# Comparison by treatment
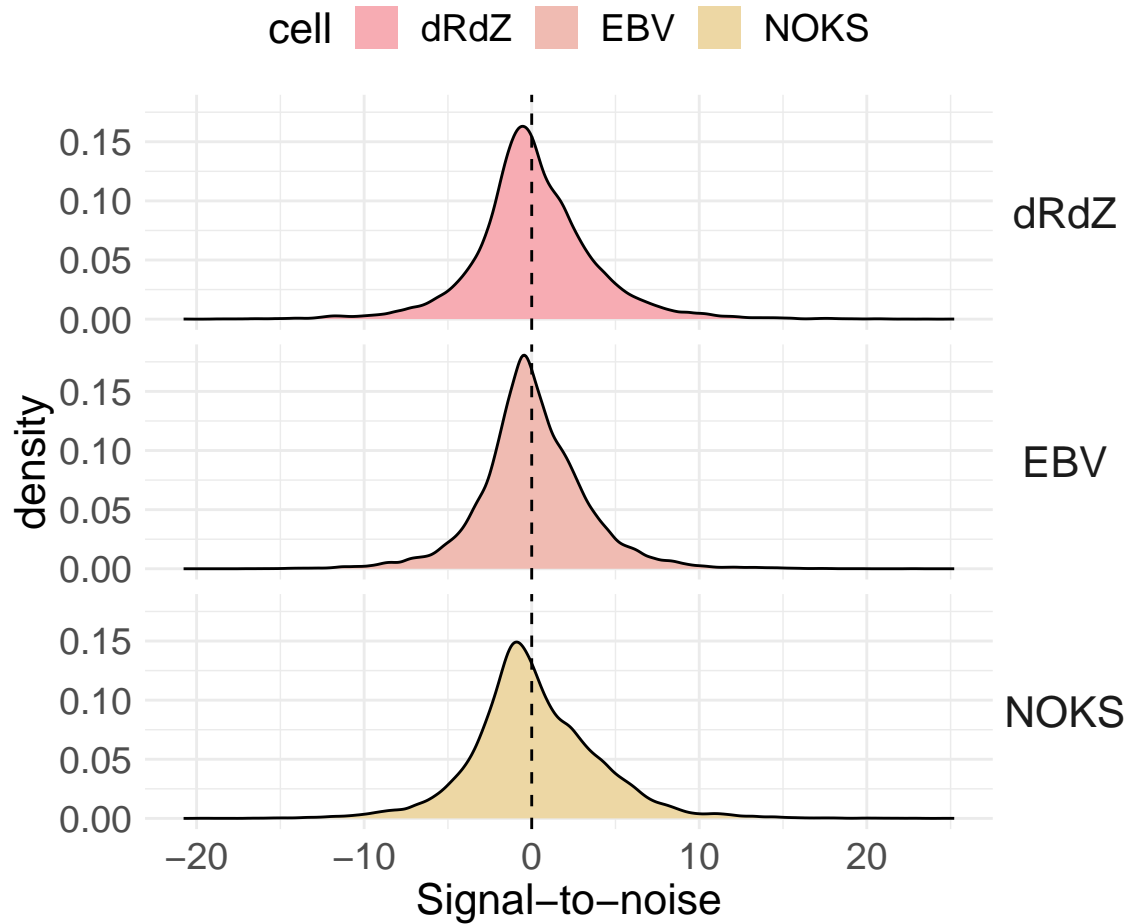
## Contents

## Intro

I guess the idea is to find genes such that show high expression under the treatment, but at the same time exhibit low expression without the treatment, i.e. we are going to focus on the genes that are on the Wald's statistic distribution's tails. For that purpose, we are going to:

1. Perform contrast of MC-treated vs untreated samples for each cell line.

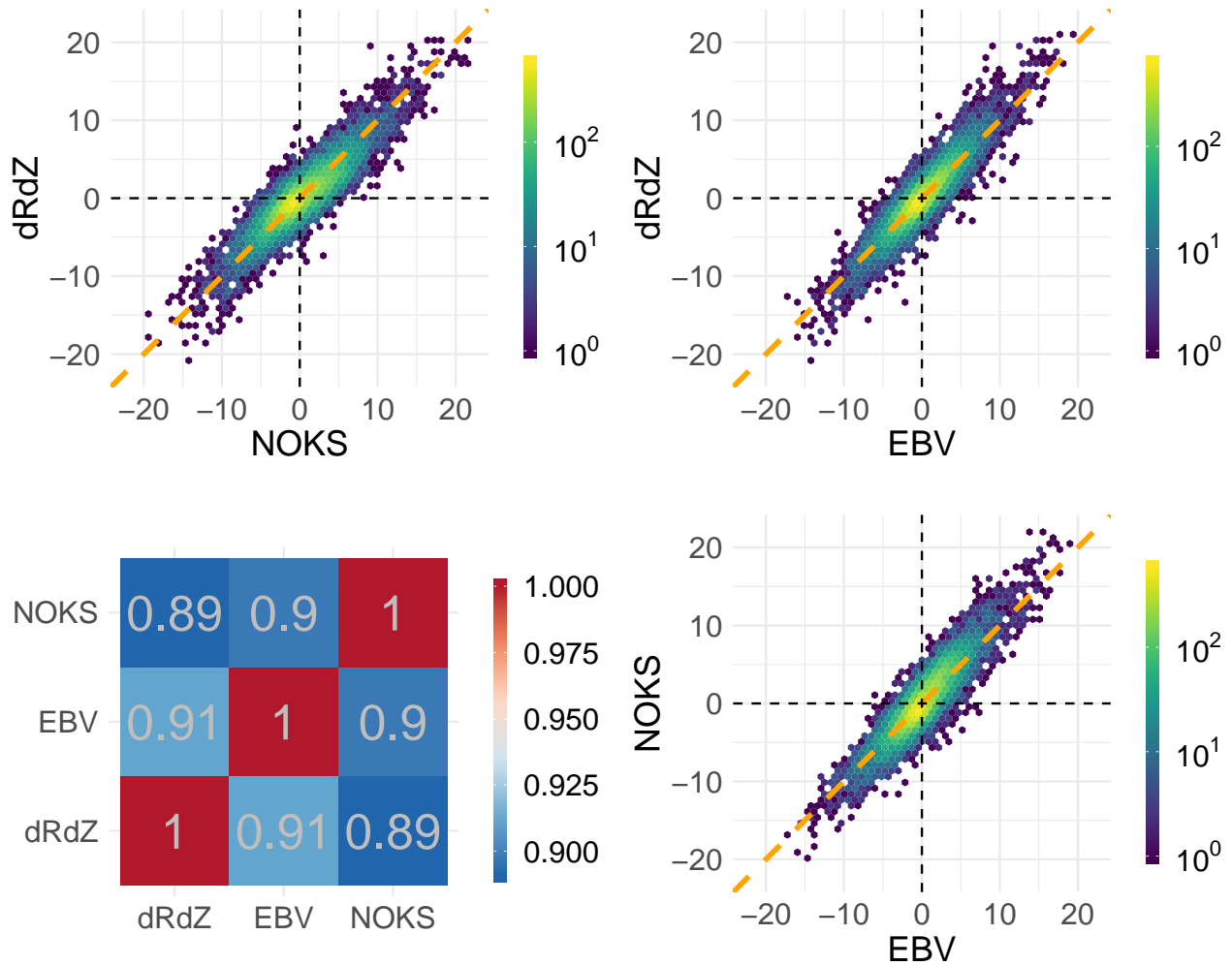2. Compare the Wald's t-statistic between the three cell lines.



## Signal to noise analysis

The figure below shows the densities of the signal-to-noise measures of the treated vs untreated contrast for all three cell lines. Clearly, all three cell lines resemble a similar pattern, with a slightly heavier tails for the NOKS case.

We then compare those three vectors, and we can notice that for most of the genes the signal-to-noise (i.e the log2FoldChange) share the same sign, which means that the MC treatment is affecting the majority of the genes in a similar fashion. This figure shows that the most differentially expressed genes are shared across cell lines, hence it is unlikely to find genes that are differentially expressed by the treatment in one cell line but not the other.

The figure above exhibits that in the case of NOKS vs dRdZ, the signal to noise metrics are relatively symmetric around the diagonal, but in the case of the other two pairs (EBV vs NOKS, and EBV vs dRdZ) we can notice a slightly tilted distribution towars NOKS, and dRdZ respectively. Searching in the bioconductor forums, I found the following answer by the creator of DESeq2, which states a framework to test ratio of ratio (which appears to be a common problem in RIP-seq / CLIP-seq), where the IP counts are normalized by the Input counts (very similar to ChIP-seq).

For that purpose, we fit a different models for each pair of cell lines, which means that we are testing the `cell:treatment` effect, which indicates that we are testing the `treatment` effect across `cell` lines.

```r
ratio_of_ratios_deseq <- function(rsem_data,thr = 20)
{
  count_matrix =  rsem_data %>%
    dplyr::select(file,rsem) %>%
    unnest() %>%
    dplyr::select(file,gene_id,expected_count) %>%
    mutate(
      expected_count = floor(expected_count)
      ) %>%
    spread(file,expected_count) %>%
    as_matrix()

  coldata = rsem_data %>%
```

```
    dplyr::select(file,cell,treatment) %>%
    mutate(interac = paste(cell,treatment, sep = ".")) %>%
    as.data.frame() %>%
    tibble::remove_rownames() %>%
    tibble::column_to_rownames("file")

  deseq = DESeqDataSetFromMatrix(
    count_matrix,colData = coldata,
    design = ~ cell + treatment + cell:treatment)

  deseq = deseq[ rowSums(assay(deseq) ) > thr,]
  deseq = DESeq(deseq, test = "LRT", reduced =  ~ cell + treatment)

  deseq

}
```
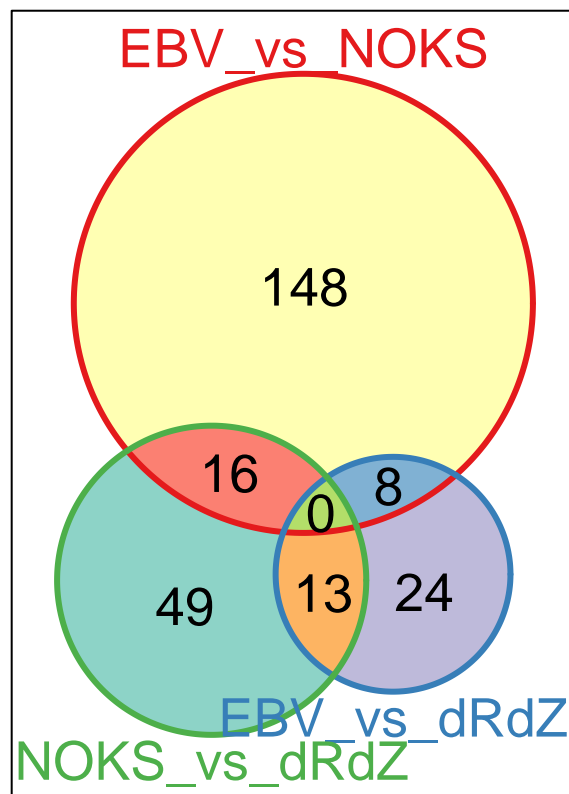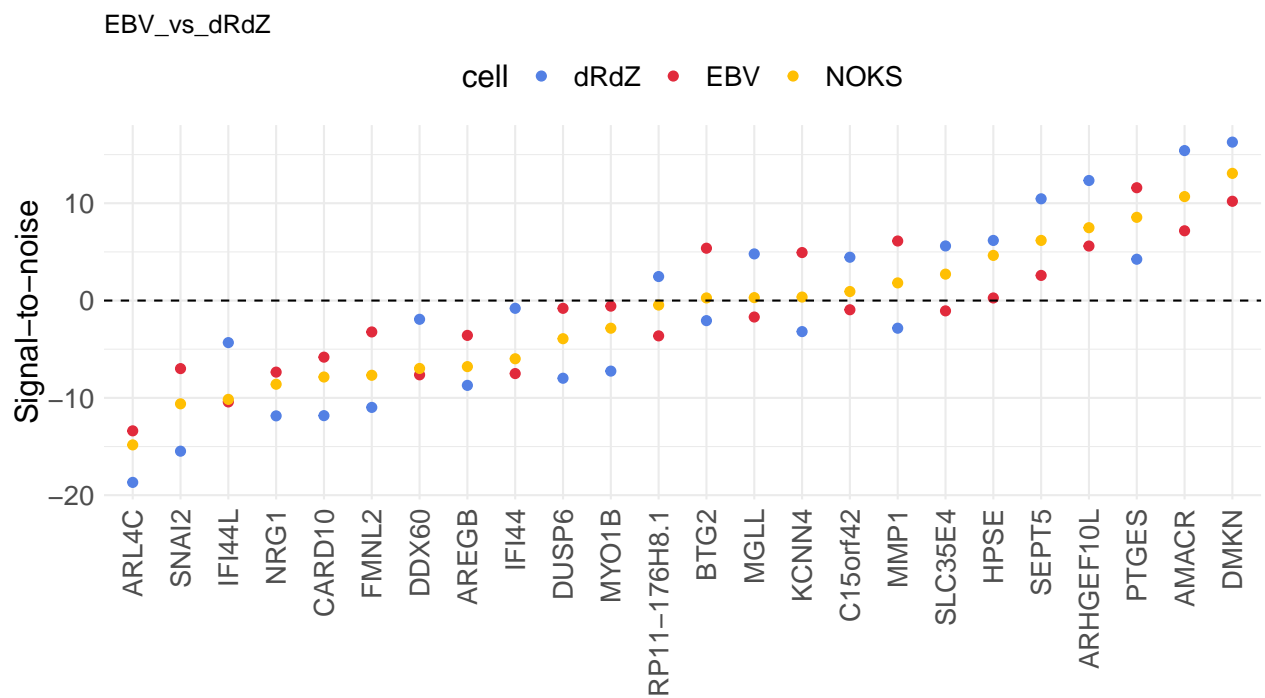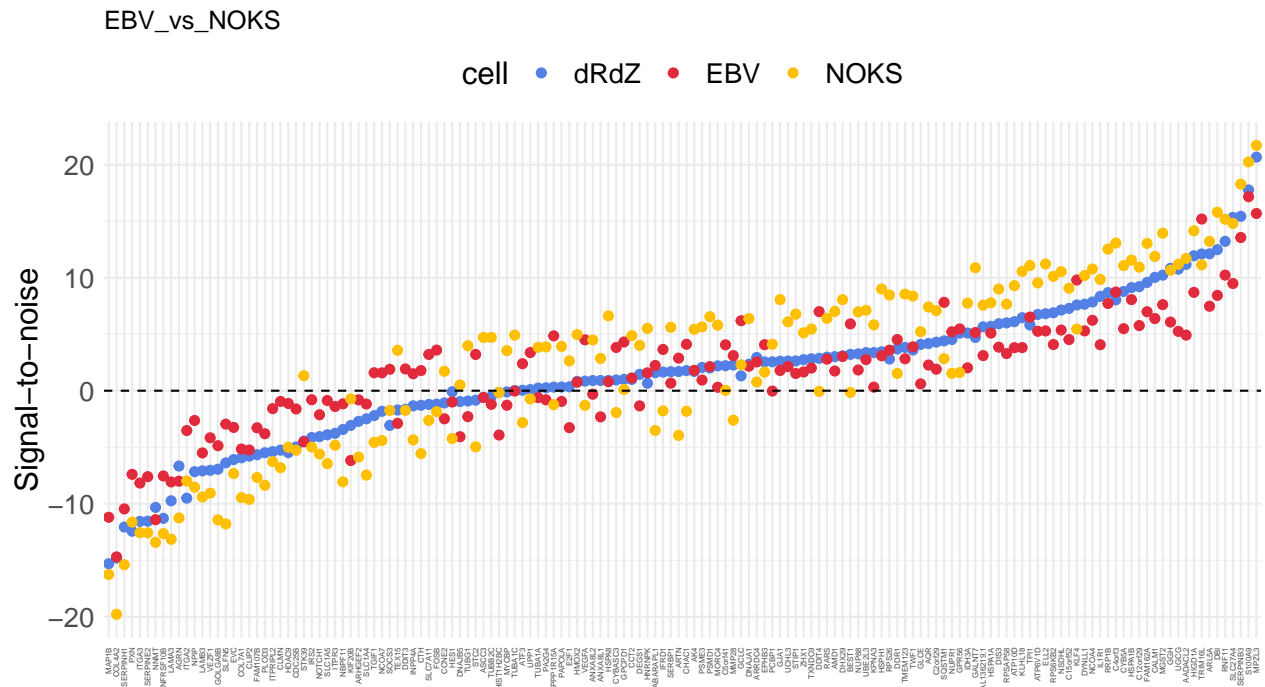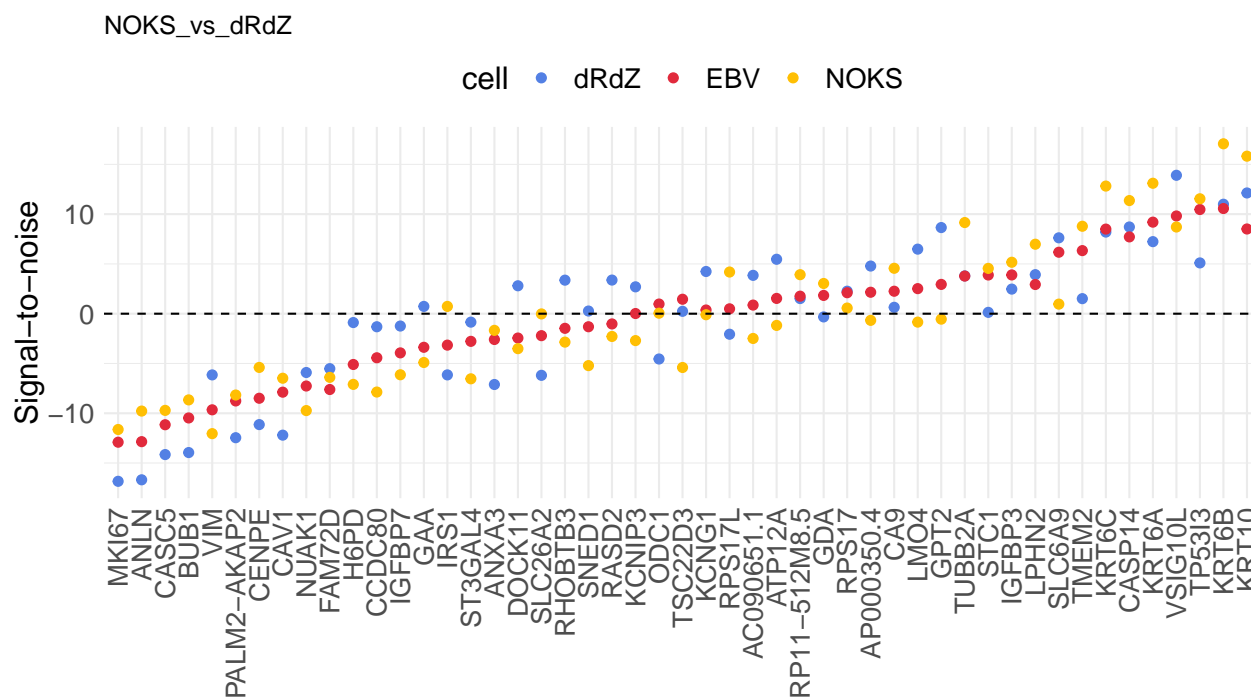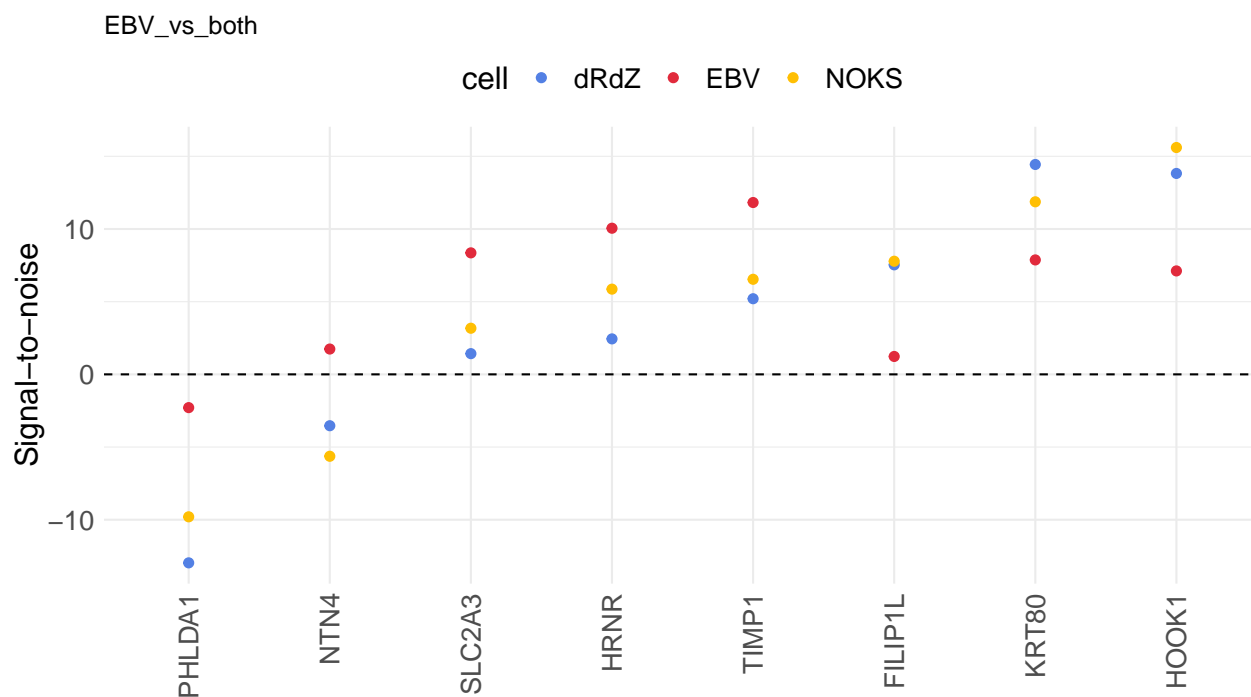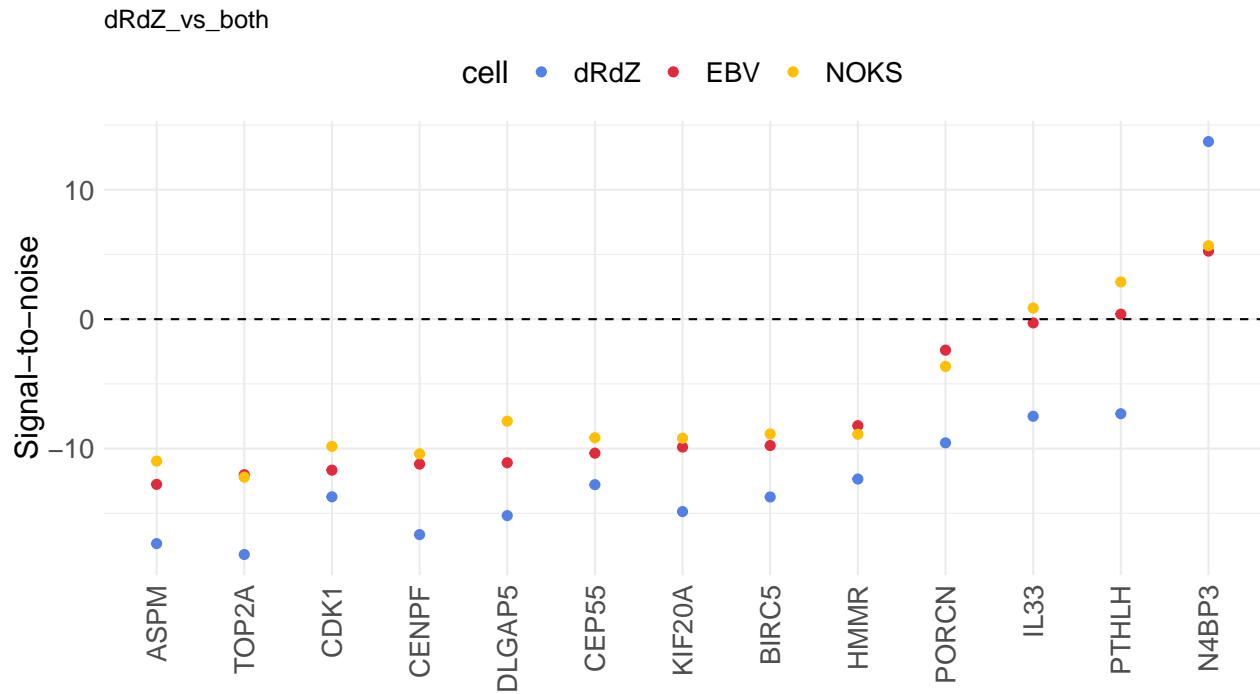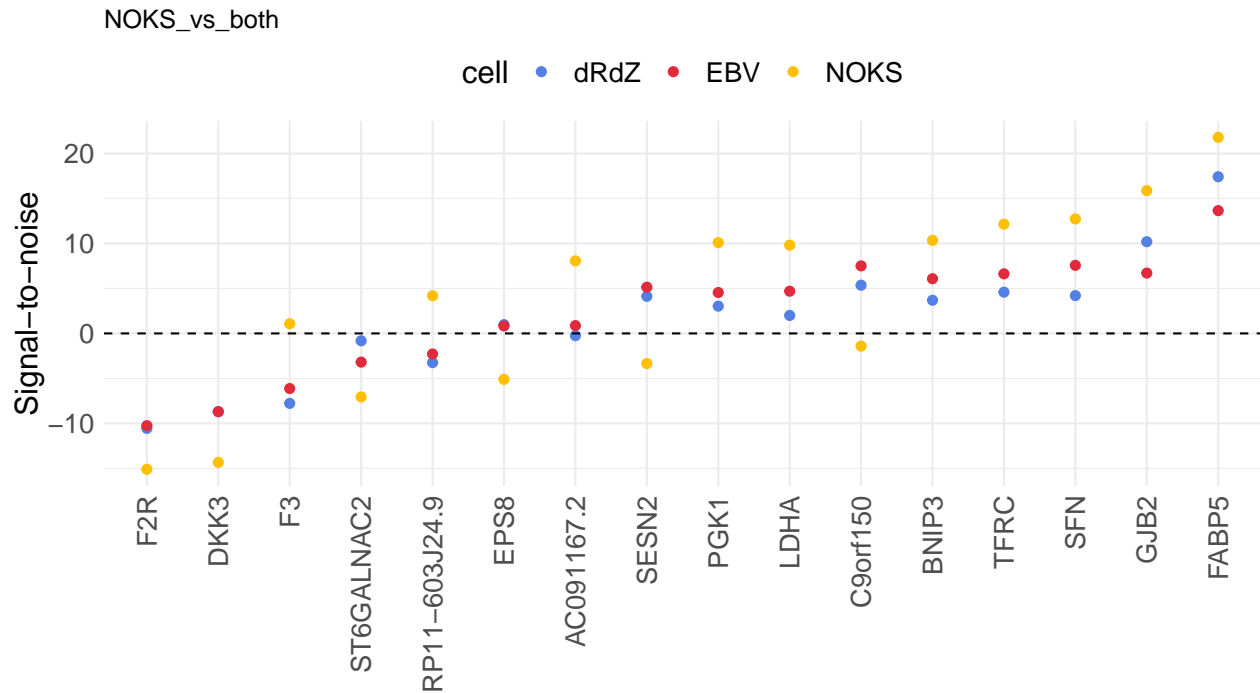
This test returns a smaller list of genes that are differentially expressed. For example, if we consider the genes that are differentially expressed with adjusted p.value ≤ 0.01, we can notice that the number of genes that are differentially expressed is much larger when testing the MC effect across the EBV and NOKS cell lines, than when any of them is the mutant type dRdZ.



Furthermore, we can examine the genes in these subgroups. For example, we can notice that the intersection of `EBV_vs_dRdZ` and `EBV_vs_NOKS` have 8, and in that group there are genes such that the distance between EBV and the other cell lines are maximized, thus it results in genes where the signal-to-noise is close between them (i.e. NOKS and dRdZ). Alternatively, in the regions of the Venn diagram where the genes are only differentiated in one category, we can observe that the cell line that was not considered is usually in the middle.

EBV_vs_NOKS



EBV_vs_dRdZ

EBV_vs_both

NOKS_vs_dRdZ

NOKS_vs_both



dRdZ_vs_both

# Pathway analysis

This analysis returns a different signal-to-noise metric for each cell line. Hence, we are capable of performing a pathway analysis too. For example, we can notice that the none of the `MYC_TARGETS_V1` or `MYC_TARGETS_V2` pathways are enriched in the `EBV_vs_dRdZ` test, but the first one is enriched in the `EBV_vs_NOKS` test and both are enriched in the `NOKS_vs_dRdZ` test.