

Exploiting convergent evolution to derive a pan-cancer cisplatin sensitivity gene expression signature

Jessica A. Scarborough^{1,2}, Steven A. Eschrich³, Javier Torres-Roca⁴, Andrew Dhawan^{5,7,*}, and Jacob G. Scott^{1,2,6,8,*}

¹Systems Biology and Bioinformatics Department, School of Medicine, Case Western Reserve University, Cleveland, OH

²Department of Translational Hematology and Oncology Research, Taussig Cancer Institute, Cleveland Clinic, Cleveland, OH

⁵Neurological Institute, Cleveland Clinic, Cleveland, OH

³Biostatistics and Bioinformatics Program, Moffitt Cancer Center, Tampa, FL

⁴Department of Radiation Oncology, Moffitt Cancer Center, Tampa, FL

⁶Department of Radiation Oncology, Taussig Cancer Institute, Cleveland Clinic, Cleveland, OH

*Both authors are equal corresponding authors.

⁷dhawana@ccf.org

⁸scottj10@ccf.org

ABSTRACT

Precision medicine offers remarkable potential for the treatment of cancer, but is largely focused on tumors that harbor actionable mutations. Gene expression signatures can expand the scope of precision medicine by predicting response to traditional (cytotoxic) chemotherapy agents without relying on changes in mutational status. We present a novel signature extraction method, inspired by the principle of convergent evolution, which states that tumors with disparate genetic backgrounds may evolve similar phenotypes independently. This evolutionary-informed method can be utilized to produce signatures predictive of response to over 200 chemotherapeutic drugs found in the Genomics of Drug Sensitivity in Cancer Database. Here, we demonstrate its use by extracting the Cisplatin Response Signature, CisSig, for use in predicting a common trait (sensitivity to cisplatin) across disparate tumor subtypes (epithelial-origin tumors). CisSig is predictive of cisplatin response within the cell lines and clinical trends in independent datasets of tumor samples. This novel methodology can be used to produce robust signatures for the prediction of traditional chemotherapeutic response, dramatically increasing the reach of personalized medicine in cancer.

1 Introduction

2 Despite rich collections of cancer “-omic” data, precision medicine research has largely focused on producing therapies that
3 target somatic mutations in proposed driver genes. These therapies have produced some inspiring successes, extending the
4 lives of patients with targetable mutations by months to years.¹⁻³ However, the reach of genome-driven care is narrow and
5 most patients without targetable mutations simply have not seen the benefits of personalized medicine. In fact, it was estimated
6 that in 2018, less than 5% of cancer patients in the United States could benefit from genome driven care.⁴ Even among the
7 patients who do benefit from genome-driven care, the costs of targeted agents are high and the clinical responses are typically
8 not durable.

9 Without an actionable mutation, patients often receive conventional cytotoxic chemotherapy. In these scenarios, there
10 are significant opportunities for expanding the reach of precision medicine. For example, gene expression signatures can be
11 used to predict response to these traditional chemotherapy agents without relying on changes in mutational status. Not only is
12 gene expression a powerful measure of phenotype, it is readily translatable to a clinical setting, as patient tumors can undergo
13 RNA-sequencing at relatively low cost and high scale.

14 Defined as a set of genes (typically fewer than 100), certain gene expression signatures have already been incorporated
15 into standard-of-care and clinical decision-making algorithms (e.g. OncotypeDx⁵, MammaPrint⁶). Additionally, signatures of
16 radiosensitivity have been developed and have achieved level 1 evidentiary status for archival tissue.⁷⁻¹⁰ Yet, a major obstacle
17 in the field is finding gene expression signatures that are robust enough to be predictive in novel datasets. And although there is
18 a great need for distilling complex gene expression data into a clinical tool, most published gene expression signatures perform
19 no better than signatures consisting of random genes.¹¹ To address this problem, we propose a novel method for the extraction
20 of chemotherapeutic response signatures, utilizing both cell line data, to add isolated drug response information, and tumor
21 sample data, to improve clinical translatability.

22 As seen in experimental evolution, a variety of evolutionary trajectories can lead to the same phenotype.¹²⁻¹⁵ **Figure 1A**
23 shows a canonical example of convergent evolution, where genomically disparate species (bats and birds) both evolved the same
24 phenotype of flight independently of one another. Just as bats and birds are genetically closer to mice and reptiles, respectively,
25 individual tumors may be genotypically similar to tumors with differing drug response phenotypes, **Figure 1B**. Under the
26 selection pressure of a chemotherapeutic agent, tumors may take a wide variety of genomic pathways when evolving drug
27 sensitivity or resistance, meaning that searching for a single genomic marker would be infeasible. In order to understand the
28 basis of chemotherapeutic response, our approach exploits the principles of convergent evolution by combining hundreds of
29 cell lines from a variety of cancer subtypes and extracting transcriptomic patterns of this phenotypic state.

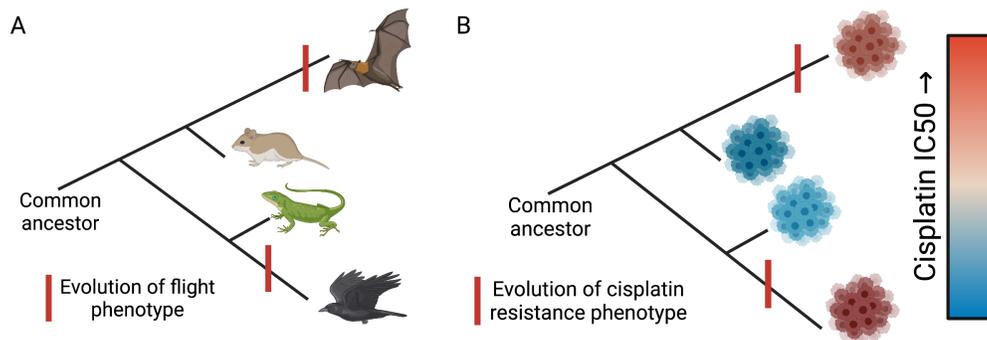


Figure 1. Visual representation of convergent evolution in animals and tumors. A. Birds and bats are genomically disparate, but both have individually evolved the ability to fly. **B.** Two tumors may evolve cisplatin resistance independently, despite being genomically distinct from one another.

30 Our work leverages a seed gene approach, as in Buffa et al., where previously identified hypoxia-regulated genes became
31 seeds in a co-expression network, and highly connected genes formed a hypoxia metagene (gene signature)¹⁶. By extracting
32 genes that are highly co-expressed with biologically significant genes, Buffa et al. produced a robust hypoxia gene signature
33 which was prognostic, even in multivariate analysis and across multiple tissue types.

34 Our approach empirically derives these seed genes using differential gene expression analysis, comparing cisplatin-sensitive
35 and -resistant cell lines from the Genomics of Drug Sensitivity in Cancer (GDSC) database. The seed genes are then trimmed
36 based on co-expression in epithelial-based tumor samples from The Cancer Genome Atlas (TCGA) ensuring that the final
37 signature contains genes that tend to be expressed together in both cell lines and clinical samples. This novel method may be
38 used to extract gene expression signatures for any quantitative or binary phenotype, and here, we will demonstrate its utility

39 with the extraction and validation of the Cisplatin Response Signature (CisSig), for use in predicting response to cisplatin in
 40 epithelial-origin tumors (carcinomas). We then show that our final signature is highly predictive of drug response within GDSC
 41 cell lines. And finally, we establish that signature expression in independent datasets of clinical tumor samples is congruent
 42 with use of cisplatin in standard of care guidelines between disease sites.

43 Results

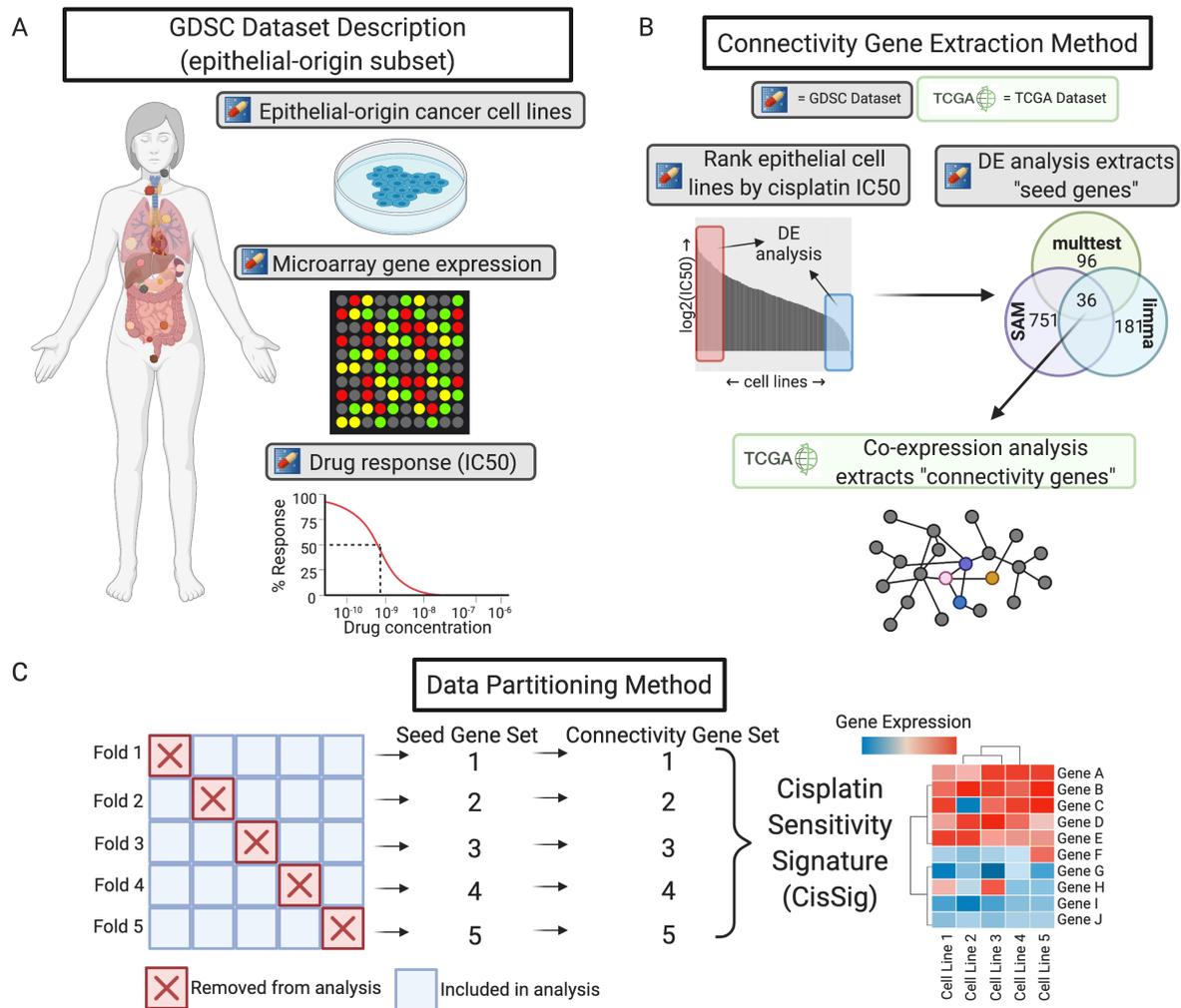


Figure 2. Schematic representation of CisSig derivation. **A.** Description of the epithelial-origin subset of the Genomics of Drug Discovery in Cancer (GDSC) dataset (denoted with the pill icon in future figures). These data include 429 epithelial-based cancer cell lines, with drug response measurements to over 200 drugs and gene expression characterization via microarray. **B.** Pipeline for extracting connectivity seeds. First, differential gene expression analysis between the top and bottom 20% of cisplatin responders found genes with significantly increased expression in a state of cisplatin sensitivity. These differentially expressed genes became "seed genes" in a co-expression network built using gene expression from clinical samples of epithelial-based tumors in The Cancer Genome Atlas (TCGA). Seed genes that were highly co-expressed with each other were denoted as "connectivity genes." **C.** Schematic of data partitioning, where GDSC epithelial-based cancer cell lines from **A.** are split into 5 folds. Each fold underwent the pipeline in **B.** Genes found in at least 3 of the 5 connectivity gene sets were included in the final signature, CisSig.

44 Convergent evolution informs Cisplatin Response Signature (CisSig) derivation

45 CisSig was derived using 429 epithelial-based cancer cell lines in the GDSC Database, each characterized for gene expression
 46 and drug response (see **Figure 2A**). The distribution of disease sites for these cell lines may be found in **Supplementary**

47 **Table 2.** GDSC gene expression consists of RMA normalized microarray data, details discussed in Methods. This database
48 reports both half-maximal inhibitory concentration (IC50) and area under the drug response curve (AUC) as measures of drug
49 response. A Spearman correlation between these two metrics demonstrated reasonable concordance ($\rho = 0.84$, $p < 0.001$) in
50 measuring cisplatin response for our cell lines of interest (**Supplementary Fig. 1**). We therefore moved forward with IC50 as
51 the metric of drug response, as it is a more commonly reported measure.

52 The GDSC epithelial cell lines were partitioned into five folds (each containing 343 or 344 cell lines) with a different
53 20% of the cell lines removed, illustrated in **Figure 2C**. Each of these folds was analyzed with a pipeline of differential gene
54 expression and co-expression analysis, visually depicted in **Figure 2B** and discussed below. This pipeline was performed across
55 multiple partitions of the data in order to find genes that are consistent between folds, reducing the chance for outlier cell lines
56 to influence the results.

57 With no pre-filtering of genes, differential gene expression (DE) analysis using limma,¹⁷ SAM,¹⁸ and multtest¹⁹ methods
58 was performed between the top and bottom 20% of responders (i.e. cell lines with the highest and lowest 20% of IC50
59 values). The distribution of disease sites found in each comparison group (resistant and sensitive) for each fold may be found
60 in **Supplementary Tables 2-6**. More details on parameters and version numbers for each DE method can be found in the
61 Methods section. For each fold, the genes found to be over-expressed in a cisplatin-sensitive state by all three DE methods were
62 termed the “seed genes,” resulting in 5 sets of seed genes, as depicted in **Figure 2C**. Using only intersecting genes between the
63 three methods is done with the goal of increasing stringency by reducing false discovery rate. Results of the DE analysis for
64 each fold are summarized in **Supplementary Table 7**, and lists of differentially expressed genes from each method, for each
65 fold can be found in Supplementary Data.

66 A co-expression network was built for each set of seed genes, as described in Methods and visually represented in the
67 bottom panel of **Figure 2B**. These networks were built using The Cancer Genome Atlas (TCGA) RNA-Seq expression data
68 from epithelial-based tumor samples, comparing the expression of each seed gene and all other genes in the dataset. Seed genes
69 that were highly co-expressed with each other are extracted from each fold, termed “connectivity seeds.” Here, we bring in
70 gene expression from tumor samples (not cell lines) to ensure that only genes that are expressed together in both cell lines and
71 tumor samples are included in the final signature. The final gene signature, CisSig, contains any gene found in at least 3 of the
72 5 sets of connectivity seeds, and the genes included in the signature are listed in **Table 1**.

73 Using the ‘sigQC’ package in R, we analyzed a suite of quality control metrics to assess the robustness of CisSig in a clinical
74 sample (TCGA) dataset.^{20,21} The signature is compared to the 5 sets of seed genes originally extracted from GDSC prior
75 to being refined by co-expression analysis. These results are visualized in a radar plot in **Supplementary Figure 2**. CisSig
76 demonstrates greater intra-signature correlation, increased correlation between mean and median, and decreased skewness
77 within RNA-expression from TCGA samples of epithelial origin. Other metrics of interest include the coefficient of variance
78 and the proportion (σ) of signature genes found in the top 10%, 25% or 50% of variable genes. These metrics can be used
79 to assess the variability of signature genes within a dataset, where it is ideal to have signature genes that vary more than the
80 background noise. Here, CisSig performs similarly to the unfiltered differential gene expression results. Finally, these
81 metrics are summarized into a score, also displayed in **Supplemental Figure 2**, where CisSig outperformed all sets of seed
82 genes.

83 **Increased CisSig expression predicts cisplatin sensitivity within GDSC dataset.**

84 **Figure 3A** demonstrates the expression of CisSig genes in cisplatin-sensitive and -resistant GDSC cell lines (top and bottom
85 IC50) quintiles. From this, we see that signature expression tends to be higher (more red) in sensitive, rather than resistant, cell
86 lines. Next, a “CisSig score,” the median normalized expression of the 19 CisSig genes, is calculated for the same sensitive and
87 resistant cell lines. The distribution of CisSig score and IC50 among all cell lines can be found in **Supplementary Figure 3**.
88 **Figure 3B** shows that sensitive cell lines tend to have higher CisSig scores than resistant cell lines. This is expected, given that
89 the seed genes were initially extracted as genes with increased expression in a cisplatin-sensitive state in the GDSC dataset.

90 **Figure 3C** compares the distribution of IC50 between cohorts of GDSC cell lines in this top and bottom quintile of CisSig
91 score. We are terming this plot a “Cell Line Persistence Curve,” which resembles a Kaplan-Meier survival curve, but uses IC50
92 in place of survival time for cell lines. Here, we assume that a cell line does not “survive” when the concentration of cisplatin is
93 greater than its IC50. For example, at 50% “survival” on the y-axis, the median IC50 of the high CisSig cohort is 2.76 μM
94 (left, vertical dashed line), while the median IC50 of the low CisSig cohort is 5.15 μM (right, vertical dashed line). In other
95 words, cell lines predicted to be resistant (low CisSig) tend to have greater IC50 values and cell lines predicted to be sensitive
96 (high CisSig) tend to have lower IC50 values.

97 As demonstrated by Venet et al, many published gene signatures do not perform significantly better when predicting survival
98 outcomes than random gene signatures of the same length¹¹. Given the large sample size of cell lines, simply testing for
99 statistical significance may not be stringent enough. Therefore, we compared the performance of CisSig’s Cell Line Persistence
100 Curve (hazard ratio) to the performance of a null distribution. This null distribution was created using 1000 random gene

Table 1. Genes included in CisSig. These genes all appear in at least 3 of the 5 sets of connectivity seeds.

HGNC Gene Symbol	Gene Name
<i>ADAT2</i>	Adenosine Deaminase tRNA Specific 2
<i>ATP1B3</i>	ATPase Na ⁺ /K ⁺ transporting subunit beta 3
<i>CDIN1</i>	CDAN1 interacting nuclease 1
<i>CIQBP</i>	Complement C1q binding protein
<i>CDC7</i>	Cell division cycle 7
<i>CDCA7</i>	Cell division cycle associated 7
<i>FKBP14</i>	FKBP prolyl isomerase 14
<i>KRT5</i>	Keratin 5
<i>LRR8C</i>	Leucine rich repeat containing 8 VRAC subunit C
<i>LY6K</i>	Lymphocyte antigen 6 family member K
<i>MMP10</i>	Matrix metalloproteinase 10
<i>NPM3</i>	Nucleophosmin 3
<i>PSAT1</i>	Phosphoserine aminotransferase 1
<i>RIOK1</i>	RIO kinase 1
<i>SLFN11</i>	Schlafen family member 11
<i>STOML2</i>	Stomatin like 2
<i>USP31</i>	Ubiquitin specific peptidase 31
<i>WDR3</i>	WD repeat domain 3
<i>ZNF750</i>	Zinc finger protein 750

101 signatures with the same length as CisSig, assessing the hazard ratio between each signature's Cell Line Persistence Curve. In
102 **Figure 3D**, we see that CisSig drastically outperforms the top 95% of this null distribution.

103 **CisSig outperforms the null distributions of drug response prediction models in the GDSC dataset.**

104 In **Figure 3C-D**, we demonstrated a novel method to show the stark difference in IC₅₀ distribution for GDSC cell lines with
105 high and low CisSig scores, but it is also important to assess CisSig's predictive power using more traditional methods. To that
106 aim, we built a variety of prediction models using CisSig to predict IC₅₀ as a continuous or binary outcome in epithelial-based
107 GDSC cell lines, described in **Table 2**. We chose to evaluate the efficacy of using a summary score (CisSig score) in addition to
108 individual gene expression in order to show the value of more "basic" statistical models (e.g. simple linear regression) for
109 producing an easier to interpret model while also gauging the power of using individual CisSig genes in accurately predicting
110 drug response (e.g. random forest). When utilizing expression of each gene individually as the input for our models, we chose
111 a penalized form of regression to prevent overfitting. Finally, for each method selected, we chose to build two models, one
112 with all epithelial-based cell lines and one with only epithelial-based cell lines with high or low signature expression (based on
113 CisSig score quintiles). In doing so, we can gauge whether more extreme expression of CisSig is related to improved drug
114 response prediction accuracy.

115 In short, simple linear regression models used CisSig score to predict a cell line's IC₅₀ as a continuous variable, while
116 elastic net, L1-, and L2-penalized linear regression models used expression of all CisSig genes to predict a cell line's IC₅₀ as a
117 continuous variable. For these linear regression models, performance was compared using the Spearman correlation coefficient
118 (ρ) between the predicted and actual IC₅₀ value for the cell lines withheld from a given fold's training dataset. The best
119 correlation coefficient between the five folds is chosen to represent each model, shown in **Table 2**. Simple logistic regression
120 models used CisSig score to predict a cell line's IC₅₀ as a binary outcome (above or below the median), while elastic net-,
121 L1-, and L2-penalized logistic regression, support vector machine (with linear and polynomial kernels), and random forest
122 models were built to use expression of each CisSig gene to predict IC₅₀ as a binary outcome. We used area under the receiver
123 operating characteristic (ROC) curve (AUC) to represent each classification model's performance, again choosing the best of
124 five folds to represent the model in **Table 2**.

125 As expected, all models demonstrate improved performance when trained and tested on only cell lines with the highest
126 and lowest signature scores. Additionally, the penalized regression models outperform the simple regression models when
127 comparing the same cell line data inputs. It is expected that including CisSig genes as individual variables would improve
128 performance in comparison to CisSig score, but it is noteworthy that something as simple as median normalized expression of
129 all CisSig genes (also known as the CisSig score) could predict IC₅₀ with the performance shown here.

130 **Figure 4** shows the performance of CisSig for each of the modeling methods described in **Table 2**. In **Figures 4A-B**, we

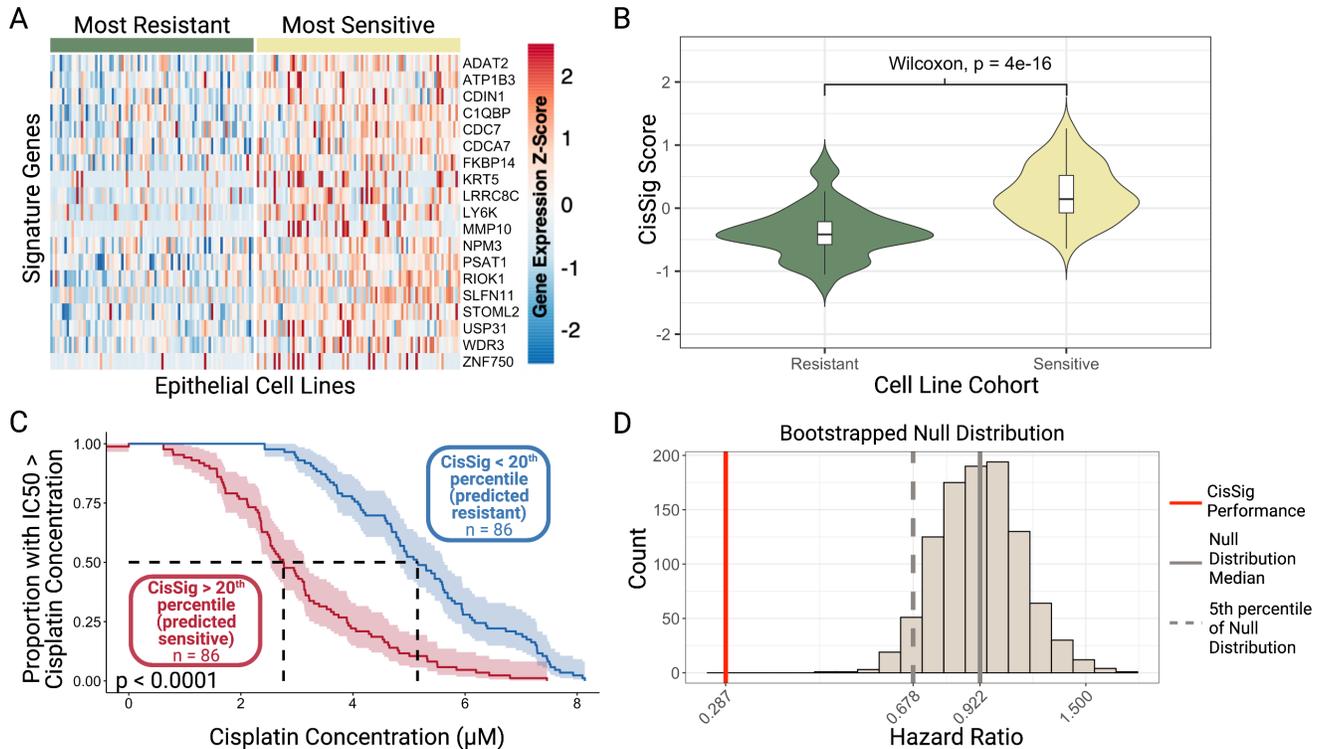


Figure 3. Visualization of CisSig expression within GDSC Dataset. **A.** An unclustered heatmap showing gene expression of the CisSig genes (rows) in cell lines (columns) from the top and bottom quintiles of cisplatin IC50. Color of the heatmap represents the Z-score of gene expression, normalized to each gene. Cell lines denoted as sensitive (right, yellow bar) tend to display higher expression of CisSig genes than cell lines denoted as resistant (left, green bar). Z-scores above 2.5 are denoted as 2.5, and Z-scores below -2.5 are denoted as -2.5. **B.** Violin plots comparing the distribution of CisSig scores between the cell lines in the highest and lowest quintile of cisplatin IC50. A Wilcoxon Rank Sum Test found that the median CisSig scores between these two cohorts was significantly different ($p < 0.001$). **C.** Comparison of the distribution cisplatin IC50 between cell lines in the highest and lowest quintile of CisSig score. Y-axis represents the proportion of the cohort with a cisplatin IC50 greater than the cisplatin concentration on the X-axis. A log-rank test between the two cohorts demonstrates significantly different drug response between the two cohorts ($p < 0.0001$). **D.** Null distribution of hazard ratio using 1000 random gene signatures with the same length as CisSig and the model described in C. CisSig’s performance is compared to the 95% confidence interval of the null distribution, where each signature’s performance (CisSig and nulls) is represented by the hazard ratio between two cohorts separated by the signature score.

131 demonstrate how each of the violin plots in **Figures 4C-D** were built. For example, in **Figure 4A**, we assess a linear regression
 132 model with CisSig score from all epithelial-based GDSC cell lines as the input and IC50 as the continuous outcome. Each
 133 model is built with five-fold cross validation, and performance is measured by comparing the predicted and actual IC50 of the
 134 testing set using a Spearman correlation. The best performance of the five-folds is used to represent CisSig’s performance,
 135 shown in **Figure 4A**. Next, a null distribution, shown in **Figure 4B**, is produced using 1000 random gene signatures with the
 136 same length as CisSig and the same modeling method. Again, the best performance of the five-folds is used to represent each
 137 null signature’s performance, and CisSig is compared to the null distribution.

138 We repeated the modeling described in **Figures 4A-B** for 10 additional modeling methods and the two versions of the
 139 dataset (one including all cell lines and another including only cell lines in the top and bottom quintile of signature expression).
 140 In **Figures 4C-D**, we show that CisSig outperforms these null distributions for each of the 11 modeling methods using both
 141 versions of the dataset, often outperforming the null distribution altogether. Finally, **Supplementary Figures 4-14** presents
 142 CisSig’s performance in each of the cross validation folds and show a detailed histogram of each model’s null distribution.

143 It is important to note that the wide variety of modeling methods shown here demonstrate that no one method is predictably
 144 superior to another, and CisSig shows strong predictive power when using any of them. Models that include only cell lines
 145 with more extreme signature expression tend to have improved performance compared to the same modeling method that includes
 146 all cell lines. This intimates that more extreme CisSig expression can more accurately predict a cell line’s response to cisplatin.

It is made available under a [CC-BY-NC 4.0 International license](https://creativecommons.org/licenses/by-nc/4.0/).

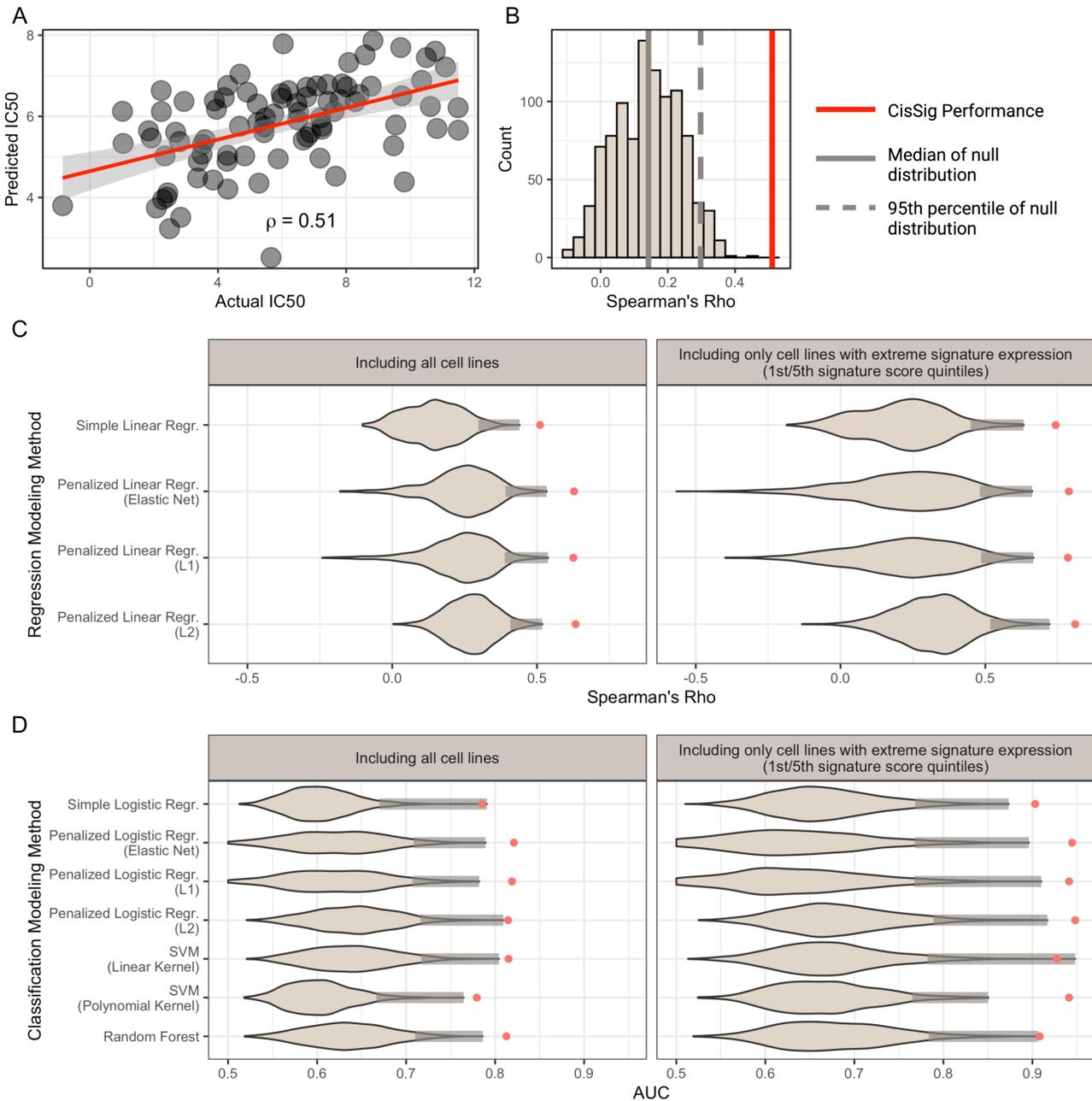


Figure 4. CisSig predicts IC50 using a variety of modeling techniques in the GDSC dataset. **A.** Scatterplot of the actual vs. predicted IC50 using CisSig score to predict IC50 with linear regression. Plot shows the best performing fold (measured by Spearman's rho) from 5-fold cross validation. **B.** Null distribution of the performance metric from **A.** (Spearman's rho), built using 1000 random gene signatures to predict IC50 as described in **A.** As with CisSig, the metric of the best performing fold is used to represent each null signature. The median of the null distribution and the cutoff for the 95th percentile of the null distribution are represented by the solid and dashed gray line, respectively. CisSig's performance, red solid line, outperforms at least 95% of the null distribution. **C.** Violin plots containing the null distribution of performance metrics for 11 modeling methods. Each distribution was created as discussed in **A-B**, where CisSig's performance is compared to the performance of 1000 random gene signatures of the same length. For each violin, a shaded gray bar represents the top 5% of each null distribution and CisSig's performance is shown with a red dot. The modeling methods, including input and output, are described in Table 2.

Table 2. Model details and validation results for the prediction of cisplatin response using CisSig in GDSC dataset.

Input	Output	Method	Included Data	Metric	Value
CisSig Score	IC50 (continuous)	Simple Linear Regression	All	Corr. Coef.	0.51
CisSig Score	IC50 (continuous)	Simple Linear Regression	Quintiles	Corr. Coef.	0.74
All gene expression	IC50 (continuous)	Elastic Net Linear Regression	All	Corr. Coef.	0.63
All gene expression	IC50 (continuous)	Elastic Net Linear Regression	Quintiles	Corr. Coef.	0.79
All gene expression	IC50 (continuous)	L1 Linear Regression	All	Corr. Coef.	0.63
All gene expression	IC50 (continuous)	L1 Linear Regression	Quintiles	Corr. Coef.	0.79
All gene expression	IC50 (continuous)	L2 Linear Regression	All	Corr. Coef.	0.63
All gene expression	IC50 (continuous)	L2 Linear Regression	Quintiles	Corr. Coef.	0.81
All gene expression	IC50 (binary)	Simple Logistic Regression	All	AUC	0.79
All gene expression	IC50 (binary)	Simple Logistic Regression	Quintiles	AUC	0.90
All gene expression	IC50 (binary)	Elastic Net Logistic Regression	All	AUC	0.82
All gene expression	IC50 (binary)	Elastic Net Logistic Regression	Quintiles	AUC	0.94
All gene expression	IC50 (binary)	L1 Logistic Regression	All	AUC	0.82
All gene expression	IC50 (binary)	L1 Logistic Regression	Quintiles	AUC	0.94
All gene expression	IC50 (binary)	L2 Logistic Regression	All	AUC	0.81
All gene expression	IC50 (binary)	L2 Logistic Regression	Quintiles	AUC	0.95
All gene expression	IC50 (binary)	SVM (linear kernel)	All	AUC	0.82
All gene expression	IC50 (binary)	SVM (linear kernel)	Quintiles	AUC	0.93
All gene expression	IC50 (binary)	SVM (polynomial kernel)	All	AUC	0.78
All gene expression	IC50 (binary)	SVM (polynomial kernel)	Quintiles	AUC	0.94
All gene expression	IC50 (binary)	Random Forest	All	AUC	0.81
All gene expression	IC50 (binary)	Random Forest	Quintiles	AUC	0.91

147 **Ranking cancer subtypes by CisSig expression is concordant with observed clinical trends.**

148 The consistently strong validation statistics displayed in **Figures 3** and **4** demonstrate that this novel signature extraction
 149 methodology is capable of selecting genes with strong predictive power within the source dataset. In other words, it is a
 150 powerful tool for feature selection. In order to assess translation into novel datasets; however, predictive power must be
 151 demonstrated in datasets that were not used to select genes of interest.

152 Using three large datasets, we assessed how expression of CisSig relates to cisplatin use across epithelial-based cancer
 153 disease sites. CisSig score was calculated for all samples (cell lines or clinical tumor samples) in GDSC, TCGA, and Total
 154 Cancer Care (TCC) databases. In order to visualize these scores on a log-transformed axis, signature score was linearly scaled,
 155 such that the lowest score became exactly 1.

156 In **Figure 5**, disease sites were ranked by the median signature score for the cohort in GDSC (left), TCGA (middle), and
 157 TCC (right) datasets. Furthermore, each disease site is labeled as utilizing cisplatin in NCCN treatment guidelines (green circle),
 158 using cisplatin in very select circumstances (yellow bars), or not having cisplatin included in NCCN treatment guidelines
 159 (red square). In all datasets, we see that disease sites with higher CisSig scores tend to have cisplatin included in treatment
 160 guidelines, while those with lower scores tend to not have cisplatin included in treatment guidelines.

161 Finally, disease site rank was compared between datasets using Spearman’s correlation. There is a strong correlation
 162 between the rank of shared disease sites of all three datasets. Between GDSC and TCGA, Spearman’s ρ is 0.77 ($p < 0.001$).
 163 Between GDSC and TCC, Spearman’s ρ is 0.92 ($p < 0.001$). And between TCGA and TCC, Spearman’s ρ is 0.93 ($p < 0.001$).
 164 This high degree of concordance between datasets signifies that CisSig displays consistent expression between a variety of data
 165 sources (including between microarray and RNA-seq methods).

166 **Discussion**

167 The principles of convergent evolution tell us that genetically distant organisms can evolve similar traits in order to become
 168 more fit under the same selection pressure. In cancer, therefore, we cannot ignore the possibility that different mutations may
 169 lead to the same phenotype. Therefore, our novel method groups convergent phenotypes and uses expression profiling to better
 170 predict drug response in cancer. In doing so, we harnessed the power of over 400 epithelial-origin cell lines in the GDSC
 171 Database to extract CisSig, a gene expression signature for use in predicting cisplatin response in epithelial-origin tumors.

172 Gene expression signatures can expand the reach of precision medicine to impact the vast majority of patients whose tumors

It is made available under a [CC-BY-NC 4.0 International license](https://creativecommons.org/licenses/by-nc/4.0/).

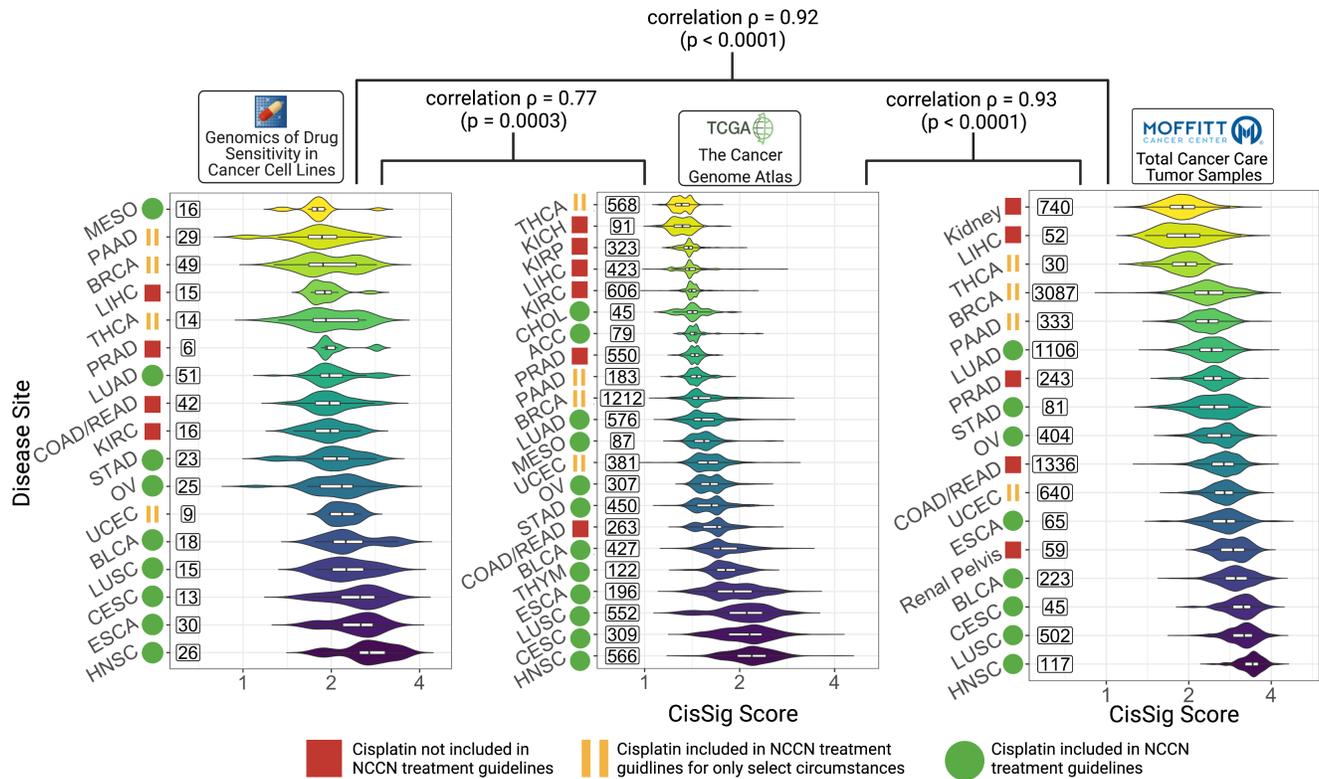


Figure 5. Cancer subtypes with greater CisSig expression tend to have cisplatin included in standard of care guidelines.

Cancer subtypes are ranked by median CisSig Score in three data sets, GDSC (left), TCGA (middle), and TCC (right). The color of each violin plot represents the rank of the cancer subtype. The ranks of intersecting subtypes between each dataset are compared with Spearman's rank correlation, reported with correlation ρ and p-value. Rank correlation ρ between GDSC and TCGA and GDSC and TCC datasets is 0.77 ($p = 0.0003$) and 0.902 ($p < 0.0001$). Rank correlation ρ between TCGA and TCC datasets is 0.93 ($p < 0.0001$). Violin plots display the distribution of CisSig scores for each cancer subtype. Within each violin, a boxplot denotes median signature score for each subtype (middle horizontal line) and 25th/75th percentile for signature scores (box edges). Numbers to the left of each violin plot represent sample size included in each cancer subtype.

do not harbor actionable mutations. Yet, finding signatures with significant translational potential remains difficult. This is in part because although cell lines are preferable for high throughput analysis, they tend to demonstrate some divergence from their tumors of origin.^{22,23}

As demonstrated by many predictive modeling methods, our gene signature is highly effective at predicting drug response in GDSC cell lines, from which it was originally derived. This initial validation is an important step, but demonstrating utility with independent clinical samples is crucial for assessing the translational potential of our signature. Unlike with cell lines, high throughput characterization of drug response (i.e. IC50, AUC, etc) in clinical tumor samples is not feasible.²² Because of this, many researchers use survival as a surrogate measure of treatment response for tumor samples. However, without a known clinical history of cisplatin treatment, we cannot use survival as a surrogate measure of cisplatin response.

As such, we chose to assess how well CisSig expression in tumor samples correlates to clinical treatment trends. With this analysis, we show that cancer subtypes frequently treated with cisplatin (e.g. head and neck, cervical) tend to have greater CisSig scores in GDSC, TCGA, and TCC datasets. GDSC was directly used in the extraction of CisSig and TCGA is used only for co-expression analysis in trimming the signature genes, but the TCC database was not used in any part of the extraction methodology. And although predictive modeling within this independent dataset is not feasible, this result does show that expression of CisSig tends to be congruent with current clinical practices, an indication that CisSig has translational potential.

This signature extraction method is, of course, not without limitations. First, a single tumor sample may not capture the intratumoral heterogeneity that is crucial for predicting the physiological response to a drug. Next, although the signature was extracted to find genes with importance across pan-cancer (epithelial-based) tumor subtypes, clinical validation must occur within individual disease sites. Given the heterogeneity between tumor subtypes, disease-site specific versions of CisSig may

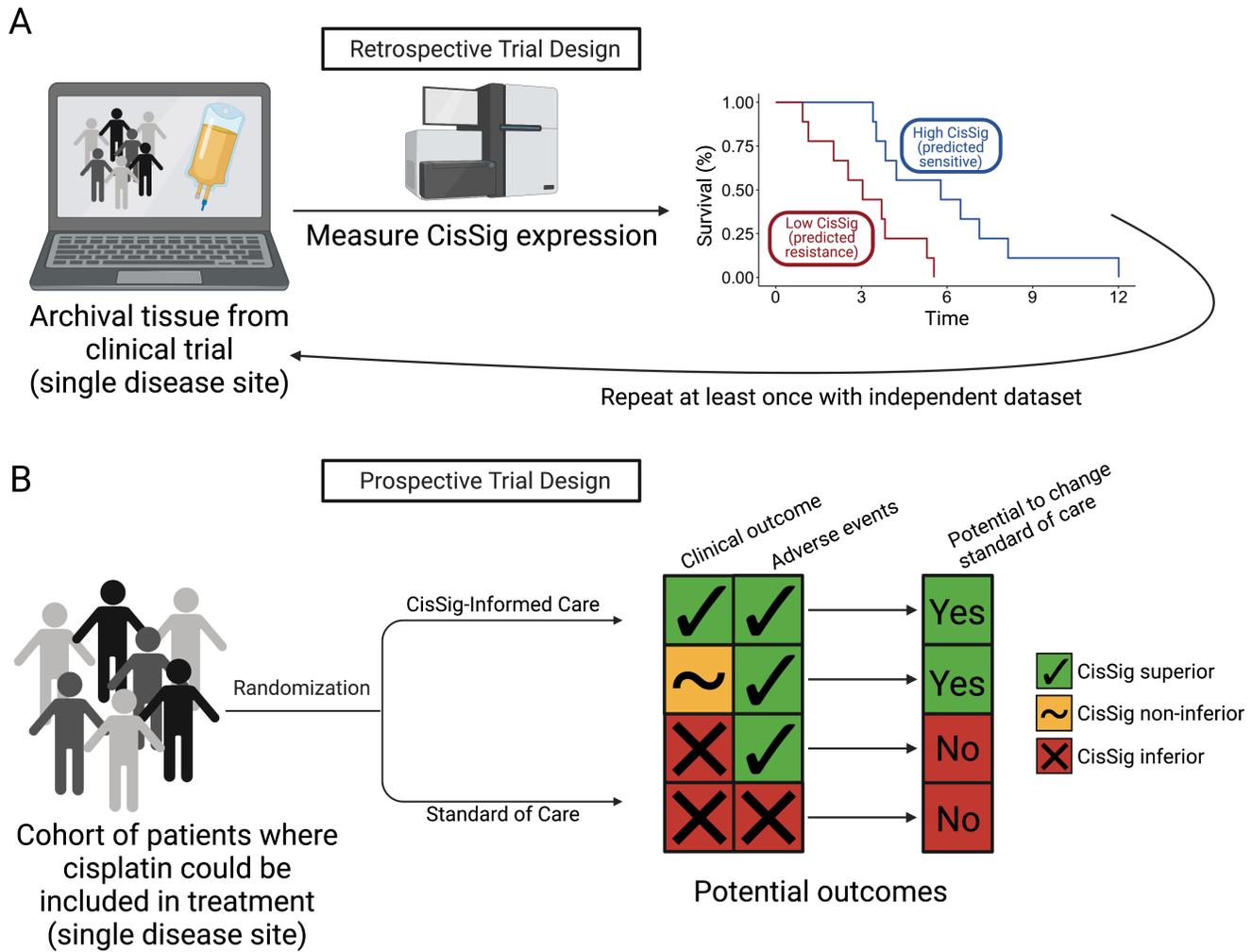


Figure 6. Pathways for CisSig to reach level 1 evidence. A. Retrospective trial design. Measure CisSig expression in archival tissue of a single disease site from a prior clinical trial. Determine if cohorts of high vs. low CisSig expression have significantly different survival trends. Repeat at least once with an independent dataset to reach level 1 evidence. **B. Prospective trial design.** Begin with a cohort of patients with a single cancer subtype where cisplatin may or may not be included in their treatment plan. Patients are randomized into a CisSig-Informed Cohort or a Standard of Care Cohort. In the CisSig-Informed Cohort, clinicians will be informed on CisSig expression and what it means regarding predicted therapeutic response. They will use this information when counseling the patient in deciding between therapeutic options. The Standard of Care Cohort will not receive any information regarding a patient's CisSig expression. The two cohorts will be compared regarding clinical outcome, adverse events, and other factors. A variety of differences between the two cohorts could lead to level 1 evidence for the use of CisSig in that disease site.

192 require trimming the genes of this pan-cancer signature even further. Additionally, as discussed previously, using cell line
 193 expression data as the basis of a clinical signature is necessary given the current limitations of high throughput databases, but it
 194 can hinder translation. Therefore, a key future direction will be testing the signature in clinical data to determine if patient
 195 response to cisplatin can be stratified by signature expression.

196 **Figure 6** shows two pathways that could demonstrate the successful clinical translation of CisSig, providing level 1
 197 evidence for its use. First, a retrospective trial design, displayed in **Figure 6A**, could take archival tissue from a clinical
 198 trial of an epithelial-based disease site (e.g. squamous lung cancer, cervical cancer, etc.) where all patients have undergone
 199 cisplatin-containing treatment. After CisSig expression is measured for all samples, patients will be separated into cohorts of
 200 high and low CisSig expression, where it is predicted that the high CisSig cohort will have improved survival. According to
 201 Burns, et al²⁴, this retrospective trial design must be completed with at least two independent datasets to reach level 1 evidence.

202 Next, **Figure 6B** presents a prospective trial design for assessing the utility of CisSig in a clinical setting. Starting with

203 a cohort of patients where cisplatin could be included in the treatment plan, patients would be randomized into the CisSig-
204 informed care or standard of care treatment cohort. With CisSig-informed care, clinicians will be informed of the patient's
205 predicted response to cisplatin and encouraged to utilize this information when determining whether cisplatin should be
206 included in the patient's treatment plan. There are many possible outcomes of this trial design, a few of which are demonstrated
207 in **Figure 6B**. Ideally, CisSig-informed care will improve both survival outcomes and adverse events. Yet, there would still
208 be level 1 evidence to support the use of CisSig even if survival outcomes were non-inferior, given that adverse events are
209 improved in the CisSig-informed care cohort.

210 Selection, like drug treatment acts on phenotype. And in this work, we demonstrate a novel gene signature extraction
211 method-informed by principles of convergent evolution-where we find shared transcriptomic markers of drug response
212 phenotype in tumors that appear genotypically disparate. By harnessing the power of a large dataset, such as the GDSC, we
213 extracted a biologically-inspired product, CisSig. Expanding this method to produce signatures for response prediction to a
214 variety of chemotherapeutic agents will lead to a monumental expansion of precision medicine in cancer.

215 **Methods**

216 **Data Collection and Pre-Processing**

217 All data cleaning, analysis, and plotting was performed using R (Version 4.0.5) with RStudio.

218 ***GDSC Gene Expression Data***

219 Microarray mRNA expression, drug response, and meta-data for 983 cell lines and 251 drugs was downloaded from the
220 Genomics in Drug Sensitivity Database (GDSC)²⁵. The expression and meta-data were last updated 4 July 2016. The GDSC
221 database can be accessed at <https://www.cancerrxgene.org/>. Documentation for the GDSC database states that the RMA
222 normalized^{26,27} expression data for all cell lines were collected via Human Genome U219 96-Array Plate using the Gene
223 Titan MC instrument (Affymetrix). Further the robust multi-array analysis (RMA) algorithm was used to normalize the data,
224 reporting intensity values for 18562 individual loci. The raw data and probe ID mappings were deposited in ArrayExpress
225 (accession number: E-MTAB-3610). The RMA processed dataset is available at <http://www.cancerrxgene.org/gdsc1000/>.

226 Epithelial-based cell lines are extracted based on the following GDSC tissue descriptors (exact labels found in database):
227 head and neck, oesophagus, breast, biliary_tract, large_intestine, liver, adrenal_gland, stomach, kidney, lung_NSCLC_adenocarcinoma,
228 lung_NSCLC_squamous-_cell_carcinoma, mesothelioma, pancreas, skin_other, thyroid, Bladder, cervix, endometrium, ovary,
229 prostate, testis, urogenital_system_other, uterus.

230 ***GDSC Drug Response Data***

231 The drug response data in the GDSC database was last updated 27 March 2018; this version is referred to as "GDSC2." Cisplatin
232 drug concentration is reported in μM . Raw viability data were processed using the R package, *gdscIC50*, where they were
233 normalized with negative controls (media alone) and positive controls (media only wells with no cells). Dose-response curves
234 were fit using a multi-level fixed effect model with a classic sigmoidal curve shape assumed. This model was fitted using all
235 cell line/drug combinations that were screened instead of fitting separate models to individual drug-response series. In this
236 approach, the shape parameter only changes between cell lines, but the position parameter is adjusted between cell lines and
237 compounds. Additional information regarding dose-response curve fitting may be found at Vis et al.²⁸. Fitting models to all
238 dose-response series leads to improved robustness for more accurate IC50 and AUC estimates.

239 ***TCGA Gene Expression Data***

240 RNA-Seq by Expectation Maximization (RSEM) normalized gene expression for epithelial-based cancers was downloaded from
241 The Cancer Genome Atlas (TCGA) database, which was accessed through the Firebrowse database using the 'RTCGAToolbox'
242 package (version 2.20.0)²⁹ in R. The following TCGA Study Abbreviations were downloaded (exact labels found in database):
243 ACC, BLCA, BRCA, CESC, CHOL, COADREAD, ESCA, HNSC, KIRC, KIRP, KICH, LIHC, LUAD, LUSC, MESO, OV,
244 PAAD, PRAD, STAD, THCA, THYM, UCEC. These values were measured through the Illumina HiSeq RNAseq V2 platform
245 and were log2 transformed.

246 ***Total Cancer Care (TCC) Gene Expression Data***

247 The Total Cancer Care Dataset is collected by the H. Lee Moffitt Cancer Center and Research Institute using protocols described
248 in Fenstermacher et al.^{30,31}. The Total Cancer Care (TCC) protocol is a prospective tissue collection protocol that has been
249 active at Moffitt Cancer Center (Tampa, FL, USA) and 17 other institutions since 2006. We assayed tumours from adult patients
250 enrolled in the TCC protocol on Affymetrix Hu-RSTA-2a520709, which contains approximately 60,000 probesets representing
251 25,000 genes. Chips were normalised using iterative rank-order normalisation.³² Batch effects were reduced using partial-least
252 squares. We extracted from the TCC database normalised, debatched expression values for 9,063 samples from 17 sites of

253 epithelial origin and the 19 CisSig genes. We excluded all metastatic duplicate samples and disease sites with fewer than 25
254 samples.

255 **Drug Response Quality Control**

256 IC50 is an imperfect measure of drug response, yet it is widely used throughout the literature. It is defined as the concentration
257 of drug at which cells experience 50% inhibitory effect. Another measure of drug response is area under the drug response
258 curve, which is defined as the integral of a drug response curve, where cellular activity is measured on the y-axis and drug
259 concentration is measured on the x-axis. IC50 and AUC values for all epithelial cell lines are compared using a Spearman
260 correlation test (see **Supplementary Figure 1**) in order to assess concordance between the two metrics.

261 **Differential Gene Expression Analysis**

262 As seen in **Figure 2C**, the GDSC dataset is split into 5-folds, where 20% of the cell lines are removed from further analysis for
263 each of the 5 runs. This leaves 343 or 344 cell lines in each of the 5 partitions. After data partitioning, the top 20% and bottom
264 20% are extracted for comparison using differential expression analysis, **Figure 2C**.

265 Differential expression analysis is performed using three algorithms: significance analysis of microarrays (SAM),
266 resampling-based multiple hypothesis testing, and linear models for microarrays (limma), which are implemented using
267 R packages ‘samr’¹⁸ (version 3.0), ‘multtest’¹⁹ (version 2.46.0), and ‘limma’¹⁷ (version version 3.46.0), respectively. Gene
268 expression was pre-normalized using RMA (discussed above) and genes were not pre-filtered before this analysis. This analysis
269 has 69 samples per group, which is appropriate given the demonstration by Baccarella et al. showing that differential expression
270 results begin to vary problematically beginning when there are as few as 8 samples per group³³.

271 A false discovery rate or p-value cutoff of 0.20 was chosen for each method. The ‘samr’ and ‘multtest’ method were both
272 set to the same seed. The ‘samr’ method used 10,000 permutations (parameter: “nperm”) and test statistic was set to “standard”
273 for t-test (parameter: “testStatistic”). The ‘limma’ method used no p-value adjustment method (parameter: “adjust.method”)
274 and a log-fold change cutoff of 0.5 (parameter: “lfc”). The ‘multtest’ method used 1,000 bootstrap iterations (parameter: “B”)
275 and single-step minP for multiple testing procedure (parameter: “method”). All other parameters for the three algorithms were
276 set to default. The intersection of the genes found to have significantly increased expression in sensitive cell lines by the three
277 algorithms is termed “seed genes” for use in future co-expression analysis. An FDR cutoff of 0.2 is a relatively non-stringent
278 FDR cutoff; it was chosen in order to include a variety of genes before taking the intersection of results between the three
279 methods.

280 **Co-Expression Network Analysis and Final Signature Derivation**

281 The co-expression network, represented in the pipeline of **Figure 2B**, is made by performing a pairwise Spearman correlation
282 between the expression of each seed gene and every other gene (including other seed genes) except itself. The correlation
283 coefficient for each pairwise comparison is termed the “affinity score.” Next, the network is transformed so that the largest
284 5% of affinity scores are transformed to 1 and all other scores become 0. This is done without squaring the scores in order to
285 extract only positive correlations. The average affinity score for each gene compared to each seed gene is then derived; this
286 value becomes known as a gene’s “connectivity score.” The intersection between the differentially expressed seed genes and
287 genes with the top 20% of the highest connectivity scores become known as the “connectivity genes.” Five sets of connectivity
288 genes are compiled, one for each data partition. The final signature (CisSig) is produced by extracting any gene that is found in
289 at least three of the five connectivity gene sets.

290 **Signature Quality Control in TCGA**

291 In order to examine how CisSig compares to the original differential gene expression results and ensure portability to novel
292 datasets, we perform a quality control analysis within the TCGA dataset using the ‘sigQC’ R package²⁰ with methodology as in
293 Dhawan et al. 2019²¹. Here, various metrics are calculated using the expression of the genes found in the gene expression
294 signature and the 5 sets of differential expression analysis results. These metrics include intra-signature correlation, correlation
295 between the mean expression and first principal component, and skewness of the signature expression. The final results of all
296 the metrics calculated for each signature are displayed in a radar plot, with a summary score of each set of genes (signature)
297 tested. This summary score is the ratio of the area within the radar plot and the full polygon if each metric was the highest
298 value possible.

299 **Predicting cell line IC50 using CisSig in GDSC**

300 A cell line or sample’s median normalized expression value of the CisSig genes is termed the CisSig score. Cell lines were
301 again organized into five folds (independent of the data partitioning used in the signature extraction, described in **Figure 2C**).
302 Predictive models were built using 80% of the cell lines (training cell lines) and tested on the 20% of the cell lines withheld
303 from the model (validation cell lines). All models were built with two versions of input—one using all of the epithelial-based

304 cell lines in the GDSC database and the other using only the cell lines in the top and bottom quintiles of CisSig score. When
305 using all the epithelial-based cell lines, training sets consist of 344-345 cell lines, while testing sets consist of 86 cell lines.
306 When using only the cell lines in the top and bottom quintiles for signature expression, training sets consist of 137 or 138 cell
307 lines and testing sets consist of 34 or 35 cell lines.

308 Simple linear and logistic regression was used to predict IC50 as a continuous variable with CisSig score as the input.
309 Elastic net, L1-, and L2-penalized linear regression methods utilized the expression of each of the 19 CisSig genes to predict
310 IC50 as a continuous variable. Elastic net, L1-, and L2-penalized logistic regression methods, support vector machine (SVM),
311 and random forest methods utilized expression of each of the 19 CisSig genes to predict IC50 as a binary variable (above
312 or below the median of the group). All linear regression models were evaluated using the Spearman correlation coefficient
313 between true and predicted IC50 values from the validation set. Classification models (logistic regression, SVM, and random
314 forest) were evaluated using area under the receiver operating characteristic (ROC) curve (AUC).

315 Elastic net, L1-, and L2-penalized linear and logistic regression models were built using the 'glmnet' package (version
316 4.1-2) in R. The alpha parameter was set to 0.5, 1, and 0 for elastic net, L1-, and L2-penalized regression, respectively. Models
317 were tuned with 10-fold cross validation to choose a value for lambda with the best predictive capabilities based on mean square
318 error for linear models and misclassification error for logistic models.

319 SVM models were built with the 'e1071' package (version 1.7-8) in R, using both a linear and polynomial kernel. Models
320 were tuned with 10-fold cross validation to choose the best value for degree (from 3, 4, 5), gamma (from 10^{-3} , 10^{-2} , 10^{-1} , 1, 10^1 ,
321 10^2 , 10^3), and cost (from 10^{-3} , 10^{-2} , 10^{-1} , 1, 10^1 , 10^2 , 10^3).

322 Random forest models were built with the 'randomForest' package (version 4.6-14), and each model grew 500 trees. All
323 other parameters in training the prediction models were default.

324 Cell Line Persistence Curves

325 Cell lines with high CisSig scores (predicting the more sensitive cell lines) and low signatures scores (predicting the more
326 resistant cell lines) are separated by quintile. A Kaplan-Meier survival model is built for the two cohorts using IC50 in lieu of
327 survival time. A log-rank test compares the two survival curves to analyze if the two cohorts of signature expression are related
328 to different "survival" of higher IC50s in each group.

329 Null distributions of cell line IC50 models

330 CisSig's performance was compared to a null distribution for all models built, including all models used to predict IC50 as a
331 continuous or binary variable and the cell line persistence models using the log-rank test to compare the two survival curves. To
332 build each null distribution, 1000 random gene signatures with the same length as CisSig were chosen. Each random gene
333 signature was selected using all genes included in the GDSC expression profiling without replacement. The performance of
334 each random signature was tested in each individual modeling method, producing a null distribution for each modeling method.

335 As discussed above, the predictive models utilize five-fold cross validation and the best summary statistic of the five folds is
336 chosen to represent the signature's performance. This remains consistent for the null models, where the best summary statistic
337 of the five folds is used to represent each random signature. Again, all code for building the testing and null models may be
338 found in the GitHub repository listed in Code and Data Availability.

339 Ranking disease sites in GDSC, TCGA, and TCC by CisSig Score

340 All epithelial-origin cell lines or tumor samples in the GDSC, TCGA, and TCC datasets had CisSig Score calculated as
341 previously described. For the purposes of plotting on a log-scale, the scores were linearly adjusted by adding the absolute value
342 of the lowest score plus 1 to each sample's score, making the lowest score now 1. For example, if the lowest signature score
343 for the dataset was -5, 6 was added to each sample's score. Disease sites within each dataset were ranked by median CisSig
344 score. For disease sites shared between datasets, a Spearman correlation was performed to assess how the rank of disease sites
345 compare between datasets.

346 Classifying disease sites by cisplatin use

347 NCCN Treatment Guidelines for each disease site were manually searched, versions listed in **Supplementary Table 8**. Disease
348 sites were classified as including cisplatin in treatment guidelines, only including cisplatin in very select circumstances, or not
349 including cisplatin in treatment guidelines. For those classified as only using cisplatin in select circumstances, details are noted
350 in **Supplementary Table 8**.

351 Data and code availability

352 The code to download all data, extract CisSig, perform validation of the signature, and reproduce all figures in the manuscript is
353 available via GitHub at

354 <https://github.com/jessicascarborough/cissig>.

355 **Acknowledgements**

356 J.G.S. would like to thank NIH (5R37CA244613-02) and the American Cancer Society (RSG-20-096-01) for their generous
357 support. J.A.S. thanks the NIH for their support through the T32GM007250 and 1F30CA257076-01 grants. The results
358 published here are in whole or part based upon data generated by the TCGA Research Network: <https://www.cancer.gov/tcga>.
359 All authors are grateful to the cancer patients who provided tissue for further study in the GDSC, TCGA, and TCC datasets.
360 This work made use of the High Performance Computing Resource in the Core Facility for Advanced Research Computing at
361 Case Western Reserve University.

362 **Author contributions**

363 J.A.S. contributed to experimental design, wrote all associated code, analyzed data, and wrote the manuscript. A.D. contributed
364 to experimental design and analyzed data. S.A.E. analyzed data. J.T.R. contributed to experimental design. J.G.S. contributed
365 to experimental design, analyzed data, and wrote the manuscript. All authors read and approved of the manuscript.

References

- 366 **1.** Hirsch, F. R. *et al.* Lung cancer: current therapies and new targeted treatments. *The Lancet* **389**, 299–311 (2017).
- 367 **2.** Solomon, B. J. *et al.* First-line crizotinib versus chemotherapy in alk-positive lung cancer. *New Engl. J. Medicine* **371**,
- 368 2167–2177 (2014).
- 369 **3.** Prasad, V., De Jesus, K. & Mailankody, S. The high price of anticancer drugs: origins, implications, barriers, solutions.
- 370 *Nat. reviews Clin. oncology* **14**, 381 (2017).
- 371 **4.** Marquart, J., Chen, E. Y. & Prasad, V. Estimation of the percentage of us patients with cancer who benefit from
- 372 genome-driven oncology. *JAMA oncology* **4**, 1093–1098 (2018). PMID: PMC6143048.
- 373 **5.** Sparano, J. A. *et al.* Adjuvant chemotherapy guided by a 21-gene expression assay in breast cancer. *New Engl. J. Medicine*
- 374 **379**, 111–121 (2018).
- 375 **6.** Soliman, H. *et al.* Mammprint guides treatment decisions in breast cancer: results of the impact trial. *BMC cancer* **20**, 81
- 376 (2020).
- 377 **7.** Scott, J. G. *et al.* A genome-based model for adjusting radiotherapy dose (gard): a retrospective, cohort-based study. *The*
- 378 *lancet oncology* **18**, 202–211 (2017).
- 379 **8.** Scott, J. G. *et al.* Pan-cancer prediction of radiotherapy benefit using genomic-adjusted radiation dose (gard): a cohort-based
- 380 pooled analysis. *The Lancet Oncol.* **22**, 1221–1229 (2021).
- 381 **9.** Eschrich, S. A. *et al.* A gene expression model of intrinsic tumor radiosensitivity: prediction of response and prognosis
- 382 after chemoradiation. *Int. J. Radiat. Oncol. Biol. Phys.* **75**, 489–496 (2009).
- 383 **10.** Torres-Roca, J. F. A molecular assay of tumor radiosensitivity: a roadmap towards biology-based personalized radiation
- 384 therapy. *Pers. medicine* **9**, 547–557 (2012).
- 385 **11.** Venet, D., Dumont, J. E. & Detours, V. Most random gene expression signatures are significantly associated with breast
- 386 cancer outcome. *PLoS computational biology* **7**, e1002240 (2011).
- 387 **12.** Nichol, D. *et al.* Antibiotic collateral sensitivity is contingent on the repeatability of evolution. *Nat. communications* **10**,
- 388 1–10 (2019). PMID: PMC6338734.
- 389 **13.** Scarborough, J. A. *et al.* Identifying states of collateral sensitivity during the evolution of therapeutic resistance in ewing’s
- 390 sarcoma. *Iscience* **23**, 101293 (2020).
- 391 **14.** Dhawan, A. *et al.* Collateral sensitivity networks reveal evolutionary instability and novel treatment strategies in alk
- 392 mutated non-small cell lung cancer. *Sci. Reports* **7**, 1–9 (2017). PMID: PMC5430816.
- 393 **15.** Blount, Z. D., Lenski, R. E. & Losos, J. B. Contingency and determinism in evolution: Replaying life’s tape. *Science* **362**
- 394 (2018).
- 395 **16.** Buffa, F., Harris, A., West, C. & Miller, C. Large meta-analysis of multiple cancers reveals a common, compact and highly
- 396 prognostic hypoxia metagene. *Br. journal cancer* **102**, 428 (2010).
- 397 **17.** Ritchie, M. E. *et al.* limma powers differential expression analyses for rna-sequencing and microarray studies. *Nucleic*
- 398 *acids research* **43**, e47–e47 (2015).
- 399 **18.** Tusher, V., Tibshirani, R. & Chu, C. Significance analysis of microarrays applied to ionizing radiation response. *Proc. Natl.*
- 400 *Acad. Sci.* **98**, 5116–5121 (2001).
- 401 **19.** Pollard, K. S., Dudoit, S. & van der Laan, M. J. Multiple testing procedures: the multtest package and applications to
- 402 genomics. In *Bioinformatics and computational biology solutions using R and bioconductor*, 249–271 (Springer, 2005).
- 403 **20.** Dhawan, A., Barberis, A., Cheng, W.-C. & Buffa, F. *sigQC: Quality Control Metrics for Gene Signatures* (2018). R
- 404 package version 0.1.21.
- 405 **21.** Dhawan, A. *et al.* Guidelines for using sigqc for systematic evaluation of gene signatures. *Nat. Protoc.* **14**, 1377 (2019).
- 406 **22.** Azuaje, F. Computational models for predicting drug responses in cancer research. *Briefings bioinformatics* **18**, 820–829
- 407 (2017).
- 408 **23.** Goodspeed, A., Heiser, L. M., Gray, J. W. & Costello, J. C. Tumor-derived cell lines as molecular models of cancer
- 409 pharmacogenomics. *Mol. Cancer Res.* **14**, 3–13 (2016).
- 410 **24.** Simon, R. M., Paik, S. & Hayes, D. F. Use of archived specimens in evaluation of prognostic and predictive biomarkers. *J.*
- 411 *Natl. Cancer Inst.* **101**, 1446–1452 (2009).
- 412

- 413 **25.** Yang, W. *et al.* Genomics of drug sensitivity in cancer (gdsc): a resource for therapeutic biomarker discovery in cancer
414 cells. *Nucleic Acids Res.* **41**, D955–D961, DOI: [10.1093/nar/gks1111](https://doi.org/10.1093/nar/gks1111) (2013). [/oup/backfile/content_public/journal/nar/41/
415 d1/10.1093/nar/gks1111/2/gks1111.pdf](https://oup/backfile/content_public/journal/nar/41/d1/10.1093/nar/gks1111/2/gks1111.pdf).
- 416 **26.** Irizarry, R. A. *et al.* Exploration, normalization, and summaries of high density oligonucleotide array probe level data.
417 *Biostatistics* **4**, 249–264 (2003).
- 418 **27.** Iorio, F. *et al.* A landscape of pharmacogenomic interactions in cancer. *Cell* **166**, 740–754 (2016).
- 419 **28.** Vis, D. J. *et al.* Multilevel models improve precision and speed of ic50 estimates. *Pharmacogenomics* **17**, 691–700 (2016).
- 420 **29.** Samur, M. K. Rtcgatoobox: a new tool for exporting tcga firehose data. *PloS one* **9**, e106397 (2014).
- 421 **30.** Fenstermacher, D. A., Wenham, R. M., Rollison, D. E. & Dalton, W. S. Implementing personalized medicine in a cancer
422 center. *Cancer journal (Sudbury, Mass.)* **17**, 528 (2011).
- 423 **31.** Dalton, W. S. The “total cancer care” concept: linking technology and health care. *Cancer Control.* **12**, 140–141 (2005).
- 424 **32.** Welsh, E. A., Eschrich, S. A., Berglund, A. E. & Fenstermacher, D. A. Iterative rank-order normalization of gene expression
425 microarray data. *BMC bioinformatics* **14**, 1–11 (2013).
- 426 **33.** Baccarella, A., Williams, C. R., Parrish, J. Z. & Kim, C. C. Empirical assessment of the impact of sample number and read
427 depth on rna-seq analysis workflow performance. *BMC bioinformatics* **19**, 423 (2018).

428 **Supplementary Tables**

Table 1. Tissue of origin for all 429 epithelial-origin GDSC cell lines.

Tissue of origin	No. cell lines.
NSCLC adenocarcinoma	53
Breast	50
Large intestine	42
Ovary	33
Esophagus	30
Pancreas	30
Head and neck	28
Stomach	25
Bladder	18
Kidney	17
Mesothelium	16
Liver	15
NSCLC squamous	15
Cervix	14
Thyroid	14
Endometrium	11
Prostate	7
Biliary tract	5
Testis	3
Uterus	2
Adrenal gland	1

Table 2. Tissue of origin for DE comparison groups for fold 1.

Drug Response	Tissue of origin	No. cell lines
Resistant	NSCLC adenocarcinoma	13
	Pancreas	9
	Breast	8
	Large intestine	7
	Mesothelium	5
	Kidney	4
	Ovary	4
	Stomach	3
	Thyroid	3
	Biliary Tract	2
	Endometrium	2
	Liver	2
	NSCLC squamous	2
	Esophagus	2
	Cervix	1
	Head and neck	1
Prostate	1	
Sensitive	Head and neck	10
	NSCLC adenocarcinoma	10
	Ovary	6
	Breast	5
	Large intestine	5
	Esophagus	5
	Bladder	4
	Kidney	4
	Pancreas	4
	Stomach	4
	Cervix	3
	Thyroid	3
	NSCLC squamous	2
	Prostate	1
Testis	1	
Uterus	1	

Table 3. Tissue of origin for DE comparison groups for fold 2.

Drug Response	Tissue of origin	No. cell lines
Resistant	NSCLC adenocarcinoma	13
	Breast	9
	Pancreas	9
	Large intestine	5
	Mesothelium	5
	Stomach	5
	Esophagus	4
	Ovary	4
	Biliary Tract	2
	Head and neck	2
	Kidney	2
	Liver	2
	NSCLC squamous	2
	Thyroid	2
	Bladder	1
	Endometrium	1
Prostate	1	
Sensitive	Head and neck	11
	NSCLC adenocarcinoma	8
	Bladder	6
	Breast	5
	Large intestine	5
	Esophagus	5
	Ovary	5
	Stomach	5
	Cervix	3
	Kidney	3
	Pancreas	3
	NSCLC squamous	2
	Testis	2
	Thyroid	2
	Endometrium	1
	Prostate	1
Uterus	1	

Table 4. Tissue of origin for DE comparison groups for fold 3.

Drug Response	Tissue of origin	No. cell lines
Resistant	NSCLC adenocarcinoma	13
	Pancreas	10
	Breast	8
	Large intestine	5
	Mesothelium	5
	Kidney	4
	Ovary	4
	Stomach	4
	Esophagus	3
	Thyroid	3
	Endometrium	2
	Head and neck	2
	Biliary Tract	1
	Bladder	1
	Cervix	1
	Liver	1
	NSCLC squamous	1
Prostate	1	
Sensitive	Ovary	10
	Head and neck	9
	Bladder	8
	NSCLC adenocarcinoma	6
	Esophagus	5
	Breast	4
	Cervix	4
	Kidney	4
	Pancreas	4
	Stomach	4
	Large intestine	3
	Thyroid	3
	Testis	2
	Endometrium	1
	NSCLC squamous	1
Uterus	1	

Table 5. Tissue of origin for DE comparison groups for fold 4.

Drug Response	Tissue of origin	No. cell lines
Resistant	NSCLC adenocarcinoma	12
	Breast	10
	Large intestine	8
	Pancreas	6
	Mesothelium	5
	Ovary	5
	Thyroid	4
	Kidney	3
	Liver	3
	Esophagus	3
	Biliary Tract	2
	Head and neck	2
	Stomach	2
	Bladder	1
	Cervix	1
	Endometrium	1
NSCLC squamous	1	
Sensitive	Head and neck	9
	NSCLC adenocarcinoma	8
	Ovary	8
	Bladder	6
	Kidney	5
	Large intestine	5
	Stomach	5
	Cervix	4
	Esophagus	4
	Pancreas	3
	Thyroid	3
	Breast	2
	Testis	2
	Endometrium	1
	NSCLC squamous	1
	Prostate	1
Uterus	1	

Table 6. Tissue of origin for DE comparison groups for fold 5.

Drug Response	Tissue of origin	No. cell lines
Resistant	NSCLC adenocarcinoma	13
	Pancreas	10
	Breast	8
	Large intestine	8
	Mesothelium	4
	Esophagus	4
	Stomach	4
	Kidney	3
	Ovary	3
	Thyroid	3
	Endometrium	2
	Biliary Tract	1
	Bladder	1
	Cervix	1
	Head and neck	1
	Liver	1
	NSCLC squamous	1
Prostate	1	
Sensitive	Head and neck	12
	NSCLC adenocarcinoma	8
	Ovary	8
	Stomach	6
	Bladder	5
	Breast	5
	Large intestine	5
	Esophagus	5
	Kidney	4
	Cervix	3
	Pancreas	2
	Endometrium	1
	NSCLC squamous	1
	Prostate	1
	Testis	1
Thyroid	1	

Table 7. DE genes by fold. The SAM method consistently extracts more genes than limma or multtest. The intersection, however, is much smaller than either limma or multtest, showing significant filtering during the intersection step.

Fold	Method	No. Up-regulated Genes	No. Down-regulated Genes
1	SAM	1979	1083
	limma	181	322
	multtest	219	150
	intersection	59	58
2	SAM	1397	853
	limma	159	302
	multtest	139	115
	intersection	32	41
3	SAM	2290	1143
	limma	176	355
	multtest	247	173
	intersection	58	73
4	SAM	1904	1069
	limma	188	263
	multtest	237	147
	intersection	61	42
5	SAM	566	636
	limma	156	221
	multtest	93	87
	intersection	34	28

Table 8. NCCN Guideline versions used for assessing disease-site specific treatment guidelines.

Disease Site	NCCN Guideline Version	Cisplatin Use	Notes for select circumstances
ACC	Neuroendocrine and Adrenal Tumors Version 3.2021	Yes	
BLCA	Bladder Cancer Version 3.2021.	Yes	
BRCA	Breast Cancer Version 5.2021	Select circumstances	Only for recurrent, unresectable triple negative BRCA with germline BRCA1/2 mutation
CECSC	Cervical Cancer Version 1.2021	Yes	
CHOL	Hepatobiliary Cancers Version 5.2021	Yes	
COAD	Colon Cancer Version 2.2021	No	
ESCA	Esophageal and Esophagogastric Junction Cancers Version 3.2021	Yes	
HNSC	Head and Neck Cancers Version 3.2021	Yes	
KICH	Kidney Cancer Version 2.2022	No	
KIRP	Kidney Cancer Version 2.2022	No	
KIRC	Kidney Cancer Version 2.2022	No	
Kidney	Kidney Cancer Version 2.2022	No	
Renal	Kidney Cancer Version 2.2022	No	
Pelvis			
LIHC	Hepatobiliary Cancers Version 3.2021	No	
LUAD	Non-Small Cell Lung Cancer Version 5.2021	Yes	
LUSC	Non-Small Cell Lung Cancer Version 5.2021	Yes	
MESO	Malignant Pleural Mesothelioma Version 2.2021	Yes	
OV	Ovarian Cancer/Fallopian Tube Cancer/ Primary Peritoneal Cancer Version 1.2021	Yes	
PAAD	Pancreatic Adenocarcinoma Version 2.2021	Select circumstances	Only for BRCA1/2 or PALB2 mutations
PRAD	Prostate Cancer Version 2.2021	No	
READ	Rectal Cancer Version 1.2021	No	
STAD	Gastric cancer Version 3.2021	Yes	
THCA	Thyroid Carcinoma Version 1.2021	Select circumstances	Only as adjuvant/radiosensitizer for anaplastic carcinoma
THYM	Thymomas and Thymic Carcinomas Version 1.2021	Yes	
UCEC	Uterine Neoplasms Version 3.2021	Yes	

429 **Supplementary Figures**

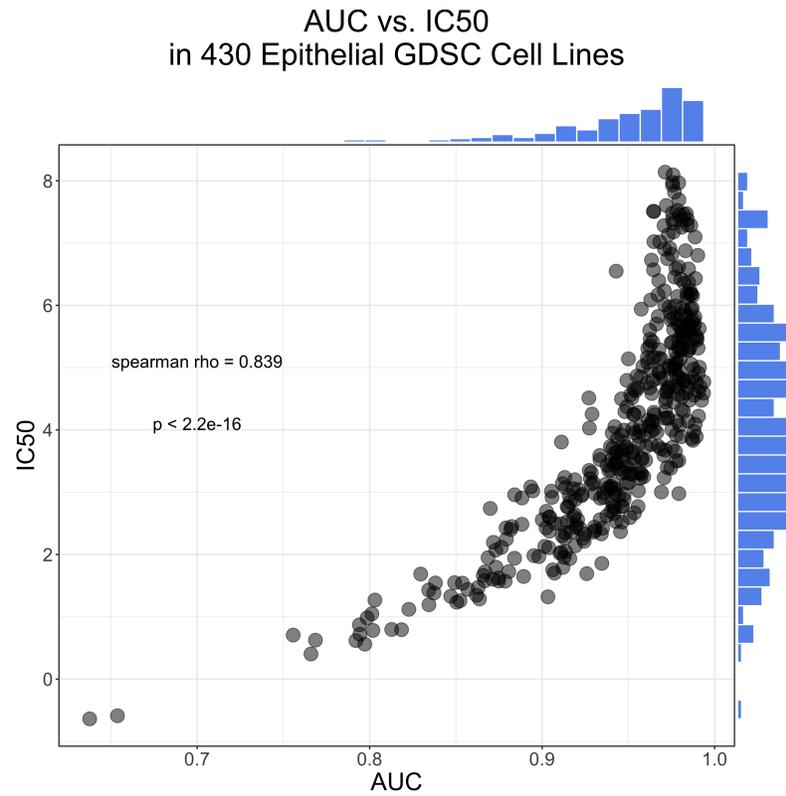


Figure 1. Correlation between AUC and IC50 drug response metrics for epithelial-based cancer cell lines in the Genomics of Drug Discovery in Cancer (GDSC) dataset.

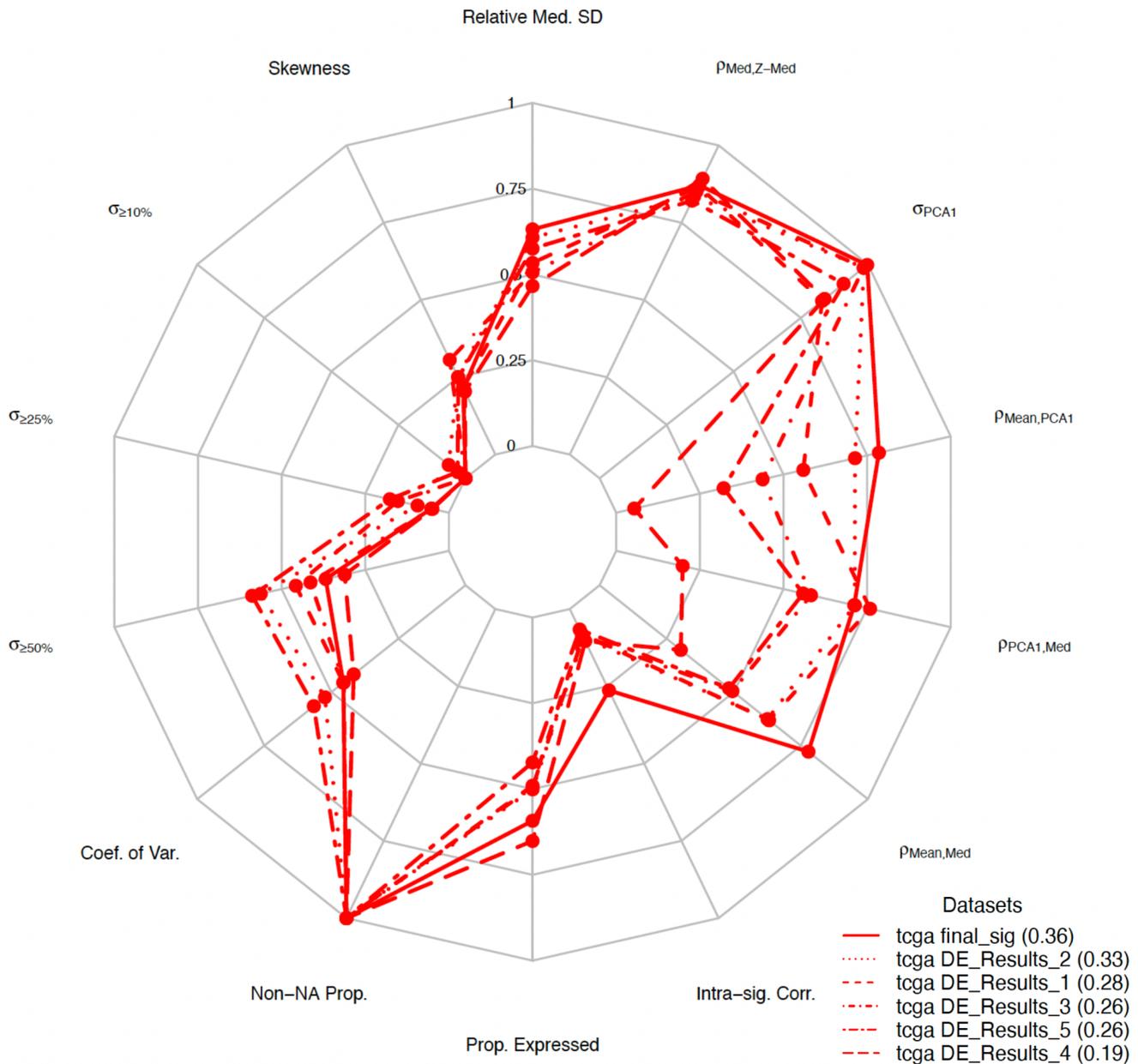


Figure 2. Quality control metrics comparing differential expression results to the final gene signature using sigQC^{20,21}. CisSig is compared to the folds of differential gene expression analysis, comparing results using a radar plot. It shows greater intra-signature correlation, higher correlation between signature mean and median, and decreased skewness within RNA-seq expression from TCGA samples of epithelial origin.

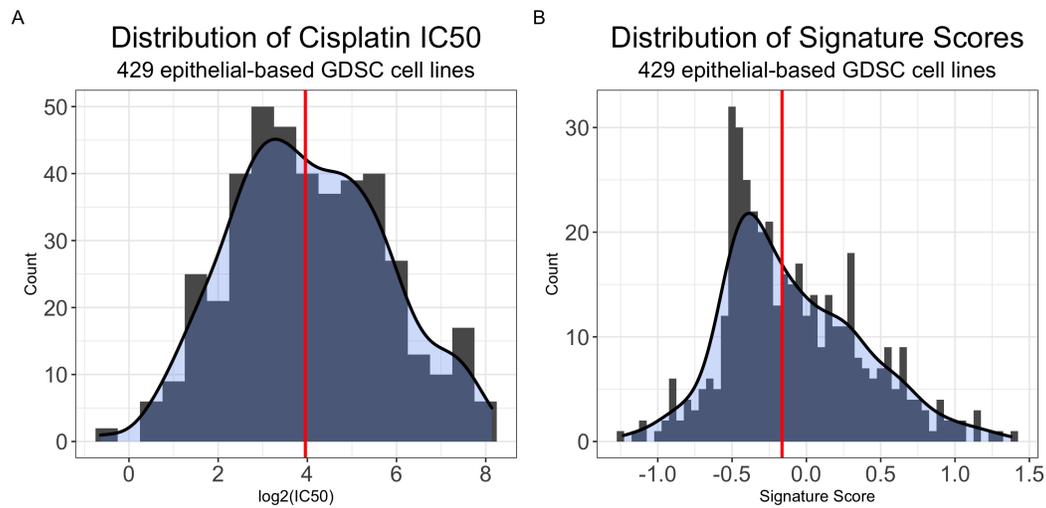


Figure 3. Cisplatin IC₅₀ (log₂-transformed) in epithelial-origin GDSC cell lines is relatively normally distributed, while CisSig Score has a right skew. A. Distribution of CisSig across 429 epithelial-based GDSC cell lines, using a histogram (gray) and kernel density estimation (blue). Median score marked by red vertical line. CisSig score is calculated as a cell line’s median normalized expression of CisSig genes listed in A. **B.** Distribution of cisplatin IC₅₀ across 429 epithelial-based GDSC cell lines, using a histogram (gray) and kernel density estimation (blue). Median IC₅₀ marked by red vertical line.

It is made available under a [CC-BY-NC 4.0 International license](https://creativecommons.org/licenses/by-nc/4.0/).

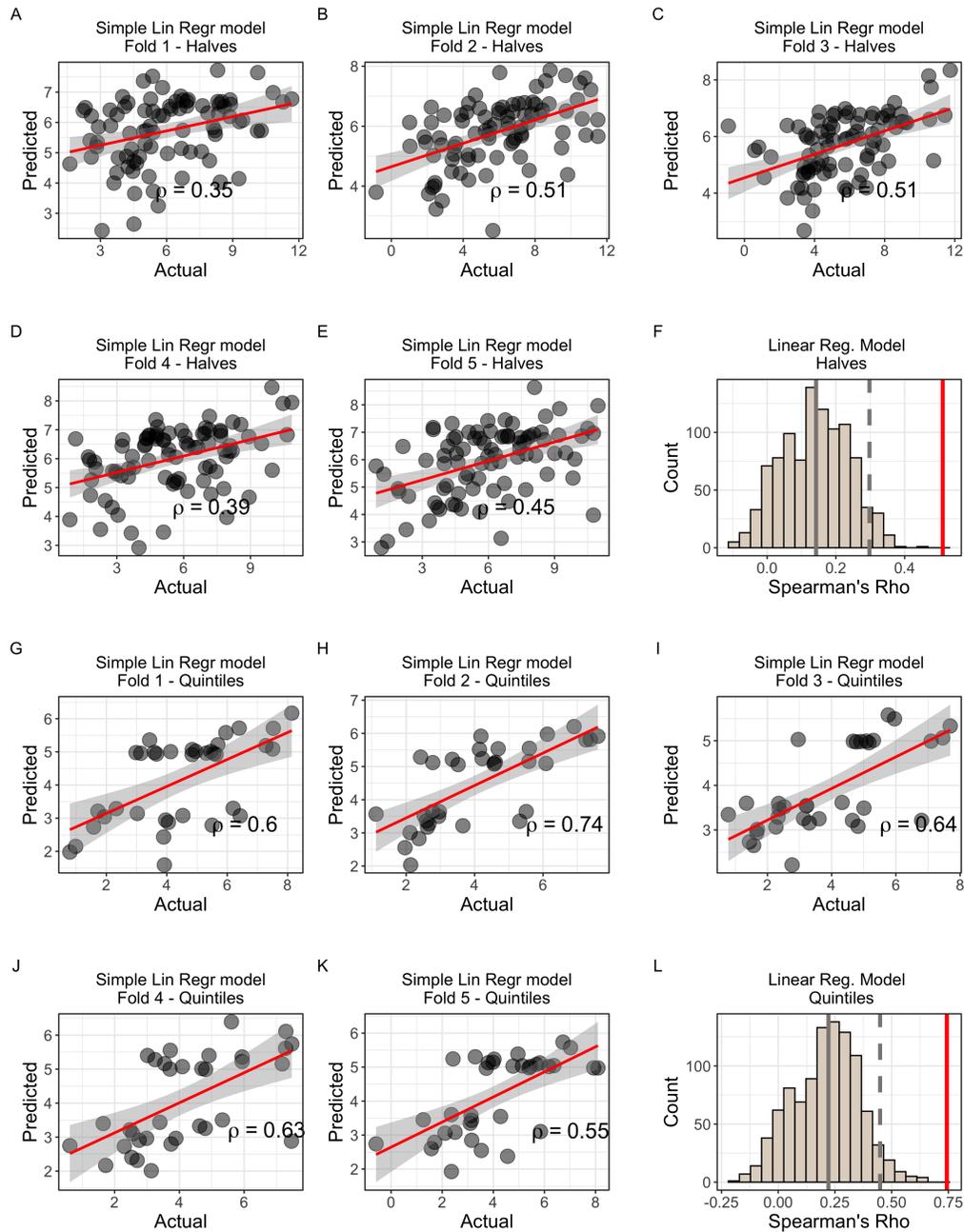


Figure 4. Modeling IC50 response using CisSig Score to predict IC50 in GDSC with simple linear regression. A-E. Predicted vs. Actual IC50 for validation sets of folds 1-5 for models built with all 429 cell lines. **F.** Null distribution of modeling metrics using 1000 random gene signatures with the same length as CisSig and the model described in A-E. CisSig's performance (red solid line) is within the top 5% of the null distribution (cutoff at gray dashed line). Gray solid line represents median of null distribution. **G-K.** Predicted vs. Actual IC50 for validation sets of folds 1-5 for models built using cell lines in the top and bottom 20% of cisplatin IC50. **F.** Null distribution of modeling metrics using 1000 random gene signatures with the same length as CisSig and the model described in G-K. CisSig's performance (red solid line) is compared to the 95% confidence interval (gray dashed line) of the null distribution. Gray solid line represents median of null distribution.

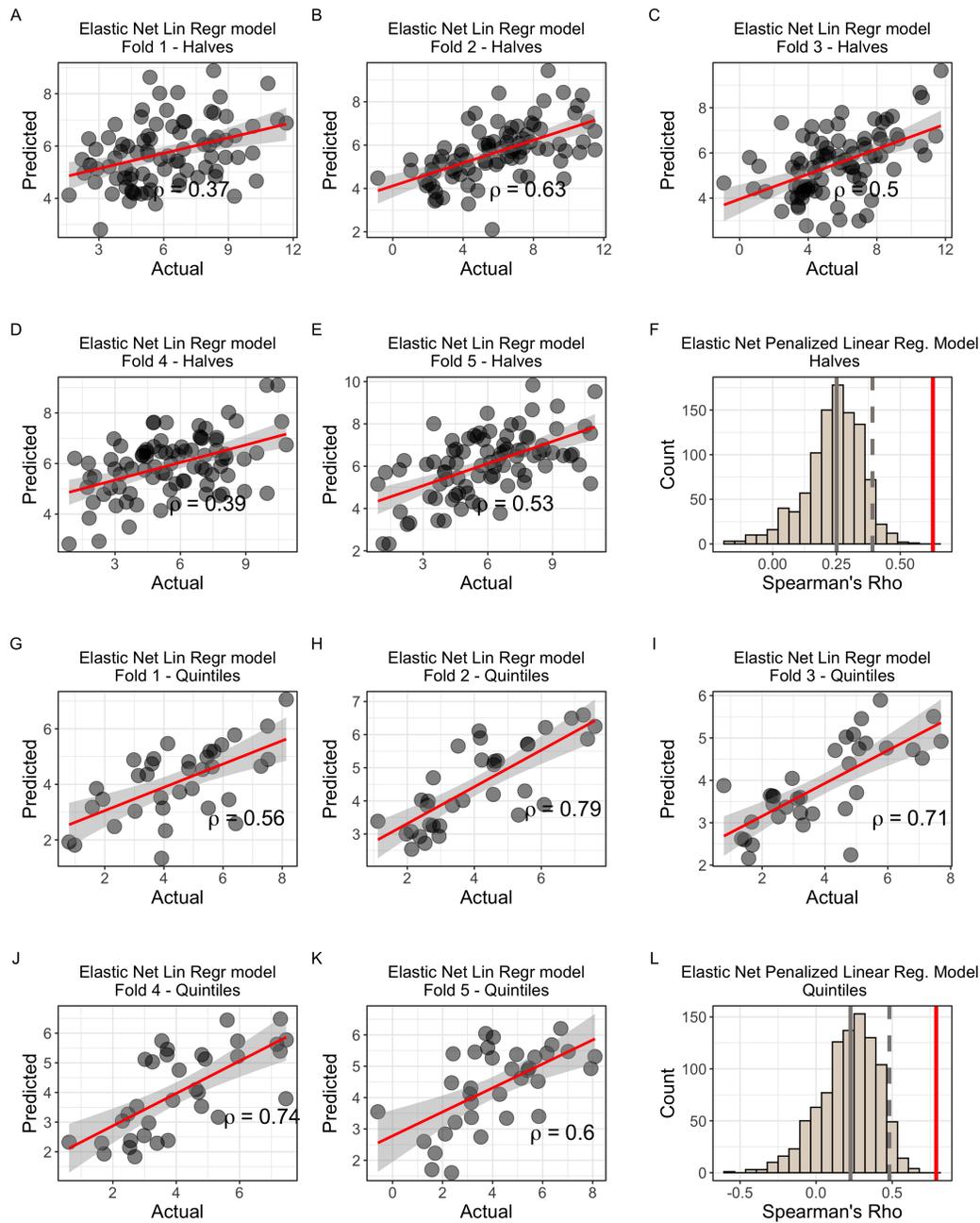


Figure 5. Modeling IC50 response using individual CisSig genes to predict IC50 in GDSC with elastic net penalized linear regression. **A-E.** Predicted vs. Actual IC50 for validation sets of folds 1-5 for models built with all 429 cell lines. **F.** Null distribution of modeling metrics using 1000 random gene signatures with the same length as CisSig and the model described in **A-E**. CisSig's performance (red solid line) is within the top 5% of the null distribution (cutoff at gray dashed line). Gray solid line represents median of null distribution. **G-K.** Predicted vs. Actual IC50 for validation sets of folds 1-5 for models built using cell lines in the top and bottom 20% of cisplatin IC50. **F.** Null distribution of modeling metrics using 1000 random gene signatures with the same length as CisSig and the model described in **G-K**. CisSig's performance (red solid line) is compared to the 95% confidence interval (gray dashed line) of the null distribution. Gray solid line represents median of null distribution.

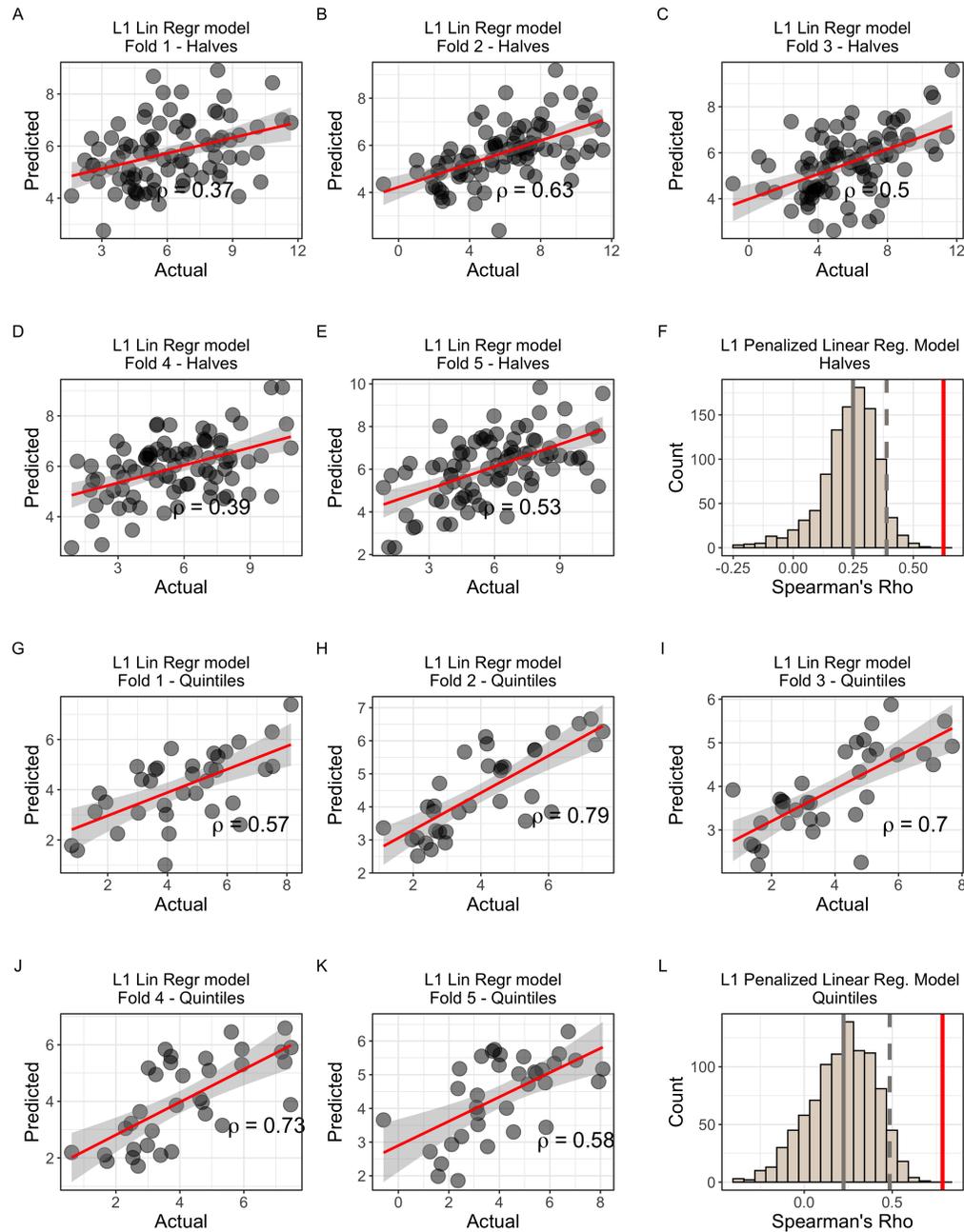


Figure 6. Modeling IC50 response using individual CisSig genes to predict IC50 in GDSC with L1 penalized linear regression. **A-E.** Predicted vs. Actual IC50 for validation sets of folds 1-5 for models built with all 429 cell lines. **F.** Null distribution of modeling metrics using 1000 random gene signatures with the same length as CisSig and the model described in **A-E**. CisSig's performance (red solid line) is within the top 5% of the null distribution (cutoff at gray dashed line). Gray solid line represents median of null distribution. **G-K.** Predicted vs. Actual IC50 for validation sets of folds 1-5 for models built using cell lines in the top and bottom 20% of cisplatin IC50. **F.** Null distribution of modeling metrics using 1000 random gene signatures with the same length as CisSig and the model described in **G-K**. CisSig's performance (red solid line) is compared to the 95% confidence interval (gray dashed line) of the null distribution. Gray solid line represents median of null distribution.

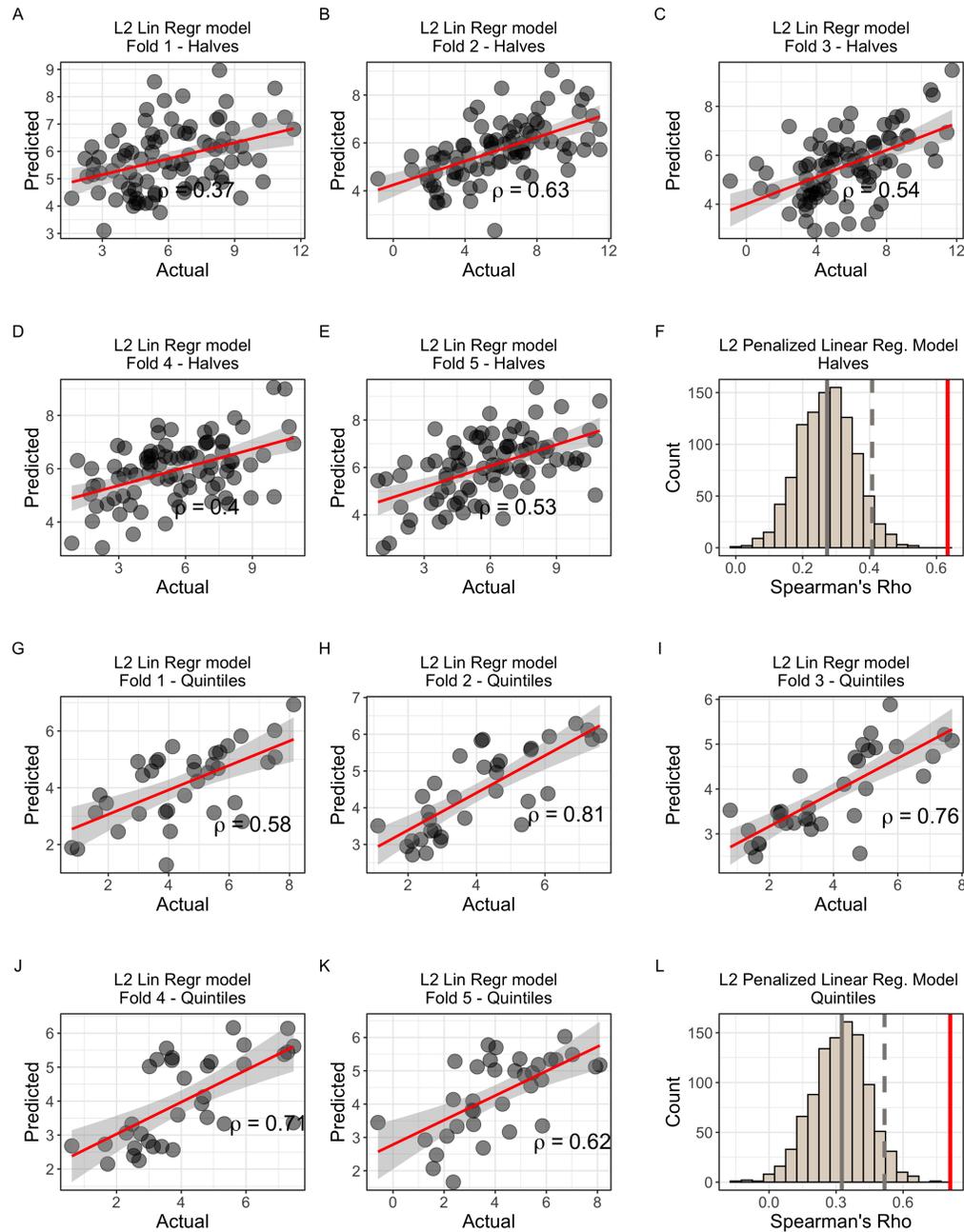


Figure 7. Modeling IC50 response using individual CisSig genes to predict IC50 in GDSC with L2 penalized linear regression. **A-E.** Predicted vs. Actual IC50 for validation sets of folds 1-5 for models built with all 429 cell lines. **F.** Null distribution of modeling metrics using 1000 random gene signatures with the same length as CisSig and the model described in **A-E**. CisSig's performance (red solid line) is within the top 5% of the null distribution (cutoff at gray dashed line). Gray solid line represents median of null distribution. **G-K.** Predicted vs. Actual IC50 for validation sets of folds 1-5 for models built using cell lines in the top and bottom 20% of cisplatin IC50. **F.** Null distribution of modeling metrics using 1000 random gene signatures with the same length as CisSig and the model described in **G-K**. CisSig's performance (red solid line) is compared to the 95% confidence interval (gray dashed line) of the null distribution. Gray solid line represents median of null distribution.

It is made available under a [CC-BY-NC 4.0 International license](https://creativecommons.org/licenses/by-nc/4.0/).

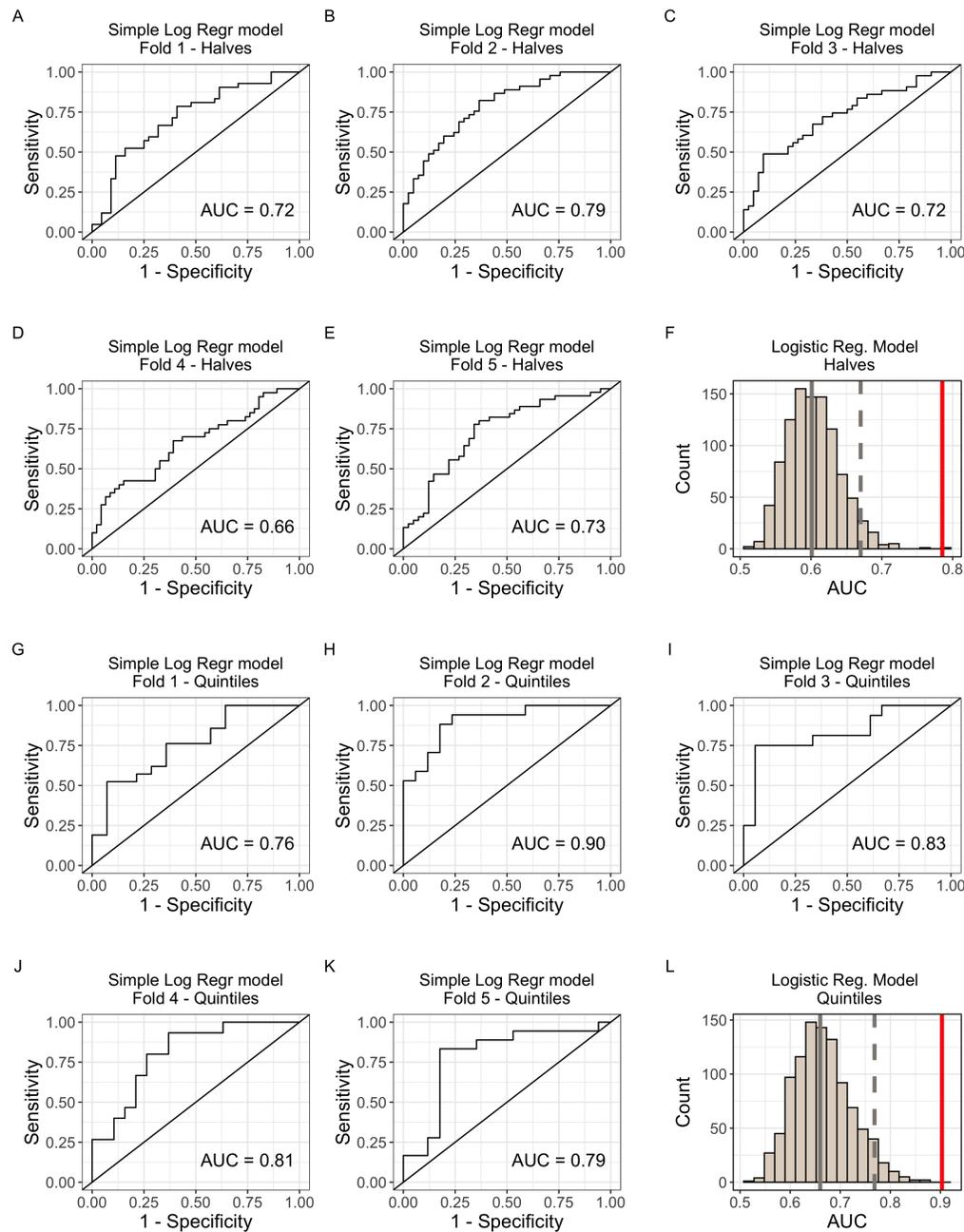


Figure 8. Modeling IC₅₀ response using CisSig score to predict IC₅₀ class in GDSC with simple logistic regression.

A-E. Predicted vs. Actual IC₅₀ for validation sets of folds 1-5 for models built with all 429 cell lines. **F.** Null distribution of modeling metrics using 1000 random gene signatures with the same length as CisSig and the model described in **A-E**. CisSig's performance (red solid line) is within the top 5% of the null distribution (cutoff at gray dashed line). Gray solid line represents median of null distribution. **G-K.** Predicted vs. Actual IC₅₀ for validation sets of folds 1-5 for models built using cell lines in the top and bottom 20% of cisplatin IC₅₀. **L.** Null distribution of modeling metrics using 1000 random gene signatures with the same length as CisSig and the model described in **G-K**. CisSig's performance (red solid line) is compared to the 95% confidence interval (gray dashed line) of the null distribution. Gray solid line represents median of null distribution.

It is made available under a [CC-BY-NC 4.0 International license](https://creativecommons.org/licenses/by-nc/4.0/).

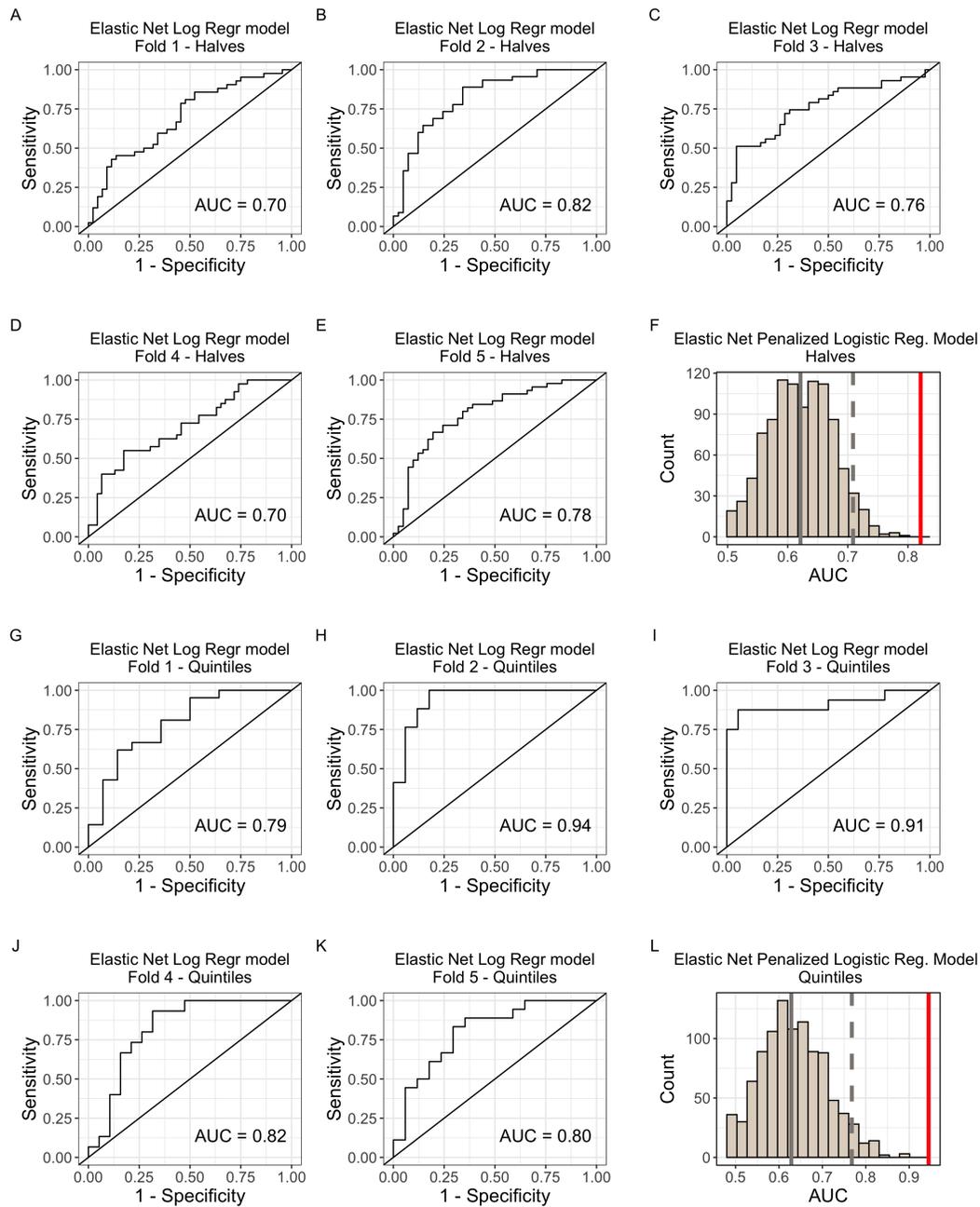


Figure 9. Modeling IC50 response using individual CisSig genes to predict IC50 class in GDSC with elastic net penalized logistic regression. **A-E.** Predicted vs. Actual IC50 for validation sets of folds 1-5 for models built with all 429 cell lines. **F.** Null distribution of modeling metrics using 1000 random gene signatures with the same length as CisSig and the model described in **A-E**. CisSig's performance (red solid line) is within the top 5% of the null distribution (cutoff at gray dashed line). Gray solid line represents median of null distribution. **G-K.** Predicted vs. Actual IC50 for validation sets of folds 1-5 for models built using cell lines in the top and bottom 20% of cisplatin IC50. **F.** Null distribution of modeling metrics using 1000 random gene signatures with the same length as CisSig and the model described in **G-K**. CisSig's performance (red solid line) is compared to the 95% confidence interval (gray dashed line) of the null distribution. Gray solid line represents median of null distribution.

It is made available under a [CC-BY-NC 4.0 International license](https://creativecommons.org/licenses/by-nc/4.0/).

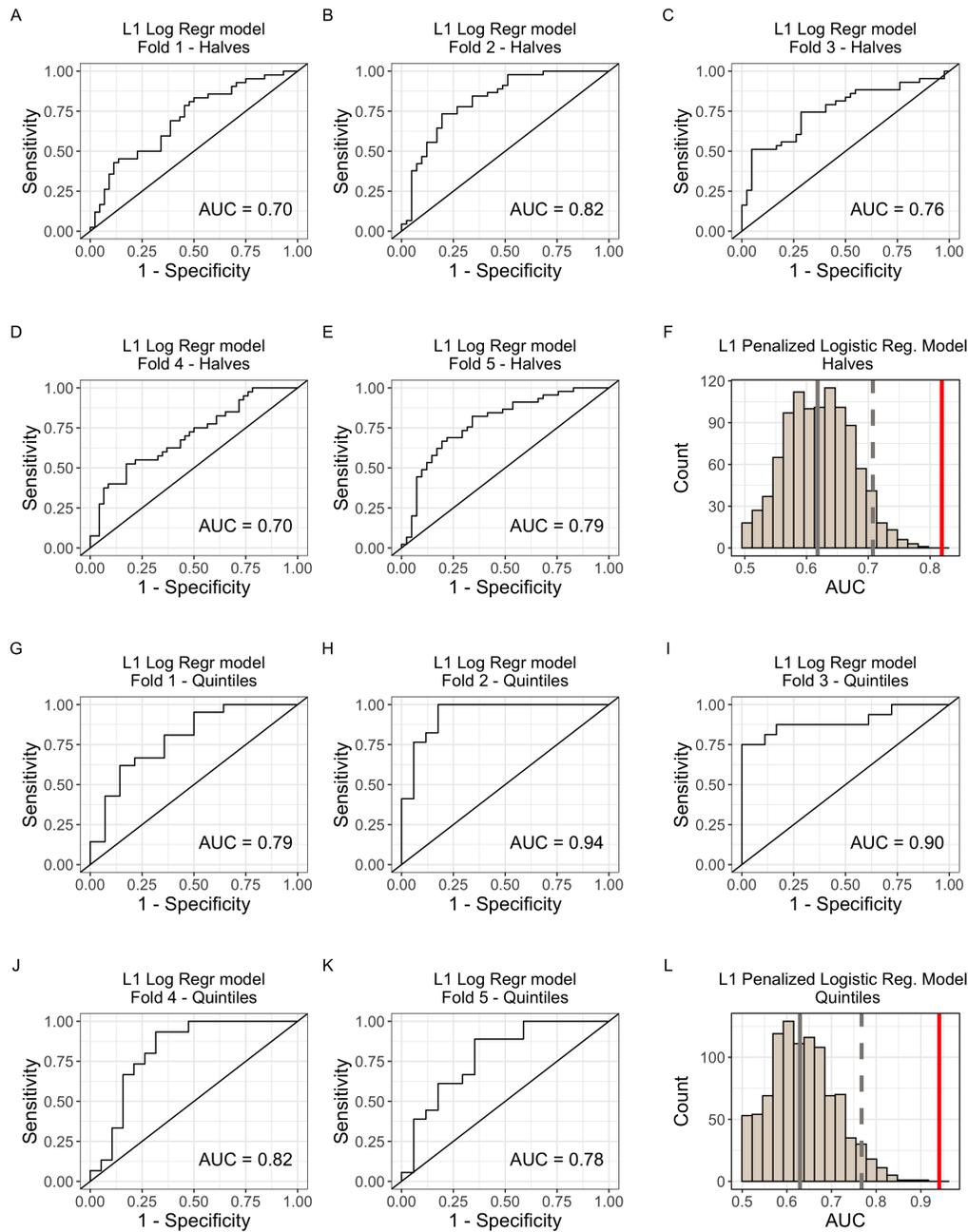
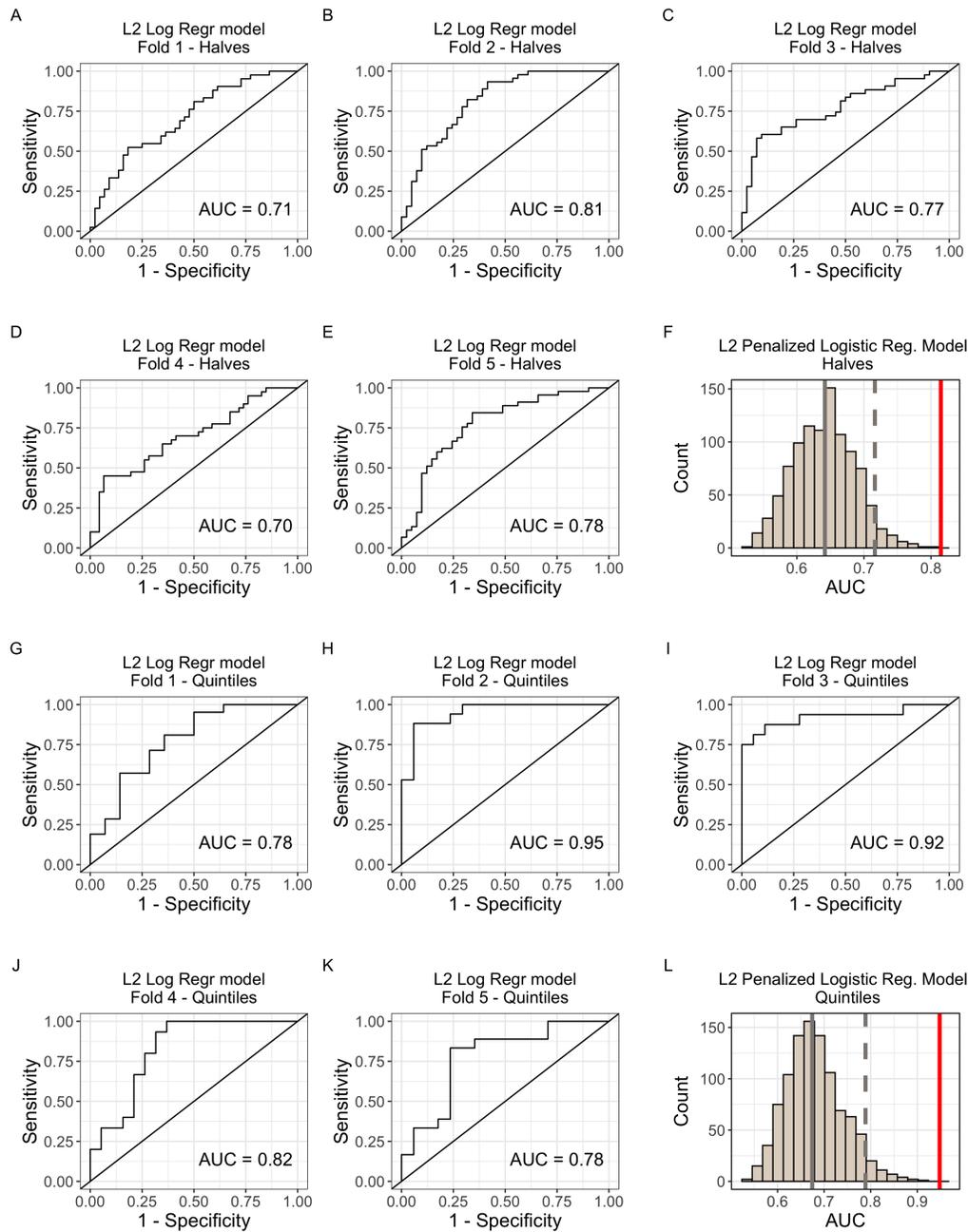


Figure 10. Modeling IC50 response using individual CisSig genes to predict IC50 class in GDSC with L1 penalized logistic regression. **A-E.** Predicted vs. Actual IC50 for validation sets of folds 1-5 for models built with all 429 cell lines. **F.** Null distribution of modeling metrics using 1000 random gene signatures with the same length as CisSig and the model described in **A-E**. CisSig's performance (red solid line) is within the top 5% of the null distribution (cutoff at gray dashed line). Gray solid line represents median of null distribution. **G-K.** Predicted vs. Actual IC50 for validation sets of folds 1-5 for models built using cell lines in the top and bottom 20% of cisplatin IC50. **F.** Null distribution of modeling metrics using 1000 random gene signatures with the same length as CisSig and the model described in **G-K**. CisSig's performance (red solid line) is compared to the 95% confidence interval (gray dashed line) of the null distribution. Gray solid line represents median of null distribution.

It is made available under a [CC-BY-NC 4.0 International license](https://creativecommons.org/licenses/by-nc/4.0/).



It is made available under a [CC-BY-NC 4.0 International license](https://creativecommons.org/licenses/by-nc/4.0/).

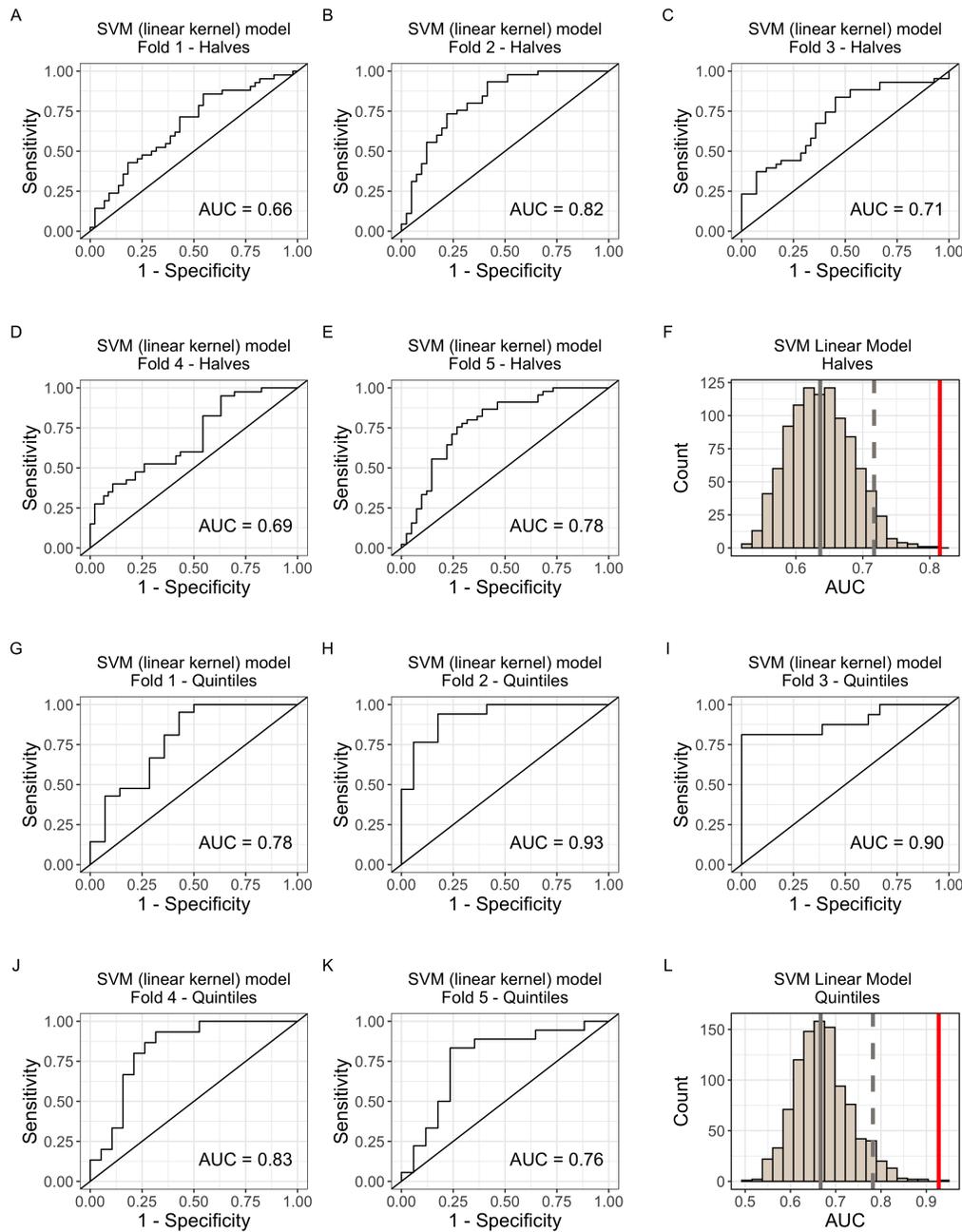


Figure 12. Modeling IC50 response using individual CisSig genes to predict IC50 class in GDSC with support vector machine modeling (linear kernel). **A-E.** Predicted vs. Actual IC50 for validation sets of folds 1-5 for models built with all 429 cell lines. **F.** Null distribution of modeling metrics using 1000 random gene signatures with the same length as CisSig and the model described in **A-E**. CisSig's performance (red solid line) is within the top 5% of the null distribution (cutoff at gray dashed line). Gray solid line represents median of null distribution. **G-K.** Predicted vs. Actual IC50 for validation sets of folds 1-5 for models built using cell lines in the top and bottom 20% of cisplatin IC50. **F.** Null distribution of modeling metrics using 1000 random gene signatures with the same length as CisSig and the model described in **G-K**. CisSig's performance (red solid line) is compared to the 95% confidence interval (gray dashed line) of the null distribution. Gray solid line represents median of null distribution.

It is made available under a [CC-BY-NC 4.0 International license](https://creativecommons.org/licenses/by-nc/4.0/).

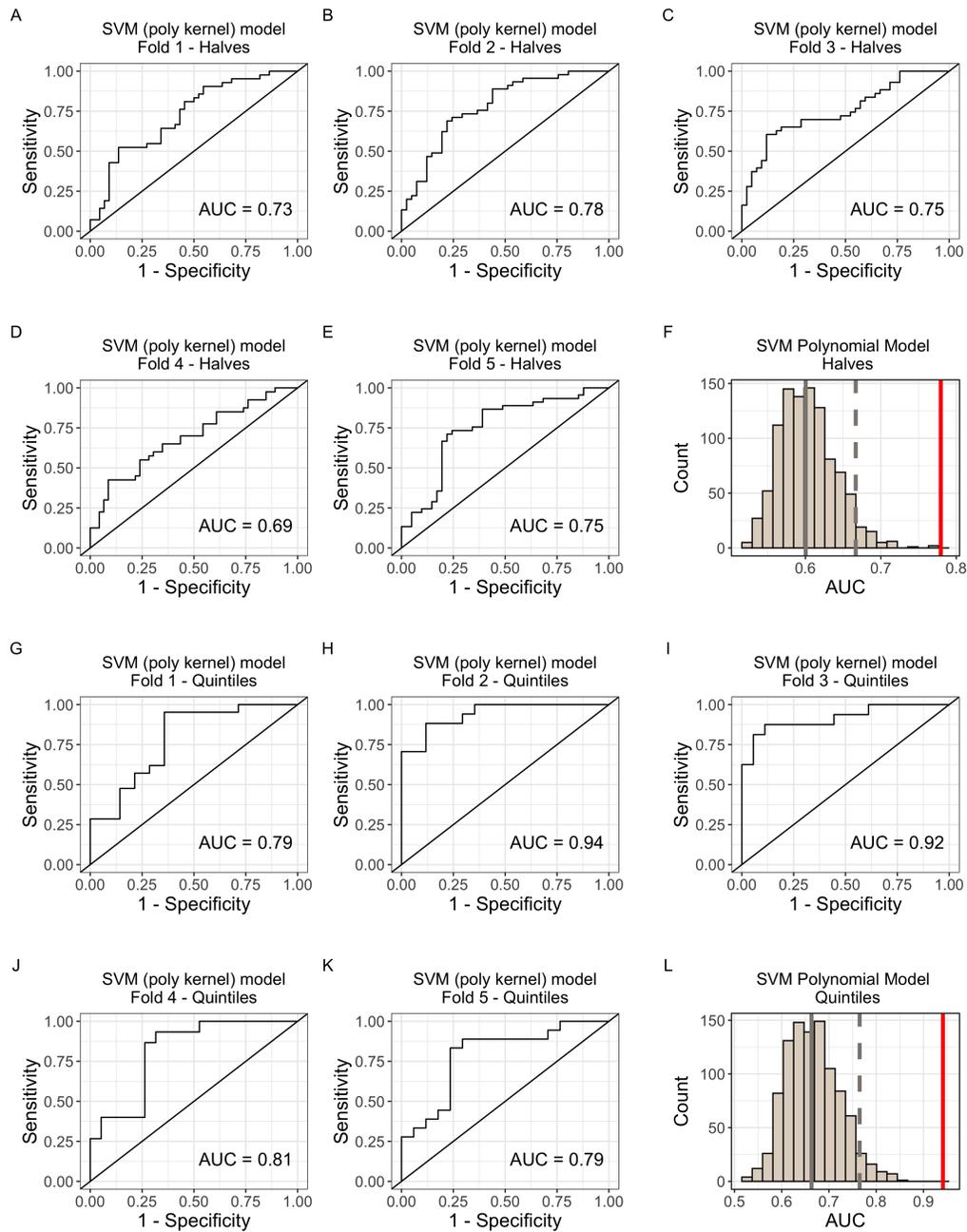


Figure 13. Modeling IC50 response using individual CisSig genes to predict IC50 class in GDSC with support vector machine modeling (polynomial kernel). A-E. Predicted vs. Actual IC50 for validation sets of folds 1-5 for models built with all 429 cell lines. F. Null distribution of modeling metrics using 1000 random gene signatures with the same length as CisSig and the model described in A-E. CisSig's performance (red solid line) is within the top 5% of the null distribution (cutoff at gray dashed line). Gray solid line represents median of null distribution. G-K. Predicted vs. Actual IC50 for validation sets of folds 1-5 for models built using cell lines in the top and bottom 20% of cisplatin IC50. F. Null distribution of modeling metrics using 1000 random gene signatures with the same length as CisSig and the model described in G-K. CisSig's performance (red solid line) is compared to the 95% confidence interval (gray dashed line) of the null distribution. Gray solid line represents median of null distribution.

It is made available under a [CC-BY-NC 4.0 International license](https://creativecommons.org/licenses/by-nc/4.0/).

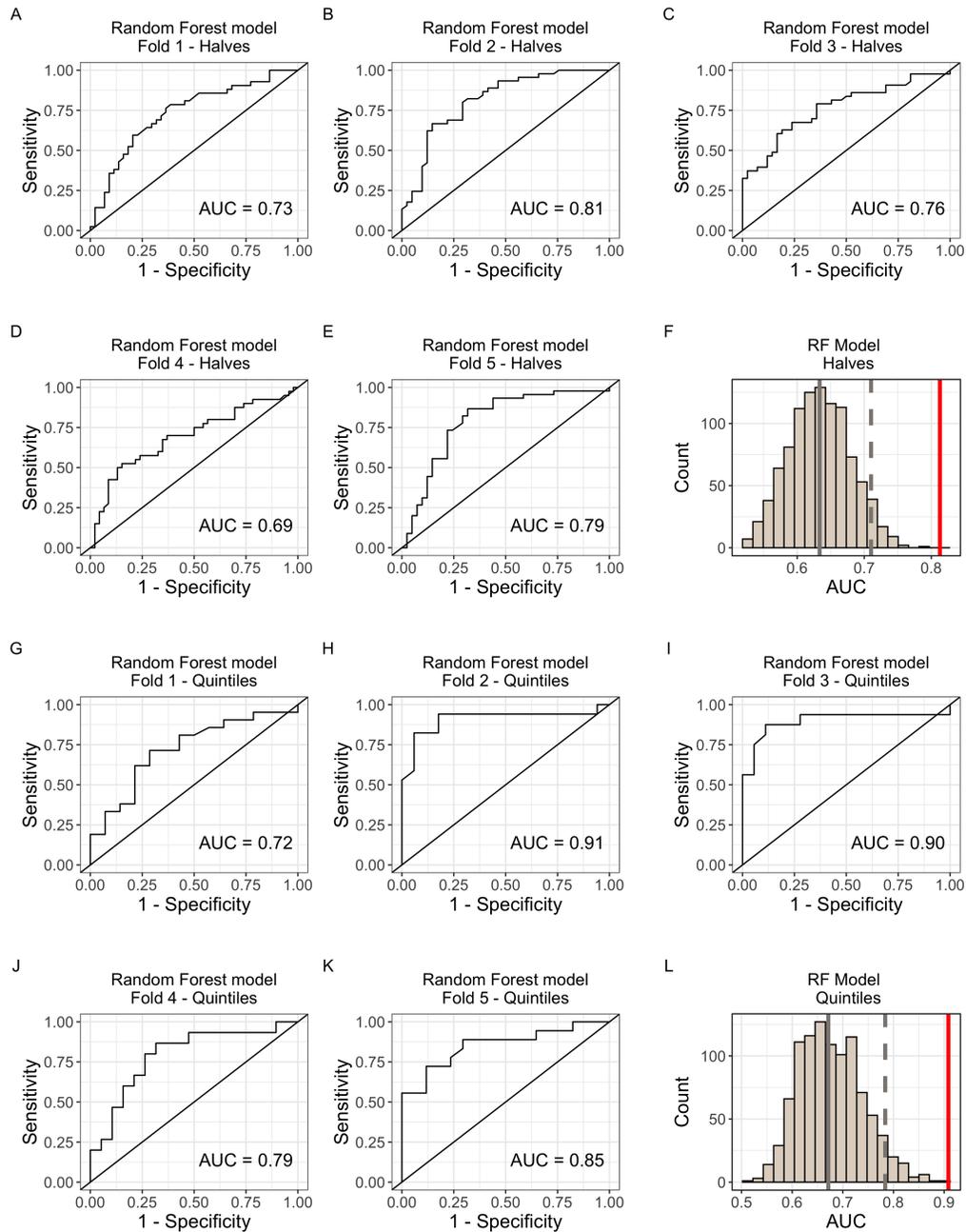


Figure 14. Modeling IC50 response using individual CisSig genes to predict IC50 class in GDSC with random forest modeling. A-E. Predicted vs. Actual IC50 for validation sets of folds 1-5 for models built with all 429 cell lines. F. Null distribution of modeling metrics using 1000 random gene signatures with the same length as CisSig and the model described in A-E. CisSig's performance (red solid line) is within the top 5% of the null distribution (cutoff at gray dashed line). Gray solid line represents median of null distribution. G-K. Predicted vs. Actual IC50 for validation sets of folds 1-5 for models built using cell lines in the top and bottom 20% of cisplatin IC50. F. Null distribution of modeling metrics using 1000 random gene signatures with the same length as CisSig and the model described in G-K. CisSig's performance (red solid line) is compared to the 95% confidence interval (gray dashed line) of the null distribution. Gray solid line represents median of null distribution.