificity} = \frac{TN}{TN+FP} \times 100 \dots\dots\dots (2)$$

$$229 \text{ Accuracy} = \frac{TP+TN}{TP+FN+FP+TN} \times 100 \dots\dots\dots (3)$$

$$230 \text{ Matthews Correlation Coefficient} = \frac{(TN \times TP) - (FN \times FP)}{\sqrt{[(FP+TP)(FN+TP)(FP+TN)(FN+TN)]}} \times 100 \dots\dots\dots (4)$$

231 The measurements obtained from the above parameters are expressed in terms of
232 TP=True Positive, FP=False Positive, TN=True Negative, FN=False Negative.

233

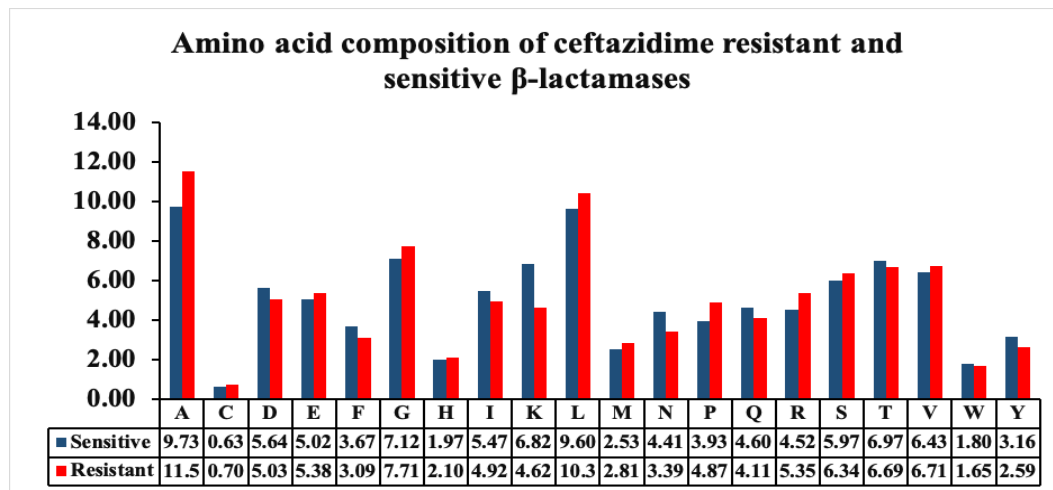
234 **Results**

235 We have used 87 beta-lactamases proteins that are ceftazidime sensitive, having MIC values
236 less than or equal to 4 and 112 ceftazidime resistant beta-lactamase protein sequences with
237 MIC values greater than 4. All analysis and model development have been done on the above
238 dataset.

239 **Analysis based on the amino acid composition**

240 We have analyzed the average amino acid composition for each residue for ceftazidime
241 resistant and sensitive beta-lactamase sequences and found out that residues such as A, G, L,
242 P, and R, is higher in ceftazidime resistant beta-lactamase sequences as compared to sensitive
243 sequences. Whereas, in the case of ceftazidime sensitive sequences, D, I, K, N, T, and Y
244 residues are higher, as shown in Figure 2.

245



246

247 Figure 2: Average amino acid composition of each amino acid residues for ceftazidime-resistant and
248 ceftazidime-sensitive beta-lactamases.

249

250 **Predictions based on machine-learning models**

251 We have implemented various machine learning classifiers such as KNN, DT, RF, GNB, LR,
252 SVC and XGB to develop the prediction model to classify the sequences of ceftazidime
253 resistant and sensitive beta-lactamases. We have calculated each protein sequence's features
254 using the composition-based module of Pfeature, which resulted in 9149 features. On
255 applying the feature selection method using the support vector classifier with L1
256 regularization, we were left with 33 most relevant features. We have ranked these 33 features
257 using feature selector python package and generated prediction models for the top 10, 20, and

258 33 features. For top 10 features, RF has obtained balanced results with AUROC of 0.78 with
 259 MCC of 0.44 for training dataset, whereas AUROC and MCC for the validation dataset are
 260 0.76 and 0.49, respectively. Performance for all the implemented classifier is exhibited in
 261 Table 2. To understand the difference between the positive and negative datasets, we
 262 calculated the average values of the top-10 features of ceftazidime-sensitive and ceftazidime-
 263 resistant beta-lactamases as represented in Table 1.
 264 Table 1: Brief description of top 10 features and their average values in ceftazidime-sensitive and
 265 ceftazidime-resistant beta-lactamases.

Name of features	Description of features	#Average Value-1	#Average Value-2
total_bonds	Bond composition of peptide	4980.345	5353.161
hydrogen_bonds	Bond composition of peptide	2747.425	2958.643
single_bond	Bond composition of peptide	4544.747	4890.884
R_ddor	Distance distribution of Arginine	47.46805	36.44679
Y_ddor	Distance distribution of Tyrosine	55.61483	69.62223
Grantham_gap1	Quasi sequence order of peptide	9422.904	9228.332
Grantham_gap3	Quasi sequence order of peptide	9815.743	9486.559
CeTD_33_SA	Number of transitions taking place from group 1 residues to group 2 residues for solvent accessibility attribute	65.29885	53.53571
CeTD_22_PC	Number of transitions taking place from group 2 residues to group 2 residues for polarizability attribute	66.2069	55.08036
CeTD_33_PO	Number of transitions taking place from group 3 residues to group 3 residues for polarity attribute	67.4023	53.71429

266 #Average Value-1: average values of ceftazidime-resistance beta-lactamases; # Average Value-2:
 267 average values of ceftazidime-sensitive beta-lactamases.

268
 269
 270

Table 2: Performance of various classifiers using top 10 features.

Classifier	Training Dataset	Validation dataset
------------	------------------	--------------------

	Sens	Spec	Acc	AUROC	MCC	Sens	Spec	Acc	AUROC	MCC
KNN	83.33	67.21	76.26	0.83	0.51	76.47	57.69	68.33	0.72	0.35
DT	60.26	72.13	65.47	0.72	0.32	76.47	65.38	71.67	0.72	0.42
RF	76.92	67.21	72.66	0.78	0.44	76.47	73.08	75.00	0.76	0.49
GNB	55.13	77.05	64.75	0.71	0.32	23.53	80.77	48.33	0.64	0.05
LR	69.23	68.85	69.06	0.73	0.38	61.76	65.38	63.33	0.67	0.27
SVC	76.92	78.69	77.70	0.78	0.55	70.59	73.08	71.67	0.74	0.43
XGB	78.21	65.57	72.66	0.75	0.44	76.47	65.38	71.67	0.78	0.42

271 # Sens: Sensitivity; Spec: Specificity; Acc: Accuracy; AUROC: Area Under Receiver Operating
 272 Curve; MCC: Matthews Correlation Coefficient

273 Similarly, RF model developed using top 20 features performed best among all the
 274 classifiers, with AUROC of 0.79 and MCC of 0.48 on training dataset, and AUROC 0.76 and
 275 MCC 0.4 on validation dataset. Performance using other classifiers is given in Table 3.

276 Table 3: Performance of various classifiers using top 20 features.

Classifier	Training Dataset					Validation dataset				
	Sens	Spec	Acc	AUROC	MCC	Sens	Spec	Acc	AUROC	MCC
KNN	78.21	68.85	74.10	0.82	0.47	73.53	57.69	66.67	0.71	0.32
DT	61.54	62.30	61.87	0.67	0.24	70.59	61.54	66.67	0.68	0.32
RF	74.36	73.77	74.10	0.79	0.48	67.65	73.08	70.00	0.76	0.40
GNB	53.85	77.05	64.03	0.70	0.31	29.41	73.08	48.33	0.63	0.03
LR	70.51	70.49	70.50	0.77	0.41	76.47	65.38	71.67	0.69	0.42
SVC	56.41	85.25	69.06	0.75	0.43	64.71	84.62	73.33	0.74	0.49
XGB	75.64	63.93	70.50	0.72	0.40	70.59	65.38	68.33	0.71	0.36

277 # Sens: Sensitivity; Spec: Specificity; Acc: Accuracy; AUROC: Area Under Receiver Operating
 278 Curve; MCC: Matthews Correlation Coefficient

279
 280 For all 33 features, RF obtained the maximum AUROC of 0.80 with 0.48 MCC on the
 281 training dataset, and AUROC of 0.79 with MCC of 0.46 on the validation dataset. We have
 282 reported performance for all classifiers using 33 features in the Table 4.

283 Table 4: Performance of various classifiers using 33 selected features on training and
 284 validation datasets.

Classifier	Training Dataset					Validation dataset				
	Sens	Spec	Acc	AUROC	MCC	Sens	Spec	Acc	AUROC	MCC
KNN	74.36	77.05	75.54	0.81	0.51	73.53	57.69	66.67	0.74	0.32

DT	65.38	70.49	67.63	0.72	0.36	79.41	69.23	75.00	0.73	0.49
RF	74.35	73.77	74.10	0.80	0.48	73.53	73.08	73.33	0.79	0.46
GNB	56.41	72.13	63.31	0.67	0.29	29.41	73.08	48.33	0.64	0.03
LR	74.36	65.57	70.50	0.77	0.40	76.47	69.23	73.33	0.78	0.46
SVC	57.69	88.52	71.22	0.74	0.47	50.00	88.46	66.67	0.74	0.40
XGB	74.36	73.77	74.10	0.77	0.48	76.47	76.92	76.67	0.79	0.53

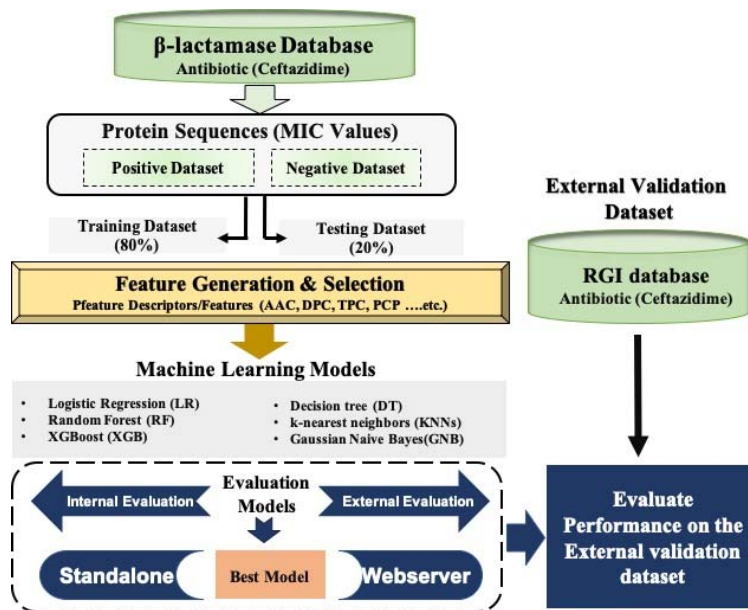
285 # Sens: Sensitivity; Spec: Specificity; Acc: Accuracy; AUROC: Area Under Receiver Operating
286 Curve; MCC: Matthews Correlation Coefficient
287

288 In order to check the robustness of our final model, we have downloaded 22 ceftazidime
289 resistant protein sequences from Resistance Gene Identifier (RGI) database and checked the
290 performance by implementing random forest based model developed on top 33 features. 19
291 out of 22 sequences were giving the correct result, with AUROC of 0.81 and MCC of 0.50 on
292 training dataset, and AUROC 0.79 and MCC 0.71 on validation dataset.

293

294 **Webserver implementation**

295 We have developed a webserver named ABCRpred
296 (<https://webs.iiitd.edu.in/raghava/abcrpred/>) using Random Forest based machine learning
297 approach to serve the scientific world. Since, we wanted to identify sensitive strains of beta-
298 lactamases therefore we developed this method to discriminate between antibiotic resistant
299 and sensitive variants of beta-lactamase strains. 87 antibiotic-sensitive and 112 antibiotic-
300 resistant beta-lactamases protein sequences data were used for training and testing, while
301 building the webserver. The complete architecture of ABCRpred is shown in figure 3.



302

303

304 Figure 3: Overall ABCRpred architecture that shows process of creating datasets, features selection,
305 model development and process of model evaluation.

306

307 The 'Predict' page on the webserver has been developed to predict resistance/susceptibility of
308 any new beta-lactamase protein sequence towards ceftazidime antibiotic. The page enables
309 the users to enter the sequence in FASTA format or upload the file with multiple peptide
310 sequences. User is required to set a random forest threshold and select physicochemical
311 properties as per their requirement. Prediction of each sequence will be carried out according
312 to the selected model. After submitting the input, the output file contains various columns of
313 sequence ID, random forest score, prediction outcome whether the input sequence is resistant
314 or susceptible and the result of selected physicochemical properties. The standalone package
315 (<https://webs.iiitd.edu.in/raghava/abcrpred/stand.php>) has also been incorporated in the
316 webserver to let the users predict the resistance/susceptibility profile of protein sequences
317 even in the absence of the internet. The standalone version incorporated our best models and
318 can work on Linux or Unix operating systems.

319 Discussion

320 The beta-lactam antibiotics are regarded as the drug of choice for the treatment of severe
321 infections caused by *Enterobacteriaceae*. Most of the beta-lactam antibiotics face resistance
322 against beta-lactamases carrying bacteria. Moreover, exposure of beta-lactamase carrying
323 bacterial strains to multitude of beta-lactams has induced active continuous production and
324 mutation of beta-lactamases expanding their activity even against the newly developed beta-

325 lactam antibiotics. In this study we used MIC data of ceftazidime (a beta-lactam antibiotic)
326 against beta-lactamase carrying bacteria for building a prediction model to predict resistance
327 and susceptibility of any newly emerged variant of beta-lactamase carrying bacterial strain. A
328 total of 199 experimental MIC data was collected from a comprehensive database of beta-
329 lactamase enzymes called as β -lactamase Database [16]. Our data of ceftazidime MIC
330 against various beta-lactamase carrying bacterial strains was divided into two sets. One set
331 have MIC values greater than 4 referred to as ceftazidime-resistant strain; other set have
332 MIC values less than or equal to 4 called as ceftazidime-sensitive strains. We obtained beta-
333 lactamase corresponding to each strain; a beta-lactamase corresponding to resistant strain is
334 called resistant beta-lactamase and a beta-lactamase corresponding to sensitive strain is called
335 sensitive beta-lactamase. Our final dataset have sensitive and resistant variants of beta-
336 lactamase.

337 Amino acid composition analysis revealed that certain residues like Alanine, Glycine,
338 Leucine, Proline and Arginine are more frequent in ceftazidime resistant beta-lactamases as
339 compared to ceftazidime sensitive ones. Similarly, in case of ceftazidime sensitive beta-
340 lactamases the residues like Aspartic acid, Isoleucine, Lysine, Asparagine, Threonine and
341 Tyrosine are more in abundance in comparison to resistant beta-lactamases. From these
342 findings it can be inferred that in ceftazidime-resistant beta-lactamases, amino acids with
343 non-polar side chains predominates. No wonder this gives these resistant beta-lactamases
344 extra stability making it hard for ceftazidime to inhibit their activity. In case of ceftazidime-
345 sensitive beta-lactamases, amino acids with polar side chains predominates. Moreover, amino
346 acid with charged entities is more in number in this case. This makes these proteins quite
347 unstable and prone to attack by the antibiotic.

348 In this study, Pfeature software has been used to compute different types of descriptors that
349 includes amino acid composition, dipeptide composition, residue entropy, repeats,
350 distribution of amino acids. In order to identify relevant features, we adopt different
351 techniques to remove useless features or descriptors. All descriptors having low variance has
352 been removed as they are not suitable for classification. Highly correlated or redundant has
353 been removed to decrease the noise. Finally algorithms in Scikit-learn has been used for
354 selecting important descriptors for developing prediction models. We employed different
355 machine learning algorithms using python-library-scikit-learn. We implemented widely used
356 machine learning classifiers, like KNN, DT, RF, GNB, LR, SVC and XGB [30]. In order to
357 our models we used internal and external validations [31] [32]. The result of the generated

358 model was analysed using various parameters called as threshold-dependent parameters and
359 threshold-independent parameters [33] [34]. We also validated the sturdiness of our model by
360 cross checking the resistance of 22 ceftazidime resistance beta-lactamases downloaded from
361 RGI database. Our model correctly predicted 19 ceftazidime resistance strains out of 22. We
362 hold an opinion that this method will be very helpful in prior prediction of ceftazidime
363 resistance/susceptibility towards any newly emerging strain of beta-lactamases. This also
364 open vista for researchers to look for alternative therapeutic options to fight continuously
365 emerging beta-lactamases. The method also has a major utility in doing prediction of
366 sensitive beta-lactamase strains in metagenomics data.

367 **Conclusion**

368 In conclusion, this is the first study of resistance/sensitivity prediction model development
369 using one particular antibiotic. The study brings about in-silico model to predict
370 resistance/susceptibility of ceftazidime antibiotic towards beta-
371 lactamases(<http://webs.iiitd.edu.in/raghava/abcrpred/>). This will help in identification of
372 ceftazidime sensitive beta-lactamases strains. Prediction can be done even when only protein
373 sequence of any beta-lactamase is known. We believe in future, researchers will build similar
374 model for other antibiotics. Prior prediction of sensitive antibiotics against a bacterial
375 infection will lead to era of strain-specific antibiotics; basically, end of present hit and trial
376 era. This will reduce time and cost of treatment as well a significant reduction in side-effects
377 due to the treatment by inappropriate antibiotics.

378 **Conflict of Interest**

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

380 **Author Contributions**

381 LM, SSU, AD, and SP collected and processed the datasets. LM, AD, SP, SSU and GPSR
382 implemented the algorithms. SP and AD developed the prediction models. LM, AD, SP, SSU
383 and GPSR analysed the results. SP, NS and AD created the back-end of the web server and
384 front-end user interface. LM, AD, SP, SSU, NS and GPSR penned the manuscript. GPSR
385 conceived and coordinated the project and gave overall supervision to the project. All authors
386 have read and approved the final manuscript.

387 **Acknowledgement**

388 Authors are thankful to J.C. Bose National Fellowship, Department of Science and
389 Technology (DST), Government of India, and DST-INSPIRE, NPDF-SERB, and DBT for
390 fellowships and the financial support.

391 **Data Availability Statement**

392 All the datasets generated for this study are either included in this article/Supplementary
393 material or available at the “ABCRpred” webserver,
394 <https://webs.iiitd.edu.in/raghava/abcrpred/download.php> as mentioned in the Materials and
395 Methods section.

396 **References**

- 397 1. Munita JM, Arias CA. Mechanisms of Antibiotic Resistance. *Microbiol. Spectr.* 2016; 4(2):481-
398 511.
- 399 2. Thakuria B, Lahon K. The Beta Lactam Antibiotics as an Empirical Therapy in a Developing
400 Country: An Update on Their Current Status and Recommendations to Counter the Resistance
401 against Them. *J. Clin. Diagn. Res.* 2013; 7(6):1207–1214.
- 402 3. Tooke CL, Hinchliffe P, Bragginton EC, et al. beta-Lactamases and beta-Lactamase Inhibitors in
403 the 21st Century. *J. Mol. Biol.* 2019; 431(18):3472–3500.
- 404 4. Bush K. Past and Present Perspectives on beta-Lactamases. *Antimicrob. Agents Chemother.* 2018;
405 62(10).
- 406 5. Naas T, Oueslati S, Bonnin RA, et al. Beta-lactamase database (BLDB) - structure and function. *J.*
407 *Enzyme Inhib. Med. Chem.* 2017; 32(1):917–919.
- 408 6. Canton R, Novais A, Valverde A, et al. Prevalence and spread of extended-spectrum beta-
409 lactamase-producing Enterobacteriaceae in Europe. *Clin. Microbiol. Infect.* 2008; 14:144–153.
- 410 7. Maryam L, Usmani SS, Raghava GPS. Computational resources in the management of antibiotic
411 resistance: Speeding up drug discovery. *Drug Discov. Today* 2021.
- 412 8. Aytan-Aktug D, Clausen PTLC, Bortolaia V, et al. Prediction of Acquired Antimicrobial
413 Resistance for Multiple Bacterial Species Using Neural Networks. *mSystems* 2020; 5(1).
- 414 9. Ransom EM, Potter RF, Dantas G, et al. Genomic Prediction of Antimicrobial Resistance: Ready
415 or Not, Here It Comes! *Clin. Chem.* 2020; 66(10):1278–1289.
- 416 10. Kim J, Greenberg DE, Pifer R, et al. VAMPr: VARIant Mapping and Prediction of antibiotic
417 resistance via explainable features and machine learning. *PLoS Comput. Biol.* 2020;
418 16(1):e1007511.
- 419 11. Nguyen M, Olson R, Shukla M, et al. Predicting antimicrobial resistance using conserved genes.
420 *PLoS Comput. Biol.* 2020; 16(10):e1008319.
- 421 12. Gupta S, Chaudhary K, Kumar R, et al. Prioritization of anticancer drugs against a cancer using
422 genomic features of cancer cells: A step towards personalized medicine. *Sci. Rep.* 2016; 6(1):1-11.

- 423 13. Dhanda SK, Vir P, Singla D, et al. A Web-Based Platform for Designing Vaccines against
424 Existing and Emerging Strains of Mycobacterium tuberculosis. *PLoS One* 2016; 11(4):e0153771
- 425 14. Gupta S, Chaudhary K, Dhanda SK, et al. A Platform for Designing Genome-Based Personalized
426 Immunotherapy or Vaccine against Cancer. *PLoS One* 2016; 11(11):e0166372.
- 427 15. Jamal S, Khubaib M, Gangwar R, et al. Artificial Intelligence and Machine learning based
428 prediction of resistant and susceptible mutations in Mycobacterium tuberculosis. *Sci. Rep.* 2020;
429 10(1):1-16.
- 430 16. Keshri V, Diene SM, Estienne A, et al. An Integrative Database of beta-Lactamase Enzymes:
431 Sequences, Structures, Functions, and Phylogenetic Trees. *Antimicrob. Agents Chemother.* 2019;
432 63(5).
- 433 17. Forest Laboratories, LLC, at 1-800-678-1605 or FDA at 1-800-FDA-1088 or
434 www.fda.gov/medwatch. Revised: 02/2015.
- 435 18. European Committee On Antimicrobial Susceptibility Testing Breakpoint Tables For
436 Interpretation Of Mics And Zone Diameters Version 10.0, valid from 2020-01-01.
- 437 19. Alcock BP, Raphenya AR, Lau TTY, et al. CARD 2020: Antibiotic Resistome Surveillance with
438 the Comprehensive Antibiotic Resistance Database. *Nucleic Acids Res.* 2020; 48(D1):D517–D525.
- 439 20. Pande A, Patiyal S, Lathwal A, et al. Computing wide range of protein/peptide features from their
440 sequence and structure. *bioRxiv* 2019;599126.
- 441 21. Pedregosa F, Varoquaux G, Gramfort A, et al. Scikit-learn: Machine Learning in Python. *J. Mach.*
442 *Learn. Res.* 2011; 12:2825–2830.
- 443 22. Hussain HM, Benkrid K, Seker H. Dynamic partial reconfiguration implementation of the
444 SVM/KNN multi- classifier on FPGA for bioinformatics application. *Annu. Int. Conf. IEEE Eng.*
445 *Med. Biol. Soc. IEEE Eng. Med. Biol. Soc. Annu. Int. Conf. 2015*; 2015:7667–7670.
- 446 23. Mushtaq M-S, Mellouk A. 2 - Methodologies for Subjective Video Streaming QoE Assessment.
447 *Quality of Experience Paradigm in Multimedia Services*, 27–57.
- 448 24. Tolles J, Meurer WJ. Logistic Regression: Relating Patient Characteristics to Outcomes. *JAMA*
449 2016; 316(5):533–534.
- 450 25. Agrawal P, Bhalla S, Chaudhary K, et al. In Silico Approach for Prediction of Antifungal
451 Peptides. *Front. Microbiol.* 2018; 9:323.
- 452 26. Dhall A, Patiyal S, Sharma N, et al. Computer-aided prediction and design of IL-6 inducing
453 peptides: IL-6 plays a crucial role in COVID-19. *Brief. Bioinform.* 2021; 22(2):936–945.
- 454 27. Bhalla S, Kaur H, Dhall A, et al. Prediction and Analysis of Skin Cancer Progression using
455 Genomics Profiles of Patients. *Sci. Rep.* 2019; 9(1):1-16.
- 456 28. Patiyal S, Agrawal P, Kumar V, et al. NAGbinder: An approach for identifying N-
457 acetylglucosamine interacting residues of a protein from its primary sequence. *Protein Sci.* 2020;
458 29(1):201–210.
- 459 29. Sharma N, Patiyal S, Dhall A, et al. AlgPred 2.0: an improved method for predicting allergenic
460 proteins and mapping of IgE epitopes. *Brief. Bioinform.* 2020.

- 461 30. Ivanescu AE, Li P, George B, et al. The importance of prediction model validation and assessment
462 in obesity and nutrition research. *Int. J. Obes. (Lond)*. 2016; 40(6):887–894.
- 463 31. Debray TPA, Vergouwe Y, Koffijberg H, et al. A new framework to enhance the interpretation of
464 external validation studies of clinical prediction models. *J. Clin. Epidemiol.* 2015; 68(3):279–289.
- 465 32. Nagpal G, Usmani SS, Dhanda SK, et al. Computer-aided designing of immunosuppressive
466 peptides based on IL-10 inducing potential. *Sci. Rep.* 2017; 7(1):1-10.
- 467 33. Kumar V, Agrawal P, Kumar R, et al. Prediction of cell-penetrating potential of modified peptides
468 containing natural and chemically modified residues. *Front. Microbiol.* 2018; 9:725.
- 469 34. Robin X, Turck N, Hainard A, et al. pROC: an open-source package for R and S+ to analyze and
470 compare ROC curves. *BMC Bioinformatics* 2011; 12(1):1-8.