

B

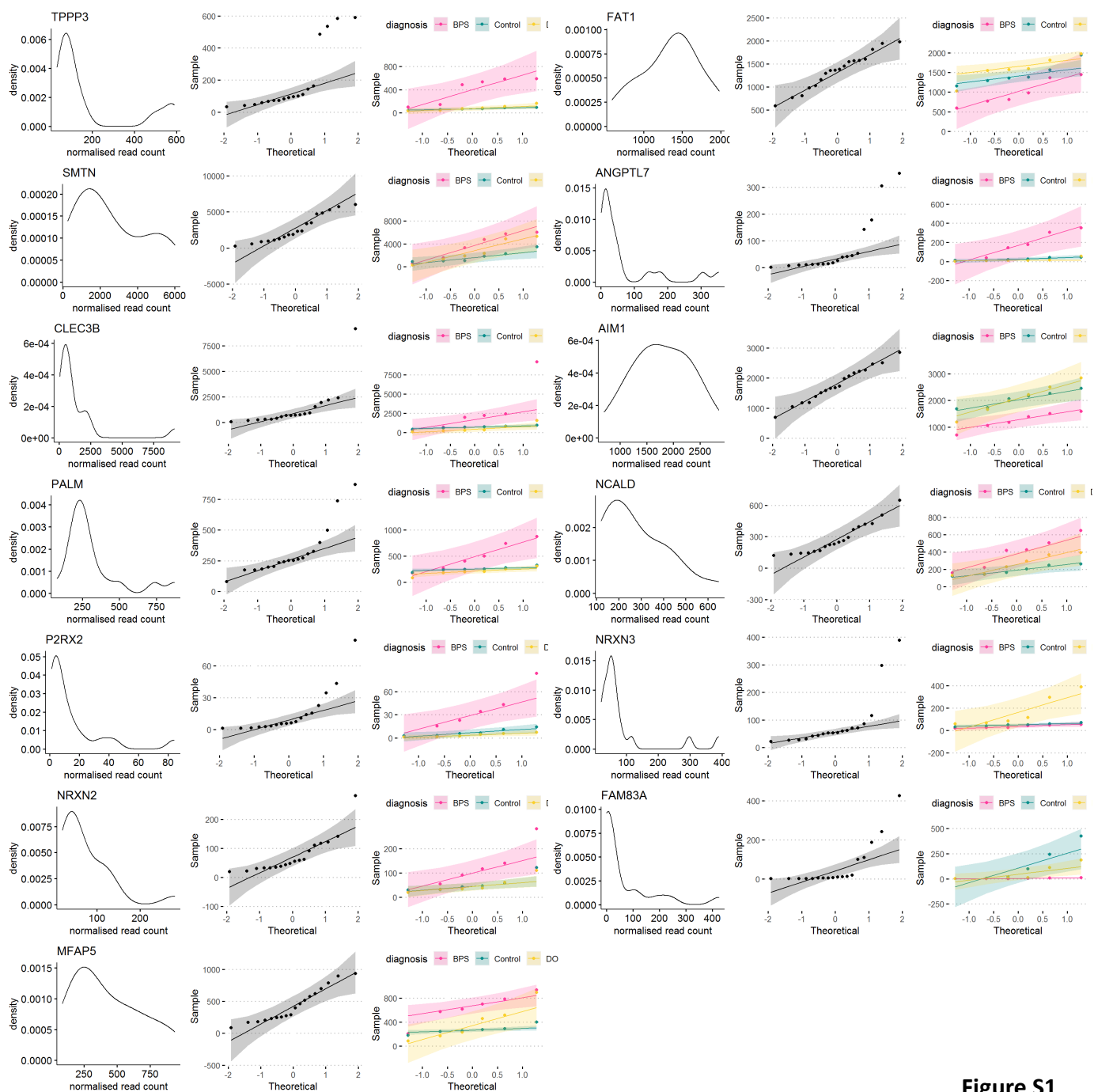


Figure S1

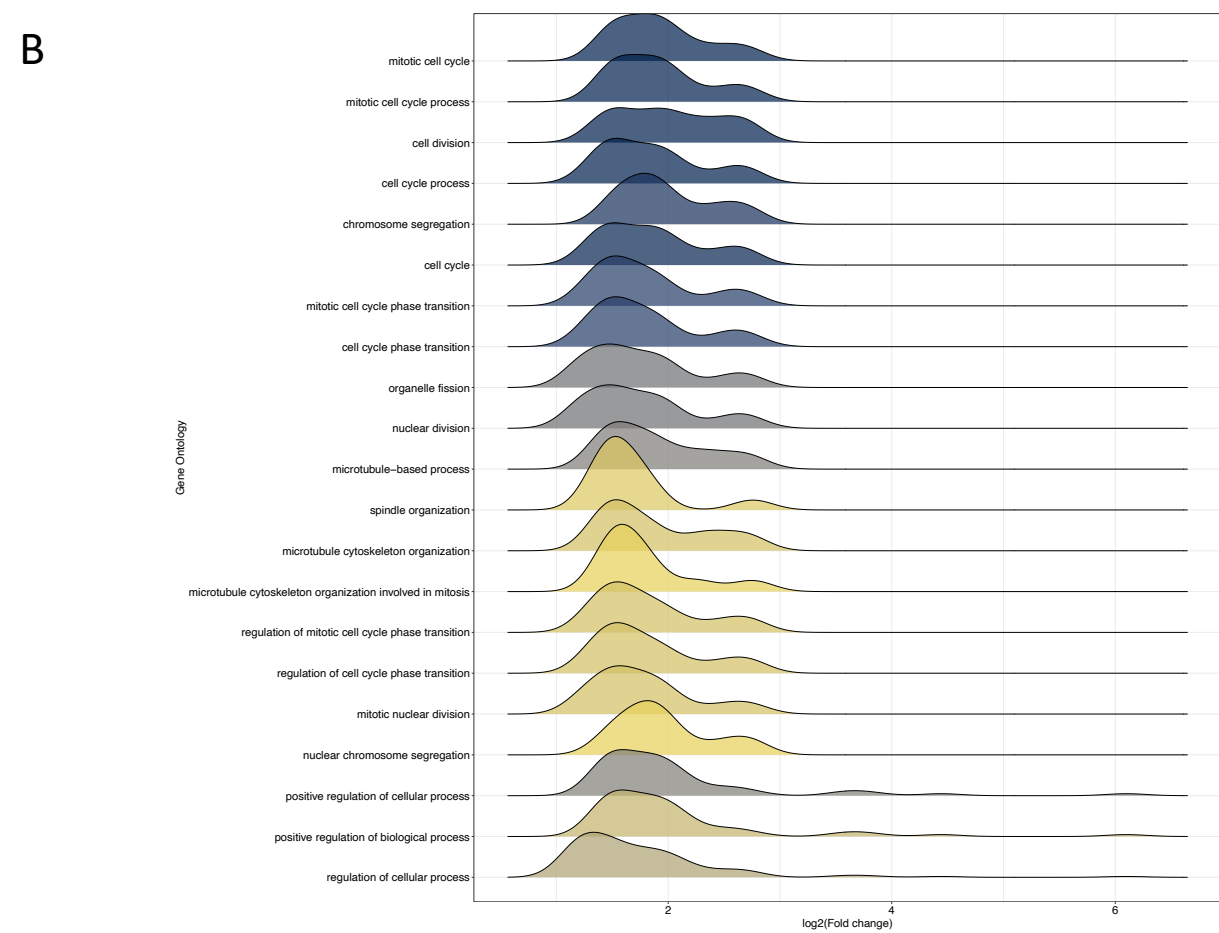
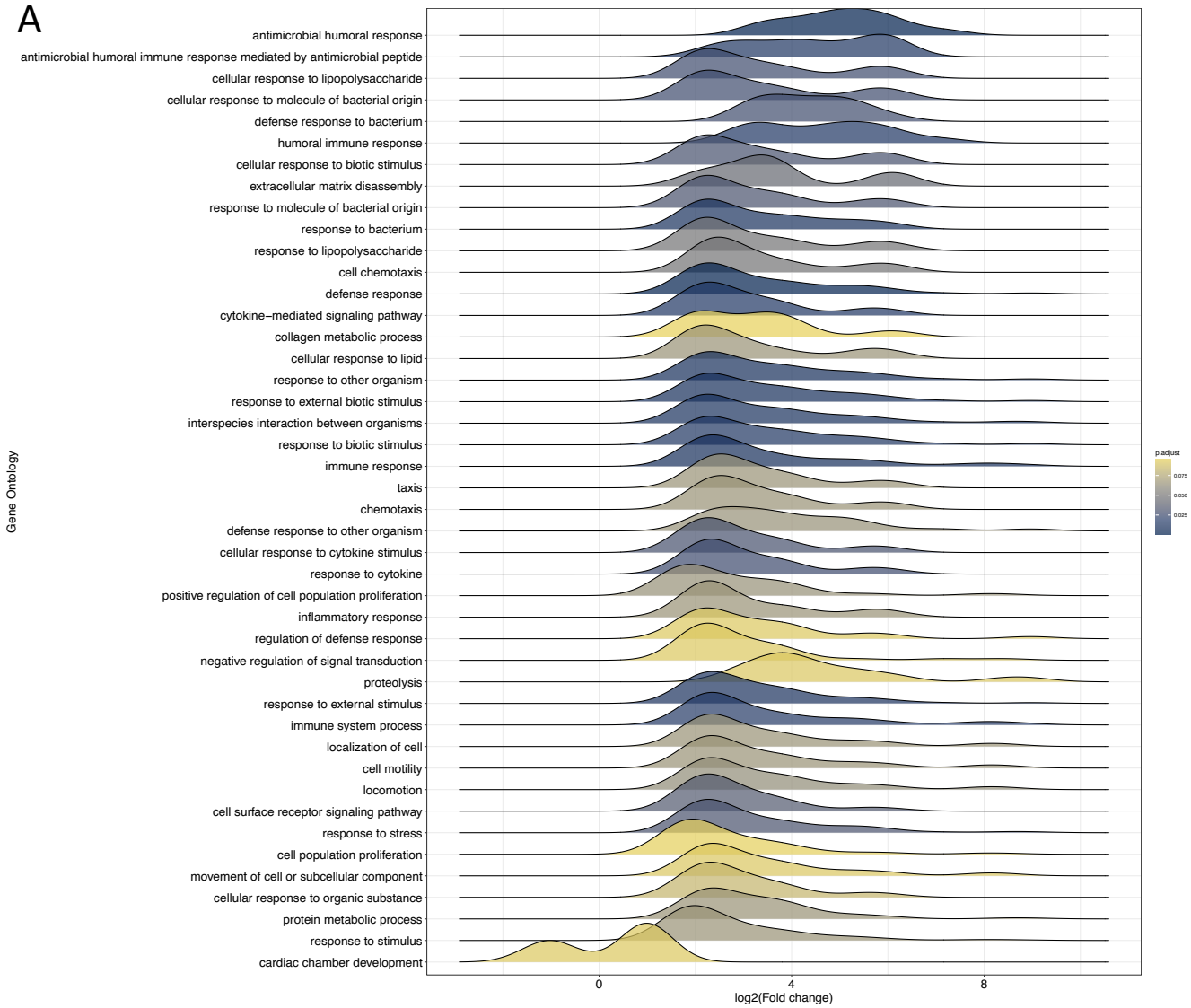
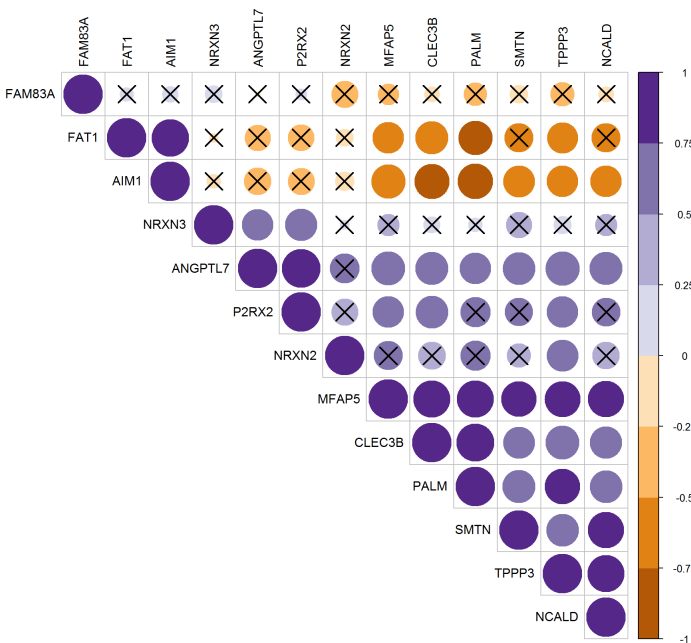


Figure S2



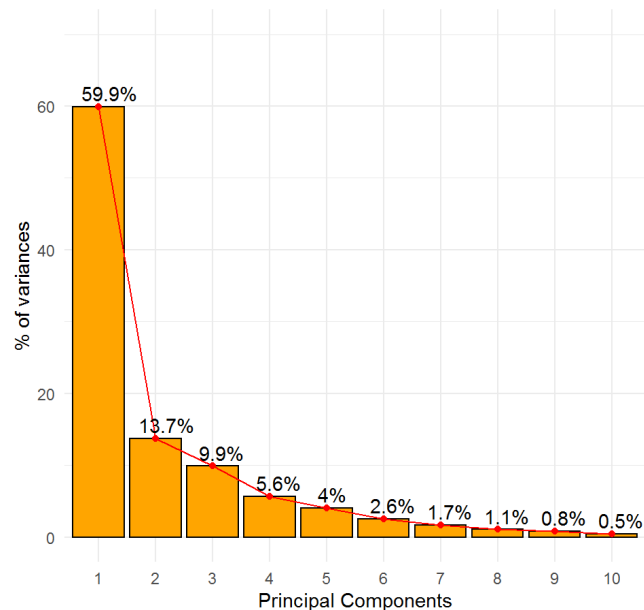
Figure S3

A



B

Variances - PCA



C

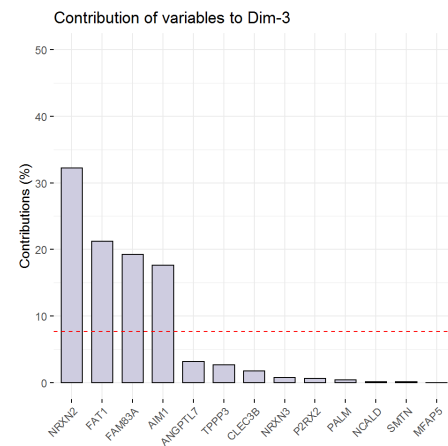
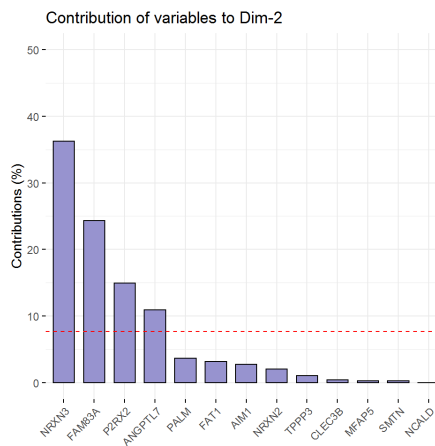
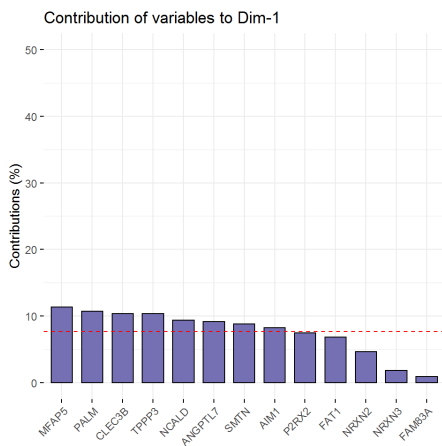


Figure S4

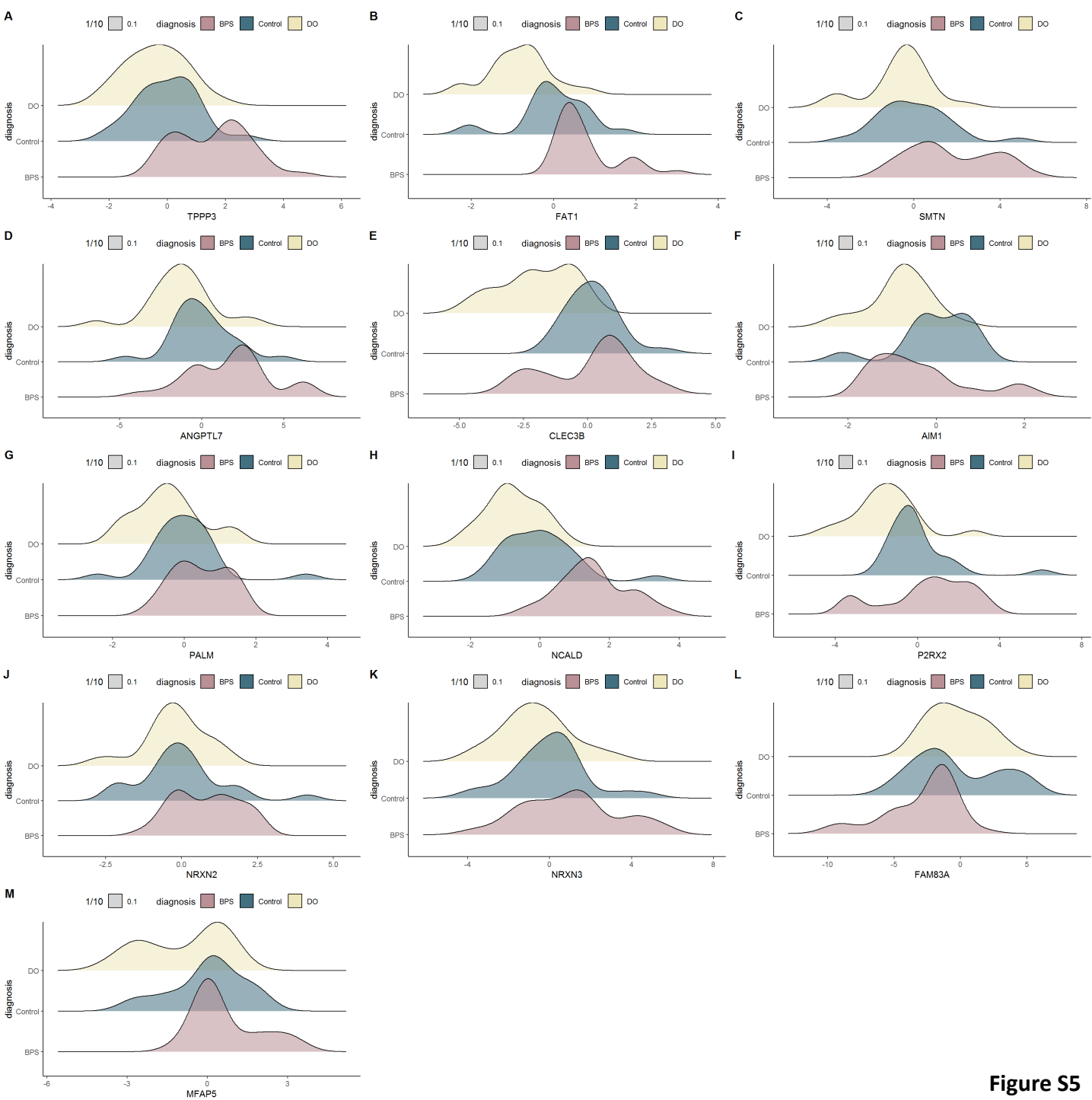


Figure S5

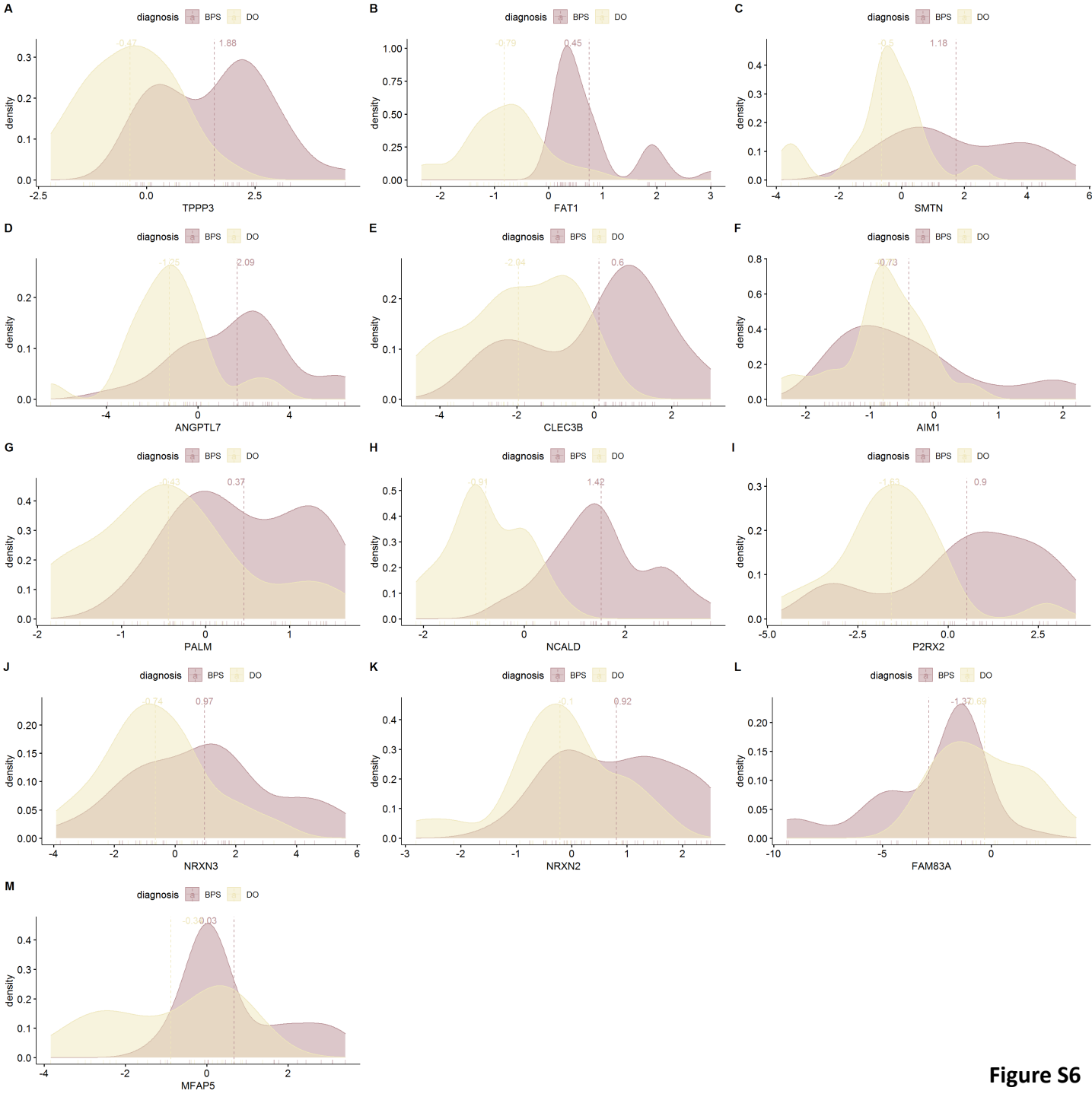


Figure S6

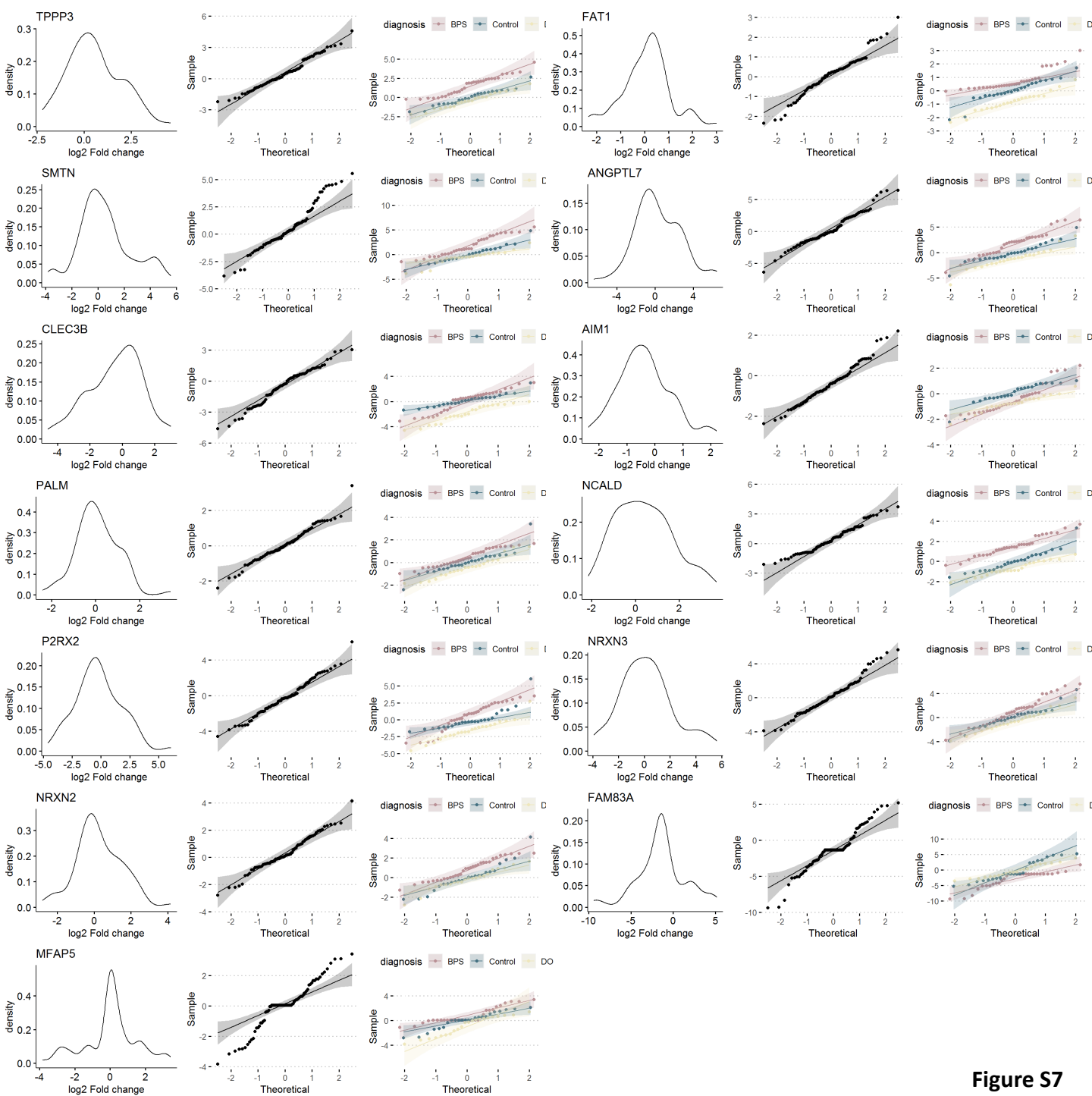


Figure S7

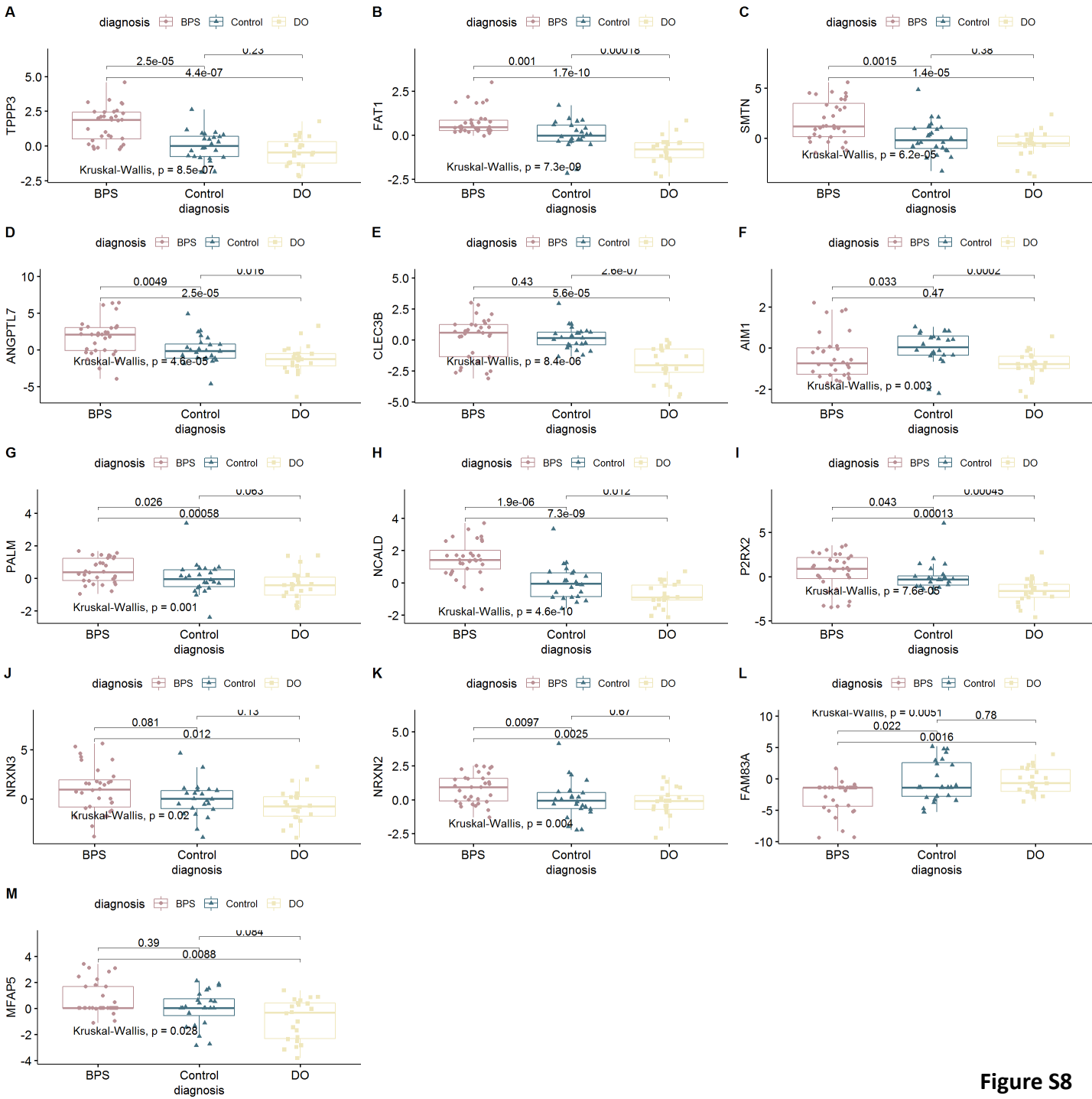


Figure S8

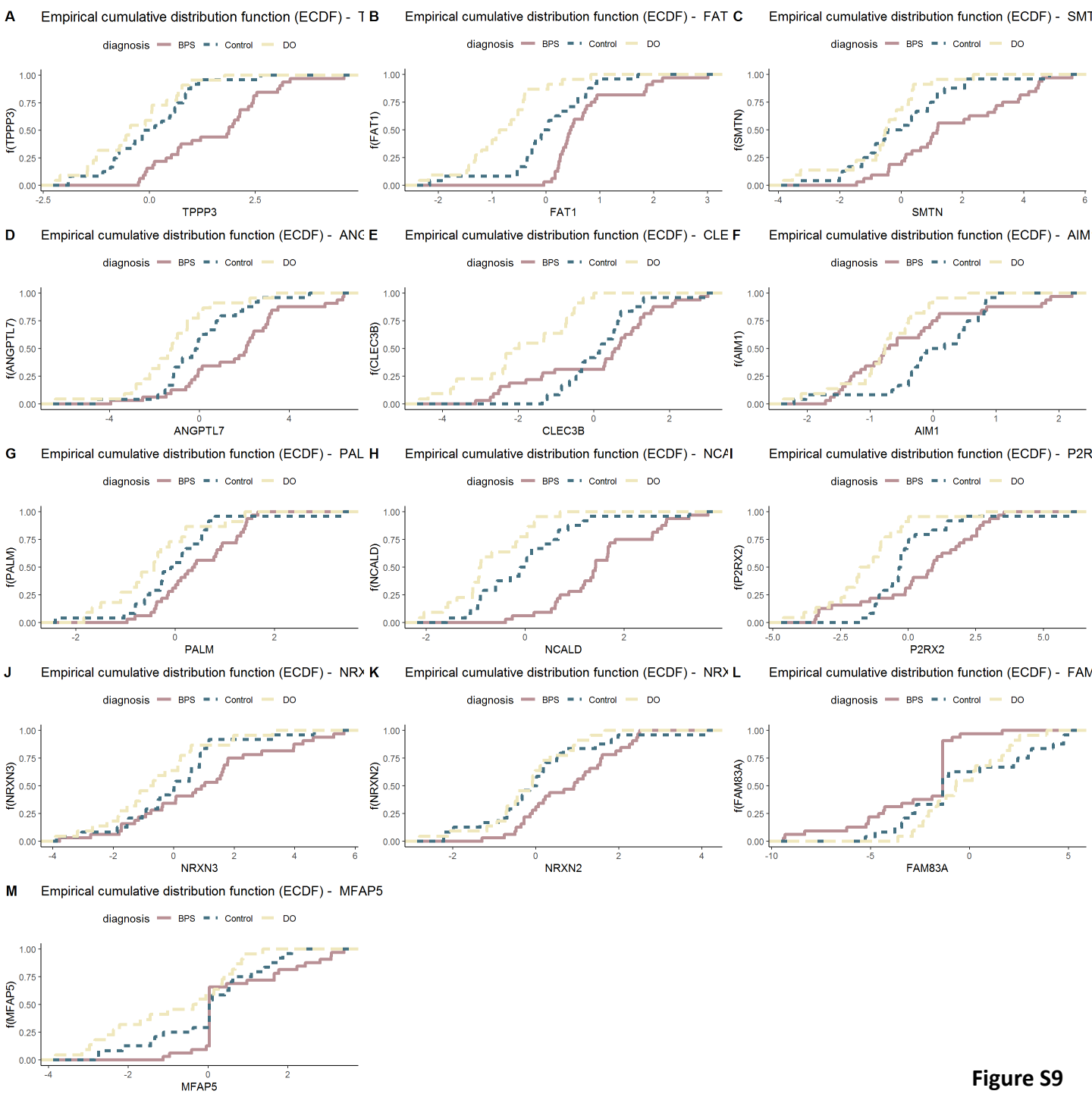
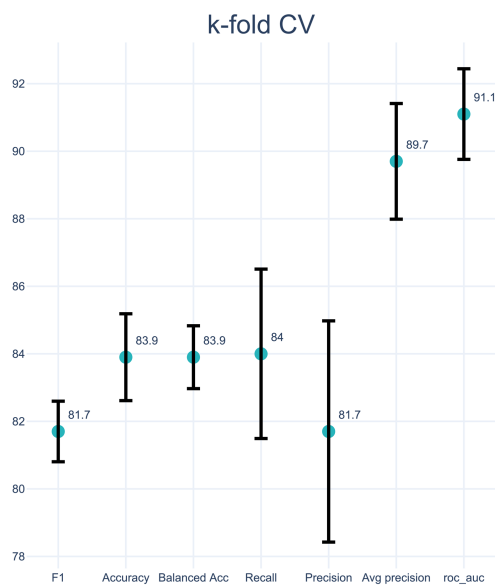




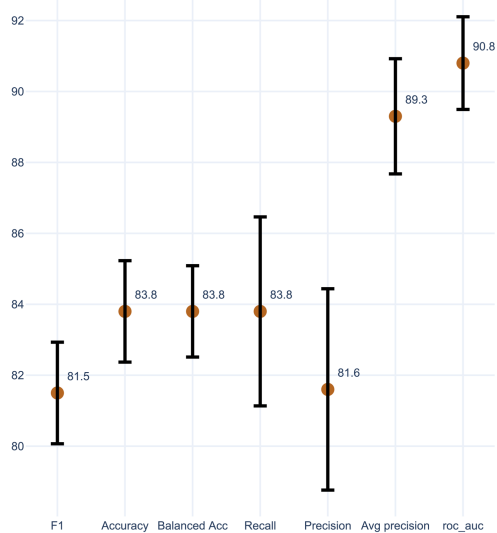
Figure S11

A

Standard Deviation in Models Performance (%)



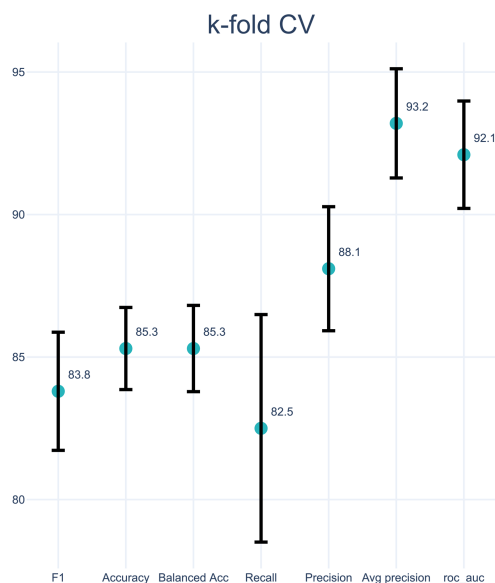
Nested CV



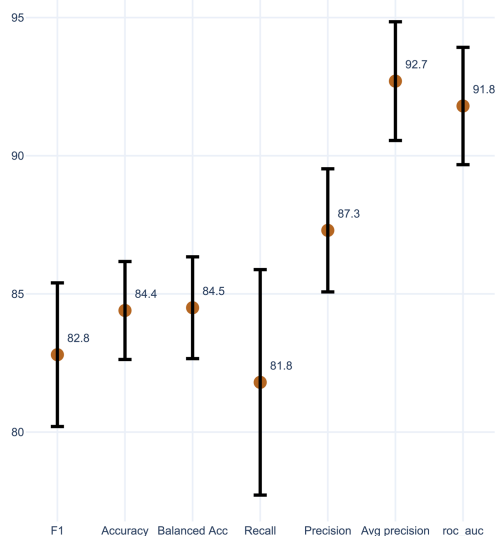
Metrics

B

Standard Deviation in Models Performance (%)



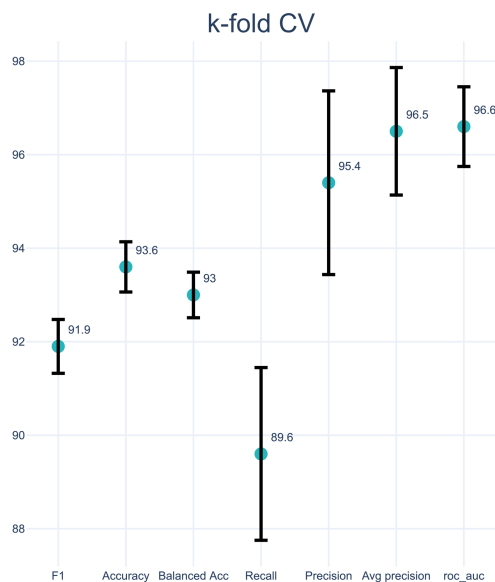
Nested CV



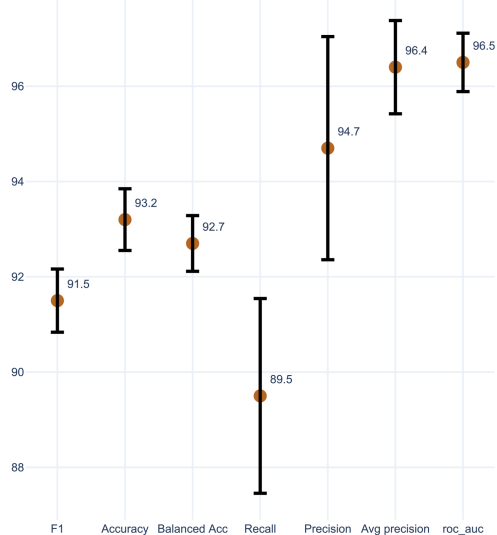
Metrics

C

Standard Deviation in Models Performance (%)



Nested CV



Metrics

Figure S12

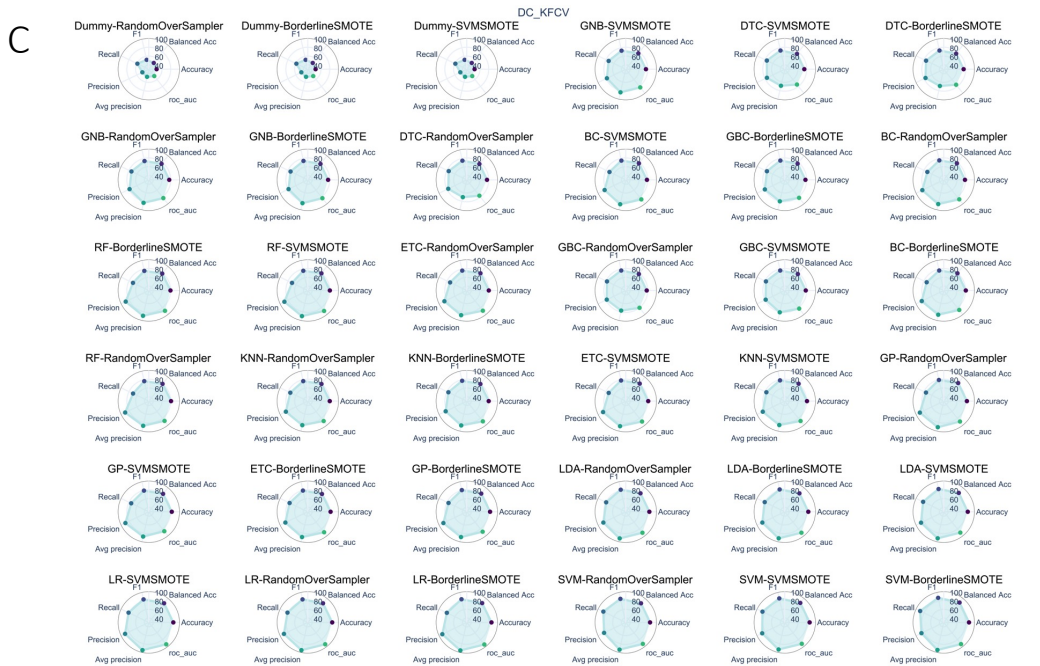
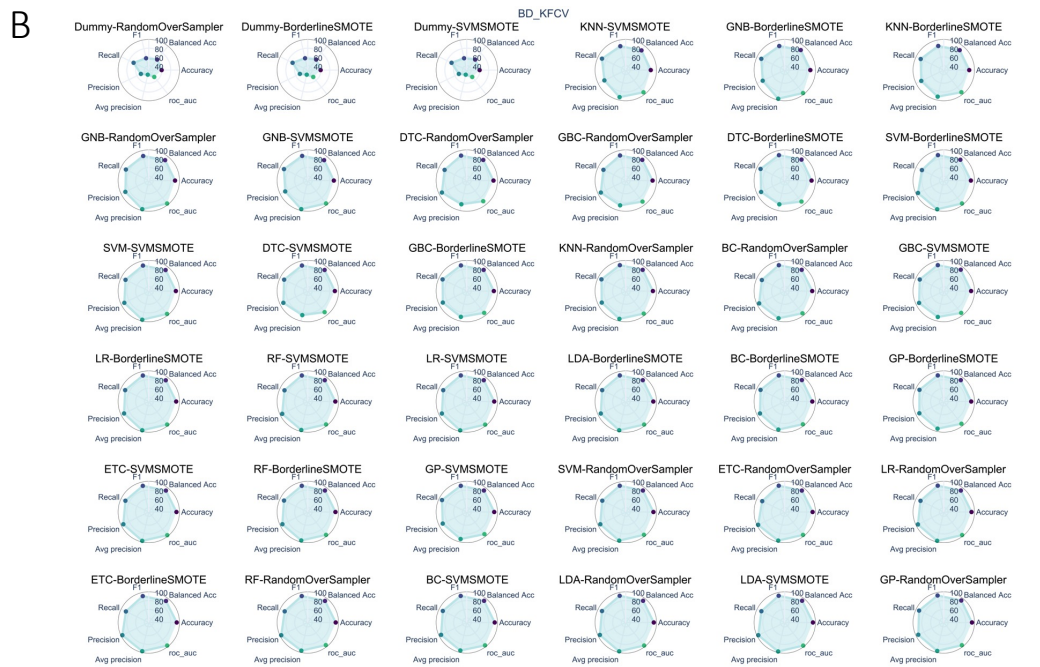
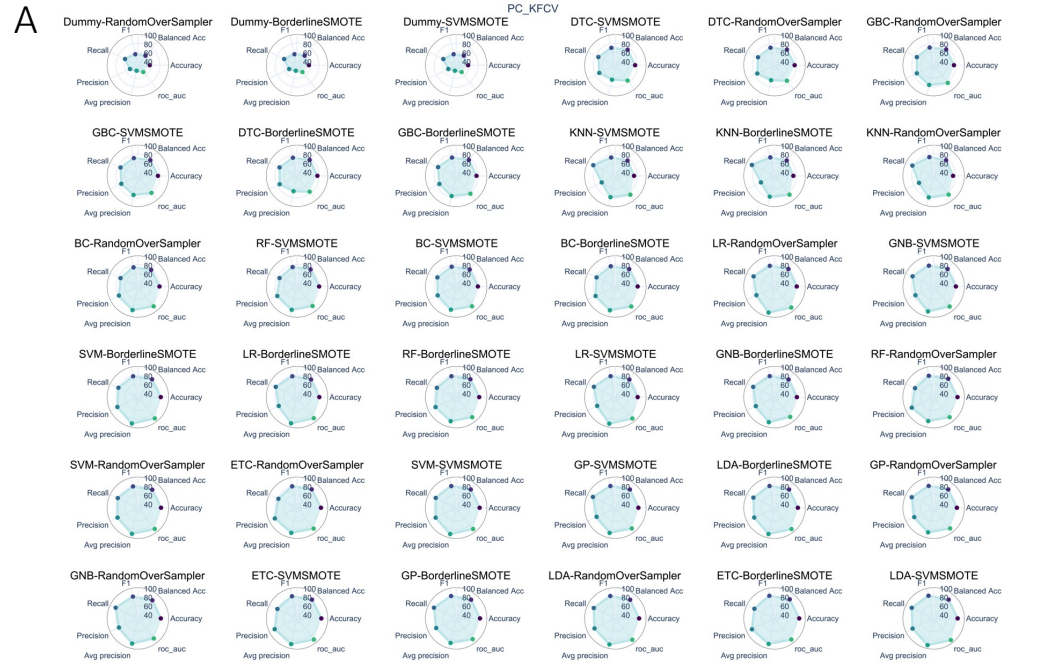
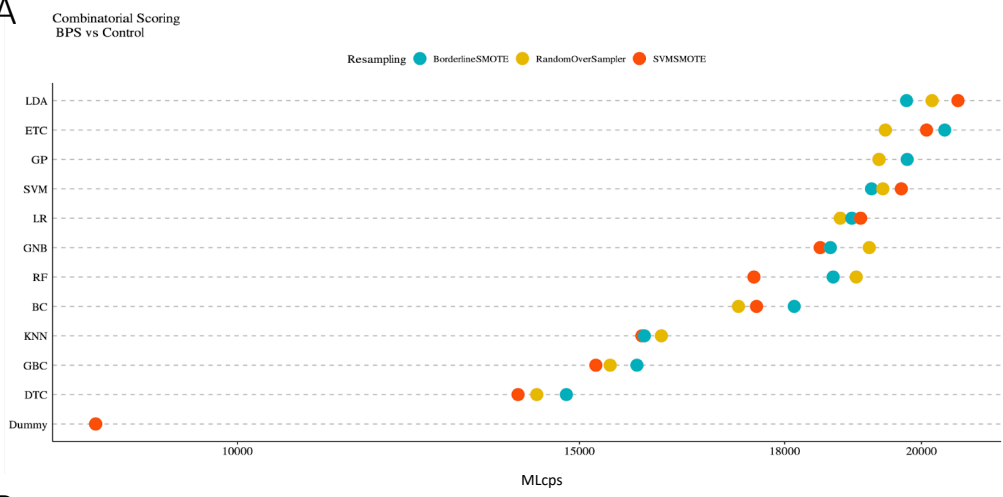


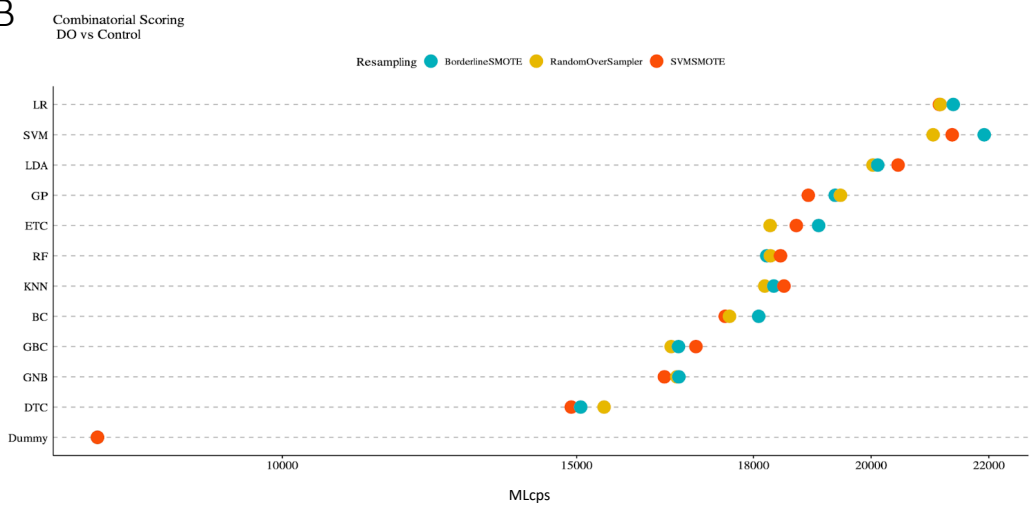
Figure S13



A



B



C

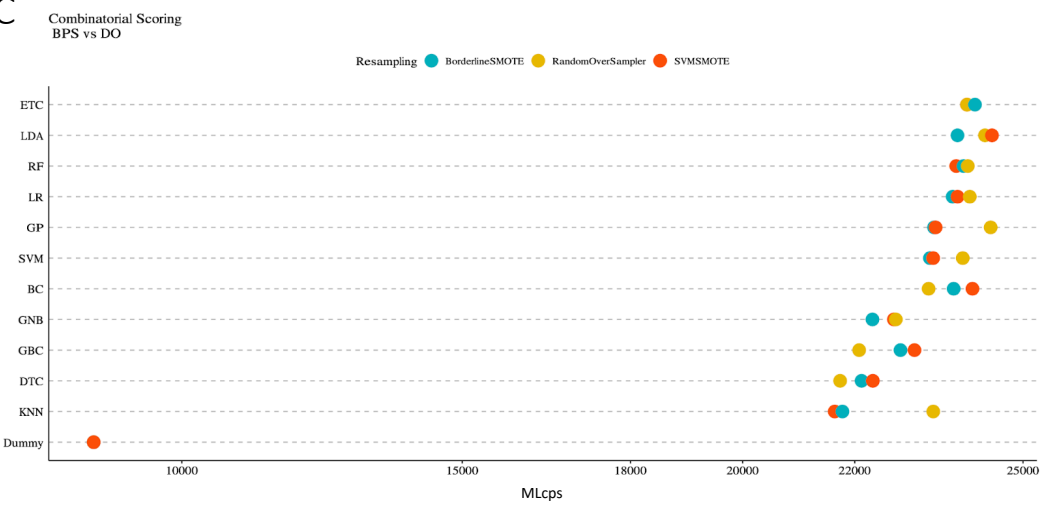


Figure S15

Group	N
Control	24
BPS	32
DO	22

qPCR data

Repeat 1, ..., Repeat 10

Group pairs

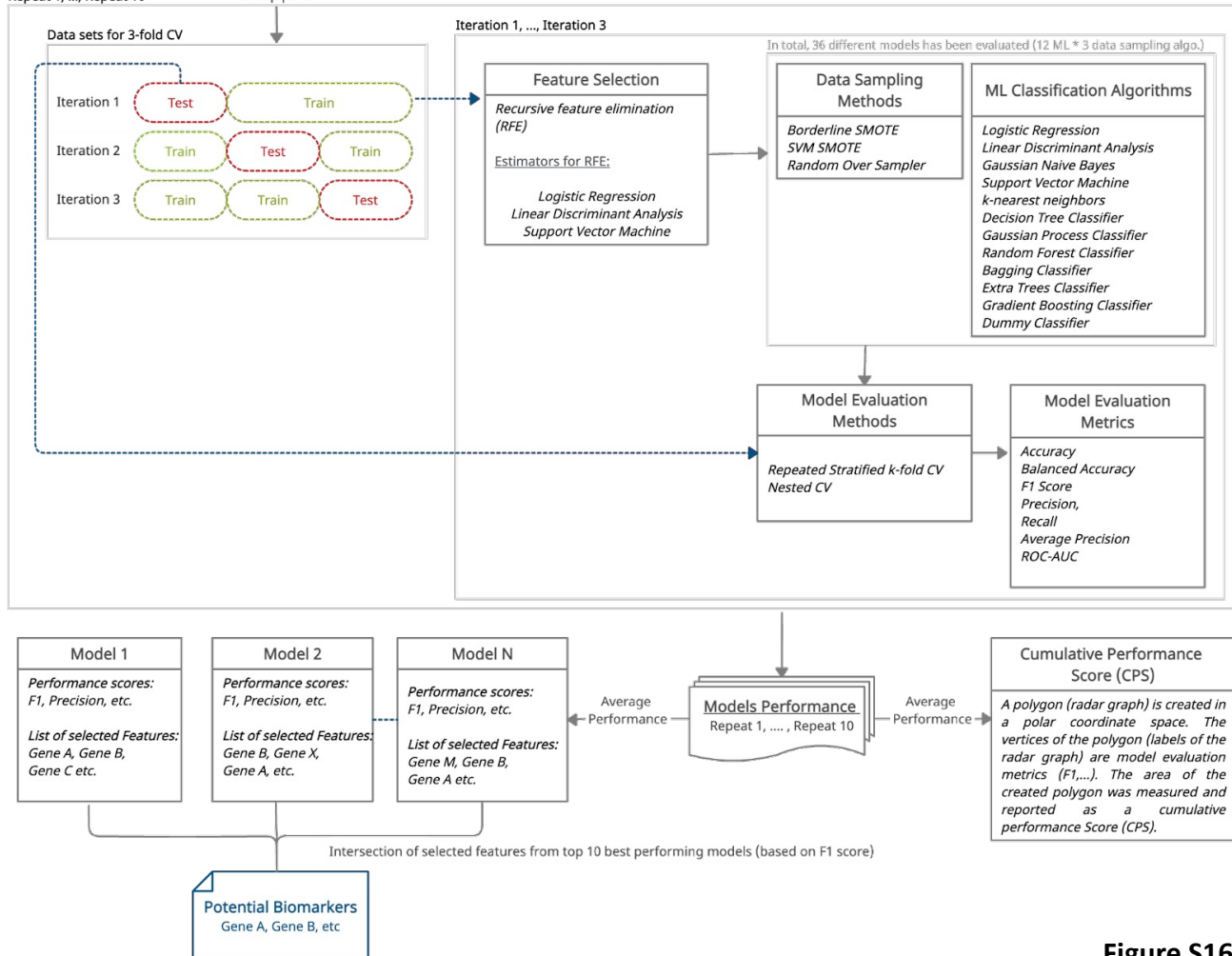


Figure S16

Supplementary Figures - Legends

Figure S1. Similarity matrix and normality test of selected genes in NGS dataset.

(A) Similarity matrix for DO and BPS patients using 1390 differentially regulated mRNAs in DO and BPS patients purple color show lower similarity and green color represents higher similarity. (B) Normality test for 13 selected markers. For each gene, the normal plots (histogram), Q-Q plot (quantile-quantile plot) of all patients in a single group and Q-Q plot of each patient group separately is calculated.

Figure S2. Functional enrichment analysis of DEGs based on NGS dataset.

(A) Enriched GO (biological processes) in differentially expressed genes (DEGs) between DO and control identified using GSEA methods. X-axis indicates LFC values. The color describes the FDR/adjusted p-value as listed in the color bar notifying whether an enriched term is significant (cutoff <0.05). (B) Enriched GO (biological processes) in DEGs between BPS and control identified using GSEA methods.

Figure S3. Deviation Graphs of selected 13 markers in NGS datasets of BPS, DO and controls. Deviation Graphs estimating the cumulative distribution of selected genes. Z-score of 1 denotes that the observation is at a distance of one standard deviation to the right from the centre. (A) TPPP3 (B) FAT1 (C) SMTN (D) ANGPTL7 (E) CLEC3B (F) AIM1 (G) PALM (H) NCALD (I) P2RX2 (J) NRXN3 (K) NRXN2 (L) FAM83A (M) MFAP5

Figure S4. Correlation plot and PCA contribution of selected genes in NGS dataset.

(A) Correlation plot visualizing the relationship between the different attributes. Similarity between observations is defined using the inter-observation distance measures or correlation-based distance measures. The relationships with no significance are visualised by a cross on the correlogram. (B) Scree plot of eigenvalues

showing the percentage of variances explained by each principal component. (C) Barplot of gene contributions in the PCA. A reference dashed line is shown on the barplot corresponding to the expected value if the contribution were uniform. For a given dimension, any gene with a contribution above the reference line could be considered as important in contributing to the dimension.

Figure S5. Log2FC distribution of 13 markers in QPCR dataset. Density ridgeline plot of distribution of BPS, DO and controls for a given gene. (A) TPPP3 (B) FAT1 (C) SMTN (D) ANGPTL7 (E) CLEC3B (F) AIM1 (G) PALM (H) NCALD (I) P2RX2 (J) NRXN3 (K) NRXN2 (L) FAM83A (M) MFAP5

Figure S6. Log2FC distribution of the mean of each 13 markers in QPCR dataset. Density ridgeline plot of distribution of BPS and DO and visualizing the mean of each population.

Figure S7. Normality test and PCA contribution of different markers in QPCR dataset. Normality test for 13 selected markers: the normal plots (histogram), Q-Q plot (quantile-quantile plot) of all patients in a single group and Q-Q plot of each patient group separately is calculated.

Figure S8. Testing the performance of selected 13 markers in QPCR datasets of BPS, DO and controls. Boxplot Statistics visualization (Based on log2FC), pairwise comparisons including p-value. (A) TPPP3 (B) FAT1 (C) SMTN (D) ANGPTL7 (E) CLEC3B (F) AIM1 (G) PALM (H) NCALD (I) P2RX2 (J) NRXN3 (K) NRXN2 (L) FAM83A (M) MFAP5

Figure S9. Empirical cumulative distribution function (ECDF) of selected 13 markers in QPCR datasets of BPS, DO and controls. ECDF reports for any given number (mRNA read count) the percentage of individuals that are below that threshold.

(A) TPPP3 (B) FAT1 (C) SMTN (D) ANGPTL7 (E) CLEC3B (F) AIM1 (G) PALM (H) NCALD (I) P2RX2 (J) NRXN3 (K) NRXN2 (L) FAM83A (M) MFAP5

Figure S10. Deviation Graphs of selected 13 markers in QPCR datasets of BPS, DO and controls. Deviation Graphs estimating the cumulative distribution of selected genes. Z-score of 1 denotes that the observation is at a distance of one standard deviation to the right from the centre. (A) TPPP3 (B) FAT1 (C) SMTN (D) ANGPTL7 (E) CLEC3B (F) AIM1 (G) PALM (H) NCALD (I) P2RX2 (J) NRXN3 (K) NRXN2 (L) FAM83A (M) MFAP5

Figure S11. Performance of selected 13 markers (QPCR datasets of BPS, DO and controls) in PCA (A) Correlogram of the relationship between the different attributes. Relationship with no significance is visualised by a cross (B) Scree plot of eigenvalues showing the percentage of variances explained by each principal component. (C) Barplot of gene contributions. For a given dimension, any gene with a contribution above the reference line could be considered as important in contributing to the dimension. (D) PCA biplot showing the separation samples using QPCR data.

Figure S12. Standard deviation of performance metrics for various models. (A) Standard deviation of metrics in BPS vs control, (B) DO vs control and (C) BPS vs DO.

Figure S13. Models' performance based on k-fold cross-validation (CV) method. Each plot represents an ML model, where the first part of the model's name corresponds to a particular ML algorithm, and the latter part indicates a data resampling method. It displays metric scores for the corresponding model, with a higher shaded area indicating better performance. The models are sorted in increasing order based on the F1 score from the k-fold CV result. (A) Thirty-six radar graphs illustrate model performances for BPS vs. Control. (B) Thirty-six radar graphs illustrate

model performances for BPS vs. DO. **(C)** Thirty-six radar graphs illustrate model performances for DO vs. Control.

Figure S14. Models' performance based on nested cross-validation (CV) method.

Each plot represents an ML model, where the first part of the model's name corresponds to a particular ML algorithm, and the latter part indicates a data resampling method. It displays metric scores for the corresponding model, with a higher shaded area indicating better performance. The models are sorted in increasing order based on the F1 score from the k-fold CV result. **(A)** Thirty-six radar graphs illustrate model performances for BPS vs. Control. **(B)** Thirty-six radar graphs illustrate model performances for BPS vs. DO. **(C)** Thirty-six radar graphs illustrate model performances for DO vs. Control.

Figure S15. MLcps Results. MLcps scores are depicted in Cleveland's dot plot for **(A)** BPS vs. control, **(B)** DO vs. control, and **(C)** BPS vs. DO.

Figure S16. Proposed ML Framework. The proposed pipeline first splits the datasets into k (3) equal-sized bins in a stratified manner, where k-1 bins will be used as training datasets and the remaining bin as a test dataset. Next, it employs the RFECV method for feature selection, three data resampling techniques, and twelve ML algorithms. Then, it evaluates the model performance on the test dataset using the k-fold and nested CV (k=3) method and calculates seven different performance metrics. As part of the k-fold CV methods, it repeats the last two steps for each unique bin. It repeats the complete process ten times and takes the average performance as the final model performance. In the end, the pipeline provides a list of the intersection of selected features from the top 10 best-performing models (based on F1 score) as the final list of selected features.