

Examining Predictions for Doxorubicin in Detail

Keith A. Baggerly

November 19, 2009

Contents

1	Executive Summary	2
1.1	Introduction	2
1.2	Methods	2
1.3	Results	3
1.4	Conclusions	3
2	Options and Libraries	4
3	Loading Data	4
3.1	Earlier Rda Files: NCI60 Cell Line Data, Lists of Cell Lines Used	4
3.2	List of Probesets	4
3.3	Doxorubicin 2007 Data: Adria_ALL.txt	5
3.4	Doxorubicin 2008 Data: Adria_ALL_data1_n95.doc	7
3.5	Holleman et al. [2] Data	7
3.6	Lugthart et al. [3] Data	8
4	Identifying the Cell Lines Used in Adria_ALL.txt	9
4.1	Getting the NCI60 Quantifications	9
4.2	Transforming the Adria_ALL.txt Data	9
4.3	Matching the First Probeset	10
5	Clustering the Adria_ALL.txt Test Data	12
6	Examining the Adria_ALL.txt Test Data	14
7	Examining Doxorubicin After the August 2008 Correction	17
7.1	Looking for Ties	17
7.2	Checking Consistency of St. Jude and Potti et al. [6] Sensitive/Resistant Calls	20
8	Looking at the Data Cited	23
9	Assembling a Summary Figure and Table	23
10	Appendix	24
10.1	File Location	24
10.2	Saves	24
10.3	SessionInfo	24

1 Executive Summary

1.1 Introduction

In late 2006, Potti et al. [4] introduced a method for combining microarray profiles of cell lines with drug sensitivity data to derive “signatures” of sensitivity to specific drugs. These signatures could then be used to predict patient response. In theory, the approach is straightforward:

- Using drug sensitivity data for a panel of cell lines, select those that are most sensitive and most resistant to the drug of interest.
- Using array profiles of the identified cell lines, identify the most differentially expressed genes.
- Using the most differentially expressed genes, build a model that takes an array profile and returns a classification.

This report is part of a series in which we try to trace the specific steps involved in order to better understand the approach. In this report, we focus on response to doxorubicin, which Potti et al. have now addressed multiple times [4,5,6].

Potti et al. [4] initially predicted adriamycin response using 122 test samples. However, the initially cited source, Holleman et al. [2], names 94 responders and 28 nonresponders but Potti et al. [4] find 23 and 99, respectively. Coombes et al. [1] suggested labels might have been reversed. In reply, Potti and Nevins [5] disputed this conclusion, commenting on “the acute lymphocytic leukemia dataset in which the labels are accurate — full details are provided on our web page.” However, in August of 2008 Potti et al. [6] posted a second correction in which they noted that

of the 122 samples assayed for sensitivity to daunorubicin for which the authors applied a predictor of adriamycin sensitivity, 27 samples were replicated owing to the fact that the same samples were included in several separate series files in the Gene Expression Omnibus generated in 2004 and 2005, which were the source of the data provided for the study.

We examine both the data posted in November 2007 and the data posted in August 2008.

1.2 Methods

We acquired the raw doxorubicin (adriamycin) data posted on the Potti et al. [4] web site in November 2007 (Adria_ALL.txt), and in August 2008 (Adria_ALL_data1_n95.doc). We also acquired raw data pertaining to the Holleman et al. [2] test set from the supplementary web site <http://www.stjuderesearch.org/data/ALL4> and from the Gene Expression Omnibus (GEO) page for GSE2351, <http://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE2351>.

Adria_ALL.txt contains 144 data columns: 22 for training data cell lines, and 122 for test data samples. Sample columns are not named, but sensitive/resistant status is indicated for each. Starting from the NCI60 cell line data and the cell lines named for doxorubicin, we first mimicked the processing applied: the data were log-transformed, the values for each row (gene) were centered and scaled to have mean zero and unit variance (separately for training and test data), the data were exponentiated to undo the log-transform, and the final results were rounded to two decimal places. We then used correlation to identify the cell lines used for training. We also used correlation to look for patterns that might drive clustering in the test set.

Adria_ALL_data1_n95.doc gives no quantifications, but rather lists 95 Gene Expression Omnibus (GEO) array ids and gives a sensitive/resistant label for each. We first checked the list of names for ties and internal consistency. Then, using the GEO identifiers, we identified the classification that the Holleman et al. [2] rules would assign to each sample: sensitive, intermediate, or resistant as the LC50 for daunorubicin is below,

in between, or above the two values (0.075, 0.114) respectively. These LC50 values are given on the GEO web site for GSE2351 for all but three of the samples; for the others the classification is available from the supplementary St. Jude site noted above. We then compared these classifications with those assigned by Potti et al. [6].

Finally, we examined the listed sources of data to clarify the rationale behind their selection. The initial Potti et al. [4] paper lists GSE650 and GSE651. A side comment in *Adria_ALL.txt* notes that “Validation data is from GSE4698, GSE649, GSE650, GSE651, and others”. The second correction to Potti et al. [6] adds two more accession numbers, GSE2351 and GSE649, to those listed in Potti et al. [4].

1.3 Results

The cell lines used in the 2007 data (*Adria_ALL.txt*) match those supplied in the list of “cell lines used in each chemo predictor” in 2008, but with the sensitive and resistant lines reversed. The 2007 labeling is incorrect; in this dataset cell line NCI/ADR-RES (“adriamycin resistant”) is called sensitive.

We computed all pairwise sample correlations for the 2007 data; values above 0.9999 are shown in the first panel of Figure 1. Symbols off the main diagonal show the presence of ties: red triangles indicate pairs where sensitive/resistant labels for the same sample differ. Only 84 of the 122 test samples are distinct; others are present 2, 3, or 4 times, and the labels for repeated samples often conflict.

We also examined all pairs of sample names for the 2008 data; matches are shown in the second panel of Figure 1. Again, symbols off of the main diagonal show the presence of ties, and red triangles indicate where sensitive/resistant labels for the same sample differ. Only 80 of the 95 “unique” samples are distinct; 15 names are present twice. Of the 15 duplicates, 6 are labeled sensitive once and resistant once.

Labels Potti et al. [6] assigned to the 80 distinct samples in 2008 are cross-tabulated against labels the Holleman et al. [2] rules assign in Table 1. Holleman et al. [2] would classify 48 as sensitive, 10 as intermediate and 22 as resistant. Potti et al. [4] list 13 as sensitive, 61 as resistant, and 6 as both. Every sample that Holleman et al. [2] call resistant or intermediate is called resistant by Potti et al. [4], but only 13 (maybe 19, depending on how we treat the inconsistent calls) of the samples that Holleman et al. [2] call sensitive are called sensitive by Potti et al. [4]. Between 29 and 35 samples change classes. This is not explainable by a simple shift in LC50 cutoffs, as values for “sensitive” and “resistant” samples overlap.

Of the 4 accession numbers listed, all use samples profiled by Holleman et al. [2]. GSE650 and GSE651 list arrays for 94 and 28 samples that are sensitive and resistant to daunorubicin, respectively. There are no ties in these two datasets: all samples are distinct. GSE649 lists arrays for samples resistant to vincristine; there is some overlap between GSE649 and the two others listed above. GSE2351 examines the same samples for cross-resistance to multiple drugs and introduces a small number of new samples; there is overlap between GSE2351 and the three others listed above. No rationale is given for including samples from the latter two sets, or for excluding any of the non-tied samples from the prior two.

1.4 Conclusions

The labels applied to the doxorubicin samples were incorrect. The 2007 cell line labels were reversed, which should negate predictive accuracy. These reversed labels were posted after the Coombes et al. [1] suggestion of reversal was denied. Ties present in the test data in both 2007 and 2008 further invalidate predictive accuracy, particularly when labels for the same samples are inconsistent. Ties and inconsistencies persist in 2008 even though the 2008 correction was specifically to address the issue of ties in the data. There were also ties in the data posted for docetaxel in 2007 (there was no data posted in 2008), though there were no ties in the GEO accession numbers there. The labels Potti et al. [6] assign conflict with the labels assigned by the suppliers of the test data, with no explanation given. The additional data sources listed by Potti et al. [6] in their second correction make no sense. As shown in *checkingCellLines*, vincristine induces a

wholly different response pattern than doxorubicin, so there is no reason that a signature for one should be predictive for the other. Similarly, it is not clear why a signature designed for doxorubicin should be predictive for cross-resistance. The doxorubicin data does not support their main thesis.

2 Options and Libraries

```
> options(width = 80)
> library(geneplotter)
```

3 Loading Data

3.1 Earlier Rda Files: NCI60 Cell Line Data, Lists of Cell Lines Used

Here, we simply load RData objects assembled earlier for the U95Av2 NCI60 quantification data (novartisA), and the lists of cell lines used (cellLinesUsed, built in enumeratingCellLines).

```
> rdaList <- c("novartisA", "cellLinesUsed")
> for (rdaFile in rdaList) {
+   rdaFullFile <- file.path("RDataObjects", paste(rdaFile, "Rda",
+     sep = "."))
+   if (file.exists(rdaFullFile)) {
+     cat("loading ", rdaFullFile, " from cache\n")
+     load(rdaFullFile)
+   }
+   else {
+     cat("building ", rdaFullFile, " from raw data\n")
+     Stangle(file.path("RNowebSource", paste("buildRda", rdaFile,
+       "Rnw", sep = ".")))
+     source(paste("buildRda", rdaFile, "R", sep = "."))
+   }
+ }
```

```
loading RDataObjects/novartisA.Rda from cache
loading RDataObjects/cellLinesUsed.Rda from cache
```

3.2 List of Probesets

For the list of genes used, we simply load the list of probesets supplied on the Potti et al. [4] website.

```
> doxorubicinGenes <- read.table(file.path("RawData", "PottiNatMed",
+   "GeneListsNov07", "Adria(final).txt"), sep = "\t", header = TRUE)
> dim(doxorubicinGenes)
```

```
[1] 80 3
```

```
> doxorubicinGenes[1:3, ]
```

```

      Probe.Set.ID                      Gene.Title
1      1051_g_at                        melan-A
2      110_at  chondroitin sulfate proteoglycan 4 (melanoma-associated)
3      1319_at                        discoidin domain receptor family, member 2
Gene.Symbol
1      MLANA
2      CSPG4
3      DDR2

```

```

> doxorubicinGenes <- as.character(doxorubicinGenes[, "Probe.Set.ID"])
> doxorubicinGenes[1:3]

```

```
[1] "1051_g_at" "110_at" "1319_at"
```

3.3 Doxorubicin 2007 Data: Adria_ALL.txt

In the first correction to Potti et al. [4], Potti and Nevins [5] comment on “the acute lymphocytic leukemia dataset in which the labels are accurate — full details are provided on our web page.” Of the files available on the Potti et al. [4] web site (<http://data.genome.duke.edu/NatureMedicine.php>) as of November 6, 2007, only one, “Adria_ALL.txt” involves the test samples for doxorubicin. This file contains processed data for 144 samples: 22 training samples and 122 validation samples. A note to the side (in position EP3 when the table is read into Excel) states that “Validation data is from GSE4698, GSE649, GSE650, GSE651, and others.”

We trimmed off this comment to get a more consistently formatted table for easier loading (the untrimmed table was saved as Adria_ALLOriginal.txt). There are two header lines. Row 1 indicated whether the column is training (Adria) or testing (Validation). Row 2 indicates whether the column is Sens or Resistant (training) or Resp or NR (testing). We load and parse these headers below.

```

> tempDoxorubicin07Header1 <- read.table(file.path("RawData", "PottiNatMed",
+ "Adria_ALL.txt"), sep = "\t", nrows = 1, header = FALSE)
> tempDoxorubicin07Header1 <- as.vector(t(tempDoxorubicin07Header1))
> tempDoxorubicin07Header2 <- read.table(file.path("RawData", "PottiNatMed",
+ "Adria_ALL.txt"), sep = "\t", skip = 1, nrows = 1, header = FALSE)
> tempDoxorubicin07Header2 <- as.vector(t(tempDoxorubicin07Header2))
> tempDoxorubicin07Header1[1:25]

```

```

[1] "Adria0"      "0"           "0"           "0"           "0"
[6] "0"           "0"           "0"           "0"           "0"
[11] "Adria1"     "1"           "1"           "1"           "1"
[16] "1"           "1"           "1"           "1"           "1"
[21] "1"           "1"           "Validation2" "2"           "2"

```

```

> tempDoxorubicin07Header2[1:25]

```

```

[1] "Resistant" "Resistant" "Resistant" "Resistant" "Resistant" "Resistant"
[7] "Resistant" "Resistant" "Resistant" "Resistant" "Sens"       "Sens"
[13] "Sens"       "Sens"       "Sens"       "Sens"       "Sens"       "Sens"
[19] "Sens"       "Sens"       "Sens"       "Sens"       "NR"         "NR"
[25] "Resp"

```

```
> table(tempDoxorubicin07Header1)

tempDoxorubicin07Header1
      0      1      2      Adria0      Adria1 Validation2
      9     11    120         1         1         2
```

```
> table(tempDoxorubicin07Header2)

tempDoxorubicin07Header2
  NR Resistant  Resp  Sens
   99         10    23   12
```

In the training set, 0 denotes resistant and 1 denotes sensitive. There are 10 lines labeled as resistant, and 12 lines labeled as sensitive. The test (validation) data is all labeled 2, regardless of status. The test data contains 99 NR (resistant) samples and 23 Resp (sensitive) samples, matching the numbers reported by Potti et al. [4]

We now name the samples and assemble the sample information.

```
> tempSampleNames <- c(paste("Training", c(1:22), sep = ""), paste("Test",
+   c(1:122), sep = ""))
> tempGroup <- c(rep("Training", 22), rep("Test", 122))
> tempStatus <- tempDoxorubicin07Header2
> tempStatus[tempStatus == "Sens"] <- "Sensitive"
> tempStatus[tempStatus == "NR"] <- "Resistant"
> tempStatus[tempStatus == "Resp"] <- "Sensitive"
> doxorubicin07Info <- data.frame(sampleGroup = tempGroup, status = tempStatus,
+   row.names = tempSampleNames)
> doxorubicin07Info[c(1:2, 22:25), ]
```

	sampleGroup	status
Training1	Training	Resistant
Training2	Training	Resistant
Training22	Training	Sensitive
Test1	Test	Resistant
Test2	Test	Resistant
Test3	Test	Sensitive

```
> rm(list = ls(pattern = "^temp"))
```

Finally, we load the numerical quantifications.

```
> doxorubicin07Numbers <- read.table(file.path("RawData", "PottiNatMed",
+   "Adria_ALL.txt"), sep = "\t", skip = 2, header = FALSE)
> colnames(doxorubicin07Numbers) <- rownames(doxorubicin07Info)
> doxorubicin07Numbers[1:4, c(1:2, 22:25)]
```

	Training1	Training2	Training22	Test1	Test2	Test3
1	1.18	1.12	1.22	0.60	3.53	2.16
2	1.75	4.02	0.63	0.71	0.63	0.30
3	0.13	0.35	2.54	0.97	0.29	1.67
4	0.19	0.42	0.94	0.71	1.86	3.11

3.4 Doxorubicin 2008 Data: Adria_ALL_data1_n95.doc

The August 2008 correction of Potti et al. [6] notes that 27 of the 122 samples used in their development of a signature for doxorubicin were replicates, not independent observations, and that they have now redone their analysis using only the 95 “unique” samples. Their revised figure shows 74 resistant samples and 21 sensitive samples. The correction also adds two more GEO identifiers (GSE649 and GSE2351) to accompany the GSE650 and GSE651 entries initially listed, in order to “more clearly identify the sources of the samples”. Checking the Potti et al. [4] web site shows that Adria_ALL.txt has been replaced with Adria_ALL_data1_n95.doc. This file states that

we have recently realized that there were replicates in the data provided by the investigators at St. Jude constituting one of the independent validations. It is clear that 27 samples had been inadvertently replicated. The duplication appears to be result from the fact that the same data was posted under several separate series files in GEO (GSE649, GSE650 and GSE651) in 2004 and 2005. The corrected sample list, with clinical response data, is listed below.

We have copied the tabular data into a csv file, Adria_ALL_data1_n95_table.csv, for easier loading.

```
> doxorubicin08Table <- read.table(file.path("RawData", "PottiNatMed",
+     "Adria_ALL_data1_n95_table.csv"), sep = ",", header = TRUE,
+     nrows = 95)
> dim(doxorubicin08Table)
```

```
[1] 95  2
```

```
> doxorubicin08Table[1:3, ]
```

```
Sample.ID Response
1  GSM44303      RES
2  GSM44304      RES
3  GSM9653       RES
```

3.5 Holleman et al. [2] Data

Of the GEO data sets now listed, three (GSE649, GSE650, and GSE651) are part of a larger experiment described by Holleman et al. [2]. In this experiment, 173 samples from children with ALL were tested for sensitivity to each of four chemotherapeutic agents (asparagine, ASP, daunorubicin, DNR, prednisone, PRED, and vincristine, VCR). As noted in their second correction, Potti et al. [6] are focusing on anthracyclines, so it is the DNR measurements that are of interest here.

Sensitivity was assessed using an MTT assay, giving LC50 concentrations for the sample/drug combinations. Holleman et al. [2] set cutoff LC50 values for each drug, which were used to split the samples into “sensitive”, “intermediate”, and “resistant” categories, which were scored as 1-3 respectively. Supplementary table 1 of Holleman et al. [2] gives the LC50 cutoff values that were used for classifying daunorubicin:

Daunorubicin Rules	
LC50 Value	Classification
$LC50 \leq 0.075$	Sensitive
$0.075 < LC50 < 0.114$	Intermediate
$0.114 \leq LC50$	Resistant

Not every sample was evaluated against every drug; e.g., there are only 147 numerical scores of 1-3 recorded for DNR.

The supplementary web site for Holleman et al. [2], <http://www.stjude.com/research/data/ALL4>, includes a “key.xls” file. This file contains two sheets, describing which samples are allocated to which categories, and how these samples were posted to GEO. We have extracted both sheets as individual csv files for easier loading.

```
> hollemanKeyGSM <- read.table(file.path("RawData", "HollemanNEJM",
+   "key_gsm.csv"), sep = ",", header = TRUE)
> dim(hollemanKeyGSM)

[1] 173  7

> hollemanKeyGSM[1:4, ]

  PharmGKBSubject.ID   gsm IPT ASP.score DNR.score PRED.score VCR.score
1          PA126710896 GSM9707  B         1         1         2         2
2          PA126710897 GSM9708  B         1         1         2         2
3          PA126710898 GSM9808  B         3         3         2         3
4          PA126710899 GSM9709  B         1         1         2         2

> hollemanKeyGSE <- read.table(file.path("RawData", "HollemanNEJM",
+   "key_gse.csv"), sep = ",", header = TRUE)
> dim(hollemanKeyGSE)

[1] 17  3

> hollemanKeyGSE[1:3, ]

      group series score
1          all 173 GSE635  NA
2 ASP sensitive all ALL GSE643  1
3 ASP resistant all ALL GSE645  3
```

3.6 Lugthart et al. [3] Data

The final GEO dataset listed, GSE2351, comes from another paper from St. Jude: Lugthart et al. [3]. The 129 samples used in this study are for the most part a subset of the 173 samples used by Holleman et al. [2] The GEO entry for GSE2351 contains a table of clinical information giving (among other things) the mapping of GSM ids to LC50 values for the samples used. We have extracted this clinical information table as a csv file for easier loading.

```
> gse2351Clinical <- read.table(file.path("RawData", "GEO", "GSE2351",
+   "gse2351.csv"), sep = ",", header = TRUE)
> dim(gse2351Clinical)

[1] 129  9

> gse2351Clinical[1:4, ]
```


	code	gsm	LC50.ASP	LC50.DNR	LC50.PRED	LC50.VCR	Comp1	Comp2
1	n1001	GSM9653	0.0617	0.1985	250.0000	0.6859	-1.6724358	-0.13943068
2	n1002	GSM9654	0.0107	0.0682	0.1322	1.7681	0.1271482	-1.21809544
3	n1003	GSM9655	0.0032	0.0297	0.0595	0.0488	1.8210632	-0.03051928
4	n1004	GSM9656	0.0032	0.0290	0.1055	0.1784	1.4034015	-0.54825000
		Comp3						
1		1.3079215						
2		-0.1763163						
3		0.8716396						
4		0.6402796						

4 Identifying the Cell Lines Used in Adria_ALL.txt

4.1 Getting the NCI60 Quantifications

We begin by extracting the NCI60 array quantifications for the cell lines listed by Potti et al. [4] for doxorubicin, using the orientation from the “cell lines in each chemo predictor” file posted on the Potti et al. [4] web site as of August 2008. First, we list the lines involved.

```
> cellLinesUsed[["doxorubicin"]][["listPotti06CorrAug08"]]

$Sensitive
[1] "SF-539"      "SNB-75"      "MDA-MB-435" "NCI-H23"     "M14"
[6] "MALME-3M"   "SK-MEL-2"    "SK-MEL-28"  "SK-MEL-5"   "UACC-62"

$Resistant
[1] "NCI/ADR-RES" "HCT-15"      "HT29"       "EKVX"       "NCI-H322M"
[6] "IGROV1"      "OVCAR-3"     "OVCAR-4"    "OVCAR-5"    "OVCAR-8"
[11] "SK-OV-3"     "CAKI-1"
```

Adria_ALL.txt lists 10 resistant and 12 sensitive lines, whereas the listing above identifies 10 sensitive and 12 resistant. We put the group of size 10 first.

```
> doxorubicinNCI60Quants <-
+   novartisA[,
+     c(cellLinesUsed[["doxorubicin"]][
+       "listPotti06CorrAug08"])[["Sensitive"]],
+     cellLinesUsed[["doxorubicin"]][
+       "listPotti06CorrAug08"])[["Resistant"]])]
```

4.2 Transforming the Adria_ALL.txt Data

The posted data for doxorubicin has been transformed, apparently using the same set of operations as were applied for docetaxel. Specifically, the data were log-transformed, the values for each row (gene) were centered and scaled to have mean zero and unit variance (separately for training and test data), the data were exponentiated to undo the log-transform, and the final results were rounded to two decimal places. In order to match the numbers more precisely, we transform the NCI60 quantifications the same way.

```
> doxorubicinNCI60Scaled <-
+   round(exp(t(scale(t(log(doxorubicinNCI60Quants))))), 2)
```

4.3 Matching the First Probeset

Matching probesets is difficult due to a difference in dimensions:

```
> dim(doxorubicin07Numbers)
[1] 8958 144
> dim(doxorubicinNCI60Scaled)
[1] 12625 22
```

The Adria_ALL.txt data matrix has 8958 rows, which is a smaller number than we've worked with before. This comes about because the doxorubicin testing data was run on U133Av2 arrays, and the training data on U95Av2 arrays, necessitating a mapping across platforms. As noted by Potti et al. [4], this mapping was accomplished using "Chip Comparer", available at <http://tenero.duhs.duke.edu/genearray/perl/chip/chipcomparer.pl>. We ran Chip Comparer to get mappings going each way (the row orderings are different depending on which platform is named first), and found that there were 8958 distinct Locus Link clusters that the data were mapped to.

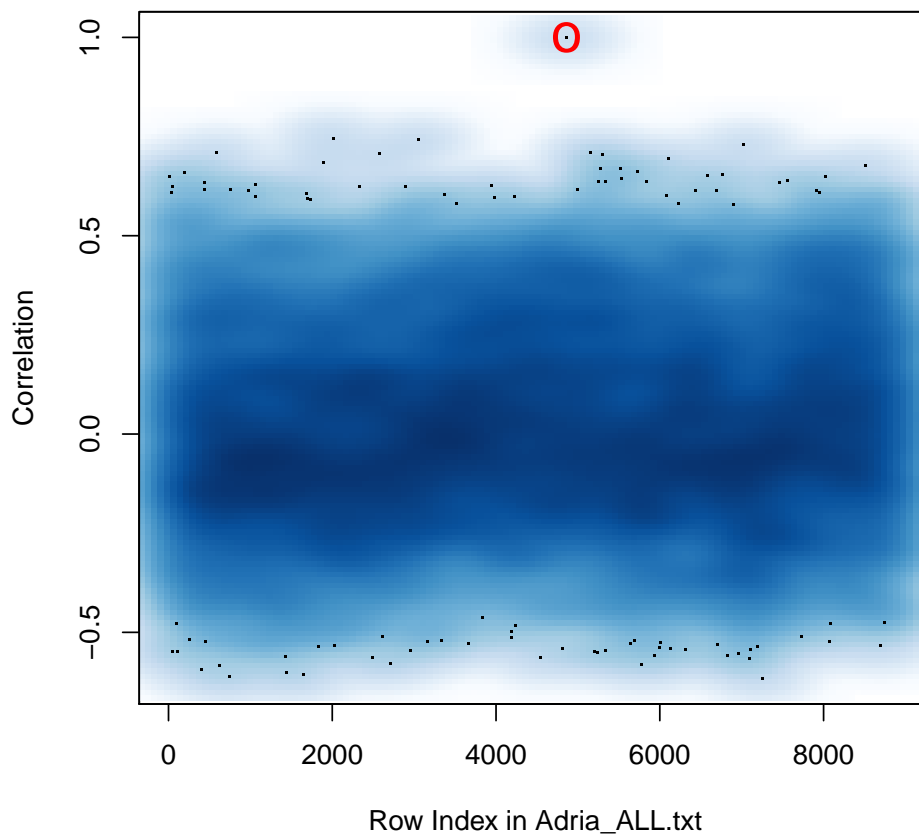
We then used correlation to try matching the first NCI60 probeset by brute force.

```
> temp <- cor(doxorubicinNCI60Scaled[1, ], t(doxorubicin07Numbers[,
+     1:22]))
> max(temp)
[1] 1
> which.max(temp)
[1] 4860
```

There is a perfect match for the first set of probeset values. To see if this stands out clearly from the rest of the data, we plot a smoothed version of the correlations obtained.

```
> smoothScatter(1:length(temp), temp, xlab = "Row Index in Adria_ALL.txt",
+   ylab = "Correlation", main = paste("Correlation with Scaled Values for",
+   rownames(doxorubicinNCI60Scaled)[1]))
> points(which.max(temp), max(temp), pch = "o", col = "red", cex = 2)
```

Correlation with Scaled Values for 36460_at



Nothing else is even close. We now do a quick survey of the other probesets to see if we can obtain similarly good matches for them.

```
> tempIndices <- apply(doxorubicin07Numbers[, 1:22], 1, function(x) {
+   which.max(cor(x, t(doxorubicinNCI60Scaled)))
+ })
> tempCors <- apply(cbind(doxorubicin07Numbers[, 1:22], doxorubicinNCI60Scaled[tempIndices,
+   ]), 1, function(x) {
+   cor(x[1:22], x[23:44])
+ })
> min(tempCors)

[1] 0.9999924
```

The correlations are all perfect modulo discrepancies that can be explained by rounding error. This shows that in this instance, Chip Comparer was applied by simply “picking one” of the U95Av2 probeset

ids mapping to a given Locus Link id. We use this maximal correlation mapping to assign names (U95Av2 probeset ids) to the rows of the Adria_ALL.txt matrix.

```
> rownames(doxorubicin07Numbers) <- rownames(doxorubicinNCI60Scaled)[tempIndices]
> doxorubicin07Numbers[1:3, 1:5]
```

	Training1	Training2	Training3	Training4	Training5
35753_at	1.18	1.12	3.46	0.65	3.07
36138_at	1.75	4.02	0.43	0.31	0.76
41765_at	0.13	0.35	1.13	1.14	0.84

```
> doxorubicinNCI60Scaled[rownames(doxorubicin07Numbers)[1:3], 1:5]
```

	SF-539	SNB-75	MDA-MB-435	NCI-H23	M14
35753_at	1.18	1.12	3.46	0.65	3.07
36138_at	1.75	4.02	0.43	0.31	0.76
41765_at	0.13	0.35	1.13	1.14	0.84

```
> rm(list = ls(pattern = "^temp"))
```

The perfect match shows that we are working with the same cell lines, though the sensitive/resistant labeling has been reversed.

5 Clustering the Adria_ALL.txt Test Data

In order to see if the expression values for the reported signature genes clearly separated the test set responders from nonresponders, we extracted and clustered the relevant expression submatrix. To do this, we first identified the subset of signature genes that were retained as chip comparer condensed the initial list of U95A probesets down to 8958.

```
> match(doxorubicinGenes, rownames(doxorubicin07Numbers))
```

```
[1] NA 2776 3087 485 888 3555 5515 1961 5367 114 2348 7081 8873 1772 7282
[16] 1820 5684 6583 5012 2216 142 1067 4167 7362 7025 6549 NA 337 6408 NA
[31] 7286 1971 651 1573 55 NA 5652 213 NA 2181 3702 5764 NA 144 2734
[46] 2276 NA 909 5788 NA 5497 1519 5803 NA 2128 376 6305 NA NA 7563
[61] 1620 1018 5501 1927 637 343 NA 6343 7779 6157 NA 7180 1937 6315 7436
[76] NA 803 2570 NA 2235
```

```
> sum(is.na(match(doxorubicinGenes, rownames(doxorubicin07Numbers))))
```

```
[1] 15
```

```
> keyRows <- match(doxorubicinGenes, rownames(doxorubicin07Numbers))
```

```
> keyRows <- keyRows[!is.na(keyRows)]
```

Next, we extract the expression submatrix for the test data (the last 122 columns) and log transform it for clustering. Since the values were rounded to two decimal places, we impose a floor of $\log(0.01)$; all values below this value are replaced with this value.

```
> doxorubicin07Submatrix <- as.matrix(log(doxorubicin07Numbers[keyRows,
+   23:144]))
> doxorubicin07Submatrix[doxorubicin07Submatrix < -4.7] <- log(0.01)
```

Now we define colors corresponding to status; red for Resistant, and blue for Sensitive.

```
> doxorubicin07TestColors <- rep("red", 122)
> doxorubicin07TestColors[doxorubicin07Info[23:144, "status"] ==
+   "Sensitive"] <- "blue"
```

Finally, we draw a heatmap using the jet colormap and adding boxes to highlight similar features. The box boundaries were found by trial and error. Box placement seems to involve some odd interactions with the graphics window, so we invoke `dev.off` first to clear any unwanted carryover,

```
> dev.off()

null device
  1

> dev.new()
```

and then we draw the heatmap and explicitly save the figure with `dev.copy2pdf` and `dev.copy2eps` as opposed to simply setting `fig=TRUE` in the Sweave invocation.

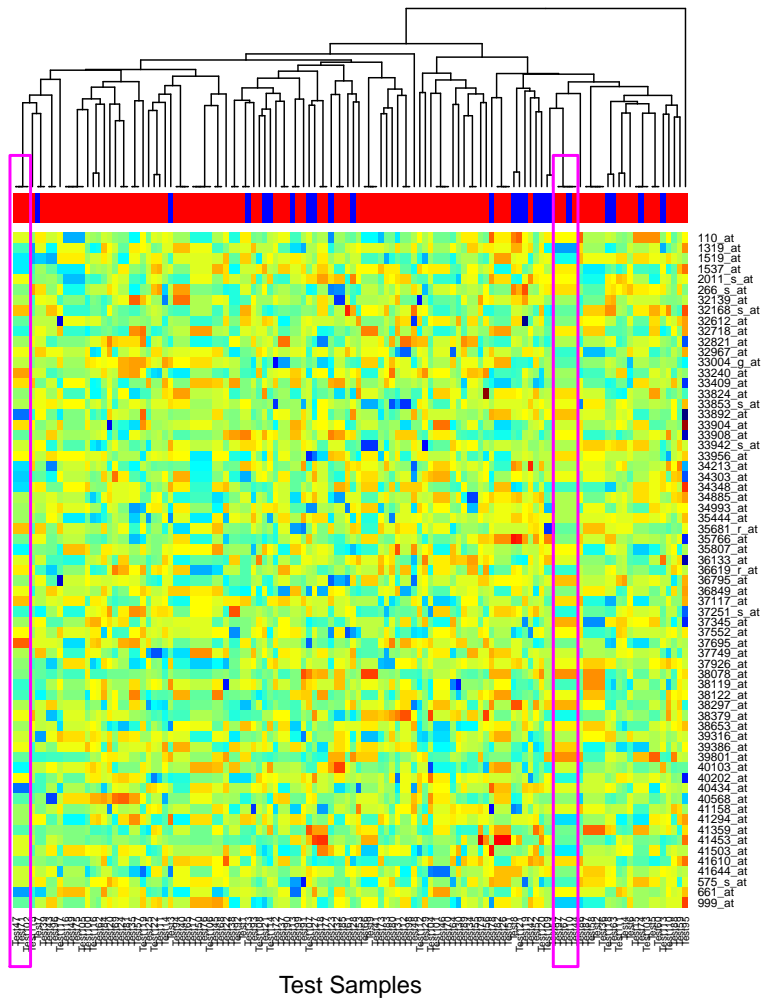
```
> source(file.path("Scripts", "jet.colors.R"))
> nSubmatrixRows <- dim(doxorubicin07Submatrix)[1]
> heatmap(doxorubicin07Submatrix[nSubmatrixRows:1, ], col = jet.colors(64),
+   ColSideColors = doxorubicin07TestColors, Rowv = NA, xlab = "Test Samples")
> par(xpd = TRUE)
> tempBlock <- c(-0.15, 0.95)
> names(tempBlock) <- c("Bottom", "Top")
> tempBlock1 <- c(0, 0.025)
> tempBlock2 <- c(0.66, 0.69)
> names(tempBlock1) <- c("Left", "Right")
> names(tempBlock2) <- c("Left", "Right")
> lines(tempBlock1[c("Left", "Left", "Right", "Right", "Left")],
+   tempBlock[c("Bottom", "Top", "Top", "Bottom", "Bottom")],
+   col = "magenta", lwd = 2)
> lines(tempBlock2[c("Left", "Left", "Right", "Right", "Left")],
+   tempBlock[c("Bottom", "Top", "Top", "Bottom", "Bottom")],
+   col = "magenta", lwd = 2)
> par(xpd = FALSE)
> dev.copy2pdf(file = file.path("Figures", "doxoSubmatrixWithBoxes.pdf"))

quartz
  2

> dev.copy2eps(file = file.path("Figures", "doxoSubmatrixWithBoxes.eps"))

quartz
  2
```

We then load the image produced.



Several blocks of columns show identical expression patterns, suggesting that the same same samples have been used multiple times.

6 Examining the Adria_ALL.txt Test Data

Since the test data columns were not named, we decided to check the correlation structure to see if there were obvious clusters.

```
> doxorubicin07Cors <- cor(doxorubicin07Numbers)
> sum(doxorubicin07Cors > 0.9999)
```

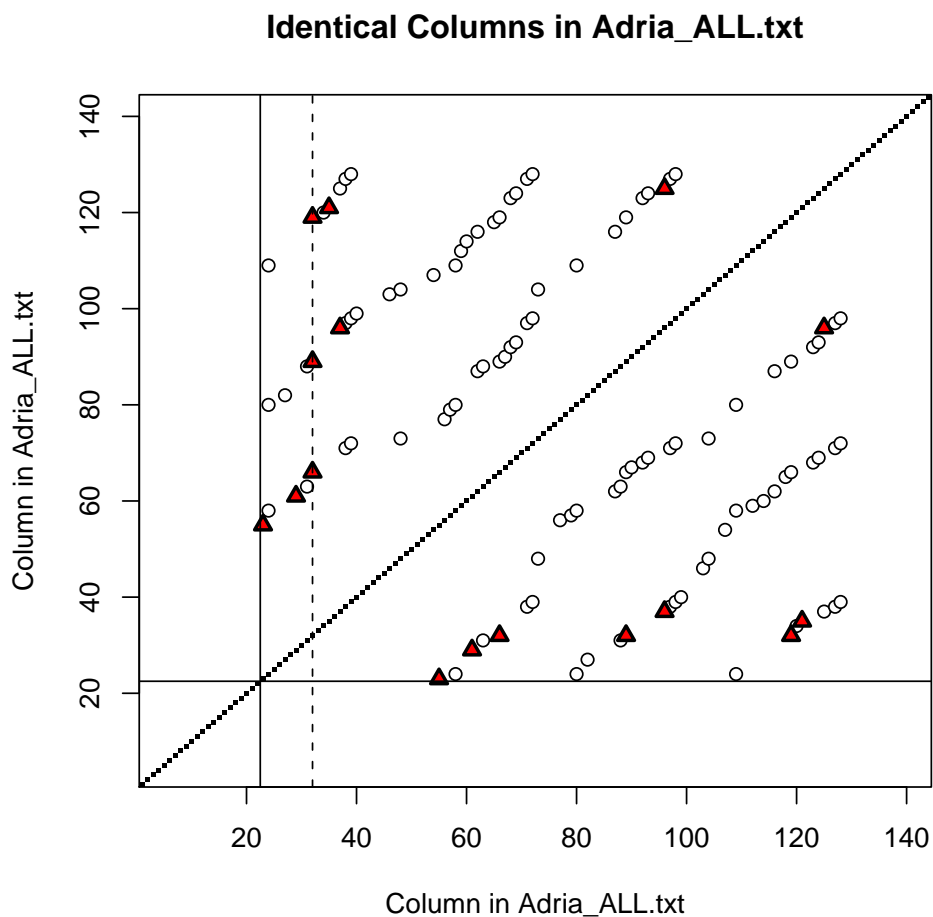
```
[1] 256
```

```
> sum(diag(doxorubicin07Cors) > 0.9999)
```

[1] 144

Some columns are not distinct – in some cases different columns are from the same sample. The structure may be more apparent graphically.

```
> doxorubicin07HighCors <- (doxorubicin07Cors > 0.9999)
> same07Status <- matrix(rep(doxorubicin07Info[, "status"], 144),
+   144, 144)
> same07Status <- (same07Status == t(same07Status))
> temp07Ties <- which(doxorubicin07HighCors & same07Status, arr.ind = TRUE)
> temp07Ties <- temp07Ties[temp07Ties[, 1] != temp07Ties[, 2],
+   ]
> temp07BadTies <- which(doxorubicin07HighCors & (!same07Status),
+   arr.ind = TRUE)
> plot(1:144, 1:144, pch = ".", cex = 3, xlim = c(0.5, 144.5),
+   ylim = c(0.5, 144.5), xaxs = "i", yaxs = "i", xlab = "Column in Adria_ALL.txt",
+   ylab = "Column in Adria_ALL.txt", main = "Identical Columns in Adria_ALL.txt",
+   )
> points(temp07Ties[, 1], temp07Ties[, 2], pch = 21, bg = "white",
+   lwd = 1)
> points(temp07BadTies[, 1], temp07BadTies[, 2], pch = 24, bg = "red",
+   lwd = 2)
> abline(h = 22.5, v = 22.5)
> abline(v = 32, lty = "dashed")
```



A very regular structure appears in the correlation heatmap, almost band-diagonal in nature. Further, some of the tied columns (red triangles in the plot) have conflicting sensitive/resistant labels – the same sample has been labeled both sensitive and resistant. The ties are clearly confined to the test data (columns 23-144), though the last 16 samples appear not to have any ties. A dashed line indicates one sample present 4 times that was labeled Sensitive once and Resistant three times; these were boxed on the right in the previous figure.

We now tabulate the number of independent samples and the structure of the ties.

```
> nTest07Copies <- apply(doxorubicin07HighCors[23:144, 23:144],
+   1, sum)
> table(nTest07Copies)

nTest07Copies
 1  2  3  4
60 28 18 16

> sum(table(nTest07Copies)/c(1:4))
```



```
[1] 84
```

```
> table(nTest07Copies, doxorubicin07Info[23:144, "status"])
```

nTest07Copies	Resistant	Sensitive
1	42	18
2	25	3
3	17	1
4	15	1

Only 84 of the 122 test samples are distinct. The cross-tabulation by status emphasizes the additional conflict noted above of the same sample being classified both ways. There are 4 samples replicated 4 times each, accounting for 16 of the columns. Of these 16, however, exactly 15 are classed as resistant. One of these 4 samples is classified as resistant 3 times, and sensitive 1 time.

7 Examining Doxorubicin After the August 2008 Correction

7.1 Looking for Ties

In the previous section, we saw only 84 distinct test samples. The second correction to Potti et al. [6], however, asserts that there are 95. Thus, our first check is whether all entries in the table supplied are unique.

```
> same08Name <- matrix(rep(doxorubicin08Table[, "Sample.ID"], 95),
+   95, 95)
> same08Name <- (same08Name == t(same08Name))
> same08Status <- matrix(rep(doxorubicin08Table[, "Response"],
+   95), 95, 95)
> same08Status <- (same08Status == t(same08Status))
> table(apply(same08Name, 1, sum))
```

```
 1  2
65 30
```

```
> sum(table(apply(same08Name, 1, sum))/c(1:2))
```

```
[1] 80
```

```
> table(apply(same08Name & (!same08Status), 1, sum))
```

```
 0  1
83 12
```

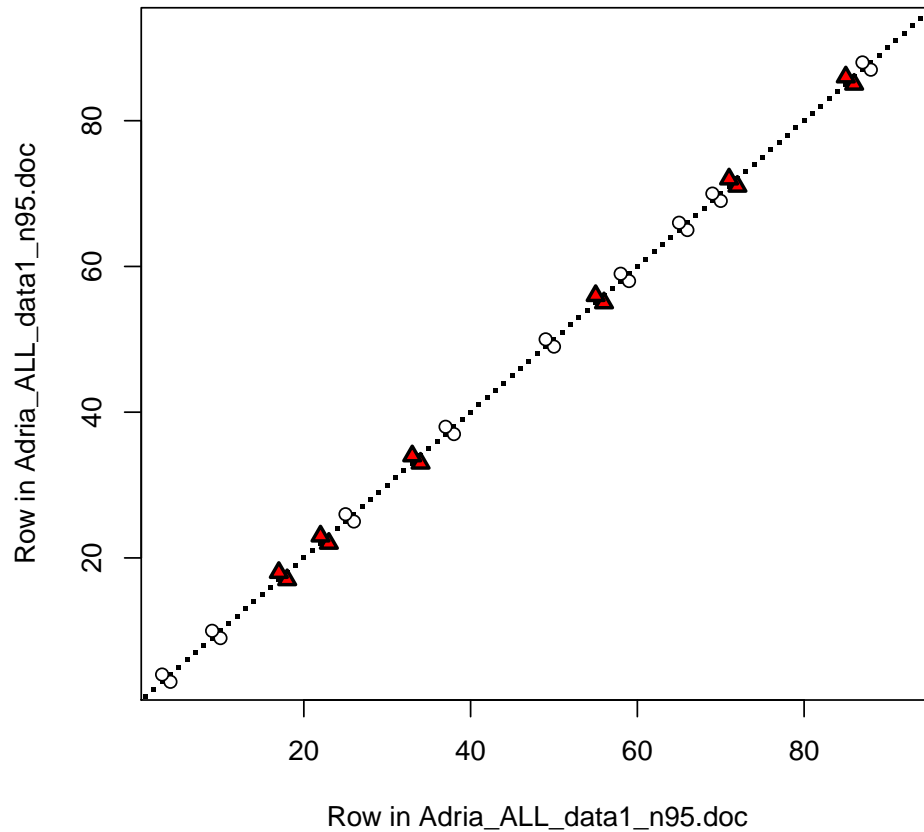
```
> sum(table(apply(same08Name & (!same08Status), 1, sum)) * c(0,
+   0.5))
```

```
[1] 6
```

The entries are not all distinct. Only 80 values are unique; there are 15 sets of pairs (but no sets of three or four). Of these 15 pairs, 6 are labeled as both sensitive and resistant. The structure may be more apparent graphically.

```
> same08Name <- matrix(rep(doxorubicin08Table[, "Sample.ID"], 95),
+   95, 95)
> same08Name <- (same08Name == t(same08Name))
> same08Status <- matrix(rep(doxorubicin08Table[, "Response"],
+   95), 95, 95)
> same08Status <- (same08Status == t(same08Status))
> temp08Ties <- which(same08Name & same08Status, arr.ind = TRUE)
> temp08Ties <- temp08Ties[temp08Ties[, 1] != temp08Ties[, 2],
+   ]
> temp08BadTies <- which(same08Name & (!same08Status), arr.ind = TRUE)
> plot(1:95, 1:95, pch = ".", cex = 3, xlim = c(0.5, 95.5), ylim = c(0.5,
+   95.5), xaxs = "i", yaxs = "i", xlab = "Row in Adria_ALL_data1_n95.doc",
+   ylab = "Row in Adria_ALL_data1_n95.doc", main = "Identical Rows in Adria_ALL_data1_n95.doc")
> points(temp08Ties[, 1], temp08Ties[, 2], pch = 21, bg = "white",
+   lwd = 1)
> points(temp08BadTies[, 1], temp08BadTies[, 2], pch = 24, bg = "red",
+   lwd = 2)
```

Identical Rows in Adria_ALL_data1_n95.doc



All of the duplicate entries are sequential, since the table supplied is sorted by Sample ID. For now, we simply identify the unique rows and tabulate the tied rows.

```
> distinct08GSMs <-
+   as.character(doxorubicin08Table$Sample.ID[
+     !duplicated(doxorubicin08Table$Sample.ID)])
> duplicate08GSMs <-
+   as.character(doxorubicin08Table$Sample.ID[
+     duplicated(doxorubicin08Table$Sample.ID)])
> duplicate08Rows <- which(same08Name & upper.tri(same08Name), arr.ind=TRUE)
> duplicate08Calls <-
+   data.frame(Row1 = duplicate08Rows[,1],
+             Row2 = duplicate08Rows[,2],
+             Call1 = as.character(doxorubicin08Table[duplicate08Rows[,1],
+             "Response"]),
+             Call2 = as.character(doxorubicin08Table[duplicate08Rows[,2],
```

```

+         "Response"]),
+         Inconsistent =
+         (doxorubicin08Table[duplicate08Rows[,1],"Response"]) !=
+         (doxorubicin08Table[duplicate08Rows[,2],"Response"]),
+         row.names = as.character(doxorubicin08Table[duplicate08Rows[,2],
+         "Sample.ID"]))
> duplicate08Calls

```

	Row1	Row2	Call1	Call2	Inconsistent
GSM9653	3	4	RES	RES	FALSE
GSM9658	9	10	SEN	SEN	FALSE
GSM9708	17	18	RES	SEN	TRUE
GSM9715	22	23	RES	SEN	TRUE
GSM9724	25	26	RES	RES	FALSE
GSM9738	33	34	RES	SEN	TRUE
GSM9745	37	38	SEN	SEN	FALSE
GSM9765	49	50	RES	RES	FALSE
GSM9772	55	56	RES	SEN	TRUE
GSM9777	58	59	RES	RES	FALSE
GSM9784	65	66	RES	RES	FALSE
GSM9790	69	70	RES	RES	FALSE
GSM9791	71	72	RES	SEN	TRUE
GSM9821	85	86	RES	SEN	TRUE
GSM9824	87	88	RES	RES	FALSE

7.2 Checking Consistency of St. Jude and Potti et al. [6] Sensitive/Resistant Calls

Since GSE2351 supplies the DNR LC50 values for the samples used by Lugthart et al. [3], and we have the rule that Holleman et al. [2] used to score samples, we can assign scores a la Holleman et al. [2] for all of the samples in either experiment. Using these, we can check the consistency between the Holleman et al. [2] calls and those of Potti et al. [6].

We first check the number of samples used by just one one of Holleman et al. [2] and Lugthart et al. [3].

```
> setdiff(distinct08GSMs, as.character(hollemanKeyGSM$gsm))
```

```
[1] "GSM44303" "GSM44304"
```

```
> setdiff(distinct08GSMs, as.character(gse2351Clinical$gsm))
```

```
[1] "GSM9756" "GSM9792" "GSM9793"
```

There are only two cases which were not assayed by Holleman et al. [2], and three cases for which we do not have LC50 values. For the two prior cases, we look at the LC50 values to score them.

```
> gse2351Clinical[gse2351Clinical$gsm == "GSM44303", c("code",
+ "gsm", "LC50.DNR")]
```

	code	gsm	LC50.DNR
128	nl176	GSM44303	0.0848

```
> gse2351Clinical[gse2351Clinical$gsm == "GSM44304", c("code",
+   "gsm", "LC50.DNR")]
      code      gsm LC50.DNR
129 nl177 GSM44304  0.0793
```

Since both LC50 values are between the cutoff values of 0.075 and 0.114, Holleman et al. [2] would have classed these two samples as Intermediate and assigned a score of 2 to each. We now assemble a table combining the Potti et al. [6] calls (SEN, RES, or Inconsistent), the Holleman et al. [2] scores (1, 2, 3), and the DNR LC50 values for the 80 distinct GSM IDs (where available).

```
> gsmTable <-
+   data.frame(HollemanScore=rep(0,80),
+             LC50.DNR=rep(0,80),
+             PottiCall=I(rep("",80)),
+             row.names=distinct08GSMs)
> for(tempGSM in 1:length(distinct08GSMs)){
+
+   tempHollemanRow <- which(hollemanKeyGSM$gsm == distinct08GSMs[tempGSM])
+   if(length(tempHollemanRow) == 1){
+     gsmTable[tempGSM,"HollemanScore"] <-
+       as.integer(as.character(hollemanKeyGSM$DNR.score[tempHollemanRow]))
+   }
+
+   tempGSE2351Row <- which(gse2351Clinical$gsm == distinct08GSMs[tempGSM])
+   if(length(tempGSE2351Row) == 1){
+     gsmTable[tempGSM,"LC50.DNR"] <-
+       gse2351Clinical[tempGSE2351Row,"LC50.DNR"]
+   }
+
+   tempPottiRow <- which(doxorubicin08Table$Sample.ID == distinct08GSMs[tempGSM])
+   if(length(tempPottiRow) == 1){
+     gsmTable[tempGSM,"PottiCall"] <-
+       as.character(doxorubicin08Table[tempPottiRow,"Response"])
+   }
+   if(length(tempPottiRow) == 2){
+     if(doxorubicin08Table$Response[tempPottiRow[1]] ==
+        doxorubicin08Table$Response[tempPottiRow[2]]){
+       gsmTable[tempGSM,"PottiCall"] <-
+         as.character(doxorubicin08Table[tempPottiRow[1],"Response"])
+     }
+     else{
+       gsmTable[tempGSM,"PottiCall"] <- "Inconsistent"
+     }
+   }
+ }
> rm(list=ls(pattern=~temp))
> gsmTable$LC50.DNR[gsmTable$LC50.DNR==0] <- NA
> gsmTable$HollemanScore[gsmTable$HollemanScore == 0] <- 2
```

Now we simply tabulate and cross-tabulate the scores.

```
> attach(gsmTable)
> table(HollemanScore)

HollemanScore
 1  2  3
48 10 22

> table(PottiCall)

PottiCall
Inconsistent      RES      SEN
           6          61          13

> table(PottiCall, HollemanScore)

           HollemanScore
PottiCall   1  2  3
Inconsistent  6  0  0
RES           29 10 22
SEN           13  0  0

> detach(gsmTable)
```

Of the 80 distinct samples, Holleman et al. [2] classified 48 as sensitive (a score of 1) and 22 as resistant (a score of 3). The other 10 samples are intermediate, including the 2 not examined by Holleman et al. [2]. These “intermediate” samples were not available in either GSE650 or GSE651. By contrast, Potti et al. [6] list only 13 as sensitive; 61 are labeled resistant and 6 are labeled both ways.

Cross-tabulating, we see that every sample that Holleman et al. [2] call resistant or intermediate is called resistant by Potti et al. [6], but only 13 (or maybe 19, depending on how we treat the inconsistent calls) of the samples that Holleman et al. [2] call sensitive are called sensitive by Potti et al. [6]: between 29 and 35 samples change classes. Further, both of the new samples (intermediates) are classed by Potti et al. [6] as resistant.

One possibility is that Potti et al. [6] are using a different LC50 cutoff value for DNR. We check this by looking at the ranges of values in each group.

```
> summary(gsmTable$LC50.DNR[gsmTable$PottiCall == "SEN"])

  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.   NA's
0.00200 0.01070 0.01900 0.02489 0.02625 0.07500 2.00000

> summary(gsmTable$LC50.DNR[gsmTable$PottiCall == "RES"])

  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.   NA's
0.0082 0.0309 0.0761 0.1160 0.1710 0.4196 1.0000

> summary(gsmTable$LC50.DNR[gsmTable$PottiCall == "RES"] < 0.075)

  Mode  FALSE  TRUE  NA's
logical   33   27    1
```

```

> sort(gsmTable$LC50.DNR[gsmTable$PottiCall == "SEN"])

[1] 0.0020 0.0063 0.0100 0.0114 0.0189 0.0190 0.0247 0.0252 0.0273 0.0540
[11] 0.0750

> sort(gsmTable$LC50.DNR[gsmTable$PottiCall == "RES"])

[1] 0.0082 0.0104 0.0110 0.0127 0.0149 0.0178 0.0181 0.0236 0.0244 0.0257
[11] 0.0257 0.0273 0.0276 0.0290 0.0297 0.0313 0.0398 0.0446 0.0461 0.0466
[21] 0.0596 0.0597 0.0671 0.0682 0.0691 0.0728 0.0732 0.0750 0.0750 0.0753
[31] 0.0769 0.0793 0.0811 0.0813 0.0848 0.0901 0.0932 0.0992 0.1078 0.1166
[41] 0.1191 0.1195 0.1292 0.1408 0.1667 0.1839 0.1897 0.1917 0.1985 0.1993
[51] 0.2115 0.2327 0.2707 0.2781 0.2981 0.3180 0.3555 0.3987 0.4185 0.4196

```

Unfortunately, this explanation doesn't work, since the value ranges overlap – the largest “sensitive” value is 0.075 (the cutoff set by Holleman et al. [2]), and almost half of the “resistant” values (27 cases) have LC50 values lower than this.

8 Looking at the Data Cited

Wholly aside from whether the calls agree, there is a question of whether the listed data sources make sense.

The initial Potti et al. [4] paper listed 122 samples, which corresponds exactly to the number of samples that Holleman et al. [2] classed as either sensitive or resistant to DNR (94 and 28 cases, respectively) and gave data for in GSE650 and GSE651. No explanation is given for not using all of these 122 samples or for including 10 new samples which would have been classed as intermediate. Almost all of the resistant cases from GSE651 are retained, but about half of the sensitive cases from GSE650 are dropped.

Listing GSE649 as a source makes no sense since the samples in that set are chosen because they were resistant to vincristine, not daunorubicin. Since this study was focused on anthracyclines (according to `Adria_ALL_data1_n95.doc`), the vincristine measurements should not be relevant. As shown in the report `checkingCellLines`, GI50 values for doxorubicin and vincristine do not track at all, so using one to assemble a signature for the other is problematic at best.

Listing GSE2351 as a source makes no sense as the Lugthart et al. [3] study was focused on assessing crossresistance to multiple therapies, and the signature was nominally focused on doxorubicin.

Since GSE650 and GSE651 contain no duplicates, the claim in the second correction to Potti et al. [6] that “27 samples were replicated owing to the fact that the same samples were included in several separate series files in the Gene Expression Omnibus generated in 2004 and 2005, which were the source of the data provided for the study” makes no sense. It also omits mention of the fact that the data posted for docetaxel at the time of the first correction also contained ties, though the GEO sources did not.

9 Assembling a Summary Figure and Table

Here, we simply condense figures and tables assembled above.

```

> ## This figure requires two panels, both identifying
> ## the presence of ties in the data.
>
> tempDoxo07Panel <- c(0.07, 0.47, 0.08, 0.88)
> tempDoxo08Panel <- c(0.57, 0.97, 0.08, 0.88)

```

```

> par(plt=tempDoxo07Panel)
> same07Status <- matrix(rep(doxorubicin07Info[, "status"], 144), 144, 144)
> same07Status <- (same07Status == t(same07Status))
> temp07Ties <- which(doxorubicin07HighCors & same07Status, arr.ind=TRUE)
> temp07Ties <- temp07Ties[temp07Ties[,1] != temp07Ties[,2],]
> temp07BadTies <- which(doxorubicin07HighCors & (!same07Status), arr.ind=TRUE)
> plot(1:144, 1:144, pch=".", cex=3,
+      xlim=c(0.5,144.5), ylim=c(0.5,144.5),
+      xaxs="i", yaxs="i",
+      xlab="Column in Adria_ALL.txt",
+      ylab="Column in Adria_ALL.txt",
+      main="Identical Columns in Adria_ALL.txt",
+      )
> points(temp07Ties[,1],temp07Ties[,2],pch=21,bg="white",lwd=1)
> points(temp07BadTies[,1],temp07BadTies[,2],pch=24,bg="red",lwd=2)
> abline(h=22.5, v=22.5)
> par(plt=tempDoxo08Panel, new=TRUE)
> same08Name <- matrix(rep(doxorubicin08Table[, "Sample.ID"], 95), 95, 95)
> same08Name <- (same08Name == t(same08Name))
> same08Status <- matrix(rep(doxorubicin08Table[, "Response"], 95), 95, 95)
> same08Status <- (same08Status == t(same08Status))
> temp08Ties <- which(same08Name & same08Status, arr.ind=TRUE)
> temp08Ties <- temp08Ties[temp08Ties[,1] != temp08Ties[,2],]
> temp08BadTies <- which(same08Name & (!same08Status), arr.ind=TRUE)
> plot(1:95, 1:95, pch=".", cex=3,
+      xlim=c(0.5,95.5), ylim=c(0.5,95.5),
+      xaxs="i", yaxs="i",
+      xlab="Row in Adria_ALL_data1_n95.doc",
+      ylab="Row in Adria_ALL_data1_n95.doc",
+      main="Identical Rows in Adria_ALL_data1_n95.doc")
> points(temp08Ties[,1],temp08Ties[,2],pch=21,bg="white",lwd=1)
> points(temp08BadTies[,1],temp08BadTies[,2],pch=24,bg="red",lwd=2)

quartz
  2

```

10 Appendix

10.1 File Location

```

> getwd()
[1] "/Users/kabagg/ReproRsch/WebSite"

```

10.2 Saves

10.3 SessionInfo

```

> sessionInfo()

```

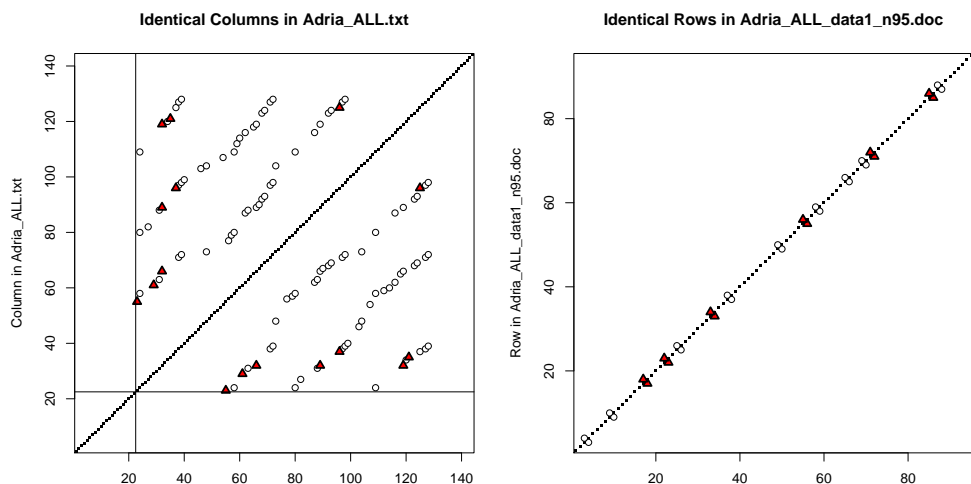



Figure 1: **A.** Locations of identical data columns in the quantifications supplied by Potti et al. [4] for the doxorubicin signature in November 2007. Lines mark the training/test boundary. There are no tied samples in the training data (lower left), but sensitive and resistant labels for the cell lines used were reversed. Of the 122 test samples (upper right), only 84 are distinct; some samples are present 2, 3, or 4 times (off-diagonal symbols). Further, some tied samples were labeled as sensitive in one column and resistant in another (red triangles). **B.** Locations of identical data rows in the table supplied by Potti et al. [6] for the same signature in August 2008. Names of 95 “unique” test data samples are given. Only 80 names are distinct; 15 are present in pairs, using the same symbols as in A. Of the 15 pairs, 6 involve one row labeled sensitive and another labeled resistant.

		Holleman et al. [2] Classifications			
		Sensitive	Intermediate	Resistant	
Potti et al. [4] Classifications	Sensitive	13	0	0	13
	Resistant	29	10	22	61
	Both	6	0	0	6
		48	10	22	

Table 1: Cross-tabulation of drug sensitivities assigned to the 80 distinct samples listed in the Potti et al. [6] doxorubicin table from August 2008, using both the rule from Holleman et al. [2], the source of this test data, and the labels attached by Potti et al. [6]. The actual Potti et al. [4] table has 95 entries; 15 names are present twice. Of those, 6 are listed as sensitive in one row and resistant in the other, and are classed as “Both” in the table above.

R version 2.9.1 (2009-06-26)

i386-apple-darwin8.11.1

locale:

en_US.UTF-8/en_US.UTF-8/C/C/en_US.UTF-8/en_US.UTF-8

attached base packages:

[1] stats graphics grDevices utils datasets methods base

other attached packages:

[1] genplotter_1.22.0 lattice_0.17-25 annotate_1.22.0

[4] AnnotationDbi_1.6.1 Biobase_2.4.1

loaded via a namespace (and not attached):

[1] DBI_0.2-4 grid_2.9.1 KernSmooth_2.23-2 RColorBrewer_1.0-2

[5] RSQLite_0.7-1 xtable_1.5-5

References

- [1] Coombes KR, Wang J, Baggerly KA: Microarrays: retracing steps. *Nat Med*, **13**:1276-7, 2007.
- [2] Holleman A, Cheek MH, den Boer ML, et al.: Gene-expression patterns in drug-resistant acute lymphoblastic leukemia cells and response to treatment. *N Engl J Med*, **351**:533-42, 2004.
- [3] Lugthart S, Cheek MH, den Boer ML, et al.: Identification of genes associated with chemotherapy crossresistance and treatment response in childhood acute lymphoblastic leukemia. *Cancer Cell*, **7**:375-86.
- [4] Potti A, Dressman HK, Bild A, et al: Genomic signatures to guide the use of chemotherapeutics. *Nat Med*, **12**:1294-1300, 2006.
- [5] Potti A, Nevins JR: Reply to Microarrays: retracing steps. *Nat Med*, **13**:1277-8, 2007.
- [6] Potti A, Dressman HK, Bild A, et al: Corrigendum to Genomic signatures to guide the use of chemotherapeutics. *Nat Med*, **14**:889, 2008.