

Elsevier required licence: © <2022>. This manuscript version is made available under the CC-BY-NC-ND 4.0 license <http://creativecommons.org/licenses/by-nc-nd/4.0/>
The definitive publisher version is available online at
[\[https://www.sciencedirect.com/science/article/pii/S0147651322001117?via%3Dihub\]](https://www.sciencedirect.com/science/article/pii/S0147651322001117?via%3Dihub)

Predicting sustainable arsenic mitigation using machine learning

techniques Sushant K. Singh^{1,2*}, Taylor W. Robert¹, Biswajeet Pradhan^{3,4}, Ataollah Shirzadi⁵,
and Binh Thai Pham⁶

¹Department of Earth and Environmental Studies, Montclair State University, New Jersey, USA

*E-mail: sushantorama@gmail.com

²Artificial Intelligence & Analytics | Health Care and Life Sciences, Virtusa Corporation, New York City, NY, USA (Current affiliation)

³Centre for Advanced Modelling and Geospatial Information Systems (CAMGIS), School of Information, Systems & Modelling Faculty of Engineering and IT, University of Technology Sydney, Australia

⁴Department of Energy and Mineral Resources Engineering, Sejong University, Choongmu-gwan, 209 Neungdong-ro Gwangjin-gu, Seoul 05006, Republic of Korea

⁵College of Natural Resources, Department of Rangeland and Watershed Management Sciences, University of Kurdistan, Sanandaj, Iran

⁶Department of Geotechnical Engineering, University of Transport Technology, 54 Trieu Khuc, Thanh Xuan, Ha Noi, VietNam

ABSTRACT: Several artificial intelligence techniques have been applied to developing an array of environmental prediction models for various environmental challenges. However, there is no such prediction models exist for sustainable arsenic mitigation technologies. This study evaluates the state-of-the-art artificial intelligence models such as linear, nonlinear, ensemble, tree-based, Naïve Bayes, and neural network machine learning classifiers in predicting the preference of the most sustainable arsenic mitigation technology and provides following insights: (a) which machine learning algorithm has the highest prediction accuracy and robustness, and (b) which machine learning models are best to fit socioeconomic-environmental data for developing prediction models of sustainable arsenic mitigation technology? We evaluated 19 machine learning models for their predictive accuracy and the robustness by comparing their overall prediction accuracy, precision, recall, and the Area Under the Curve (AUC) of the Receiver Operating Characteristic (ROC) curve. A Gaussian distribution-based Naïve Bayes classifier outperformed the rest of the algorithms with the highest AUC of 0.825 on test data. The second two best models were Nu

30 Support Vector Classification (NuSVC) (AUC=0.800) (a radial basis function kernel-based
31 support vector machine algorithm) and K-Neighbors (AUC=0.790). All the ensemble classifiers
32 scored higher than 70% AUC, Random Forest being the top performer (AUC=0.769). We used
33 only one tree-based classifier Decision Tree, and it produced promising results (AUC=0.769) after
34 the three top classifiers. The neural network-based multilayer perceptron model, although ranked
35 9th position, also had a considerably good performance (AUC=0.748). Most linear classifiers did
36 not perform well with the Ridge classifier at the top (AUC=0.727) and perceptron at the bottom
37 (AUC=0.567). A Naïve Bayes-based classifier with Bernoulli distribution was the worst model
38 (AUC=0.500). Socioeconomic, demographic, and psychological data may not be linearly
39 associated with each other or with the outcomes. Therefore, nonlinear or ensemble classifiers could
40 better understand these complex relationships and help develop the most accurate and robust
41 prediction models. Gaussian NB is the best option for developing such prediction models on
42 socioeconomic and psychological data with small sample size. The proposed methodological
43 framework and the outcomes of the 19 machine learning models will help develop informed and
44 intelligent research methods as well as in targeting the population who are ready to adopt
45 sustainable arsenic mitigation technology.

46 Keywords: *Arsenic; Arsenic mitigation technologies; Machine learning; Linear classifier;*
47 *Nonlinear classifier; Ensemble*

48

49

50

51

52 1. INTRODUCTION

53 Socioeconomic, demographic, psychological, and cultural aspects of the communities exposed to
54 environmental contaminants, arsenic in this case, play a significant role in adopting and sustainably
55 practicing mitigation technologies [1-4]. Groundwater arsenic contamination is a global
56 environmental as well as a social threat to nearly 296 million individuals' lives in more than 100
57 countries, India and Bangladesh being the foremost victims [5, 6]. Arsenic is a human carcinogen,
58 can adversely impact dermal, cardiovascular, respiratory, neurological, and genetic systems, and
59 could lead to incurable varieties of cancers, if consumed for a prolonged period [5]. It is a well-
60 known fact that arsenic concentration in those contaminated countries has increased multiple times
61 of the standard norm of 10 µg/L stipulated by the World Health Organization (WHO) and the
62 United Nations Food and Agricultural Organization's (FAO) standard of 100 µg/L for irrigation
63 water [7]. Arsenic has also entered to the human, animal, and aquatic food chain evidenced from
64 high concentrations of arsenic in rice, beans, pulses, vegetable, cereals, fruits, poultry, egg, fresh
65 milk, milk powder, mother's milk, fish, shellfish, and algae [7-10]. There are studies that report a
66 significant concentration of arsenic in the human urine, blood, hair, and nail samples, evidence of
67 the arsenic exposure and accumulation in the human body [5, 11].

68 Low cost and simple arsenic mitigation techniques, such as arsenic treatment (filtration) units,
69 deep tube wells, piped water supply system, and rainwater harvesting system have been the
70 primary ways of providing arsenic-free water in the arsenic-contaminated areas globally [4, 12,
71 13]. However, because of the technical [4, 14-18], social, economic, and cultural challenges [1,
72 12-15, 19-22], these interventions could not be achieved sustainability. There are a fair number of
73 studies that highlight the technical challenges of arsenic mitigation technologies, however,
74 research on the socioeconomic, psychological, and cultural aspects of arsenic mitigation is still in
75 the rudimentary stages [3, 19, 23-25]. Based on these handfuls of studies, the authors found that
76 there is a lack of arsenic awareness and ownership of the implemented arsenic mitigation
77 technologies; low willingness to pay for arsenic mitigation technologies; complicated operation
78 and maintenance of arsenic mitigation technologies manuals; expensive technologies; long
79 distance between the households and the arsenic-free water sources; and social resistance by a
80 group of people to not let access the arsenic-free sources [1, 12-15, 19-22] have negatively
81 impacted the sustainable adoption of arsenic mitigation technologies.

82 In some recent studies [19, 22], it was discovered that communities' trust in the local agencies and
83 institutions as well as their social capital played a crucial role in their decision-making to adopt
84 arsenic mitigation technologies. In other studies [19, 26, 27], the authors highlighted that people's
85 perceived risk of health, income, and social discrimination due to arsenic contamination
86 significantly impact their decision-making process to adopt arsenic mitigation technologies. The
87 cost-effectiveness of a proposed arsenic mitigation technology also ensures their sustainable use
88 by the beneficiaries [21].

89 Accurately capturing the socioeconomic, demographic, psychological, and cultural information of
90 arsenic-affected communities is a challenging work for researchers, developing prediction models
91 on these data is even more daunting [1]. The reasons are lack of empirical data, a complex
92 relationship between the variables, socioeconomic, psychological, and cultural data are prone to
93 multicollinearity, and lack of successful case studies [1]. These all may affect the selection of the
94 most important predictors as in statistical analysis the model will only select the significant
95 variables unless we enforce expert opinion and include the variables known to be important but
96 not statistically significant [28]. In recent studies, the authors captured information on the
97 socioeconomic, demographic, social trust and capital aspects from an arsenic-exposed community
98 located in the middle-Ganga Plain of Bihar, India [19, 22]. While developing a logistic regression
99 prediction model of the adoption of arsenic treatment units, the authors started with 19 statistically
100 significant variables but end-up having eight variables in the final model with both significant and
101 nonsignificant variables. The model accomplished an overall prediction accuracy of 80.2%, which
102 looks promising [22]. However, since this model was not compared with other state-of-the art
103 modeling techniques such as machine learning models, we cannot say this is the best
104 socioeconomic model of predicting sustainable arsenic mitigation technologies. Also, considering
105 the lack of such data, how a robust machine learning model can be developed that could help
106 predict sustainable arsenic mitigation technologies in arsenic contaminated areas.

107 Several artificial intelligence techniques have been used in developing various prediction models
108 on environmental data including landslide susceptibility [29-33], groundwater potential [34, 35],
109 groundwater vulnerability [36], and groundwater contaminations [37, 38]. Pertaining to arsenic
110 research, various machine learning algorithms have been used in predicting arsenic contamination
111 in groundwater using the physical, chemical, hydrogeological, and topographical data [39-43].

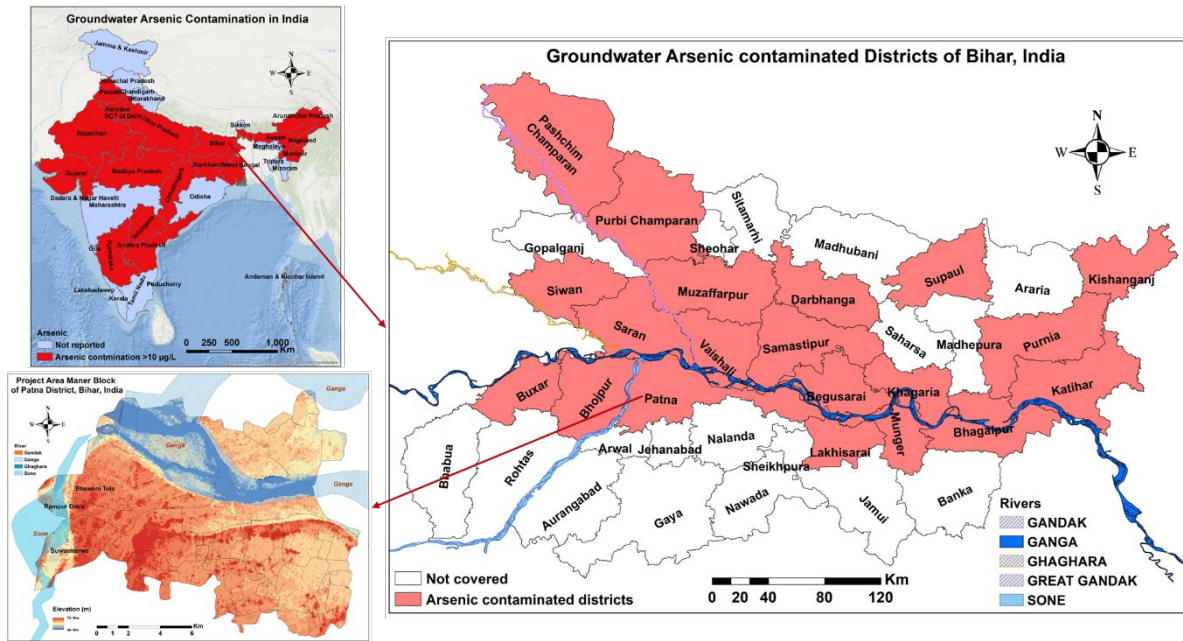
112 However, to the best of our knowledge, these state-of-the-art technologies have never been applied
113 to develop a socioeconomic model of arsenic mitigation. In a recent study by Singh et al. (2018),
114 the authors have used various machine learning models including logistic regression (LR), support
115 vector machine (SVM), decision trees (DT), k-nearest neighbor (k-NN), naïve Bayes (NB), and
116 random forests (RF) to predict arsenic awareness as a function of various socioeconomic,
117 sanitation, socio-behavioral, and social trust factors captured through an empirical study. In this
118 study, the authors discovered that arsenic awareness is a nonlinear classification problem and the
119 SVM and RF appeared to be the most appropriate machine learning algorithms in correctly
120 classifying arsenic awareness [1]. The authors further suggested that survey-based complex
121 environmental data may require advanced computational techniques opposed to traditional
122 statistical approach for developing accurate and robust prediction models.

123 Therefore, this study is a founding step in filling the above research gaps through answering
124 following questions: (a) which machine learning algorithm can achieve the highest prediction
125 accuracy and robustness in predicting the preference of sustainable arsenic mitigation technology
126 and (b) whether prediction of the preference of arsenic mitigation technology is a linear or a non-
127 linear classification challenge?

128 **2. METHODS**

129 **2.1. Study Area**

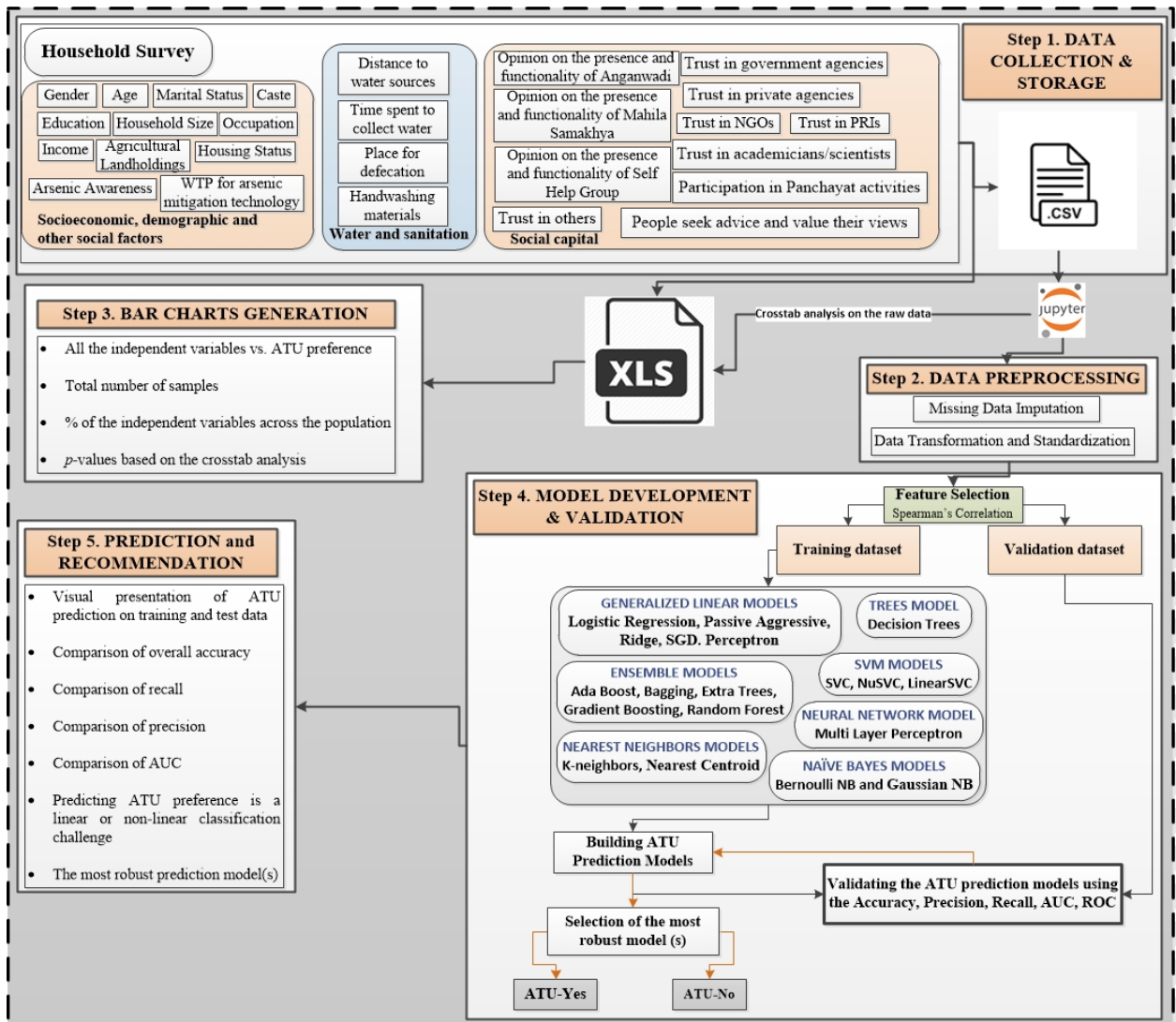
130 The state Bihar, the study area, is the second worst arsenic affected states of India after West
131 Bengal, which shares its geographical boundaries with other arsenic impacted regions including
132 Uttar Pradesh state of India, Bangladesh, Nepal, and Tibet [7]. The groundwater used for drinking
133 purposes is contaminated with elevated levels of arsenic in over 50% of the districts of Bihar.
134 Groundwater used for irrigation is also found to be contaminated with arsenic in some areas along
135 with a considerable amount of arsenic in agricultural soils and food materials [44-47]. Elevated
136 levels of arsenic in urine, blood, hair, and nail samples are also detected and several arsenicosis
137 victims are also diagnosed in the state [48-51]. Socioeconomic, health, and psychological aspects
138 of arsenic in the study area are also investigated, but still confined in a few geographical regions
139 of the state [19, 24, 25, 44, 48, 50, 51].



140

141 **Figure 1.** A map showing the arsenic affected districts of Bihar and the three villages selected
 142 for this study with their elevations.

143 In this study, we developed a five steps methodological framework to achieve our goals including:
 144 (1) data collection and storage, (2) data pre-processing, (3) data visualization, (4) model
 145 development and validation, and (5) prediction and recommendation.



146
147 **Figure 2.** Methodological framework adopted for this study.

148 **2.2. Data**

149 The data was captured by interviewing 340 households, randomly selected and stratified by their
 150 caste, through a structured questionnaire in three villages Suarmarwa, Rampur Diara, and Bhawani
 151 Tola: all located in the severely arsenic contaminated block Maner of Patna district of Bihar, India.
 152 Survey methodological details are explained in Singh (2015). The socioeconomic and
 153 demographic information were captured through asking questions on gender, age, marital status,
 154 caste, education, household size, occupation, income, agricultural landholdings, and housing
 155 status. Water and sanitation behaviors were captured through asking questions on a number of
 156 households involved in water collection, distance travelled and time spent to collect water, place

157 for defecation, and materials used for hand washing after defecation. Social capital and trust were
158 captured through asking questions on communities' opinion on the presence and functionality of
159 Anganwadi, Mahila Samakhya, Self-Help Group; trust in others, government agencies, NGOs,
160 Panchayat Raj Institutions, private agencies, academicians and scientists; participation in PRI's
161 activities, and whether people seek advice and value their views. Arsenic awareness was captured
162 through 10 questions converted to arsenic awareness index (low awareness vs. high awareness). A
163 detailed analysis on arsenic awareness is explained in [1]. Willingness to pay for arsenic mitigation
164 technology was also captured through a structured question. Communities' preference for
165 sustainable arsenic mitigation technologies was recorded through a structured questionnaire with
166 options of arsenic treatment unit (ATU), piped water supply systems, deep tube wells, dug
167 wells/open wells, and rainwater harvesting system. An in-depth analysis on these technological
168 preferences is available in [22].

169 This study provides a comprehensive analysis of the most preferred sustainable arsenic mitigation
170 technology (ATU) and investigates how the state-of-the-art machine learning technologies can
171 efficiently predict communities' preference of sustainable mitigation technology. The survey data
172 was transferred to a .excel file for frequency graphs generation and to a .csv file for data pre-
173 processing, statistical analysis, and machine learning model development using Jupyter Notebook
174 version 6.0.1 web application and Python 3 [52].

175 **2.3. Data pre-processing**

176 A majority of the variables were categorical and captured at the Likert-scale of five. Because of
177 their imbalanced frequency distribution across the responses (strongly disagree to strongly agree),
178 we reconstructed and recoded them for further analysis and model development. The
179 transformation of the original categories of the variables to new categories is explained in SI-1.
180 After screening the data, we found that one household did not answer the question on preference
181 of sustainable arsenic mitigation technology; therefore, we had a total 339 samples. Using Pandas
182 Python library, we imported the data to Jupyter Notebook and put it into a data frame for further
183 analysis [53]. Scikit-learn, an open access machine learning library in Python programming
184 language, was used for further analysis [54]. We also found 17 missing data that was imputed
185 using the mode of each feature. A contingency analysis was performed between all the independent

186 variables and ATU preference. Scipy library was used for all statistical analysis including
187 Spearman's correlation [55].

188 2.4. Data visualization

189 The results of contingency analysis were graphically presented in Figure 2-4 in the result section
190 using Excel Spreadsheet [56]. The graphs contain four important information including the number
191 of data points, categories of each feature wherever applicable, percentage of different categories
192 of features, and *p-value* to determine whether the responses across the categories were significantly
193 different from each other.

194 2.5. Machine learning algorithm selection

195 Applying machine learning to efficiently predict communities' preference of sustainable arsenic
196 mitigation technology is inspired by two recent researches where the first study [22] models the
197 preferences using a traditional statistical technique logistic regression, but did not provide a
198 comparison of how other statistical or machine learning techniques may fit the data to the various
199 models. The second study [1] provides a great deal of insights that developing prediction models
200 in the context of arsenic using complex socioeconomic, demographic, and other social factors may
201 need very specific type of algorithms or a hybrid model. Applying all the algorithms together and
202 comparing them in one study is not feasible therefore, we decided to select the state-of-the-art
203 linear, nonlinear, ensemble, Naïve Bayes, and tree-based classifiers to develop the models and
204 compare them for their prediction accuracy and robustness. A brief description of each algorithm
205 applied in this study is described below.

206 2.5.1. Generalized Linear Models

207 *Logistic Regression (LR)*

208 The LR is a multivariate regression that provides the probability of the presence of an event at
209 each response according to the predictors [57, 58]. It has some advantages that environmental
210 researchers have encouraged to apply it, including; (1) the LR does not need to set normality for
211 independent feature, (2) predictors can either be continuous or discrete or any combination of these
212 types of data, (3) it is easy to implement in most statistical packages such as SPSS, SAS, STATA,
213 R and so on [59-61]. The dependent variable in the LR should be binary (present/occurrence and
214 absent/non-occurrence of an event) to achieve the probability values. In this study, we aim to

215 predict the probability of adoption of arsenic treatment units (ATU) using LR model and some
216 predictors. The LR can be formulated in its simplest form as follows:

$$217 \quad P_{LR} = \frac{e^z}{1+e^z} \quad (1)$$

218 where, P_{LR} is the probability of present/occurrence of an event that varies from 0 to 1 as s-shaped
219 curve, Z is a linear combination that varies from $-\infty$ to $+\infty$ and can be computed as bellow:

$$220 \quad Z = c_0 + a_1x_1 + a_2x_2 + \dots + a_nx_n \quad (2)$$

221 where, c_0 is the constant coefficient/intercept of the LR model, n is the number of predictors,
222 a_i ($i = 1, 2, 3, \dots, n$) is the independent variables as input to the model, x_i ($i = 1, 2, 3, \dots, n$) is
223 the coefficients of each predictors as input to the model [62].

224 *Passive Aggressive (PA)*

225 The passive-aggressive algorithm is similar to Perceptron as there is no learning rate required.
226 However, it contain a regularization parameter “C” [54, 63]. Scikit learn machine learning library
227 in python offers a passive aggressive classifier for binary classification such as in this case for
228 classifying preference of ATU as a most sustainable arsenic mitigation technology. This technique
229 is less explored though.

230 *Ridge*

231 The Ridge classifier is another linear classifier and is also known as least squares support vector
232 machines where in Scikit learn this classifier first converts binary targets ATU-No and ATU-Yes
233 to, respectively -1 and 1. The algorithm then regress the dependent variable against the independent
234 variables, and the predicted class resembles to the sign of the regressor’s prediction [54].

235 *Stochastic Gradient Descent (SGD)*

236 The SGD is another linear classifier and commonly used with large sample size. The SGD can fit
237 both logistic regression model and support vector machine by selecting appropriate loss functions,
238 respectively “log” and “hinge” [54]. It requires a learning rate, and the loss is estimated for each
239 data point at a time. For a normally distributed data, the SGD provides a better result. Although
240 SGD classifier is known for its efficiency on large datasets, it is known to be highly sensitive to
241 feature scaling [54].

242 *Perceptron*

243 The Perceptron is another linear classifier that doesn't require learning rate, that's why it is a
244 favorable classifier for large scale learning. Additionally, it doesn't penalize the learning and
245 updates the model only on mistakes [\[54\]](#).

246 2.5.2. Trees Models

247 *Decision Trees (DT)*

248 The DT is a tree-base non-parametric classifier that learns decision rules from the data features.
249 This is a very popular classifier for developing binary classification models as it is simple to
250 understand and to interpret. Although it doesn't work well with missing data, it is a good classifier
251 that efficiently handles both numerical and categorical data, and requires less assumptions [\[54\]](#).

252 2.5.3. Ensemble Models

253 *Ada Boost*

254 The AdaBoost is an ensemble classifier, first introduced in 1995 by Freund and Schapire. It fits a
255 sequence of weak learners on modified versions of the data and produces a combined predicted
256 class. In this entire process, the classifier makes sure that none of the data point is left in the training
257 phase [\[54, 64\]](#).

258 *Bagging*

259 Bagging is another ensemble classifier that consolidates a final prediction based on the previous
260 predictions on randomly selected subsets of the original training dataset. It works well with strong
261 and complex models [\[54, 65\]](#).

262 *Extra Trees (ET)*

263 The ET is another ensemble model, well known to control over-fitting but less explored. In scikit-
264 learn, this classifier fits several randomized decision trees i.e. "extra-trees" on several sub-samples
265 of the dataset [\[54\]](#).

266 *Gradient Boosting (GB)*

267 The GB is another ensemble classifier offered by Scikit-learn library and is very popular among
268 scientific computation community [54]. It builds an additive model in a forward stage-wise fashion
269 and allows for the optimization of arbitrary differentiable loss functions.

270 *Random Forest (RF)*

271 The RF is another popular ensemble classifier available through Scikit-learn that fits several
272 decision tree classifiers on many sub-samples of the dataset. The RF averages the probabilistic
273 prediction value of each decision-tree and uses to improve the classification accuracy. It is also
274 known to be prone to over-fitting [54, 66].

275 2.5.4. Support Vector Machines Models

276 Support vector machines (SVM) is known for efficiently classifying linearly separable data as well
277 as non-linearly separable data by using a kernel function, such as sigmoid, radial, or polynomial.
278 It is advantageous if the data is clearly separable and the ratio between the number of dimensions
279 and the number of samples is greater. It is also memory efficient. However, the SVM takes more
280 time in training the model therefore; it is not feasible for large data set. Likewise, with noisy data,
281 the performance is poor.

282 *SVC*

283 The SVC is known as C-Support Vector Classification, a non-linear SVM classifier [54].

284 *Nu-Support Vector Classification (Nu-SVC)*

285 The Nu-SVC is another non-linear SVM classifier that uses a parameter to control the number of
286 support vectors [54].

287 *Linear Support Vector Classification (LinearSVC)*

288 The LinearSVC uses 'linear' kernel and Scikit-learn library offers suppleness in choosing loss
289 functions and penalties [54].

290 2.5.5. Nearest Neighbors Models

291 *K-nearest neighbors*

292 *k*-Nearest Neighbors (k-NN) is a nonparametric non-linear classifier that categorizes the intended
293 event by a majority vote of its *k* nearest neighbors, which is a positive integer and can be derived
294 through elbow-test [54, 67]. k-NN is advantageous in several ways, as it does not require
295 assumptions, easily interpretable, good for nonlinear data, and can produce comparatively better
296 accuracy [54, 67]. However, k-NN is very sensitive to irrelevant features, the scale of the variables,
297 the dimensions of the dataset, and class imbalance. In addition, it can be highly computation
298 intensive as it stores all the training data.

299 *Nearest Centroid (NC)*

300 In Scikit-learn, the NC classifier belongs to the nearest neighbor algorithms where each class is
301 characterized by the centroid of its members [54].

302 2.5.6. Neural Network Model

303 *Multi-Layer Perceptron*

304 The ANN is organized and structured based on the skills of human brain cells to extract knowledge
305 from the input dataset [32, 68]. It has some advantages that it has been a strong and promising
306 technique to prediction environmental problems including, (i) it can efficiently detect a different
307 subset of data within a whole dataset, (ii) it do not need to any experience and pre-knowledge
308 process, and (iii) It do not need to a given statistical model in the training dataset [69]. The
309 multilayer perceptron (MLP) and radial base function (RBF) are two popular and well-known as
310 functions of ANN. Although the capability of these two function are different from one case study
311 to another, in general the MLP is more popular and general than the RBF kernel function [70]. The
312 MLP is more successfully and flexibility in modeling, especially on non-linear, imprecise and
313 imperfect data so that it can extract the reliable results [71]. Therefore, in this study, we used of
314 MLP function to construct a network for determining the relationship between the ATU and
315 predictors. The MLP has a structure with three layers including a n input, an output and one or
316 more hidden layers between them [72, 73]. In this study, ATU predictors are taking into
317 consideration as inputs (neurons) and the weights for each predictor is output. In a simplest form,
318 let x_i and w_i are input predictors and they obtained weights during the modeling process. In hidden

319 layer, they are multiplied and then summed up to extract the final output or the final weights (y_i)
320 by a non-linear activation function as follows:

$$321 \quad \text{net} = \sum_{i=0}^n w_i x_i \quad (3)$$

$$322 \quad y_i = f(\text{net}) \quad (4)$$

323 2.5.7. Naïve Bayes Models

324 Naïve Bayes (NB) is another state-of-the-art nonlinear classifier that works on the Bayes theorem,
325 has a strong assumption that all the predictors are independent and not correlated to each other,
326 and can be mathematically presented as below:

$$327 \quad \text{Posterior Probability}[Y = P(c|x)] =$$
$$328 \quad \frac{[\text{Likelihood: } P(x_1|C) \times P(x_2|C) \times \dots \times P(x_n|C)] \times [\text{Class Prior Probability: } P(c)]}{\text{Predictor Prior Probability: } P(x)} \quad (5)$$

329 NB is known to be outperforming other state-of-the-art classifiers, such as logistic regression [54].
330 It requires less training dataset and can quickly predict on test dataset. It is also known to be a good
331 classifier for categorical variables [54]. However, the strong assumption of independence could be
332 a challenge while applying NB on the dataset with multicollinearity.

333 *Bernoulli Naïve Bayes*

334 This is one of the NB classifiers that assumes the data has multivariate Bernoulli distributions.
335 Therefore, this class requires samples to be represented as binary-valued feature vectors; if handed
336 any other kind of data, a BernoulliNB instance may binarize its input (depending on
337 the binarize parameter) [54].

338 The decision rule for Bernoulli naive Bayes is based on

$$339 \quad P(x_i | y) = P(i | y)x_i + (1 - P(i | y))(1 - x_i) \quad (6)$$

340 which differs from multinomial NB's rule in that it explicitly penalizes the non-occurrence of a
341 feature i that is an indicator for class y , where the multinomial variant would simply ignore a non-
342 occurring feature [54].

343 *Gaussian Naïve Bayes*

344 This is another NB classifier that assumes that the data has a Gaussian distribution, and the
345 Gaussian distribution can be presented in Eq. 7.

346
$$P(x_i | y) = \frac{1}{\sqrt{2\pi\sigma_y^2}} \exp\left(-\frac{(x_i - \mu_y)^2}{2\sigma_y^2}\right) \quad (7)$$

347 The parameters σ_y and μ_y are estimated using maximum likelihood [54].

348 2.5.8. Nearest Shrunken Centroids (NSC)

349 The NSC [74, 75] in classification issues calculates a standardized centroid for each class using
350 the average value of each feature of each class divided by its class standard deviation of that
351 feature. In the next step, the feature vector of a new conditioning factor as input is compared to the
352 centroids of each of these classes. Consequently, the class with the closest centroid (in squared
353 distance) is the predicted class for that new conditioning factor as input data [76]. In this model,
354 using a threshold, each of the class centroids of the features shrinks toward the overall centroid.
355 Mathematically, first a threshold value is assigned to the class centroids of the features, and if it is
356 small for all classes, it is set to zero. Consequently, when shrinking the centroids for all classes is
357 completed the new sample of the feature is classified by the usual nearest centroid rule [76].

358 2.6. Model Development and Validation

359 After imputing the missing values, except the response variable ATU, the data was scaled and
360 centered, and 75% of the data was used for model development and 25% for testing. The data was
361 split using model selection function of Scikit-learn [54].

362 2.6.1. Feature selection

363 Feature selection is an important step in developing any machine-learning model, and there are
364 various ways of selecting the most appropriate predictors. Some machine learning models can
365 handle feature selection, but not most of them. Therefore, we decided to apply Spearman's
366 correlation to select all the predictors with a significant correlation with ATU [54].

367 2.6.2. Training machine learning models

368 Using `train_test_split` function of Scikit-learn the data was split into 75% for training the model
369 and 25% of the data for model validation [54]. To train all 19 models together, we created a
370 function by listing ensemble methods including `ensemble.AdaBoostClassifier`,
371 `ensemble.BaggingClassifier`, `ensemble.ExtraTreesClassifier`,
372 `ensemble.GradientBoostingClassifier`, `ensemble.RandomForestClassifier`; generalized linear
373 models including `linear_model.LogisticRegressionCV`,
374 `linear_model.PassiveAggressiveClassifier`, `linear_model.RidgeClassifierCV`,
375 `linear_model.SGDClassifier`, `linear_model.Perceptron`; Navies Bayes methods including
376 `naive_bayes.BernoulliNB`, `naive_bayes.GaussianNB`; Nearest Neighbor methods including
377 `neighbors.KNeighborsClassifier` and `neighbors.NearestCentroid`; SVM techniques including
378 `svm.SVC`, `svm.NuSVC`, `svm.LinearSVC`; Trees-based methods including
379 `tree.DecisionTreeClassifier`, `tree.ExtraTreeClassifier()`, and Neural Network methods including
380 `neural_network.MLPClassifier` [54].

381 2.6.3. Validation of machine learning models

382 All models were validated on 25% of the data using accuracy, precision, recall, and AUC score.

383 Predicted number of people preferred ATU (TP); Predicted number of people not preferred ATU
384 (TN); incorrectly predicted number of people preferred ATU (FP); incorrectly predicted number
385 of people not preferred ATU (FN)

$$386 \text{ Accuracy} = \frac{TP+TN}{TP+TN+FP+FN} \quad (8)$$

$$387 \text{ Sensitivity} = \frac{TP}{TP+FN} \quad (9)$$

$$388 \text{ Sensitivity} = \frac{\text{predicted number of people preferred ATU}}{\text{total number of people preferred ATU in the population}} \quad (10)$$

389 Sensitivity is also known as recall, hit rate, and true positive rate.

$$390 \text{ Specificity} = \frac{TN}{TN+FP} \quad (11)$$

$$391 \text{ Specificity} = \frac{\text{predicted number of people not preferred ATU}}{\text{total number of people not preferred ATU in the population}} \quad (12)$$

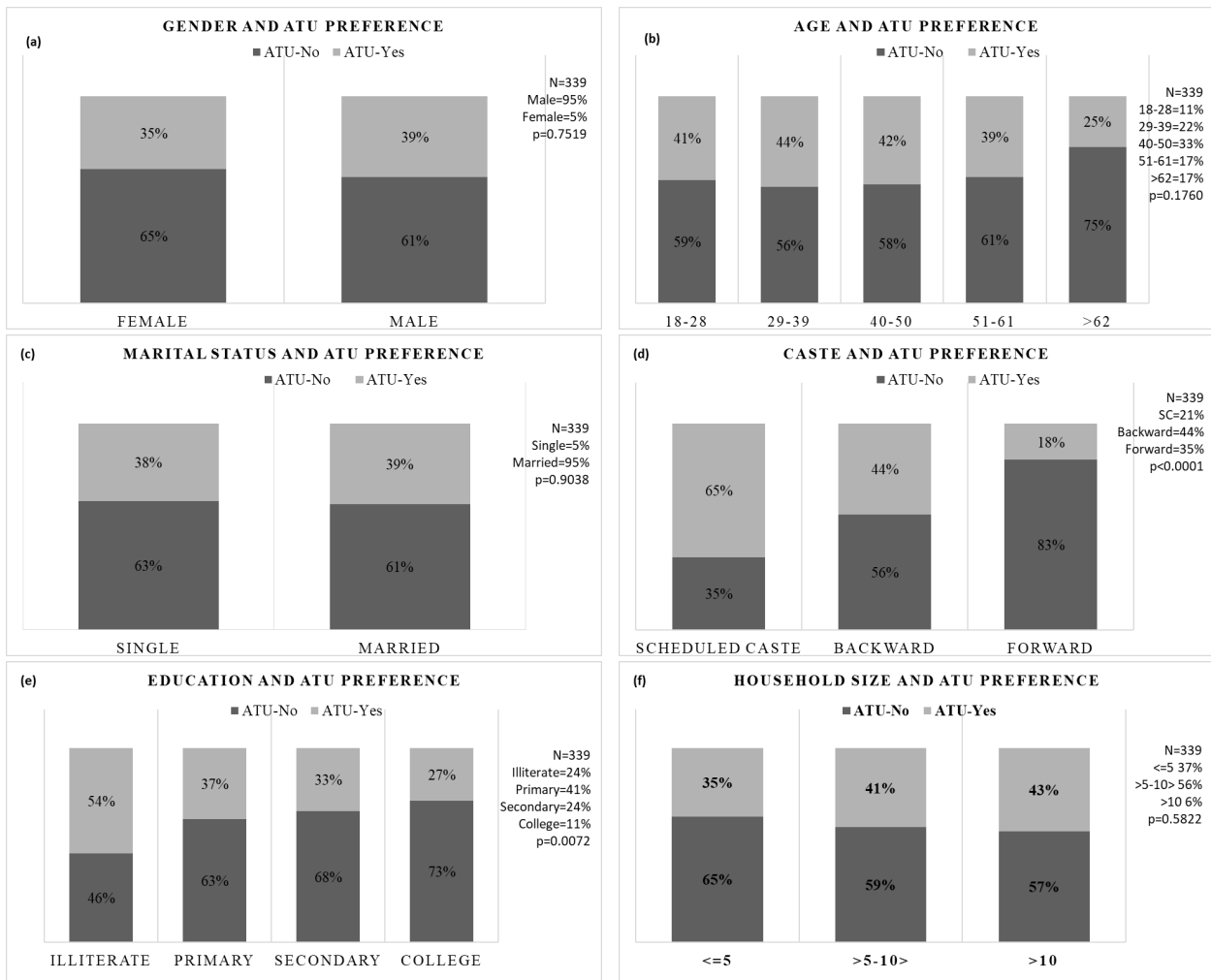
392 Specificity is also known as true negative rate and selectivity.

393 **3. RESULTS**

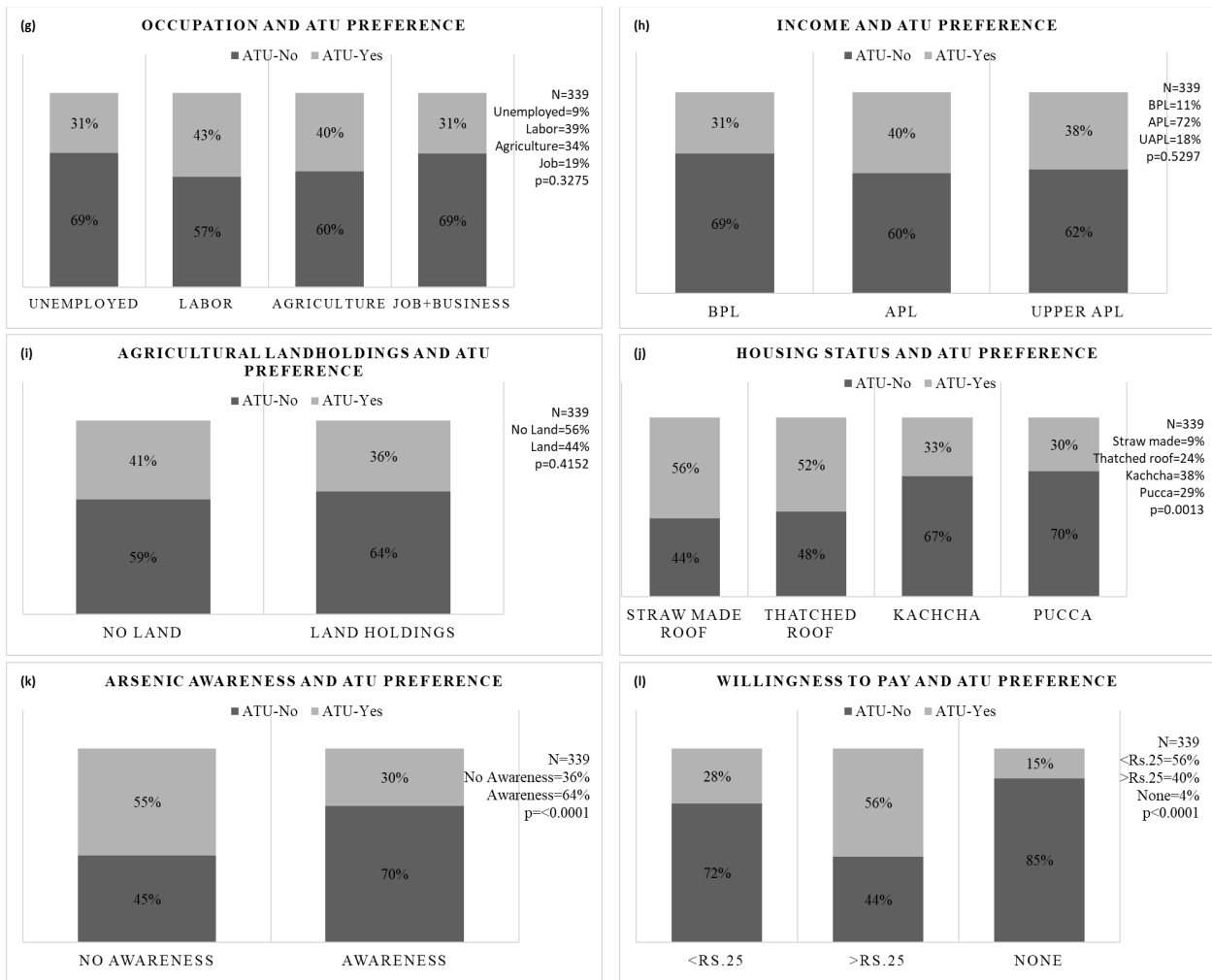
394 **3.1. Socioeconomic, demographic and other social factors**

395 With 39% of the population, ATU was the most preferred sustainable arsenic mitigation
 396 technology and significantly different from the other options including piped water supply system,
 397 deep tube wells, dug wells/open wells, and rainwater harvesting system [22]. Although there was
 398 a less participation of females in the survey than males, we did not find any significant difference
 399 ($p=0.7519$) among them preferring ATU as the most sustainable arsenic mitigation technology
 400 (Figure 3-a). People of age group of 40-50 (Figure 3-b) were more interested in adopting ATU as
 401 the sustainable arsenic mitigation technology than the other age groups, but not significantly
 402 different from each other ($p=0.1760$).

403



405



406

407

408

409

410 **Figure 3.** (a) gender and ATU preference; (b) age and ATU preference; (c) marital status and ATU
 411 preference; (d) caste and ATU preference; (e) education and ATU preference; (f) household size and ATU
 412 preference; (g) occupation and ATU preference; (h) income and ATU preference; (i) agricultural
 413 landholdings and ATU preference; (j) housing status and ATU preference; (k) arsenic awareness and ATU
 414 preference; (l) willingness to pay for arsenic mitigation and ATU preference.

415

416 A similar trend was observed with marital status (Figure 3-c) of the respondents, where there was
 417 no difference in the response of preferring ATU ($p=0.9038$) among single and married people.

418 This could be also because of the low participation of unmarried people (5%) in the survey. Caste
 419 was appeared to be one of the most important features of the respondents that distinguish the

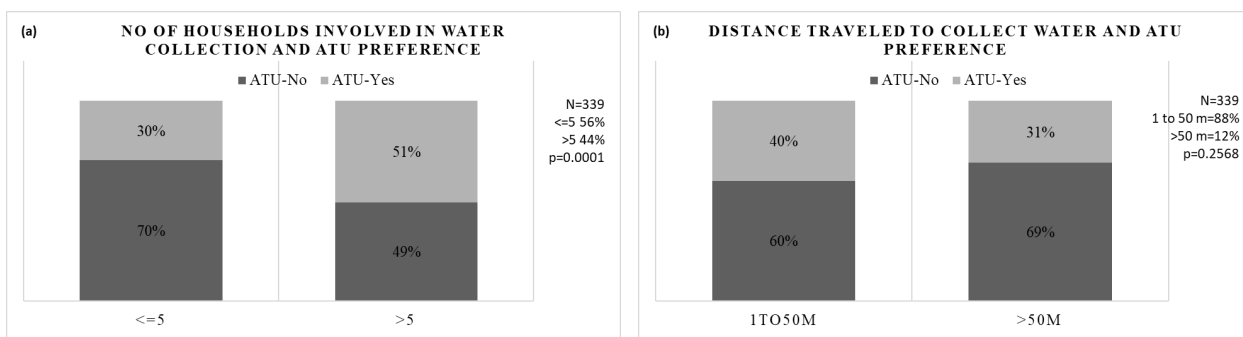
420 preference of ATU among various castes. A majority of SC (65%) preferred ATU (Figure 3-d)
 421 followed by BC (44%) and FC (18%) and their preferences were significantly different ($p<0.0001$)
 422 from each other. A similar trend was observed across various education levels of the respondents

423 (Figure 3-e) where people with no education found to be more likely (54%) to adopt ATU than

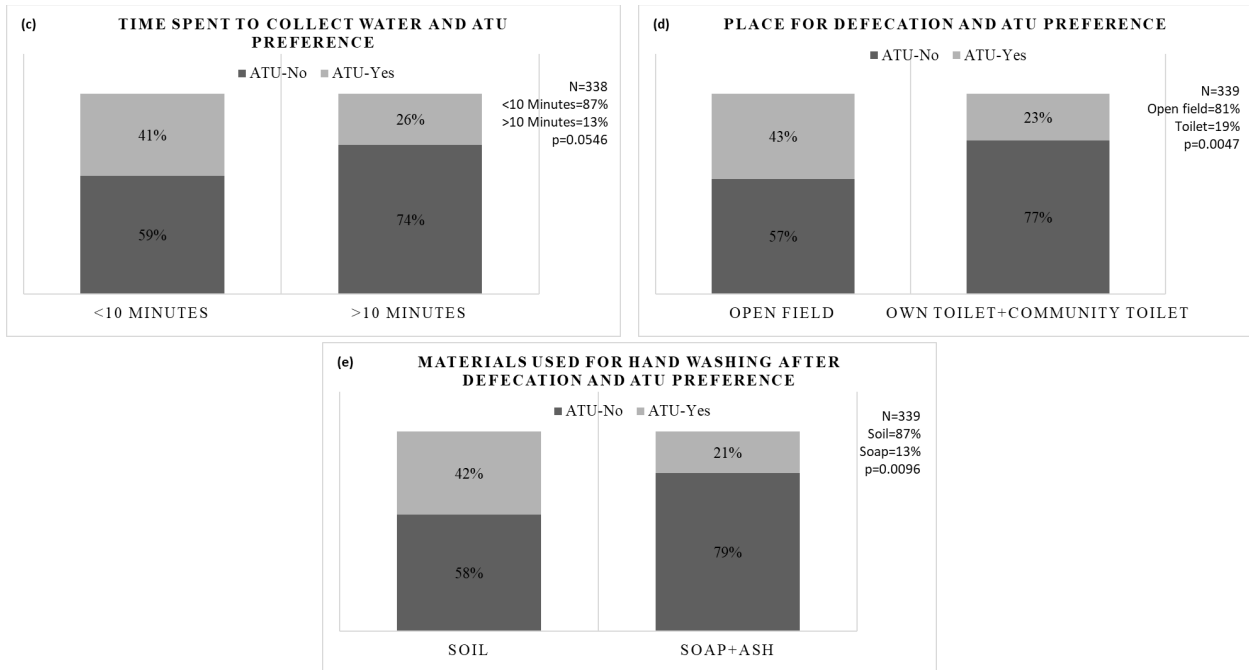
424 people with various levels of education and was significantly different from each other (0.0072).
 425 There was no association between household size and preference of ATU ($p=0.5822$) (Figure 3-f).
 426 A similar trend was observed with the occupation ($p=0.3275$) (Figure 2-g), income ($p=0.5297$)
 427 (Figure 3-h), and agricultural landholdings ($p=0.4152$) (Figure 3-i) of the respondents where their
 428 preference of ATU as the sustainable arsenic mitigation technology was not different. People live
 429 in straw-made roofed houses were more likely (56%) (Figure 3-j) to prefer ATU than the people
 430 live in better housing structures and the responses were significantly different ($p=0.0013$) from
 431 each other. The respondents less aware of arsenic (55%) were more likely to prefer ATU than the
 432 respondents with arsenic awareness (Figure 3-k) and their responses were significantly different
 433 ($p<0.001$). This further interprets that the people perceive technology/filters as a better solution to
 434 purify any water contaminants. The respondents with a WTP >Rs. 25 were more likely (56%) to
 435 prefer ATU and their responses were different ($p<0.001$) across various WTP levels (Figure 3-l).

436 3.2. Water and sanitation factors

437 When it comes to respondents' water and sanitation behaviors, the number of households involved
 438 in water collection was significantly ($p<0.0001$) associated with ATU preference (Figure 4-a). The
 439 households with more than five members involved in water collection were more likely to prefer
 440 ATU, which may indicate their concern about dependency on many people for water collection
 441 and collective time loss. The distance travelled (Figure 4-b) and the time spent to collect water
 442 (Figure 4-c) were not significantly associated with ATU preference. However, people's sanitation
 443 habits were significantly associated with ATU preference. The respondents who defecate in the
 444 open field were more likely (43%) to prefer than who uses their own toilets (23%). The similar
 445 trend was observed for the materials used for hand washing after the defecation (Figure 4-c).



446



447

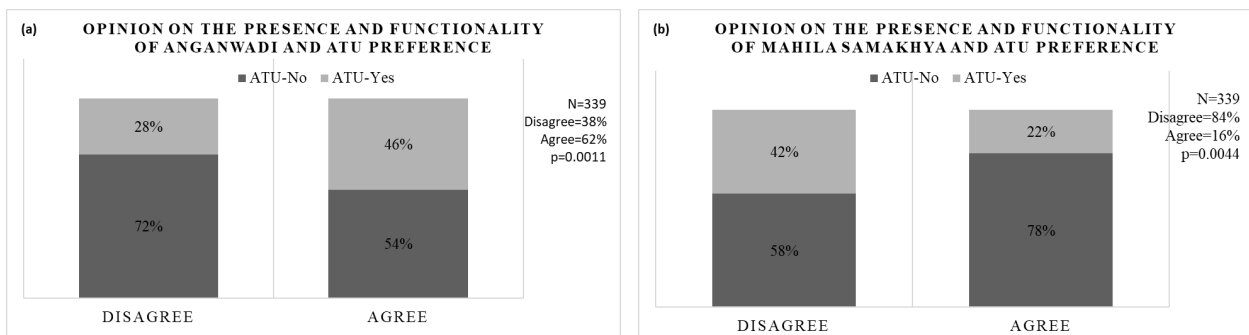
448

449

450 **Figure 4.** (a) number of households involved in water collection and ATU preference; (b) distance travelled
 451 to collect water and ATU preference; (c) time spent to collect water and ATU preference; (d) materials used
 452 for hand washing after defecation and ATU preference.

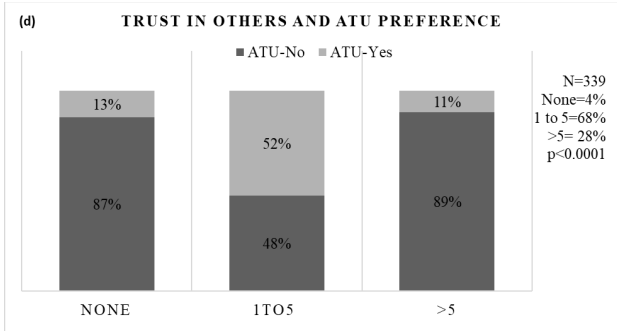
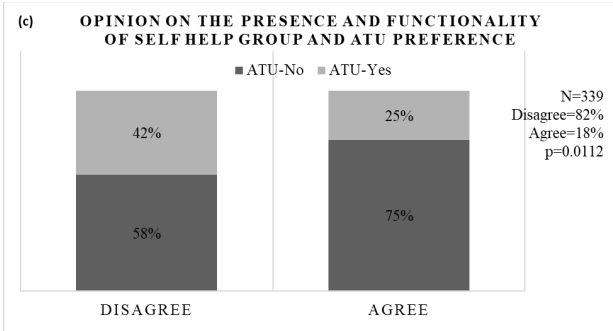
453 3.3. Social capital

454 Social capital found to be a major player in guiding the decision making to adopt arsenic mitigation
 455 technology. The respondents, who agreed on the presence and functionality of Anganwadi, were
 456 more likely to prefer ATU with a significant difference ($p=0.0011$) from who disagreed (Figure 5-
 457 a). A contrasting pattern was found with the respondents' response on the presence and
 458 functionality of Mahila Samakhya (Figure 5-b) and Self-Help Groups (Figure 5-c) with the
 459 respondents who disagree were more likely to prefer ATU than who agreed on this question.

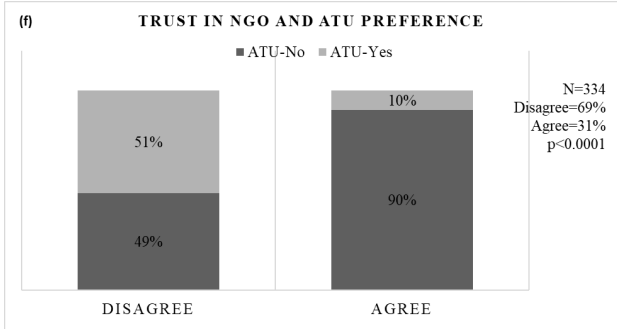
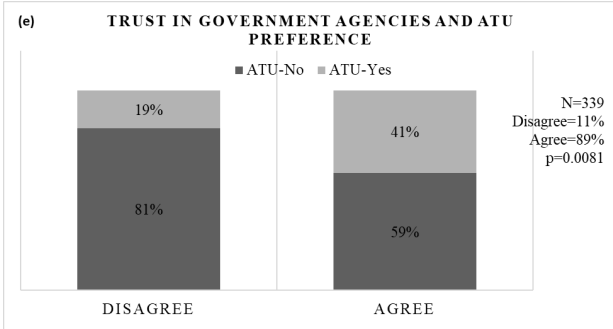


460

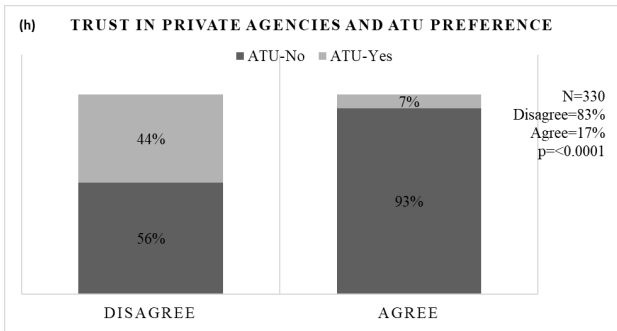
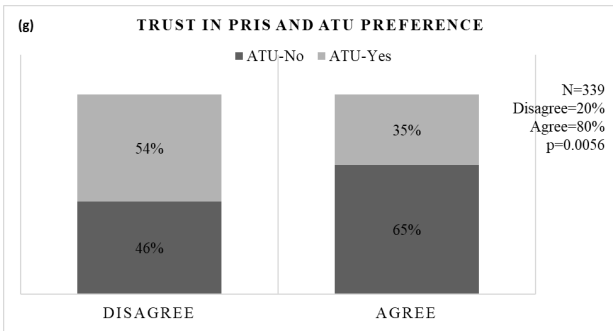
461



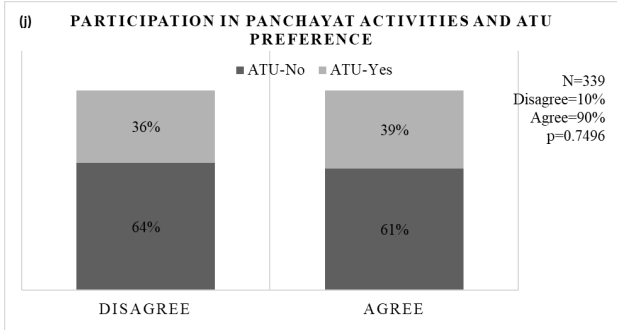
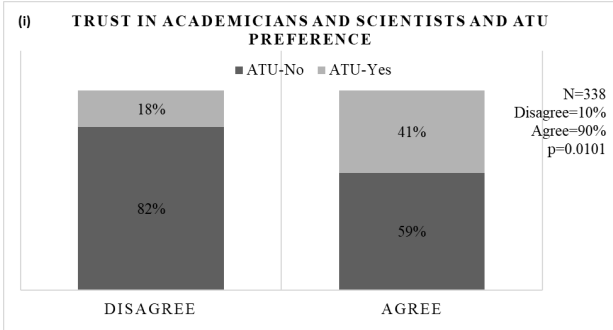
462



463

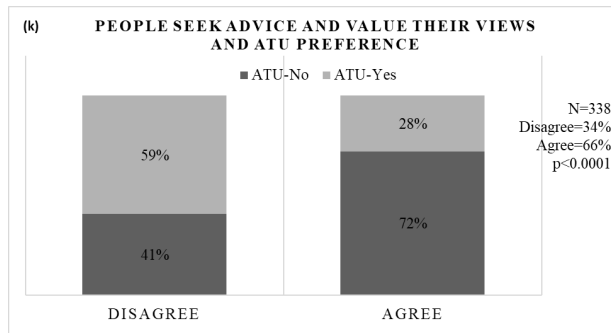


464



465

466

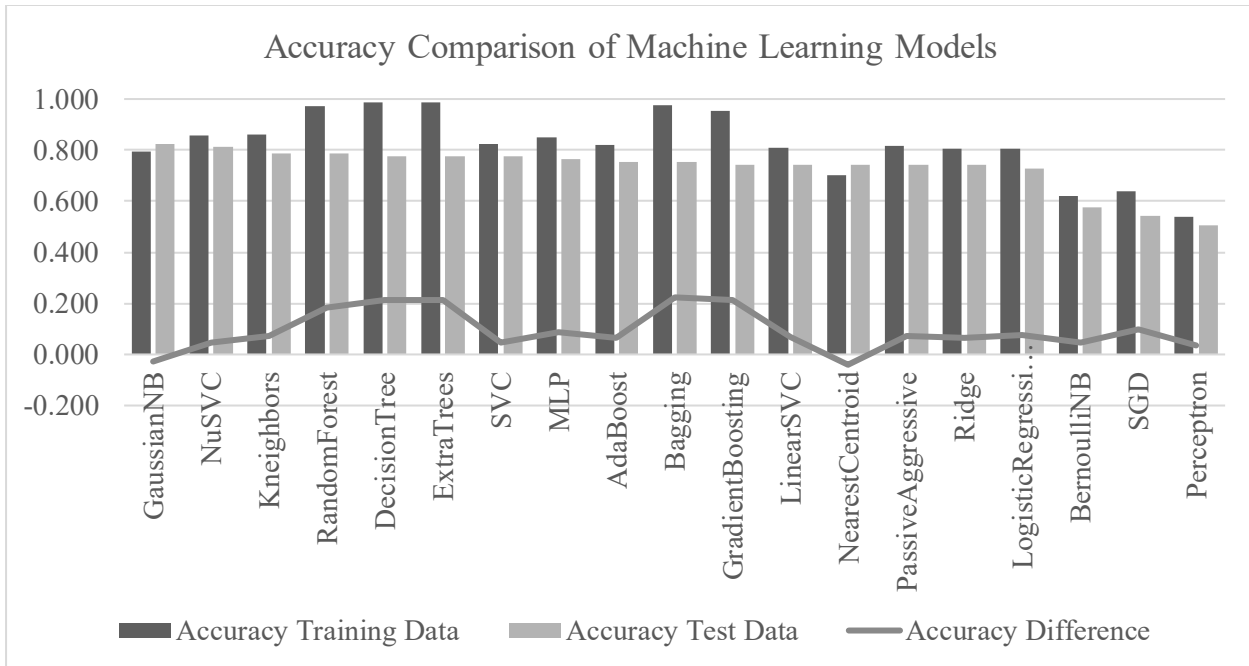


467 **Figure 5.** (a) opinion of the presence and functionality of Anganwadi and ATU preference; (b) opinion of
468 the presence and functionality of Mahila Samakhya and ATU preference; (c) opinion of the presence and
469 functionality of Self Help Group and ATU preference; (d) trust in others and ATU preference; (e) trust in
470 government agencies and ATU; (f) trust in NGOs and ATU; (g) trust in PRIs and ATU preference; (h) trust
471 in private agencies and ATU preference; (i) trust in academicians and scientists and ATU preference; (j)
472 participation in panchayat activities and ATU preference; (k) people seek advice and value their views and
473 ATU preference.

474
475 The respondents who trust between 1 and 5 people were more likely (Figure 5-d) to prefer ATU
476 and significantly different ($p < 0.001$) than other trust categories. The people who trust in
477 government agencies (Figure 5-e) and academicians (Figure 4-i) were more likely to prefer ATU
478 and their responses were significantly different ($p = 0.0081$ and $p = 0.0101$ respectively). In contrast,
479 the respondents who do not trust NGO (Figure 5-f, $p < 0.0001$), PRIs (Figure 5-g, $p = 0.0056$), and
480 private agencies (Figure 5-h, $p < 0.0001$) were more likely to prefer ATU. People's participation in
481 panchayat activities had no significant association ($p = 0.7496$) with ATU preference (Figure 5-j).
482 On the other hand, people who reported that other people in the society do not seek advice and
483 value their views were more likely (59%, $p < 0.0001$) to prefer ATU than who agreed on this
484 question.

485 3.4. Comparing machine learning models

486 We developed 19 different models and all of them were trained on the 75% of the data set (Figure
487 6). Among these 19 models, Decision Tree (accuracy=0.988), Extra Trees (accuracy=0.988),
488 Bagging (accuracy=0.967), Random Forests (accuracy=0.972), and Gradient Boosting
489 (accuracy=0.953) achieved the highest accuracy, all above 95%, Decision Tree being the top
490 performer (Figure 6).

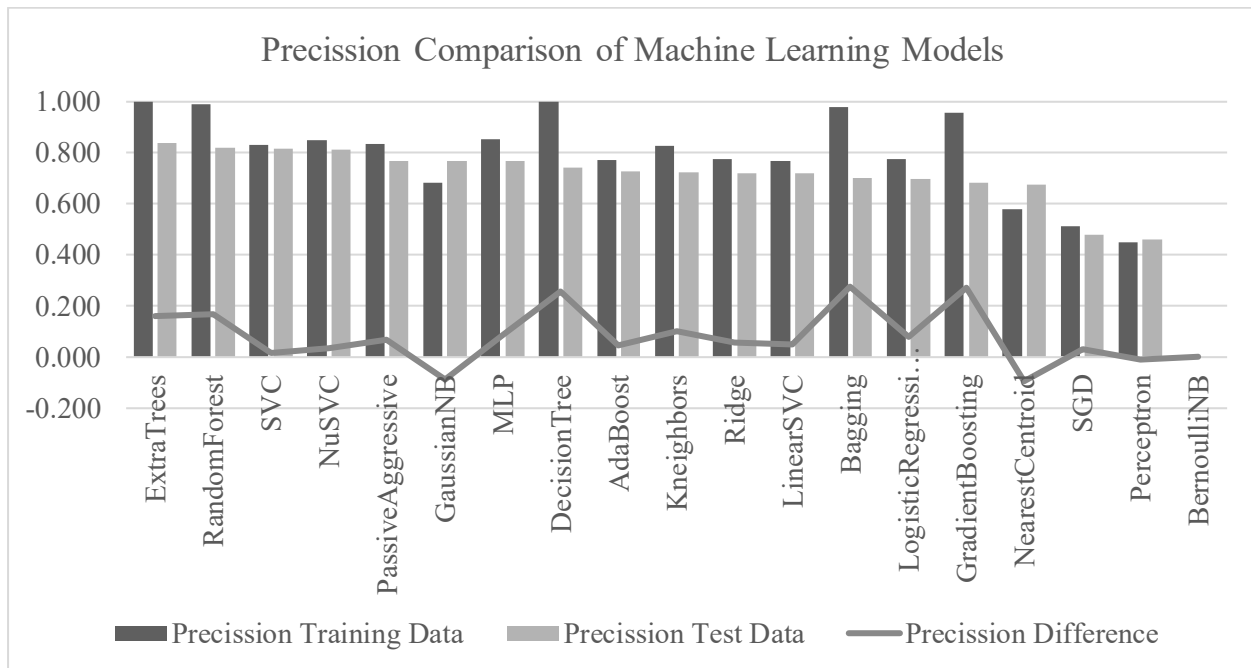


491
 492 **Figure 6.** Model comparison for accuracy: the graph is in descending order based on the accuracy on test
 493 data.

494 Nine models, including K-neighbors (accuracy=0.862), NuSVC (accuracy=0.858), MLP
 495 (accuracy=0.850), SVC (accuracy=0.823), AdaBoost (accuracy=0.819), PassiveAggressive
 496 (accuracy=0.815), LinearSVC (accuracy=0.811), LogisticRegression (0.807), and RidgeClassifier
 497 (accuracy=0.807) had the accuracy above 80%. Gaussian NB and Nearest Centroid also had the
 498 satisfactory accuracy of 0.795 and 0.704, respectively (Figure 6). However, SGD
 499 (accuracy=0.638), Bernouli NB (accuracy=0.622), and Perceptron (accuracy=0.539) had the
 500 poorest performance (Figure 6).

501 When the accuracy of 19 models on test data was compared it was apparent that a majority of the
 502 models had high variance in the accuracy on training and test datasets (Figure 6). Considering the
 503 accuracy as a model performance criterion, Gaussian NB model was found to have less variance
 504 in the accuracy where the overall accuracy on test data was 0.824, 0.028 greater than the accuracy
 505 on the training data. The second most stable model was NuSVC (accuracy=0.812) with a
 506 difference between training and testing dataset of 0.047. Other models with good performance on
 507 test data after Gaussian NB and NuSVC were, respectively, K-neighbors (accuracy=0.788),
 508 Random Forests (accuracy=0.788), Decision Trees (accuracy=0.777), Extra Trees
 509 (accuracy=0.777), SVC (accuracy=0.777), MLP (accuracy=0.765), AdaBoost (accuracy=0.753),
 510 Bagging (accuracy=0.753), Gradient Boosting (accuracy=0.741), LinearSVC (accuracy=0.741),

511 the Nearest Centroid (accuracy=0.741), Passive Aggressive (accuracy=0.741), Ridge
 512 (accuracy=0.741), and Logistic Regression (accuracy=0.729). Bernouli NB (accuracy=0.577),
 513 SGD (accuracy=0.541), and Perceptron (accuracy=0.506) models were the poorest performers.
 514 After comparing the accuracy of all the models of training and test data, Gaussian NB was found
 515 to be the clear winner and Perceptron was as good as an intuition, a random model.
 516 As mentioned in the methodology section, Precision measures the type-I error i.e., false positive.
 517 It is also known as how sensitive the model is correctly predicting the positive events. While
 518 comparing the precision score on the training data, we found that Extra Trees, Random Forests,
 519 and Decision Trees have the perfect precision of “1.” Bagging (precision=0.979) and Gradient
 520 Boosting (precision=0.957) almost equally precise after Extra Trees, Random Forests and
 521 Decision Trees. MLP (precision=0.854), NuSVC (precision=0.849), Passive Aggressive
 522 (precision=0.836), SVC (precision=0.831), and K-neighbors (precision=0.828) achieved good
 523 precision scores all above 80%. Ridge (precision=0.776), Logistic regression (precision=0.776),
 524 Ada Boost (precision=0.772), Linear SVC (precision=0.767) and Gaussian NB (precision=0.683)
 525 also had good precision score ranging between 70% and 80%. The Nearest Centroid
 526 (precision=0.581), SGD (precision=0.511), and Perceptron (precision=0.451) scored less
 527 precision. Bernouli NB (precision=0) couldn’t produce any precision score. It is clear that
 528 Bernouli failed to precisely predict the positive outcome, i.e. the preference of ATU as a
 529 sustainable arsenic mitigation technology (Figure 7).

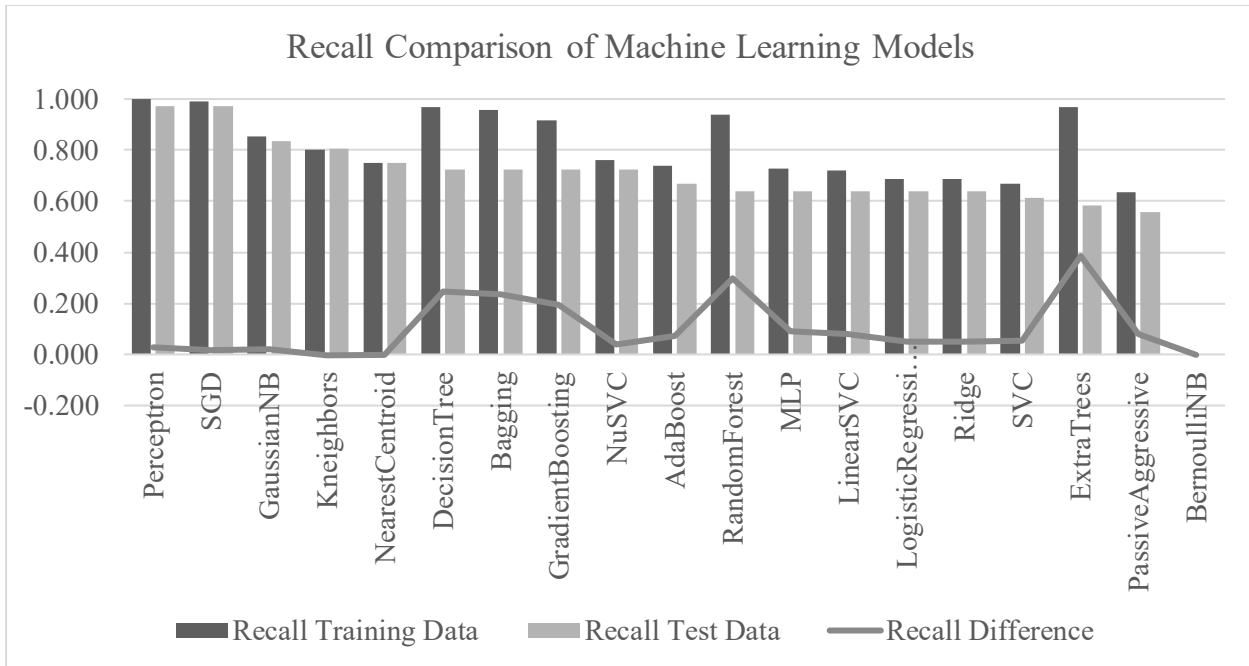


530

531 **Figure 7.** Model comparison for precision: the graph is in descending order based on the precision on test
532 data.

533 The precision on test data for Extra Trees got reduced from 1 to 0.840, still the highest among rest
534 of the models. Random Forest model had 0.821 precision value on test data, followed by SVC
535 (precision=0.815), NuSVC (precision=0.813), Passive Aggressive (precision=0.769),
536 GaussianNB (precision=0.769), MLP (precision=0.767), Decision Tree (precision=0.743), Ada
537 Boost (precision=0.727), K-neighbors (precision=0.725), Ridge (precision=0.719), Linear SVC
538 (precision=0.719), and Bagging (precision=0.703), all above 70%. Logistic regression
539 (precision=0.697), Gradient Boosting (precision=0.684), and Nearest Centroid (precision=0.675)
540 scored precision between 60% and 70%. SGD (precision=0.479), Perceptron (precision=0.461),
541 and Bernoulli NB (precision=0.000) had the lower precision scores, Bernoulli NB being the lowest
542 similar to its precision score on training data. That further states that when Bernoulli NB predicts
543 ATU as a sustainable arsenic mitigation technology, it is correct 0% of the time (Figure 7).

544 As opposed to Precision, Recall measures the type-II error i.e., false negative. It is also known as
545 how specific the model correctly predicts the negative events. While comparing the precision score
546 on the test data, we found that Perceptron had the perfect Recall of 1, followed by SGD
547 (recall=0.990), Decision Tree (recall=0.969), Extra Trees (recall=0.969), Bagging (recall=0.958),
548 Random Forest (recall=0.938), and Gradient Boosting (recall=0.917), all above 90%. Gaussian
549 NB (recall=0.854), K-neighbors (recall=0.802), NuSVC (recall=0.760), the Nearest Centroid
550 (recall=0.750), Ada Boost (recall=0.740), MLP (recall=0.729), and Linear SVC (recall=0.719) had
551 recall score above 70% but below 90%. Logistic Regression (recall=0.688), Ridge (recall=0.688),
552 SVC (recall=0.667), and Passive Aggressive (recall=0.635) had comparatively lower recall value
553 among all 19 models, and Bernoulli NB (recall=0.000) being at the bottom of this sequence.
554 (Figure 8).



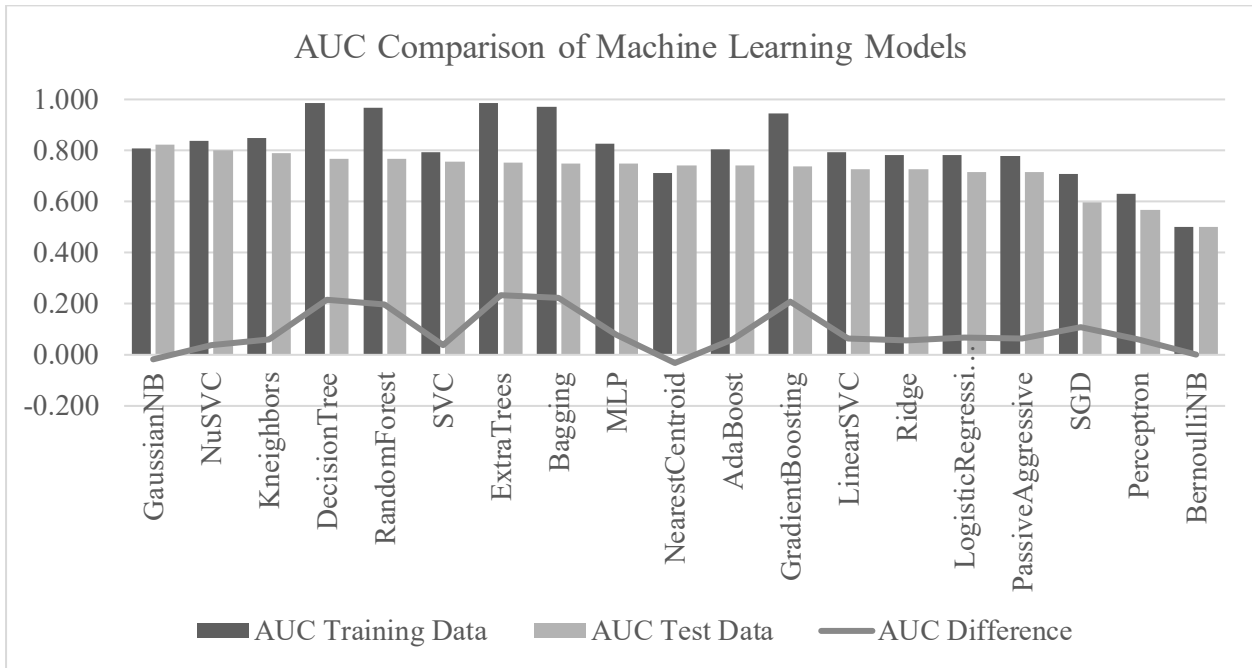
555

556 **Figure 8.** Model comparison for recall: the graph is in descending order based on the recall on test data.

557 The recall on test data for Perceptron got reduced from 1 to 0.972, still the highest among rest of
 558 the models. SGD achieved 0.972 recall score on the test data, followed by Gaussian NB
 559 (recall=0.833), K-neighbors (recall=0.806), the Nearest Centroid (recall=0.750), Decision Tree
 560 (recall=0.722), Bagging (recall=0.722), Gradient Boosting (0.722), Nu SVC (recall=0.722), Ada
 561 Boost (recall=0.667), Random Forest (recall=0.639), MLP (recall=0.639), Linear SVC
 562 (recall=0.639), Logistic regression (recall=0.639), Ridge (recall=0.639), SVC (recall=0.611),
 563 Extra Trees (recall=0.583), Passive Aggressive (recall=0.556), and Bernoulli NB (recall=0.000)
 564 (Figure 8). Perceptron and SGD had the highest recall on both training and testing data. On the
 565 other hand, Bernoulli NB achieved “0” recall score (Figure 8).

566 AUC is considered as the most preferable machine learning model performance metrics to evaluate
 567 the accuracy and the robustness of any machine learning models. We found that Decision Tree
 568 (AUC=0.984) and Extra Trees (AUC=0.984) achieved the highest AUC score on training data.
 569 Bagging (AUC=0.973) secured the third place, followed by Random Forest (AUC=0.966),
 570 Gradient Boosting (AUC=0.946), K-neighbors (AUC=0.850), NuSVC (AUC=0.839), MLP
 571 (AUC=0.827), Gaussian NB (AUC=0.807), Ada Boost (AUC=0.803), Linear SVC (AUC=0.793),
 572 SVC (AUC=0.792), Ridge (AUC=0.784), Logistic Regression (AUC=0.784), Passive Aggressive

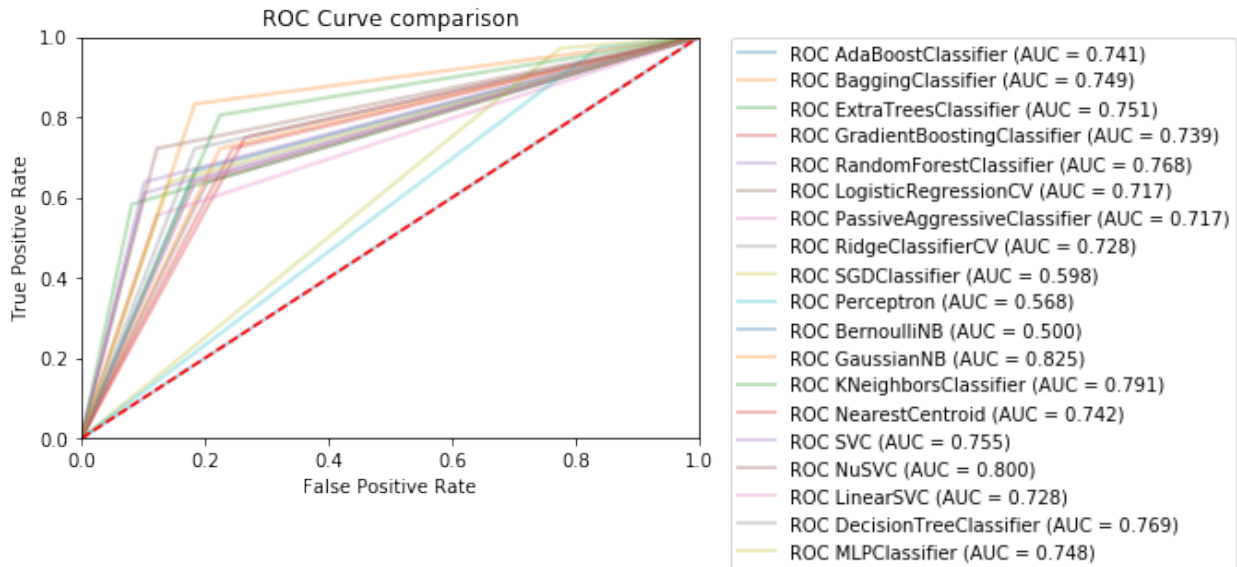
573 (AUC=0.780), Nearest Centroid (AUC=0.710), SGD (AUC=0.707), Perceptron (AUC=0.630),
 574 and Bernoulli NB (AUC=0.500) (Figure 9).



575
 576 **Figure 9.** Model comparison for AUC: the graph is in descending order based on the AUC on test data.

577 Gaussian NB had the best AUC of 0.825 on test data, followed by NuSVC (AUC=0.800), K-
 578 neighbors (AUC=0.791), Decision Trees (AUC=0.769), Random Forest (AUC=0.768), SVC
 579 (AUC=0.755), Extra Trees (AUC=0.751), Bagging (AUC=0.749), MLP (AUC=0.748), Nearest
 580 Centroid (AUC=0.742), Ada Boost (AUC=0.741), Gradient Boosting (AUC=0.739), Linear SVC
 581 (AUC=0.728), Ridge (AUC=0.728), Logistic Regression (AUC=0.717), Passive Aggressive
 582 (AUC=0.717), SGD (AUC=0.589), Perceptron (AUC=0.568), and Bernoulli NB (AUC=0.500)
 583 (Figure 10).

584



585

586 **Figure 10.** The AUC of all models on test data.

587 With the highest AUC of 0.825 on test data, GaussianNB can be declared as the best model in
 588 predicting ATU as a sustainable arsenic mitigation technology (Figure10).

589 4. DISCUSSION

590 The DT (21.4%), ExtraTrees (21.4%), Bagging (22.9%), Gradient Boosting (22.2%), and Random
 591 Forest (18.9%) models found to be over fitted on the training data while comparing on the test data
 592 (Figure 6). There were other ten models that overfitted above 5% on the training data. It appeared
 593 that these models learned the noise in the training data as a natural trend in the data that
 594 unfortunately could not be applied to test data. Therefore, the accuracy on test data went down by
 595 as high as 22.9%. That further indicates that these models cannot be generalized to unseen dataset.
 596 On the other hand, only two models including GaussianNB (-3.5%) and NearestCentroid (-5.8%)
 597 seem to be under-fitted on the training data, being very conservative in learning only the
 598 meaningful trend in the data and predicted with slightly less accuracy on test data. Each algorithm
 599 has its own advantages and disadvantages and have been developed to address specific scenario.
 600 Tree-based and ensemble algorithms are known to be prone to over-fitting if the sample size is not
 601 appropriate, assumptions are not met, or features are varied in nature [77, 78]. Among all,
 602 GaussianNB was found to be the most robust model. The reason could be that NB is a fast, highly
 603 scalable algorithm and is a good choice for binary classification problems [79, 80]. It can easily
 604 be updated on new data. In this study, we only had 339 samples, which is not a large sample size

605 and NB is known to be easily get trained on a small dataset that reflected from its highest AUC
606 score (Figure 10). Since it is highly scalable, this technique can easily be applied to other similar
607 arsenic contaminated areas [79, 80]. In a recent study, the authors found that the GaussianNB
608 outperformed the SVM as it is statistically robust, neutrally reasonable, and could reproduce across
609 unseen datasets [79]. Socioeconomic and psychological data comprises both continuous data (age
610 and income) and discrete data (gender, education level, marital status, etc.) and Gaussian NB works
611 pretty well with multidimensional data. In reality, survey-based data also suffers from missingness,
612 GaussianNB is also prone to missing data and over-fitting, and it also ignores irrelevant features
613 in the model. A majority of machine learning outcomes are difficult to interpret. Gaussian NB
614 provides predictive ability to users as it can make probabilistic predictions [80, 81].

615 5. CONCLUSION

616 Application of various techniques of artificial intelligence in environmental data modeling and
617 prediction is a recent phenomenon. Its application on building prediction models on socioeconomic
618 and psychological data collected from communities living in environmental contaminated regions
619 has just began. This study is a founding step in providing insights on how various state-of-the-art
620 artificial intelligence can be used for developing accurate prediction models of sustainable arsenic
621 mitigation technologies. Selecting an appropriate method is key in developing meaningful
622 prediction models as the machine (computational instruments) only understands the data as
623 numbers. In this study, we have evaluated several cutting edges linear, nonlinear, ensemble, tree-
624 based, and Naïve Bayes-based machine learning algorithms for predicting sustainable arsenic
625 mitigation technology. Gaussian NB found to be the best model to fit to such multidimensional
626 data. Achieving greater than 70% of AUC on test data by other 15 models is also promising. The
627 top three models Gaussian NB, NuSVC, and K-neighbors are nonlinear classifiers that considers
628 a nonlinear association between independent and dependent variables. The bottom performers
629 including logistic regression, passive aggressive, SGD, and perceptron are all linear classifiers and
630 could not do a better justice with this data. Bernoulli NB was clearly a bad choice of model as it
631 assumes the features should be binary and that's the reason it failed to develop a meaning
632 prediction model. Nonlinear and ensemble models are the better choice of models for
633 multidimensional data where the association between the features are not linear, but complex. We
634 also understand that if a few top linear and ensemble models can be further explored and fine-

635 tuned, we could probably enhance the model performance. Where funding to generate such data
636 in a large number is a challenge, Gaussian NB model is like a life-saving method that requires less
637 data to learn, can handle missing data, prone to over-fitting, and easy to interpret. We evaluated a
638 few less common algorithms that provided hope for exploring more for developing prediction
639 models on socioeconomic-environmental data including NuSVC, Extra Tree, and Nearest
640 Centroid. A larger sample size, careful feature selection, feature engineering may also help
641 improve the performance of these models. Some models do require hyper parameter tuning, thus
642 selecting the optimum hyper parameters for such models will also help improve the model
643 performance.

644 **References**

- 645 [1] S.K. Singh, R.W. Taylor, M.M. Rahman, B. Pradhan, Developing robust arsenic awareness
646 prediction models using machine learning algorithms, *Journal of environmental management*,
647 211 (2018) 125-137.
- 648 [2] S.K. Singh, E.A. Stern, *Global Arsenic Contamination: Living with the Poison Nectar*,
649 *Environment: Science and Policy for Sustainable Development*, 59 (2017) 24-28.
- 650 [3] S.K. Singh, R.W. Taylor, Assessing the role of risk perception in ensuring sustainable arsenic
651 mitigation, *Groundwater for Sustainable Development*, (2019) 100241.
- 652 [4] J. Bundschuh, M. Litter, V.S. Ciminelli, M.E. Morgada, L. Cornejo, S.G. Hoyos, J. Hoinkis,
653 M.T. Alarcon-Herrera, M.A. Armienta, P. Bhattacharya, Emerging mitigation needs and
654 sustainable options for solving the arsenic problems of rural and isolated urban areas in Latin
655 America—A critical analysis, *water research*, 44 (2010) 5828-5845.
- 656 [5] D. Chakraborti, S.K. Singh, H.M. Rashid, M.M. Rahman, *Arsenic: Occurrence in*
657 *Groundwater*, *Encyclopedia of Environmental Health*. Burlington: Elsevier, 2 (2017) 1-17.
- 658 [6] S. Murcott, *Arsenic contamination in the world*, IWA publishing, 2012.
- 659 [7] D. Chakraborti, S.K. Singh, M.M. Rahman, R.N. Dutta, S.C. Mukherjee, S. Pati, P.B. Kar,
660 *Groundwater Arsenic Contamination in the Ganga River Basin: A Future Health Danger*,
661 *International Journal of Environmental Research and Public Health*, 15 (2018) 180.
- 662 [8] A. Heikens, *Arsenic contamination of irrigation water, soil and crops in Bangladesh: Risk*
663 *implications for sustainable agriculture and food safety in Asia*, RAP Publication (FAO),
664 (2006).
- 665 [9] M.M. Moriarty, I. Koch, R.A. Gordon, K.J. Reimer, *Arsenic speciation of terrestrial*
666 *invertebrates*, *Environmental Science & Technology*, 43 (2009) 4818-4823.
- 667 [10] M. Bassil, F. Daou, H. Hassan, O. Yamani, J.A. Kharma, Z. Attieh, J. Elaridi, *Lead, cadmium*
668 *and arsenic in human milk and their socio-demographic and lifestyle determinants in*
669 *Lebanon*, *Chemosphere*, 191 (2018) 911-921.
- 670 [11] M. Molin, S.M. Ulven, H.M. Meltzer, J. Alexander, *Arsenic in the human food chain,*
671 *biotransformation and toxicology—Review focusing on seafood arsenic*, *Journal of Trace*
672 *Elements in Medicine and Biology*, 31 (2015) 249-259.
- 673 [12] M. Hossain, S.N. Rahman, P. Bhattacharya, G. Jacks, R. Saha, M. Rahman, *Sustainability of*
674 *arsenic mitigation interventions—an evaluation of different alternative safe drinking water*
675 *options provided in MATLAB, an arsenic hot spot in Bangladesh*, *Frontiers in Environmental*
676 *Science*, 3 (2015) 30.
- 677 [13] M. Hossain, P. Bhattacharya, G. Jacks, M. von Brömssen, K.M. Ahmed, M.A. Hasan, S.K.
678 *Frape, Sustainable arsenic mitigation—from field trials to implementation for control of*
679 *arsenic in drinking water supplies in Bangladesh*, in: *Best Practice Guide on the Control of*
680 *Arsenic in Drinking Water*, IWA Publishing UK, 2017, pp. 99-116.
- 681 [14] A. Kabir, G. Howard, *Sustainability of arsenic mitigation in Bangladesh: Results of a*
682 *functionality survey*, *International Journal of Environmental Health Research*, 17 (2007) 207-
683 218.
- 684 [15] M. Shafiquzzaman, M.S. Azam, I. Mishima, J. Nakajima, *Technical and social evaluation of*
685 *arsenic mitigation in rural Bangladesh*, *Journal of health, population, and nutrition*, 27 (2009)
686 674.

- 687 [16] N. Shibasaki, P. Lei, A. Kamata, Evaluation of deep groundwater development for arsenic
688 mitigation in western Bangladesh, *Journal of Environmental Science and Health Part A*, 42
689 (2007) 1919-1932.
- 690 [17] L.H. Winkel, P.T.K. Trang, V.M. Lan, C. Stengel, M. Amini, N.T. Ha, P.H. Viet, M. Berg,
691 Arsenic pollution of groundwater in Vietnam exacerbated by deep aquifer exploitation for
692 more than a century, *Proceedings of the National Academy of Sciences*, 108 (2011) 1246-
693 1251.
- 694 [18] M.M. Hira-Smith, Y. Yuan, X. Savarimuthu, J. Liaw, A. Hira, C. Green, T. Hore, P.
695 Chakraborty, O.S. Von Ehrenstein, A.H. Smith, Arsenic concentrations and bacterial
696 contamination in a pilot shallow dugwell program in West Bengal, India, *Journal of*
697 *Environmental Science and Health Part A*, 42 (2007) 89-95.
- 698 [19] S.K. Singh, Assessing and mapping vulnerability and risk perceptions to groundwater arsenic
699 contamination: Towards developing sustainable arsenic mitigation models (Order No.
700 3701365). Available from ProQuest Dissertations & Theses Full Text. (1681668682). in:
701 Earth and Environmental Studies, Montclair State University, USA, 2015, pp. 392.
- 702 [20] S.K. Singh, N. Vedwan, Mapping Composite Vulnerability to Groundwater Arsenic
703 Contamination: An Analytical Framework and a Case Study in India, *Natural Hazards*, 75
704 (2015) 1883-1908.
- 705 [21] S.K. Singh, An Analysis of the Cost-Effectiveness of Arsenic Mitigation Technologies:
706 Implications for Public Policy, *International Journal of Sustainable Built Environment*, 6
707 (2017) 522-535.
- 708 [22] S.K. Singh, R.W. Taylor, H. Su, Developing Sustainable Models of Arsenic-Mitigation
709 Technologies in the Middle-Ganga Plain in India, *Current Science*, 113 (2017) 80-93.
- 710 [23] S.K. Singh, R.W. Taylor, Assessing and Mapping Human Health Risks Due to Arsenic and
711 Socioeconomic Correlates for Proactive Arsenic Mitigation, in: *Arsenic Water Resources*
712 *Contamination*, Springer, 2020, pp. 231-256.
- 713 [24] B.K. Thakur, V. Gupta, Arsenic concentration in drinking water of Bihar: health issues and
714 socio-economic problems, *Journal of Water, Sanitation and Hygiene for Development*, 6
715 (2016) 331-341.
- 716 [25] B.K. Thakur, V. Gupta, Arsenic-Contaminated Drinking Water and the Associated Health
717 Effects in the Shahpur Block of Bihar: A Case Study From Five Villages, in: *Arsenic Water*
718 *Resources Contamination*, Springer, 2020, pp. 257-271.
- 719 [26] S. Priyadarshini, How the arsenic-affected perceive risk, in, 2014.
- 720 [27] S.K. Singh, R.W. Taylor, Likelihood of adoption of arsenic-mitigation technologies under
721 perceived risks to health, income, and social discrimination to arsenic contamination, in: Y.
722 Zhu, H. Guo, P. Bhattacharya, A. Ahmad, J. Bundschuh, R. Naidu (Eds.) *Environmental*
723 *Arsenic in a Changing World: Proceedings of the 7th International Congress and Exhibition*
724 *on Arsenic in the Environment (AS 2018)*, CRC Press, Beijing, P.R. China, 2018, pp. 700.
- 725 [28] R.M. Warner, *Applied statistics: From bivariate through multivariate techniques*, Sage
726 Publications, 2012.
- 727 [29] B.T. Pham, B. Pradhan, D.T. Bui, I. Prakash, M. Dholakia, A comparative study of different
728 machine learning methods for landslide susceptibility assessment: a case study of Uttarakhand
729 area (India), *Environmental Modelling & Software*, 84 (2016) 240-250.
- 730 [30] B.T. Pham, I. Prakash, J. Dou, S.K. Singh, P.T. Trinh, H. Trung Tran, T. Minh Le, V.P. Tran,
731 D. Kim Khoi, A. Shirzadi, A Novel Hybrid Approach of Landslide Susceptibility Modeling

- 732 Using Rotation Forest Ensemble and Different Base Classifiers, Geocarto International,
733 (2018) 1-38.
- 734 [31] B.T. Pham, I. Prakash, S.K. Singh, A. Shirzadi, H. Shahabi, D.T. Bui, Landslide susceptibility
735 modeling using Reduced Error Pruning Trees and different ensemble techniques: Hybrid
736 machine learning approaches, CATENA, 175 (2019) 203-218.
- 737 [32] T.V. Phong, T.T. Phan, I. Prakash, S.K. Singh, A. Shirzadi, K. Chapi, H.-B. Ly, L.S. Ho, N.K.
738 Quoc, B.T. Pham, Landslide susceptibility modeling using different artificial intelligence
739 methods: a case study at Muong Lay district, Vietnam, Geocarto International, (2019) 1-24.
- 740 [33] B. Pradhan, A comparative study on the predictive ability of the decision tree, support vector
741 machine and neuro-fuzzy models in landslide susceptibility mapping using GIS, Computers
742 & Geosciences, 51 (2013) 350-365.
- 743 [34] B.T. Pham, A. Jaafari, I. Prakash, S.K. Singh, N.K. Quoc, D.T. Bui, Hybrid computational
744 intelligence models for groundwater potential mapping, Catena, 182 (2019) 104101.
- 745 [35] W. Chen, B. Pradhan, S. Li, H. Shahabi, H.M. Rizeei, E. Hou, S. Wang, Novel hybrid
746 integration approach of bagging-based fisher's linear discriminant function for groundwater
747 potential analysis, Natural Resources Research, 28 (2019) 1239-1258.
- 748 [36] W.S. Jang, B. Engel, C.M. Yeum, Integrated environmental modeling for efficient aquifer
749 vulnerability assessment using machine learning, Environmental Modelling & Software, 124
750 (2020) 104602.
- 751 [37] L. Knoll, L. Breuer, M. Bach, Large scale prediction of groundwater nitrate concentrations
752 from spatial data using machine learning, Science of the total environment, 668 (2019) 1317-
753 1327.
- 754 [38] H.M. Rizeei, O.S. Azeez, B. Pradhan, H.H. Khamees, Assessment of groundwater nitrate
755 contamination hazard in a semi-arid region by using integrated parametric IPNOA and data-
756 driven logistic regression models, Environmental monitoring and assessment, 190 (2018) 633.
- 757 [39] F.-J. Chang, L.-s. Kao, Y.-M. Kuo, C.-W. Liu, Artificial neural networks for estimating
758 regional arsenic concentrations in a blackfoot disease area in Taiwan, Journal of hydrology,
759 388 (2010) 65-76.
- 760 [40] B. Purkait, S. Kadam, S. Das, Application of Artificial Neural Network Model to Study
761 Arsenic Contamination in Groundwater of Malda District, Eastern India, Journal of
762 Environmental Informatics, 12 (2008).
- 763 [41] K.H. Cho, S. Sthiannopkao, Y.A. Pachepsky, K.-W. Kim, J.H. Kim, Prediction of
764 contamination potential of groundwater arsenic in Cambodia, Laos, and Thailand using
765 artificial neural network, Water research, 45 (2011) 5535-5544.
- 766 [42] J.D. Ayotte, B.T. Nolan, J.A. Gronberg, Predicting arsenic in drinking water wells of the
767 Central Valley, California, Environmental science & technology, 50 (2016) 7555-7563.
- 768 [43] Y. Park, M. Ligaray, Y.M. Kim, J.H. Kim, K.H. Cho, S. Sthiannopkao, Development of
769 enhanced groundwater arsenic prediction model using machine learning approaches in
770 Southeast Asian countries, Desalination and Water Treatment, 57 (2016) 12227-12236.
- 771 [44] S. Singh, Arsenic contamination in water, soil, and food materials in Bihar, Lambert
772 Academic Publishing, Germany, 2011.
- 773 [45] S.K. Singh, A.K. Ghosh, Entry of arsenic into food material—a case study, World Appl Sci J,
774 13 (2011) 385-390.
- 775 [46] S.K. Singh, A.K. Ghosh, Health Risk Assessment due to Groundwater Arsenic
776 Contamination: Children are at High Risk, Human and Ecological Risk Assessment: An
777 International Journal, 18 (2012) 751-766.

- 778 [47] S.K. Singh, A. Ghosh, A. Kumar, K. Kislai, C. Kumar, R. Tiwari, R. Parwez, N. Kumar, M.
779 Imam, Groundwater Arsenic Contamination and Associated Health Risks in Bihar, India,
780 International Journal of Environmental Research, 8 (2014) 49-60.
- 781 [48] D. Chakraborti, S.C. Mukherjee, S. Pati, M.K. Sengupta, M.M. Rahman, U.K. Chowdhury,
782 D. Lodh, C.R. Chanda, A.K. Chakraborti, G.K. Basu, Arsenic groundwater contamination in
783 Middle Ganga Plain, Bihar, India: a future danger?, Environmental Health Perspectives, 111
784 (2003) 1194-1198.
- 785 [49] S. Ahamed, D. Chakraborti, Groundwater Arsenic contamination and Health Effects in Bihar
786 and UP, LAP Lambert Academic Publishing, Germany, 2012.
- 787 [50] D. Chakraborti, M.M. Rahman, S. Ahamed, R.N. Dutta, S. Pati, S.C. Mukherjee, Arsenic
788 contamination of groundwater and its induced health effects in Shahpur block, Bhojpur
789 district, Bihar state, India: risk evaluation, Environmental Science and Pollution Research, 23
790 (2016) 9492-9504.
- 791 [51] D. Chakraborti, M.M. Rahman, S. Ahamed, R.N. Dutta, S. Pati, S.C. Mukherjee, Arsenic
792 groundwater contamination and its health effects in Patna district (capital of Bihar) in the
793 middle Ganga plain, India, Chemosphere, 152 (2016) 520-529.
- 794 [52] T. Kluyver, B. Ragan-Kelley, F. Pérez, B.E. Granger, M. Bussonnier, J. Frederic, K. Kelley,
795 J.B. Hamrick, J. Grout, S. Corlay, Jupyter Notebooks-a publishing format for reproducible
796 computational workflows, in: ELPUB, 2016, pp. 87-90.
- 797 [53] W. McKinney, Python for data analysis: Data wrangling with Pandas, NumPy, and IPython,
798 " O'Reilly Media, Inc.", 2012.
- 799 [54] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P.
800 Prettenhofer, R. Weiss, V. Dubourg, Scikit-learn: Machine learning in Python, Journal of
801 machine learning research, 12 (2011) 2825-2830.
- 802 [55] P. Virtanen, R. Gommers, T.E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau, E.
803 Burovski, P. Peterson, W. Weckesser, J. Bright, SciPy 1.0: fundamental algorithms for
804 scientific computing in Python, Nature methods, (2020) 1-12.
- 805 [56] W. Winston, Microsoft Excel 2010 Data Analysis and Business Modeling: Data Analysis and
806 Business Modeling, Pearson Education, 2011.
- 807 [57] P.M. Atkinson, R. Massari, Generalised linear modelling of susceptibility to landsliding in
808 the central Apennines, Italy, Computers & Geosciences, 24 (1998) 373-385.
- 809 [58] W. Chen, X. Zhao, H. Shahabi, A. Shirzadi, K. Khosravi, H. Chai, S. Zhang, L. Zhang, J. Ma,
810 Y. Chen, Spatial prediction of landslide susceptibility by combining evidential belief function,
811 logistic regression and logistic model tree, Geocarto International, (2019) 1-25.
- 812 [59] D. Westreich, J. Lessler, M.J. Funk, Propensity score estimation: neural networks, support
813 vector machines, decision trees (CART), and meta-classifiers as alternatives to logistic
814 regression, Journal of clinical epidemiology, 63 (2010) 826-833.
- 815 [60] S.-B. Bai, J. Wang, G.-N. Lü, P.-G. Zhou, S.-S. Hou, S.-N. Xu, GIS-based logistic regression
816 for landslide susceptibility mapping of the Zhongxian segment in the Three Gorges area,
817 China, Geomorphology, 115 (2010) 23-31.
- 818 [61] A. Shirzadi, L. Saro, O.H. Joo, K. Chapi, A GIS-based logistic regression model in rock-fall
819 susceptibility mapping along a mountainous road: Salavat Abad case study, Kurdistan, Iran,
820 Natural hazards, 64 (2012) 1639-1656.
- 821 [62] S.Z. Mousavi, A. Kavian, K. Soleimani, S.R. Mousavi, A. Shirzadi, GIS-based spatial
822 prediction of landslide susceptibility using logistic regression model, Geomatics, Natural
823 Hazards and Risk, 2 (2011) 33-50.

- 824 [63] K. Crammer, O. Dekel, J. Keshet, S. Shalev-Shwartz, Y. Singer, Online passive-aggressive
825 algorithms, *Journal of Machine Learning Research*, 7 (2006) 551-585.
- 826 [64] Y. Freund, R.E. Schapire, A decision-theoretic generalization of on-line learning and an
827 application to boosting, in: *European conference on computational learning theory*, Springer,
828 1995, pp. 23-37.
- 829 [65] L. Breiman, Bagging predictors, *Machine learning*, 24 (1996) 123-140.
- 830 [66] L. Breiman, Random forests, *Machine learning*, 45 (2001) 5-32.
- 831 [67] N.S. Altman, An introduction to kernel and nearest-neighbor nonparametric regression, *The*
832 *American Statistician*, 46 (1992) 175-185.
- 833 [68] S. Shanmuganathan, Artificial neural network modelling: An introduction, in: *Artificial*
834 *Neural Network Modelling*, Springer, 2016, pp. 1-14.
- 835 [69] L. Jing, J. Hudson, Numerical method in rock engineering, *International Journal of Rock*
836 *Mechanics and Mining Sciences*, 39 (2002) 409-427.
- 837 [70] I. Yilmaz, O. Kaynar, Multiple regression, ANN (RBF, MLP) and ANFIS models for
838 prediction of swell potential of clayey soils, *Expert systems with applications*, 38 (2011)
839 5958-5966.
- 840 [71] D. Kanungo, M. Arora, S. Sarkar, R. Gupta, A comparative study of conventional, ANN black
841 box, fuzzy and combined neural and fuzzy weighting procedures for landslide susceptibility
842 zonation in Darjeeling Himalayas, *Engineering Geology*, 85 (2006) 347-366.
- 843 [72] C. Polykretis, C. Chalkias, Comparison and evaluation of landslide susceptibility maps
844 obtained from weight of evidence, logistic regression, and artificial neural network models,
845 *Natural hazards*, 93 (2018) 249-274.
- 846 [73] A. Shirzadi, H. Shahabi, K. Chapi, D.T. Bui, B.T. Pham, K. Shahedi, B.B. Ahmad, A
847 comparative study between popular statistical and machine learning methods for simulating
848 volume of landslides, *Catena*, 157 (2017) 213-226.
- 849 [74] R. Tibshirani, T. Hastie, B. Narasimhan, G. Chu, Diagnosis of multiple cancer types by
850 shrunken centroids of gene expression, *Proceedings of the National Academy of Sciences*, 99
851 (2002) 6567-6572.
- 852 [75] R. Tibshirani, T. Hastie, B. Narasimhan, G. Chu, Class prediction by nearest shrunken
853 centroids, with applications to DNA microarrays, *Statistical Science*, 18 (2003) 104-117.
- 854 [76] M. Pardo, G. Sberveglieri, Random forests and nearest shrunken centroids for the
855 classification of sensor array data, *Sensors and Actuators B: Chemical*, 131 (2008) 93-99.
- 856 [77] M. LeBlanc, J. Crowley, A review of tree-based prognostic models, in: *Recent advances in*
857 *clinical trial design and analysis*, Springer, 1995, pp. 113-124.
- 858 [78] M. Re, G. Valentini, 1 Ensemble methods: a review 3, (2012).
- 859 [79] R.D. Raizada, Y.-S. Lee, Smoothness without smoothing: why Gaussian naive Bayes is not
860 naive for multi-subject searchlight studies, *PloS one*, 8 (2013).
- 861 [80] F. Pereira, M. Botvinick, Information mapping with pattern classifiers: a comparative study,
862 *Neuroimage*, 56 (2011) 476-496.
- 863 [81] K.P. Murphy, Naive bayes classifiers, *University of British Columbia*, 18 (2006) 60.
- 864