

This is the peer reviewed version of the following article:

Chowdhury, A. Sayed, E. Khaledian, and S.L. Broschat, (2019). Capreomycin resistance prediction in two species of Mycobacterium using a stacked ensemble method. Journal of Applied Microbiology, Vol.127, No.6, which has been published in final form at <https://sfamjournals.onlinelibrary.wiley.com/doi/full/10.1111/jam.14413>. This article may be used for non-commercial purposes in accordance with Wiley Terms and Conditions for Use of Self-Archived Versions

Capreomycin Resistance Prediction in Two Species of *Mycobacterium* Using a Stacked Ensemble Method

Abu Sayed Chowdhury¹, Ehdieh Khaledian¹, and Shira L. Broschat^{1,2,3}

¹School of Electrical Engineering and Computer Science,

²Paul G. Allen School for Global Animal Health,

³Department of Veterinary Microbiology and Pathology,

Washington State University,

Pullman, Washington, United States

E-mail: {abu.chowdhury, ehdieh.khaledian, shira}@wsu.edu

Corresponding author: Abu Sayed Chowdhury

Email: abu.chowdhury@wsu.edu

Phone number: +1 509-715-9463

Address: School of Electrical Engineering and Computer Science, Washington State University,

P.O. Box 642752, Pullman, Washington 99164-2752, United States

1 **Abstract**

2 **Aims:** Predicting bacterial resistance provides valuable information that can assist in clinical
3 decisions. With recent advances in whole genome sequencing technology, the detection of
4 antibiotic resistance (AR) proteins directly from genomic data is becoming feasible. AR
5 genes/proteins can be identified using best-hit methods that work by comparing candidate
6 sequences with known AR genes in public databases. However, these approaches may fail to
7 detect resistance genes with sequences that differ significantly from known sequences. Our goal
8 is to develop a machine learning technique to accurately predict capreomycin resistance in
9 *Mycobacteria* with low false discovery rates.

10 **Methods and Results:** We present a stacked ensemble learning model as an alternative to
11 traditional DNA sequence alignment-based methods using optimal features generated from the
12 physicochemical, evolutionary, and secondary structure properties of protein sequences. We
13 train logistic regression, C5.0, and support vector machine (SVM) algorithms as our base
14 classifiers, and our stacked ensemble predictors combine the results from the base classifiers to
15 achieve higher accuracy. Compared with our most accurate base classifier (SVM), our most
16 accurate stacked ensemble predictor increases training accuracy by 2.43%. Our stacked
17 ensemble predictors achieve test accuracy up to 81.25%.

18 **Conclusions:** We developed a stacked ensemble model to predict capreomycin resistance for
19 *Mycobacteria* with an accuracy greater than 80% using protein sequences with sequence
20 similarity ranging between 10% and 70%. This performance cannot be achieved with best-hit
21 methods due to differences in sequence similarity.

22 **Significance and Impact of the Study:** Today an estimated one-half million cases of multidrug-

23 resistant (MDR) and extensively-drug resistant (XDR) tuberculosis (TB) occur annually worldwide
24 at a great cost. Because capreomycin is a second-line drug used to treat drug resistant TB, the
25 ability to use a machine learning approach to classify capreomycin-resistant TB in a timely
26 manner is crucial for the successful treatment of MDR or XDR TB.

27 **Keywords:** Capreomycin resistance; antibiotic resistance; tuberculosis; physicochemical
28 features; secondary structure features; feature selection; ensemble learning; machine learning

29

30 **Introduction**

31 The introduction of antibiotics has saved millions of lives over the past decades. However, their
32 widespread use has resulted in the rapid emergence of antibiotic resistance (AR) (Aminov 2009,
33 Dijkstra et al. 2018). AR impacts the safety and efficacy of treatment options that require the use
34 of antibiotics and leads to the risk of fatal outcomes. Therefore, various strategies such as
35 reduction in antibiotic usage or alternative second-line antibiotics have been suggested. One of
36 the major AR bacteria that is a global health threat is *Mycobacterium tuberculosis* (MTB). MTB
37 bacteria cause the life-threatening, airborne infectious disease, tuberculosis (TB). TB has infected
38 almost one-third of the world population. It usually attacks lungs, but it can affect the brain,
39 kidneys, and spine as well. It can be multidrug resistant (MDR) due to resistance to first-line TB
40 antibiotics such as isoniazid and rifampin, and it also can be extensively drug resistant (XDR) to
41 second-line antibiotics such as capreomycin, amikacin, and kanamycin (Gygli et al. 2017). A
42 species known as *M. yongonense* TTK-01-0059 (originally *M. sp.* TTK-0100059), which causes
43 pulmonary infections in humans, has been found to be resistant to capreomycin. This species has
44 related strains *M. yongonense* 05-1390(T), Asan 36912, and Asan 36527 which may also be
45 resistant to capreomycin, but their genomes are not available for inspection (Mnyambwa et al.

46 2018).

47 AR detection in *Mycobacterium* spp. can provide valuable information for improving clinical
48 decisions in the treatment of the infectious diseases they cause, such as identifying drugs with
49 the potential for success or eliminating a **drug that is ineffective due to antibiotic resistance**. As
50 such, AR prediction has become a fundamental challenge for improving health care. Traditional
51 *in vitro* culture-based tests for AR phenotype detection take many days to provide susceptibility
52 test results (Chen et al. 2018), and phenotyping using the gold-standard proportion method can
53 take up to two months (Niehaus et al. 2014). A molecular approach is an alternative to
54 conventional culture-based detection; however, it often fails to detect rare gene variants and
55 provides limited biomarkers for AR detection. Whole genome sequencing is currently a feasible
56 approach for the detection of target mutations, and sequence alignment techniques such as best-
57 hit approaches can be applied to detect AR genes using sequence identity in existing online
58 databases (Kleinheinz et al. 2014, Davis et al. 2016, Yang et al. 2016). These methods work well
59 for finding known and highly conserved AR sequences and have low probability of producing false
60 positives, *i.e.*, predicting non-AR genes as AR genes (Forsberg et al. 2014). However, they may
61 produce high false negative rates when they fail to identify AR genes that have lower sequence
62 identity with known AR genes (Xavier et al. 2016, Yang et al. 2016, McArthur et al. 2017). Machine
63 learning can be a useful alternative computational framework for identifying AR phenotypes
64 accurately using features, *i.e.*, characteristics of known AR protein sequences, found in the
65 genomic data regardless of sequence similarity (Niehaus et al. 2014, Santerre et al. 2016).
66 Machine learning models require a large number of genomes with AR metadata to make a strong
67 prediction. Databases such as the Pathosystems Resource Integration Center (PATRIC) (Wattam

68 et al. 2013, Davis et al. 2016, Wattam et al. 2016) and the Antibiotic Resistance Genes Database
69 (ARDB) (Lio and Pop 2008) currently provide genomes with AR metadata. In some recent studies
70 (Her and Wu 2018; Kavvas et al. 2018; Moradigaravand et al. 2018), machine learning techniques
71 were applied to find resistance to varieties of antibiotics from whole genome sequences using
72 databases of known AR genes. These studies first constructed a pan-genome with identifying core
73 and accessory gene clusters, then computed features such as presence-absence patterns and the
74 population structures of gene clusters, and finally applied a machine learning algorithm to
75 identify putative AR genes using the extracted features. However, these approaches used a small
76 number of genetic features for AR gene prediction.

77 In this work, we developed a machine learning approach to detect capreomycin resistance in
78 two species of *Mycobacterium* which may be extended to other species. We collected protein
79 sequences for both antibiotic resistance and antibiotic susceptibility with sequence similarity
80 ranging between 10% and 70% from the two databases mentioned above. We then extracted
81 potential features including physicochemical, evolutionary, and structural properties and
82 evaluated these features to find the optimal ones for prediction. The advantage of using protein
83 sequences rather than gene sequences is that a greater number of features are discernible. We
84 applied a stacked ensemble learning approach to predict capreomycin resistance. To our
85 knowledge, this is the first time stacked ensemble learning has been used for this task. Our
86 stacked ensemble predictor trained a learning algorithm to combine the results from logistic
87 regression, C5.0, and SVM algorithms. Finally, we divided our data into training and test sets to
88 validate the performance of our model. We also compared the performance of our classifier with
89 BLASTp results.

90

91 **Materials and Methods**

92 In this section, we discuss the details of data collection, feature generation from protein
93 sequences, and feature evaluation to find an optimal feature set.

94 **Data Collection**

95 To collect capreomycin resistance and susceptibility protein sequences, we searched the PATRIC
96 (Davis et al. 2016) database and ARDB (Lio and Pop 2008). The Pathosystems Resource
97 Integration Center (PATRIC) is a resource center designed to store and integrate a variety of data
98 types, e.g., genomic sequences, three-dimensional protein structures, and sequence typing data.
99 The Antibiotic Resistance Genes Database (ARDB) is a manually curated database that provides a
100 centralized collection of information on antibiotic resistance. We retrieved 3,366 resistance and
101 5,411 susceptibility protein sequences. The vast majority of these were duplicate sequences. We
102 applied CD-HIT (Li and Godzik 2006, Fu et al. 2012) to remove sequences with similarity greater
103 than or equal to 70% from both sets of protein sequences. Our final result was 44 resistance and
104 36 susceptibility sequences (Supplementary Tables: Table S1, and Table S2). We used this dataset
105 to train and test our classification model.

106 **Feature Generation**

107 We performed a thorough literature search to determine potential protein features to consider
108 for use with our machine learning model. A number of researchers have suggested techniques
109 for extracting features from proteins based on composition, physical and chemical
110 characteristics, and secondary structure (Ding and Dubchak 2001, Cai et al. 2003, Zhang et al.

111 2011, Liu et al. 2013, Wei et al. 2015, Li et al. 2016, Eslami et al. 2018). In this work, we first
112 computed the amino acid composition for each sequence. Because there are 20 different amino
113 acids, this gives 20 features, each representing a fraction of the amino acids within a protein
114 sequence of length n . Next, we applied the composition, transition, and distribution (CTD) model
115 to obtain global physicochemical information from the protein sequences (Dubchak et al. 1995,
116 Dubchak et al. 1999). Amino acids are grouped into three classes using hydrophobicity properties.
117 These classes are polar, neutral, and hydrophobic. Arginine, lysine, glutamic acid, aspartic acid,
118 glutamine, and asparagine belong to the polar class. The neutral group contains glycine, alanine,
119 serine, threonine, proline, histidine, and tyrosine, and the remaining amino acids cysteine,
120 leucine, valine, isoleucine, methionine, phenylalanine, and tryptophan are in the hydrophobic
121 group. Composition C represents the number of amino acids of each class in a protein sequence
122 divided by the total number of amino acids. Transition T gives the percent frequency of one
123 amino acid class followed by amino acids of another class. The distribution D measures the
124 fraction of the whole sequence where the first, 25%, 50%, 75%, and 100% of the amino acids of
125 a class are located.

126 The CTD model considers eight physicochemical amino acid properties, namely, hydrophobicity,
127 normalized van der Waals volume, polarity, polarizability, charge, surface tension, secondary
128 structure, and solvent accessibility. Because there are three classes for amino acids, we obtain
129 three features for composition, three features for transition, and $3 \times 5 = 15$ features for
130 distribution for each of the eight physicochemical amino acid properties. Thus, in total, the CTD
131 model gives $(3 + 3 + 15) \times 8 = 168$ features for each protein sequence.

132 We also considered the position-specific scoring matrix (PSSM) to obtain evolutionary-

133 relevant features from protein sequences. To accomplish this, we generated the PSSM from
 134 protein sequences using PSI-BLAST (Altschul et al. 1997) and then calculated transition scores
 135 between neighboring amino acids to generate a 400D feature vector for each sequence. Let $M =$
 136 $(m_{i,j}) \in \mathbb{R}^{n \times 20}$ be the PSSM matrix where $m_{i,j}$ is the entry in the i th row and j th column of
 137 M . The transition $t_{x,y}$ between neighboring amino acids x and y can be computed using Eq.
 138 (1).

$$139 \quad t_{x,y} = \sum_{k=1}^{n-1} m_{k,x} m_{k+1,y} \quad (1)$$

140 Finally, we generated features from the secondary structure of a protein sequence which has
 141 proven to be very efficient in protein fold recognition. We applied PSIPRED (Jones 1999) which
 142 identifies the secondary structure of a protein sequence based on PSI-BLAST. The output of
 143 PSIPRED is a secondary structure sequence composed of $H(\alpha\text{-helix})$, $E(\beta\text{-strand})$, $C(\gamma\text{-coil})$,
 144 and an $n \times 3$ structure probability matrix where n is the number of amino acids in a protein
 145 sequence. Suppose $SPM = (p_{i,j}) \in \mathbb{R}^{n \times 3}$ is the structure probability matrix where each entry
 146 in the matrix represents the probability of an amino acid to be in one of the states H , E , or C .
 147 Then we extracted features from the secondary structure sequence and SPM as described in
 148 existing works (Kurgan and Homaeian 2006, Kurgan et al. 2008, Liu and Jia 2010, Zhang et al.
 149 2011, Wei et al. 2015). We generated location-oriented features CMV_H , CMV_E , and CMV_C by
 150 measuring the spatial arrangements of the three states— H , E , and C —using Eqs. (2), (3), and
 151 (4), respectively.

$$152 \quad CMV_H = \frac{\sum_{j=1}^{T_H} I_{Hj}}{n(n-1)} \quad (2)$$

$$153 \quad CMV_E = \frac{\sum_{j=1}^{T_E} I_{Ej}}{n(n-1)} \quad (3)$$

$$154 \quad CMV_C = \frac{\sum_{j=1}^{T_C} I_{C_j}}{n(n-1)} \quad (4)$$

$$155 \quad Max_H = \frac{\max\{n_H\}}{n} \quad (5)$$

$$156 \quad Max_E = \frac{\max\{n_E\}}{n} \quad (6)$$

157 where I_{H_j} , I_{E_j} , and I_{C_j} are the position indexes, and T_H , T_E , and T_C are the total number of
 158 states for H , E , and C , respectively. We computed the normalized maximum spatial
 159 consecutive E and H states in the secondary structure sequence using Eqs. (5) and (6) where
 160 Max_H and Max_E are the maximum lengths of consecutive H and E states in the secondary
 161 structure sequence, respectively. We also computed the segmented sequences by ignoring the
 162 coil segments in the secondary structure. For example, let $S =$
 163 $CCCEEEEEHHHCCEECCCHHHHCCEEE$ be the secondary structure sequence. If we ignore
 164 the coil segments C and combine consecutive H 's and E 's as H' and E' , then we obtain
 165 $S' = E'H'E'H'E'$. As α and β are usually separated and interspersed in α/β proteins and
 166 $\alpha + \beta$ proteins, respectively (Zhang et al. 2011, Wei et al. 2015), we measured the frequency of
 167 $E'H'E'$ (denoted as $f_{E'H'E'}$) using Eq. (7) where $T_{E'H'E'}$ is the total number of occurrences of
 168 $E'H'E'$ in S' , and the length of S' is $L_{S'}$.

$$169 \quad f_{E'H'E'} = \frac{T_{E'H'E'}}{L_{S'} - 2} \quad (7)$$

170 After extracting features from the secondary structure sequence, we retrieved features from the
 171 structure probability matrix SPM . We obtained three global information features from the
 172 SPM using Eqs. (8), (9), and (10).

$$173 \quad F_{g_1} = \frac{1}{n} \sum_{j=1}^n p_{j,1} \quad (8)$$

$$174 \quad F_{g_2} = \frac{1}{n} \sum_{j=1}^n p_{j,2} \quad (9)$$

$$175 \quad F_{g_3} = \frac{1}{n} \sum_{j=1}^n p_{j,3} \quad (10)$$

176 Here, $p_{j,1}$, $p_{j,2}$, and $p_{j,3}$ are the j th probability values of the *SPM* for H , E , and C ,
 177 respectively. We also captured local information features by subdividing the *SPM* into smaller
 178 δ matrices each with $\lfloor \frac{n}{\delta} \rfloor$ rows and three columns. We chose $\delta = 8$ and computed 3D features
 179 for a particular sub-matrix q as shown in Eq. (11). Note that $F_{local_q}(1)$, $F_{local_q}(2)$, and
 180 $F_{local_q}(3)$ are computed using the same formulas given in Eqs. (8), (9), and (10), respectively. In
 181 this way, we found a total of 24 features with local information.

$$182 \quad F_q = \{F_{local_q}(1), F_{local_q}(2), F_{local_q}(3)\}, 1 \leq q \leq \delta \quad (11)$$

183 In summary, we obtained a 20D feature vector based on amino acid composition, a 168D
 184 feature vector of physicochemical properties using a CTD model, a 400D feature vector using
 185 evolutionary information from the PSSM, a 6D feature vector having location oriented features
 186 of secondary structure, and a 27D feature vector having global and local information from the
 187 structure probability matrix (SPM) of a protein sequence. We merged all of these features to
 188 obtain a 621D high-dimensional feature vector as summarized in Table 1. In the next section, we
 189 explain how we removed redundant and noisy information to obtain our final feature set.

190 Dimensionality Reduction

191 To reduce the dimension of our features, we first calculated Pearson's correlation coefficient
 192 between features, that is, we measured the linear correlation between two feature vectors u
 193 and v . Pearson's correlation coefficient $\rho_{u,v}$ is calculated using

$$194 \quad \rho_{u,v} = \frac{E[(u-\mu_u)(v-\mu_v)]}{\sigma_u \sigma_v} \quad (12)$$

195 where E is the expectation, μ_u and μ_v are the mean values, and σ_u and σ_v are the
196 standard deviations of u and v , respectively. The $\rho_{u,v}$ value can be $[-1, +1]$ where the signs
197 '+' and '-' indicate positive and negative correlations between features. We considered the
198 absolute value of $|\rho_{u,v}|$ to estimate the correlation between features. When the value of
199 $|\rho_{u,v}|$ is large, then two features are highly correlated, and we can consider either of the
200 features for classification purposes. We set a threshold value of 0.95, keeping one feature and
201 eliminating the other. In this way, we removed redundant information in the feature set and
202 reduced the number of features from 621 to 392.

203 Next we performed a standard t -test using our dataset of resistant and susceptible protein
204 sequences together with our feature set. The t -test is used to check whether two populations
205 are significantly different based on a statistical measurement. In the t -test, the null hypothesis
206 is that there is no difference between the two populations (resistance and susceptibility protein
207 sequences), and we measure the p -values for all the features to see whether we can reject the
208 null hypothesis. A low p -value for a feature means that the feature is important for predicting
209 capreomycin resistance proteins, and we can reject the null hypothesis. To further reduce the
210 dimension of the feature set, we set a threshold p -value of 0.05, eliminating features with $p >$
211 0.05. Figure 1 depicts the p -values for selected features both in linear and logarithmic scales.
212 We found that the top selected features were from the distribution measurements of the CTD
213 model which clearly indicates the importance of this feature group in our prediction model.

214 We utilized the R *stats* package (version 3.4.3) to measure correlation and to perform t -tests
215 for all the features. Following the filtering of features through this t -test step, our feature vector
216 was further reduced from 392D features to 336D. We considered this 336D feature vector as our

217 optimal feature set. In the following section, we describe how we used our optimal features to
218 develop a classification model for recognizing capreomycin resistance proteins.

219 **Stacked Ensemble Model**

220 For our machine learning model, we applied stacking (also called stacked generalization) to
221 increase accuracy. We measured the accuracy of a classifier using Eq. (13) where FP , FN , TP ,
222 and TN are false positive, false negative, true positive, and true negative, respectively.

$$223 \quad Accuracy = \frac{TP+TN}{TP+FP+TN+FN} \quad (13)$$

224 Stacking is an efficient ensemble machine learning technique that takes the predictions of
225 different machine learning algorithms in the first layer as inputs to a second layer (Wolpert 1992,
226 Breiman 1996, Van der Laan et al. 2007, Chowdhury et al. 2015). It combines the predictions of
227 the first-layer models together using another classifier as depicted in Figure 2. In the stacked
228 ensemble model, the second layer classifier (also called a meta-classifier or meta-learner) is
229 trained using the predictions of the first-layer models (also called base classifiers) to make an
230 optimal combination of the predictions.

231 As shown in Figure 2, we used three different models as base classifiers in our first layer:
232 logistic regression via a generalized linear model (GLM), a C5.0 model, and a support vector
233 machine (SVM). The GLM is a statistical linear model which generalizes multiple linear regression
234 for more than one dependent variable (Madsen and Thyregod 2010). C5.0 is an algorithm for
235 creating decision trees, which supports boosting (Golmah 2014). An SVM is a supervised learning
236 model formally characterized by a separating hyperplane that divides binary data (Scholkopf
237 2001). These three classifiers can accommodate multiple variables, including mixtures of
238 categorical and continuous variables. They accept the same dataset as input and provide

239 individual prediction outputs PG , PC , and PS , respectively. The meta-classifier in the second
240 layer delivers the final prediction PM based on the predictions of the base classifiers PG , PC ,
241 and PS . We tested several different meta-classifiers: a generalized linear model (GLM), linear
242 discriminate analysis (LDA), and a random forest (RF). LDA is a classification method that searches
243 for a linear combination of variables that best divides two classes. It is mathematically robust and
244 often produces models with accuracy as high as more complex classifiers (Balakrishnama and
245 Ganapathiraju, 1998). An RF is an ensemble of decision trees. It uses a bagging method to
246 combine the ensemble learning methods in order to obtain a more accurate and stable prediction
247 and to improve the overall result (Breiman 2001).

248 Note that in order to combine the results of the base classifiers, the pairwise correlations
249 between their predicted outputs should be low. If they are low (< 0.75), the stacked ensemble
250 model is likely to have better accuracy than a base classifier alone. If the correlations are not low,
251 the base models will provide the same or very similar predictions making the ensemble of
252 classifiers worthless.

253 To build the stacked ensemble model, we used six R packages: *mlbench* (version 2.1-1), *caret*
254 (version 6.0-79), *caretEnsemble* (version 2.0.0), *C50* (version 0.1.2), *plyr* (version 1.8.4), and
255 *randomForest* (version 4.6-14). In the next section, we provide test results and show that the
256 stacked ensemble technique is superior in predicting capreomycin resistance in *Mycobacterium*
257 spp. rather than applying our machine learning algorithms individually.

258 Results

259 For classification, we divided the dataset into training and test sets. We considered stratified
260 sampling where 80% of the dataset is used as the training set and the remaining 20% is used as

261 the test set. First we trained our base classifiers—the generalized linear model (GLM), C5.0
262 algorithm, and SVM—using the training set and then fed the predictions obtained from this step
263 to the meta-classifier. Next we applied this stacked model to predict the resistance and
264 susceptibility of our test dataset. We tuned the C5.0 and SVM models based on our training
265 dataset to find the best parameters and selected *trials* = 1, *model* = rules, *winnnow* = true for the
266 C5.0 model, and the radial basis kernel function in the SVM with C and σ parameter values of
267 1 and 0.002, respectively. We considered 10 repeats of 10-fold cross validation to train the model
268 using 100 resamples for the training dataset.

269 Figure 3 depicts the accuracy statistics for all simulation runs for the three base classifiers.
270 The box-and-whisker plots in Figure 3 show the accuracy statistics for the training dataset where
271 each box represents the median (indicated by a solid vertical line) and lower and upper quartiles,
272 and each whisker gives the lowest and highest observations. Note that a small circle beyond the
273 left whisker in Figure 3 indicates an outlier that is far from the lowest value. By calculating the
274 average training accuracy values, we found that the SVM achieved the most accurate result
275 (78.95%), much more accurate than those of the GLM and C5.0, 50.25% and 54.20%, respectively.
276 Based on the training accuracy, the SVM is identified as a strong base classifier for our dataset.
277 As stated earlier, we can combine the predictions of all three models to increase the classification
278 accuracy if pairwise predictions of the classifiers have low correlation. The histograms in Figure
279 4 represent the accuracies of all simulation runs for each classifier, and the scatter plots show
280 pairwise comparisons of the accuracies obtained from the three classifiers. The fitted line in each
281 scatter plot gives the positive or negative correlation of all accuracies (indicated by small circles)
282 between any two classifiers.

AR Prediction Using Ensemble Learning

283 It is clear from the scatter plots in Figure 4 that the predicted outputs of the base classifiers
284 have low correlation. The correlation between the GLM and C5.0 results is 0.094 (positive
285 correlation). The GLM and SVM results have a correlation of 0.046 (negative correlation), and the
286 C5.0 and SVM results have a correlation value of 0.024 (negative correlation). As all the
287 correlations are much less than 0.75, combining the predictions using a meta-classifier is likely to
288 increase the classification accuracy.

289 The training and test accuracies obtained using the stacked ensemble model with the three
290 different meta-classifiers are listed in Table 2. We tuned the RF meta-classifier and selected an
291 *mtry* parameter value of 2.

292 The meta-classifiers provide significant improvement in training accuracy. In particular, when
293 we combine the predictions of the base classifiers using the GLM meta-classifier, the training
294 accuracy is improved to 81.38%, an increase in training accuracy of approximately 2.43% from
295 the best base model alone. We applied our stacked ensemble model to make predictions on our
296 test dataset and achieved test accuracies up to 81.25%. The confusion matrices for the test
297 dataset for the meta-classifiers are given in supplementary tables (Table S3-Table S5).

298 We also compared our stacked ensemble approach with the performance of BLASTp using
299 default parameter settings (<https://blast.ncbi.nlm.nih.gov/Blast.cgi?PAGE=Proteins>). The results
300 depicted in Figure 5 are the number of matched AR protein sequences as a function of percent
301 identity threshold in test data using the AR samples from the training data. For example, if we
302 consider percent identity $\geq 50\%$ for a test protein sequence to be matched (true positive) with a
303 training AR sequence, we obtain five of nine true positives. Importantly, in order to achieve the
304 same number of true positives as for our GLM meta-classifier, the percent identity threshold of

305 **BLASTp must be 33%**, but the “cost” comes in the form of increasing false positives. For example,
306 when we set the percent identity threshold to 33%, all non-AR test sequences are incorrectly
307 identified as AR sequences, resulting in a high false positive rate. Thus, there is a big tradeoff with
308 the threshold setting when using a best-hit method such as BLASTp.

309 **Discussion**

310 Prediction of capreomycin resistance in *Mycobacterium* spp. is a challenging but important
311 problem. In this paper we introduced the use of machine learning with protein sequences which
312 allows a rich feature set to be used to predict capreomycin resistance in *Mycobacterium* spp. We
313 based our classifier on a stacked ensemble machine learning technique which uses two layers of
314 classifiers to achieve higher accuracy. We gathered 44 unique resistance and 36 unique
315 susceptibility protein sequences and filtered 621 original features, removing redundant and less
316 informative features, to obtain a set of 336 features. We considered amino acid composition,
317 physicochemical properties, evolutionary-relevant properties, and the secondary structure of
318 proteins to generate features.

319 To validate the effectiveness of our approach, we built a stacked ensemble classifier using a
320 generalized linear model (GLM), a C5.0 model, and an *SVM* as base classifiers and GLM, linear
321 discriminant analysis, and a random forest as individual meta-classifiers. Our repeated ten-fold
322 cross validation results showed that stacking can increase the training accuracy up to 2.43%
323 compared to the best base classifier and can achieve test accuracy up to 81.25%. The identity
324 threshold must be quite low to obtain similar classification results with BLASTp, and this leads to
325 a marked increase in false positives, i.e., sequences identified as AR that are not. Although we
326 considered only capreomycin resistance, our **approach could potentially** be applied to identify

327 other antibiotic resistance such as amikacin and kanamycin resistance in *Mycobacterium* spp. as
328 well as antibiotic resistance in other bacteria. The stacked ensemble model with the GLM meta-
329 classifier gives superior performance for our training and test datasets. While the accuracy of our
330 classifier is reasonable, i.e., greater than 80%, given the availability of more sequences to use for
331 training, the accuracy would be higher. In the future, when a greater number of AR sequences
332 are available, we will retrain our algorithm.

333 **Conflict of Interest**

334 We declare that no conflict of interest exists.

335

336 **Acknowledgements**

337

338 We wish to thank the reviewers for their careful reading of the manuscript and for their
339 thoughtful comments and suggestions. Their contributions have resulted in a stronger
340 publication.

341

342 **References**

343 Altschul, S.F., Madden, T.L., Schäffer, A.A., Zhang, J., Zhang, Z., Miller, W. and Lipman, D.J.,
344 1997. Gapped BLAST and PSI-BLAST: a new generation of protein database search programs.
345 *Nucleic Acids Res*, 25(17), pp.3389-3402.

346 Aminov, R.I., 2009. The role of antibiotics and antibiotic resistance in nature.
347 *Environ Microbiol*, 11(12), pp.2970-2988.

348 Balakrishnama, S., and Ganapathiraju, A., 1998. Linear discriminant analysis-a brief
349 tutorial. *Institute for Signal and information Processing*, 18, 1-8.

350 Breiman, L., 1996. Stacked regressions. *Machine learning*, 24(1), pp.49-64.

351 Breiman, L. 2001. Random forests. *Machine learning*, 45(1), 5-32.

352 Cai, C.Z., Han, L.Y., Ji, Z.L., Chen, X. and Chen, Y.Z., 2003. SVM-Prot: web-based support vector
353 machine software for functional classification of a protein from its primary sequence. *Nucleic*
354 *Acids Res*, 31(13), pp.3692-3697.

355 Chen, M.L., Doddi, A., Royer, J., Freschi, L., Schito, M., Ezewudo, M., Kohane, I.S., Beam, A.
356 and Farhat, M., 2018. Deep learning predicts tuberculosis drug resistance status from genome

357 sequencing data. bioRxiv, p.275628.

358 Chowdhury, A.S., Alam, M.M. and Zhang, Y., 2015, August. A biomarker ensemble ranking
359 framework for prioritizing depression candidate genes. In 2015 IEEE Conference on
360 Computational Intelligence in Bioinformatics and Computational Biology (CIBCB) (pp. 1-6). IEEE.

361 Davis, J.J., Boisvert, S., Brettin, T., Kenyon, R.W., Mao, C., Olson, R., Overbeek, R., Santerre, J.,
362 Shukla, M., Wattam, A.R. and Will, R., 2016. Antimicrobial resistance prediction in PATRIC and
363 RAST. *Sci Rep*, 6, p.27930.

364 Dijkstra, J.A., van der Laan, T., Akkerman, O.W., Bolhuis, M.S., de Lange, W.C.M., Kosterink,
365 J.G., van der Werf, T.S., Alffenaar, J.W.C. and van Soolingen, D., 2018. In vitro susceptibility of
366 *Mycobacterium tuberculosis* to amikacin, kanamycin, and capreomycin. *Antimicrob Agents*
367 *Chemother*, 62(3), pp.e01724-17.

368 Ding, C.H. and Dubchak, I., 2001. Multi-class protein fold recognition using support vector
369 machines and neural networks. *Bioinformatics*, 17(4), pp.349-358.

370 Dubchak, I., Muchnik, I., Holbrook, S.R. and Kim, S.H., 1995. Prediction of protein folding class
371 using global description of amino acid sequence. *Proceedings of the National Academy of*
372 *Sciences*, 92(19), pp.8700-8704.

373 Dubchak, I., Muchnik, I., Mayor, C., Dralyuk, I. and Kim, S.H., 1999. Recognition of a protein
374 fold in the context of the SCOP classification. *Proteins: Structure, Function, and Bioinformatics*,
375 35(4), pp.401-407.

376 Eslami Manoochehri, H., Kadiyala, S. S., Birjandtalab, J., and Nourani, M., 2018. Feature
377 Selection to Predict Compound's Effect on Aging. *Proceedings of the 2018 ACM International*
378 *Conference on Bioinformatics, Computational Biology, and Health Informatics* (pp. 419-427).
379 ACM.

380 Forsberg, K.J., Patel, S., Gibson, M.K., Lauber, C.L., Knight, R., Fierer, N. and Dantas, G., 2014.
381 Bacterial phylogeny structures soil resistomes across habitats. *Nature*, 509(7502), p.612.

382 Fu, L., Niu, B., Zhu, Z., Wu, S. and Li, W., 2012. CD-HIT: accelerated for clustering the next-
383 generation sequencing data. *Bioinformatics*, 28(23), pp.3150-3152.

384 Golmah, V., 2014. An efficient hybrid intrusion detection system based on C5. 0 and
385 SVM. *International Journal of Database Theory and Application*, 7(2), 59-70.

386 Gygli, S.M., Borrell, S., Trauner, A. and Gagneux, S., 2017. Antimicrobial resistance in
387 *Mycobacterium tuberculosis*: mechanistic and evolutionary perspectives. *FEMS Microbiol Rev*,
388 41(3), pp.354-373.

389 Her, H.L. and Wu, Y.W., 2018. A pan-genome-based machine learning approach for predicting
390 antimicrobial resistance activities of the *Escherichia coli* strains. *Bioinformatics*, 34(13), pp.i89-

391 i95.

392 Jones, D.T., 1999. Protein secondary structure prediction based on position-specific scoring
393 matrices. *J Mol Biol*, 292(2), pp.195-202.

394 Kavvas, E.S., Catoi, E., Mih, N., Yurkovich, J.T., Seif, Y., Dillon, N., Heckmann, D., Anand, A.,
395 Yang, L., Nizet, V. and Monk, J.M., 2018. Machine learning and structural analysis of
396 *Mycobacterium tuberculosis* pan-genome identifies genetic signatures of antibiotic resistance.
397 *Nature communications*, 9(1), p.4306.

398 Kleinheinz, K.A., Joensen, K.G. and Larsen, M.V., 2014. Applying the ResFinder and
399 VirulenceFinder web-services for easy identification of acquired antibiotic resistance and *E. coli*
400 virulence genes in bacteriophage and prophage nucleotide sequences. *Bacteriophage*, 4(2),
401 p.e27943.

402 Kurgan, L., Cios, K. and Chen, K., 2008. SCPRED: accurate prediction of protein structural class
403 for sequences of twilight-zone similarity with predicting sequences. *BMC bioinformatics*, 9(1),
404 p.226.

405 Kurgan, L.A. and Homaeian, L., 2006. Prediction of structural classes for protein sequences
406 and domains—impact of prediction algorithms, sequence representation and homology, and test
407 procedures on accuracy. *Pattern Recognition*, 39(12), pp.2323-2343.

408 Li, W. and Godzik, A., 2006. Cd-hit: a fast program for clustering and comparing large sets of
409 protein or nucleotide sequences. *Bioinformatics*, 22(13), pp.1658-1659.

410 Li, Y.H., Xu, J.Y., Tao, L., Li, X.F., Li, S., Zeng, X., Chen, S.Y., Zhang, P., Qin, C., Zhang, C. and
411 Chen, Z., 2016. SVM-Prot 2016: a web-server for machine learning prediction of protein
412 functional families from sequence irrespective of similarity. *PloS one*, 11(8), p.e0155290.

413 Liu, B. and Pop, M., 2008. ARDB—antibiotic resistance genes database. *Nucleic Acids Res*,
414 37(suppl_1), pp.D443-D447.

415 Liu, B., Zhang, D., Xu, R., Xu, J., Wang, X., Chen, Q., Dong, Q. and Chou, K.C., 2013. Combining
416 evolutionary information extracted from frequency profiles with sequence-based kernels for
417 protein remote homology detection. *Bioinformatics*, 30(4), pp.472-479.

418 Liu, T. and Jia, C., 2010. A high-accuracy protein structural class prediction algorithm using
419 predicted secondary structural information. *J Theor Biol*, 267(3), pp.272-275.

420 Madsen, H., and Thyregod, P., 2010. Introduction to general and generalized linear models.
421 CRC Press.

422 McArthur, A.G. and Tsang, K.K., 2017. Antimicrobial resistance surveillance in the genomic
423 age. *Annals of the New York Academy of Sciences*, 1388(1), pp.78-91.

424 Mnyambwa, N.P., Kim, D.J., Ngadaya, E., Chun, J., Ha, S.M., Petrucka, P., Addo, K.K., Kazwala,
425 R.R. and Mfinanga, S.G., 2018. Genome sequence of Mycobacterium yongonense RT 955-2015
426 isolate from a patient misdiagnosed with multidrug-resistant tuberculosis: First clinical detection
427 in Tanzania. *Int J Infect Dis*, 71, pp.82-88.

428 Moradigaravand, D., Palm, M., Farewell, A., Mustonen, V., Warringer, J. and Parts, L., 2018.
429 Prediction of antibiotic resistance in *Escherichia coli* from large-scale pan-genome data. *PLoS*
430 *computational biology*, 14(12), p.e1006258.

431 Niehaus, K.E., Walker, T.M., Crook, D.W., Peto, T.E. and Clifton, D.A., 2014, June. Machine
432 learning for the prediction of antibacterial susceptibility in *Mycobacterium tuberculosis*. In *IEEE-*
433 *EMBS International Conference on Biomedical and Health Informatics (BHI)* (pp. 618-621). IEEE.

434 Santerre, J.W., Davis, J.J., Xia, F. and Stevens, R., 2016. Machine learning for antimicrobial
435 resistance. arXiv preprint arXiv:1607.01224.

436 Scholkopf, B., and Smola, A. J., 2001. Learning with kernels: support vector machines,
437 regularization, optimization, and beyond. MIT press.

438 Van der Laan, M.J., Polley, E.C. and Hubbard, A.E., 2007. Super learner. *Stat Appl Genet Mol*
439 *Biol*, 6(1).

440 Wattam, A.R., Abraham, D., Dalay, O., Disz, T.L., Driscoll, T., Gabbard, J.L., Gillespie, J.J.,
441 Gough, R., Hix, D., Kenyon, R. and Machi, D., 2013. PATRIC, the bacterial bioinformatics database
442 and analysis resource. *Nucleic Acids Res*, 42(D1), pp.D581-D591.

443 Wattam, A.R., Davis, J.J., Assaf, R., Boisvert, S., Brettin, T., Bun, C., Conrad, N., Dietrich, E.M.,
444 Disz, T., Gabbard, J.L. and Gerdes, S., 2016. Improvements to PATRIC, the all-bacterial
445 bioinformatics database and analysis resource center. *Nucleic Acids Res*, 45(D1), pp.D535-D542.

446 Wei, L., Liao, M., Gao, X. and Zou, Q., 2015. Enhanced protein fold prediction method through
447 a novel feature extraction technique. *IEEE Trans Nanobioscience*, 14(6), pp.649-659.

448 Wolpert, D.H., 1992. Stacked generalization. *Neural Netw*, 5(2), pp.241-259.

449 Xavier, B.B., Das, A.J., Cochrane, G., De Ganck, S., Kumar-Singh, S., Aarestrup, F.M., Goossens,
450 H. and Malhotra-Kumar, S., 2016. Consolidating and exploring antibiotic resistance gene data
451 resources. *Journal of clinical microbiology*, 54(4), pp.851-859.

452 Yang, Y., Jiang, X., Chai, B., Ma, L., Li, B., Zhang, A., Cole, J.R., Tiedje, J.M. and Zhang, T., 2016.
453 ARGs-OAP: online analysis pipeline for antibiotic resistance genes detection from metagenomic
454 data using an integrated structured ARG-database. *Bioinformatics*, 32(15), pp.2346-2351.

455 Zhang, S., Ding, S. and Wang, T., 2011. High-accuracy prediction of protein structural class for
456 low-similarity sequences based on predicted secondary structure. *Biochimie*, 93(4), pp.710-714.

List of figure legends

- Figure 1. Illustrating changes of p -values with selected features, (a) p -values vs. features and (b) $-\log_{10}(p\text{-value})$ vs. features.
- Figure 2. Stacked ensemble model of generalized linear model (GLM), C5.0, and SVM algorithms using a meta-classifier. The meta-classifier can be a GLM, linear discriminant analysis (LDA), or a random forest classifier (RF).
- Figure 3. Whisker plots illustrating the training accuracy of the three base classifiers.
- Figure 4. Correlations between predictions of base classifiers. Low correlations indicate that an ensemble classifier can be used.
- Figure 5. Identification of AR sequences in test data using BLASTp as a function of percent identity for AR samples from training data. In order to identify all AR sequences, an identity threshold of 33% must be used, and this leads to a high rate of false positives.

List of supplementary information

- Table S1: List of 44 Resistance Sequences
- Table S2: List of 36 Susceptibility Sequences
- Table S3: Confusion matrix for GLM meta-classifier
- Table S4: Confusion matrix for LDA meta-classifier
- Table S5: Confusion matrix for RF meta-classifier

Tables

Table 3. List of 621 Protein Features.

Feature	Feature dimension
Amino acid composition	20D
Hydropathy	21D
Normalized van der Waals volume	21D
Polarity	21D
Polarizability	21D
Charge	21D
Secondary structure	21D
Solvent accessibility	21D
Surface tension	21D
PSSM-based bigrams	400D
Structure-sequence	6D
Structure-probability	27D

Table 4. Performance Comparison of Stacked Ensemble Models.

Meta Classifier	Training Accuracy	Test Accuracy
GLM	81.38%	81.25%
LDA	80.30%	75.0%
RF	80.94%	81.25%

Figures

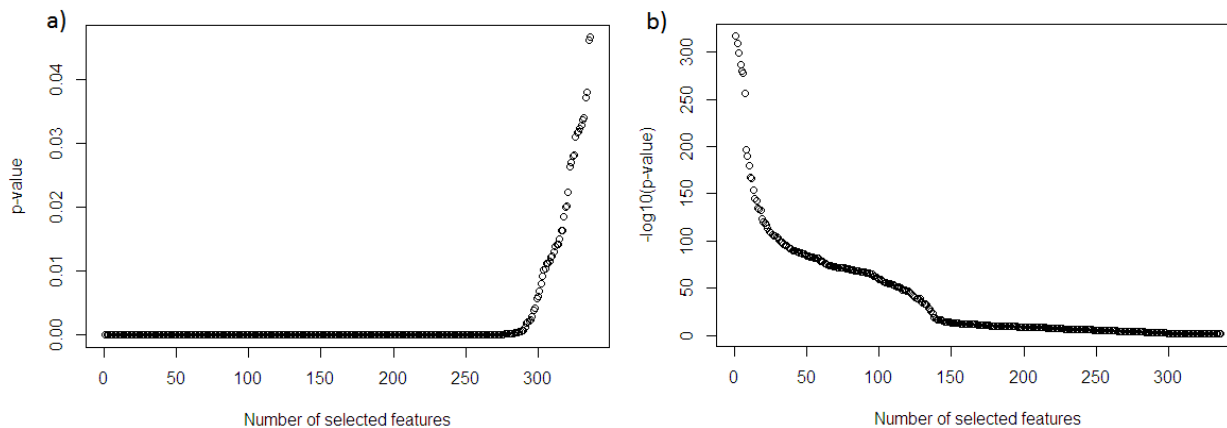


Figure 1. Illustrating changes of p-values with selected features: (a) p-values vs. features and (b) $-\log_{10}(\text{p-value})$ vs. features

AR Prediction Using Ensemble Learning

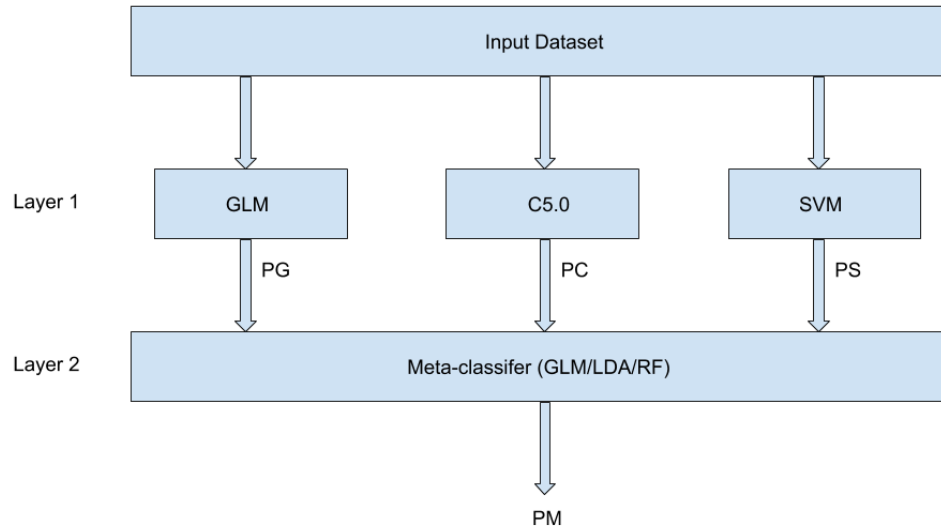


Figure 2. Stacked ensemble model of generalized linear model (GLM), C5.0, and SVM algorithms using a meta-classifier. The meta-classifier can be a GLM, linear discriminant analysis (LDA), or a random forest classifier (RF).

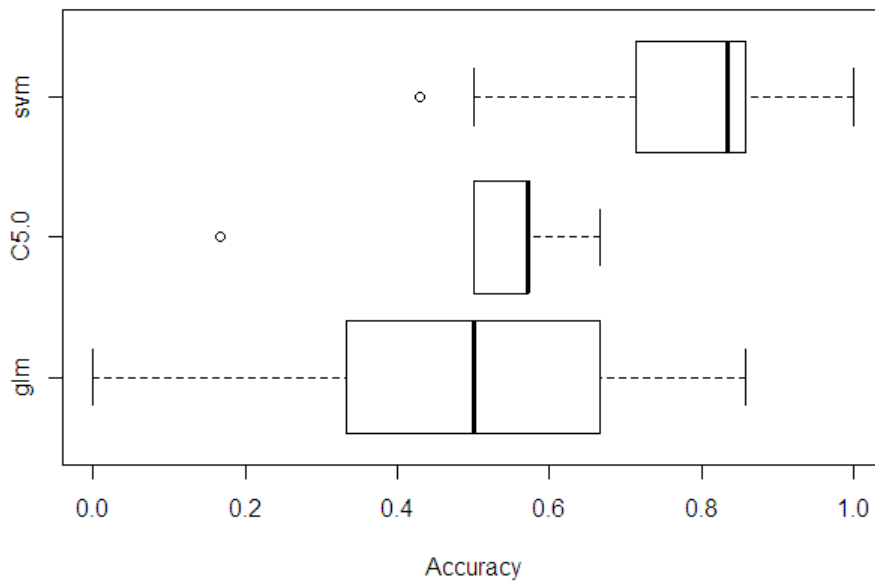


Figure 3. Whisker plots illustrating the training accuracy of the three base classifiers.

AR Prediction Using Ensemble Learning

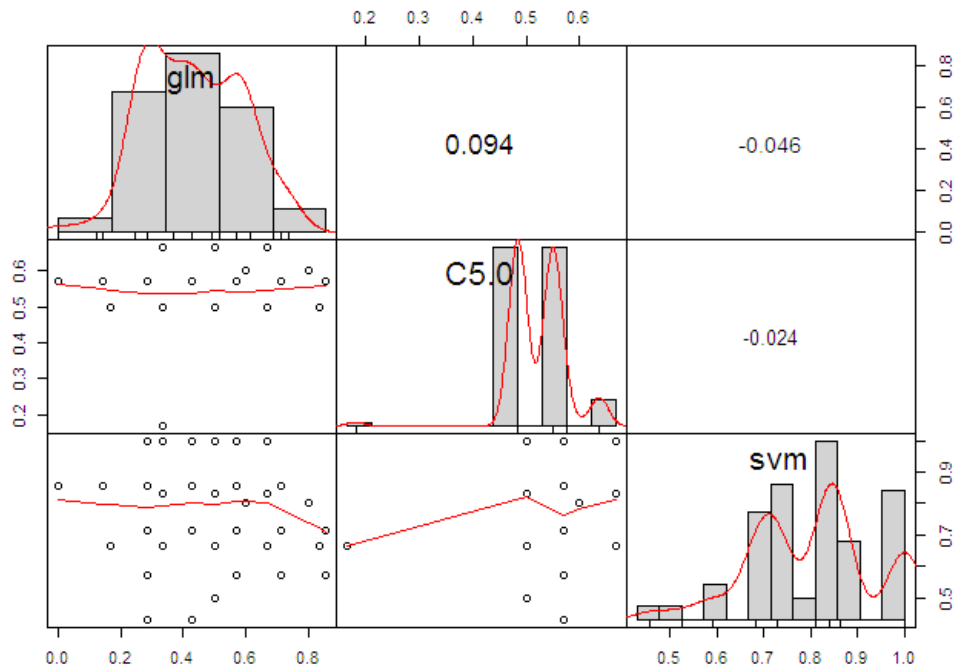


Figure 4. Correlations between predictions of base classifiers. Low correlations indicate that an ensemble classifier can be used.

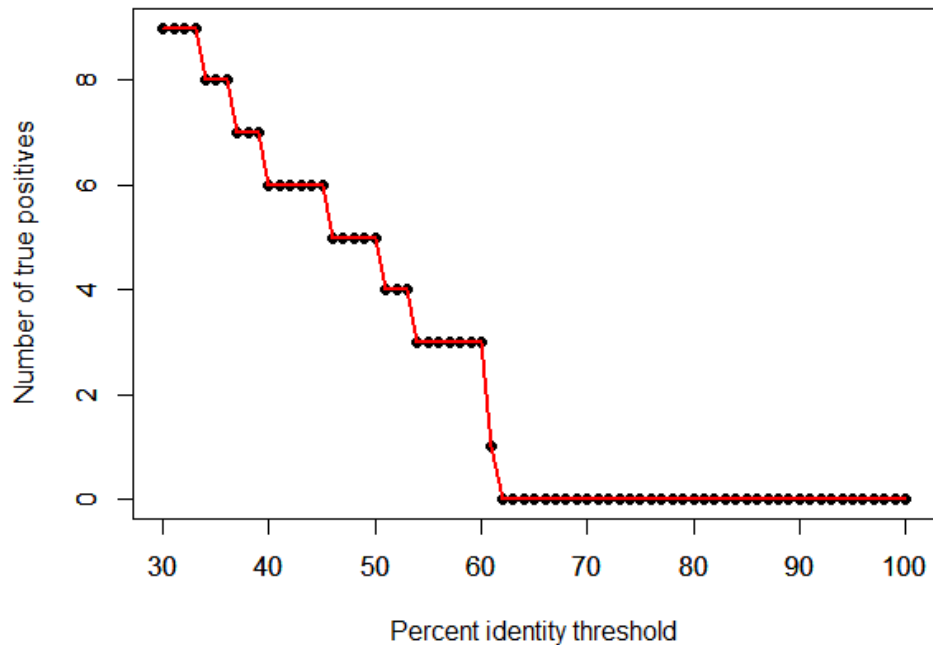


Figure 5. Identification of AR sequences in test data using BLASTp as a function of percent identity for AR samples from training data. In order to identify all AR sequences, an identity threshold of 33% must be used, and this leads to a high rate of false positives.