

Supplementary Figure 7. RT-PCR versus microarray for the 30 malignancy-risk genes: prediction performance, PCA comparison, and correlation plot.

1. Prediction performance:

- a. Both microarray expression and RT-PCR expression were standardized by centered at mean and scaled by standard deviation (SD). That is, for each gene, we calculate the mean and SD among the 21 samples, then generate a standardized score, $(x - \text{mean}(x)) / \text{sd}(x)$ where x is an expression value.
- b. Classifier by support vector machine (SVM): We used the SVM to build a classifier using the 30 genes microarray expression. Leave-One-Out-Cross-Validation (LOOCV) was used to assess the quality of the classifier. The accuracy rate was 86%.
- c. Prediction: The classifier was used to evaluate PR-PCR platform. That is, the classifier converted the 30 genes RT-PCR expressions into a predicted value for each sample to determine which group each tissue belongs to. The prediction results yielded a 90% accuracy rate.

| | | Predicted Outcome | | |
|--------------|-----------------|-------------------|-----------------|-----|
| | | Normal | IDC-like normal | IDC |
| True Outcome | Normal | 7 | 1 | 0 |
| | IDC-like normal | 0 | 4 | 1 |
| | IDC | 0 | 0 | 8 |

2. Comparison of various PCAs: We performed two ways to compare PCAs. The first one is to compare correlation of microarray and RT-PCR at each PCA. Results in Figure A showed a very high correlation in PCA1 ($r=0.95$). In contrast, the other PCAs had a weak correlation ($r<0.5$). The second approach is to test each PCA (PCA1, PCA2, and PCA3) to see if it is associated with cancer status (i.e., normal, IDC-like normal, to, IDC). Results in Figure B showed that an increasing trend of the risk score (from normal to IDC tissues) by PCA1 in the first column panel. On the other hand, the PCA2 and PCA3 did not yield their association with cancer status.

Figure A: Correlation of microarray and RT-PCR at each PCA

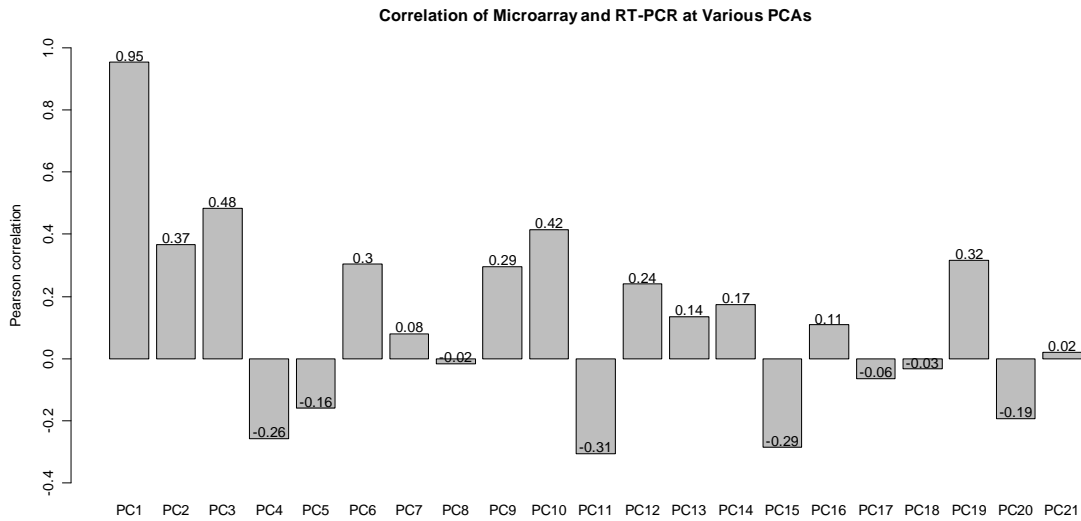
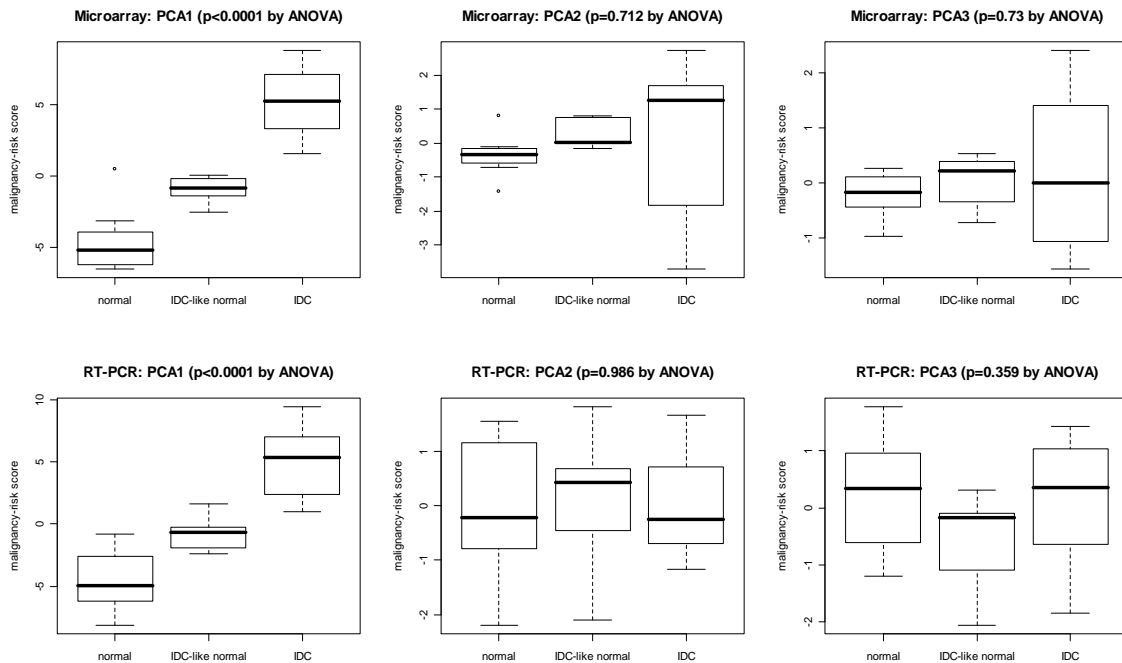
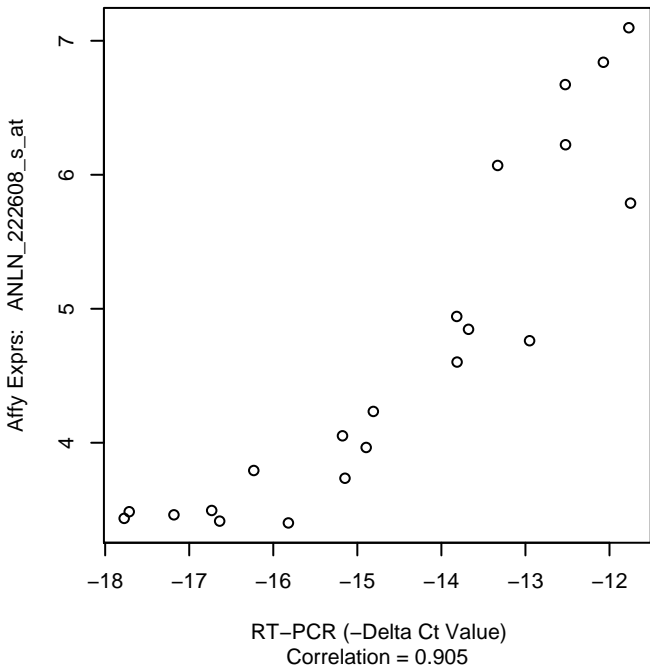
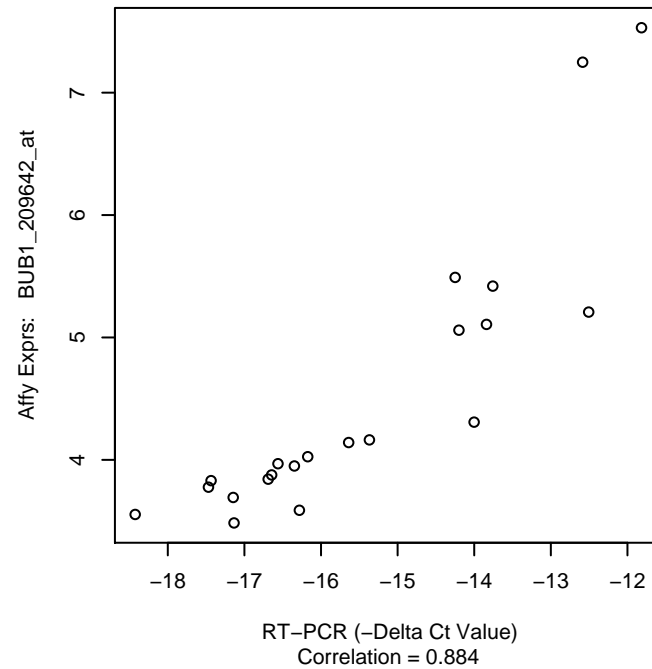
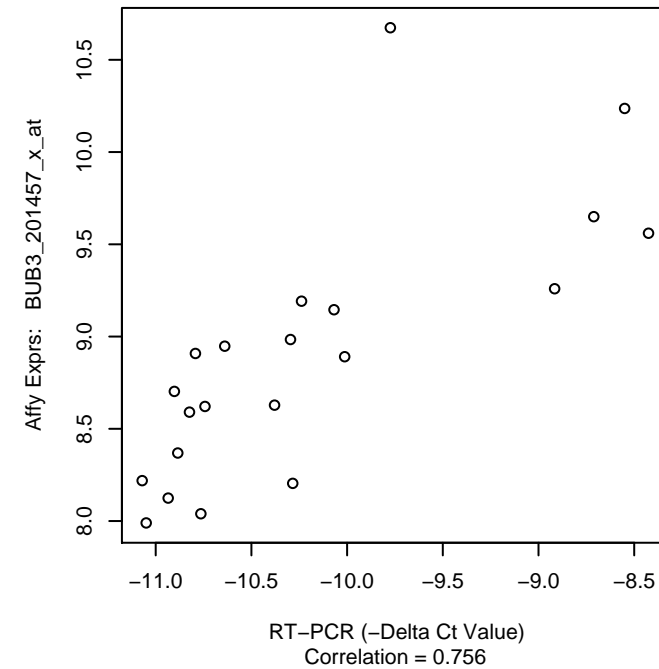
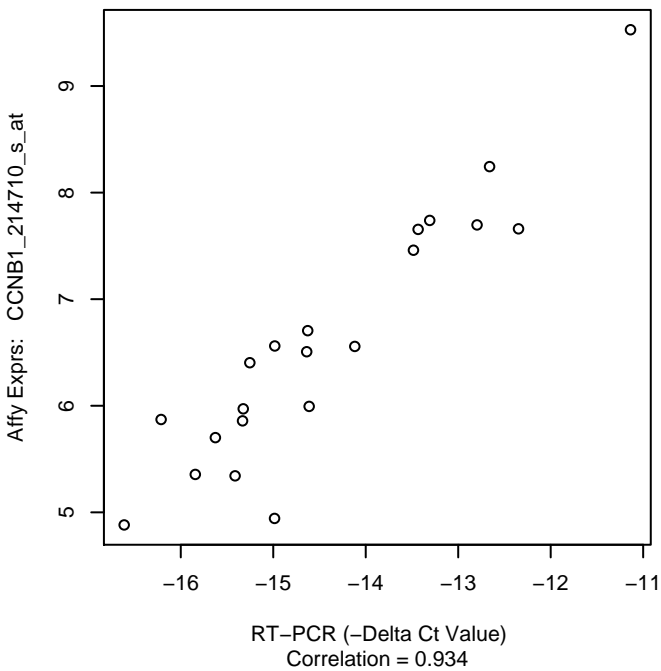
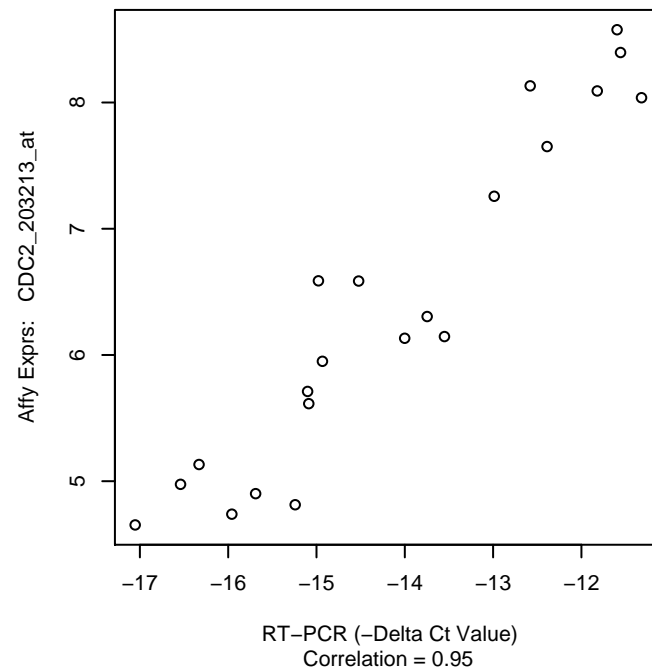
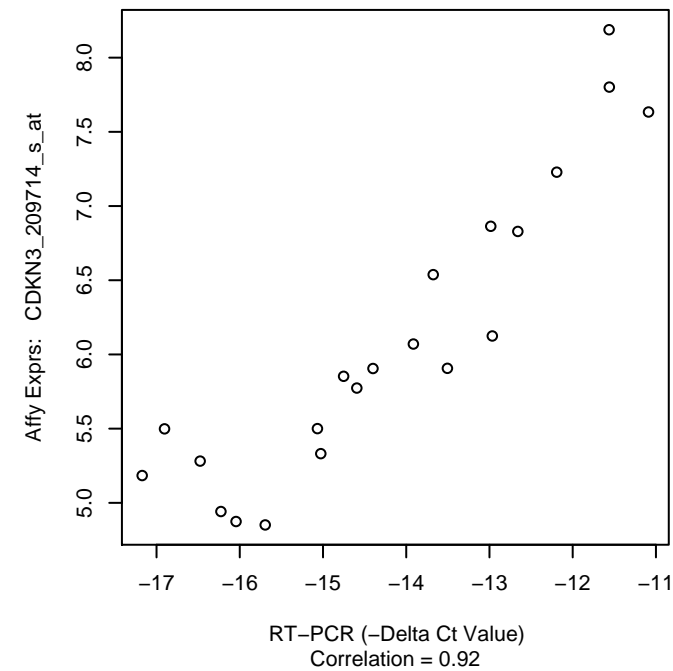
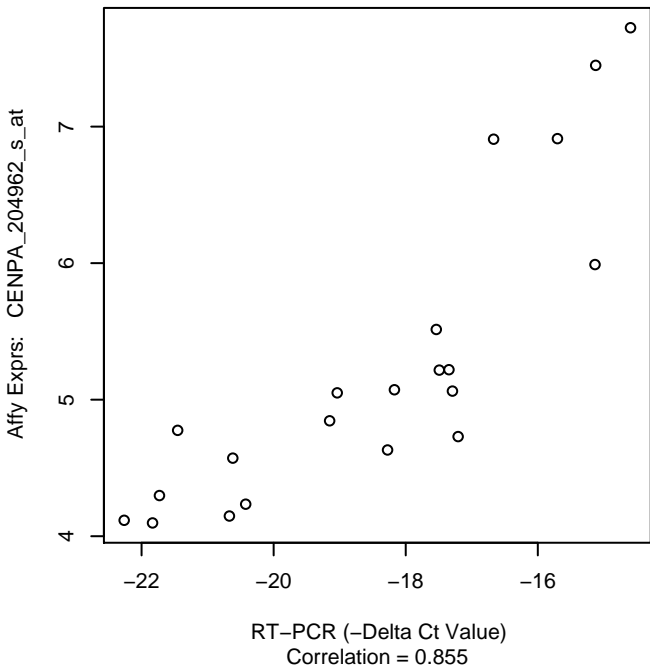
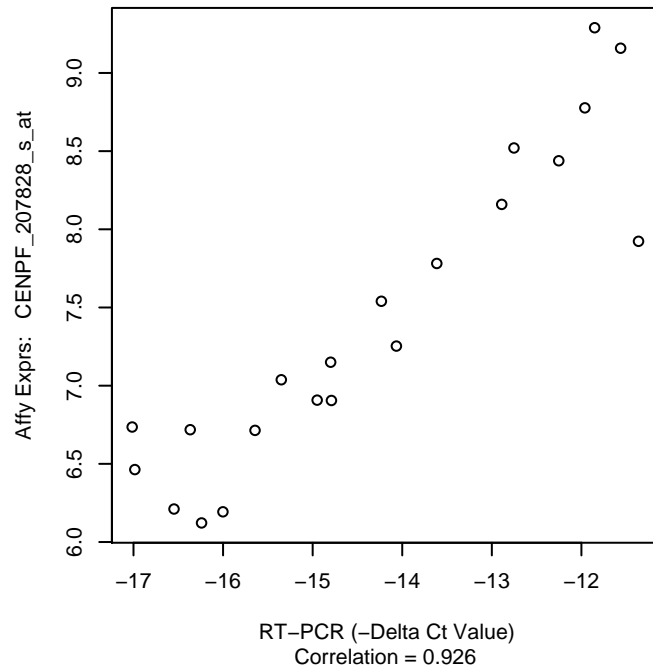
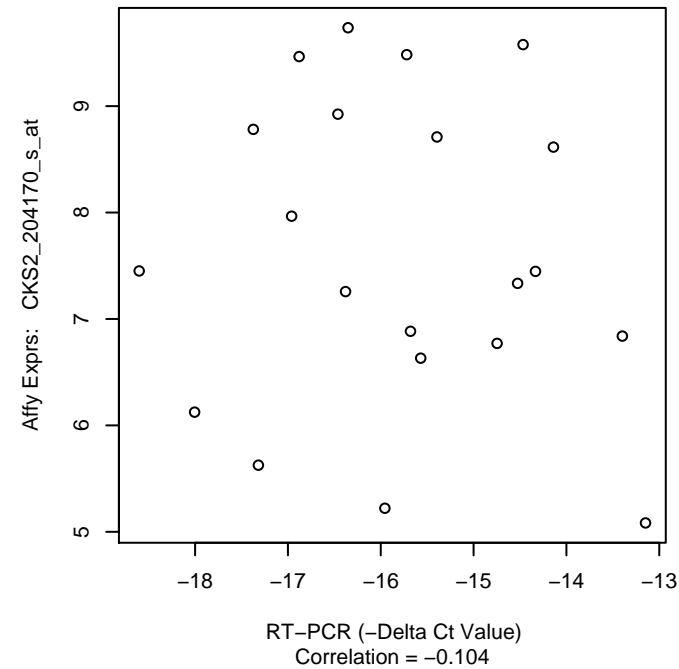
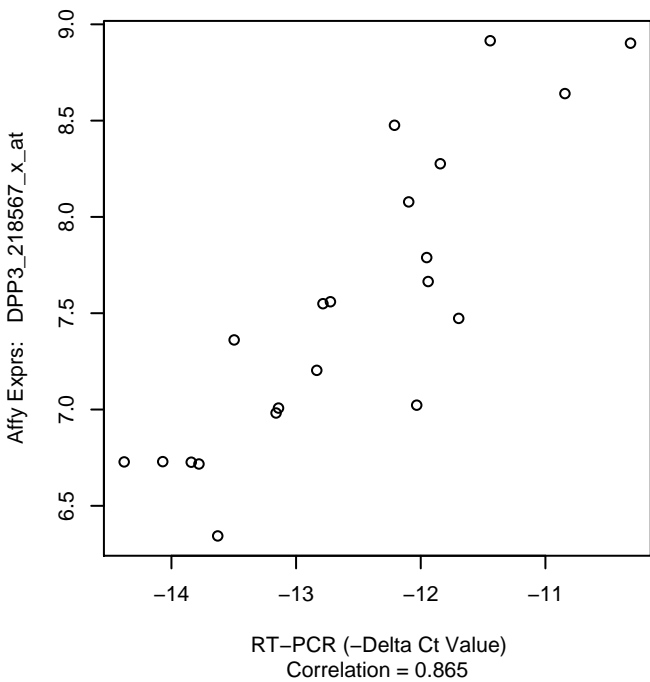
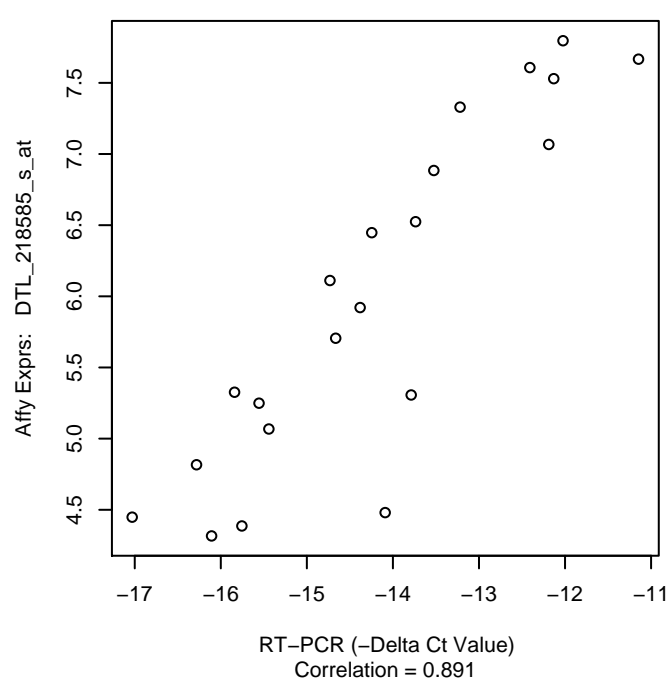
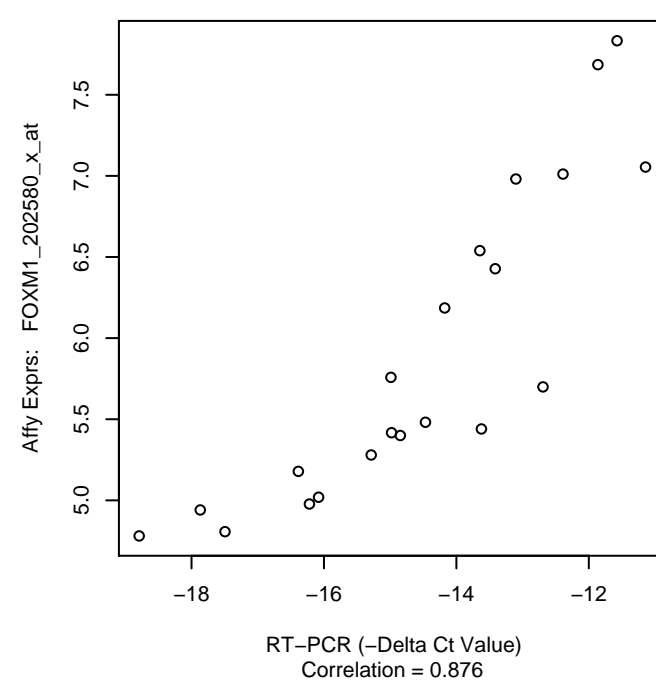


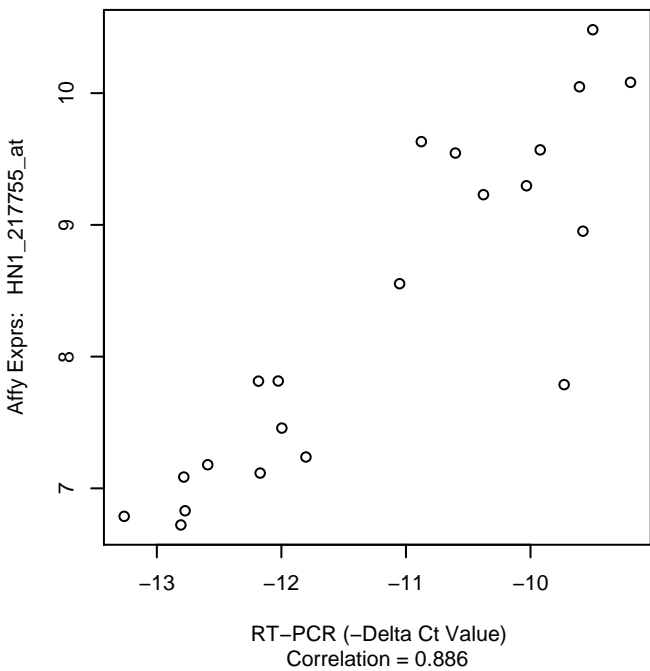
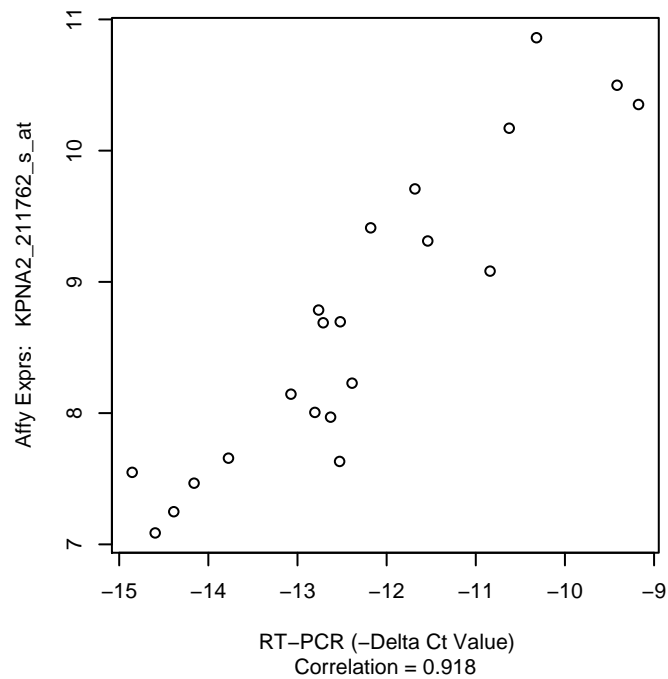
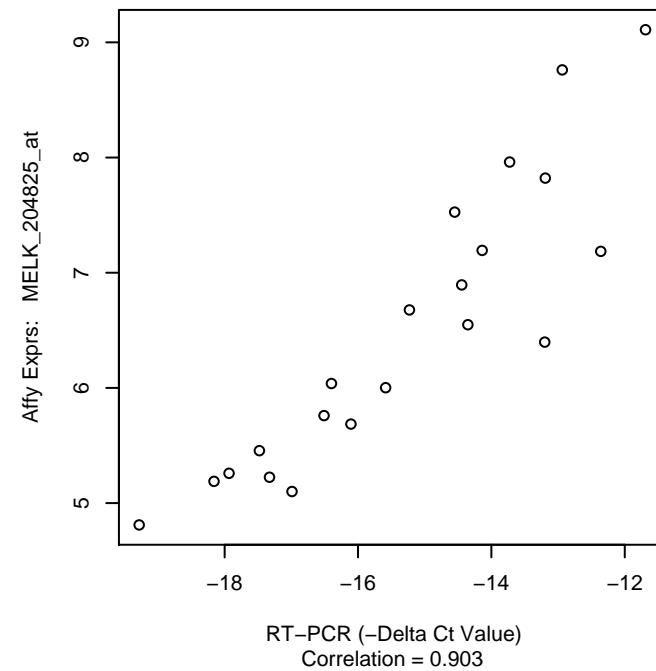
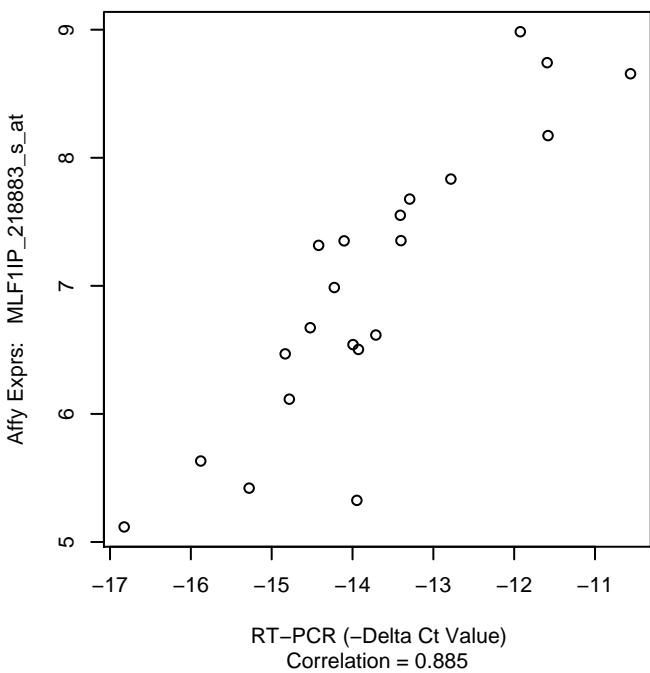
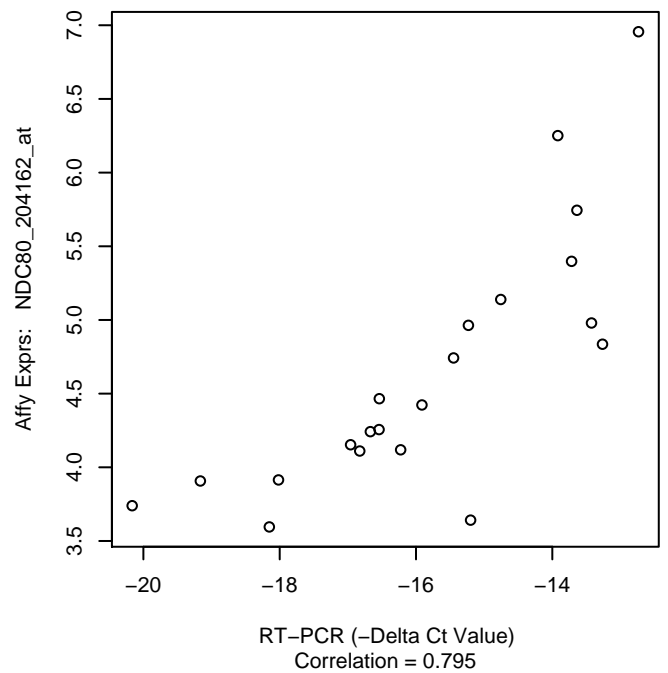
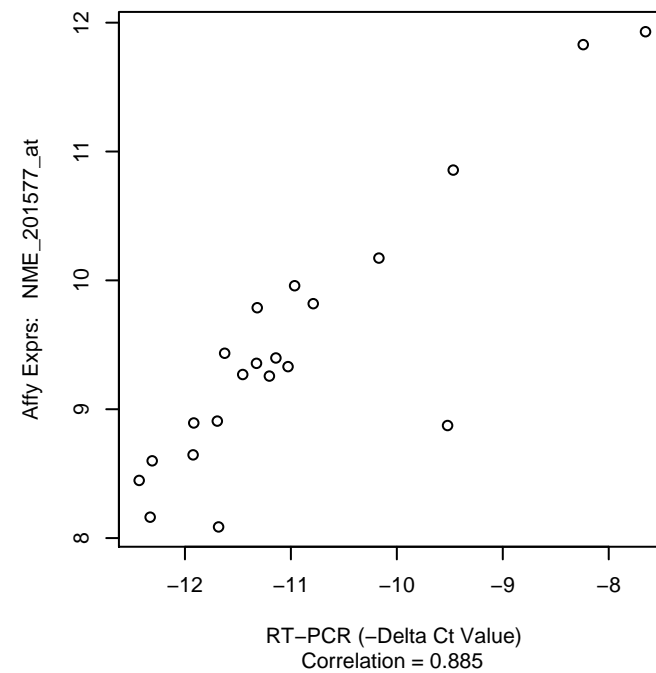
Figure B: Comparison of Various Principle Components (PCA): PCA 1 means the malignancy-risk score was generated by the 1st PCA. We did the same way for PCA2 and PCA3 in microarray and RT-PCR, respectively. The PCA1 (microarray and RT-PCR) demonstrated a trend of the risk score from normal to IDC tissues in the first column panel. In contrast, the PCA2 and PCA3 did not show its association with cancer status.

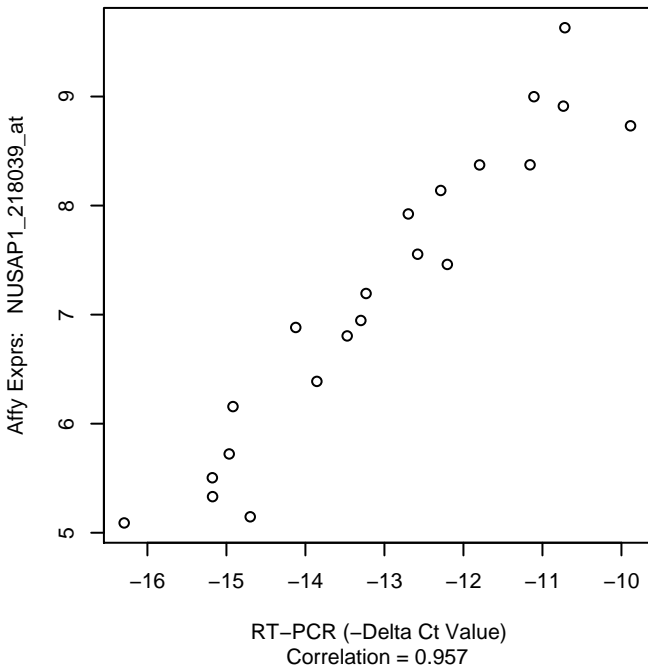
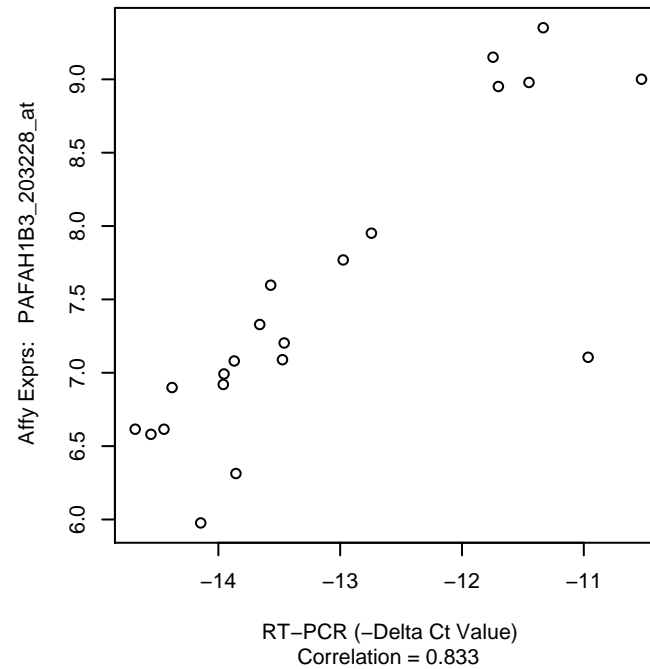
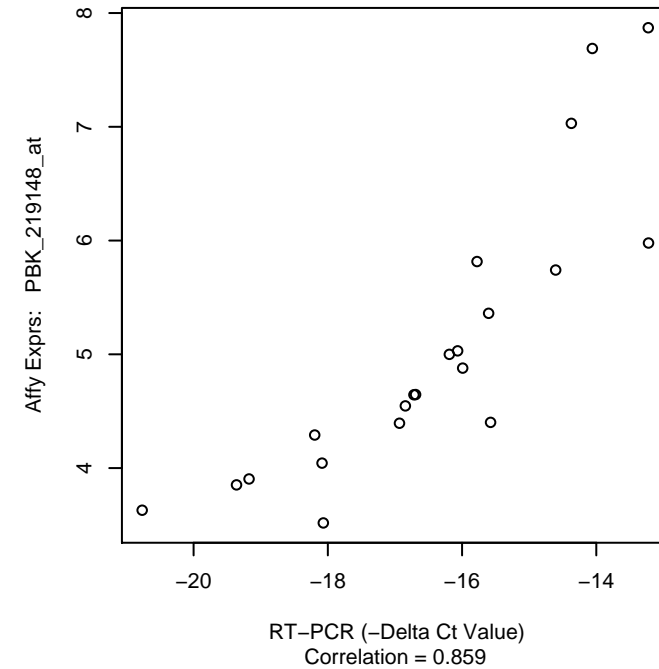
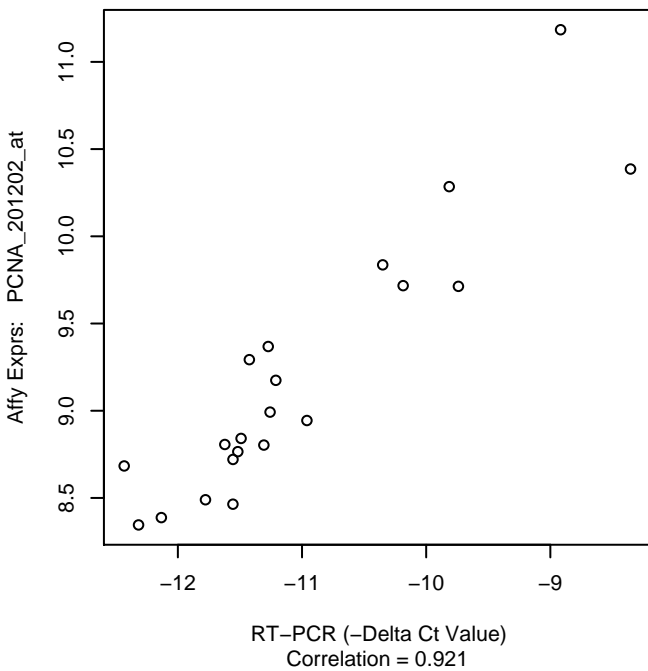
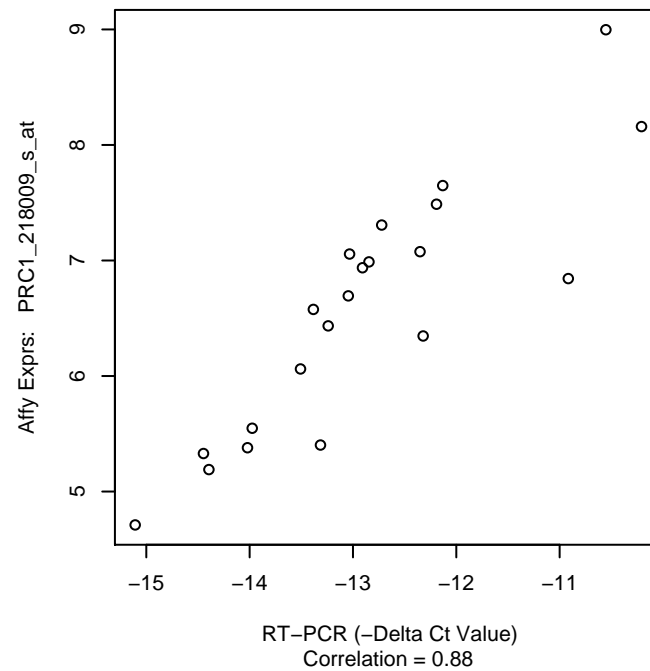
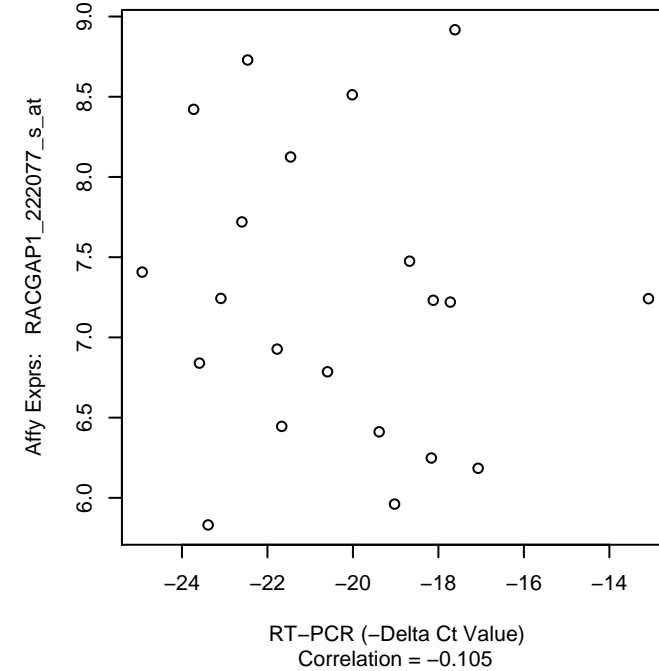


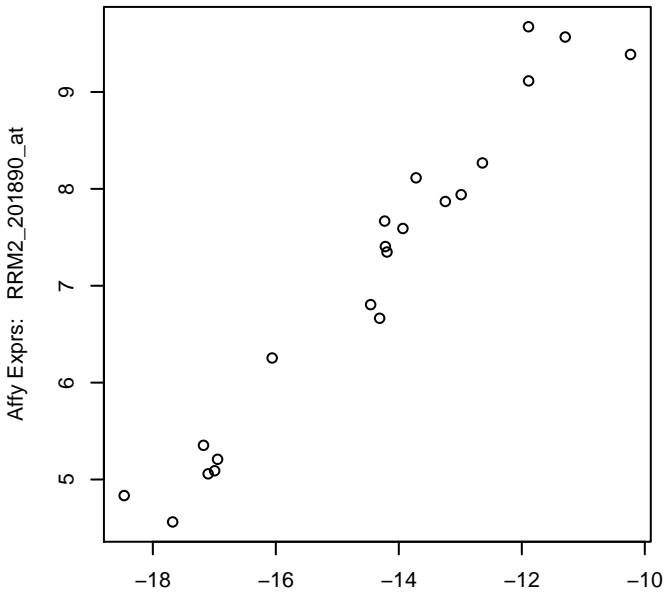
3. Correlation plot:

ANLN**BUB1****BUB3****CCNB1****CDC2****CDKN3**

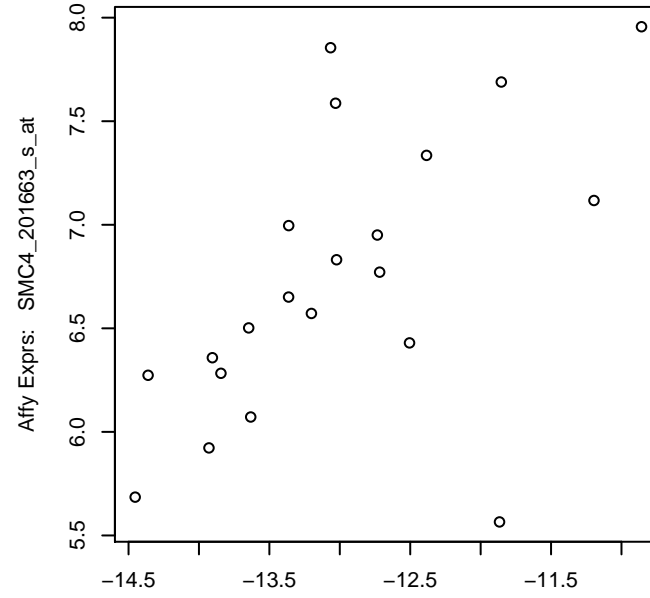
CENPA**CENPF****CKS2****DPP3****DTL****FOXM1**

HN1**KPNA2****MELK****MLF1IP****NDC80****NME**

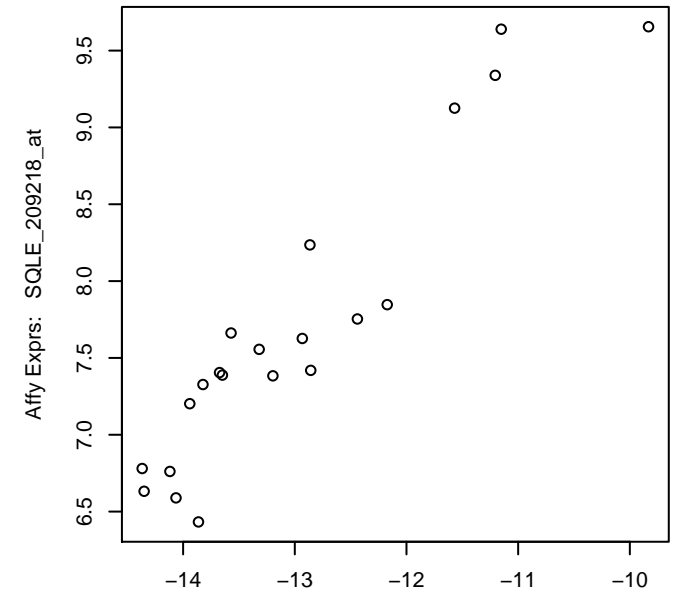
NUSAP1**PAFAH1B3****PBK****PCNA****PRC1****RACGAP1**

RRM2

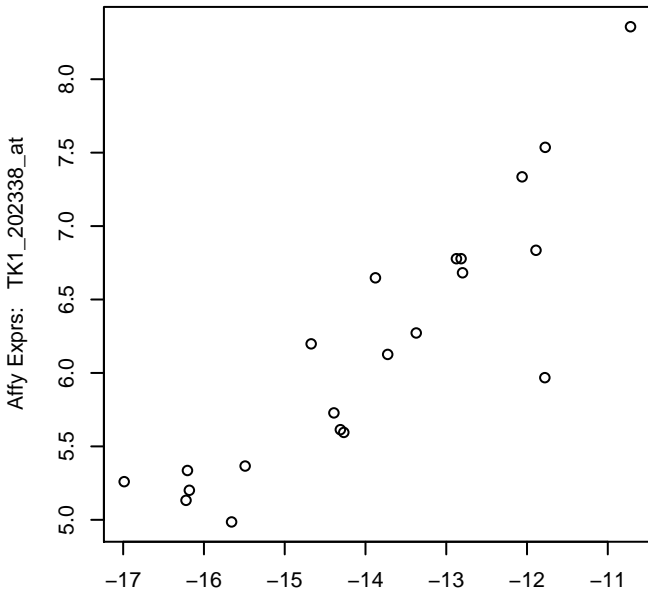
Correlation = 0.975

SMC4

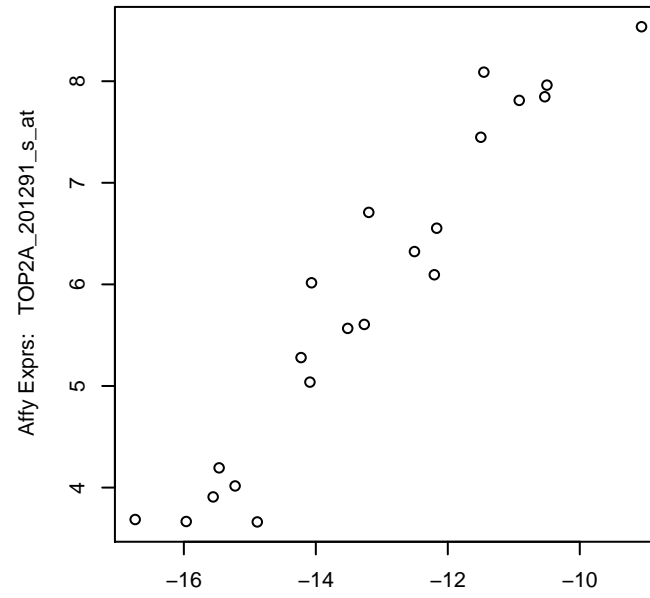
Correlation = 0.562

SQLE

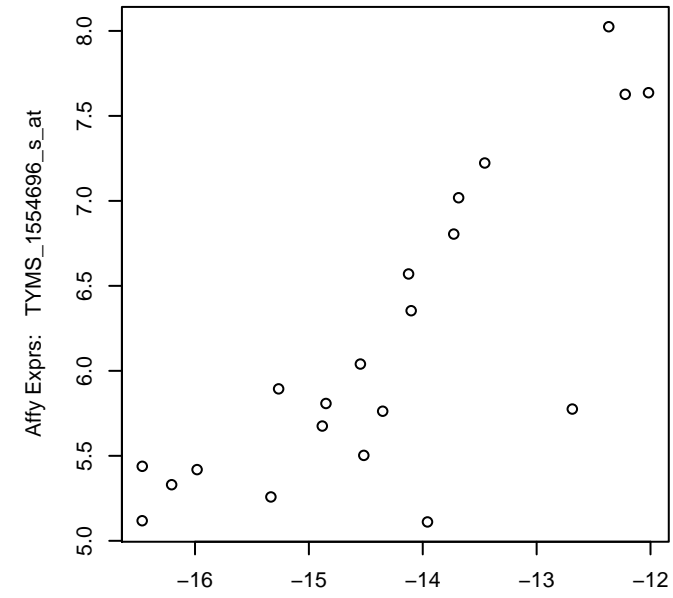
Correlation = 0.938

TK1

Correlation = 0.886

TOP2A

Correlation = 0.961

TYMS

Correlation = 0.804