1  **Predicting Response to Platin Chemotherapy Agents with Biochemically-inspired**
2  **Machine Learning**

3  Eliseos J. Mucaki[1], Jonathan Z.L. Zhao[1], Dan Lizotte[2,3], and [§]Peter K. Rogan[1,2,3,4,5]

4

5  **Running Title:**

6  Predicting Responses to Platin Drugs by Machine Learning

7

8  **Author Affiliations**

9  [1]Department of Biochemistry, Schulich School of Medicine and Dentistry, Western University,

10  London, Canada, N6A 2C1

11  [2]Department of Computer Science, Faculty of Science, Western University, London, Canada, N6A

12  2C1

13  [3]Department of Epidemiology & Biostatistics, Faculty of Science, Western University, London,

14  Canada, N6A 2C1

15  [4]Cytognomix Inc. London, Canada N5X 3X5

16  [5]Department of Oncology, Schulich School of Medicine and Dentistry, Western University,

17  London, Canada, N6A 2C1

18

19  Author Emails: emucaki@uwo.ca, jzhao293@uwo.ca, dlizotte@uwo.ca, progan@uwo.ca

20  [§]**Correspondence to:** Peter K. Rogan (progan@uwo.ca), Department of Biochemistry, Schulich
21  School of Medicine and Dentistry, Western University, London, Ontario, Canada, N6A 2C1. 1
22  (519) 661-4255.

23 **ABSTRACT.**

24 Selection of effective genes that accurately predict chemotherapy response could

25 improve cancer outcomes. We compare optimized gene signatures for cisplatin,

26 carboplatin, and oxaliplatin response in the same cell lines, and respectively validate

27 each with cancer patient data. Supervised support vector machine learning was used to

28 derive gene sets whose expression was related to cell line $GI_{50}$ values by backwards

29 feature selection with cross-validation. Specific genes and functional pathways

30 distinguishing sensitive from resistant cell lines are identified by contrasting signatures

31 obtained at extreme vs. median $GI_{50}$ thresholds. Ensembles of gene signatures at

32 different thresholds are combined to reduce dependence on specific $GI_{50}$ values for

33 predicting drug response. The most accurate models for each platin are: cisplatin:

34 *BARD1*, *BCL2*, *BCL2L1*, *CDKN2C*, *FAAP24*, *FEN1*, *MAP3K1*, *MAPK13*, *MAPK3*,

35 *NFKB1*, *NFKB2*, *SLC22A5*, *SLC31A2*, *TLR4*, *TWIST1*; carboplatin: *AKT1*, *EIF3K*,

36 *ERCC1*, *GNGT1*, *GSR*, *MTHFR*, *NEDD4L*, *NLRP1*, *NRAS*, *RAF1*, *SGK1*, *TIGD1*, *TP53*,

37 *VEGFB*, *VEGFC;* oxaliplatin: *BRAF*, *FCGR2A*, *IGF1*, *MSH2*, *NAGK*, *NFE2L2*, *NQO1*,

38 *PANK3*, *SLC47A1*, *SLCO1B1*, *UGT1A1*. TCGA bladder, ovarian and colorectal cancer

39 patients were used to test cisplatin, carboplatin and oxaliplatin signatures (respectively),

40 resulting in 71.0%, 60.2% and 54.5% accuracy in predicting disease recurrence and

41 59%, 61% and 72% accuracy in predicting remission. One cisplatin signature predicted

42 100% of recurrence in non-smoking bladder cancer patients (57% disease-free; N=19),

43 and 79% recurrence in smokers (62% disease-free; N=35). This approach should be

44 adaptable to other studies of chemotherapy response, independent of drug or cancer

45 types.

46

47     **KEY WORDS**. Chemotherapy response, support vector machines, gene signatures,

48     cancer, cisplatin, oxaliplatin, carboplatin, machine learning, bladder cancer, breast

49     cancer, ovarian cancer

50    **INTRODUCTION**

51    Chemotherapy regimens are selected based on overall outcomes for specific

52    types and subtypes of cancer pathology, progression to metastasis, other high-risk

53    indications, and prognosis[1,2], and variability in tumor resistance has led to tiered

54    sequential strategies for selection of agents based on their overall efficacy[3]. We and

55    others have developed machine learning (ML)-based gene signatures aimed at

56    predicting response to specific chemotherapeutic agents and minimizing

57    chemoresistance based on inhibition of growth or drug targets ($GI_{50}$ or $IC_{50}$)[4–6]. In this

58    study, we present integrated ML models of platin drug responses (cis-, carbo- and

59    oxaliplatin). Previous studies have reviewed the genes[7], gene products[8] and specific

60    individual pathways that are activated and repressed by drugs[9], but lack comprehensive

61    models of the global cellular response to drugs. We use integrated ML-based signatures

62    based on expression of multiple genes to predict key responses to each of these platin

63    agents, for the first time, at different resistance levels.

64    Cisplatin, carboplatin and oxaliplatin are each widely prescribed compounds for

65    their antineoplastic effects. While each contains platinum to form adducts with tumour

66    DNA, their effectiveness differs for specific types of cancers, such as bladder (cisplatin),

67    ovarian (cisplatin and carboplatin) and colorectal cancer (oxaliplatin). Carboplatin differs

68    in structure from cisplatin, exchanging the latter's dichloride ligands with a CBDCA

69    (cyclobutane dicarboxylic acid) group, while oxaliplatin is paired with both a DACH

70    (diaminocyclohexane) ligand and a bidentate oxalate group. These chelating ligands

71    have greater stability and solubility to aqueous solutions, which lead to differences in

72    drug toxicity compared to cisplatin[10]. Oxaliplatin can be up to two times as cytotoxic as

73    cisplatin, but it forms fewer DNA adducts[11]. The large hydrophobic DACH ligand which

74    overlaps the major groove is thought to prevent binding of certain DNA repair enzymes

4

75  such as the POL polymerases, and may contribute to the low cross-resistance between

76  oxaliplatin and cisplatin and carboplatin[10]. While all three drugs can enter the cell via

77  copper transporters, organic cation transporters are oxaliplatin-specific and likely play a

78  role in its efficacy in colorectal cancer (CRC) cells where these transporters are

79  commonly overexpressed[7]. Oxaliplatin specifically plays a role in interfering with both

80  DNA and RNA synthesis, unlike cisplatin which only infers with DNA[12]. It is these

81  intrinsic properties between the platinum drugs which lead to differences in their activity

82  and resistance profiles, despite their similar mode of action.

83      We derived gene signatures to predict drug response at different sensitivity and

84  resistance levels for each of these agents. We and others have used supervised

85  learning algorithms, including random forest models[13]; support vector machine (SVM)

86  models[6]; neural networks[14]; and linear regression models[5] to make these predictions.

87  Pathway and network analysis of gene expression have been used to indicate hundreds

88  of genes potentially up- and down-regulated upon cisplatin treatment[15]. Cisplatin-specific

89  gene signatures have been developed with integrative approaches such as elastic net

90  regression using inferred pathway activity of bladder cancer cell line data[16]. These

91  methods have implicated genes that have not been described previously. Supervised ML

92  with biochemically-relevant genes has also been useful for predicting drug response[6]. A

93  concern with each of these ML approaches is that an insufficient number of samples

94  coupled to a large number of features, i.e. gene expression changes, in each sample

95  can result in overfitting of the model affecting its generalizability with other sources of

96  data[17]. We therefore reduce the number of dimensions by selecting genes biologically

97  relevant to the drugs under observation[6,17]. Additional selection criteria are necessary

98  when the number of genes implicated in peer-reviewed reports is still prohibitively large

99  compared to sample size.

100    Biochemically-inspired gene signatures have shown good performance in

101    predicting treatment response. A paclitaxel ML signature based on tumor gene

102    expression had a higher success predicting the pathological complete response rate

103    (pCR [18]) for sensitive patients (84% of patients with no / minimal residual disease) than

104    models based on differential gene expression (GE) analysis[6]. For gemcitabine, a

105    signature derived from both expression and copy number (CN) data from breast cancer

106    cell lines was derived, and subsequently applied to analysis of nucleic acids from patient

107    archival material. Multiple other outcome measures used to validate gene signatures

108    include prognosis[5], Miller-Payne response[19], and disease recurrence. Binary SVM

109    classifiers based on discrete time thresholds have been used to classify continuous

110    outcome measures such as prognosis and recurrence. By contrast, pCR is simpler to

111    interpret with binary SVM models. Nevertheless, differences in clinical recurrence have

112    been noted between patients demonstrated with pCR and those who do not exhibit

113    disease pathology[18]. This source of variability in defining patient response can confound

114    transferability of SVM models between different datasets.

115    We apply biochemically-inspired ML to predict and compare the cellular and

116    patient responses to cisplatin, carboplatin and oxaliplatin. We train models for

117    classification of platin resistance with cancer cell line data and validate with patient GE

118    and outcome data. Our previous gene signatures were based on median $GI_{50}$ for each

119    drug[6]. This has been a necessary compromise, however in this study we consider

120    signatures that differ at the highest vs. the lowest levels of drug resistance. A series of

121    gene signatures are derived by shifting the $GI_{50}$ thresholds that distinguish sensitivity

122    from resistance. The frequency of genes selected at median vs. extreme thresholds

123    highlights pathways that define these responses among different patient subsets.

124

125    **RESULTS**

126    **Selection of Platin Drug Related Genes**

127    We documented genes in the peer-reviewed literature associated with drug

128    effectiveness or response (Supplemental References). For cisplatin, carboplatin and

129    oxaliplatin, this implicated 178, 90, and 288 genes, respectively (Suppl. Table S1).

130    Multiple factor analysis (MFA) was used to determine which genes were correlated to

131    $GI_{50}$ in breast cancer cell lines through either GE and/or CN[13], significantly reducing the

132    sizes of the gene sets for cisplatin (N=39), carboplatin (N=28), and oxaliplatin (N=55).

133    Genes with significant relationships to $GI_{50}$ and direction of correlation (positive or

134    inverse) are indicated in Figure 1. The diverse functions of these genes included

135    apoptosis, DNA repair, transcription, cell growth, metabolism, immune system, signal

136    transduction and membrane transport. Analysis of $IC_{50}$ and gene expression levels for

137    cisplatin-treated bladder cancer cell lines confirmed these relationships evident from $GI_{50}$

138    values of different breast cancer lines. $IC_{50}$ values were related to GE for *CFLAR*, *FEN1*,

139    *MAPK3*, *MSH2*, *NFKB1*, *PNKP*, *PRKAA2*, and *PRKCA*[20]. Similarly, separate bladder cell

140    line $IC_{50}$ values from the Genomics of Drug Sensitivity in Cancer project

141    (http://www.cancerrxgene.org; N=17) were correlated with GE for *CFLAR*, *FEN1*, and

142    *NFKB1*, in addition to *ATP7B*, *BARD1*, *MAP3K1*, *NFKB2*, *SLC31A2* and *SNAI1*.

143          We performed MFA on the $GI_{50}$ values for cisplatin, carboplatin and oxaliplatin,

144    without consideration of either GE or CN. Responses to cis- and carboplatin were

145    directly correlated (a 6.2º separation between vectors), but neither was related to the

146    oxaliplatin response (Figure 2). Previous studies have shown that cisplatin-resistant cell

147    lines are generally sensitive to oxaliplatin[21–23].

7

148    SVM-based signatures were initially derived for each platin drug from breast

149    cancer cell line GE data. A 13-gene signature for cisplatin at the median $GI_{50}$ threshold

150    (5.2% misclassification rate) consisted of *BARD1, BCL2L1, FAAP24, CFLAR, MAP3K1,*

151    *MAPK3, NFKB1, POLQ, PRKAA2, SLC22A5, SLC31A2, TLR4,* and *TWIST1*. A similarly

152    derived carboplatin signature included *AKT1*, *ATP7B*, *EGF*, *EIF3I*, *ERCC1*, *GNGT1*,

153    *HRAS*, *MTR*, *NRAS*, *OPRM1*, *RAD50*, *RAF1*, *SCN10A*, *SGK1*, *TIGD1*, *TP53*, and

154    *VEGFB* (10.4% misclassification). For oxaliplatin, the final SVM model consisted of

155    *AGXT*, *APOBEC2*, *BRAF*, *CLCN6*, *FCGR2A*, *IGF1*, *MPO*, *MSH2*, *NAGK*, *NAT2*,

156    *NFE2L2*, *NOTCH1*, *PANK3*, *PRSS1*, and *UGT1A1* (2.1% misclassification). A cisplatin

157    SVM generated from 17 bladder cancer cell lines in cancerRxgene resulted in 2 equally

158    accurate signatures (with 11.8% misclassification) consisting of either *PNKP* and

159    *PRKCA* or *ATP7B*, *CFLAR*, *FEN1*, *MAPK3*, *NFKB1* and *SLC22A11*. These models were

160    not useful for predicting patient outcomes due to the limited size of the training set.

161    **$GI_{50}$-Threshold Independent Modeling**

162    In our previous studies, we set median $GI_{50}$ value as the threshold to

163    distinguished drug resistance and sensitivity[5,6]. An important question is whether the

164    genes contributing to drug response are consistent among different cell lines, each with

165    their own unique $GI_{50}$ values. Different ML models were obtained by shifting the $GI_{50}$

166    threshold, which changed the labels of resistant vs. sensitive cell lines. After feature

167    selection, the compositions of the corresponding gene signatures for each threshold

168    were compared. Finally, ensemble averaging of all of these optimized Gaussian SVM

169    models derived for different $GI_{50}$ thresholds was used to create a threshold-independent

170    ML-based signature.

8

171    Kinase (*MAPK3*, *MAP3K1*) genes and apoptotic family members (*BCL2*,

172    *BCL2L1*) were most the common in the cisplatin signatures at different $GI_{50}$ thresholds,

173    with consistent representation of error-prone and base-excision DNA repair genes as

174    well (Figure 3A; Supplementary Table S2A). The kinases are more concentrated in

175    signatures with lower drug sensitivity thresholds, whereas *BCL2* and *BCL2L1* are more

176    ubiquitous at all levels. The error prone polymerases, *POLD1* and *POLQ,* are more

177    frequent in models with lower sensitivity thresholds, while the flap endonuclease *FEN1*

178    tends to be present at high levels of resistance. Thresholded models for carboplatin-

179    related genes commonly contained the apoptotic family member *AKT1*, transcription

180    regulation genes *ETS2* and *TP53*, as well as cell growth factors *VEGFB* and *VEGFC*,

181    although the latter was less common at lower sensitivity thresholds (Figure 3B).

182    Common oxaliplatin-related genes included transporters *SLCO1B1* and *GRTP1* (but not

183    *SLC47A1*), transcription genes *NFE2L2*, *PARP15* and *CLCN6*, as well as multiple

184    metabolism-related genes (Figure 3C).

185    SVM models were also derived using the cisplatin and/or carboplatin-treated

186    TCGA (The Cancer Genome Atlas) bladder urothelial carcinoma patients, using post-

187    treatment time to relapse as a surrogate criterion for different $GI_{50}$ resistance thresholds

188    (as performed in Mucaki *et al.* [2017][24]; Supplementary Table S3). Similar trends to cell

189    line SVMs are apparent: *POLQ* is frequently included in models with recurrence

190    threshold of longer duration, while *FEN1* is a marker of resistance, when time to relapse

191    is shorter. However *BCL2*, which is present in a majority of breast cancer cell line SVMs,

192    is present in only one model derived from TCGA data. Similarly, *MSH2* was rarely

193    selected using cell lines, yet appears in nearly all patient derived SVMs with > 1 year

194    recurrence.

195          $GI_{50}$ thresholded ML models were also derived using the log-loss function, which

196    penalizes false classifications, whose value ranges from zero (or completely accurate),

197    to 1 (or completely inaccurate; Supplementary Table S4). The overall distribution of

198    genes across $GI_{50}$ thresholds has many distinct similarities with the models derived by

199    misclassification. For both sets of cisplatin models, *BCL2*, *BCL2L1* and *FEN1* are

200    common in low-to-moderate $GI_{50}$ thresholds, while *NFKB1* is enriched at high thresholds

201    (Figure 3A; Suppl. Figure 1A). For carboplatin, *AKT1*, *VEGFB* and *VEGFC* are similarly

202    distributed across $GI_{50}$ thresholds with both methods, although *VEGFB* is less dense

203    with log-loss models at low $GI_{50}$ values (Figure 3B; Suppl. Figure 1B). In both sets of

204    models for oxaliplatin, *SIAE* and *SLC47A1* show a high density across all $GI_{50}$

205    thresholds, while *ABCG2* shows low density across each (<50% inclusion; Figure 3C

206    and Suppl. Figure 1C). There are some distinct differences. *EGF* and *ERCC1* were

207    selected at a greater frequency at a moderate carboplatin $GI_{50}$ with log-loss rather than

208    misclassification. Similarly, the following oxaliplatin genes were selected considerably

209    more often when using log-loss: *APOBEC2*, *HLA-B*, *LTA*, and *MPO*. Therefore, while the

210    misclassification and log-loss based models are not interchangeable, the models are

211    overall quite similar.

212          Log-loss models were initially constructed either by (a) a modified version of the

213    misclassification-based method, or (b) using the BFS software described in Zhao *et al.*

214    (2018)[25]. Multiple signatures with low log-loss values can have different compositions,

215    consistent with the possibility that there may be various diverse gene combinations that

216    can give rise to signatures with satisfactory performance. However, these signatures

217    often contain a larger number of gene features than the misclassification based

218    signatures, and raised concerns that they might be more prone to overfitting. The log-

219    loss minimized models generated by both methods had comparable compositions. The

220    median $GI_{50}$ thresholded cisplatin model generated by the log-loss modified software

221    [*ATP7B*, *BCL2L1*, *CDKN2C*, *CFLAR*, *ERCC2*, *ERCC6*, *FAAP24*, *FOS*, *GSTO1*, *GSTP1*,

222    *MAP3K1*, *MAPK13*, *MAPK3*, *MSH2*, *MT2A*, *PNKP*, *POLD1*, *POLQ*, *PRKAA2*, *PRKCA*,

223    *PRKCB*, *SLC22A5*, *SLC31A2*, *SNAI1*, *TLR4*, *TP63*] shares 15/19 genes with the

224    signature generated by the BFS software[25] [*ATP7B*, *BARD1*, *BCL2*, *BCL2L1*, *ERCC2*,

225    *FAAP24*, *FEN1*, *FOS*, *MAP3K1*, *MAPK13*, *MAPK3*, *MSH2*, *MT2A*, *NFKB1*, *PNKP*,

226    *POLQ*, *PRKCB*, *SLC22A5*, *SNAI1*]).

227    **Traditional Model Validation against Cancer Patient Data**

228    $GI_{50}$-thresholded models for each platin drug, generated with the breast cancer

229    cell line data, produced 70 cisplatin, 83 carboplatin, and 83 oxaliplatin SVM models,

230    respectively. Each model was validated using available platin-treated patient datasets[26–

231    30]. The chemotherapy response metadata differed between studies. Als *et al.*[29] reported

232    survival post-treatment, whereas Tsuji *et al.*[30] categorized patients as responders and

233    non-responders. TCGA provided two different measures which were used to assess

234    predictive accuracy in our models – chemotherapy response and disease-free survival.

235    Accuracy is similar using either measure (Supplementary Table S5A); however

236    recurrence and disease-free survival was used as the primary measure of response as it

237    was more often recorded in the TCGA data sets tested. Patients from Als *et al.* with a ≥

238    5 year survival post-treatment were labeled as sensitive to treatment. The differences

239    between these metadata may, in part, account for the differences in the prediction

240    accuracy of the thresholded SVM models.

241    At higher resistance thresholds for any platin drug (low $GI_{50}$), where more cell

242    lines are labeled sensitive, the positive class (disease-free survival) is correctly

243    classified, while the negative class (recurrence) is highly misclassified (Suppl. Figures 2

244    and 3). The reverse is true for models built using lower resistance thresholds (high $GI_{50}$).

245    We therefore state SVMs generated at these extreme thresholds are not very useful at

246    predicting patient data. When used to predict recurrence in the TCGA datasets,

247    sensitivity and specificity appears to be maximized in models where the $GI_{50}$ threshold

248    for resistance was set near (but not necessarily at) the median (Suppl. Figure 2; Suppl.

249    Tables S5A to 5C). While this pattern holds true for Tsuji *et al.*[30], oxaliplatin models

250    where $GI_{50}$ thresholds were set above the median could better separate primary and

251    metastatic CRC patients (best model predicting 92.6% metastatic and 60.7% primary

252    cancers; Suppl. Table S5C). While less consistent, cisplatin models generated with

253    thresholds above median $GI_{50}$ performed better when evaluating the Als *et al.*[29] patient

254    dataset (Suppl. Figure 3).

255        Models were further evaluated for their accuracy in TCGA patients using various

256    recurrence times post-treatment to classify resistant and sensitive patients (0.5 - 5 years;

257    Supplemental Table S6A-C). The best performing cisplatin model (hereby identified as

258    **Cis1**; Table 1) was able to accurately predict 71.0% of bladder cancer patients who

259    recurred after 18 mo. (N=31; 58.5% accurate for disease-free patients [N=41]).

260    Response of TCGA bladder patients treated with carboplatin (without cisplatin; N=19)

261    were best predicted by **Cis12** two years post-treatment (80% accurate for responding

262    patients [N=5]; 93% for recurrent patients [N=14]). The best performing carboplatin

263    model (designated **Car1** [Table 1]) predicted recurrence of ovarian cancer after 4 years

264    at an accuracy of 60.2% (N=302; 61.0% accurate for disease-free patients [N=108]).

265    These models were also used to test TCGA bladder cancer patients treated carboplatin

266    but not cisplatin (N=19), of which the best performing model (**Car73**) was 84% accurate

267    for patients after 1 year of treatment (100% for responding patients [N=11]; 62.5%

268    accuracy for recurrent [N=8]). Two additional carboplatin models are tied for overall

12

269    accuracy (84%; **Car9** and **Car51**), but more successfully predict non-responsive patients

270    (87.5%; 82% accuracy for responding patients). These three models share four genes:

271    *AKT1, ETS2, GNGT1*, and *VEGFB.* For oxaliplatin, the best performing model

272    (designated **Oxa1** [Table 1]) accurately predicted 71.6% of the disease-free TCGA CRC

273    patients after one year (N=88; 54.5% accuracy predicting recurrence [N=11]). These

274    models (based on gene expression measured by Affymetrix Gene Chip Human Exon 1.0

275    ST arrays), TCGA sample expression data, as well as SVMs based on bladder cell line

276    data (based on expression measured by Affymetrix U133A microarray), were added to

277    the online web-based SVM calculator (http://chemotherapy.cytognomix.com; introduced

278    in Dorman *et al.* [2016][6]) to predict platin response.

279        To evaluate the consistency in the response prediction of TCGA bladder cancer

280    patients treated with cisplatin, the distance from the hyperplane for all SVMs generated

281    were plotted for each patient with a short recurrence time (<6 mo., N=10; Supplementary

282    Figure 4). Despite showing similar levels of resistance to treatment, patterns differed

283    between patients. While these patients would be expected to be indicated as highly

284    cisplatin resistant (hyperplane distance < 0), two patients (TCGA-XF-A9SU and TCGA-

285    FJ-A871) were predicted sensitive across nearly all SVM models. Similar variation was

286    also seen in patients with either a long recurrence time (>4 years) or no recurrence at all

287    after 6 years (Suppl. Figure 5).

288        Threshold independent models were generated for each individual platin drug at

289    different $GI_{50}$ thresholds through ensemble ML, which involves the averaging of

290    hyperplane distances for each model to generate a composite score for each TCGA

291    patient tested. Hyperplane distances across all 70 cisplatin models were similar, with a

292    mean score of -0.22 and a standard deviation of 3.5 hyperplane units (hu) across the set

293    of patient data. The ensemble model classified disease-free bladder cancer patients

13

294    with 59% accuracy and those with recurrent disease with 47% accuracy. Limiting

295    ensemble averaging to only cisplatin models generated at a moderate $GI_{50}$ threshold

296    (ranging from 5.10 to 5.50) did not significantly improve accuracy (44% for disease-free

297    and 66% for recurrent patients; Suppl. Table S7A). For carboplatin, ensemble ML did not

298    produce significantly better predictions than random, regardless of the $GI_{50}$ threshold

299    interval selected (Suppl. Table S7B) or the similar mean hyperplane distances (-0.11 +/-

300    3.9 hu). For oxaliplatin, the ensemble ML model (mean = -0.12 +/- 2.7 hu) was most

301    accurate after 1 year (60% accuracy for disease-free and 73% for recurrent patients;

302    Suppl. Table S7C). As in cisplatin, limiting this analysis to oxaliplatin SVM models with

303    moderate $GI_{50}$ thresholds did not significantly increase accuracy.

304    To determine the impact of individual genes on overall model accuracy, each

305    gene within every SVM model was excluded, and model accuracy was reassessed

306    (Supplementary Tables S2A; S2B and S2C contain cis-, carbo- and oxaliplatin models,

307    respectively). Genes which consistently significantly increase misclassification

308    (averaging > 16% increase) in moderate threshold SVMs ($GI_{50}$ thresholds set from 5.1 to

309    5.5) include *ERCC2*, *POLD1*, *BARD1*, *BCL2*, *PRKCA* and *PRKCB*. *ERCC2* and *POLD1*

310    perform critical functions in nucleotide and base excision repair, respectively. *PRKCA*

311    and *PRKCB* are paralogs with significant roles in signal transduction. *BARD1* has been

312    shown to reduce apoptotic *BCL2* in the mitochondria[31], and has a key role in genomic

313    stability through its association with *BRCA1*. Genes with a high variance in increased

314    misclassification between different models include *NFKB1*, *NFKB2*, *TWIST1*, *TP63*,

315    *PRKAA2*, and *MSH2*. The variance of these genes may be due to epistatic interactions

316    with other biological components, including the other genes in the SVM. For example,

317    *NFKB1* and *NFKB2* are jointly included in 7 SVMs generated at a moderate $GI_{50}$

318    threshold. There is evidence of possible epistasis in that the removal of either of these

14

319   genes, but not necessary both, will have a large impact in model misclassification rates

320   (≥ 18.0% increase). The misclassification variance of *NFKB1* with *NFKB2*, is significantly

321   lower than in SVM models lacking *NFKB2*.

322        To further evaluate the predictive capability of the misclassification-based gene

323   signatures, k-fold cross-validation of the cisplatin, carboplatin and oxaliplatin models

324   were performed on TCGA bladder, ovarian and colorectal cancer patient data,

325   respectively. Patients were evenly distributed in 5 groups with an equal (or near-equal)

326   ratio of disease-free and recurrent patients. The majority of the cisplatin models showed

327   an overall accuracy > 50%. The cisplatin model which performed best under the k-fold

328   analysis (6-resistance level; *BARD1*, *BCL2*, *BCL2L1*, *PRKAA2*, *PRKCA*, *PRKCB*,

329   *TWIST1*) showed an overall accuracy of 71.2% (84.4% accurate for sensitive and 53.9%

330   accurate for resistant patients). The accuracy of the carboplatin and oxaliplatin models

331   did not exceed 60%. In general, traditional validation outperformed the k-fold validation

332   results.

333   **Predicting cisplatin response in patients based on smoking history**

334        Tobacco smoking is known as the highest risk factor for the development of

335   bladder cancer[32]. We therefore subdivided the patients based on their smoking history

336   and tested the thresholded models (Supplementary Tables S8 and S9). When testing

337   patients who were lifelong non-smokers, the prediction accuracy of **Cis1** predicted all

338   non-smoking patients who were recurrent after 18 months as cisplatin-resistant (N=5).

339   Prediction accuracy for disease-free patients was 57.1% (N=14). Another model (**Cis18**;

340   Suppl. Table S8) had performed equally as well for non-smokers, and these two models

341   share 7 genes: *BCL2*, *BCL2L1*, *FAAP24*, *MAP3K1*, *MAPK13*, *MAPK3*, and *SLC31A2*.

342   Threshold independent analysis predicted disease-free equally well, but recurrence was

15

343     less accurate (66.7%). Note that non-smokers make up a small subset of the patients

344     tested (N=19). Threshold-independent prediction of recurrence in patients with a

345     smoking history was 46% accurate (N=13), while disease-free patients were correctly

346     predicted at a rate of 58% (N=19). Recurrence in these patients was best predicted by a

347     model built at the median $GI_{50}$ threshold (**Cis2**). Accuracy improved for both disease-free

348     (57.7% -> 61.9%) and recurrent patients (76.0% -> 78.6%) when excluding patients who

349     quit smoking more than 15 years before diagnosis. Genes in this SVM which are not

350     present in the two models which performed well for non-smokers include *CFLAR* and

351     *PRKAA2.*

352     Tobacco smoking has a significant impact on cytosine methylation levels in the

353     genome[33]. CpG island methylation has been associated with smoking pack years in a

354     subset of the TCGA bladder urothelial carcinoma patients[26]. We suspected that the level

355     of methylation measured in the SVMs which performed best for smoking and non-

356     smoking patients might differ, and with possible concomitant effects on GE. When

357     ranking each gene from **Cis1** by highest methylation and GE, 88 of 1080 patient: gene

358     combinations showed the expected inverse correlation between methylation levels and

359     GE (i.e. high methylation and low GE). Inverse correlation of methylation and GE was

360     more common than direct correlation (i.e. high methylation and high GE; N=17).

361     However, direct correlation was more common in patients with a recent smoking history

362     (70.5%). This pattern was also observed for **Cis2**, which best predicted recurrence in

363     smokers. In cases where methylation and GE are directly correlated, we propose that

364     smoking may alter expression by other effects, e.g. mutagenic, rather solely than by

365     epigenetic inactivation through methylation.

366     To determine which genes in these models led to discordant predictions of

367     patient outcome, we conducted a bioinformatic analysis in which the expression of each

16

368  signature gene was gradually altered until the misclassification was corrected. If the GE

369  value required to cross this threshold exceeded ≥ 3-fold the highest/lowest expression of

370  that gene, it was interpreted as a minor contributor to the prediction. Genes which could

371  not correct a discordant prediction included *PRKAA2*, *NFKB1*, *NFKB2* and *TWIST1*.

372  Significant genes which, when altered, corrected discordant predictions included

373  *MAP3K1*, *MAPK3*, *SLC22A5* and *SLC31A2*. Altering *BCL2L1* expression was more

374  likely to correct the discordant predictions of **Cis1** (4 out of 5) than with **Cis2** (2 out of 4).

375  **DISCUSSION**

376  Using gene expression signatures, we derived both $GI_{50}$ threshold-dependent

377  and -independent ML models which predict the chemotherapy responses for cisplatin,

378  carboplatin and oxaliplatin, respectively. The cisplatin model **Cis1** (Supplementary Table

379  S6A) most accurately predicted response in bladder cancer patients after 18 months,

380  and **Car1** (Suppl. Table S6B) best predicted response in ovarian cancer patients after 4

381  years. **Oxa1** (Suppl. Table S6C) more accurately predicted disease-free patients than

382  recurrent disease at the one year treatment threshold. The thresholds which best

383  represented time-to-recurrence differed between the platin drugs in each cancer type.

384  Cisplatin gene signatures had noticeably improved performance when smoking history

385  was taken into account.

386  The three platin drugs produce distinctly different gene signature models. Initial

387  gene sets exhibited some overlap between platin drugs (N=67 between any two platins),

388  but very few of these were correlated by MFA of $GI_{50}$ with multiple platin drugs (*ATP7B*,

389  *BCL2* and *MSH2*). Signature genes common to multiple platin drugs whose expression

390  was correlated with cisplatin $GI_{50}$ values but not with carboplatin and/or oxaliplatin values

391  include  *BCL2L1, GSTP1, MAP3K1, MAPK3, MT1A,* and *MT2*. Similarly,  genes

17

392   correlating only to carboplatin $GI_{50}$ included *AKT1, EGF, ERCC1, KRAS, LIG3, MTHFR,*

393   *MTR, RAD50, TP53*, while genes correlating to only oxaliplatin $GI_{50}$ included *ATM,*

394   *BCL2, CLCN6, ERCC2, ERCC6,* and *UGT1A1*. Despite the close similarity between

395   cisplatin and carboplatin $GI_{50}$ response (see Figure 2), only one gene (*ATP7B*) was

396   related by MFA to $GI_{50}$ levels of both drugs. *BCL2* and *MSH2* correlated with both

397   cisplatin and oxaliplatin $GI_{50}$ (*BCL2* did not correlate with carboplatin $GI_{50}$). The increase

398   in misclassification caused by the elimination of *MSH2* from any SVM model in which it

399   was present was significant; for example, misclassification of **Cis14** and **Oxa21** (Table

400   1) were increased by 28.2% and 19.1%, respectively (Suppl. Tables S2A and S2C).

401   These differences may reflect the spectrum of activity, sensitivity, and toxicity of these

402   signature genes[21–23,34,35].

403       Previous validation of patient data for other drugs validated with other datasets[6,24]

404   using biochemically inspired machine learning have had better performance than those

405   reported here. We investigated the possibility that disease and molecular heterogeneity

406   in platin-treated patients may have affected the accuracy of our results. Model

407   predictions were reevaluated after stratifying clinical features such as time-to-disease

408   recurrence, cancer stage, and metastatic lymph node count. Breast cancer patients with

409   advanced disease (stage III and IV) were analyzed separately from those with earlier

410   stage diagnoses (stage I and II). Cisplatin model **Cis1** performed best on stage IV

411   patients (overall accuracy 72.4% at a 2 year recurrence threshold), while **Oxa1** similarly

412   performed best in predicting late stage cancers (74.5% accurate for stage III and 71.4%

413   accurate for stage IV at a 2 year recurrence threshold). **Cis5** was also more accurate for

414   later stage cancer patients (72.4% overall accuracy at 18 months). The accuracies of

415   models were similar across all stages (e.g. **Car1** ranged from 58-74%). Cisplatin-treated,

416   TCGA bladder cancer patients and oxaliplatin-treated TCGA colorectal cancer patients

417  were also stratified by Lymph Node status (N0, N1, and N2 [bladder cancer patient data

418  set comprised of only two N3 patients, which were included with the analysis of N2

419  patients; N3 was not represented in colorectal cancer]). In TCGA bladder cancer

420  patients, **Cis1** exhibited ~60% accuracy across all categories; however it performed

421  better in sensitive N0 and N1 patients relative to N2. **Cis2** was less accurate for N2

422  patients than for N0 and N1. Sensitive N2 patients were more likely to be misclassified

423  (<40%) than relapsed N2 patients. In TCGA colorectal cancer patients, **Oxa1** was 88%

424  accurate in N2 patients (95% accurate for sensitive N2 patients [n=19], and 67%

425  accurate for relapsed N2 patients [n=6]). Oxaliplatin models were less accurate for N1

426  patients compared to N0 and N2. Thus, heterogeneity in disease stage as well as

427  metastatic phenotypes adversely confounds the overall accuracies of our predictions.

428  Gene signature models derived from cell lines and tested on patients differ in

429  their outcome measures. The exact $GI_{50}$ cell line threshold that is most predictive of

430  patient outcome is not known, and different groups use different methods to discretize

431  $GI_{50}$ values[36,37]. Therefore, we developed ML models for platin drugs which predict drug

432  response without relying on arbitrary $GI_{50}$ thresholds. For cisplatin, SVM ensemble

433  averaging generated on different resistance thresholds shows a small increase in

434  accuracy over most models, better representing the sensitive, disease-free class (59%

435  accuracy). Interestingly, ensemble averaging of only the models built using a moderate

436  $GI_{50}$ thresholds yielded results which better represented the resistance class. This result

437  closer matches the accuracy of **Cis1**, and may be due to **Cis1** having a greater overall

438  impact on the ensemble prediction. When limiting ensemble averaging to only those

439  models with the highest area under the curve (AUC) at each resistance threshold,

440  differences in predictions were negligible. Ensemble ML can potentially avoid problems

441  with poor performance and overfitting by combining models that individually perform

442  slightly better than chance[38].

443    It is difficult to reconcile gene signatures without features known to be related to

444    chemoresistance with tumor biology. Our thresholding approach may reveal potentially

445    important genes and pathways associated with platin resistance. It would be preferable

446    to explore pathways related to signature genes to improve accuracy, identify potential

447    targets for further study of chemoresistance, and expand the model parameters to take

448    into account alternate states besides those captured in the original signature[39].

449    Signatures for resistance may be useful for developing targeted intervention to re-

450    sensitize tumours. For example, the mismatch repair (MMR) gene *MSH2* is commonly

451    present in gene signatures at high resistance levels for oxaliplatin, which is of interest,

452    as MMR deficiency has been shown to be predictive for oxaliplatin resistance[35]. Indeed,

453    *MLH1*, *MSH2* and *MSH6*-deficient cells are more susceptible to oxaliplatin, despite

454    MMR-deficiency being associated with cisplatin resistance[34]. The autoimmune disease-

455    associated gene *SIAE*, which has been previously shown to have a strong negative

456    correlation to oxaliplatin response in advanced CRC patients[40], was selected in the

457    majority of thresholded oxaliplatin models (Supplementary Table S2C). The gene *BCL2*,

458    which was commonly selected for cisplatin (Figure 3A), was rarely selected for

459    oxaliplatin (Figure 3C). At the highest levels of resistance to cisplatin, models were

460    enriched for genes belonging to DNA repair, anti-oxidative response, apoptotic pathways

461    and drug transporters (Figure 3A). These gene pathways are known to be involved in

462    cisplatin resistance[41,42] and these specific genes may be explored in subsequent work to

463    identify the contribution to chemotherapy response in a biochemical context.

464    Log-loss evaluates the accuracy of a classifier by penalizing erroneous

465    classifications, and is relevant in cases where data is imbalanced and/or have an

466    unequally distributed error cost. We assessed whether ML models based on log-loss

467    minimization could improve model accuracy in patient data (Supplementary Table S4)

468    and compared these to models generated by minimizing cell line misclassification. When

20

469    models generated by both methods were highly similar (generated at the same $GI_{50}$

470    threshold, consist of a similar number of genes and consist of ≥ 80% shared genes),

471    prediction accuracy of TCGA cancer patient outcomes were nearly indistinguishable, as

472    accuracy can vary over different relapse thresholds. Where significant differences in

473    predictions were seen, the misclassification-based models were more accurate overall

474    (**Cis1**, **Cis17** and the "12-Resistant" carboplatin model were +8.3%, +5.6% and +3.9%

475    more accurate compared to the log-loss model, respectively). Oxaliplatin models were

476    dissimilar across all $GI_{50}$ thresholds, as the log-loss minimized ML models often contain

477    increased numbers of genes compared to the misclassification-based models. Many of

478    these larger models were less accurate in patients compared to models which minimized

479    misclassification rates consistent that this evaluation and model selection method is

480    more prone to overfitting. This pattern was also noted for models generated at extreme

481    $GI_{50}$ thresholds for all three platin drugs in which response was, by definition, somewhat

482    imbalanced.

483        It may be feasible to predict responses to combination chemotherapy with the

484    models described here. Not included in the present analysis were signatures for

485    methotrexate, vinblastine, and doxorubicin, which comprise the MVAC cocktail used to

486    treat bladder cancer. This was due primarily to a lack of patients treated with this drug

487    combination in the TCGA bladder dataset (N=11). Individual signatures for several of

488    these drugs have been derived and analyzed using the patient data from METABRIC

489    (Molecular Taxonomy of Breast Cancer International Consortium)[24]. A reasonable

490    approach to predicting combination chemotherapy would first determine the probability

491    of sensitivity or resistance to individual drugs, accounting for the misclassification rate by

492    each (defined as $d_1$, …, $d_k$). The ML classifiers output these probabilities, analogous to

493    their misclassification rates in a set of patients treated identically. If the model predicts

494    that the patient is sensitive to drug $d_1$ with 90% probability, and sensitive to drug $d_2$ with

21

495   5% probability, the probability of sensitivity to the combination is 1 - (1 - 0.9)*(1 - 0.05) =

496   90.5%, and the probability of resistance is 9.5%. The correlated responses could be

497   estimated for drug pairs, $d_1$ and $d_2$, and then adjusted for the combined probability of the

498   pair to $d_{12}$, based on the features that are shared by the signatures of both drugs. The

499   probability of sensitivity would then be given by 1 - (1 - $d_{12}$)*(1 - $d_3$)*...*(1 - $d_k$).

500   The predictive accuracy for the same model could differentiate highly between

501   the two datasets. **Cis3** (Supplemental Table S6A) had a high predictive accuracy and

502   AUC for TCGA bladder cancer patients (AUC=0.64). However, the AUC was lower when

503   applied to the Als *et al.*[29] dataset (AUC=0.18). Patient metadata in the latter study only

504   indicated patient survival times, while we base the expected TCGA patient outcome on

505   time to disease recurrence. As the basis of our expected outcome differs between

506   datasets, these differences may be acting as a confounding factor to determine accuracy

507   of gene signatures. The datasets also differ in how expression was measured

508   (microarray vs. RNA-seq). The relevance of models based on training and testing data

509   from different platforms can affect the accuracy of validation, which might not be

510   improved by data normalization. In this study, datasets were subjected to z-score

511   normalization. In subsequent studies, other techniques to correct for some of these

512   effects have been described and could be applied[43].

513   In summary, we describe $GI_{50}$- or $IC_{50}$-threshold-independent ML models to

514   predict chemotherapy response to platin agents in cancer patients. Ensemble machine

515   learning produced combined signatures that were more accurate than most individual

516   models generated with different thresholds. Genes associated cisplatin response

517   included those which exacerbate resistance in patients with a history of smoking. The

518   methodology described here should be adaptable to other drugs and cancer types. With

519   a range of models for multiple drugs, it may be possible to improve the efficacy of

22

520    treatment by tailoring treatment to a patient's specific tumour biology, and reduce

521    treatment duration by limiting the number of different therapeutic regimens prescribed

522    before achieving a successful response[44].

523    **MATERIALS AND METHODS**

524    **Data and preprocessing**

525    Microarray GE and data were from breast cancer cell lines were used to train

526    ML-based gene signatures of drug response based on respective growth or target

527    inhibition data ($GI_{50}$ or $IC_{50}$). Cell lines were treated with either cisplatin (N=39),

528    carboplatin (N=46), or oxaliplatin (N=47)[13]. Bladder cancer cell line GE and $IC_{50}$

529    measurements for cisplatin were obtained from cancerRxgene (N=17). However, all

530    testing was performed on breast cancer cell line data as the number of bladder cancer

531    cell lines was insufficient to produce accurate signatures. RNA-seq GE and survival

532    measurements were downloaded from TCGA for bladder urothelial carcinoma (N=72

533    patients treated with cisplatin)[26], ovarian epithelial tumor (N=410 treated with

534    carboplatin)[27] and colorectal adenocarcinoma (N=99 treated with oxaliplatin)[28]. GE of

535    cisplatin-treated patients of cell carcinoma of the urothelium (N=30)[29] and for oxaliplatin-

536    treated CRC patients (N=83)[30] were obtained from the Gene Expression Omnibus.

537    Clinical metadata and GE for TCGA patients were obtained from Genomic Data

538    Commons (https://gdc.cancer.gov/), while methylation HM450 (Illumina) data for these

539    patients was downloaded from cBioPortal[45].

540    Initial gene sets for developing signatures for each drug were identified from

541    previously published literature (see Supplemental References) and databases, such as

542    PharmGKB and DrugBank[46,47]. The final gene sets were chosen using MFA to analyze

543    interactions between GE, CN, and $GI_{50}$ data for the drug of interest[48]. Genes whose GE

23

544    and/or CN showed a direct or inverse correlation with $GI_{50}$ were selected for SVM

545    training. As the number of genes related to $GI_{50}$ for oxaliplatin exceeded the number of

546    cell lines available for training, we limited the input to the ML model oxaliplatin to those

547    genes whose GE were related to $GI_{50}$. Similarly, the number correlating genes in

548    cisplatin treated cells exceeded the number of cell lines. For cisplatin, genes whose

549    expression correlated with $GI_{50}$ were eliminated if they showed no or little expression in

550    TCGA bladder cancer patients (i.e. RNA-seq counts by Expectation Maximization

551    [RSEM] were < 5.0 for majority of individuals). This reduces the overall number of genes

552    for SVM analysis, and thus helps to avoid a data to size sample imbalance. For cisplatin,

553    MFA was repeated using $IC_{50}$ values for 17 bladder cancer cell lines; however, the

554    available CN data generally showed a lack of variation in the cell lines for these genes.

555    Instead, the available $IC_{50}$ values for three other cancer drugs (doxorubicin,

556    methotrexate and vinblastine) were compared with the $IC_{50}$ of cisplatin by MFA.

557    Applying an SVM model directly to patient data without a normalization approach

558    is imprecise when training and testing data are not obtained using similar methodology

559    (i.e. different microarray platforms). To compare the cell line GE microarray data and the

560    patient RNA-seq GE datasets, expression values were normalized by conversion to z-

561    scores using MATLAB[49]. Although Log2 intensity values from microarray data were not

562    available for TCGA samples, RNA-seq based GE and $log_2$ intensities from microarray

563    data are highly correlated[50].

**Machine Learning**

565    SVMs were trained with breast cancer cell line GE datasets[13] with the Statistics

566    Toolbox in MATLAB[49] similar to Dorman et al (2016)[6]. Rather than a linear kernel, we

567    used a Gaussian kernel function (fitcsvm), and then tested with leave-one-out cross-

24

568    validation (using the options '*crossval'* and '*leaveout'*). A greedy backwards feature

569    selection algorithm was used to improve classification accuracy[51]. BFS leaves out

570    individual genes from the initial MFA-qualified gene set, then trains a cross validated

571    Gaussian kernel SVM on the training samples, removing the gene with the highest

572    misclassification rate. The procedure is repeated until all genes have been evaluated.

573    The gene subset with the lowest misclassification rate[6] or log-loss statistic[25] based on

574    cross-validation is selected as the model for subsequent testing with patient GE and

575    clinical data. K-fold cross validation of the misclassification-based models was

576    performed using MATLAB software described in Zhao et al. (2018)[25].

577         SVMs minimized according to the log-loss classification function were also

578    generated with both software described in Zhao *et al.* (2018; uses multiclass compatible

579    'fitcecoc' function)[25], and with a modified version of the software described above (using

580    'fitSVMPosterior' to compute posterior probabilities). Computed probabilities differ

581    between 'fitSVMPosterior' and 'fitcecoc' (range: 0.02-0.04), thus the resultant models will

582    differ between the two programs. When given unbalanced data (e.g. lower resistance

583    thresholds), 'fitSVMPosterior' will warn that some classes are not represented, and thus

584    those folds will not predict the labels for those missing classes. The log-loss models

585    described in this manuscript were generated with the multiclass compatible 'fitcecoc'

586    function software[25].

587    *Derivation of gene signatures for different drug resistance thresholds*

588         We have previously set a conventional $GI_{50}$ threshold distinguishing sensitivity

589    from resistance at the *median* of the range of drug concentrations that inhibited cell

590    growth by 50%[6]. We hypothesized that different gene signatures could be derived for

591    different levels of drug resistance by varying this threshold. ML experiments for

592   classifying resistance or sensitivity at $GI_{50}$ values generated a series of optimized

593   Gaussian SVM models whose performance were assessed with patient expression data

594   for each signature. A heat map which illustrates the frequencies of genes appearing in

595   these models was created with the R language *hist2d* function.

596       A composite gene signature was created by ensemble averaging of all models

597   generated at each resistance threshold. Ensemble averaging combines signatures

598   through averaging the weighted accuracy of a set of related models[38]. The decision

599   function for the ensemble classifier is the mean of the decision function scores of the

600   component classifiers, weighted by the AUC.

601   *Significance of cell line-derived models*

602       The potential for models to overfit data during training and/or feature selection

603   was first assessed by permutation analysis with randomized cell line labels and with

604   random sets of genes, as described previously[6]. Using the median cisplatin $GI_{50}$ as the

605   resistance threshold, 10,000 models based on random gene selection (15 genes) had

606   higher rates of misclassification than the best median SVM models (2 signatures with

607   7.7% misclassification). Cisplatin, carboplatin and oxaliplatin GE data for random cell

608   line label combinations (n=10,000) generated only 8, 1 and 1 signatures, respectively,

609   with lower error rates than the best biochemically-inspired signatures. When minimizing

610   for log-loss (rather than misclassification), random gene analysis (10,000 iterations;

611   median cisplatin $GI_{50}$ threshold) resulted only in models with a higher log-loss than the

612   model generated with the initial cisplatin gene set. Log-loss based random label analysis

613   (n=2000 combinations) resulted in 3.4% of random label models resulted in a lower log-

614   loss than the cisplatin model at the same $GI_{50}$ threshold (5.27).This was not entirely

615   surprising, since this result depends on the $GI_{50}$ threshold used for labeling. The

26

616    differences between $GI_{50}$ values for cell lines close to the median $GI_{50}$ used in this

617    analysis are almost negligible (e.g. 5.11 vs 5.12) and likely within the measurement error

618    for these values.

619    Cell-line based model accuracies in predicting outcomes of platin-treated TCGA

620    bladder cancer patients were compared with results from participants who did not

621    receive these treatments (using an 18 months post-treatment threshold). In non-platin

622    treated patients, 36.5% of those who were disease-free were predicted accurately with

623    the **Cis1** signature (N=178; 22% less accurate than platin treated patients), and 62.9%

624    accurate for those with recurrent disease (N=70; 8.1% less accurate). **Cis2** was 43.8%

625    accurate for disease-free non-platin treated patients (N=178; 12.3% lower accuracy),

626    and 60.0% of those who relapsed (N=70; 2.9% less accurate). Gene expression

627    changes in patients treated with platin drugs are better modeled by cancer cell-line

628    based predictors than in patients receiving other treatments.

629    **ACKNOWLEDGEMENTS**

634    **CONFLICTS OF INTEREST**

635    PKR cofounded CytoGnomix Inc., which hosts the interactive resource described in this

636    study for prediction of responses to chemotherapy agents. The other authors have no

637    conflicts of interest.

638    **AUTHOR CONTRIBUTIONS**

639     PKR and DL designed the methodology. EJM and JZ performed analyses. EJM and

640     PKR wrote the manuscript.

641     **FUNDING**

644 **REFERENCES**

645 1. Cardoso, F. *et al.* Locally recurrent or metastatic breast cancer: ESMO Clinical Practice

646    Guidelines for diagnosis, treatment and follow-up. *Ann. Oncol.* **23,** vii11–vii19 (2012).

647 2. Oostendorp, L. J., Stalmeier, P. F., Donders, A. R. T., van der Graaf, W. T. & Ottevanger, P. B.

648    Efficacy and safety of palliative chemotherapy for patients with advanced breast cancer

649    pretreated with anthracyclines and taxanes: a systematic review. *Lancet Oncol.* **12,** 1053–

650    1061 (2011).

651 3. Alfarouk, K. O. *et al.* Resistance to cancer chemotherapy: failure in drug response from

652    ADME to P-gp. *Cancer Cell Int.* **15,** 71 (2015).

653 4. Gąsowska-Bodnar, A. *et al.* Survivin Expression as a Prognostic Factor in Patients With

654    Epithelial Ovarian Cancer or Primary Peritoneal Cancer Treated With Neoadjuvant

655    Chemotherapy: *Int. J. Gynecol. Cancer* **24,** 687–696 (2014).

656 5. Hatzis, C. *et al.* A genomic predictor of response and survival following taxane-anthracycline

657    chemotherapy for invasive breast cancer. *JAMA* **305,** 1873–1881 (2011).

658 6. Dorman, S. N. *et al.* Genomic signatures for paclitaxel and gemcitabine resistance in breast

659    cancer derived by machine learning. *Mol. Oncol.* **10,** 85–100 (2016).

660 7. Zhang, S. *et al.* Organic Cation Transporters Are Determinants of Oxaliplatin Cytotoxicity.

661    *Cancer Res.* **66,** 8847–8857 (2006).

662 8. Poisson, L. M. *et al.* A metabolomic approach to identifying platinum resistance in ovarian

663    cancer. *J. Ovarian Res.* **8,** (2015).

664 9. Cadoná, F. C. *et al.* Guaraná a Caffeine-Rich Food Increases Oxaliplatin Sensitivity of

665    Colorectal HT-29 Cells by Apoptosis Pathway Modulation. *Anticancer Agents Med. Chem.* **16,**

666    1055–1065 (2016).

667     10. Kasparkova, J., Vojtiskova, M., Natile, G. & Brabec, V. Unique Properties of DNA Interstrand

668         Cross-Links of Antitumor Oxaliplatin and the Effect of Chirality of the Carrier Ligand. *Chem. –*

669         *Eur. J.* **14,** 1330–1341 (2008).

670     11. Woynarowski, J. M. *et al.* Oxaliplatin-Induced Damage of Cellular DNA. *Mol. Pharmacol.* **58,**

671         920–927 (2000).

672     12. Tashiro, T., Kawada, Y., Sakurai, Y. & Kidani, Y. Antitumor activity of a new platinum

673         complex, oxalato (trans-l-1,2-diaminocyclohexane)platinum (II): new experimental data.

674         *Biomed. Pharmacother.* **43,** 251–260 (1989).

675     13. Daemen, A. *et al.* Modeling precision treatment of breast cancer. *Genome Biol.* **14,** R110

676         (2013).

677     14. Yuan, Y. *et al.* Identification of the biomarkers for the prediction of efficacy in first-line

678         chemotherapy of metastatic colorectal cancer patients using SELDI-TOF-MS and artificial

679         neural networks. *Hepatogastroenterology.* **59,** 2461–2465 (2012).

680     15. L'Espérance, S., Bachvarova, M., Tetu, B., Mes-Masson, A.-M. & Bachvarov, D. Global gene

681         expression analysis of early response to chemotherapy treatment in ovarian cancer

682         spheroids. *BMC Genomics* **9,** 99 (2008).

683     16. Nickerson, M. L. *et al.* Molecular analysis of urothelial cancer cell lines for modeling tumor

684         biology and drug response. *Oncogene* (2016).

685     17. Yuryev, A. Gene expression profiling for targeted cancer treatment. *Expert Opin. Drug*

686         *Discov.* **10,** 91–99 (2015).

687     18. Sataloff, D. M. *et al.* Pathologic response to induction chemotherapy in locally advanced

688         carcinoma of the breast: a determinant of outcome. *J. Am. Coll. Surg.* **180,** 297–306 (1995).

689    19. Ogston, K. N. *et al.* A new histological grading system to assess response of breast cancers

690        to primary chemotherapy: prognostic significance and survival. *Breast Edinb. Scotl.* **12,** 320–

691        327 (2003).

692    20. Earl, J. *et al.* The UBC-40 Urothelial Bladder Cancer cell line index: a genomic resource for

693        functional studies. *BMC Genomics* **16,** 403 (2015).

694    21. Rixe, O. *et al.* Oxaliplatin, tetraplatin, cisplatin, and carboplatin: Spectrum of activity in drug-

695        resistant cell lines and in the cell lines of the national cancer institute's anticancer drug

696        screen panel. *Biochem. Pharmacol.* **52,** 1855–1865 (1996).

697    22. Mehmood, R. K. Review of Cisplatin and oxaliplatin in current immunogenic and monoclonal

698        antibody treatments. *Oncol. Rev.* **8,** 256 (2014).

699    23. Kweekel, D. M., Gelderblom, H. & Guchelaar, H.-J. Pharmacology of oxaliplatin and the use

700        of pharmacogenomics to individualize therapy. *Cancer Treat. Rev.* **31,** 90–105 (2005).

701    24. Mucaki, E. J. *et al.* Predicting Outcomes of Hormone and Chemotherapy in the Molecular

702        Taxonomy of Breast Cancer International Consortium (METABRIC) Study by Biochemically-

703        inspired Machine Learning. *F1000Research* **5,** 2124 (2017).

704    25. Zhao, J. Z. L., Mucaki, E. J. & Rogan, P. K. Predicting ionizing radiation exposure using

705        biochemically-inspired genomic machine learning. *F1000Research* **7,** 233 (2018).

706    26. Robertson, A. G. *et al.* Comprehensive Molecular Characterization of Muscle-Invasive

707        Bladder Cancer. *Cell* **171,** 540-556.e25 (2017).

708    27. Cancer Genome Atlas Research Network. Integrated genomic analyses of ovarian

709        carcinoma. *Nature* **474,** 609–615 (2011).

710    28. Cancer Genome Atlas Network. Comprehensive molecular characterization of human colon

711        and rectal cancer. *Nature* **487,** 330–337 (2012).

712     29. Als, A. B. *et al.* Emmprin and survivin predict response and survival following cisplatin-

713         containing chemotherapy in patients with advanced bladder cancer. *Clin. Cancer Res. Off. J.*

714         *Am. Assoc. Cancer Res.* **13,** 4407–4414 (2007).

715     30. Tsuji, S. *et al.* Potential responders to FOLFOX therapy for colorectal cancer by Random

716         Forests analysis. *Br. J. Cancer* **106,** 126–132 (2012).

717     31. Tembe, V. *et al.* The BARD1 BRCT domain contributes to p53 binding, cytoplasmic and

718         mitochondrial localization, and apoptotic function. *Cell. Signal.* **27,** 1763–1771 (2015).

719     32. Freedman, N. D., Silverman, D. T., Hollenbeck, A. R., Schatzkin, A. & Abnet, C. C. Association

720         between smoking and risk of bladder cancer among men and women. *JAMA* **306,** 737–745

721         (2011).

722     33. Joehanes, R. *et al.* Epigenetic Signatures of Cigarette Smoking. *Circ. Cardiovasc. Genet.*

723         (2016).

724     34. Raymond, E., Faivre, S., Chaney, S., Woynarowski, J. & Cvitkovic, E. Cellular and Molecular

725         Pharmacology of Oxaliplatin. *Mol. Cancer Ther.* **1,** 227–235 (2002).

726     35. Alex, A. K. *et al.* Response to Chemotherapy and Prognosis in Metastatic Colorectal Cancer

727         With DNA Deficient Mismatch Repair. *Clin. Colorectal Cancer* (2016).

728     36. Sos, M. L. *et al.* Predicting drug susceptibility of non-small cell lung cancers based on genetic

729         lesions. *J. Clin. Invest.* **119,** 1727–1740 (2009).

730     37. Laderas, T. G., Heiser, L. M. & Sönmez, K. A Network-Based Model of Oncogenic

731         Collaboration for Prediction of Drug Sensitivity. *Front. Genet.* **6,** 341 (2015).

732     38. Clemen, R. T. Combining forecasts: A review and annotated bibliography. *Int. J. Forecast.* **5,**

733         559–583 (1989).

734     39. Airley, R. *Cancer chemotherapy*. (Wiley-Blackwell, 2009).

735   40. Li, X.-X. *et al.* RNA-seq identifies determinants of oxaliplatin sensitivity in colorectal cancer
736        cell lines. *Int. J. Clin. Exp. Pathol.* **7,** 3763–3770 (2014).
737   41. Borst, P., Rottenberg, S. & Jonkers, J. How do real tumors become resistant to cisplatin? *Cell*
738        *Cycle Georget. Tex* **7,** 1353–1359 (2008).
739   42. Wernyj, R. & Morin, P. Molecular mechanisms of platinum resistance: still searching for the
740        Achilles? heel. *Drug Resist. Updat.* **7,** 227–232 (2004).
741   43. Johnson, W. E., Li, C. & Rabinovic, A. Adjusting batch effects in microarray expression data
742        using empirical Bayes methods. *Biostat. Oxf. Engl.* **8,** 118–127 (2007).
743   44. Akamatsu, N., Nakajima, H., Ono, M. & Miura, Y. Increase in acetyl CoA synthetase activity
744        after phenobarbital treatment. *Biochem. Pharmacol.* **24,** 1725–1727 (1975).
745   45. Gao, J. *et al.* Integrative analysis of complex cancer genomics and clinical profiles using the
746        cBioPortal. *Sci. Signal.* **6,** pl1 (2013).
747   46. Whirl-Carrillo, M. *et al.* Pharmacogenomics Knowledge for Personalized Medicine. *Clin.*
748        *Pharmacol. Ther.* **92,** 414–417 (2012).
749   47. Law, V. *et al.* DrugBank 4.0: shedding new light on drug metabolism. *Nucleic Acids Res.* **42,**
750        D1091–D1097 (2014).
751   48. Abdi, H. & Williams, L. J. Principal component analysis. *Wiley Interdiscip. Rev. Comput. Stat.*
752        **2,** 433–459 (2010).
753   49. MATLAB and Statistics Toolbox Release 2012b, The MathWorks, Inc., Natick, Massachusetts,
754        United States.
755   50. Marioni, J. C., Mason, C. E., Mane, S. M., Stephens, M. & Gilad, Y. RNA-seq: an assessment
756        of technical reproducibility and comparison with gene expression arrays. *Genome Res.* **18,**
757        1509–1517 (2008).

758    51.  Bermingham, M. L. *et al.* Application of high-dimensional feature selection: evaluation for

759            genomic prediction in man. *Sci. Rep.* **5,** 10312 (2015).

760

761    **Tables**

762    **Table 1: Models Which Best Predicted Response in TCGA Cancer Patients**

| Model ID | Cancer Type Tested | GI50 Threshold | Signature (C, σ*) |
|---|---|---|---|
| **Cis1** (Cisplatin) | Bladder | 5.11 | BARD1, BCL2, BCL2L1, CDKN2C, FAAP24, FEN1, MAP3K1, MAPK13, MAPK3, NFKB1, NFKB2, SLC22A5, SLC31A2, TLR4, TWIST1 (100000, 100) |
| Cis2 (Cisplatin) | Bladder | 5.12 | BARD1, BCL2L1, CFLAR, FAAP24, MAP3K1, MAPK3, NFKB1, POLQ, PRKAA2, SLC22A5, SLC31A2, TLR4, TWIST1 (10000, 100) |
| Cis3 (Cisplatin) | Bladder | 5.60 | BCL2, CFLAR, ERCC2, ERCC6, FAAP24, FEN1, MAP3K1, NFKB1, NFKB2, PNKP, POLQ, PRKCB, SLC22A5, SNAI1, TLR4 (100000, 100) |
| Cis12 (Cisplatin) | Bladder | 5.40 | ATP7B, BCL2, BCL2L1, CDKN2C, ERCC2, FAAP24, GSTO1, MAP3K1, MAPK3, MT2A, NFKB1, NFKB2, POLD1, POLQ, PRKCB, SNAI1, TLR4, TP63 (10000, 100) |
| Cis14 (Cisplatin) | Bladder | 5.16 | BARD1, BCL2, BCL2L1, CDKN2C, FAAP24, FEN1, FOS, GSTP1, MAP3K1, MAPK13, MAPK3, MSH2, NFKB1, POLD1, POLQ, PRKAA2, PRKCB, SLC22A5, SLC31A2, SNAI1, TWIST1 (10000, 100) |
| Cis17 (Cisplatin) | Bladder | 5.10 | ATP7B, BCL2, BCL2L1, FEN1, GSTP1, MAP3K1, MAPK3, MT2A, NFKB1, PNKP, POLQ, PRKAA2, PRKCB, SLC31A2, TLR4, TP63 (100000, 100) |
| **Car1** (Carboplatin) | Ovarian | 4.22 | AKT1, EIF3K, ERCC1, GNGT1, GSR, MTHFR, NEDD4L, NLRP1, NRAS, RAF1, SGK1, TIGD1, TP53, VEGFB, VEGFC (100000, 100) |
| Car9 (Carboplatin) | Ovarian | 4.32 | AKT1, ATP7B, EIF3I, ETS2, GNGT1, HRAS, KRAS, LIG3, MTHFR, MTR, NRAS, RAD50, SCN10A, TIGD1, TP53, VEGFB (10000, 100) |
| Car51 (Carboplatin) | Ovarian | 4.34 | AKT1, EGF, EIF3I, ERCC1, ETS2, GNGT1, KRAS, MTHFR, MTR, NEDD4L, NLRP1, NRAS, RAD50, RAF1, SGK1, TIGD1, TP53, VEGFB, VEGFC (10000, 100) |
| Car73 (Carboplatin) | Ovarian | 4.09 | AKT1, ATP7B, ETS2, GNGT1, HRAS, NLRP1, SCN10A, VEGFB (100000, 1000) |
| **Oxa1** (Oxaliplatin) | Colorectal | 5.10 | BRAF, FCGR2A, IGF1, MSH2, NAGK, NFE2L2, NQO1, PANK3, SLC47A1, SLCO1B1, UGT1A1 (10, 10) |
| Oxa21 (Oxaliplatin) | Colorectal | 5.10 | BRAF, IGF1, IGF1R, KLF3, MSH2, NAT2, NFE2L2, NQO1, PANK3, PRSS1, SIAE, SLC47A1, SLCO1B1, UGT1A1 (1000, 100) |

*C - The box-constraint. $\sigma$ – the kernel-scale ("sigma"). Bolded models are those that best overall performance against TCGA cancer patient gene expression data.

763

764

**FIGURE LEGENDS**

**Figure 1.** Schematic of platinum drug sensitivity and resistance genes which showed MFA correlation for $GI_{50}$ of A) cisplatin, B) carboplatin, and C) oxaliplatin. The genes used to derive the SVM are shown in context of their effect in the cell and role in cisplatin mechanisms of action. GE and CN correlation with inhibitory drug concentration by MFA of breast ($GI_{50}$) and bladder ($IC_{50}$) cancer cell line data.

**Figure 2:** $GI_{50}$ values for cell lines treated with the three platin drugs were plotted in order of ascending oxaliplatin $GI_{50}$. For most cell lines, there is a visible trend between the $GI_{50}$ for cisplatin and carboplatin, reflecting the correlation between the two drugs seen by MFA. Despite this correlation, carboplatin shows a much smaller variance (0.22) compared to cisplatin (0.37; oxaliplatin variance is 0.34).

**Figure 3.** The variation in gene composition of misclassification-based SVMs at different $GI_{50}$ thresholds for A) cisplatin, B) carboplatin, and C) oxaliplatin. $GI_{50}$ intervals are indicated on the left, with the number of cell lines with $GI_{50}$ values within said intervals in brackets. Each box represents the density of genes appearing in optimized Gaussian SVM models in those functional categories, with darker grey indicating frequent genes in indicated $GI_{50}$ threshold intervals, while lighter grey indicates less commonly selected genes. The number of thresholded models used to derive the density plot within each interval is equal (or greater, in the case of multiple equally performing models) to the number of cell lines within that $GI_{50}$ interval.

**Supplementary Figure 1.** The variation in gene composition of log-loss based SVMs at different $GI_{50}$ thresholds for A) cisplatin, B) carboplatin, and C) oxaliplatin. Each box represents the density of genes appearing in optimized Gaussian log-loss SVM models

36

788    in those functional categories, with darker grey indicating frequent genes in indicated

789    $GI_{50}$ threshold intervals, while lighter grey indicates less commonly selected genes.

790    **Supplementary Figure 2.** Classification accuracy of models on TCGA bladder cancer

791    patients treated with cisplatin and/or carboplatin as the resistance threshold is varied.

792    Recurrence and disease-free survival are used as a binary measure to assess

793    performance. The x-axis indicates movement of the resistance threshold, with more cell

794    lines labeled sensitive on the left and more labeled resistant on the right. Maximal AUC

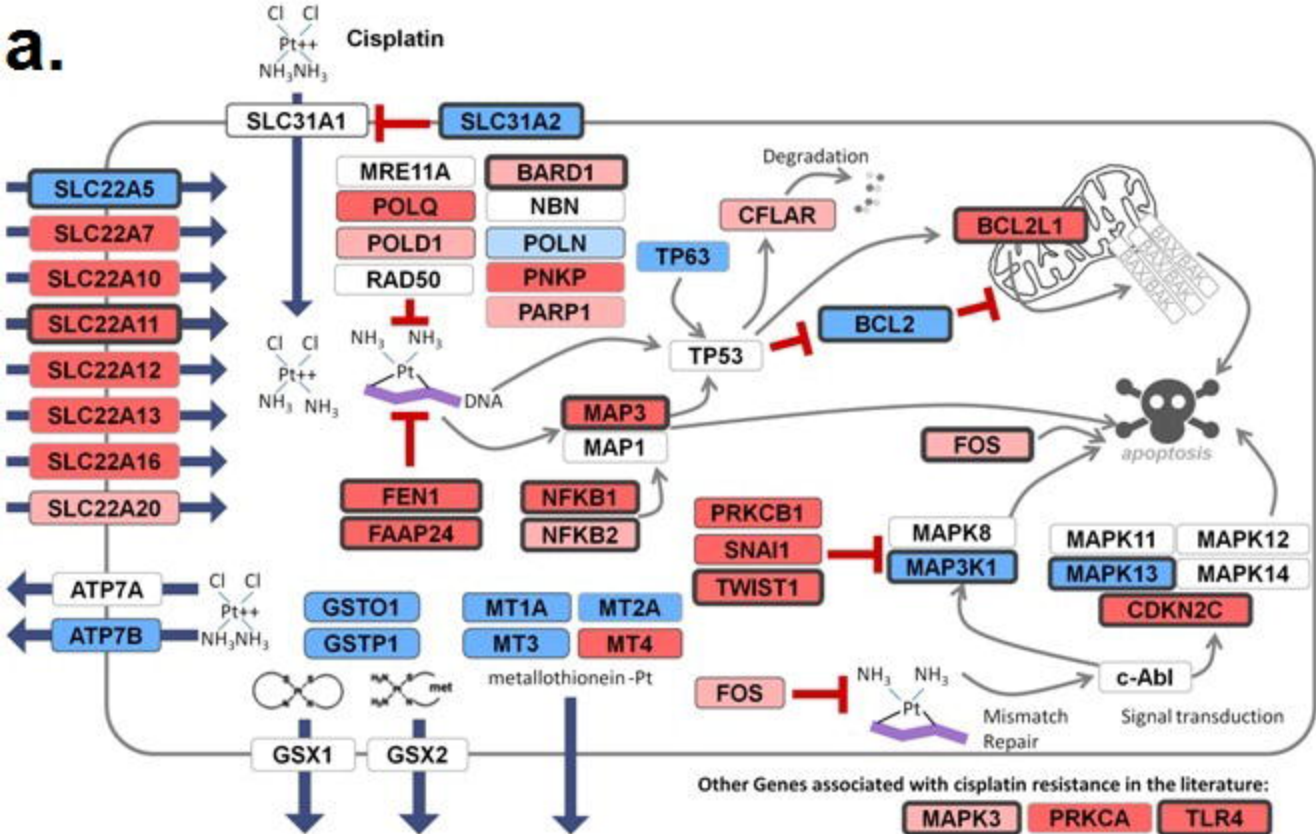795    is indicated by the downward arrows.

796    **Supplementary Figure 3.** Classification accuracy of SVM models for cisplatin, at a

797    range of response thresholds, were assessed using gene expression data for cisplatin-

798    treated bladder cancer patients from Als *et al.* [29]. Patients with a ≥ 5 year survival post-

799    treatment were labeled sensitive. Red arrows indicate the SVM models with the highest

800    positive predictive value (PPV) in the accuracy of classification of patient outcome.

801    **Supplementary Figure 4.** Hyperplane distance calculated by all thresholded SVMs for

802    recurrent (<6 months) TCGA patients. Each diagram represents the predictions of all

803    SVMs for all patients who had recurrence less than 6 months after treatment (N=10).

804    Each point represents an SVM, where the x-axis represents the number of cell lines set

805    to resistant (in order of lowest to highest $GI_{50}$), and the y-axis represents the calculated

806    hyperplane distance. A negative hyperplane distance would represent a prediction of

807    resistance to cisplatin. Despite this, some patients show a strong preference towards

808    predictions of sensitivity (i.e. TCGA-XF-A9SU).

809    **Supplementary Figure 5.** Hyperplane distance calculated by all thresholded SVMs for

810    sensitive TCGA patients. Each diagram represents the predictions of all SVMs for all

811    patients who had recurrence > 4 years after treatment (top; N=3), or patients who
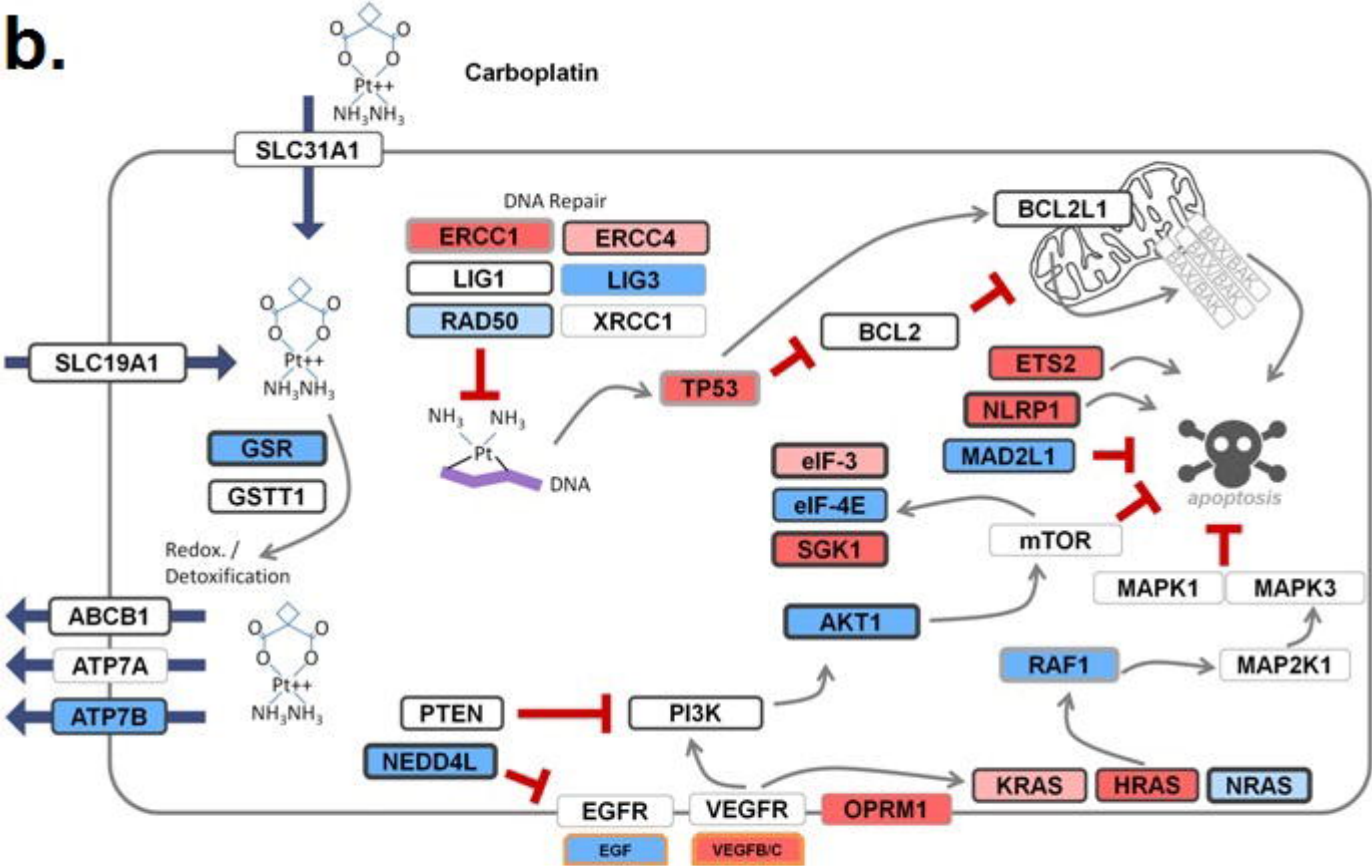
812     showed no recurrence after 6 years (bottom; N=6). Each point represents an SVM,

813     where the x-axis represents the number of cell lines set to resistant (in order of lowest to

814     highest $GI_{50}$), and the y-axis represents the calculated hyperplane distance. A positive

815     hyperplane distance would represent a prediction of sensitivity to cisplatin.

**a.**

**b.**

Carboplatin

SLC31A1

SLC19A1

DNA Repair

| ERCC1 | ERCC4 |
| LIG1 | LIG3 |
| RAD50 | XRCC1 |

GSR

GSTT1

Redox. / Detoxification

ABCB1

ATP7A

ATP7B

PTEN

NEDD4L

TP53

BCL2

BCL2L1

ETS2

NLRP1

MAD2L1

eIF-3

eIF-4E

SGK1

mTOR

MAPK1  MAPK3

AKT1

RAF1

MAP2K1

PI3K

EGFR  VEGFR  OPRM1

EGF

VEGFB/C

KRAS  HRAS  NRAS

*apoptosis*
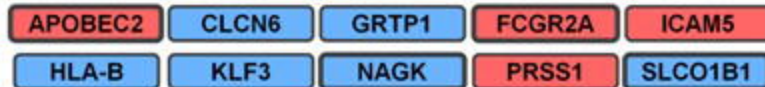
Other Genes associated with carboplatin resistance in the literature:

| GNGT | MTR | MTHFR |
| OPRM1 | SCN10A | TIGD1 |

**C.**

# a.



GI$_{50}$ (# cell lines in interval)

Resistant

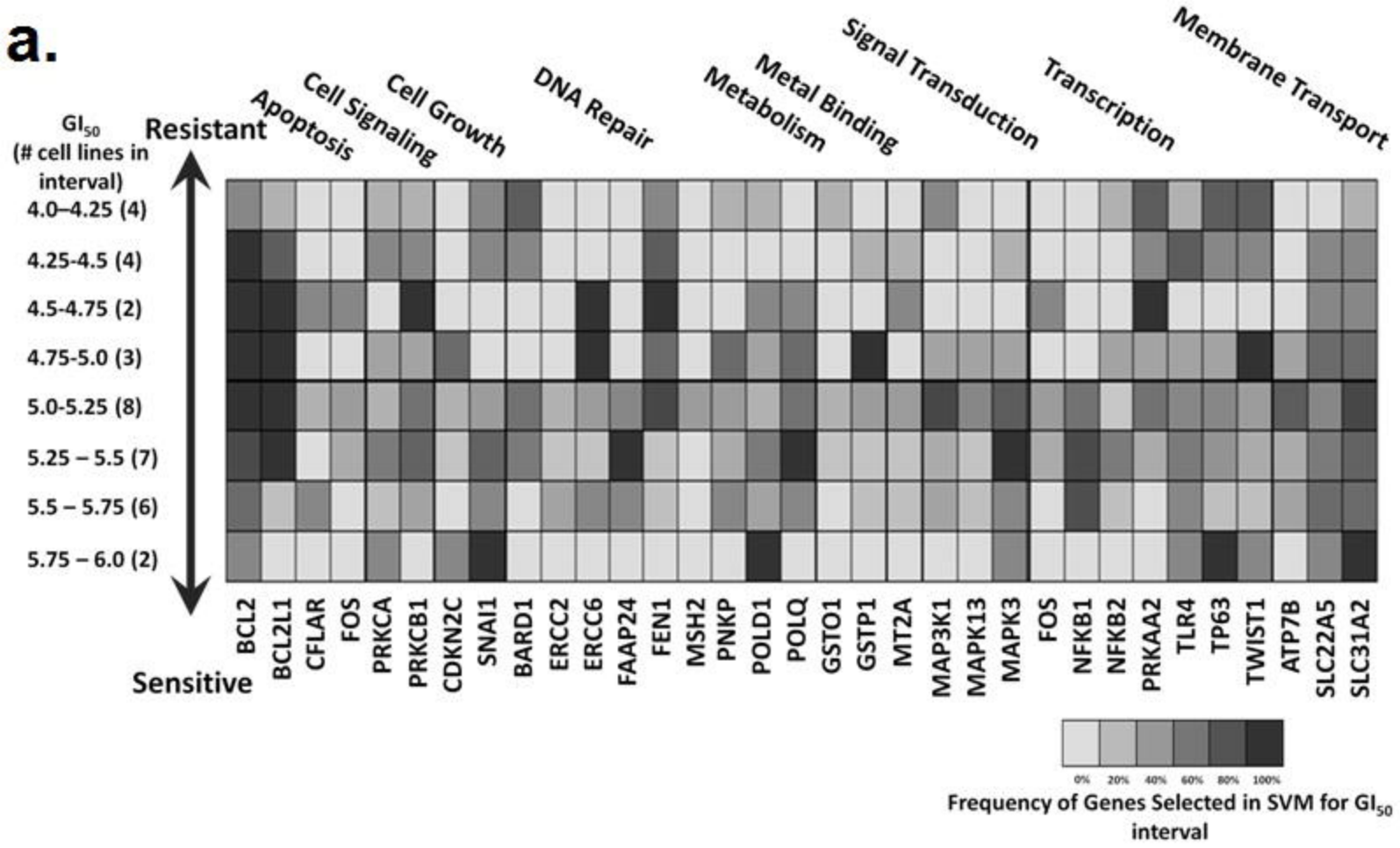4.0–4.25 (4)

4.25-4.5 (4)

4.5-4.75 (2)

4.75-5.0 (3)

5.0-5.25 (8)

5.25 – 5.5 (7)

5.5 – 5.75 (6)

5.75 – 6.0 (2)

Sensitive

Apoptosis · Cell Signaling · Cell Growth · DNA Repair · Metabolism · Metal Binding · Signal Transduction · Transcription · Membrane Transport

BCL2, BCL2L1, CFLAR, FOS, PRKCA, PRKCB1, CDKN2C, SNAI1, BARD1, ERCC2, ERCC6, FAAP24, FEN1, MSH2, PNKP, POLD1, POLQ, GSTO1, GSTP1, MT2A, MAP3K1, MAPK13, MAPK3, FOS, NFKB1, NFKB2, PRKAA2, TLR4, TP63, TWIST1, ATP7B, SLC22A5, SLC31A2

0%  20%  40%  60%  80%  100%

Frequency of Genes Selected in SVM for GI$_{50}$ interval

**b.**

Frequency of Genes Selected in SVM for GI$_{50}$ interval

**C.**

GI$_{50}$ (# cell lines in interval)

Resistant ↑

3.50 – 4.50 (4)
4.51 – 4.80 (7)
4.81 – 5.10 (7)
5.11 – 5.40 (7)
5.41 – 5.70 (14)
5.71 – 6.00 (4)

Sensitive ↓

Apoptosis · Immune System · Metabolism · Signal Transduction · Transcription / DNA Repair · Membrane Transport · Uncharacterized

BCL2, CSMD1, ICAM5, FCGR2A, HLA-B, LTA, MPO, SIAE, PRSS1, AGXT, NAT2, NQO1, PANK3, SGPP2, UGT1A1, NAGK, PROC, APOBEC2, KIT, KLC3, BRAF, IGF1R, KISS1, NOTCH1, IGF1, KLF3, HIF1A, NFE2L2, PARP15, MSH2, SLC47A1, CLCN6, ABCG2, SLCO1B1, GRTP1

0%  20%  40%  60%  80%  100%

**Frequency of Genes Selected in SVM for GI$_{50}$ interval**