

Detecting fabrication in large-scale molecular omics data

Michael Bradshaw (✉ michael.bradshawiii@colorado.edu)

University of Colorado Boulder <https://orcid.org/0000-0002-1451-8105>

Samuel H Payne

Brigham Young University

Research article

Keywords: data integrity, machine learning, omics

Posted Date: April 21st, 2020

DOI: <https://doi.org/10.21203/rs.3.rs-22923/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published at PLOS ONE on November 30th, 2021. See the published version at <https://doi.org/10.1371/journal.pone.0260395>.

1 Detecting fabrication in large-scale 2 molecular omics data

3

4 Michael S. Bradshaw^{1,2}, Samuel H. Payne¹

5 1. Biology Department, Brigham Young University, Provo UT 84602 USA

6 2. Computer Science Department, University of Colorado Boulder, Boulder CO 80309 USA

7 Corresponding Author: Michael S. Bradshaw, michael.bradshawiii@colorado.edu

8 Abstract

9 **Background:** Fraud is a pervasive problem and can occur as fabrication, falsification,
10 plagiarism or theft. The scientific community is not exempt from this universal problem and
11 several studies have recently been caught manipulating or fabricating data. Current measures
12 to prevent and deter scientific misconduct come in the form of the peer-review process and on-
13 site clinical trial auditors. As recent advances in high-throughput omics technologies have
14 moved biology into the realm of big-data, fraud detection methods must be updated for
15 sophisticated computational fraud. In the financial sector, machine learning and digit-preference
16 are successfully used to detect fraud.

17 **Results:** Drawing from these sources, we develop methods of fabrication detection in
18 biomedical research and show that machine learning can be used to detect fraud in large-scale
19 omic experiments. Using the raw data as input, the best machine learning models correctly
20 predicted fraud with 84-95% accuracy. With digit frequency as input features, the best models
21 detected fraud with 98%-100% accuracy. All of the data and analysis scripts used in this project
22 are available at <https://github.com/MSBradshaw/FakeData>.

23 **Conclusions:** Using digit frequencies as a generalized representation of the data, multiple
24 machine learning methods were able to identify fabricated data with near perfect accuracy.

25

26 **Keywords:**

27 data integrity, machine learning, omics

28 Introduction

29 Fraud is a pervasive problem and can occur as fabrication, falsification, plagiarism or theft.
30 Examples of fraud are found in virtually every field, such as education, commerce and
31 technology. With the rise of electronic crimes, specific criminal justice and regulatory bodies
32 have been formed to detect sophisticated fraud, creating an arms-race between methods to
33 deceive and methods to detect deception. The scientific community is not exempt from the
34 universal problem of fraud, and several studies have recently been caught manipulating or
35 fabricating data (George and Buyse, 2015; Kupferschmidt, 2018) or are suspected of it (Al-
36 Marzouki *et al.*, 2005). More than two million scientific articles are published yearly and ~2% of
37 authors admit to data fabrication (Fanelli, 2009). When asked if their colleagues had fabricated
38 data, positive response rates rose to 14-19% (Fanelli, 2009; George and Buyse, 2015). This
39 potentially means tens to hundreds of thousands of articles are published each year with
40 manipulated data.

41

42 Fraud in biological data represents a growing problem, as the scale of datasets can make it
43 easier to hide data manipulation. Recent advances in high-throughput omics technologies have
44 moved biology into the realm of big-data. Many diseases are now characterized in populations,
45 with thousands of individuals characterized for cancer (Blum *et al.*, 2018), diabetes (TEDDY
46 Study Group, 2007), bone strength (Orwoll *et al.*, 2005), and health care services for the general

47 populace (Bycroft *et al.*, 2018). Large-scale characterization studies are also done for cell lines
48 and drug responses (Barretina *et al.*, 2012; Subramanian *et al.*, 2017). This transition in biology,
49 where deep molecular characterization of biological samples is now routinely available, points to
50 a future where clinical trial requirements might include omics data collection.

51
52 Current trial monitoring methods include auditing, site monitoring, data reviews and central
53 monitoring (Knepper *et al.*, 2016; Baigent *et al.*, 2008). The decision to use these forms of
54 oversight and at what frequency is not driven by empirical data but rather is determined by
55 clinics' usual practice (Morrison *et al.*, 2011). The emerging data deluge challenges the
56 effectiveness of traditional auditing practices to detect fraud, and several studies have
57 suggested addressing the issue with improved centralized and independent statistical
58 monitoring (Baigent *et al.*, 2008; George and Buyse, 2015; Calis *et al.*, 2017). However, these
59 recommendations are given chiefly to help ensure the safety and efficacy of the study, not data
60 integrity.

61
62 In 1937, physicist Frank Benford observed in a compilation of 20,000 numbers that the first digit
63 did not follow a uniform distribution as one may anticipate (Benford, 1938). This pattern holds
64 true in most large collections of numbers, including scientific data. Comparing a distribution of
65 first digits to a Benford distribution can be used to identify deviations from the expected
66 frequency, often because of fraud. Recently Benford's law has been used to identify fraud in
67 financial records of international trade (Barabesi *et al.*, 2018) and money laundering (Badal-
68 Valero *et al.*, 2018). It has also been used in smaller scale to reaffirm suspicions of fraud in
69 clinical trials (Al-Marzouki *et al.*, 2005) and medical studies (Hein *et al.*, 2012).

70
71 The distinction between fraud and honest error is important to make. Fraud is the intent to cheat
72 (George and Buyse, 2015). This is the definition used throughout this paper. An honest error

73 might be, forgetting to include a few samples, while intentionally excluding samples would be
74 fraud. Copying and pasting values from one table to another incorrectly is an honest error but
75 intentionally changing the values is fraud. In these examples the results maybe the same but
76 the intent behind them differs wildly. In application, identifying the difference in intent behind
77 something like this may well be impossible. This paper focuses on a specific type of fraud free
78 of this ambiguity: data fabrication. Data fabrication is “making up data or results and recording
79 or reporting them” (George and Buyse, 2015).

80
81 We explore methods of data fabrication and detection in molecular omics data using supervised
82 machine learning and Benford digit preferences. The data used in this study comes from the
83 Clinical Proteomic Tumor Analysis Consortium (CPTAC) cohort for endometrial carcinoma,
84 which contains copy number alteration (CNA) measurements from 100 tumor samples. We
85 created 50 additional fake samples, or fake patients, for these datasets. Three different methods
86 of varying sophistication are used for fabrication: random number generation, resampling with
87 replacement and imputation. We show that machine learning and digit-preference can be used
88 to detect fraud with near perfect accuracy.

89

90 Methods

91 Real Data. The real data used in this publication originated from the genomic analysis of uterine
92 endometrial cancer. As part of the Clinical Proteomics Tumor Analysis Consortium (CPTAC),
93 100 tumor samples underwent whole genome and whole exome sequencing and subsequent
94 copy number analysis. We used the results of the copy number analysis as *is*, which is stored in
95 our GitHub repository at <https://github.com/MSBradshaw/FakeData>.

96 Fake Data. Fake data used in this study was generated using three different methods. In each
97 method, we created 50 fake samples which were combined with the 100 real samples to form a
98 mixed dataset. The first method to generate fake data was random number generation. For
99 every gene locus, we first find the maximum and minimum values observed in the original data.
100 A new sample is then fabricated by randomly picking a value within this gene specific range.
101 The second method to generate fake data was sampling with replacement. For this, we create
102 lists of all observed values across the cohort for each gene. A fake sample is created by
103 randomly sampling from these lists with replacement. The third method to generate fake data
104 was imputation. The R package missForrest (Stekhoven and Bühlmann, 2012) was repurposed
105 for data fabrication. A fake sample was generated by first creating a copy of a real sample. Then
106 we iteratively nullified 10% of the data and imputed these NAs with missForrest until every value
107 has been imputed. See Supplemental Figure 1.

108

109 Machine Learning Training. With a mixed dataset containing 100 real samples and 50 fake
110 samples, we proceeded to create and evaluate machine learning models which predict whether
111 a sample is real or fabricated (Supplemental Figure 2). The 100 real and 50 fake samples were
112 both randomly split in half, one portion added to a training set and the other held out for testing.
113 Using Python's SciKitLearn library, we evaluated multiple machine learning models, gradient
114 boosting (GBD), Naïve Bayes, Random Forest, K-Nearest Neighbor (KNN), Multi-layer
115 Perceptron (MLP) and Support Vector Machine (SVM). Training validation was done using 10-
116 fold cross validation. We note explicitly that the training routine was never able to use testing
117 data. After all training was complete, the held-out test set was then fed to each model for
118 prediction and scoring. We used simple accuracy as a metric. For each sample in the test set,
119 ML models would predict whether it was real or fabricated. Model accuracy was calculated as
120 the number of correct predictions divided by the number of total predictions. The entire process
121 of fake data generation and ML training/testing was repeated 50 times. Different random seeds

122 were used when generating each set of fake data. Thus, fake samples in all 50 iterations are
123 distinct from each other. All of the data and analysis scripts used in this project are available at
124 <https://github.com/MSBradshaw/FakeData>.

125
126 Benford-Like Digit Preferences. Benford's Law or the first digit law has been instrumental at
127 catching fraud in various financial situations (Barabesi *et al.*, 2018; Badal-Valero *et al.*, 2018)
128 and in small scale clinical trials (Al-Marzouki *et al.*, 2005). The method presented here is
129 designed with the potential to generalize and be applied to multiple sets of data of varying types
130 and configurations (e.i. different measured variables (features) and different quantities of
131 variables). Machine learning typically cannot handle data where the features are not consistent
132 in number and type. Converting all measured variables to digit frequencies circumvents this
133 problem. Digit frequencies are calculated as the number of occurrences of a single digit (0-9)
134 divided by the total number of features. In the method described in this paper, a sample's
135 features are all converted to digit frequencies of the first and second digit after the decimal.
136 Thus, for each sample the features are converted from ~17,000 copy number alterations to 20
137 digit preferences. Using this approach, whether a sample has 100 or 17,000 features it can still
138 be trained on and classified by the same model.

139 Results

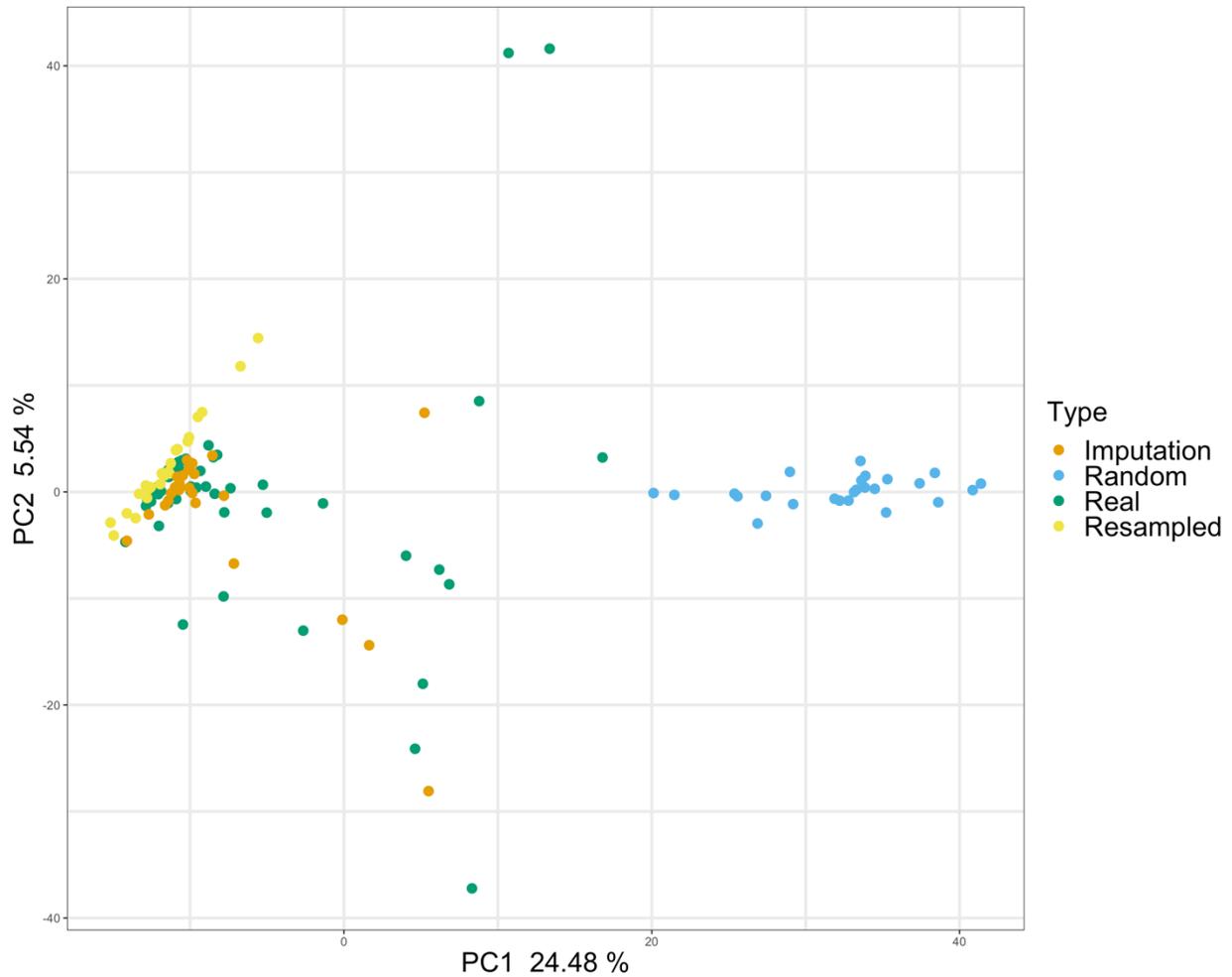
140 Our goal is to explore the ability of machine learning methods to identify fabricated data hidden
141 within large datasets. Although there are many situations where data fabrication might be
142 present, we chose a scenario wherein researchers are trying to obtain enough subjects to have
143 a sufficiently powered study. For example, if the power calculation estimates that 75 subjects
144 are required to observe a given effect, but financial constraints only allow for a cohort of 50

145 subjects. In this scenario, researchers might be tempted to fabricate data for an additional 25
146 subjects to meet regulatory requirements.

147 Fake Data

148 The real data used in this study comes from the Clinical Proteomic Tumor Analysis Consortium
149 (CPTAC) cohort for endometrial carcinoma, specifically the copy number alteration (CNA) data.
150 This real data was paired with fabricated data and used as an input to machine learning
151 classification models (see Methods). Three different methods of data fabrication were used in
152 this study: random number generation, resampling with replacement, and imputation
153 (Supplemental Figure 1). The three methods represent three realistic ways that an unscrupulous
154 scientist might create novel data. Each method has benefits and disadvantages, with imputation
155 being both the most sophisticated and also the most computationally intense and complex. As
156 seen in Figure 1, the random data clusters far from the real data. Both the resampled and
157 imputed data cluster tightly with the real data in a PCA plot, with the imputed data also
158 generating a few reasonable outlier samples.

159



160

161 **Figure 1 - Principal Component Analysis of real and fake samples.** Copy number data for
 162 the real and fabricated samples are shown. The fabricated data created via random number
 163 generation is clearly distinct from all other data. Fabricated data created via resampling or
 164 imputation appears to cluster very closely with the real data.

165

166 To look further into the fabricated data, we examined whether fake data preserved correlative
 167 relationships present in the original data (Supplemental Figure 3). This is exemplified by two
 168 pairs of genes. PLEKHN1 and HES4 are adjacent genes found on chromosome 1p36 separated
 169 by ~30,000 bp. Because they are so closely located on the chromosome, it is expected that
 170 most copy number events like large scale duplications and deletions would include both genes.

171 As expected, their CNA data has a Spearman correlation coefficient of 1.0 in the original data, a
172 perfect correlation. The second pair of genes, DFFB and OR4F5, are also on chromosome 1,
173 but are separated by 3.8 Mbp. As somewhat closely located genes, we would expect a modest
174 correlation between CNA measurements, but not as highly correlated as the adjacent gene pair.
175 Consistent with this expectation, their CNA data has a Spearman correlation coefficient of 0.27.
176 Depending on the method of fabrication, fake data for these two gene pairs may preserved
177 these correlative relationships. When we look at the random and resampled data for these two
178 genes, all correlation is lost (Supplemental Figure 3 C, D, E and F). Imputation, however,
179 produces data that closely matches the original correlations, PLEKHN1 and HES4 $R^2 = 0.97$;
180 DFFB and OR4F5 $R^2 = 0.32$ (Supplemental Figure 3 G and H).

181 Machine learning with quantitative data

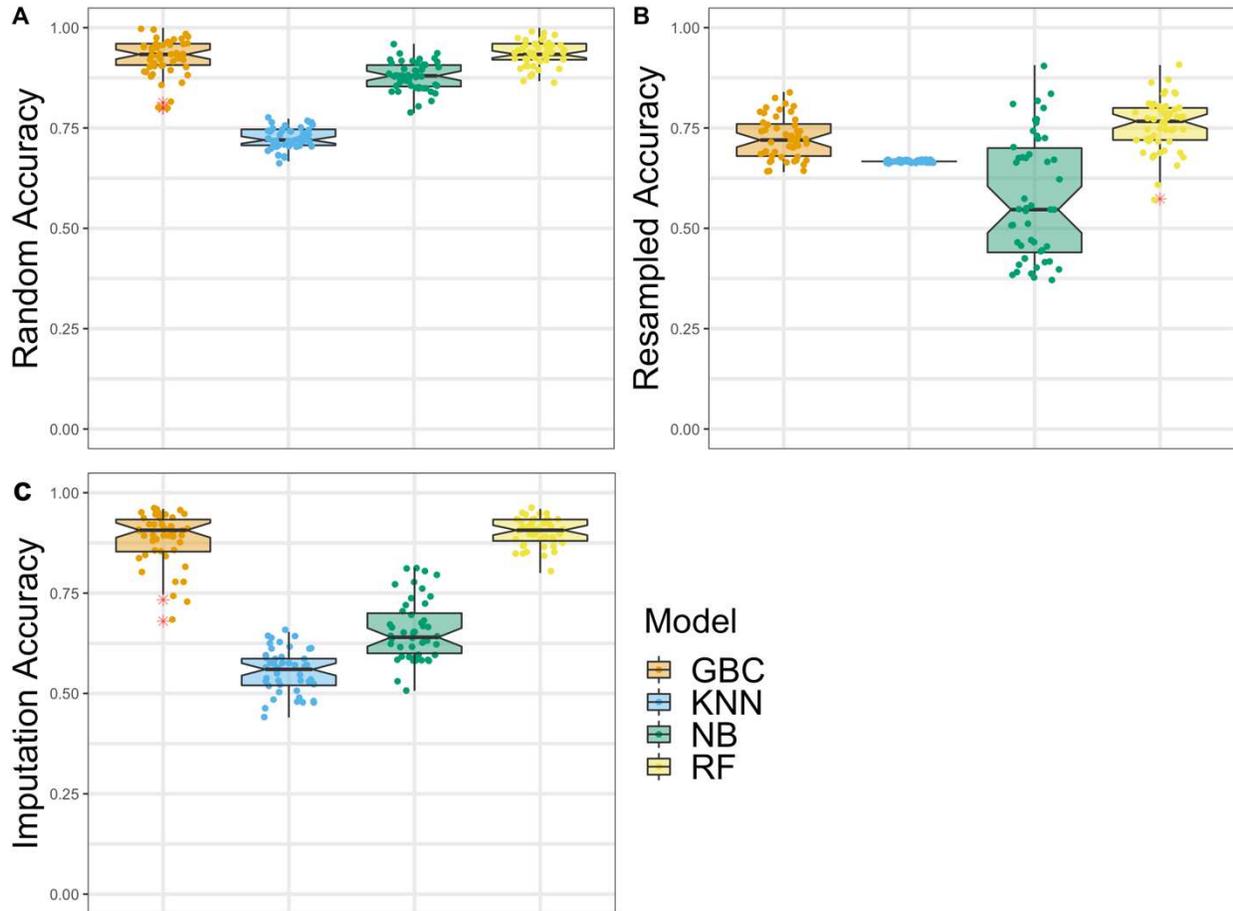
182 We tested six different methods for machine learning to create a model capable of detecting
183 fabricated data: Gradient Boosting (GBC), Naïve Bayes, Random Forest, K-Nearest Neighbor
184 (KNN), Multi-layer Perceptron (MLP) and Support Vector Machine (SVM). Models were given as
185 features the quantitative data table containing copy number data on 75 labeled samples, 50 real
186 and 25 fake. In the copy number data, each sample had measurements for ~17,000 genes,
187 meaning that each sample had ~17,000 features. After training, the model was asked to classify
188 held-out testing data containing 75 samples, 50 real and 25 fake. We evaluated the model on
189 simple accuracy, whether the predicted label was correct or incorrect. To ensure that our results
190 represent robust performance, model training and evaluation was repeated 50 times with 50
191 different fabricated datasets (see Methods). Reported results represent the average accuracy of
192 these 50 trials. We note that two methods, SVM and MLP, performed poorly compared to other
193 classification methods. Testing data was comprised of 2/3 real data and 1/3 fake data; therefore,
194 baseline accuracy (the accuracy achieved if the model predicting all test samples as the

195 majority class) is 66%. Both SVM and MLP had an average accuracy at or below this baseline
196 for classification of the simplest fabrication method (random), and were excluded from further
197 analysis.

198

199 The remaining four models performed relatively well on the classification task for data fabricated
200 with the random approach. The average accuracy of 50 trials was: Random Forest 94%, GBC
201 92%, Naïve Bayes 88%, and KNN 72% (Figure 2A). Mean classification accuracies were lower
202 for data created with the resampling method, with most models losing ~10% accuracy (Random
203 Forest 84%, GBC 83%, Naïve Bayes 73%, and KNN 70%). We also note that the variability in
204 model performance was much higher for classification of the resampled data (Figure 2B). As the
205 resampling method uses data values from the real data, it is possible that fake samples
206 sometimes more closely resemble real samples. Imputation classification results fluctuated
207 (Random Forest 90%, GBC 89%, Naïve Bayes 66%, and KNN 56%). While Random Forest and
208 GBC both increased in accuracy compared to the resampled data, Naïve Bayes and KNN both
209 now perform at or below the baseline accuracy (Figure 2C).

210



211

212

213 **Figure 2 - Classification accuracy using copy number data.** Fabricated data was mixed with

214 real data and given to four machine learning models for classification. Data shown represents

215 50 trials for 50 different fabricated dataset mixes. Features in this dataset are the copy number

216 values for each sample. **A.** Results for data fabricated with the random method, mean

217 classification accuracy: Random Forest 94% (+/- 3.1%), GBC 92% (+/- 4.5%), Naïve Bayes

218 88% (+/- 3.5%), and KNN 72% (+/- 2.6%). **B.** Results for data fabricated with the resampling

219 method, mean classification accuracy: Random Forest 84% (+/- 6.5%), GBC 83% (+/- 5.2%),

220 Naïve Bayes 73% (+/- 15.2%), and KNN 70% (+/- 0%). **C.** Results for data fabricated with the

221 imputation method, mean classification accuracy: Random Forest 90% (+/- 3.4%), GBC 89%

222 (+/- 6.4%), Naïve Bayes 66% (+/- 7.4%), and KNN 56% (+/- 5.3%).

223

224 Machine learning with digit preference

225 We were unsatisfied with the classification accuracy of the above models. One challenge for
226 machine learning in our data is that the number of features (~17,000) far exceeds the number of
227 samples (75). We therefore explored ways to reduce or transform the feature set, and also to
228 make the feature set more general and broadly applicable. Intrigued by the success of digit
229 frequency methods in the identification of financial fraud (Badal-Valero *et al.*, 2018), we
230 evaluated whether this type of data representation could work for bioinformatics data as well.
231 Therefore, all copy number data was transformed into 20 features, representing the digits 0-9 in
232 the first and second place after the decimal of each gene expression value. While Benford's
233 Law describes the frequency of the first digit, genomics and proteomics data are frequently
234 normalized or scaled and so the first digit may not be as characteristic. For this reason, our
235 method may be accurately referred to as Benford's Law inspired or Benford-like. These features
236 were tabulated for each sample to create a new data representation and fed into the exact same
237 machine learning training and testing routine described above. Each of these 20 new features
238 contain decimal values ranging from 0.0 to 1.0 representative of the proportional frequency that
239 digit occurs. For example, one sample's value in the feature column for the digit 1 may contain
240 the value 0.3. This means that in this sample's original data the digit 1 occurred in the first
241 position after the decimal place 30% of the time.

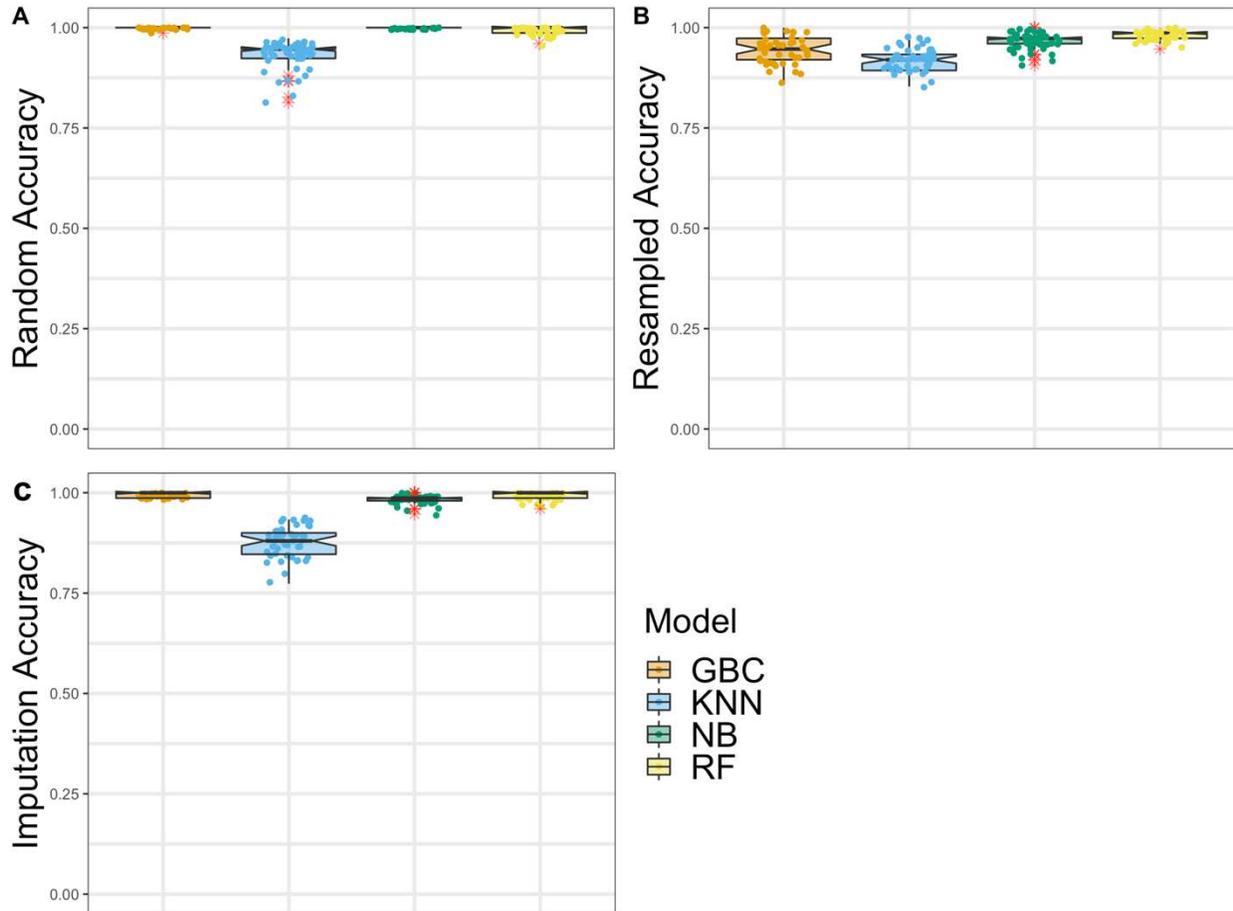
242

243 In addition to reducing the number of features, converting all features into digit frequencies
244 improves the model's generality. Machine learning typically cannot handle data where the
245 features are not consistent in number and type. Converting all measured variables to digit
246 frequencies circumvents this problem. For instance, if you had a data set of CNA and
247 transcriptomic data a machine learning model could not train and test on both of these. The

248 features in these datasets would differ in the number of features and what these features
249 represent. But once all information has been converted into digit frequencies the number and
250 type of features are standardized, enabling the model to work any number of different datasets.
251 Utilizing digit frequencies also allows the model to more easily handle missing data, which is
252 common in some types of omics measurements, such as single cell measurements.

253

254 In sharp contrast to the models built on the quantitative copy number data, machine learning
255 models which utilized the digit frequencies were highly accurate and showed little variability over
256 the 50 trials (Figure 3). When examining the results of the data fabricated via imputation (both
257 the most sophisticated and most realistic), the models achieved impressively high accuracy. As
258 an average accuracy for the 50 trials, both random forest and the gradient boosting models
259 achieved 100% accuracy. The naïve Bayes model was highly successful with a mean
260 classification accuracy 97%.



261

262

263 **Figure 3 - Classifications accuracy using digit frequency data.** Fabricated data was mixed

264 with real data and given to four machine learning models for classification. Data shown

265 represents 50 trials for 50 different fabricated dataset mixes. Features in this dataset are the

266 digit frequencies for each sample. **A.** Results for data fabricated with the random method, mean

267 classification accuracy: Random Forest 99% (+/- 1.0%), GBC 100% (+/- 0.2%), Naïve Bayes

268 100% (+/- 0.0%), and KNN 93% (+/- 3.4%). **B.** Results for data fabricated with the resampling

269 method, mean classification accuracy: Random Forest 98% (+/- 1.3%), GBC 94% (+/- 3.5%),

270 Naïve Bayes 97% (+/- 2.1%), and KNN 92% (+/- 2.8%). **C.** Results for data fabricated with the

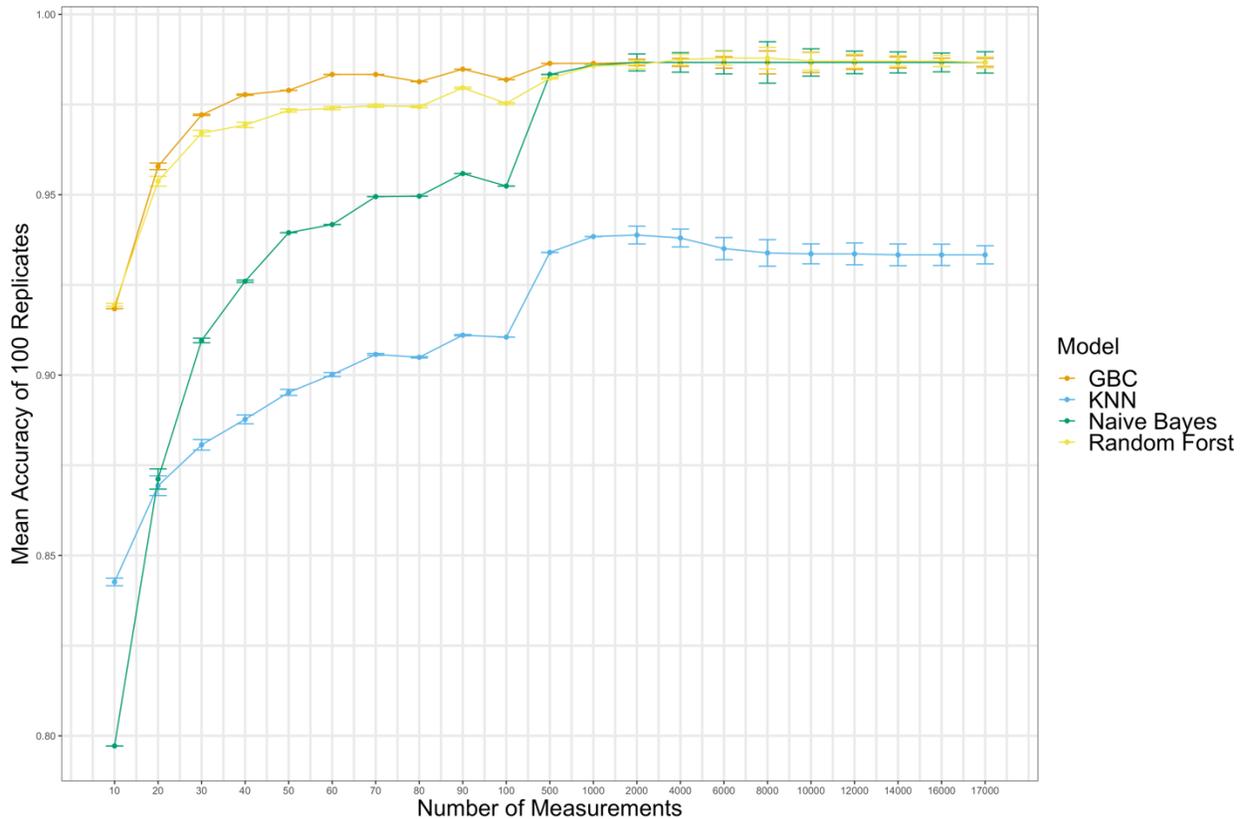
271 imputation method, mean classification accuracy: Random Forest 100% (+/- 1.0%), GBC 100%

272 (+/- 0.7%), Naïve Bayes 97% (+/- 1.1%), and KNN 89% (+/- 3.8%).

273

274 Machine learning with limited data

275 With 17,000 CNA gene measurements, the digit frequencies represent a well sampled
276 distribution. Theoretically, we realize that if one had an extremely limited dataset with CNA
277 measurements for only 10 genes, the sampling of the frequencies for the 10 digits will be poor.
278 To understand how much data is required for a good sampling of the digit-frequencies, we
279 iteratively downsampled our measurements from 17,000 to 10. With the gene-features
280 remaining in each downsample, the digit frequencies were re-calculated. Downsampling was
281 performed uniformly at random without replacement. For each measurement size 100 replicates
282 were run, all with different permutations of the downsamples. Results from this experiment can
283 be seen in Figure 4. The number of gene-features used to calculate digit frequencies does not
284 appear to make a difference at $n > 500$. In the 100 gene-feature trial, both Naive Bayes and
285 KNN have a significant drop in performance, while the Random Forest and Gradient Boosting
286 model remained relatively unaffected down to approximately 40 features. Surprisingly, these top
287 performing models (GBC and Random Forest) do not drop below 95% accuracy until they have
288 less than 20 gene-features.
289



290

291 **Figure 4 - Classifications accuracy vs number of features.** The original 17,000 CNA
 292 measurements were randomly downsampled incrementally to 10 and converted to digit-
 293 frequency training and test features for machine learning models. When 500+ measurements
 294 are used in the create of digit-preference features, there appears to be little to no effect on
 295 mean accuracy. Below 500, Naive Bayes and KNN models begin to lose accuracy quickly. GBC
 296 and Random Forest do suffer in accuracy as the number measurements used to generate
 297 features lowers but remain above 95% accurate until less than 20 measurements are included.

298

299 One hesitation for using machine learning with smaller datasets (i.e. fewer gene-features per
 300 data point) is the perceived susceptibility to large variation in performance. As noted, these
 301 downsampling experiments were performed 100 times, and error bars representing the standard
 302 error are shown in Figure 4. We note that even for the smallest datasets, performance does not
 303 noticeably vary between the 100 trials. In fact, the standard error for small datasets (e.g. 20 or

304 30 gene-features) is lower than when there were thousands. Thus, we believe that the digit-
305 frequency based models will perform well on both large-scale omics data and also on smaller
306 'targeted' data acquisition paradigms like multiplexed PCR or MRM proteomics.

307

308 Discussion

309 We present here a proof of concept method for detecting fabrication in biomedical data. Just as
310 has been previously shown in the financial sector, digit frequencies are a powerful data
311 representation when used in combination with machine learning to classify whether data is real
312 or falsified. While multiple methods of fabrication were used, we acknowledge there are more
313 subtle or sophisticated methods that could be used. We believe that fraud detection methods,
314 like the models presented herein, could be refined and generalized for broad use in monitoring
315 and oversight.

316

317 The scenario simulated in this study (50 real and 25 fabricated samples), we hope is not
318 representative of the scale of fraud in publications. Supervised machine learning for binary
319 classification requires this magnitude of fabrication for training and independent validation. After
320 training, however, that the models used here make predictions based on individual samples, not
321 whole datasets.

322

323 Regulatory bodies at multiple levels could enforce scientific integrity through the application of
324 these methods. For example, the government bodies charged with evaluating the efficacy of
325 new medicine could employ such techniques to screen large datasets that are submitted as
326 evidence for the approval of new drugs. Within individual clinical trials, data monitoring
327 committees could add a data consistency check to the statistical monitoring currently utilized.

328 For fundamental research, publishers could mandate the submission of all data to fraud
329 monitoring. Although journals commonly use software tools to detect plagiarism, a generalized
330 computational tool focused on data could make data fraud detection equally simple.

331

332 **Declarations**

333 **Ethics approval and consent to participate**

334 All data used in this publication is publicly available from the CPTAC consortium, which enrolls
335 participants and receives patient consent.

336 **Consent for publication**

337 All data used in this publication is publicly available from the CPTAC consortium, which enrolls
338 participants and receives patient consent.

339 **Availability of data and materials**

340 The datasets supporting the conclusions of this article are available in a git repository,
341 <https://github.com/MSBradshaw/FakeData>.

342 **Competing interests**

343 The authors declare that they have no competing interests.

344 **Funding**

345 This work was supported by the National Cancer Institute (NCI) CPTAC award [U24
346 CA210972].

347 **Author Contributions**

348 MSB and SHP designed the project, analyzed data and results, and wrote the publication.

349 **Acknowledgments**

350 Not applicable

351

352

353

354

355 References

356

357 Al-Marzouki,S. *et al.* (2005) Are these data real? Statistical methods for the detection of data
358 fabrication in clinical trials. *BMJ*, **331**, 267–270.

359 Badal-Valero,E. *et al.* (2018) Combining Benford’s Law and machine learning to detect money
360 laundering. An actual Spanish court case. *Forensic Sci. Int.*, **282**, 24–34.

361 Baigent,C. *et al.* (2008) Ensuring trial validity by data quality assurance and diversification of
362 monitoring methods. *Clin Trials*, **5**, 49–55.

363 Barabesi,L. *et al.* (2018) Goodness-of-Fit Testing for the Newcomb-Benford Law With
364 Application to the Detection of Customs Fraud. *Journal of Business & Economic
365 Statistics*, **36**, 346–358.

366 Barretina,J. *et al.* (2012) The Cancer Cell Line Encyclopedia enables predictive modelling of
367 anticancer drug sensitivity. *Nature*, **483**, 603–607.

368 Benford,F. (1938) The Law of Anomalous Numbers. *Proceedings of the American Philosophical
369 Society*, **78**, 551–572.

370 Blum,A. *et al.* (2018) SnapShot: TCGA-Analyzed Tumors. *Cell*, **173**, 530.

371 Bycroft,C. *et al.* (2018) The UK Biobank resource with deep phenotyping and genomic data.
372 *Nature*, **562**, 203–209.

373 Calis,K.A. *et al.* (2017) Recommendations for data monitoring committees from the Clinical
374 Trials Transformation Initiative. *Clin Trials*, **14**, 342–348.

375 Fanelli,D. (2009) How many scientists fabricate and falsify research? A systematic review and
376 meta-analysis of survey data. *PLoS ONE*, **4**, e5738.

377 George,S.L. and Buyse,M. (2015) Data fraud in clinical trials. *Clin Investig (Lond)*, **5**, 161–173.

378 Hein,J. *et al.* (2012) Scientific fraud in 20 falsified anesthesia papers : detection using financial
379 auditing methods. *Anaesthetist*, **61**, 543–549.

380 Knepper,D. *et al.* (2016) Detecting Data Quality Issues in Clinical Trials: Current Practices and
381 Recommendations. *Ther Innov Regul Sci*, **50**, 15–21.

382 Kupferschmidt,K. (2018) Tide of lies. *Science*, **361**, 636–641.

383 Morrison,B.W. *et al.* (2011) Monitoring the quality of conduct of clinical trials: a survey of current
384 practices. *Clin Trials*, **8**, 342–349.

385 Orwoll,E. *et al.* (2005) Design and baseline characteristics of the osteoporotic fractures in men
386 (MrOS) study--a large observational study of the determinants of fracture in older men.
387 *Contemp Clin Trials*, **26**, 569–585.

388 Stekhoven,D.J. and Bühlmann,P. (2012) MissForest--non-parametric missing value imputation
389 for mixed-type data. *Bioinformatics*, **28**, 112–118.

390 Subramanian,A. *et al.* (2017) A Next Generation Connectivity Map: L1000 Platform and the First
391 1,000,000 Profiles. *Cell*, **171**, 1437-1452.e17.

392 TEDDY Study Group (2007) The Environmental Determinants of Diabetes in the Young
393 (TEDDY) study: study design. *Pediatr Diabetes*, **8**, 286–298.

394

Figures

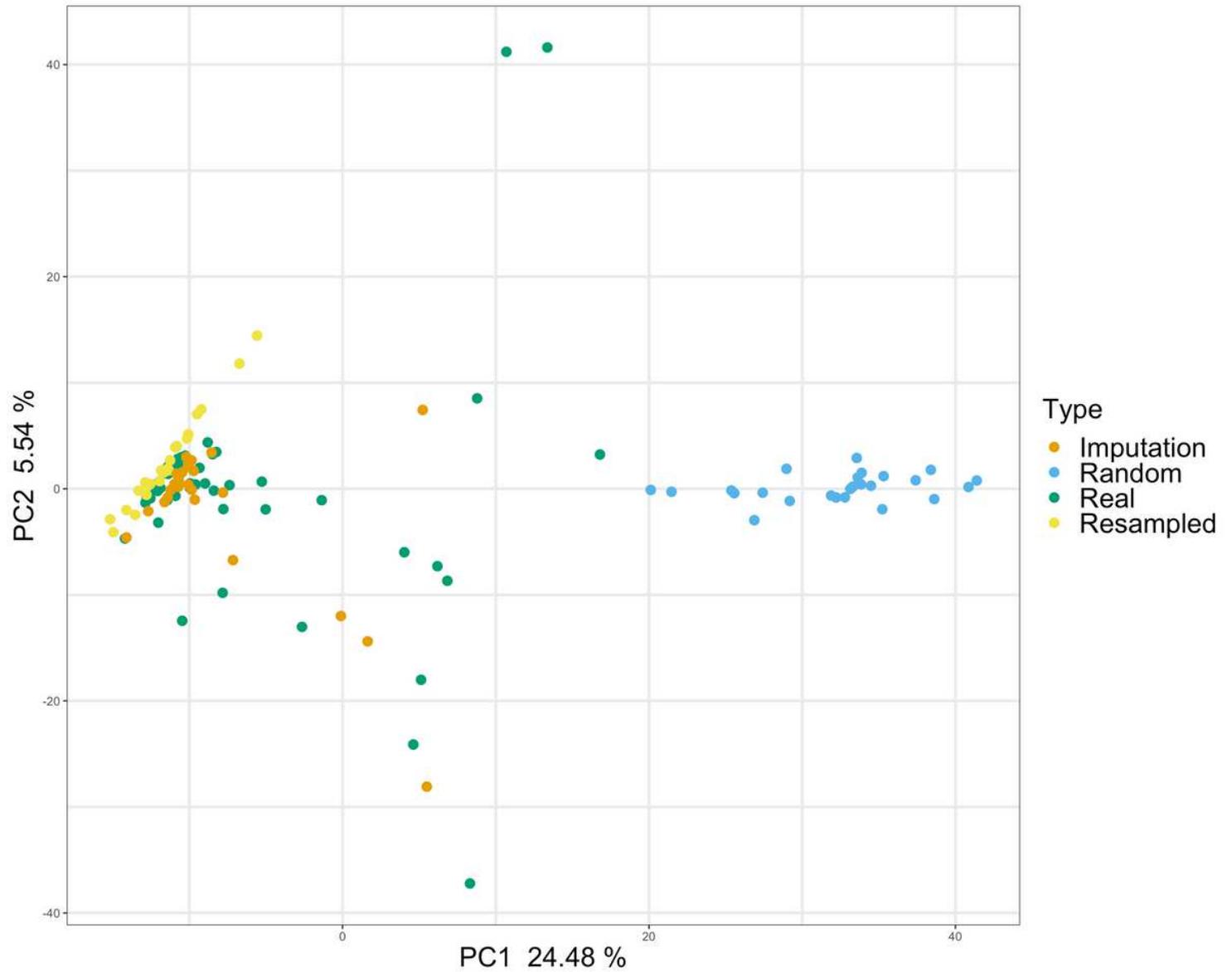


Figure 1

Principal Component Analysis of real and fake samples. Copy number data for the real and fabricated samples are shown. The fabricated data created via random number generation is clearly distinct from all other data. Fabricated data created via resampling or imputation appears to cluster very closely with the real data.

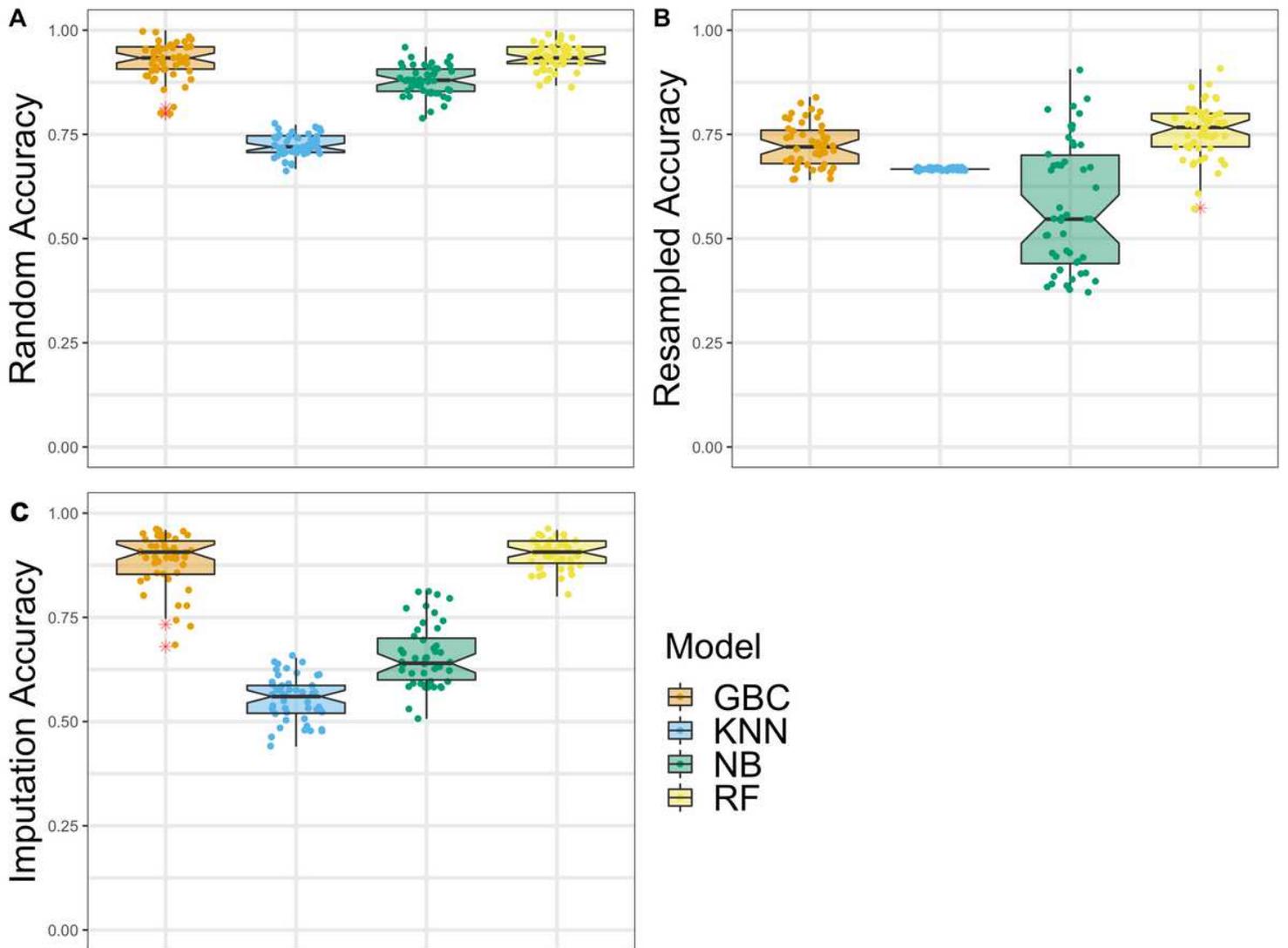


Figure 2

Classification accuracy using copy number data. Fabricated data was mixed with real data and given to four machine learning models for classification. Data shown represents 50 trials for 50 different fabricated dataset mixes. Features in this dataset are the copy number values for each sample. A. Results for data fabricated with the random method, mean classification accuracy: Random Forest 94% (+/- 3.1%), GBC 92% (+/- 4.5%), Naïve Bayes 88% (+/- 3.5%), and KNN 72% (+/- 2.6%). B. Results for data fabricated with the resampling method, mean classification accuracy: Random Forest 84% (+/- 6.5%), GBC 83% (+/- 5.2%), Naïve Bayes 73% (+/- 15.2%), and KNN 70% (+/- 0%). C. Results for data fabricated with the imputation method, mean classification accuracy: Random Forest 90% (+/- 3.4%), GBC 89% (+/- 6.4%), Naïve Bayes 66% (+/- 7.4%), and KNN 56% (+/- 5.3%).

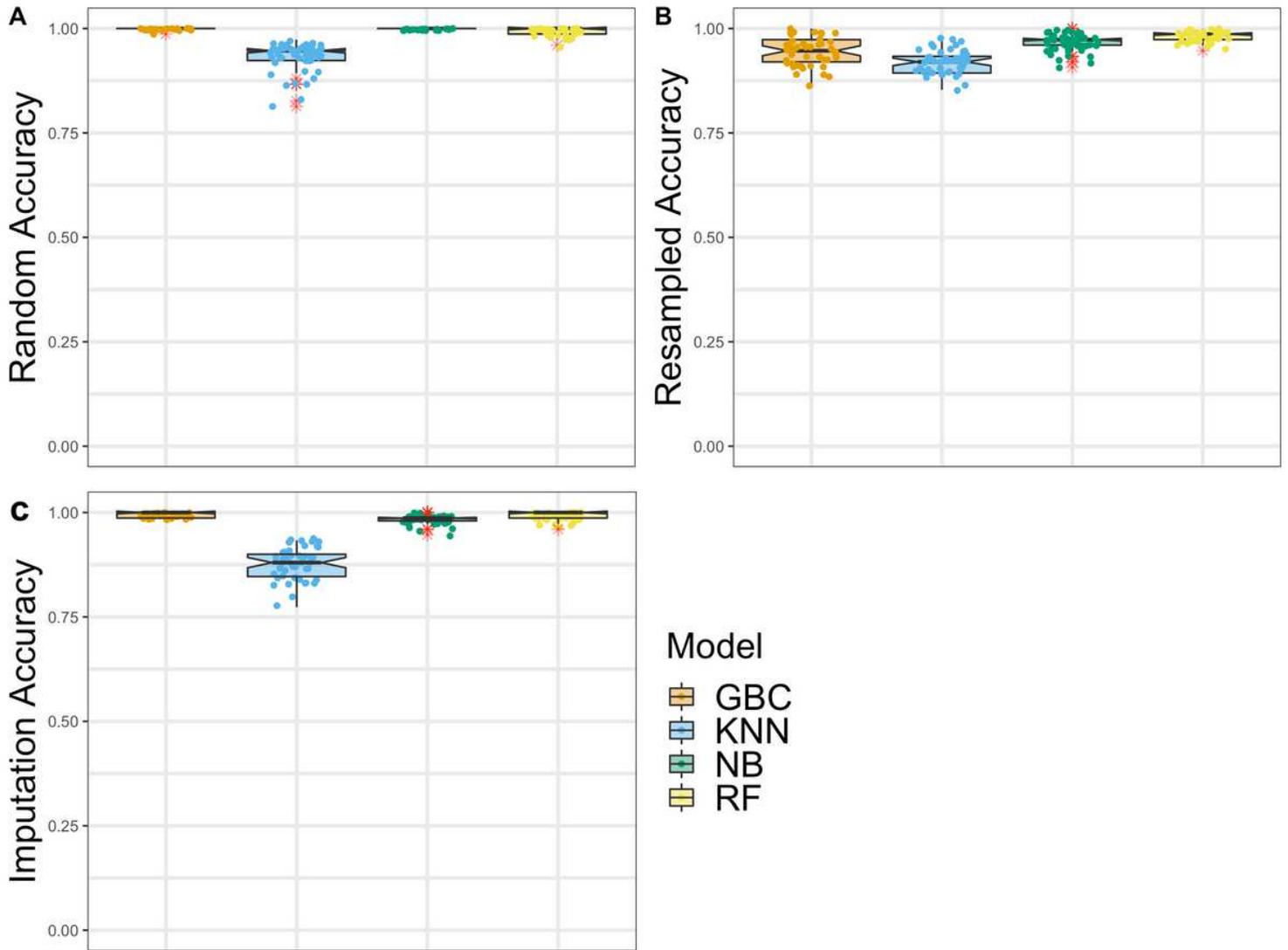


Figure 3

Classifications accuracy using digit frequency data. Fabricated data was mixed with real data and given to four machine learning models for classification. Data shown represents 50 trials for 50 different fabricated dataset mixes. Features in this dataset are the digit frequencies for each sample. A. Results for data fabricated with the random method, mean classification accuracy: Random Forest 99% (+/- 1.0%), GBC 100% (+/- 0.2%), Naïve Bayes 100% (+/- 0.0%), and KNN 93% (+/- 3.4%). B. Results for data fabricated with the resampling method, mean classification accuracy: Random Forest 98% (+/- 1.3%), GBC 94% (+/- 3.5%), Naïve Bayes 97% (+/- 2.1%), and KNN 92% (+/- 2.8%). C. Results for data fabricated with the imputation method, mean classification accuracy: Random Forest 100% (+/- 1.0%), GBC 100% (+/- 0.7%), Naïve Bayes 97% (+/- 1.1%), and KNN 89% (+/- 3.8%).

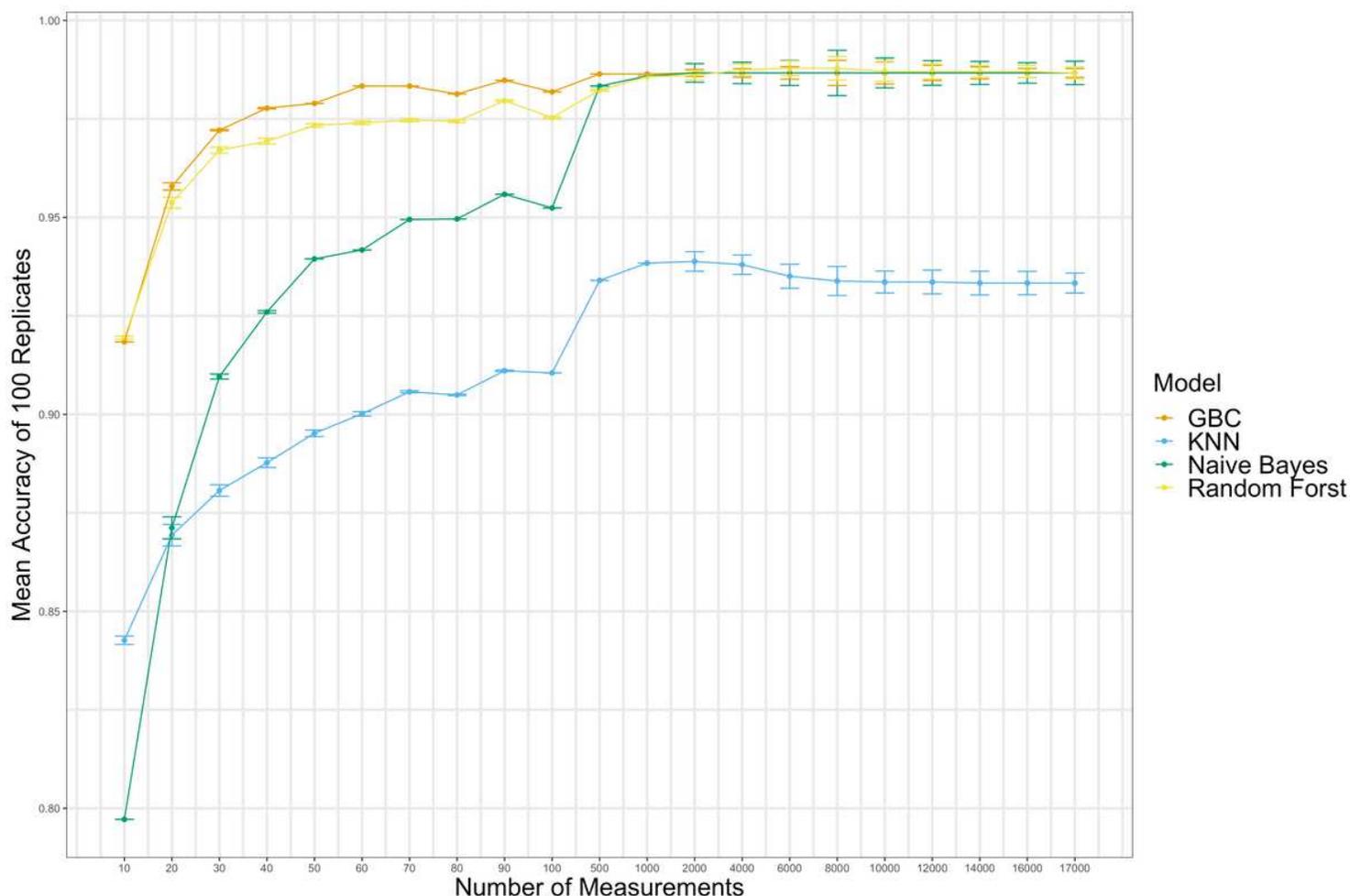


Figure 4

Classifications accuracy vs number of features. The original 17,000 CNA measurements were randomly downsampled incrementally to 10 and converted to digit-frequency training and test features for machine learning models. When 500+ measurements are used in the create of digit-preference features, there appears to be little to no effect on mean accuracy. Below 500, Naive Bayes and KNN models begin to lose accuracy quickly. GBC and Random Forest do suffer in accuracy as the number measurements used to generate features lowers but remain above 95% accurate until less than 20 measurements are included.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [FakeDatasupplement.pdf](#)