

 Open access • Posted Content • DOI:10.1101/757070

## Detecting fabrication in large-scale molecular omics data — [Source link](#)

Michael Bradshaw, Samuel H. Payne

**Institutions:** University of Colorado Boulder, Brigham Young University

**Published on:** 21 Oct 2020 - bioRxiv (Cold Spring Harbor Laboratory)

Related papers:

- [Fraud detection](#)
- [Parameters of automated fraud detection techniques during online transactions](#)
- [Detecting fraud in cellular telephone networks](#)
- [Using Analytics to Detect Possible Fraud: Tools and Techniques](#)
- [Machine learning forensics to gauge the likelihood of fraud in emails](#)

Share this paper:    

View more about this paper here: <https://typeset.io/papers/detecting-fabrication-in-large-scale-molecular-omics-data-4np2e93apy>

# 1 Detecting fabrication in large-scale 2 molecular omics data

3 Michael S. Bradshaw<sup>1</sup>, Samuel H. Payne<sup>2</sup>

4 1. Computer Science Department, University of Colorado Boulder, Boulder CO 80309 USA

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

6 **Contact:**

7 [michael.bradshawiii@colorado.edu](mailto:michael.bradshawiii@colorado.edu)

## 8 Abstract

9 Fraud is a pervasive problem and can occur as fabrication, falsification, plagiarism or theft. The  
10 scientific community is not exempt from this universal problem and several studies have  
11 recently been caught manipulating or fabricating data. Current measures to prevent and deter  
12 scientific misconduct come in the form of the peer-review process and on-site clinical trial  
13 auditors. As recent advances in high-throughput omics technologies have moved biology into  
14 the realm of big-data, fraud detection methods must be updated for sophisticated computational  
15 fraud. In the financial sector, machine learning and digit-preference are successfully used to  
16 detect fraud. Drawing from these sources, we develop methods of fabrication detection in  
17 biomedical research and show that machine learning can be used to detect fraud in large-scale  
18 omic experiments. Using the raw data as input, the best machine learning models correctly  
19 predicted fraud with 84-95% accuracy. With digit frequency as input features, the best models  
20 detected fraud with 98%-100% accuracy. All of the data and analysis scripts used in this project  
21 are available at <https://github.com/MSBradshaw/FakeData>.

## 22 Introduction

23 Fraud is a pervasive problem and can occur as fabrication, falsification, plagiarism or theft.  
24 Examples of fraud are found in virtually every field, such as: education, commerce and  
25 technology. With the rise of electronic crimes, specific criminal justice and regulatory bodies  
26 have been formed to detect sophisticated fraud, creating an arms-race between methods to  
27 deceive and methods to detect deception. The scientific community is not exempt from the  
28 universal problem of fraud, and several studies have recently been caught manipulating or  
29 fabricating data [1,2] or are suspected of it [3]. More than two million scientific articles are  
30 published yearly and ~2% of authors admit to data fabrication [4]. When asked if their  
31 colleagues had fabricated data, positive response rates rose to 14-19% [4,5]. Some domains or

32 locales have somewhat higher rates of data fabrication; in a recent survey of researchers at  
33 Chinese hospitals, 7.37% of researchers admitted to fabricating data [6]. Overall, these rates of  
34 data fabrication potentially means tens to hundreds of thousands of articles are published each  
35 year with manipulated data.

36  
37 Data in the biological sciences is particularly vulnerable to fraud given its size - which makes it  
38 easier to hide data manipulation - and researcher's dependence on freely available public data.  
39 Recent advances in high-throughput omics technologies have moved biology into the realm of  
40 big-data. Many diseases are now characterized in populations, with thousands of individuals  
41 characterized for cancer [7], diabetes [8] , bone strength [9], and health care services for the  
42 general populace [10]. Large-scale characterization studies are also done for cell lines and drug  
43 responses [11,12]. With the rise of importance of these large datasets, it becomes imperative  
44 that they remain free of errors both unintentional and intentional [13].

45  
46 Current methods for ensuring the validity of research is largely limited to the peer-review  
47 process which as of late has proven to be insufficient at spotting blatant duplication of images  
48 [14], let alone subtleties hidden in large scale data. Data for clinical trials can be subject to  
49 reviews and central monitoring [15,16]. However, the decision regarding oversight methodology  
50 and frequency is not driven by empirical data, but rather is determined by clinics' usual practice  
51 [17]. The emerging data deluge challenges the effectiveness of traditional auditing practices to  
52 detect fraud, and several studies have suggested addressing the issue with improved  
53 centralized and independent statistical monitoring [5,6,16,18]. However, these  
54 recommendations are given chiefly to help ensure the safety and efficacy of the study, not data  
55 integrity.

56

57 In 1937, physicist Frank Benford observed in a compilation of 20,000 numbers that the first digit  
58 did not follow a uniform distribution as one may anticipate [19]. This pattern holds true in most  
59 large collections of numbers, including scientific data. Comparing a distribution of first digits to a  
60 Benford distribution can be used to identify deviations from the expected frequency, often  
61 because of fraud. Recently Benford's law has been used to identify fraud in financial records of  
62 international trade [20] and money laundering [21]. It has also been used on a smaller scale to  
63 reaffirm suspicions of fraud in clinical trials [3].

64  
65 The distinction between fraud and honest error is important to make. Fraud is the intent to cheat  
66 [5]. This is the definition used throughout this paper. An honest error might be, forgetting to  
67 include a few samples, while intentionally excluding samples would be fraud. Copying and  
68 pasting values from one table to another incorrectly is an honest error but intentionally changing  
69 the values is fraud. In these examples the results may be the same but the intent behind them  
70 differs wildly. In efforts to maintain data integrity, identifying the intent of the misconduct may be  
71 impossible, and is also a secondary consideration after suspect data has been identified.

72  
73 Data fabrication is "making up data or results and recording or reporting them" [5]. This type of  
74 data manipulation is free from the above ambiguity relating to the author's intent. Making up  
75 data is always wrong. We explore methods of data fabrication and detection in molecular omics  
76 data using supervised machine learning and Benford-like digit-frequencies. We do not attempt  
77 to explain why someone may choose to fabricate their data - as other study have done [6,22];  
78 our only goal is to evaluate the utility of digit-frequencies to differentiate real from fake data. The  
79 data used in this study comes from the Clinical Proteomic Tumor Analysis Consortium (CPTAC)  
80 cohort for endometrial carcinoma, which contains copy number alteration (CNA) measurements  
81 from 100 tumor samples. We created 50 additional fake samples for these datasets. Three  
82 different methods of varying sophistication are used for fabrication: random number generation,

83 resampling with replacement and imputation. We show that machine learning and digit-  
84 preference can be used to detect fraud with near perfect accuracy.

## 85 Methods

### 86 Real Data

87 The real data used in this publication originated from the genomic analysis of uterine  
88 endometrial cancer. As part of the Clinical Proteomics Tumor Analysis Consortium (CPTAC),  
89 100 tumor samples underwent whole genome and whole exome sequencing and subsequent  
90 copy number analysis. We used the results of the copy number analysis *as is*, which is stored in  
91 our GitHub repository at <https://github.com/MSBradshaw/FakeData>.

92

### 93 Fake Data

94 Fake data used in this study was generated using three different methods. In each method, we  
95 created 50 fake samples which were combined with the 100 real samples to form a mixed  
96 dataset. The first method to generate fake data was random number generation. For every gene  
97 locus, we first find the maximum and minimum values observed in the original data. A new  
98 sample is then fabricated by randomly picking a value within this gene specific range. The  
99 second method to generate fake data was sampling with replacement. For this, we create lists  
100 of all observed values across the cohort for each gene. A fake sample is created by randomly  
101 sampling from these lists with replacement. The third method to generate fake data was  
102 imputation. The R package missForrest [23] was repurposed for data fabrication. A fake sample  
103 was generated by first creating a copy of a real sample. Then we iteratively nullified 10% of the

104 data and imputed these NAs with missForrest until every value has been imputed. See  
105 Supplemental Figure 1.

106

## 107 Machine Learning Training

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

124

## 125 Benford-Like Digit Preferences

126 Benford's Law or the first digit law has been instrumental at catching fraud in various financial  
127 situations [20,21] and in small scale clinical trials [3]. The method presented here is designed  
128 with the potential to generalize and be applied to multiple sets of data of varying types and

129 configurations (e.i. different measured variables (features) and different quantities of variables).  
130 Machine learning typically cannot handle data where the features are not consistent in number  
131 and type. Converting all measured variables to digit frequencies circumvents this problem. Digit  
132 frequencies are calculated as the number of occurrences of a single digit (0-9) divided by the  
133 total number of features. In the method described in this paper, a sample's features are all  
134 converted to digit frequencies of the first and second digit after the decimal. Thus for each  
135 sample the features are converted from ~17,000 copy number alterations to 20 digit  
136 preferences. Using this approach, whether a sample has 100 or 17,000 features it can still be  
137 trained on and classified by the same model.

## 138 Results

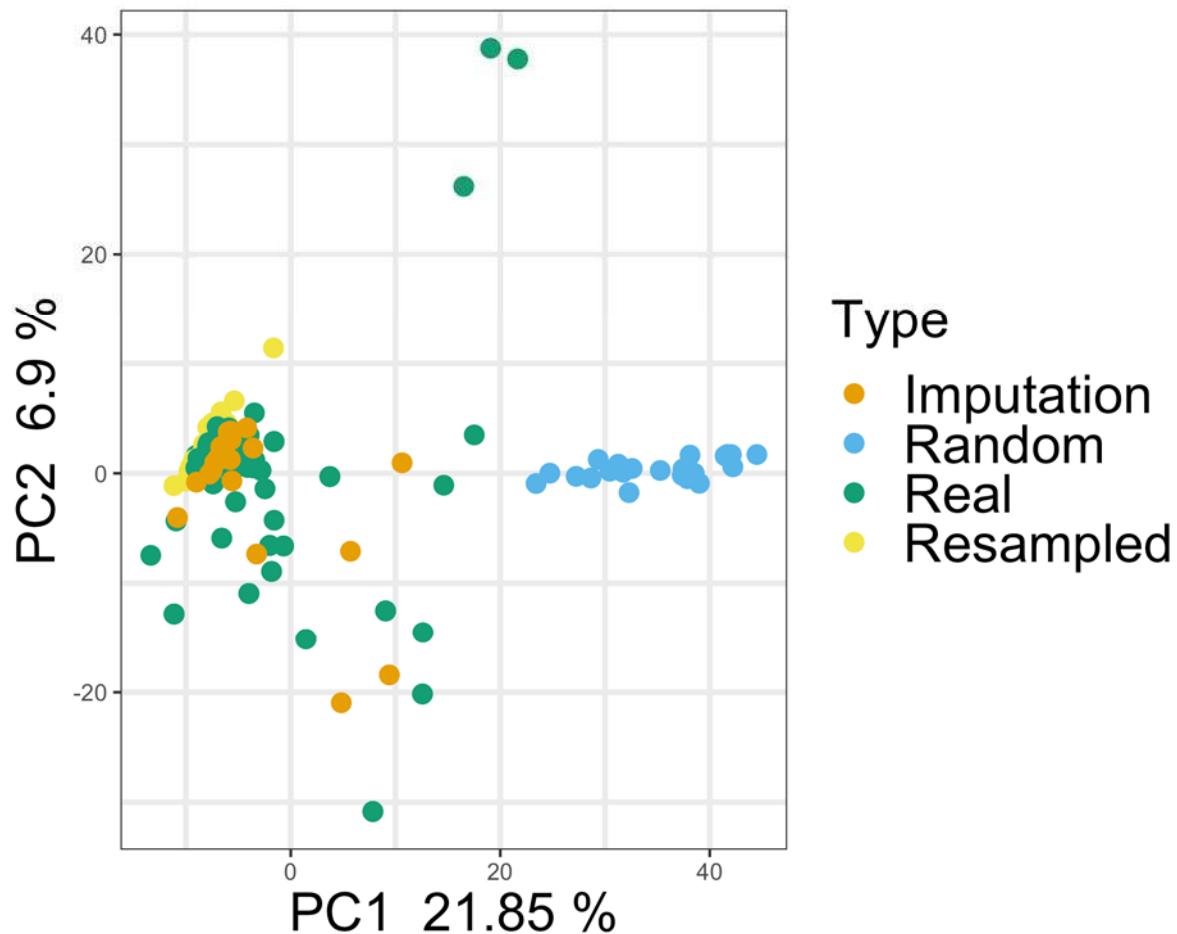
139 Our goal is to explore the ability of machine learning methods to identify fabricated data hidden  
140 within large datasets. Our results do not focus on the motivations to fabricate data, nor do they  
141 explore in depth the infinite methodological ways to do so. Our study focuses on whether  
142 machine learning can be trained to correctly identify fabricated data. Our general workflow is to  
143 take real data and mix in fabricated data. When training, the machine learning model is given  
144 access to the label (i.e. real or fabricated); the model is tested or evaluated by predicting the  
145 label of data which was held back from training (see Methods).

## 146 Fake Data

147 The real data used in this study comes from the Clinical Proteomic Tumor Analysis Consortium  
148 (CPTAC) cohort for endometrial carcinoma, specifically the copy number alteration (CNA) data.  
149 The form of this real data is a large table of floating point values. Rows represent individual  
150 tumor samples and columns represent genes; values in the cells are thus the copy number



151 quantification for a single gene in an individual tumor sample. This real data was paired with  
152 fabricated data and used as an input to machine learning classification models (see Methods).  
153 Three different methods of data fabrication were used in this study: random number generation,  
154 resampling with replacement, and imputation (Supplemental Figure 1). The three methods  
155 represent three realistic ways that an unscrupulous scientist might create novel data. Each  
156 method has benefits and disadvantages, with imputation being both the most sophisticated and  
157 also the most computationally intense and complex. As seen in Figure 1, the random data  
158 clusters far from the real data. Both the resampled and imputed data cluster tightly with the real  
159 data in a PCA plot, with the imputed data also generating a few reasonable outlier samples.



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

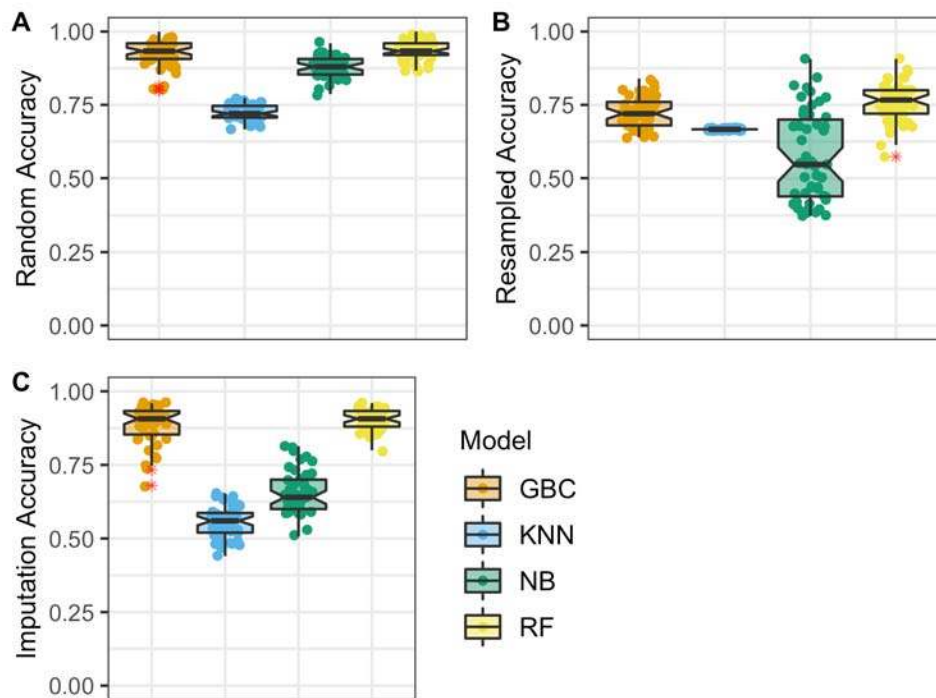
163 generation is clearly distinct from all other data. Fabricated data created via resampling or  
164 imputation appears to cluster very closely with the real data.  
165  
166 To look further into the fabricated data, we examined whether fake data preserved correlative  
167 relationships present 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 preserve these  
177 correlative relationships. When we look at the random and resampled data for these two genes,  
178 all correlation is lost (Supplemental Figure 3 C, D, E and F). Imputation, however, produces  
179 data that closely matches the original correlations, *PLEKHN1* and *HES4*  $R^2 = 0.97$ ; *DFFB* and  
180 *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. The classification task  
189 considers each sample separately, meaning that the declaration of real or fake is made only  
190 from data of a single sample. We evaluated the model on simple accuracy, whether the  
191 predicted label was correct or incorrect. To ensure that our results represent robust  
192 performance, model training and evaluation was performed 50 times; each time a completely  
193 new set of 25 fabricated samples were made (see Methods). Reported results represent the  
194 average accuracy of these 50 trials. We note that two methods, SVM and MLP, performed  
195 poorly compared to other classification methods. Testing data consisted of 2/3 real data and 1/3  
196 fake data; therefore, baseline accuracy (the accuracy achieved if the model predicting all test  
197 samples as the majority class) is 66%. Both SVM and MLP had an average accuracy at or  
198 below this baseline for classification of the simplest fabrication method (random), and were  
199 excluded from further analysis.

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



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

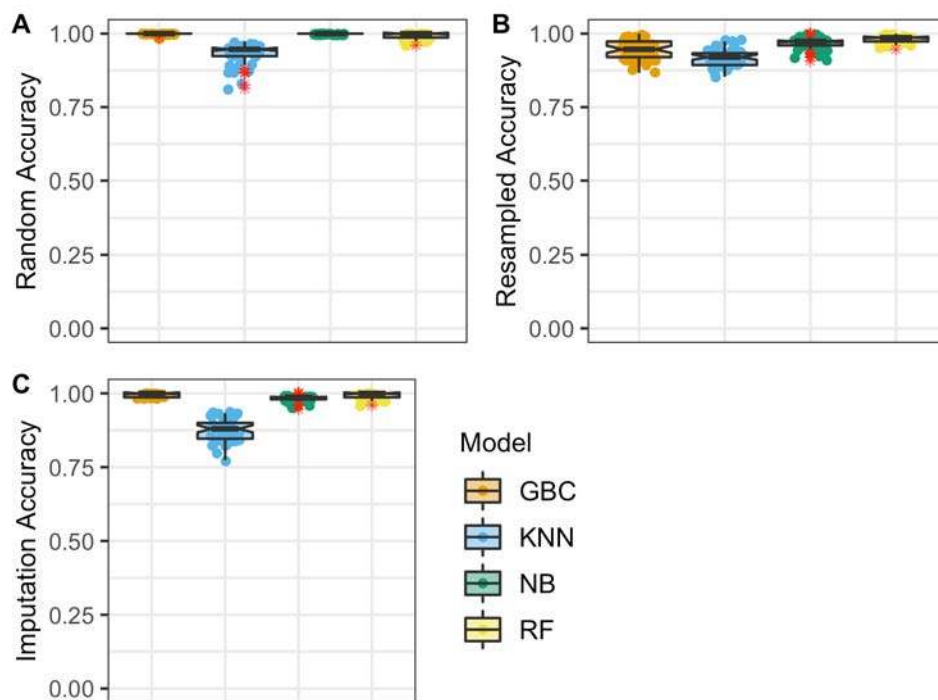
## 224 Machine learning with digit preference

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 [21], we evaluated whether this type of  
230 data representation could work for bioinformatics data as well. Therefore, all copy number data  
231 was transformed into 20 features, representing the digits 0-9 in the first and second place after  
232 the decimal of each gene expression value. While Benford's Law describes the frequency of the  
233 first digit, genomics and proteomics data are frequently normalized or scaled and so the first  
234 digit may not be as characteristic. For this reason, our method may be accurately referred to as  
235 Benford's Law inspired or Benford-like. These features were tabulated for each sample to create  
236 a new data representation and fed into the exact same machine learning training and testing  
237 routine described above. Each of these 20 new features contain decimal values ranging from  
238 0.0 to 1.0 representative of the proportional frequency that digit occurs. For example, one  
239 sample's value in the feature column for the digit 1 may contain the value 0.3. This means that  
240 in this sample's original data the digit 1 occurred in the first position after the decimal place 30%  
241 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  
252 In sharp contrast to the models built on the quantitative copy number data, machine learning  
253 models which utilized the digit frequencies were highly accurate and showed little variability over  
254 the 50 trials (Figure 3). When examining the results of the data fabricated via imputation (both  
255 the most sophisticated and most realistic), the models achieved impressively high accuracy. As  
256 an average accuracy for the 50 trials, both random forest and the gradient boosting models  
257 achieved 100% accuracy. The naïve Bayes model was highly successful with a mean  
258 classification accuracy 97%.



259  
260 **Figure 3 - Classifications accuracy using digit frequency data.** Fabricated data was mixed  
261 with real data and given to four machine learning models for classification. Data shown  
262 represents 50 trials for 50 different fabricated dataset mixes. Features in this dataset are the

263 digit frequencies for each sample. **A.** Results for data fabricated with the random method, mean  
264 classification accuracy: Random Forest 99% (+/- 1.0%), GBC 100% (+/- 0.2%), Naïve Bayes  
265 100% (+/- 0.0%), and KNN 93% (+/- 3.4%). **B.** Results for data fabricated with the resampling  
266 method, mean classification accuracy: Random Forest 98% (+/- 1.3%), GBC 94% (+/- 3.5%),  
267 Naïve Bayes 97% (+/- 2.1%), and KNN 92% (+/- 2.8%). **C.** Results for data fabricated with the  
268 imputation method, mean classification accuracy: Random Forest 100% (+/- 1.0%), GBC 100%  
269 (+/- 0.7%), Naïve Bayes 97% (+/- 1.1%), and KNN 89% (+/- 3.8%).

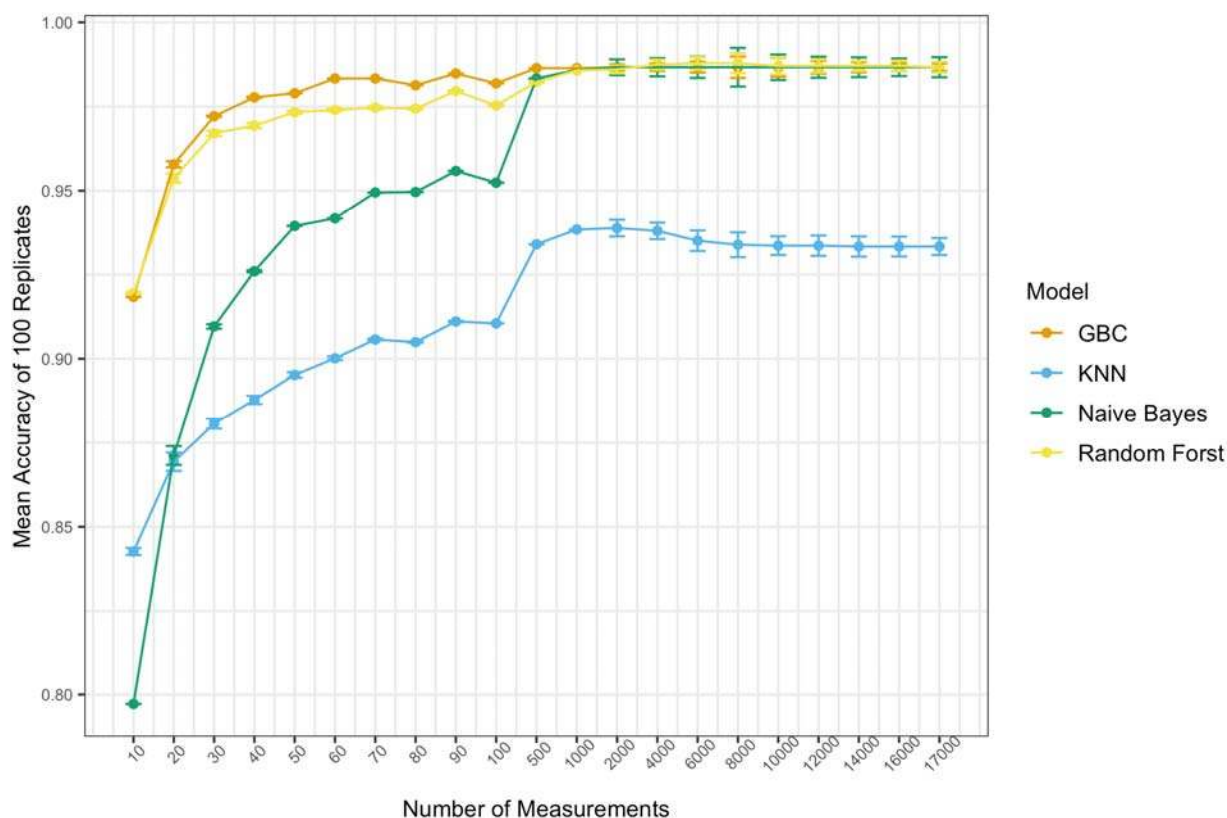
270

## 271 Machine learning with limited data

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

286

287 One hesitation for using machine learning with smaller datasets (i.e. fewer gene-features per  
288 data point) is the perceived susceptibility to large variation in performance. As noted, these  
289 downsampling experiments were performed 100 times, and error bars representing the standard  
290 error are shown in Figure 4. We note that even for the smallest datasets, performance does not  
291 noticeably vary between the 100 trials. In fact the standard error for small datasets (e.g. 20 or  
292 30 gene-features) is lower than when there were thousands. Thus we believe that the digit-  
293 frequency based models will perform well on both large-scale omics data and also on smaller  
294 'targeted' data acquisition paradigms like multiplexed PCR or MRM proteomics.



295  
296 **Figure 4 - Classifications accuracy vs number of features.** The original 17,000 CNA  
297 measurements were randomly downsampled incrementally to 10 and converted to digit-  
298 frequency training and test features for machine learning models. When 1,000+ measurements  
299 are used in the creation of digit-preference features, there appears to be little to no effect on  
300 mean accuracy. Below 1,000 Naive Bayes and KNN models begin to lose accuracy quickly.



301 GBC and Random Forest do suffer in accuracy as the number measurements used to generate  
302 features lowers but remain above 95% accurate until less than 20 measurements are included.

## 303 Discussion

304 We present here a proof of concept method for detecting fabrication in biomedical data. Just as  
305 has been previously shown in the financial sector, digit frequencies are a powerful data  
306 representation when used in combination with machine learning to predict the authenticity of  
307 data. Although the data used herein is copy number variation from a cancer cohort, we believe  
308 that the Benford-like digit frequency method can be generalized to any tabular numeric data.  
309 While multiple methods of fabrication were used, we acknowledge there are more subtle or  
310 sophisticated methods. We believe that fraud detection methods, like the models presented  
311 herein, could be refined and generalized for broad use in monitoring and oversight.

312  
313 There is an increasing call for improved oversight and review of scientific data[5,6,16,18], and  
314 various regulatory bodies or funding agencies could enforce scientific integrity through the  
315 application of these or similar methods. For example, the government bodies charged with  
316 evaluating the efficacy of new medicine could employ such techniques to screen large datasets  
317 that are submitted as evidence for the approval of new drugs. For fundamental research,  
318 publishers could mandate the submission of all data to fraud monitoring. Although journals  
319 commonly use software tools to detect plagiarism in the written text, a generalized  
320 computational tool focused on data could make data fraud detection equally simple.

## 321 Acknowledgments

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

324

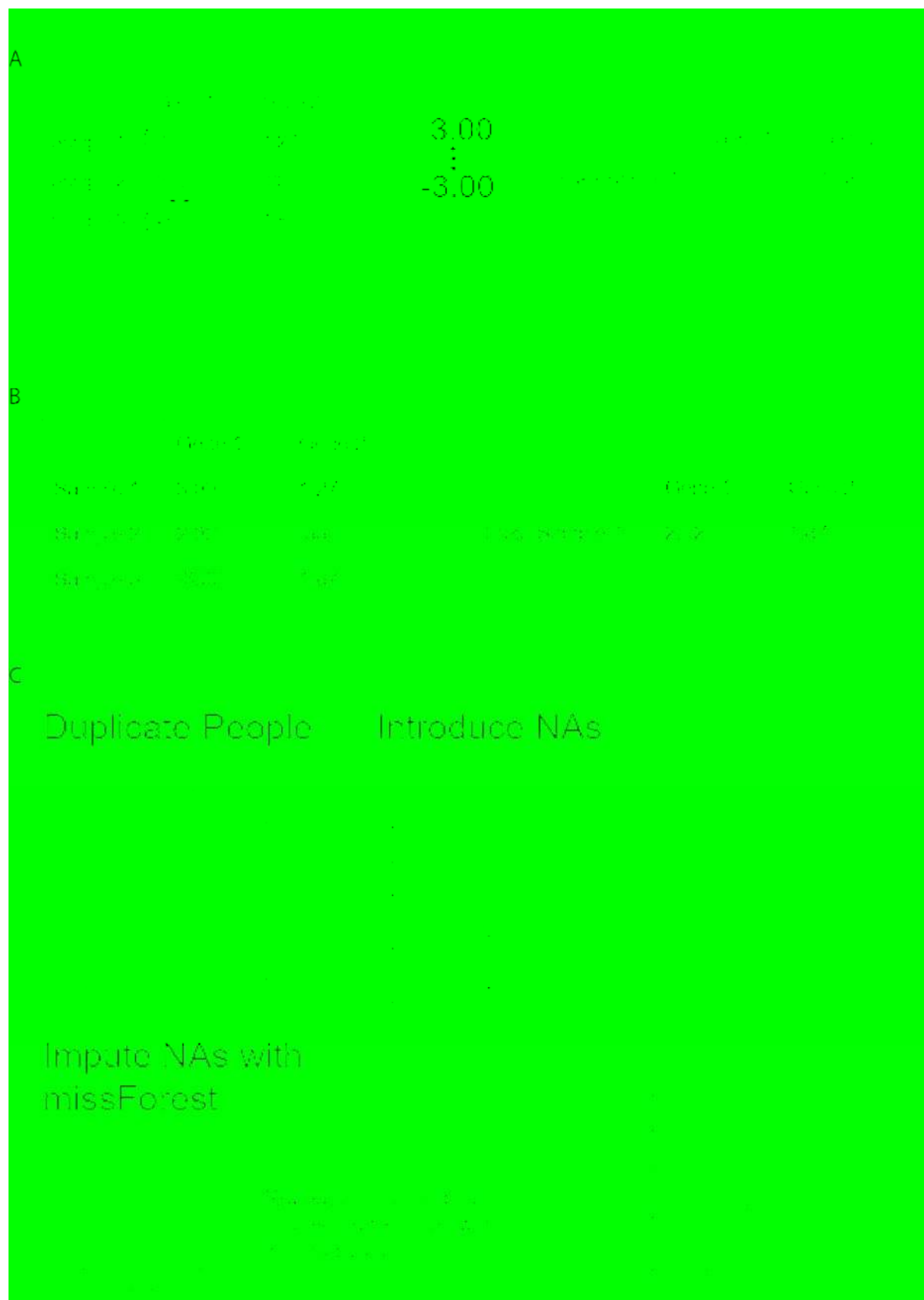
## 325 References

- 326 1. Burton F. The acquired immunodeficiency syndrome and mosquitoes. *Med J Aust.*  
327 1989;151: 539–540.
- 328 2. Kupferschmidt K. Tide of lies. *Science.* 2018;361: 636–641.
- 329 3. Al-Marzouki S, Evans S, Marshall T, Roberts I. Are these data real? Statistical methods for  
330 the detection of data fabrication in clinical trials. *BMJ.* 2005;331: 267–270.
- 331 4. Fanelli D. How many scientists fabricate and falsify research? A systematic review and  
332 meta-analysis of survey data. *PLoS One.* 2009;4: e5738.
- 333 5. George SL, Buyse M. Data fraud in clinical trials. *Clin Investig .* 2015;5: 161–173.
- 334 6. Yu L, Miao M, Liu W, Zhang B, Zhang P. Scientific Misconduct and Associated Factors: A  
335 Survey of Researchers in Three Chinese Tertiary Hospitals. *Account Res.* 2020.  
336 doi:10.1080/08989621.2020.1809386
- 337 7. Blum A, Wang P, Zenklusen JC. SnapShot: TCGA-Analyzed Tumors. *Cell.* 2018;173: 530.
- 338 8. TEDDY Study Group. The Environmental Determinants of Diabetes in the Young (TEDDY)  
339 study: study design. *Pediatr Diabetes.* 2007;8: 286–298.
- 340 9. Orwoll E, Blank JB, Barrett-Connor E, Cauley J, Cummings S, Ensrud K, et al. Design and

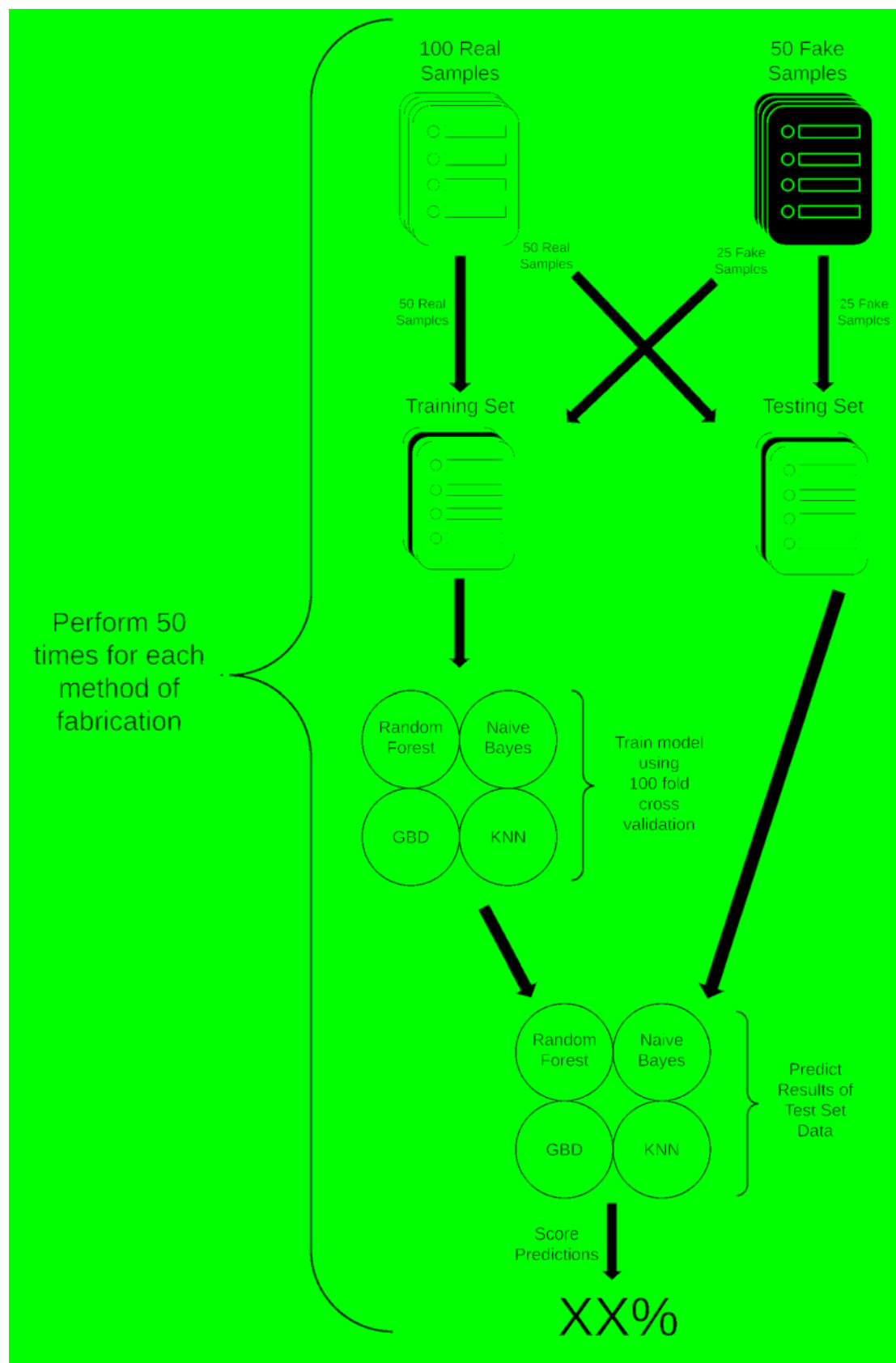
- 341 baseline characteristics of the osteoporotic fractures in men (MrOS) study--a large  
342 observational study of the determinants of fracture in older men. *Contemp Clin Trials*.  
343 2005;26: 569–585.
- 344 10. Bycroft C, Freeman C, Petkova D, Band G, Elliott LT, Sharp K, et al. The UK Biobank  
345 resource with deep phenotyping and genomic data. *Nature*. 2018;562: 203–209.
- 346 11. Barretina J, Caponigro G, Stransky N, Venkatesan K, Margolin AA, Kim S, et al. The  
347 Cancer Cell Line Encyclopedia enables predictive modelling of anticancer drug sensitivity.  
348 *Nature*. 2012;483: 603–607.
- 349 12. Subramanian A, Narayan R, Corsello SM, Peck DD, Natoli TE, Lu X, et al. A Next  
350 Generation Connectivity Map: L1000 Platform and the First 1,000,000 Profiles. *Cell*.  
351 2017;171: 1437–1452.e17.
- 352 13. Caswell J, Gans JD, Generous N, Hudson CM, Merkley E, Johnson C, et al. Defending Our  
353 Public Biological Databases as a Global Critical Infrastructure. *Front Bioeng Biotechnol*.  
354 2019;7: 58.
- 355 14. Bik EM, Casadevall A, Fang FC. The Prevalence of Inappropriate Image Duplication in  
356 Biomedical Research Publications. *MBio*. 2016;7. doi:10.1128/mBio.00809-16
- 357 15. Knepper D, Fenske C, Nadolny P, Bedding A, Gribkova E, Polzer J, et al. Detecting Data  
358 Quality Issues in Clinical Trials: Current Practices and Recommendations. *Ther Innov  
359 Regul Sci*. 2016;50: 15–21.
- 360 16. Baigent C, Harrell FE, Buyse M, Emberson JR, Altman DG. Ensuring trial validity by data  
361 quality assurance and diversification of monitoring methods. *Clin Trials*. 2008;5: 49–55.
- 362 17. Morrison BW, Cochran CJ, White JG, Harley J, Kleppinger CF, Liu A, et al. Monitoring the

- 363 quality of conduct of clinical trials: a survey of current practices. Clin Trials. 2011;8: 342–  
364 349.
- 365 18. Calis KA, Archdeacon P, Bain R, DeMets D, Donohue M, Elzarrad MK, et al.  
366 Recommendations for data monitoring committees from the Clinical Trials Transformation  
367 Initiative. Clin Trials. 2017;14: 342–348.
- 368 19. Benford F, Langmuir I. The Law of Anomalous Numbers. American Philosophical Society;  
369 1938.
- 370 20. Cerioli A, Barabesi L, Cerasa A, Menegatti M, Perrotta D. Newcomb-Benford law and the  
371 detection of frauds in international trade. Proc Natl Acad Sci U S A. 2019;116: 106–115.
- 372 21. Badal-Valero E, Alvarez-Jareño JA, Pavía JM. Combining Benford's Law and machine  
373 learning to detect money laundering. An actual Spanish court case. Forensic Sci Int.  
374 2018;282: 24–34.
- 375 22. George SL. Research misconduct and data fraud in clinical trials: prevalence and causal  
376 factors. Int J Clin Oncol. 2016;21: 15–21.
- 377 23. Stekhoven DJ, Bühlmann P. MissForest--non-parametric missing value imputation for  
378 mixed-type data. Bioinformatics. 2012;28: 112–118.

## 379 Supplemental Figures:

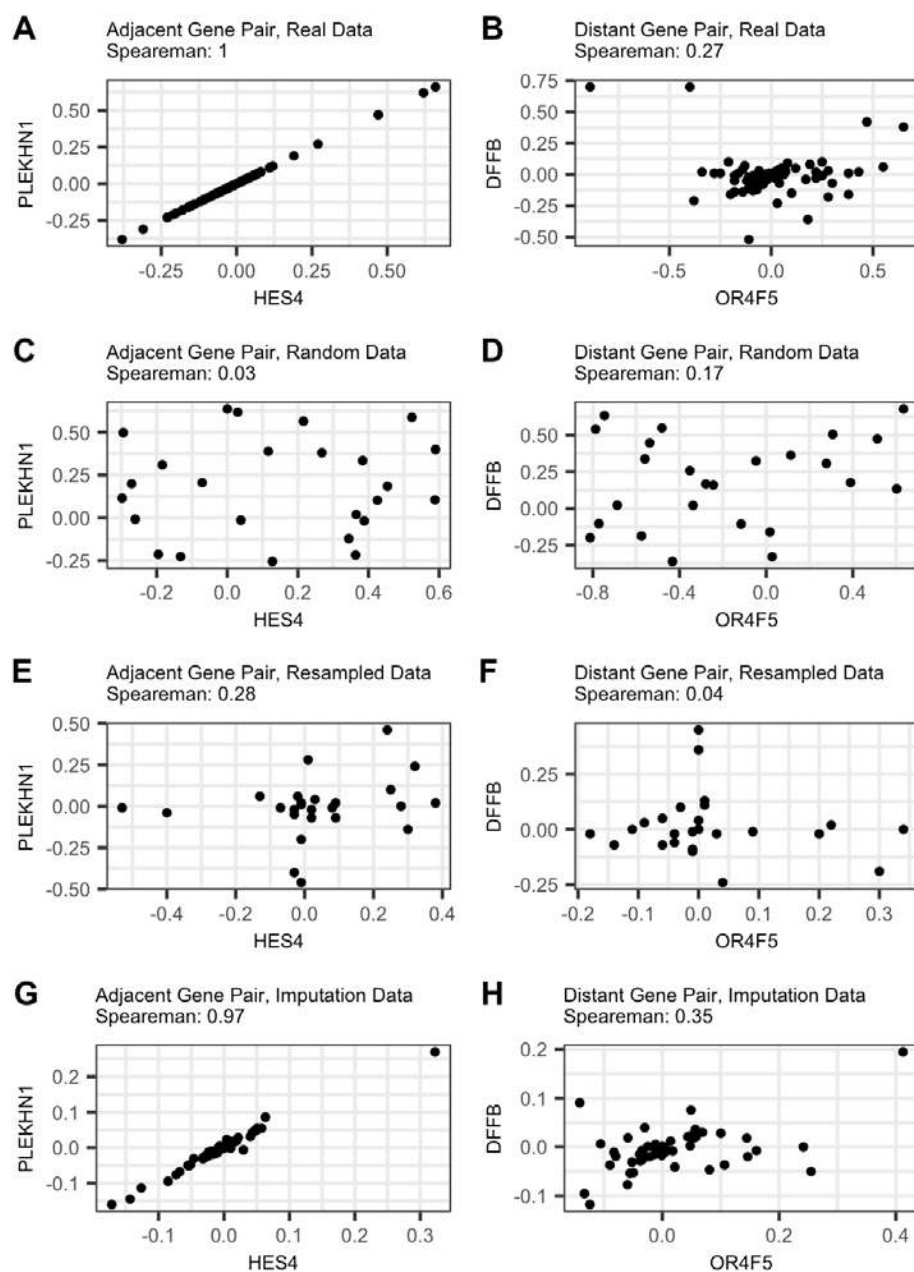


381 **Supplemental Figure 1** - Methods of Data fabrication. (A) The random method of data  
382 fabrication identifies the range of observation for a specific locus and then randomly chooses a  
383 number in that range. (B) The resampling method chooses values present in the original data.  
384 (C) The imputation method iteratively nullifies and then imputes data points from a real sample.





386 **Supplemental Figure 2** - Training and testing overview. After creating 50 fake samples using  
387 any one of the three methods of fabrication, the 100 real samples and 50 fake samples were  
388 randomly split into a train and test set of equal size and proportions (50 real and 25 fake in each  
389 set). The training sets were then used to train various machine learning models using 10-fold  
390 cross validation. Next, trained models were used to make predictions on the testing data.  
391 Predictions were then scored with total accuracy.



392

393 **Supplemental Figure 3** - Data relationships in fabricated data. The correlation between pairs of  
394 genes is evaluated to determine whether fabrication methods can replicate inter-gene patterns.  
395 Plots on the left hand side (A,C,E, and G) display data from two correlated genes *PLEKHN1*  
396 and *HES4*, adjacent genes found on 1p36. Plots on the right hand side (B,D,F, and H) display  
397 genes *DFFB* and *OR4F5* gene with marginal Spearman correlation in the real data (0.27). The  
398 plots reveal that random and resample data have little to no correlation between related genes.  
399 Imputation produces data with correlation values that are similar to the original data (0.97 and  
400 0.35, respectively).