

1 **Genetic risk score for ovarian cancer based on**
2 **chromosomal-scale length variation.**

3 **Authors:** Chris Toh¹ and James P. Brody^{1*}

4 **Affiliations:**

5 ¹Department of Biomedical Engineering, University of California, Irvine.

6 *Correspondence to: jpbrody@uci.edu.

7 **Abstract:**

8 **Introduction.** Twin studies indicate that a substantial fraction of ovarian cancers should be
9 predictable from genetic testing. Genetic risk scores can stratify women into different classes
10 of risk. Higher risk women can be treated or screened for ovarian cancer, which should
11 reduce overall death rates due to ovarian cancer. However, current ovarian cancer genetic
12 risk scores, based on SNPs, do not work that well. We developed a genetic risk score based
13 on structural variation, quantified by variations in the length of chromosomes.

14 **Methods.** We evaluated this genetic risk score using data collected by The Cancer Genome
15 Atlas. From this dataset, we synthesized a dataset of 414 women who had ovarian serous
16 carcinoma and 4225 women who had no form of ovarian cancer. We characterized each
17 woman by 22 numbers, representing the length of each chromosome in their germ line DNA.
18 We used a gradient boosting machine, a machine learning algorithm, to build a classifier that
19 can predict whether a woman had been diagnosed with ovarian cancer in this dataset.

20 **Results.** The genetic risk score based on chromosomal-scale length variation could stratify
21 women such that the highest 20% had a 160x risk (95% confidence interval 50x-450x)
22 compared to the lowest 20%. The genetic risk score we developed had an area under the
23 curve of the receiver operating characteristic curve of 0.88 (estimated 95% confidence
24 interval 0.86-0.91).

25 **Conclusion.** A genetic risk score based on chromosomal-scale length variation of germ line
26 DNA provides an effective means of predicting whether or not a woman will develop ovarian
27 cancer.
28

29

30 **Introduction:**

31 Ovarian cancer kills about 150,000 women per year worldwide[1]. The most common
32 form of ovarian cancer, ovarian serous carcinoma is often diagnosed late (stage III (51%) or IV
33 (29%)) and has a relatively bleak 5-year survival rate [2]. If women with an elevated risk of
34 developing ovarian cancers could be identified, interventions could be taken that would reduce
35 the number of women who die from ovarian cancer. These interventions include prophylactic
36 oophorectomies, which would completely avoid ovarian cancer, and more targeted screening,
37 which could identify ovarian cancers in earlier stages, where surgery is an effective cure[3–6].
38 These interventions could both increase 5-year survival times and reduce the overall number of
39 deaths due to ovarian cancer.

40 A substantial fraction of ovarian cancers should be predictable by genetic testing. The
41 heritability of ovarian cancer has been measured at about 40% (95% confidence interval 23%-
42 55%) by the Nordic Twin Study[7]. The maximum discriminative accuracy of a genetic risk test
43 is a function of both the heritability and the prevalence of the disease [8,9]. Based on the
44 measured heritability (about 40%) and prevalence (about 0.1%) of ovarian cancer, the maximum
45 accuracy, measured by the area under the receiver operating characteristic curve (AUC), should
46 be greater than 0.95, where 1.0 indicates a perfect test. Current genetic risk scores do not
47 approach that level of accuracy.

48 Most current genetic risk scores are derived from single nucleotide polymorphisms
49 (SNPs) identified by genome wide association studies[10–15]. These tests, called polygenic risk
50 scores, construct a score based on a linear combination of the value of a collection of SNPs. This
51 strategy has been moderately successful with ovarian cancer. One study followed this strategy to

52 construct a polygenic risk score where women who scored in the top 20% had a 3.4-fold
53 increased risk compared to women who scored in the bottom 20% [16].

54 We developed an alternative strategy to compute genetic risk scores. Our strategy is
55 based on structural variation rather than SNPs and uses machine learning algorithms, which
56 include non-linear effects, rather than linear combinations.

57 **Methods:**

58 We tested this strategy with data from the Cancer Genome Atlas (TCGA) project. TCGA
59 was a project sponsored by the National Cancer Institute to characterize the molecular
60 differences in 33 different human cancers [17–19]. The project collected samples from about
61 11,000 different patients, all of whom were being treated for one of 33 different types of tumors.
62 The samples collected usually included tissue samples of the tumor, tissue samples of normal
63 tissue adjacent to the tumor and normal blood samples. (Normal blood samples were not
64 available from patients diagnosed with leukemias.)

65 Most of the patient normal blood samples were processed to extract and characterize
66 germline DNA. All germline DNA samples were processed by a single laboratory, the
67 Biospecimen Core Resource at Nationwide Children’s Hospital. Single nucleotide
68 polymorphisms (SNPs) were measured from the patient samples with an Affymetrix SNP 6.0
69 array. This SNP data was then processed (by the TCGA project) through a bioinformatics
70 pipeline [20], which included the packages Birdsuite [21] and DNACopy [22]. The result of this
71 pipeline is, for each sample, a listing of a chromosomal region (characterized by the chromosome
72 number, a starting location, and an ending location) and the associated value given as the
73 “segmented mean value.” The segmented mean value is defined as the logarithm, base 2 of one-

74 half the copy number. A normal diploid region with two copies will have a segmented mean
75 value of zero.

76 NCI has provided most of the TCGA data on the Genomic Data Commons [23]. The
77 copy number variation is called the masked copy number variation on the Genomic Data
78 Commons. The masking process removes “Y chromosome and probe sets that were previously
79 indicated to have frequent germline copy-number variation.” [20].

80 This research uses de-identified coded datasets produced by TCGA. Therefore it is not
81 considered human subjects research.

82 We accessed the TCGA data through Google’s BigQuery, a cloud-based database. This
83 resource is hosted and maintained by the Institute of Systems Biology [24]. We used the copy
84 number segment (masked) table extracted from the Genomic Data Commons in February 2017.
85 We also used information from the Biospecimen (extracted April 2017) and Clinical (extracted
86 June 2018) tables. The copy number table contained all the information for the chromosome
87 scale length variation data. The Biospecimen table was used to identify which samples were
88 from normal blood (representing germ line DNA). The Clinical table provided information on
89 the individual patient’s gender, race, and ovarian cancer status. Information in the different
90 tables was tied together by the sample barcode parameter.

91 We used the statistical computer language R to query the BigQuery database, collect the
92 data and manipulate it into different forms. We took extensive care to avoid typical problems
93 that lead to falsely high AUCs in machine learning. For instance, we ensured that no data leakage
94 occurred, which can lead deceptively high AUCs when copies of a sample appear in both the
95 training and test sets.

96 We used the H2O machine learning package in R to create machine learning models.
97 H2O takes care of setting many of the proper default values, depending on whether the goal of
98 the model is classification or regression. For the gradient boosting machine (GBM) models,
99 H2O performs preprocessing, randomization, encoding categorical variables, and other data
100 processing steps appropriate for the chosen model.

101 H2O has an automated machine learning algorithm, named AutoML[25]. Given a
102 spreadsheet like- dataset, AutoML will run through four different machine learning algorithms
103 and evaluate which provides the best models for the given problem. For each of the machine
104 learning algorithms, it will evaluate several different hyperparameters. The process is limited by
105 the amount of time devoted to it. After the allotted time, AutoML reports a scoreboard ranking
106 the best algorithms. For the gradient boosting machine algorithm, we started with the default
107 H2O settings. These default settings build trees to a maximum depth of five trees with a sample
108 rate of 1 [26]. For the results reported in Table 2, we used an allotted time of one hour. In tests,
109 we found that the results do not change substantially with times up to 10 hours.

110 We used 5-fold cross validation with the GBM algorithm to produce Table 3 and Figure
111 2. Cross validation uses repeated model runs with non-overlapping data. This approach allows
112 one to use of all samples in the limited dataset. For Table 3 and Figure 2, we estimated 95%
113 confidence intervals for the odds ratios following the method described in [27].

114 Figure 3 was produced with a single model run by splitting the dataset into a training set
115 containing 80% of the data and a test set containing 20% of the data.

116 **Results:**

117 Using the TCGA dataset, we identified a measure that we call *chromosome-scale length*
118 *variation*. Taken together, structural variations like insertions, deletions, translocations and copy

119 number variations slightly alter the overall length of an individual's chromosome. Thus, the
120 lengths of the set of chromosomes can be used to characterize a person. A histogram showing the
121 distribution of relative chromosome lengths taken from germ line DNA samples in the TCGA
122 dataset is shown in Figure 1. By convention, these lengths are reported in units of log base 2. A
123 value of "0" represents the consensus, average, chromosome length.

124 **Figure 1.** This figure shows a histogram of chromosome scale
125 length variation for most of chromosomes 1,6,13, and 17. For most
126 patients in the TCGA dataset, a normal blood sample was taken,
127 genomic DNA was extracted from that sample and analyzed with
128 an Affymetrix SNP 6.0 array. The data from this array was
129 processed by the TCGA project through a bioinformatic pipeline
130 that resulted in a segment mean value, which is a number equal to
131 the log base two of one half the copy number value. This
132 histogram indicates that most people have a nominal value of 0,
133 indicating exactly two copies of the diploid chromosome. A value
134 of 0.02 would indicate the person has on average 2.028 copies of
135 the chromosome, or about 1.4% longer than the average length of
136 the chromosome.

137
138 From the TCGA dataset, we synthesized a case-control study to test whether
139 chromosome-scale length variation data can construct a genetic risk score. We identified 4225
140 women who had not been diagnosed with any form of ovarian cancer and 414 women who had
141 been diagnosed with ovarian serous carcinoma. Statistical descriptions of the two populations are
142 shown in Table 1.

143
144 **Table 1.** From the TCGA dataset, we constructed two groups, both
145 solely composed of women. The first group, containing 414
146 women, all had been diagnosed with ovarian serous carcinoma.
147 None of the second group, with 4225 women, had been diagnosed
148 with any form of ovarian cancer. This table compares the two
149 populations.

Diagnosed with Ovarian

Not diagnosed with Ovarian

	Serous Carcinoma	cancer
Total	414	4225
Mean age	58.3 years	59.7 years
% Black	25/414 = 6%	492/4225 = 12 %
% White	352/414= 85%	3064/4225= 73%
% Asian	14/414 = 3%	259/4225 6%

150
151 Next, we evaluated the effectiveness of several different machine learning algorithms.
152 We measured how well these algorithms could classify a woman, based solely on the set of 23
153 chromosome-scale length variation measurements, into either the class with ovarian cancer or
154 without. The measurement of success we used was the area under the curve (AUC) of the
155 receiver operating characteristic curve. The results of these measurements are shown in Table 2.

156 Table 2. This table lists five different machine learning algorithms
157 we evaluated for predicting ovarian cancer from chromosome-
158 scale length variation data using the H2O package in R. The
159 algorithms are ranked by the best AUC it achieved using 5-fold
160 cross validation.

Algorithm	AUC
Gradient Boosting Machine	0.88
Distributed Random Forest	0.87
Extremely Randomized Trees	0.86
Deep learning	0.82
Generalized Linear Model	0.68

161
162 Based on the results in Table 2, we used the Gradient Boosting Machine algorithm
163 throughout the rest of this manuscript. In the next step, we sought to classify the 4669 women in
164 the dataset. We used a k -fold cross validation procedure, with $k=5$. The dataset was randomly
165 partitioned into five equal groups. The first group was held out (to be the test set), while the
166 other four groups were used to train a model to distinguish the two classes (women with ovarian
167 cancer and women without ovarian cancer). The trained model assigned a numerical score to

168 each of the women in the first group (test set) quantifying how likely that woman was a member
169 of the ovarian cancer class. The process was repeated 5 times, with a different group held out
170 each time. The result is a numerical score for each of the 4669 women.

171 The predictions were compared to the known ovarian cancer status of each of the 4669
172 women. First, all 4669 women were ranked by their score, representing the likelihood that they
173 were from the ovarian cancer class. By comparing this ranking with their known ovarian cancer
174 status, we can evaluate how well the model classified the women.

175 The comparison is presented in two different forms. Table 3 provides a tabular form of
176 relative risk for the population segmented into five different groups. Figure 2 shows similar
177 information in graphical form, where the population is segmented into 50 groups.

178 Finally, we took the dataset of 4669 women and split it into a training set (80%) and a
179 test set (20%). Using H2O, we trained a Gradient Boosting Machine model to predict whether a
180 woman was in the group with ovarian cancer, or not. The results are presented in Figure 3, which
181 shows a classic receiver operating characteristic curve of the model's predictions.

182 .

183 Table 3. Using 5-fold cross validation, each woman in the dataset
184 received a score from the model built to predict ovarian cancer.
185 The women were ranked by score from lowest to highest and then
186 partitioned into five quintiles. This table presents the number of
187 women with and without ovarian cancer in each quintile along with
188 the odds ratio (relative to the entire group) and the 95% confidence
189 interval for the odds ratio.

Quintile	Number of women without ovarian cancer	Number of women with ovarian cancer	Total number of women	Odds ratio	95% confidence interval
1	925	3	928	0.03	0.01--0.09
2	925	3	928	0.03	0.01--0.09
3	901	27	928	0.30	0.21--0.45
4	842	86	928	1.04	0.82--1.33
5	632	295	927	4.76	4.01--5.65

190

191

192

193

194

195

196

197

198

Figure 2. This figure shows that women ranked higher by the predictive model are significantly more likely to have ovarian cancer. The predictive model ranked all 4669 women in the dataset based on their likelihood of having ovarian cancer, based solely on germ line DNA data. This ranking was then split into 50 equal partitions, each with about 93 women. This plot shows the odds ratio (relative to 414 ovarian cases out of 4669 total) of each of the 50 equal partitions along with the 95% confidence intervals.

199

200

201

202

Figure 3. This figure presents a receiver operating characteristic curve of the model's predictions. The area under the curve for this model was 0.88.

203 Discussion:

204

205

206

207

208

The results presented here compare favorably to other genetic risk scores for ovarian cancer. For instance, a previous study found that a polygenic risk score in the top 20% conferred a 3.4-fold risk increase compared to women in the bottom 20% [16]. As seen in Table 3, the top 20% in our results had an increase of over 100-fold risk over women who scored in the bottom 20%.

209

210

211

212

213

Table 2 quantifies different algorithms applied to this problem. These results are illustrative, but not conclusive. Tuning machine learning models is an art, and it might be possible, for instance, to tune a deep learning network to obtain superior results. In similar work on TCGA colon cancer data, we found that a pairwise neuron network algorithm performs equal to a gradient boosting machine[28]. The gradient boosting machine generally runs faster and is

214 easier to tune. Others have evaluated different machine learning algorithms for different
215 bioinformatic problems and found that no one algorithm is superior[29]. They also found that a
216 gradient boosting machine algorithm does perform well on many different types of datasets,
217 consistent with our findings.

218 A disadvantage of this approach, compared to more conventional SNP-based genetic risk
219 scores, is that the results are difficult to understand and extract biological meaning. The
220 Gradient Boosting Machine computational model is complex, consisting of dozens of decisions
221 trees. Furthermore, the data that is used to traverse the decision tree is also complex. The data
222 consists of chromosome scale length variation, which is the result of many different insertions,
223 deletions, translocations, and other structural changes. Polygenic risk scores based on SNPs are
224 easy to interpret. One can identify how much each SNP contributes to the score and one can
225 locate this SNP in the genome and understand the function of nearby genes that might change.
226 Although this approach is lacking in explanatory power, its ultimate goal is predictive power.

227 We considered whether the results were due to two common problems faced by GWAS
228 studies: batch effects or population stratification. We found it unlikely that our model is
229 identifying batch effects rather than real effects. First, all samples were collected from the same
230 tissue, blood. This eliminates one common source of batch effects, since the DNA extraction
231 process is the same for each sample. Second, all samples were processed by the same laboratory,
232 the Nationwide Children's Hospital Biospecimen Core Resource, with the same type of
233 instrument. This laboratory followed the same protocol throughout their processing phase.
234 Finally, we looked up the batch history of each sample. The 424 ovarian cancer samples were
235 processed in 15 separate batches. The non-ovarian samples were processed in several hundred
236 different batches. For these reasons, we do not believe the results are due to batch effects.

237 Population stratification occurs in case/control studies when the cases and controls
238 contain substantially different proportions of genetically discernable subclasses. Most TCGA
239 samples were collected in the United States from a racially diverse group. For instance, over half
240 the ovarian cancer samples were collected at five locations in the United States: Memorial Sloan
241 Kettering, Washington University, University of Pittsburgh, Duke, and Mayo Clinic- Rochester.
242 Table 1 lists demographic information about the two populations. Although the table does
243 indicate slightly different proportions, by race, in the case and control groups, it does not seem to
244 be different enough to account for the AUC observed.

245 This study has several weaknesses. First, the control population in this analysis is not
246 randomly drawn from the general population, but instead consists of women who were part of
247 the study because they were diagnosed with another form of cancer. Second, the results rely on a
248 single dataset. The general applicability of this method would be better established if we were
249 able to show that a model trained on one dataset would perform well on a second dataset that was
250 collected independently. Demonstrating that a model is transferrable is a longer-term goal of
251 ours.

252 Future work could refine this method to improve the predictive ability of this method.
253 The AUC might be improved through several strategies, including feature engineering, for
254 instance using sub-chromosomes rather than complete chromosomes, data augmentation
255 strategies, and the inclusion of SNP data. Further work can also establish how robust the model
256 is: can a model trained with the TCGA data be successfully applied to a person not in the TCGA
257 dataset.

258 **Conclusion:**

259 A genetic risk score based on chromosomal-scale length variation of germ line DNA
260 provides an effective means of predicting whether or not a woman will develop ovarian cancer.
261 Several avenues are open to further improve the AUC of this genetic risk score test.

262 **Competing Interests:**

263 None of the authors have any competing interests.

264 **Acknowledgements:**

265 The results published here are in whole or part based upon data generated by the TCGA
266 Research Network: <http://cancergenome.nih.gov/>.

267

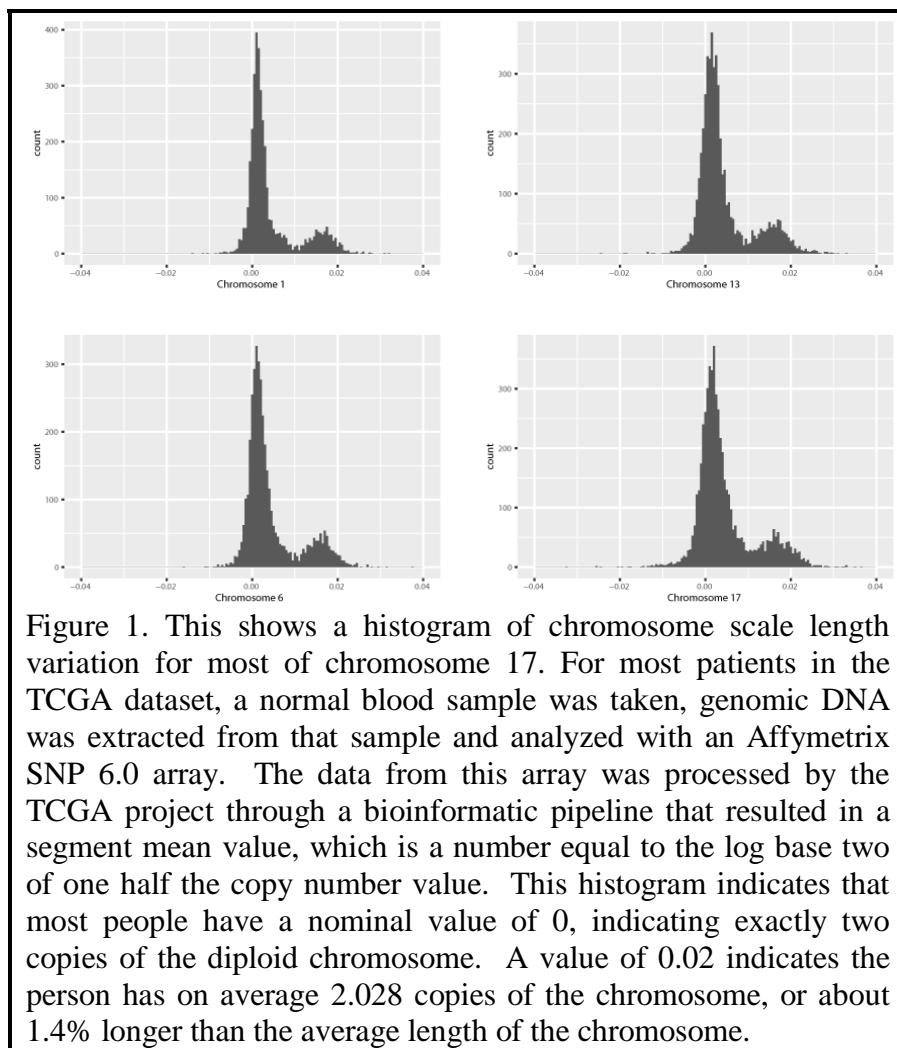
268 **References:**

- 269 1. Bray F, Ferlay J, Soerjomataram I, Siegel RL, Torre LA, Jemal A. Global cancer statistics 2018:
270 GLOBOCAN estimates of incidence and mortality worldwide for 36 cancers in 185 countries. *CA:
271 A Cancer Journal for Clinicians*. 2018;68: 394–424. doi:10.3322/caac.21492
- 272 2. Torre LA, Trabert B, DeSantis CE, Miller KD, Samimi G, Runowicz CD, et al. Ovarian cancer
273 statistics, 2018. *CA: A Cancer Journal for Clinicians*. 2018;68: 284–296. doi:10.3322/caac.21456
- 274 3. Bast RC. Status of Tumor Markers in Ovarian Cancer Screening. *Journal of Clinical Oncology*.
275 2003;21: 200s–220s. doi:10.1200/JCO.2003.01.068
- 276 4. Andrews L, Mutch DG. Hereditary Ovarian Cancer and Risk Reduction. *Best Practice & Research
277 Clinical Obstetrics & Gynaecology*. 2017;41: 31–48. doi:10.1016/j.bpobgyn.2016.10.017
- 278 5. Grossman DC, Curry SJ, Owens DK, Barry MJ, Davidson KW, Doubeni CA, et al. Screening for
279 ovarian cancer US preventive services task force recommendation statement. *JAMA - Journal of
280 the American Medical Association*. 2018. doi:10.1001/jama.2017.21926
- 281 6. Trimbos JB. Surgical treatment of early-stage ovarian cancer. *Best Practice and Research: Clinical
282 Obstetrics and Gynaecology*. 2017. doi:10.1016/j.bpobgyn.2016.10.001
- 283 7. Mucci LA, Hjelmborg JB, Harris JR, Czene K, Havelick DJ, Scheike T, et al. Familial Risk and
284 Heritability of Cancer Among Twins in Nordic Countries. *JAMA*. 2016;315: 68–76.
285 doi:10.1001/jama.2015.17703
- 286 8. Janssens ACJW, Aulchenko YS, Elefante S, Borsboom GJJM, Steyerberg EW, van Duijn CM.
287 Predictive testing for complex diseases using multiple genes: Fact or fiction? *Genetics in
288 Medicine*. 2006;8: 395–400. doi:10.1097/01.gim.0000229689.18263.f4
- 289 9. Janssens ACJW, van Duijn CM. Genome-based prediction of common diseases: advances and
290 prospects. *Human Molecular Genetics*. 2008;17: R166–R173. doi:10.1093/hmg/ddn250
- 291 10. Torkamani A, Wineinger NE, Topol EJ. The personal and clinical utility of polygenic risk scores.
292 *Nature Reviews Genetics*. 2018. doi:10.1038/s41576-018-0018-x
- 293 11. Lambert SA, Abraham G, Inouye M. Towards clinical utility of polygenic risk scores. *Human
294 Molecular Genetics*. 2019. doi:10.1093/hmg/ddz187
- 295 12. Khera A v., Chaffin M, Aragam KG, Haas ME, Roselli C, Choi SH, et al. Genome-wide polygenic
296 scores for common diseases identify individuals with risk equivalent to monogenic mutations.
297 *Nature Genetics*. 2018;50: 1219–1224. doi:10.1038/s41588-018-0183-z
- 298 13. Pharoah PDP, Tsai Y-Y, Ramus SJ, Phelan CM, Goode EL, Lawrenson K, et al. GWAS meta-analysis
299 and replication identifies three new susceptibility loci for ovarian cancer. *Nature Genetics*.
300 2013;45: 362–370. doi:10.1038/ng.2564
- 301 14. Kuchenbaecker KB, Ramus SJ, Tyrer J, Lee A, Shen HC, Beesley J, et al. Identification of six new
302 susceptibility loci for invasive epithelial ovarian cancer. *Nature Genetics*. 2015;47: 164–171.
303 doi:10.1038/ng.3185

- 304 15. Lewis CM, Vassos E. Polygenic risk scores: from research tools to clinical instruments. *Genome*
305 *Medicine*. 2020;12: 44. doi:10.1186/s13073-020-00742-5
- 306 16. Goode EL, Chenevix-Trench G, Song H, Ramus SJ, Notaridou M, Lawrenson K, et al. A genome-
307 wide association study identifies susceptibility loci for ovarian cancer at 2q31 and 8q24. *Nature*
308 *Genetics*. 2010. doi:10.1038/ng.668
- 309 17. Weinstein JN, Collisson EA, Mills GB, Shaw KRM, Ozenberger BA, Ellrott K, et al. The Cancer
310 Genome Atlas Pan-Cancer analysis project. *Nature Genetics*. 2013;45: 1113–1120.
311 doi:10.1038/ng.2764
- 312 18. Bell D, Berchuck A, Birrer M, Chien J, Cramer DW, Dao F, et al. Integrated genomic analyses of
313 ovarian carcinoma. *Nature*. 2011;474: 609–615. doi:10.1038/nature10166
- 314 19. Hutter C, Zenklusen JC. The Cancer Genome Atlas: Creating Lasting Value beyond Its Data. *Cell*.
315 2018;173: 283–285. doi:10.1016/j.cell.2018.03.042
- 316 20. Copy Number Variation Analysis Pipeline. [cited 18 Jan 2018]. Available:
317 https://docs.gdc.cancer.gov/Data/Bioinformatics_Pipelines/CNV_Pipeline/
- 318 21. Korn JM, Kuruvilla FG, McCarroll SA, Wysoker A, Nemesh J, Cawley S, et al. Integrated genotype
319 calling and association analysis of SNPs, common copy number polymorphisms and rare CNVs.
320 *Nature Genetics*. 2008;40: 1253–1260. doi:10.1038/ng.237
- 321 22. Olshen AB, Venkatraman ES, Lucito R, Wigler M. Circular binary segmentation for the analysis of
322 array-based DNA copy number data. *Biostatistics*. 2004;5: 557–572.
323 doi:10.1093/biostatistics/kxh008
- 324 23. National Cancer Institute Genomic Data Commons. [cited 18 Jan 2018]. Available:
325 <https://gdc.cancer.gov/>
- 326 24. Reynolds SM, Miller M, Lee P, Leinonen K, Paquette SM, Rodebaugh Z, et al. The ISB Cancer
327 Genomics Cloud: A Flexible Cloud-Based Platform for Cancer Genomics Research. *Cancer*
328 *Research*. 2017;77: e7–e10. doi:10.1158/0008-5472.CAN-17-0617
- 329 25. Gijsbers P, LeDell E, Thomas J, Poirier S, Bischl B, Vanschoren J. An Open Source AutoML
330 Benchmark. 6th ICML Workshop on Automated Machine Learning. 2019. Available:
331 <https://arxiv.org/pdf/1907.00909.pdf>
- 332 26. Friedman JH. Stochastic gradient boosting. *Computational Statistics and Data Analysis*. 2002;38:
333 367–378. doi:10.1016/S0167-9473(01)00065-2
- 334 27. Tenny S, Hoffman MR. Odds Ratio (OR). *StatPearls*. StatPearls Publishing; 2020. Available:
335 <http://www.ncbi.nlm.nih.gov/pubmed/28613750>
- 336 28. Zhang B. Colorectal cancer predictive test using germ-line DNA data and multiple machine
337 learning methods. 2019. Available: <https://escholarship.org/uc/item/44f3f487>
- 338 29. Olson RS, Cava W la, Mustahsan Z, Varik A, Moore JH. Data-driven advice for applying machine
339 learning to bioinformatics problems. *Pacific Symposium on Biocomputing Pacific Symposium on*
340 *Biocomputing*. 2018;23: 192–203. Available: <http://www.ncbi.nlm.nih.gov/pubmed/29218881>

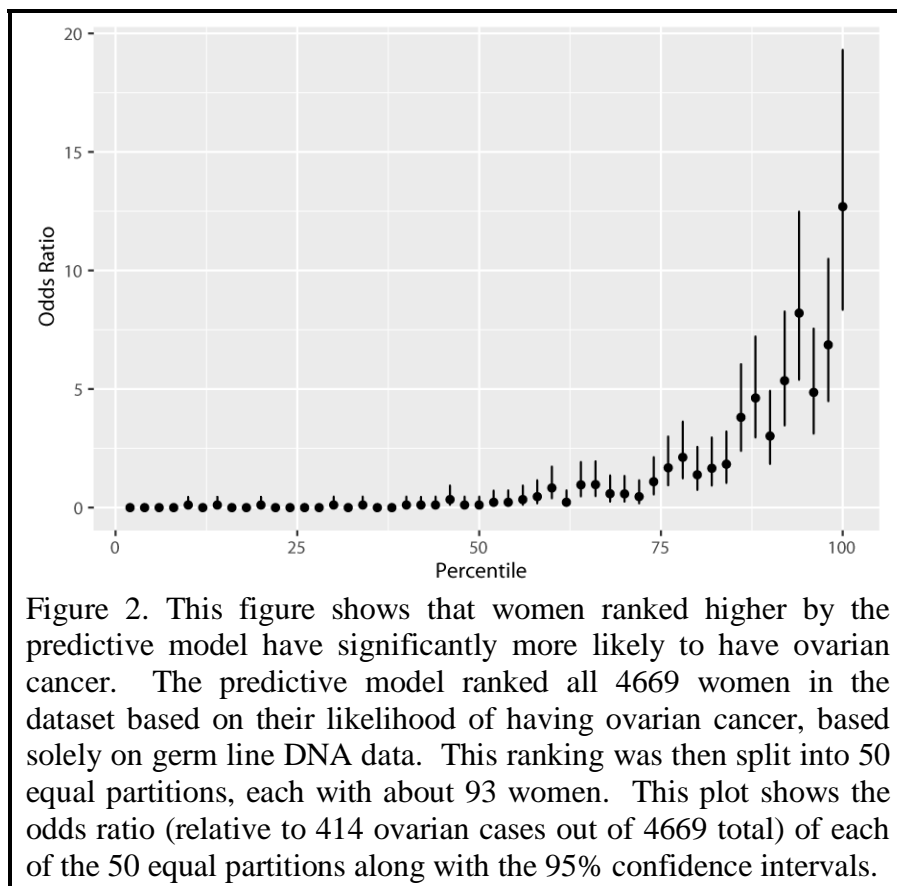
341

342



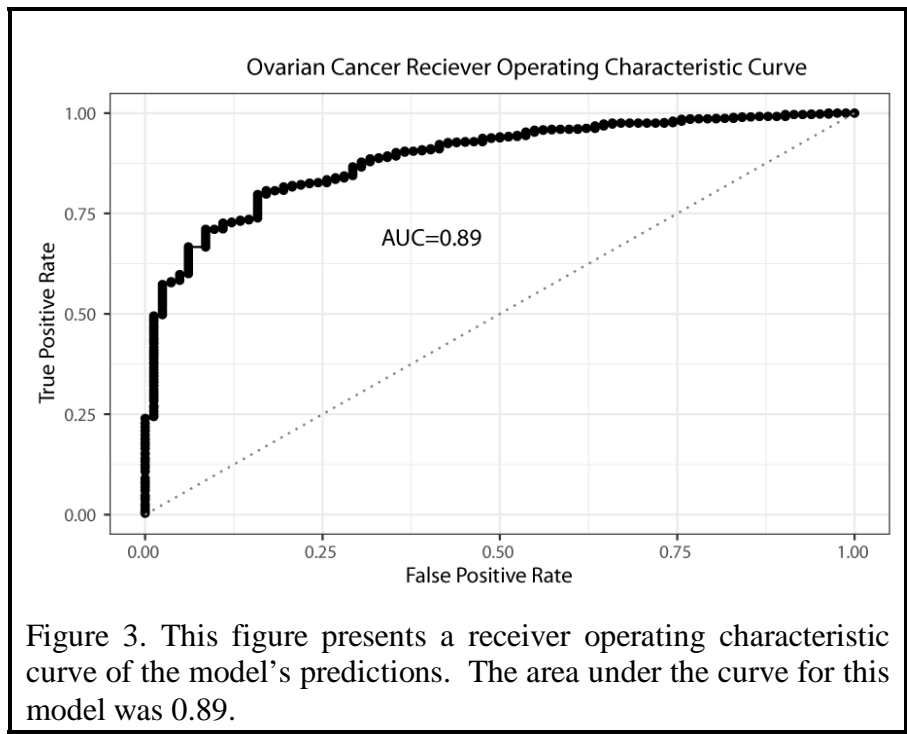
343
344
345
346
347
348
349
350
351
352
353
354
355

356



357
358
359
360
361
362
363
364
365

366



367
368
369
370
371

Figure 3. This figure presents a receiver operating characteristic curve of the model's predictions. The area under the curve for this model was 0.89.