

Supplemental File 7. Comparison with public expression data on NSCLC

A number of microarray-based studies have been carried out on lung cancer. A logical question would be how our data compares with these previous studies. From public gene expression data on lung cancer in GEO and ArrayExpress databases, we identified 12 data sets each with both tumor and normal specimens.

To avoid the complications from meta-analysis (different platforms, labs, protocols, etc.), comparisons were done at the gene set level, not on expression values of individual genes. All array data were normalized with the *RMA* method, and DEGs were obtained by using the *limma* package in R with the adjusted p-value of 0.001. The DEG sets from various experiments were clustered by hierarchical method. The distance metric was the hypergeometric p-value of testing co-membership between two gene sets.

Clustering of public experiments (see below) indicates that our result is closest to that of GSE19804 which profiled never-smoker female lung adenocarcinoma in Taiwan. The fact that a study on a patient with virtually identical demographic characteristics yielded results most similar to ours provides an additional validation of the current study.

Patient grouping according to the mutation subtype is critical for selecting treatment methods. GSE31210 contains mutation information for EGFR, KRAS, ALK genes that are established as key biomarkers for subtyping lung adenocarcinoma. Comparing up and down DEG sets, our patient group was shown to be closest to the no-mutation group in GSE31210. Again, this provides further validation of our study given that all of our patients belong to no-mutation group. In sum, comparison with publicly available data not only confirmed the validity of our data set, but also raises an important and testable hypothesis that never-smoker female adenocarcinoma patients constitute a distinctive group from those with EGFR, KRAS, and ALK mutations.

