

1 **Supplemental Article S1**

2 **Validating the automated method of detecting data availability**

3 Our method of identifying which articles create gene expression microarray data made a
4 nontrivial number of errors: about 10% of the articles it identified as creating gene expression
5 microarray data do not in fact create gene expression datasets.

6 The papers that are erroneously included in our subset to not create gene expression data, so they
7 certainly don't have associated archived datasets: all erroneously included papers were
8 automatically classified in the "no archived data" group.

9 If it were true that these erroneously-included articles received many more or many fewer
10 citations than other articles in the group, their inclusion could influence the findings of this study.
11 To verify our assumption that the influence of these mistakenly-included articles is in fact small,
12 we manually reviewed a random 226 of the 11k (get exact number) articles. Of these manually
13 reviewed articles, 206 did indeed create gene expression microarray data, and 20 did not (but
14 satisfied the boolean-search query for other reasons).

15 Examining the citations of the 20 articles that did not create gene expression data revealed that
16 these studies were cited less often than those that did create data: a mean of 26 citations compared
17 to a mean of 32 citations. The overall distribution of citations for articles that did not create gene
18 expression data is closer to zero than the distribution of citations for articles that did create gene
19 expression data.

20 We took steps to verify our assumption that the influence of articles erroneously identified these
21 mistakenly-included articles is in fact small. We began by manually reviewed a random 226 of
22 the 11k (get exact number) articles to identify those which we were assuming had created gene
23 expression microarray data but in fact had not.

24 We compared the distribution of those with errors to those without, calculated whether they were
25 statistically different, and ran a regression with the known-correct sample only.

```
26 nCitedBy      N=226, 4 Missing
27
28 +-----+-----+-----+-----+
29 |          |                               | N|nCitedBy|
30 +-----+-----+-----+-----+
31 |isCreated|   created-microarray-data|206|   31.86|
32 |          | created-microarray-data-not| 20|   26.30|
33 +-----+-----+-----+-----+
34 | Overall|                               |226|   31.37|
35 +-----+-----+-----+-----+
36
```

37

38 This difference, however, was found to be not statistically significantly different at the $p < 0.05$
39 level, using either a t-test on the log of the citation counts or a Wilcoxon rank sum test on the raw
40 citation counts.

```

1 print(ttest_citedby)
2   Welch Two Sample t-test
3
4 data: nCitedBy by isCreated
5 t = 0.5747, df = 22.61, p-value = 0.5712
6 alternative hypothesis: true difference in means is not equal to 0
7 95 percent confidence interval:
8  -14.47  25.59
9 sample estimates:
10  mean in group created-microarray-data
11                                31.86
12  mean in group created-microarray-data-not
13                                26.30
14
15 print(ttest_log_citedby)
16   Welch Two Sample t-test
17
18 data: log(1 + nCitedBy) by isCreated
19 t = 1.331, df = 21.77, p-value = 0.1968
20 alternative hypothesis: true difference in means is not equal to 0
21 95 percent confidence interval:
22  -0.2003  0.9175
23 sample estimates:
24  mean in group created-microarray-data
25                                2.991
26  mean in group created-microarray-data-not
27                                2.632
28
29 print(wilcox_citedby)
30   Wilcoxon rank sum test with continuity correction
31
32 data: nCitedBy by isCreated
33 W = 2440, p-value = 0.1733
34 alternative hypothesis: true location shift is not equal to 0

```

35 To confirm that the erroneously-included articles were not driving the findings about the citation
36 relationship with data availability, we ran a multivariate regression analysis on the subsample of
37 206 articles that we manually determined did in fact generate gene expression microarray data.
38 The estimated effect is statistically significant and similar to the findings from the whole sample.

```

39 gfm_table(anova(annotated_merged_created))
40 |-----| Df | Sum Sq | Mean Sq | F value | Pr(>F) |
41 |-----|-----|-----|-----|-----|-----|
42 | rcs(pubmed.date.in.pubmed, 3) | 2.00 | 83.82 | 41.91 | 73.91 | 0.00 |
43 | rcs(journal.impact.factor.tr, 3) | 2.00 | 18.69 | 9.35 | 16.48 | 0.00 |
44 | rcs(num.authors.tr, 3) | 2.00 | 4.03 | 2.01 | 3.55 | 0.03 |
45 | rcs(last.author.num.prev.pmc.cites.tr, 3) | 2.00 | 4.79 | 2.40 | 4.22 | 0.02 |
46 | factor(country.usa) | 1.00 | 0.05 | 0.05 | 0.09 | 0.77 |
47 | factor(dataset.in.geo.or.ae) | 1.00 | 5.68 | 5.68 | 10.03 | 0.00 |
48 | Residuals | 177.00 | 100.37 | 0.57 | | |
49 calcCI.exp(annotated_merged_created, "factor(dataset.in.geo.or.ae).L")
50 param est ciLow ciHigh p
51 Estimate factor(dataset.in.geo.or.ae).L 1.32 1.11 1.57 0.002

```