1      **Novel ratio-metric features enable the identification of new driver genes across cancer types**

2      Malvika Sudhakar[1,2,3], Raghunathan Rengaswamy[2,3,4*], Karthik Raman[1,2,3*]

3

4      [1]Department of Biotechnology, Bhupat Jyoti Mehta School of Biosciences

5      [2]Initiative for Biological Systems Engineering (IBSE)

6      [3]Robert Bosch Centre for Data Science and Artificial Intelligence (RBC-DSAI)

7      [4]Department of Chemical Engineering, Indian Institute of Technology Madras,

8      Chennai – 600 036, India

9

10      *Corresponding authors: Email: kraman@iitm.ac.in, raghur@iitm.ac.in

11

12   **ABSTRACT**

13   An emergent area of cancer genomics has been the identification of driver genes. Driver

14   genes confer a selective growth advantage to the cell and push it towards tumorigenesis.

15   Functionally, driver genes can be divided into two categories, tumour suppressor genes

16   (TSGs) and oncogenes (OGs), which have distinct mutation type profiles. While several

17   driver genes have been discovered, many remain undiscovered, especially those that are

18   mutated at a low frequency across samples. The current methods are not sufficient to

19   predict all driver genes because the underlying characteristics of these genes are not yet

20   well understood. Thus, to predict novel genes, we need to define new features and models

21   that are not biased and identify genes that might otherwise be overshadowed by mutation

22   profiles of recurrent driver genes. In this study, we define new features and build a model to

23   identify novel driver genes. We overcome overfitting and show that certain mutation types

24    such as nonsense mutations are more important for classification. Some known cancer

25    driver genes, which are predicted by the model as TSGs with high probability are ARID1A,

26    TP53, and RB1. In addition to these known genes, potential driver genes predicted are CD36,

27    ZNF750 and ARHGAP35 as TSGs and TAB3 as an oncogene. Overall, our approach surmounts

28    the issue of low recall and bias towards genes with high mutation rates and predicts

29    potential novel driver genes for further experimental screening.

30    **Keywords**: Driver genes, random forest, cancer genomics, tumour suppressor genes,

31    oncogenes, machine learning

32    **BACKGROUND**

33    Cancer is one of the leading causes of morbidity globally, with more than 18.1 million cases

34    reported in the year 2018 [1]. A major focus of cancer research has been the understanding

35    of molecular mechanisms that govern tumorigenesis and the targets that can be used for

36    treatment. Cancer cells are distinct because of their genomes, which give these cells the

37    ability to divide and metastasize to other tissues in the body. It has been observed that

38    mutations in some genes [2, 3] confer the ability of oncogenesis to these cells. The term

39    "driver" was coined to refer to mutations in the genome that pushed the cell to oncogenesis

40    [4]. Of all the mutations present in a cancer cell, not all are involved in giving a cellular

41    advantage to the cell to divide uncontrollably. Driver mutations [4, 5] are those that were

42    advantageous for tumour development and metastasis during the clonal evolution [6, 7]. On

43    the other hand, *passenger* mutations [4, 5] are mutations that are accumulated during

44    normal cell division or due to high mutational rates in cancer cells, but their presence or

45    absence does not affect the progression and establishment of tumours.

46    Driver genes are effectively those genes that harbour mutations that provide them with a

47    selective advantage to divide and grow unchecked. These driver genes not only help the

48    cells bypass the cell cycle checkpoints to divide in an uncontrolled fashion but also give

49    added functionality, such as bypassing the immune system [8, 9] and angiogenesis [10, 11],

50    which lead to their persistence in the body. While certain cancers with well-understood

51    mechanisms show that the presence of driver mutations is recurrent in most samples of a

52    cancer type [2], others seem to have mutations that occur at a lower frequency. Driver

53    genes that contain  lower frequency of mutations are difficult to identify [12] because most

54    likely these genes work in combination with other genes to confer a selective advantage to

55    the cell.

56    Driver genes can be of two kinds depending on the role of the gene in a normal cell type. A

57    tumour suppressor genes (TSG), as the name suggests, is the cell's defence mechanism from

58    becoming a cancer cell. When such a gene loses its function due to say, frameshift

59    mutations or nonsense mutations, a selective growth advantage is conferred to the cell.

60    Proto-oncogenes undergo gain of function mutations to become into an oncogene (OG).

61    Mutations in both TSGs and OGs tip the balance of a normal cell into becoming a cancer cell.

62    While many TSGs and OGs have been discovered for different cancer types, most of them

63    are highly potent and recurring in different patients. A pan-cancer model will help in

64    identifying patterns which might be lost while studying a cohort or specific cancer type,

65    owing to low sample sizes or mutation frequencies. A key aim of this study is to find low-

66    frequency driver genes by classifying them into TSGs and OGs.

67    There are broadly two classes of methods for identifying driver genes based on mutational

68    data. The first class of methods [13–15] rely on the rate of mutations in genes for a set of

3

69    patients, to identify driver genes. In these studies, the background mutation rate is

70    estimated, and genes that show statistically different mutation rates are identified as driver

71    genes. The rate of different types of mutations is used to calculate the background mutation

72    rate [14, 15]. The methods of identification differ in the statistical method used [14]. The

73    rate of cell division and length of the gene needs to be taken into account as the mutation

74    rate may change depending on cell type and length and position of the genes [15].

75    Among the different methods that exist for identifying driver genes, when validated using

76    the Cancer Gene Census (CGC) [16], it was observed that while the precision of identifying

77    these genes was high, they had a very low recall [12]. Furthermore, genes identified through

78    these approaches have a high recurrence of being mutated across different tumour

79    samples. We now know that the rate of mutation is not sufficient for the identification of

80    driver genes; instead, genes with low mutation rate can be driver genes if a mutation occurs

81    at functionally important positions.

82    The second class of methods use a ratio-metric approach, where not only the repeated

83    occurrence of mutations is taken into consideration, but also the functional impact of the

84    mutations. Ratio-metric algorithms [17–19] capture the proportion at which the different

85    mutation types occur. The type of mutations and their ratios vary and are distinct for TSGs

86    and OGs. For instance, TSGs are more likely to have indels (insertions and deletions), more

87    specifically frameshift mutations, that lead to loss of function of the protein. On the other

88    hand, OGs tend to accumulate missense mutations that confer the protein with a "gain of

89    function" [5, 20]. These features are then used for differentiating between these two types

90    of driver genes.

91    While these methods do capture some mutation patterns observed across samples, low

92    recall shows that our understanding of the characteristics that define TSGs and OGs is far

93    from complete. In this study, we define new features that calculate entropy and frequency

94    of different mutation types along with other ratio-metric features.  Our aim is to identify

95    important features for TSGs and OGs that can help classify a given gene as a TSG or an OG.

96    Since the ratio-metric approach is based on the type of mutations and these differ for TSGs

97    and OGs, genes were classified into two classes. Further, classification problems are prone

98    to overfitting resulting in high classification scores in the training set, but the model can turn

99    out to be unreliable for predictions using new data. We outline a method for estimating

100   parameters for the given classification algorithm and avoid overfitting. We use the final

101   model to predict novel driver genes by classifying a list of unlabelled genes; we validated

102   our predictions by illustrating the presence of known TSGs and OGs among our predictions

103   and through functional analysis of the predicted novel genes. We calculated the mutation

104   rates and compared our results with the widely used tool MutSigCV and show that our

105   method is able to pick out many driver genes that have very low mutation rates. Further, we

106   used a pan-cancer model to predict driver genes that were tissue-specific.


107   **RESULTS**

108   We define novel features and a method to estimate parameters and build a classifier using

109   pan-cancer data to predict TSGs and OGs. The classifier is further used to predict labels for

110   unlabelled genes, at pan-cancer and tissue-specific levels, which are analysed for functional

111   enrichment.

### 112    Novel features used for classification of TSGs and OGs

113    We trained multiple random forest models using a subset (80%) of 136 TSGs and 76 OGs for

114    each fold of the cross-validation. We performed a five-fold cross-validation while estimating

115    hyper-parameters for the model followed by multiple random iterations to estimate stable

116    hyper-parameters and avoid overfitting (as defined in Methods). It is important to carefully

117    consider overfitting as the initial training set is not very large. The accuracy for the test set

118    reduces compared to the training set, but this difference is not substantial. We note that

119    TSGs can be predicted with higher accuracy than OGs; it is probable that the features are

120    biased at capturing information regarding TSGs better than OGs. Across the multiple

121    models, an average accuracy of 0.76 ± 0.03 was achieved. These models were further used

122    for the identification of novel genes as well as tissue-specific analyses. Our model presents a

123    significant improvement in recall for TSGs. For OGs, the recall is similar to those observed in

124    other tools. Nevertheless, an average recall of driver genes (comprising both classes) shows

125    an improvement over the tools reported earlier [12].

126    To identify features important for the classification of TSGs and OGs, we calculated the

127    average rank of each feature, across all models. We observe that the top-ranking features

128    contain LOF and missense mutations (Supplementary Table S1). The new features that

129    replace old features in the top 18 ranks are Nonsense entropy, High missense frequency,

130    Compound/benign,    High    Frameshift    Frequency,    Damaging/kb,    Compound/kB,

131    Damaging/LoFl and HiFl/benign.  Further, we used the training set genes to compare the

132    distribution of feature values in TSG and OGs, and observed that our top-ranking features

133    show the highest differences between the two distributions (Fig 1). While it is common

134    knowledge that LOF mutations accumulate in TSG and recurrent missense mutations in OGs,

135    we formally show that the feature distribution is different for these two functional classes.

136    **Iterative hyper-parameter estimation avoids overfitting**

137    Initial analysis for a large number of *n_estimator* for random forest and using

138    BalancedBagging to manage class imbalance gave higher accuracy score for training sets

139    comparable to Davoli *et al.,* (2013). However, these showed very low accuracy for the test

140    set (Table 2), indicating overfitting. Additionally, we observed that changing the random

141    seed showed substantial variation in results. This variation is unexpected and could perhaps

142    stem from non-optimum parameters used for classification or the small size of the data. To

143    avoid this variation, we re-estimated the random forest parameters, *n_estimator*,

144    *max_features*, *max_depth* and *criterion*. Changing the *n_estimator* had a major effect on

145    classification, and grid search with cross-validation did not help in removing overfitting.

146    We overcame this by multiple iterations of hyper-parameter estimation by changing the

147    random seed, which helps us identify more stable hyper-parameters. This gave lower

148    accuracy for training sets but improved the accuracy of the test set considerably.  When

149    varying sets of random seeds (10, 20, 40, 80, 160, 320) were used, the results were

150    consistent across all cross-validation folds (test set accuracy 0.76 and standard deviation

151    0.03) implying the increasing number of random seed iterations do not decrease or improve

152    accuracy. We observe that for a given data fold, the hyper-parameters selected are more

153    stable for varying sets of random seeds. While different parameter sets dominate as the

154    data is changed, the overall results on the test set do not vary.

155 **Model identified novel TSGs and OGs along with known driver genes**

156 All genes that were not used for training the models were classified into TSGs and OGs. This

157 list also contained genes that are known driver genes present in CGC but not used for

158 training. The labels were predicted for the unlabelled genes, of which 126 genes or

159 transcripts showed consensus across all models. CGC known driver genes contributed to

160 40.5% of these predictions which included genes such as ARID1A, ATRX, NF1, TP53, RB1, and

161 STAG1 and their transcripts. Some novel genes predicted consistently are SIN3A, ZNF750,

162 IWS1, CD36, ARHGAP35, MGA, and RASA1 as TSGs. The model tends to be biased towards

163 TSGs with 699 genes with consistent predictions for 3 or more models out of which only 9

164 are predicted as OGs. The top OGs predicted are U2AF1, BCL2L10, KRAS, MAP1LC3B,

165 C11orf68, TAB3, MED12, MAX, and BRAF. Further, we show not all transcripts of a gene

166 behave like a driver gene, for e.g. ATRX transcript ENST00000373344 is labelled as TSG but

167 not ENST00000400866, ENST00000373341. The presence of known driver genes among top

168 TSG and OG shows the validity of the model and those other genes in the list are potential

169 driver genes.

170 Enrichment analysis of genes for various KEGG and BIOCARTA pathways revealed genes

171 involved in different cancer pathways such as myeloid leukaemia, and pancreatic cancer.

172 Genes are also enriched for various signalling pathways associated with cell growth, such as

173 EGF and PDGF signalling pathways. Further, to validate, a similar analysis was conducted

174 using genes used for training the model. We find GO terms related to cell cycle, regulation

175 of transcription, signalling and cell cycle arrest to be common for both results. These

176 keywords were further clustered with top clusters associated with genes involved in zinc-

177    finger proteins, helicases, ATP-binding, ARID binding and cancer pathways. The analysis

178    shows known driver genes and predicted driver genes enrich for similar pathways.

179    **Our approach identifies genes with low mutation frequency**

180    We analysed the mutation frequencies of the predicted genes. Mutation rates were

181    calculated using MutSigCV, a well-known driver gene predictor, which calculates mutation

182    rates to identify driver genes. MutSigCV ranks all genes of which a total of 602 driver genes

183    were identified above the threshold (p <=0.005, q <= 0.01). Training data labels were used

184    to compare the two methods. MutSigCV identified 40% for our training gene set with 85

185    genes predicted as driver, while our model did better by predicting 85% of genes. The

186    mutation rates of the genes predicted by the two models were compared. Since MutSigCV

187    ranks all genes, we picked top genes equal in size to our model predictions (>=5 model

188    consensus) and calculated KS statistic against training set and plotted the fraction of genes

189    below mutation rate of each gene. We observe that distribution of mutation rates is similar

190    to training set for our predicted genes, while MutSigCV tends to be biased towards genes

191    with higher mutation rates (Fig 2). The minimum mutation rate predicted for our model was

192    0.35 while for MutSigCV was 0.90. The KS (Kolmogorov-Smirnov) statistic for both models

193    when compared to training set shows the difference is far lesser for our model when

194    compared to MutSigCV (Table 4), which shows that the distribution of mutation rates is

195    similar to what is expected.

196    **Driver genes are tissue-specific**

197    Cohort studies tend to be specific to a cancer type. The usefulness of a pan-cancer model is

198    further elucidated when it can be used to identify tissue-specific driver genes. The objective

199    of predicting genes using a subset of data specific to tumour primary tissue source was to

9

200    identify genes specific to a cancer type. This helped in identifying genes which might

201    otherwise be lost in biological noise (Table 5). We observe TP53 predicted as TSG across the

202    different tissues. Other known driver genes that weren't identified by the pan-cancer

203    analysis were identified such as CBFB, CDH1, PTEN in breast cancer and APOB in liver. Genes

204    such FAM182A, SOX9, AHNAK2, ENSG00000121031, FLT3LG, PMEPA1, ZFP36L2 in the large

205    intestine, ALB, KRTAP19-1, APOB, CD200, CRYGD, KRTAP24-1, OR6N2 in the liver are novel

206    predictions, and their functions in these cancers can further be studied. We used the pan-

207    cancer models to predict tissue-specific driver genes and identified new genes not reported

208    by the pan-cancer analysis.

209    Genes identified for breast cancer was validated by supporting literature. CBFB [21] and

210    PTEN [22, 23] is a known TSG in breast cancer. PTEN is found to be under-expressed in

211    breast cancer [24, 25]. While CDH1 mutations are found mostly in stomach cancer, they are

212    also shown to be frequently occurring in lobular breast cancer [26, 27]. Pathway analysis of

213    breast cancer genes shows enrichment of pathways involved in gene expression regulation

214    governed by TP53, RUNX1 and PTEN which includes pathways that regulates estrogen-

215    mediated transcription. CBFB deletion leads to expression loss of RUNX1[21], which can no

216    longer regulate NOTCH signalling by repression, which is confirmed by pathway analysis.

217    Some apoptosis pathways are enriched that include CDH1 and TP53 genes. The genes

218    identified by the pan-cancer model for breast cancer samples identify genes functionally

219    important in breast tumour cells.

220    Predictions made for liver cancer were mostly novel, which made literature validation

221    difficult. RNA expression levels of genes APOB, ALB and CD200 were higher compared to all

222    other tissues (as reported by The Human Protein Atlas). Higher albumin levels are known to

10

223    decrease the risk of HCC (Hepatocellular carcinoma) [28]. APOB mutational signatures are

224    shown computationally to be significant to predict prognosis, by loss of regulation of genes

225    such as TP53, PTEN, HGF [29]. While role of other genes is difficult to elucidate, our method

226    helps identify research gaps which can be filled by studying these potential driver genes.

227    **DISCUSSION**

228    Identification of driver genes has been an important focus area of cancer research because

229    these genes are potential targets for therapy and biomarkers. Different approaches have

230    been used for identification using mutational information [17, 18, 30], gene expression

231    levels [31], protein structural information [32], network analysis [33, 34] or using multiple

232    data sources [31]. Advances in sequencing technologies have made mutational information

233    easily available, and different tools have been developed to identify driver genes. Driver

234    genes are further classified into TSGs and OGs based on the functional impact of the

235    mutations they harbour.

236    We adopt a classification approach that is able to predict TSGs and OGs by leveraging a set

237    of ratio-metric and other new features. Traditional methods identify genes based on the

238    mutation rate. Compared to previous approaches, we ascribe a higher significance to

239    functional impact along with the position of the mutations, as the genes might contain

240    mutations in functionally important regions even though the mutation rate may not be very

241    different from the background mutation rate. Features like nonsense entropy, frameshift

242    frequency captures the recurrence of a mutation when multiple samples are considered,

243    thus taking into account the position at which the mutation occurs.

244    For classification, many different algorithms are available, but the performance of the

245    algorithm is dependent on the data and estimation of parameters. It is especially important

11

246     while solving biological problems, where the training data might be small, to build robust

247     models. We tried the classification of genes using support vector machines (SVM), logistic

248     regression, balanced bagging as well as random forest and found that random forest

249     performed better in this case. Further, high performance on a given data might also be due

250     to overfitting. We sought to avoid overfitting by performing a standard 5-fold cross-

251     validation while estimating random forest parameters as well as multiple iterations for

252     estimation of stable parameters. We developed a procedure to verify that the predictions

253     are reasonably stable. An ensemble of models is used to make final predictions.

254     It is important that the estimated parameters are robust to changes in data. For random

255     forest, we estimated four parameters out of which *n_estimator* seemed to have a large

256     effect on the classification. For large values of *n_estimator*, we were able to show high

257     accuracy scores similar to Davoli et al., (2013) but the accuracy scores for test set were

258     much lower. We were not able to compare our performance on the test set with that of

259     Davoli *et al* (2013), as their test set results have not been published. To build a better model

260     that is not biased to data, we needed a more robust classifier, that is sufficiently generalized

261     and not dependent on the training data.

262     The models generated were used to find which of the new features are important for

263     classification. To evaluate the model, we used 5-fold cross-validation with 20% test dataset

264     while maintaining the ratio between TSGs and OGs and calculated metrics such as accuracy

265     and F1 score. Instead of AUROC (Area under Receiver Operating Characteristic), we chose to

266     show accuracy and F1 score, as AUROC only helps in estimating if the model can separate

267     the given classes but tells us very little about the classification power for each of these

268     classes. The F1 score is calculated for each of the given classes and helps understand if the

269    model is biased towards any one of the classes. The accuracy score on the test set shows

270    that mere accuracy is not sufficient for judging a model. The models perform slightly better

271    for TSGs, though it is far poorer at classifying OGs.

272    While assessing the model, it is important to use metrics such as F1 score, as it scores

273    predictions for each of the classes. Studies reporting only AUROC statistics present an

274    incomplete picture and are not effective in estimating the performance of the model,

275    especially in datasets having a class imbalance [35]. This is evident when we compare

276    AUROC of Balanced bagging model (0.76 ± 0.07) with our model (0.54 ± 0.07). AUROC gives

277    measures the models ability to separate the classes and not the prediction power. By

278    reporting both accuracy as well as F1 score, we show that the model does not perform

279    equally for both classes but tends to be better at classifying TSG than OG. This indicates that

280    the chosen features are not sufficient to classify oncogenes.

281    Feature ranking shows that features containing information about LOF, nonsense,

282    frameshift and missense mutations are important. Nonsense and frameshift mutations are

283    frequently seen in TSGs while recurrent missense mutations are characteristic of OGs as

284    they lead to "gain of function".

285    The list of genes classified contained known driver genes and other transcript data for genes

286    present in training and test set. We found that TSGs such as ATRX, PTCH1, and STAG2 were

287    classified as TSGs with high probability. KDM6A gene and its transcripts (ENST00000377967,

288    ENST00000382899) feature among the top, which shows that the model can also help

289    classify a particular transcript of a gene. Similarly, TP53 and its six transcripts were all

290    classified as TSGs.  Genes U2AF1, KRAS, BRAF, MED12 and MAX were classified as OGs

291    among the top genes identified as OGs. As the probability scores for OGs tend to be lesser

292    than TSGs, relatively fewer OGs make the cut-off for the top 5 percentile.

293    Among the top TSGs identified, CD36 (previously known as FAT) is receptor protein for fatty

294    acids. CD36 is also a prognostic marker for different cancer types [36, 37] and found in

295    metastatic cells [36, 38]. While the expression of a gene is markedly different from normal

296    cells, the molecular mechanism that enables metastasis is not well understood. Another

297    gene, ARHGAP35, is a glucocorticoid receptor DNA binding factor, which has also been

298    previously identified as a potential driver gene by other methods [39, 40]. ZNF750, zinc

299    finger protein 750 has been established as a tumour suppressor in oesophageal squamous

300    cell carcinoma [41–43] though it is absent from the CGC diver gene list. Some other

301    potential TSGs not present in the CGC list are MBD6 and RASA1. In the human protein atlas,

302    MAP1LC3B is labelled as a prognostic marker for renal and stomach cancer among the three

303    shortlisted OGs.

304    Our model does have some limitations. We have used binary classification for identification

305    of TSGs and OGs which, classifies all genes as either TSG or OG. All genes containing

306    mutations are not driver genes, and thus, a majority of genes are neutral. We overcome this

307    by taking consensus across the five models built. It may be possible to improve on this

308    classification by solving a multi-class problem where each gene is identified as TSG, an OG or

309    neutral gene. The difficulty in this problem stems from the huge class imbalance in the data

310    as well as the definition of neutral genes. While there are studies showing the importance of

311    a gene in tumour evolution, it is difficult to define genes that are not involved in cancer

312    progression. Most methods use a list of genes that do not contain cancer driver genes and

313    genes involved in cancer pathways, but this does not exclude potential driver genes.

314 Additionally, it has been seen that mutations are not always the reason for the change in

315 functionality and regulation might also lead to change in expression at transcriptomic and

316 proteomic levels. Other than adding new features to the analysis, including transcriptomic

317 and proteomic data along with genomic mutation data might further improve the

318 classification of genes.

319 **CONCLUSION**

320 In summary, we see two main contributions of our paper. First, we developed a classifier,

321 which enabled an improved recall of TSGs and OGs compared to previously proposed

322 methods in the literature. We carefully avoided overfitting for achieving consistent and high

323 confidence results. Second, we predicted many potential TSGs and OGs at both the pan-

324 cancer and tissue-specific level, which form a ready short-list for further experimental

325 investigation. Some of the top predictions were already well-known cancer drivers while

326 others are reported in multiple cancer studies though their role in tumorigenesis is not yet

327 well understood. Our approach is also readily amenable to the integration of other omic

328 datasets, as they become available.

329 **METHODS**

330 **Data**

331 We downloaded somatic mutation data from Catalogue of Somatic Mutations in Cancer

332 (COSMIC) (v79) [44]. These data were pre-processed to exclude hyper-mutated samples

333 (samples containing more than 2000 mutations) Known SNPs were retained only if they

334 were "confirmed somatic mutations". The final processed data consist of 2,145,044

335 mutations from 20,667 samples across 37 primary tissues. COSMIC also contains transcript

336     information, where different transcripts of a gene are saved as "gene_transcript" and are

337     handled as separate genes. Splice site mutations were identified as mutations at 1 or 2 bps

338     after the end of the exon border or 1 or 2 bps before the start of exon border. We used the

339     popular tool Polyphen2 [45] to predict the phenotypic impact of missense mutations. For

340     some mutations, Polyphen2 returns multiple scores, which we averaged for the purpose of

341     our analyses.

342     TSGs and OGs for training and test were taken from the CGC [16] gene list. Only those genes

343     that were labelled "TSG" or "OG" and not "Fusion" were used for this analysis. A total of 213

344     driver genes were used, of which 136 were TSGs and 77 were OGs. The TSG:OG ratio was

345     maintained during all cross-validation steps and in both training and test sets.

346     **Ratio-metric features**

347     Mutations were divided into 11 different categories [17, 45]: silent, missense, splicing, High

348     Functional Impact (HiFI), Mid Functional Impact (MiFI), Low Functional Impact (LoFI),

349     nonsense, frameshift, in-frame, nonstop or complex. Not all missense mutations are equally

350     deleterious — labelling them into HiFI, MiFI and LoFI categories helps differentiate genes

351     that have a large number of mutations with low impact, from genes that have relatively

352     fewer mutations but with larger functional impact. We use PolyPhen2 scores to categorise

353     mutations as HiFI ($\geq$ 0.85), LoFI ($\leq$ 0.15) and MiFI (between 0.15 and 0.85), to differentiate

354     between high confidence pathogenic mutation predictions.

355     Additionally, other mutation categories were defined, which clubbed multiple mutations

356     into one, such as 'compound' and 'damaging'. Compound mutations are included because

357     mutations types such as in-frame, nonsense and complex occur at a lower frequency than

358     single nucleotide missense mutations, which might lead to patterns and impact of these

16

359    mutations to be masked. Since the functional impact is similar to missense mutations,

360    combining similar mutation types might help in capturing information of these less

361    frequently observed mutation types. Loss of function (LOF) mutations introduce large

362    changes into proteins, causing disruption of function. Damaging mutations are the sum of

363    HiFI and MiFI mutations; these capture impact of multiple MiFI and sparse HiFI mutations.

364    Many features compute a ratio of mutation types, as outlined in Table 6. We defined 37

365    features in all, with 18 of them being similar to those defined as Davoli *et al.*, (2013).

366    Entropy and Frequency features

367    Entropy and frequency features were defined for four mutation types. A mutation (M$_i$) in a

368    given gene *i* is represented by its location. For missense mutations, M$_i$ is represented as a

369    tuple (*loc*, *wt*, *mt*) where *loc* is the location of the mutation, *wt* is the wild type nucleotide,

370    and *mt* is the mutated nucleotide. If *k* unique mutations are present in a gene, $f_i$ gives the

371    frequency for each of the mutations.

$$f_i = \frac{n_M}{n}$$

372    where $n_M$ is the number of occurrences of mutation $M$ and $n$ is the number of mutations in

373    gene $i$.

$$S = \sum_{i=1}^{k} f_i \log f_i$$

$$Entropy = \log k - S$$

374    **Classification of genes**

375    Different machine learning algorithms such as random forest, support vector machines and

376    logistic regression were used, among which random forest gave the highest accuracy.

17

377    Random forest was used for building a robust model and classifying TSGs and OGs. We used

378    five-fold cross-validation to split data into training to test set ratio of 8:2; where each fold

379    acts as a test set. We used the implementation of Random forest from the Python package

380    Sci-Kit Learn [46]. We tuned the parameters using a five-fold cross-validation grid search

381    along with multiple random iterations of random seed (described later). The parameters

382    tuned are *n_estimator* (from 5-40), *max_features* ('sqrt' or 'log2'), *max_depth* (2-4) and

383    *criterion* ('gini' or 'entropy'). The number of maximum features each decision tree considers

384    is given by the parameter *max_features*, which can be calculated in two ways, as either the

385    square root or $\log_2$ of the total number of features.

386    **Tuning hyperparameters and estimating the robustness of the classifier**

387    Our initial results showed variation in classification depending on the random seed that was

388    selected for classifying, even though cross-validation was used while estimating parameters.

389    We used balanced bagging classifier to take into consideration the class imbalance and

390    estimated parameters using cross-validation, which is the standard method. Poor results for

391    this model led us to estimate hyper-parameters differently.

392    To avoid this variation, classification and parameter selection were done for multiple

393    random seeds (Fig. 3 block B). Grid search with five-fold cross-validation was done for

394    multiple different random seeds. Optimum parameters were selected by first estimating

395    parameter '*n_estimator*' and using it to estimate other parameters. Recurrence of

396    'n_estimator' across different random seeds was counted, and the maximum count was

397    considered as the best 'n_estimator' to be given to the model. If multiple estimators were

398    chosen, maximum accuracy during cross-validation was used to select one estimator.

399    Maximum accuracy was used to find other parameters for the given best '*n_estimator*'. The

18

400    classification was rerun using the given parameters, and features were ranked. The model

401    was used to predict the classification of test set genes.

402    To estimate the effect of the number of random iterations on parameter estimation, the

403    classifier was built on a varying number of iterations of random seeds (10, 20, 40, 80, 160,

404    320). The stability of hyper-parameters selected was analysed based on the variation in the

405    accuracy of the test dataset.

406    **Feature comparison and ranking**

407    All features defined were used for classification and ranked depending on their contribution

408    to the model. Average rank was calculated across the five validation sets. The features are

409    given in Table 7.

410    **Identification and functional analysis of novel TSGs and OGs**

411    We used the model built on the combined set of 37 features to classify unlabelled genes

412    into TSGs and OGs. In total, 26,866 genes were classified as TSGs or OGs and ranked using

413    their probabilities for each class. The genes given for classification contains different

414    transcripts of the same gene symbol as different genes. In all, the gene list contained 18,951

415    unique gene symbols. Genes were labelled TSG and OG depending on their presence in the

416    top 5 percentile and consensus across models built during cross-validation. Since not all

417    genes are necessarily TSGs or OGs, genes which didn't fulfil these criteria remained

418    unlabelled. Novel TSG and OG gene list predicted by greater than four models were further

419    used for functional analysis to find the major pathways and gene ontologies these genes are

420    enriched for. Functional analysis was carried out using DAVID [47, 48] for both, the genes

421    above the threshold as well as training set genes, and the results were compared.

422    Further, the pan-cancer classifier was used to predict genes in different cancer types based

423    on the primary tissue where the tumour is formed. The data were filtered based on primary

424    tissue, and the feature matrix was generated for tissues with >1000 samples. The data was

425    then standardized and run using pan-cancer models described earlier.

426    We compared and calculated mutation rates using MutSigCV. Since the ground truth is not

427    known for these predicted genes, we compared the genes used for training and calculated

428    recall of these genes. Since MutSigCV does not classify genes as TSG or OG, the classes

429    considered were Driver and Passenger. Further, we were interested in looking at the

430    mutation rate distribution across the genes predicted. Since the distribution of mutation

431    rates is unknown, we compared the similarity of the distribution of the predicted genes with

432    the genes used for training (Kolmogorov-Smirnov statistic). Similarly, the similarity was

433    compared for genes predicted by MutSigCV.

434    **LIST OF ABBREVIATIONS**

435    AUROC: Area under Receiver Operating Characteristic

436    CGC: Cancer Gene Census

437    COSMIC: Catalogue of Somatic Mutations in Cancer

438    GO: Gene ontology

439    HCC: Hepatocellular carcinoma

440    HiFI: High Functional Impact

441    Indels: insertions and deletions

442    KS statistic: Kolmogorov-Simrnov statistic

443    LOF: Loss of function

444    LoFI: Low Functional Impact

445    MiFI: Mid Functional Impact

446    OG: oncogenes

447    TSG: tumour suppressor gene

448    **REFERENCES**

449    1. Ferlay J, Colombet M, Soerjomataram I, Mathers C, Parkin DM, Piñeros M, et al.

450    Estimating the global cancer incidence and mortality in 2018: GLOBOCAN sources and

451    methods. International Journal of Cancer. 2019.

452    2. Miki Y, Swensen J, Shattuck-Eidens D, Futreal PA, Harshman K, Tavtigian S, et al. A strong

453    candidate for the breast and ovarian cancer susceptibility gene BRCA1. Science.

454    1994;266:66–71. doi:10.1126/science.7545954.

455    3. Wooster R, Bignell G, Lancaster J, Swift S, Seal S, Mangion J, et al. Identification of the

456    breast cancer susceptibility gene BRCA2. Nature. 1995;378:789–92. doi:10.1038/378789a0.

457    4. Stratton M, Campbell P, Futreal P. The cancer genome. Nature. 2009;458:719–24.

458    doi:10.1038/nature07943.The.

459    5. Vogelstein B, Papadopoulos N, Velculescu VE, Zhou S, Diaz L a, Kinzler KW. Cancer

460    Genome Landscapes. Science (80- ). 2013;339:1546–58. doi:10.1126/science.1235122.

461    6. Greaves M, Maley CC. Clonal evolution in cancer. Nature. 2012;481:306–13.

462    doi:10.1038/nature10762.

463    7. Burrell R a, McGranahan N, Bartek J, Swanton C. The causes and consequences of genetic

464    heterogeneity in cancer evolution. Nature. 2013;501:338–45. doi:10.1038/nature12625.

465    8. Beishline K, Azizkhan-Clifford J. Sp1 and the "hallmarks of cancer." FEBS Journal.

466    2015;282:224–58.

467    9. Cavallo F, De Giovanni C, Nanni P, Forni G, Lollini PL. 2011: The immune hallmarks of

468    cancer. In: Cancer Immunology, Immunotherapy. 2011. p. 319–26.

469    10. Shahmarvand N, Nagy A, Shahryari J, Ohgami RS. Mutations in the signal transducer and

470    activator of transcription family of genes in cancer. Cancer Sci. 2018; December 2017:1–8.

471    11. Zhang E, Feng X, Liu F, Zhang P, Liang J, Tang X. Roles of PI3K/Akt and c-Jun signaling

472    pathways in human papillomavirus type 16 oncoprotein-induced HIF-1alpha, VEGF, and IL-8

473    expression and in vitro angiogenesis in non-small cell lung cancer cells. PLoS One.

474    2014;9:e103440.

475    12. Hofree M, Carter H, Kreisberg JF, Bandyopadhyay S, Mischel PS, Friend S, et al.

476    Challenges in identifying cancer genes by analysis of exome sequencing data. Nat Commun.

477    2016;7 May:12096. doi:10.1038/ncomms12096.

478    13. Tamborero D, Gonzalez-Perez A, Lopez-Bigas N. OncodriveCLUST: exploiting the

479    positional clustering of somatic mutations to identify cancer genes. Bioinformatics.

480    2013;29:2238–44. doi:10.1093/bioinformatics/btt395.

481    14. Dees ND, Zhang Q, Kandoth C, Wendl MC, Schierding W, Koboldt DC, et al. MuSiC:

482    Identifying mutational significance in cancer genomes. Genome Res. 2012;22:1589–98.

483    15. Lawrence MS, Stojanov P, Polak P, Kryukov G V., Cibulskis K, Sivachenko A, et al.

484    Mutational heterogeneity in cancer and the search for new cancer-associated genes.

485    Nature. 2013;499:214–8.

486    16. Futreal PA, Coin L, Marshall M, Down T, Hubbard T, Wooster R, et al. A census of human

487    cancer genes. Nat Rev Cancer. 2004;4:177–83. doi:10.1038/nrc1299.

488    17. Davoli T, Xu AW, Mengwasser KE, Sack LM, Yoon JC, Park PJ, et al. Cumulative

489    Haploinsufficiency and Triplosensitivity Drive Aneuploidy Patterns and Shape the Cancer

490    Genome. Cell. 2013;155:948–62. doi:10.1016/j.cell.2013.10.011.

491    18. Melloni GE, Ogier AG, de Pretis S, Mazzarella L, Pelizzola M, Pelicci P, et al. DOTS-Finder:

492    a comprehensive tool for assessing driver genes in cancer genomes. Genome Med.

493    2014;6:44. doi:10.1186/gm563.

494    19. Tokheim CJ, Papadopoulos N, Kinzler KW, Vogelstein B, Karchin R. Evaluating the

495    evaluation of cancer driver genes. Proc Natl Acad Sci. 2016;113:14330–5.

496    doi:10.1073/pnas.1616440113.

497    20. Hanahan D, Weinberg RA. The Hallmarks of Cancer. Cell. 2000;100:57–70.

498    doi:10.1016/S0092-8674(00)81683-9.

499    21. Malik N, Yan H, Moshkovich N, Palangat M, Yang H, Sanchez V, et al. The transcription

500    factor CBFB suppresses breast cancer through orchestrating translation and transcription.

501    Nat Commun. 2019.

502    22. Lu Y, Lin YZ, LaPushin R, Cuevas B, Fang X, Yu SX, et al. The PTEN/MMAC1/TEP tumor

503    suppressor gene decreases cell growth and induces apoptosis and anoikis in breast cancer

504    cells. Oncogene. 1999.

505    23. Weng L-P. PTEN coordinates G1 arrest by down-regulating cyclin D1 via its protein

506    phosphatase activity and up-regulating p27 via its lipid phosphatase activity in a breast

507    cancer model. Hum Mol Genet. 2001.

508    24. Li S, Shen Y, Wang M, Yang J, Lv M, Li P, et al. Loss of PTEN expression in breast cancer:

509    Association with clinicopathological characteristics and prognosis. Oncotarget. 2017.

510    25. Zhang HY, Liang F, Jia ZL, Song ST, Jiang ZF. PTEN mutation, methylation and expression

511     in breast cancer patients. Oncol Lett. 2013.

512     26. Hansford S, Kaurah P, Li-Chang H, Woo M, Senz J, Pinheiro H, et al. Hereditary diffuse

513     gastric cancer syndrome: CDH1 mutations and beyond. JAMA Oncol. 2015.

514     27. Schrader KA, Masciari S, Boyd N, Wiyrick S, Kaurah P, Senz J, et al. Hereditary diffuse

515     gastric cancer: Association with lobular breast cancer. In: Familial Cancer. 2008.

516     28. Nojiri S, Joh T. Albumin suppresses human hepatocellular carcinoma proliferation and

517     the cell cycle. Int J Mol Sci. 2014.

518     29. Lee G, Jeong YS, Kim DW, Kwak MJ, Koh J, Joo EW, et al. Clinical significance of APOB

519     inactivation in hepatocellular carcinoma. Exp Mol Med. 2018;50:147. doi:10.1038/s12276-

520     018-0174-2.

521     30. Kumar RD, Searleman AC, Swamidass SJ, Griffith OL, Bose R. Statistically identifying

522     tumor    suppressors    and    oncogenes    from    pan-cancer    genome-sequencing    data.

523     Bioinformatics. 2015;31:3561–8.

524     31. Sanchez-Garcia F, Villagrasa P, Matsui J, Kotliar D, Castro V, Akavia UD, et al. Integration

525     of Genomic Data Enables Selective Discovery of Breast Cancer Drivers. Cell. 2014;159:1461–

526     75.

527     32. Fujimoto A, Okada Y, Boroevich KA, Tsunoda T, Taniguchi H, Nakagawa H. Systematic

528     analysis of mutation distribution in three dimensional protein structures identifies cancer

529     driver genes. Sci Rep. 2016;6 May:26483. doi:10.1038/srep26483.

530     33. Ramsahai E, Walkins K, Tripathi V, John M. The use of gene interaction networks to

531     improve    the    identification    of    cancer    driver    genes.    PeerJ.    2017;5:e2568.

532     doi:10.7717/peerj.2568.

533     34. Chen Y, Hao J, Jiang W, He T, Zhang X, Jiang T, et al. Identifying potential cancer driver

24

534   genes by genomic data integration. Sci Rep. 2013;3:3538. doi:10.1038/srep03538.

535   35. Jeni LA, Cohn JF, De La Torre F. Facing Imbalanced Data--Recommendations for the Use

536   of Performance Metrics. In: 2013 Humaine Association Conference on Affective Computing

537   and Intelligent Interaction. 2013. p. 245–51. doi:10.1109/ACII.2013.47.

538   36. Ladanyi A, Mukherjee A, Kenny HA, Johnson A, Mitra AK, Sundaresan S, et al. Adipocyte-

539   induced CD36 expression drives ovarian cancer progression and metastasis. Oncogene.

540   2018.

541   37. Hale JS, Otvos B, Sinyuk M, Alvarado AG, Hitomi M, Stoltz K, et al. Cancer stem cell-

542   specific scavenger receptor 36 drives glioblastoma progression. Stem Cells. 2014.

543   38. Pascual G, Avgustinova A, Mejetta S, Martín M, Castellanos A, Attolini CSO, et al.

544   Targeting metastasis-initiating cells through the fatty acid receptor CD36. Nature. 2017.

545   39. Lawrence MS, Stojanov P, Mermel CH, Robinson JT, Garraway LA, Golub TR, et al.

546   Discovery and saturation analysis of cancer genes across 21 tumour types. Nature.

547   2014;505:495–501. doi:10.1038/nature12912.

548   40. Zhang Y, Zhang L, Li R, Chang DW, Ye Y, Minna JD, et al. Genetic variations in cancer-

549   related significantly mutated genes and lung cancer susceptibility. Ann Oncol.

550   2017;28:1625–30. doi:10.1093/annonc/mdx161.

551   41. Lin DC, Hao JJ, Nagata Y, Xu L, Shang L, Meng X, et al. Genomic and molecular

552   characterization of esophageal squamous cell carcinoma. Nat Genet. 2014.

553   42. Otsuka R, Akutsu Y, Sakata H, Hanari N, Murakami K, Kano M, et al. ZNF750 Expression Is

554   a Potential Prognostic Biomarker in Esophageal Squamous Cell Carcinoma. Oncology.

555   2018;94:142–8. doi:10.1159/000484932.

556   43. Hazawa M, Lin DC, Handral H, Xu L, Chen Y, Jiang YY, et al. ZNF750 is a lineage-specific

557 tumour suppressor in squamous cell carcinoma. Oncogene. 2017;36:2243–54.

558 44. Forbes SA, Beare D, Boutselakis H, Bamford S, Bindal N, Tate J, et al. COSMIC: Somatic

559 cancer genetics at high-resolution. Nucleic Acids Res. 2017;45:D777–83.

560 45. Adzhubei IA, Schmidt S, Peshkin L, Ramensky VE, Gerasimova A, Bork P, et al. A method

561 and server for predicting damaging missense mutations. Nature Methods. 2010;7:248–9.

562 46. Pedregosa F, Varoquaux G. Scikit-learn: Machine learning in Python. 2011.

563 doi:10.1007/s13398-014-0173-7.2.

564 47. Huang DW, Sherman BT, Lempicki RA. Systematic and integrative analysis of large gene

565 lists using DAVID bioinformatics resources. Nat Protoc. 2009;4:44–57.

566 48. Huang DW, Sherman BT, Lempicki RA. Bioinformatics enrichment tools: Paths toward

567 the comprehensive functional analysis of large gene lists. Nucleic Acids Res. 2009.

568 **FIGURE LEGENDS**

569 *Figure 1 Distribution of top features identified by the classifier for TSG and OG.* Training

570 genes were used to study the differences between the distributions of features (kernel

571 density) in TSG and OG. Kolmogorov-Smirnov statistic and the p-value is given for each

572 feature. Higher value of KS statistic shows magnitude of difference of the two distributions.

573 *Figure 2 Fraction of genes predicted plotted against log transformed mutation*

574 *rates.* Genes predicted by a given method were sorted based on their mutation rate and

575 plotted against the fraction of genes predicted below the given mutation rate

576 *Figure 3 Methodology for identifying novel driver genes.* The figure presents an overview

577 of the different steps involved in our study. Block A (light green frame) shows how our

26

578    classifier is built and is repeated 5 times. Block B (light blue frame) shows random iterations

579    for estimation of hyper-parameters and is repeated 10 times.

580

581

582

583    **TABLES**

584    **Table 1. Classification metrics for training and test set**. Numbers in bold indicate best

585    performances for each metric between TSG and OG. The metrics are standard, and are

586    defined as follows (T stands for True, F for false, P for positives and N for negatives):

587    Accuracy = (TP + TN)/(TP + FP + TN + FN); Precision = TP/(TP+FP); Recall = TP/(TP+FN); F1

588    score is the harmonic mean of Precision and Recall.

| | | Accuracy | F1 score | Precision | Recall |
|---|---|---|---|---|---|
| Training set | OG | 0.86 ± 0.04 | 0.77 ± 0.07 | **0.93 ± 0.04** | 0.67 ± 0.09 |
| | TSG | | **0.90 ± 0.03** | 0.84 ± 0.04 | **0.97 ± 0.01** |
| Test set | OG | 0.76 ± 0.03 | 0.59 ± 0.10 | **0.79 ± 0.12** | 0.50 ± 0.19 |
| | TSG | | **0.83 ± 0.02** | 0.77 ± 0.07 | **0.91 ± 0.07** |

589

590

591

592    *Table 2.* Classification metrics for training and test set using BalancedBagging.

| | | Accuracy | F1 score | Precision | Recall |
|---|---|---|---|---|---|
| Training set | OG | **0.93 ± 0.05** | 0.92 ± 0.06 | 0.86 ± 0.09 | 0.99 ± 0.01 |
| | TSG | | 0.94 ± 0.04 | 1.00 ± 0.01 | 0.90 ± 0.07 |
| Test set | OG | 0.69 ± 0.06 | 0.64 ± 0.06 | 0.56 ± 0.07 | 0.75 ± 0.06 |
| | TSG | | 0.73 ± 0.06 | 0.82 ± 0.04 | 0.65 ± 0.09 |

593    *Table 3.* **Hyper-parameters for each of the folds.** For each cross-validation fold, the most

594    frequent hyper-parameter set is reported. The average accuracy and F1-scores across the

595    different random seed iterations (10, 20, 40, 80, 160, 320) are given along with the standard

596    deviation.

| CV fold | N estimator | Max features | Max depth | Criterion | Accuracy | F1 score | |
|---------|-------------|--------------|-----------|-----------|----------|----------|---|
| | | | | | | OG | TSG |
| 1 | 6 | log2 | 2 | entropy | 0.74 ± 0.02 | 0.55 ± 0.02 | 0.82 ± 0.02 |
| 2 | 5 | log2 | 2 | gini | 0.76 ± 0.03 | 0.60 ± 0.06 | 0.83 ± 0.02 |
| 3 | 10 | log2 | 2 | gini | 0.75 ± 0.03 | 0.60 ± 0.04 | 0.82 ± 0.03 |
| 4 | 5 | log2 | 4 | entropy | 0.76 ± 0.01 | 0.52 ± 0.04 | 0.84 ± 0.01 |
| 5 | 20 | log2 | 4 | gini | 0.79 ± 0.01 | 0.72 ± 0.02 | 0.84 ± 0.01 |

597

598

599

600    *Table 4.* **Kolmogorov-Smirnov statistic for mutation rate distribution of predicted genes**

601    **when compared to training set genes.** KS statistic for the top 60 predicted genes when

602    compared with 208 genes in the training set.

| Method | KS statistic | p-value |
|--------|--------------|---------|
| MutSigCV | 0.774 | <<0.001 |
| Our model | 0.193 | 0.054 |

603     **Table 5. *Driver genes predicted for each of the cancer types.*** *The genes reported showed*

604     *consensus for >4 CV models. Genes in bold did not find similar consensus in the pan-cancer*

605     *predictions. Novel genes are underlined.*

| Primary Tissue | Genes |
|---|---|
| Breast cancer | TP53, **CBFB**, RUNX1, **CDH1**, GATA3, **PTEN**, TBX3 |
| Central nervous system | TP53, **HCN1** |
| Endometrium | **KRAS**, PIK3R1, **PTEN** |
| Hematopoietic | TP53, B2M, **CCND3**, HLA-A |
| Kidney | PBRM1, **VHL**, TP53 |
| Large intestine | TP53, FBXW7, **FAM182A**, SOX9, **AHNAK2**, TCF7L2, ENSG00000121031, **FLT3LG**, **PMEPA1**, **ZFP36L2** |
| Liver | TP53, **ALB**, **KRTAP19-1**, **APOB**, **CD200**, **CRYGD**, **KRTAP24-1**, **OR6N2** |

606

607     **Table 6. *Definitions of mutation categories and the ratio of mutation categories.***

| |
|---|
| Compound mutations = missense + complex + inframe + nonstop − LoFI |
| Loss of Function (LOF) = nonsense + frameshift |
| Damaging = HiFI + MiFI |
| Benign = silent + LoFI |
| $ratio\left(A/B\right)_g = \begin{cases} \dfrac{A_g}{B_g} & if\ B_g \neq 0 \\ 2 * \max(A) & if\ B_g = 0 \end{cases}$ |

608

609     **Table 7. *The ratio-metric* features used in this study for classification**.

| Previously defined in the literature (18 features) | Silent/kb, Total Missense, Total Splicing, Total LOF, Missense/kb, LOF/kb, LOF/Silent, Splicing/Silent, Missense/Silent, LOF/Benign, Splicing/Benign, Missense/Benign, average Polyphen2 score, LOF/Total, Missense/Total, Splicing/Total, LOF/Missense, Missense entropy |
|---|---|
| Defined in this paper (19 features) | HiFI/LoFI, HiFI/Benign, MiFI/kb, Nonstop/kb, Inframe/kb, Complex/kb, Compound/Benign, Compound/kB, Damaging/kb, Damaging/Benign, Damaging/LoFI, High Missense frequency, Frameshift entropy, High Frameshift frequency, Splicing entropy, High Splicing frequency, Nonsense entropy, High Nonsense frequency, Total MiFI |

610

611  **DECLARATIONS**

612  **Ethics approval and consent to participate**

613  Not applicable

614  **Consent for publication**

615  Not applicable

616  **Availability of data and materials**

617  Data for this analysis was downloaded from COSMIC (v79)

618  The processed data and codes are available in GitHub.

619  (https://github.com/RamanLab/IdentifyTSGOG)

620  **Competing interests**

621  The authors declare that they have no competing interests

622  **Funding**

623  This work was supported by Department of Biotechnology, Government of India (DBT)

624  (BT/PR16710/BID/7/680/2016), IIT Madras, Initiative for Biological Systems Engineering

625  (IBSE) and Robert Bosch Center for Data Science and Artificial Intelligence (RBC-DSAI).

626 **Authors' contributions**

627 MS, RR and KR conceived and designed the study. MS, RR, and KR were involved in the

628 analysis and interpretation of data. MS, RR and KR drafted the manuscript. The study was

629 supervised by RR and KR. All authors read and approved the final manuscript.

630 **Acknowledgements**

631 Not applicable

632 **ADDITIONAL FILE INFORMATION**

633 Additional file 1: Table S1.

634 List of features and their ranks for each of the models and the calculated average rank.

635 (XLSX 10.8 kB)

636 Additional file 2: Figure S1.

637 Distribution of features across the two classes for all the other features not included in

638 Figure 1. (PDF 1.39 MB)

639

Distribution of features

Legend:
- Training
- Our predictions
- MutSigCV

x-axis: Log mutation rate

y-axis: Fraction of genes predicted below mutation rate

Coding mutation data

Feature matrix

Labelled genes  ⟷  Unlabelled genes

5-fold cross validation

| 1 | 2 | 3 | 4 | 5 |

**A**

| 1 | 2 | 3 | 4 | 5 |

Repeat 10x

**B**

Assign random seed

Grid search with 5-fold cross validation for parameter estimation

Get hyper-parameter set

Select most frequent hyper-parameter set

Build classifier

Training and test metrics

Repeat 5x

Average training and test metrics

Classifiers

| Genes predicted by 5 models | Genes predicted by 4 models | Genes predicted by 3 models |