

# **DIVAN: accurate identification of non-coding disease-specific risk variants using multi-omics profiles**

## **All Supplemental Figures**

**Figure S1: Precision-recall curves of five-fold cross-validation on four diseases: carotid artery disease, macular degeneration, ulcerative colitis, and multiple sclerosis.**

**Figure S2: ROC curves of five-fold cross-validation for 41 diseases.**

**Figure S3: Cross-disease prediction shows the necessity for disease specificity of variant prioritization by using four diseases from four different disease classes: carotid artery disease (cardiovascular disease), Alzheimer's disease (mental disease), multiple sclerosis (immune disease), and macular degeneration (eye disease). (A) ROC curves showing that disease-specific prediction outperforms cross-disease prediction. (B) Heatmap showing that disease-specific prediction outperforms cross-disease prediction.**

**Figure S4: Precision-recall curves showing the effectiveness of the feature selection and ensemble method by comparing feature selection and ensemble combined, feature selection only, ensemble only, and the**

**baseline case: neither feature selection nor ensemble.**

**Figure S5: Contribution of three different feature groups, TF binding/histone modification/open chromatin, on prediction in four diseases: carotid artery disease, macular degeneration, ulcerative colitis, and multiple sclerosis. (A)** ROC curves showing the predictive performance using each of the three feature groups of epigenomic features in the predictive model: TF binding, histone modification, and open chromatin, with read as the continuous feature for the four diseases. **(B)** ROC curves showing the predictive performance using each of the three feature groups of epigenomic features in the predictive model, TF binding, histone modification, and open chromatin, with peak as the binary feature for the four diseases.

**Figure S6: ROC curves showing predictive performance of disease-class variant prioritization for immune diseases, including rheumatoid arthritis, asthma, type 1 diabetes mellitus, systemic lupus erythematosus, and multiple sclerosis.** The ROC curves are generated by using a “leave-one-disease-out” approach; that is, the predictive model is built using variants of all other diseases within the disease class, tested on the variants of the disease being left out.

**Figure S7: Bar chart showing the number of informative features for three feature categories (TF binding/histone modification/open chromatin) for 45 diseases when using peak as the feature.**

**Figure S8: Distribution of distances between non-coding SNPs associated with the 45 diseases in ARB to their nearest TSS (ignoring stand).**

**Figure S9: Distribution of enrichment proportions of H3K9me3 features and open chromatin features when using peak and read as the feature. The enrichment proportion is defined as the ratio between average read counts of each H3K9me3/open chromatin feature on risk variants and benign variants respectively.**

**Figure S10: Distributions of test statistics and p-values calculated between risk variants and benign variants for four diseases: carotid artery disease, macular degeneration, ulcerative colitis, and multiple sclerosis. (A) Distribution of t-test statistics using read as the feature for the four diseases. (B) Distribution of p-values of t-test for the four diseases. (C) Distribution of p-values of Fisher's exact test using peak as the feature for the four diseases.**

**Figure S11: ROC curves and precision-recall curves of five-fold**

**cross-validation using three base learners, decision tree, SVM, and Lasso, for four diseases: carotid artery disease, macular degeneration, ulcerative colitis, and multiple sclerosis. (A)** ROC curves of five-fold cross-validation for three base learners: decision tree, SVM, and Lasso on the four diseases **(B)** Precision-recall curves of five-fold cross-validation for three base learners, decision tree, SVM, and Lasso, on the four diseases.

**Figure S1**

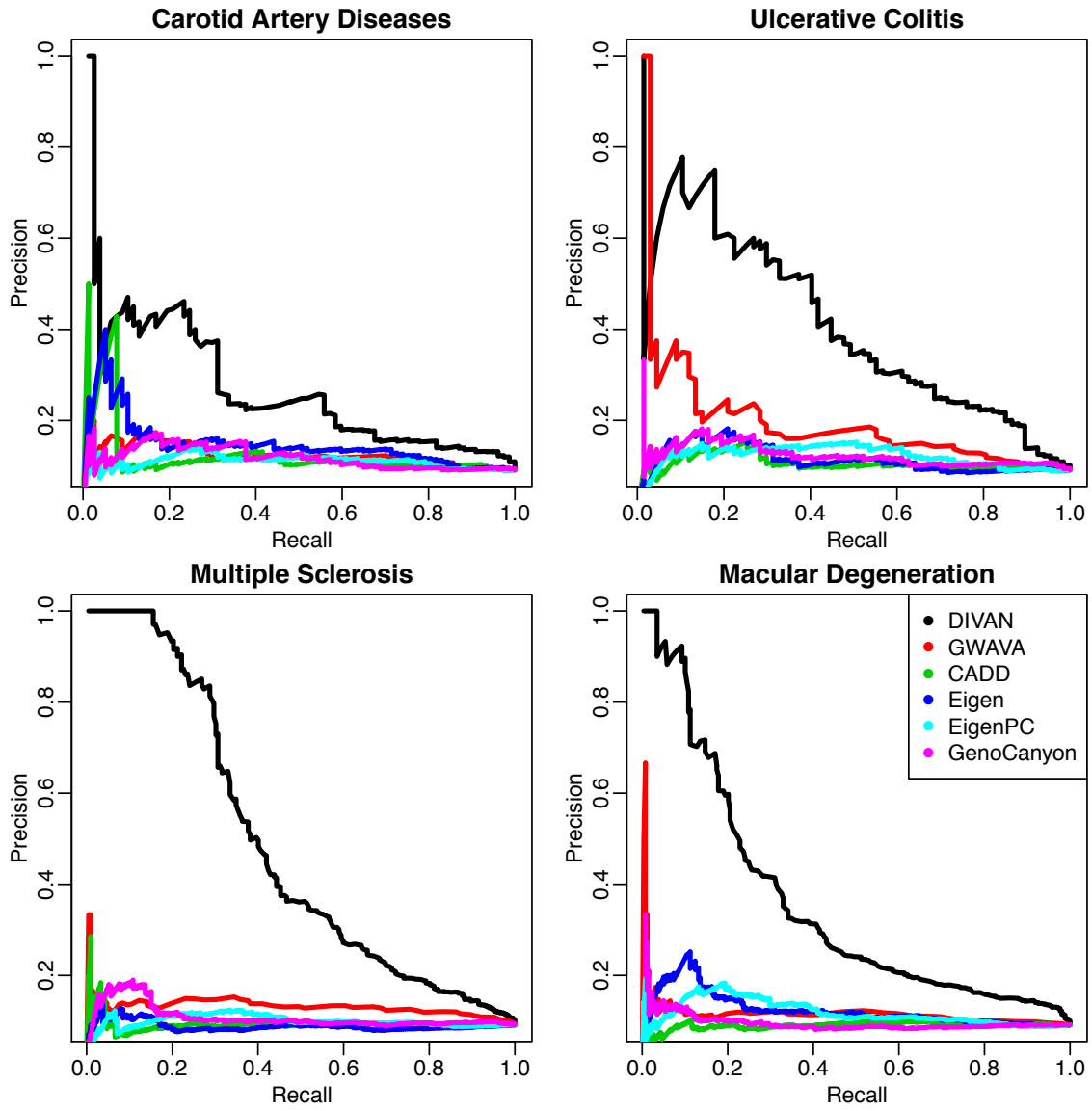
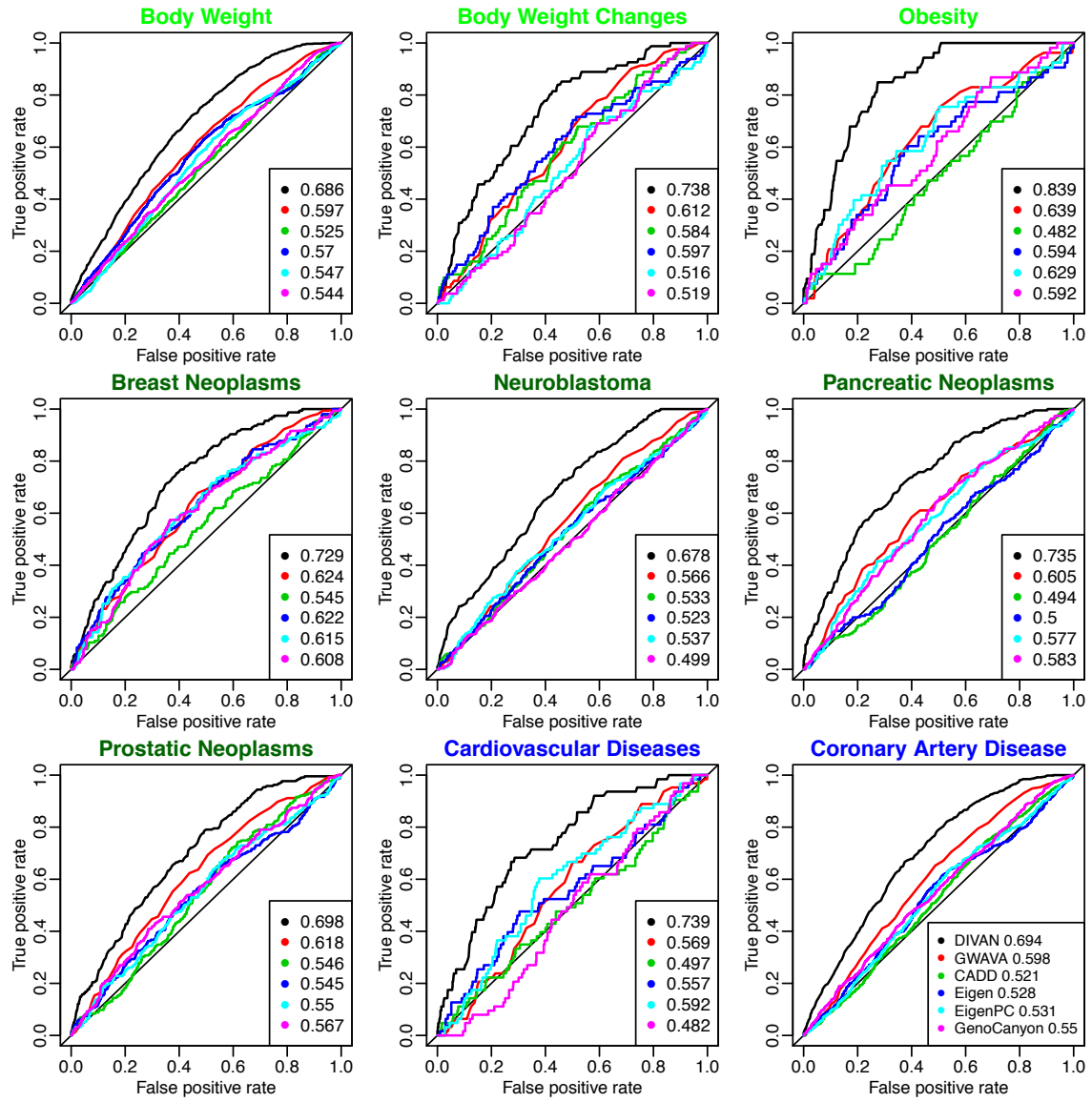
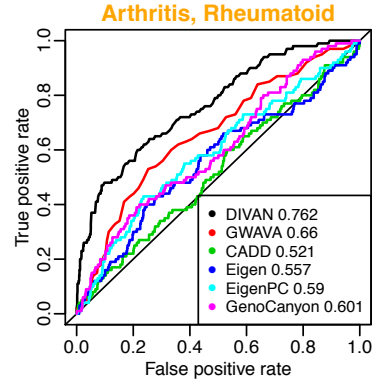
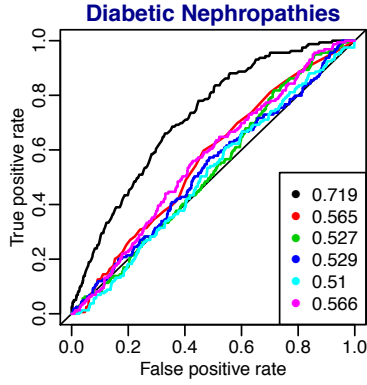
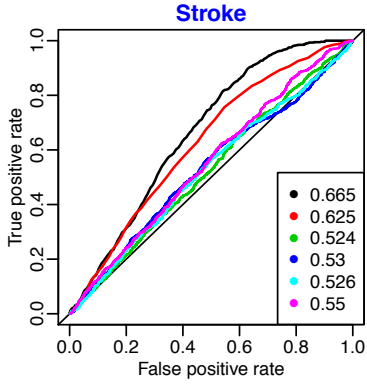
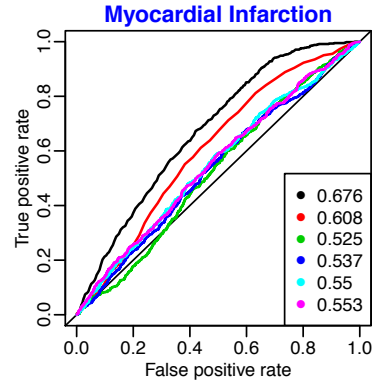
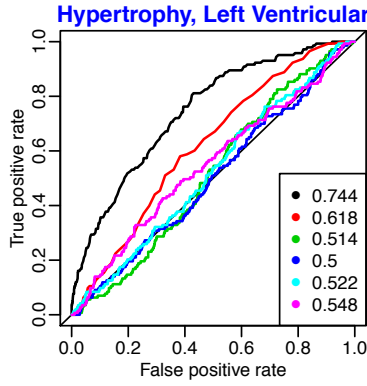
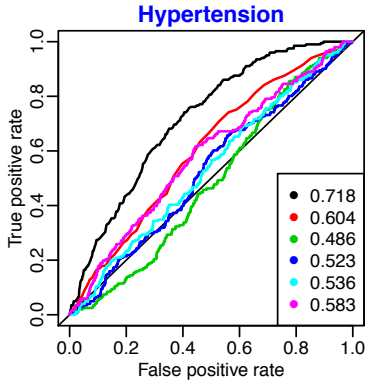
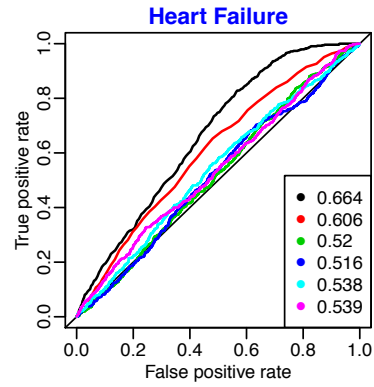
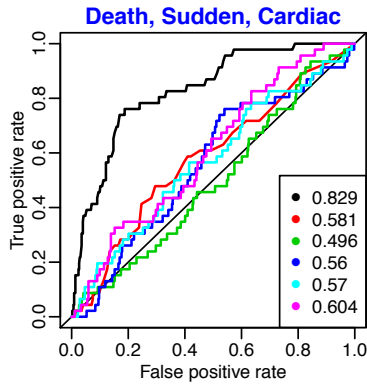
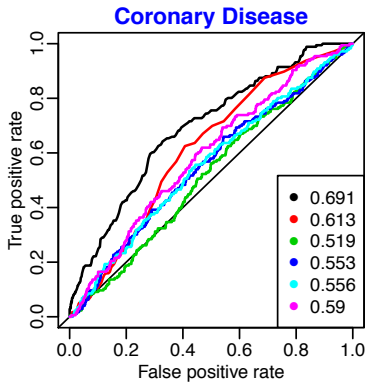
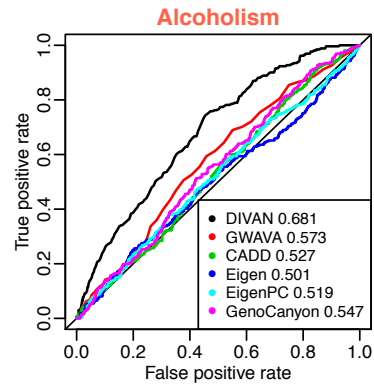
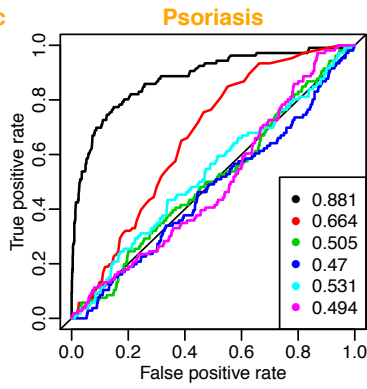
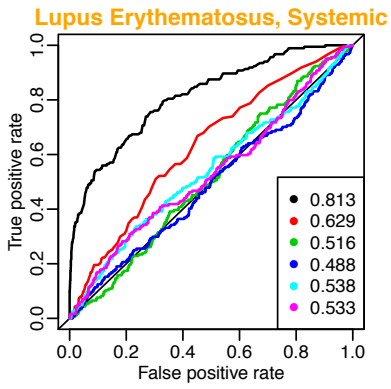
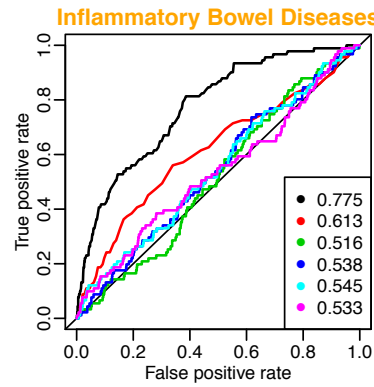
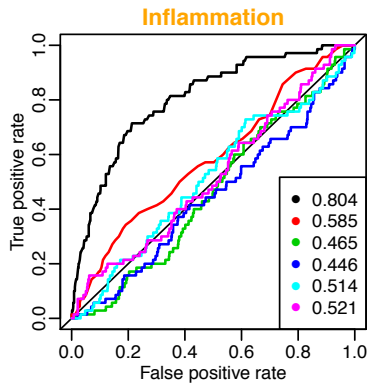
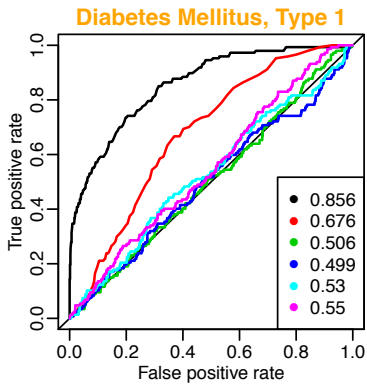
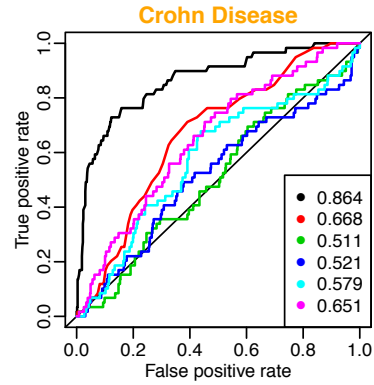
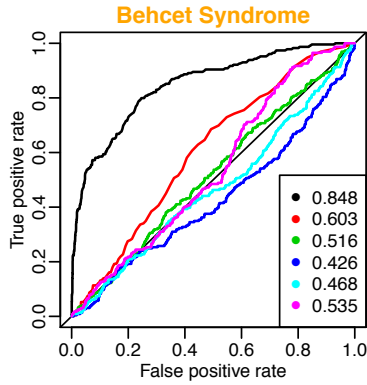
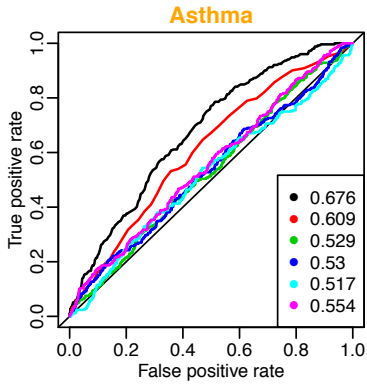


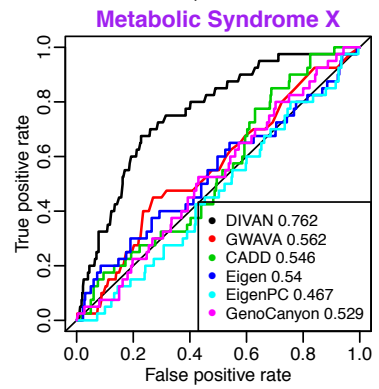
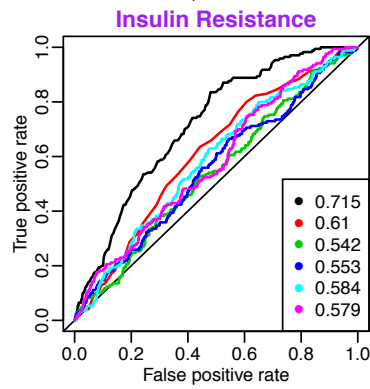
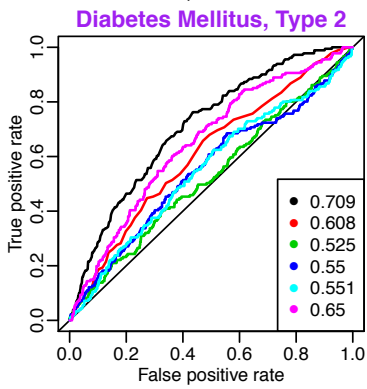
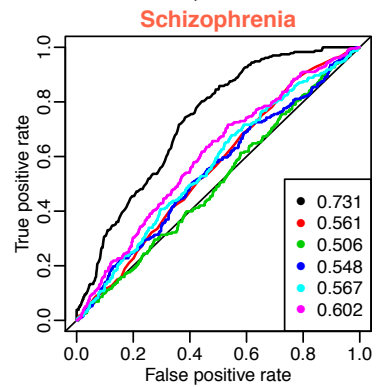
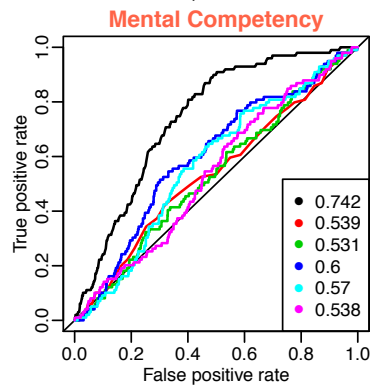
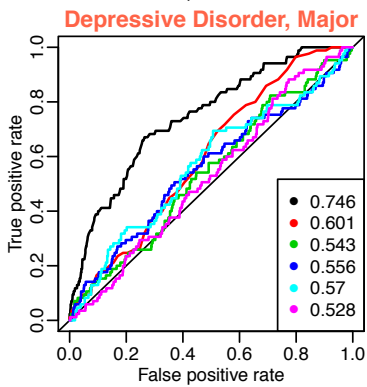
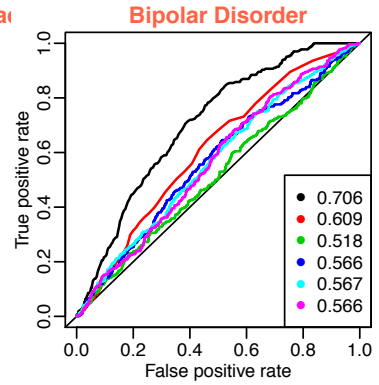
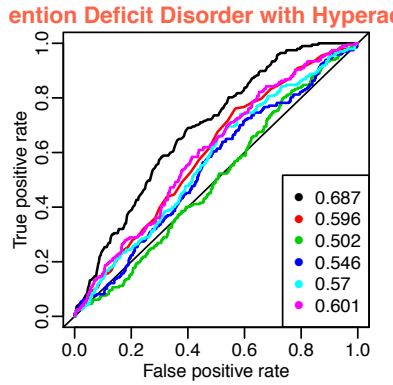
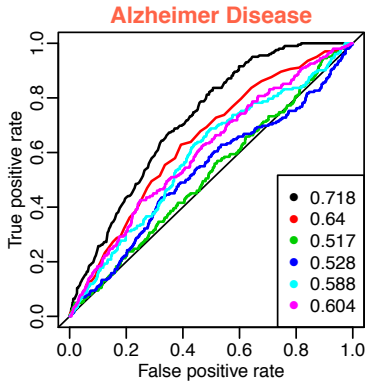
Figure S2











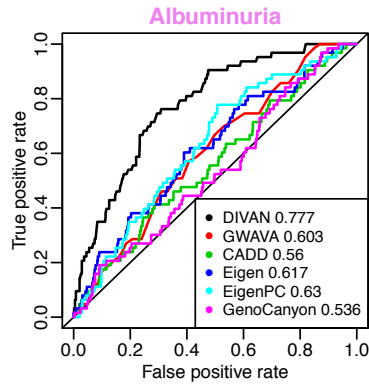
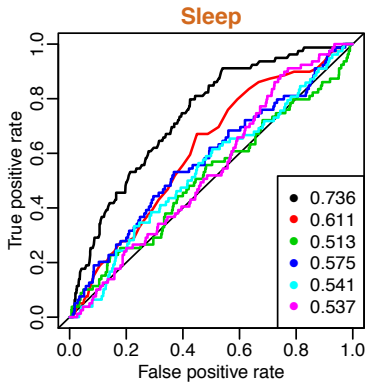
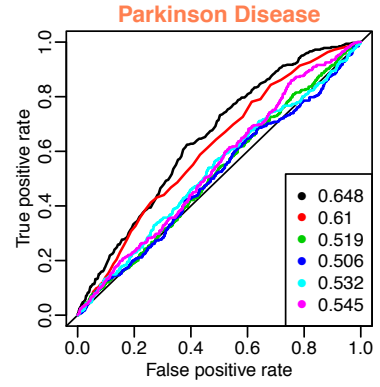
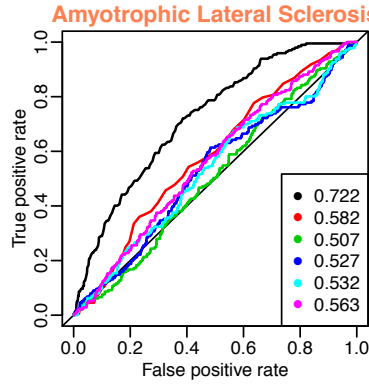
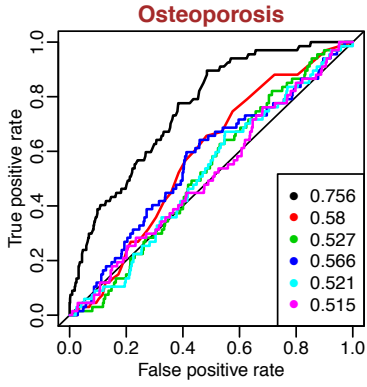
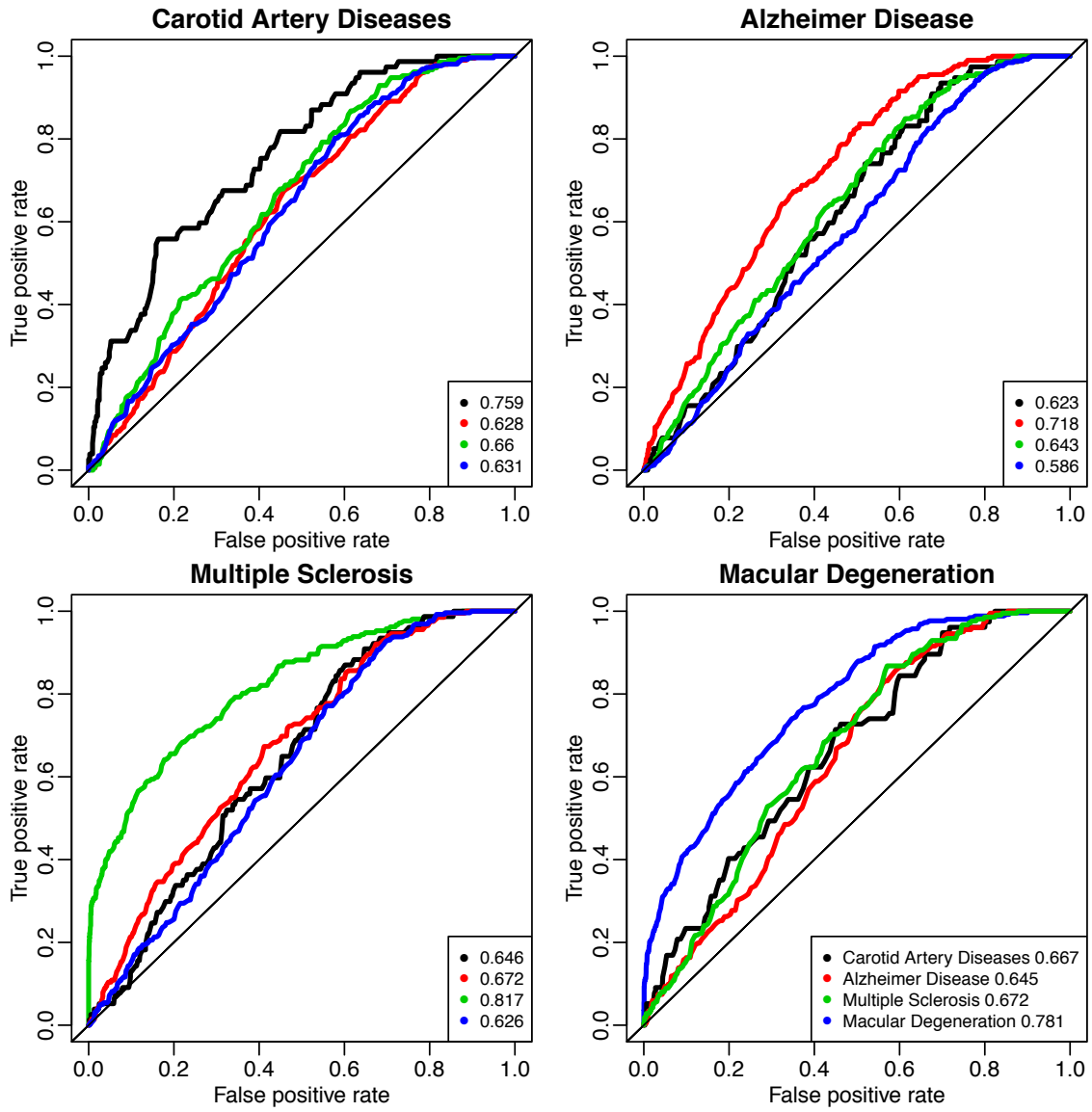
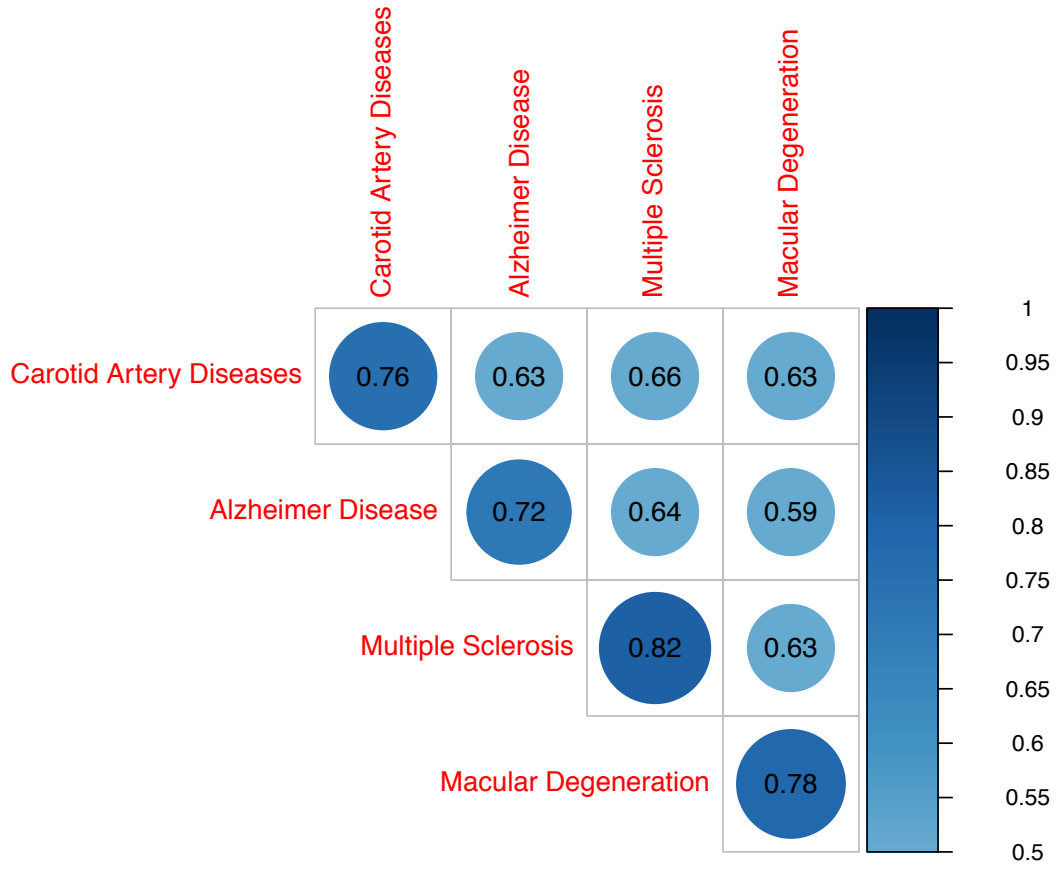


Figure S3-A



**Figure S3-B**



**Figure S4**

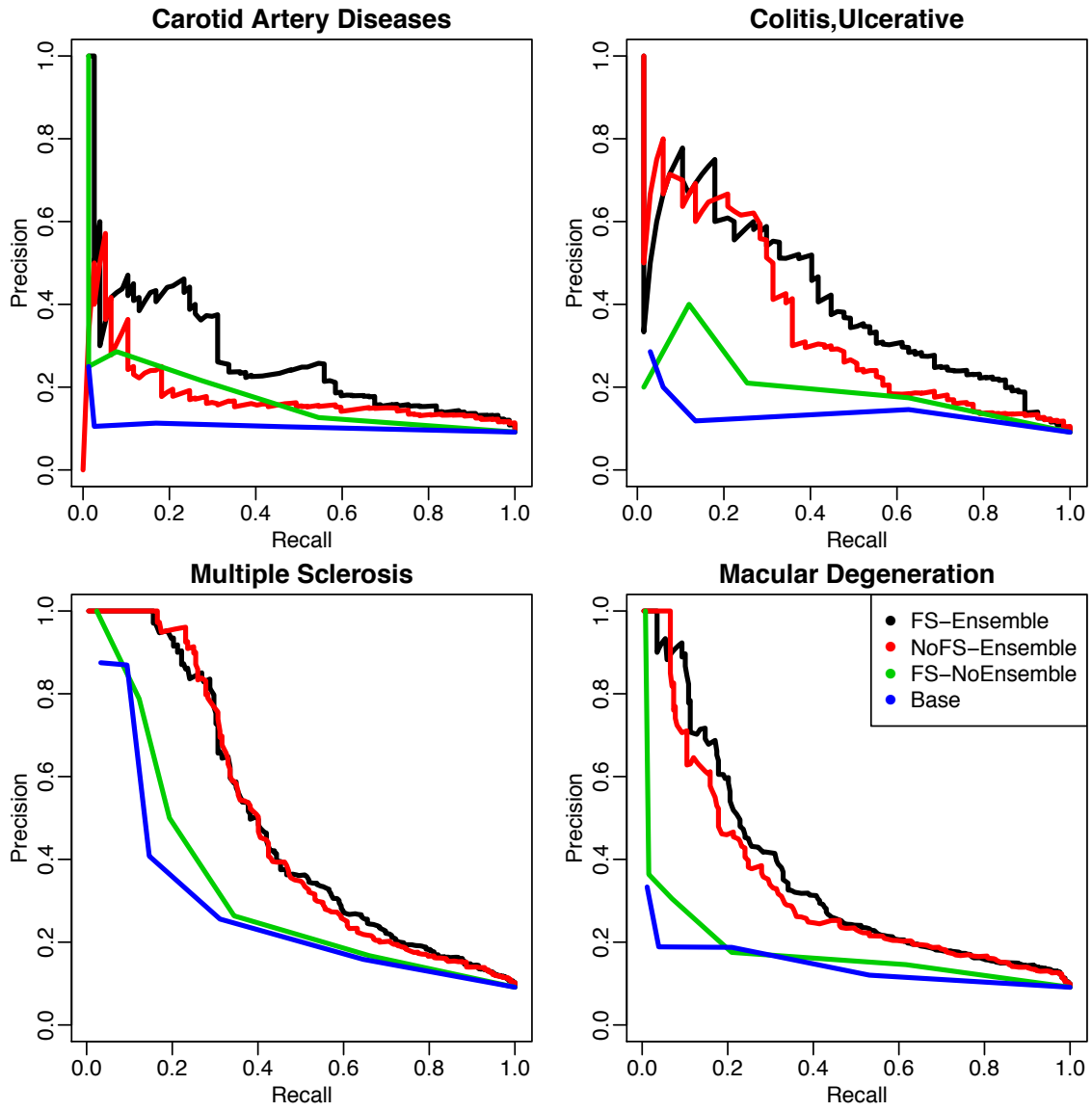
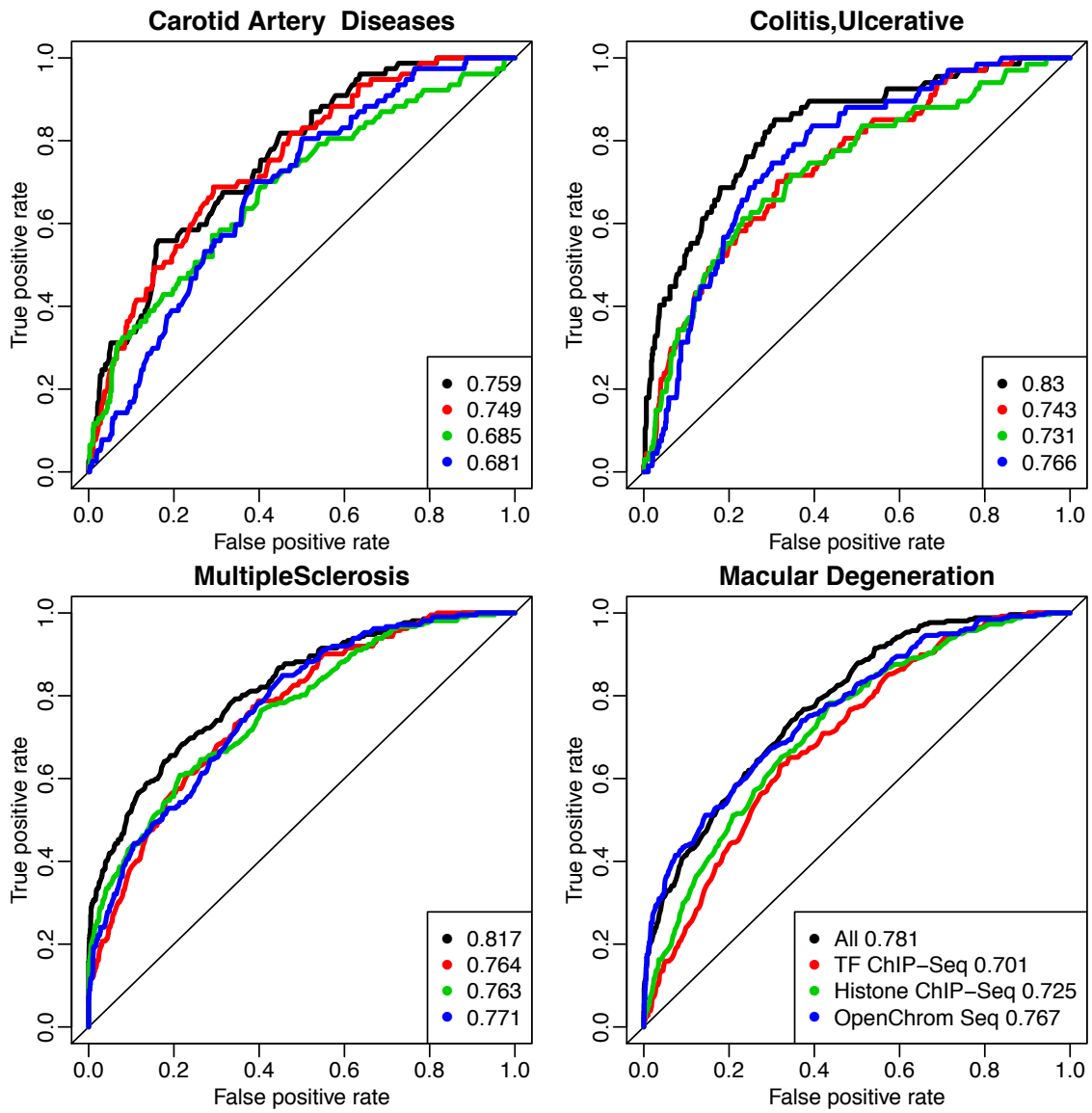


Figure S5-A



**Figure S5-B**

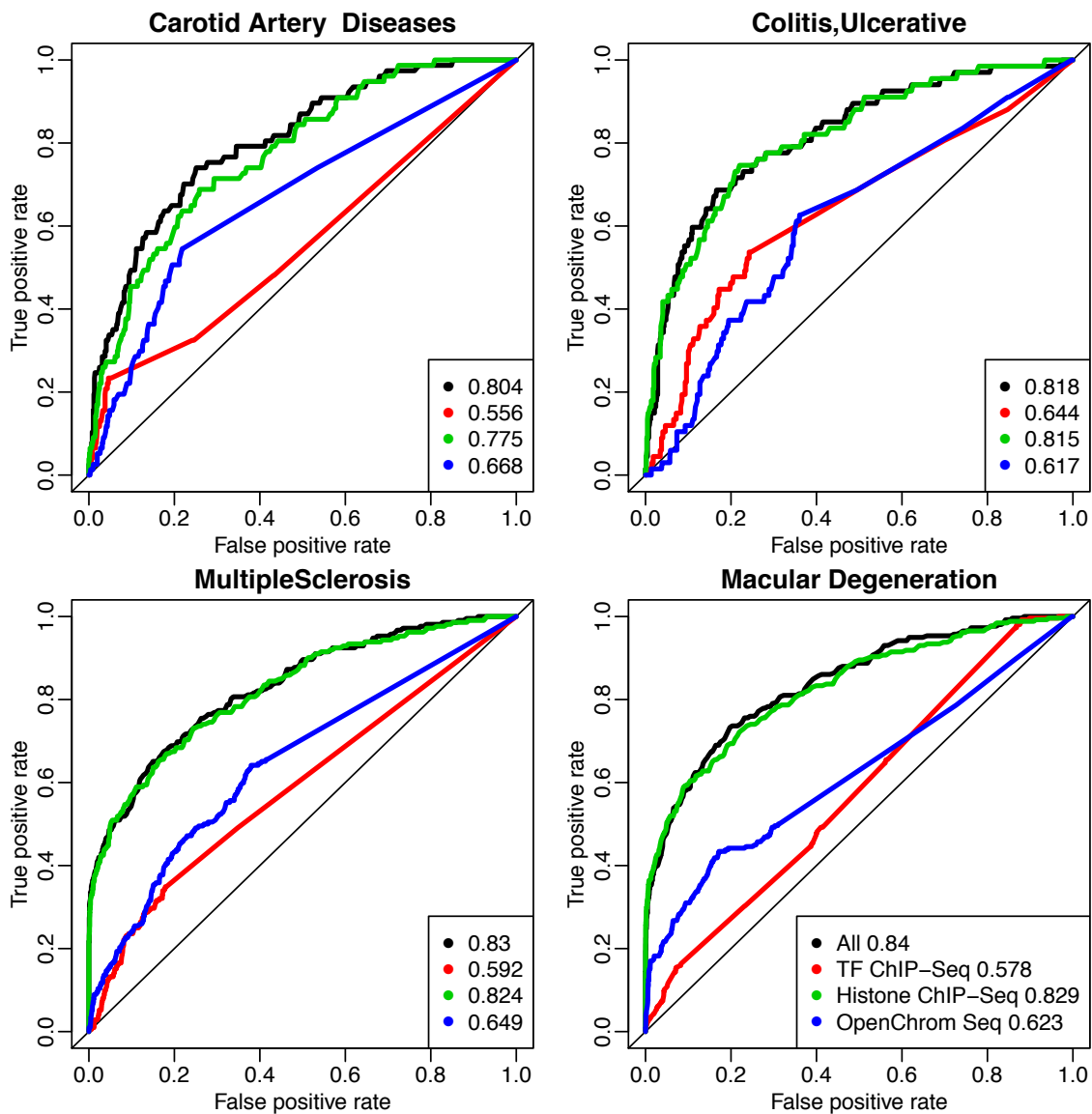


Figure S6

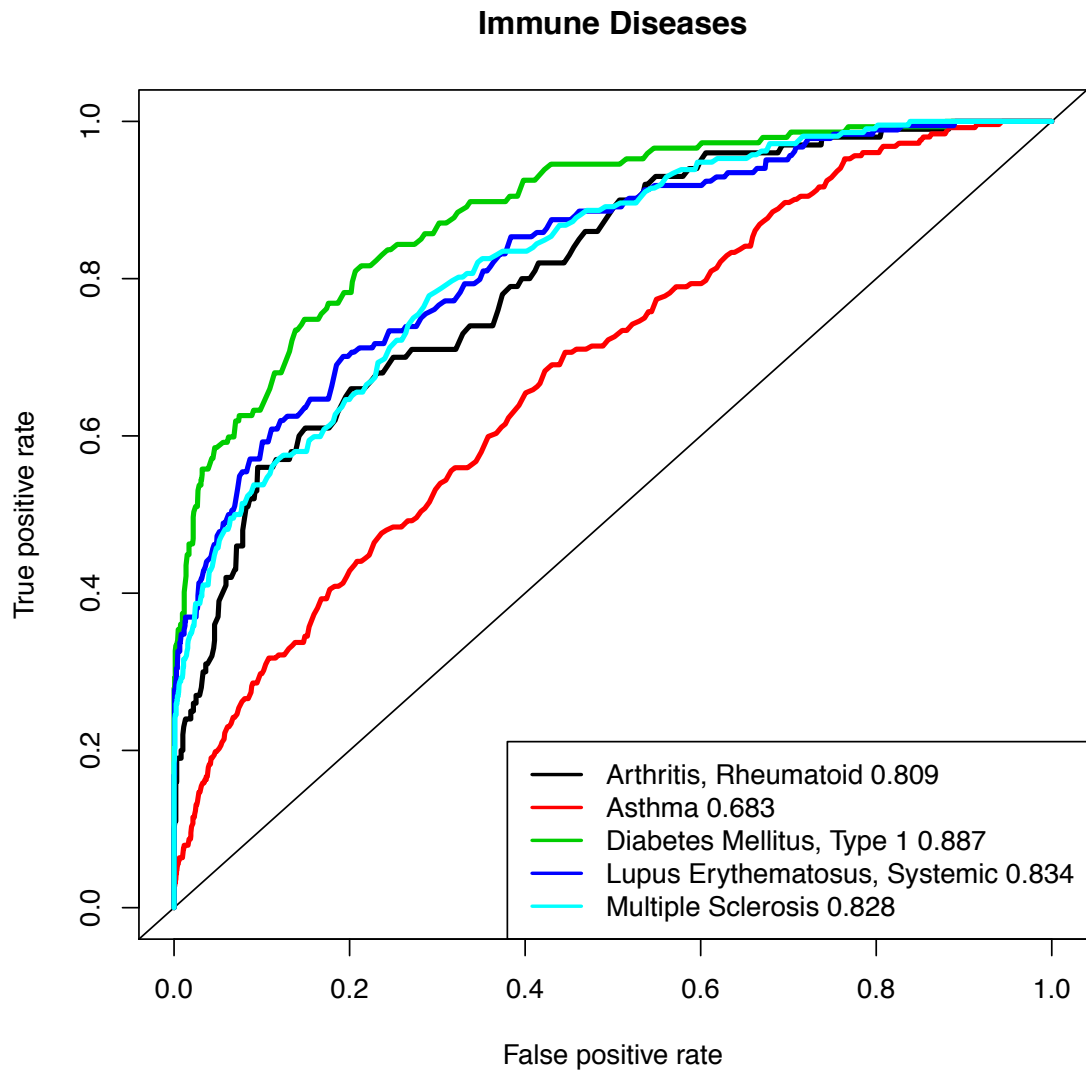
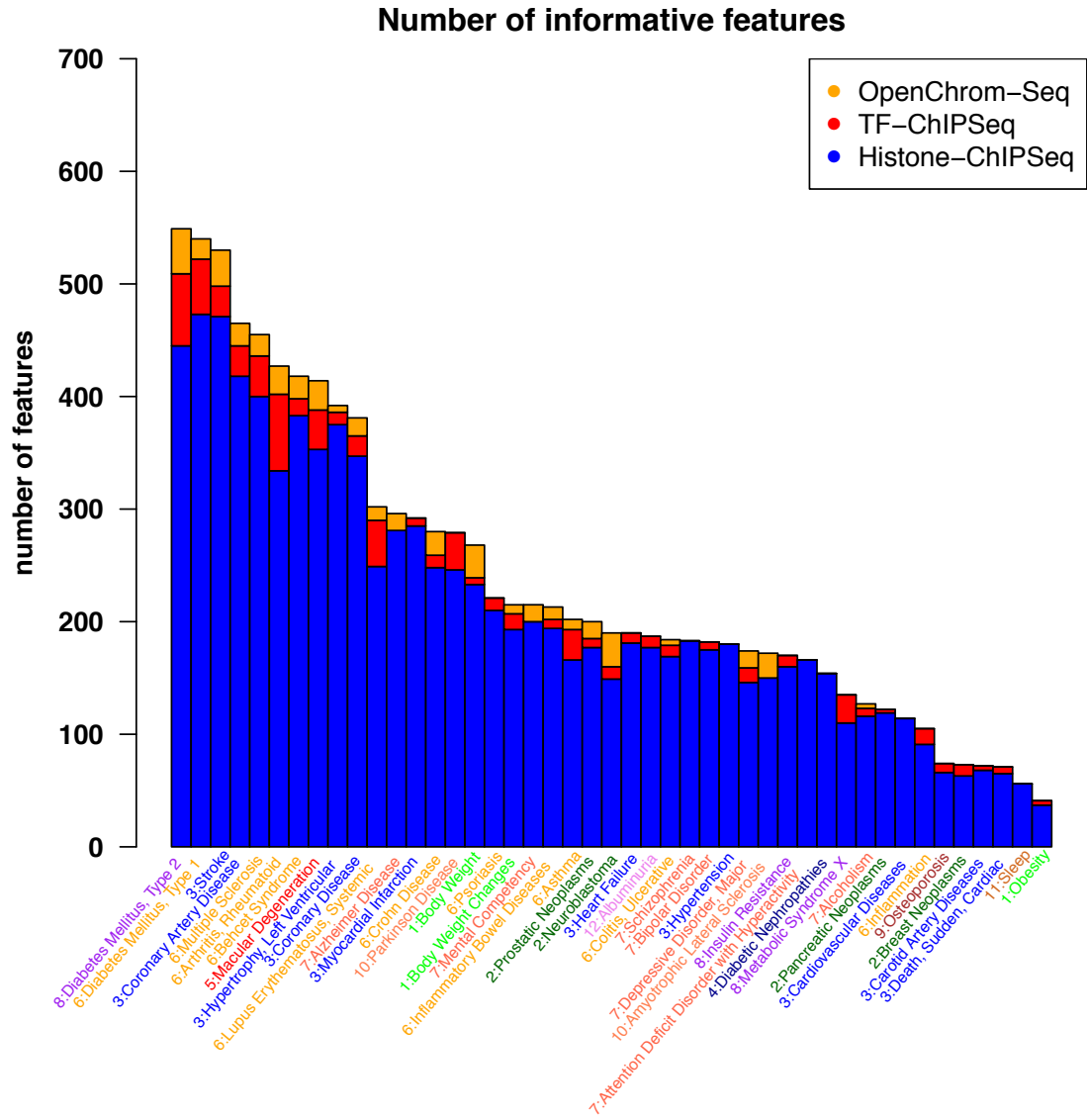




Figure S7



**Figure S8**

**Distances from non-coding SNPs(ARB) to nearest TSS(ignore strand)**

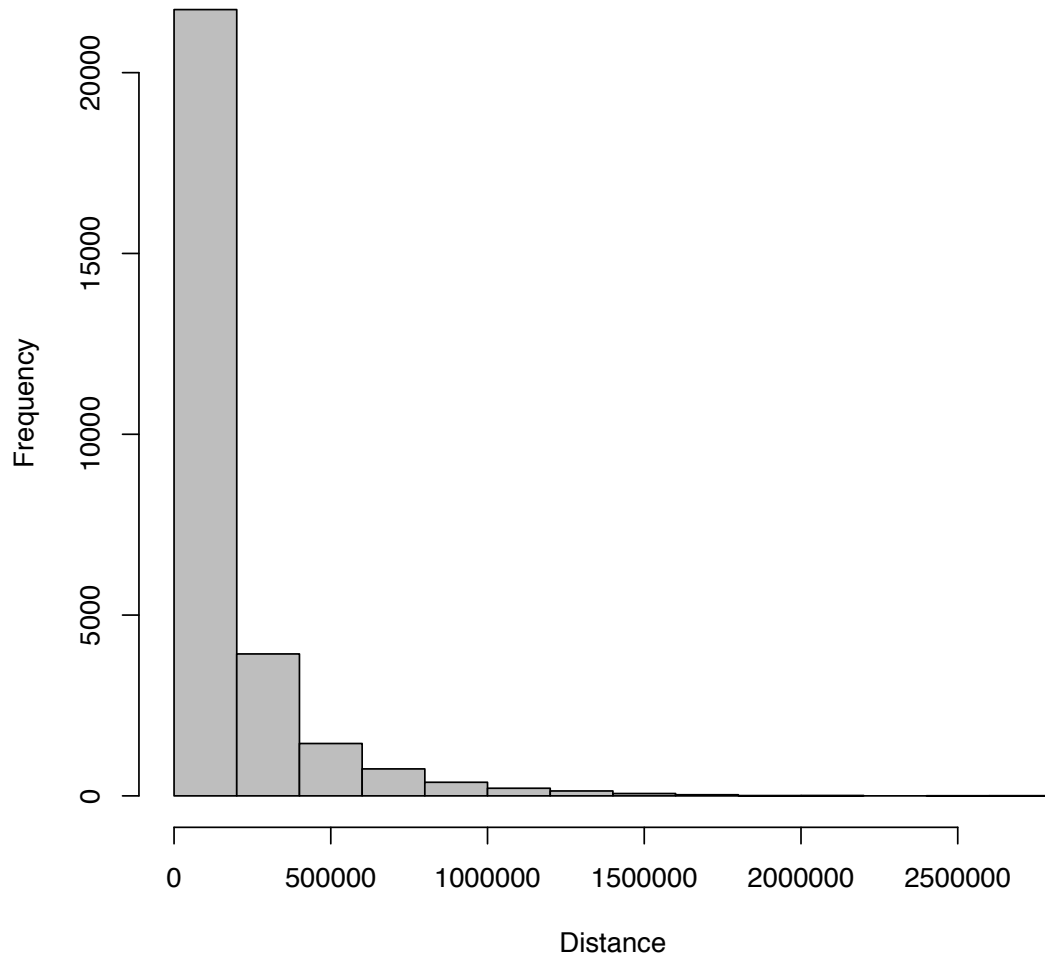
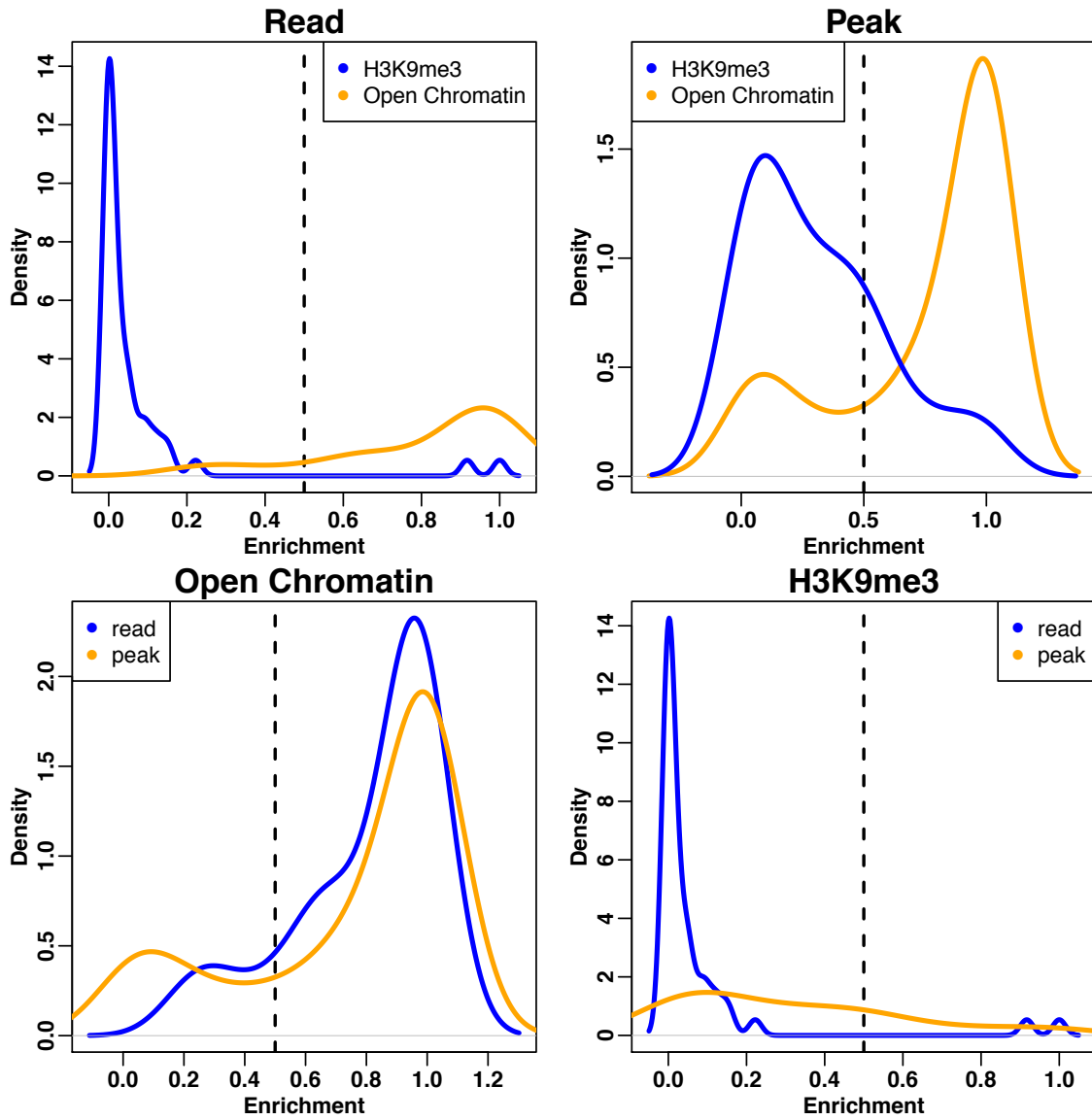
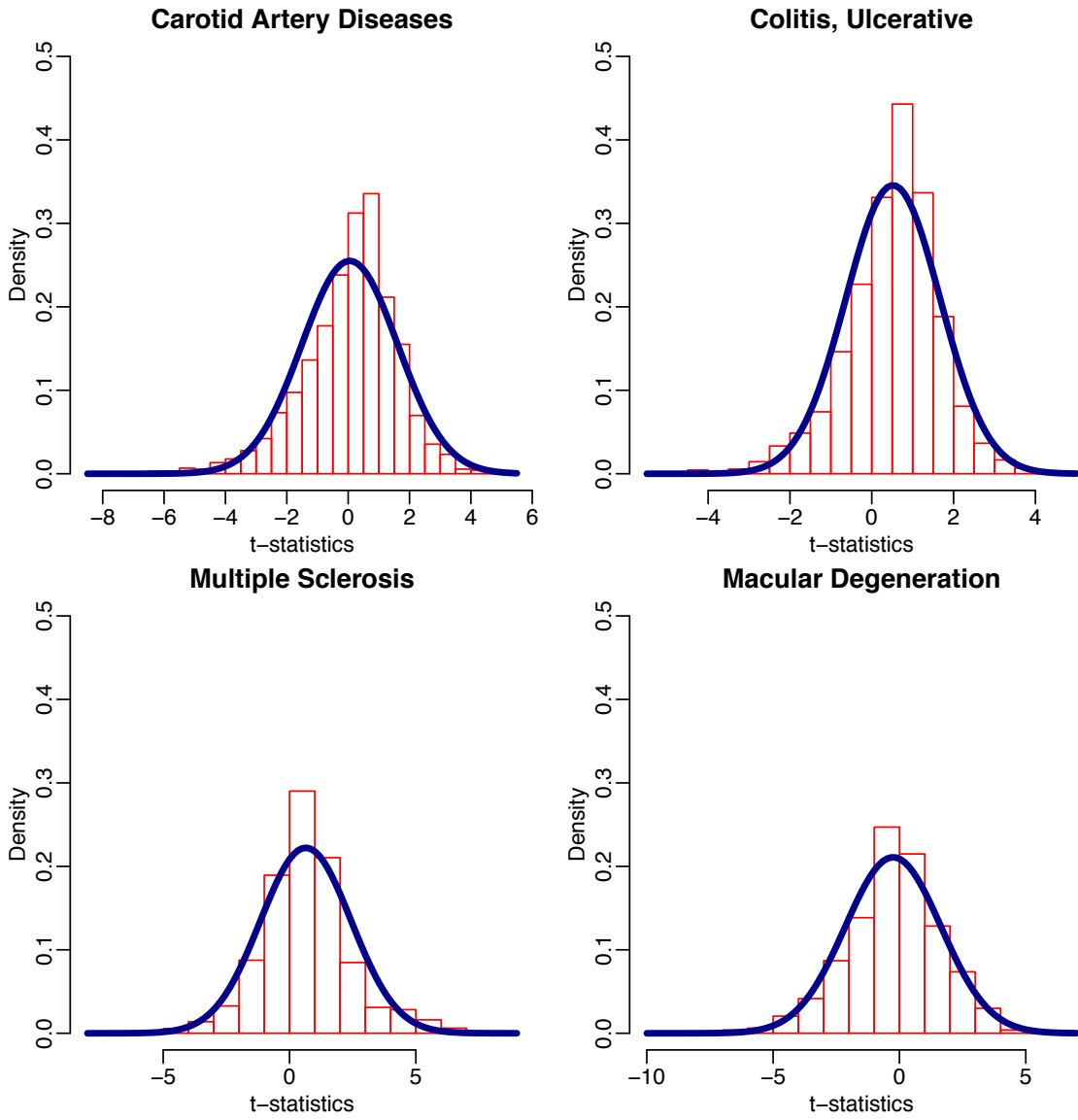


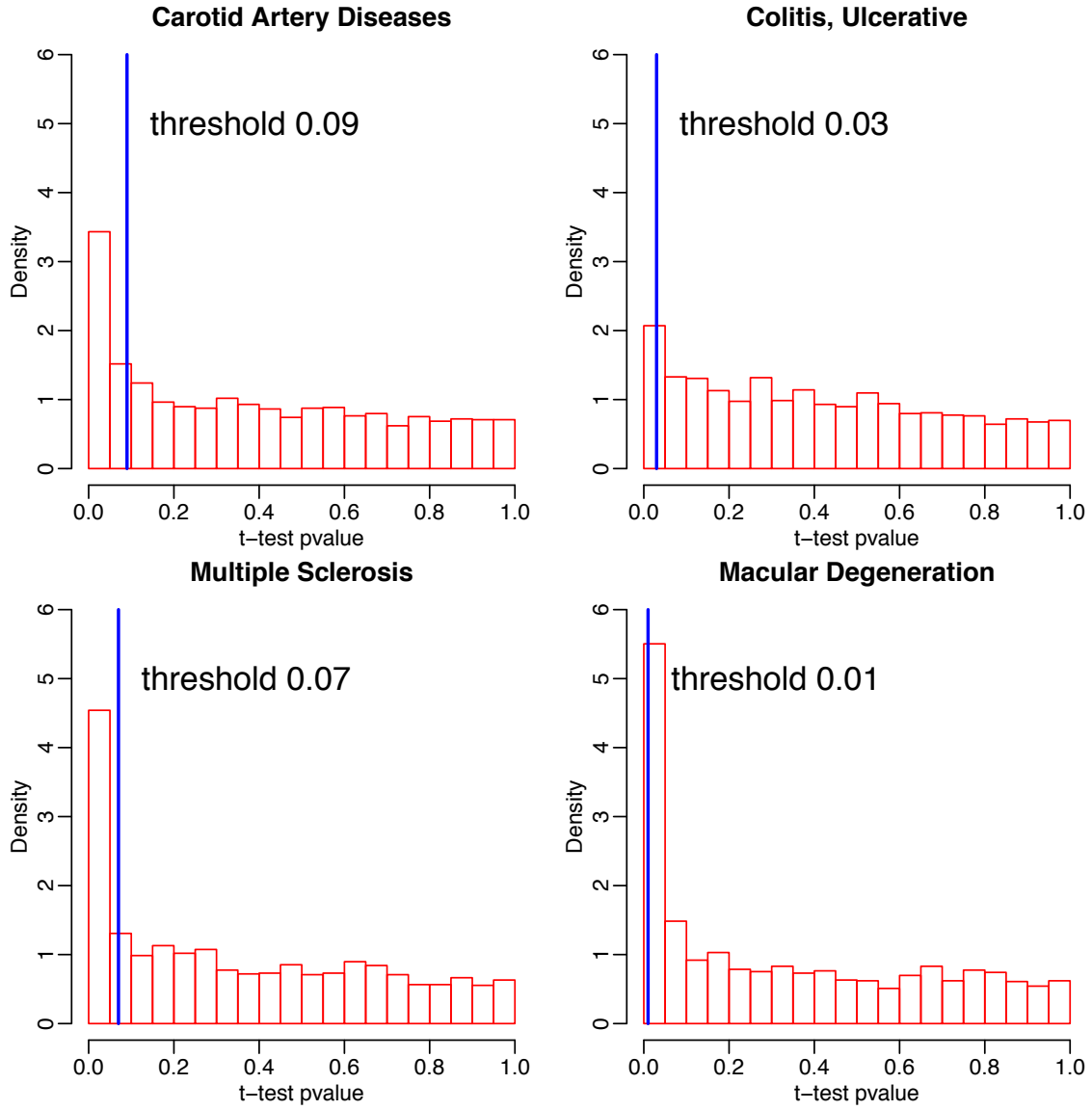
Figure S9



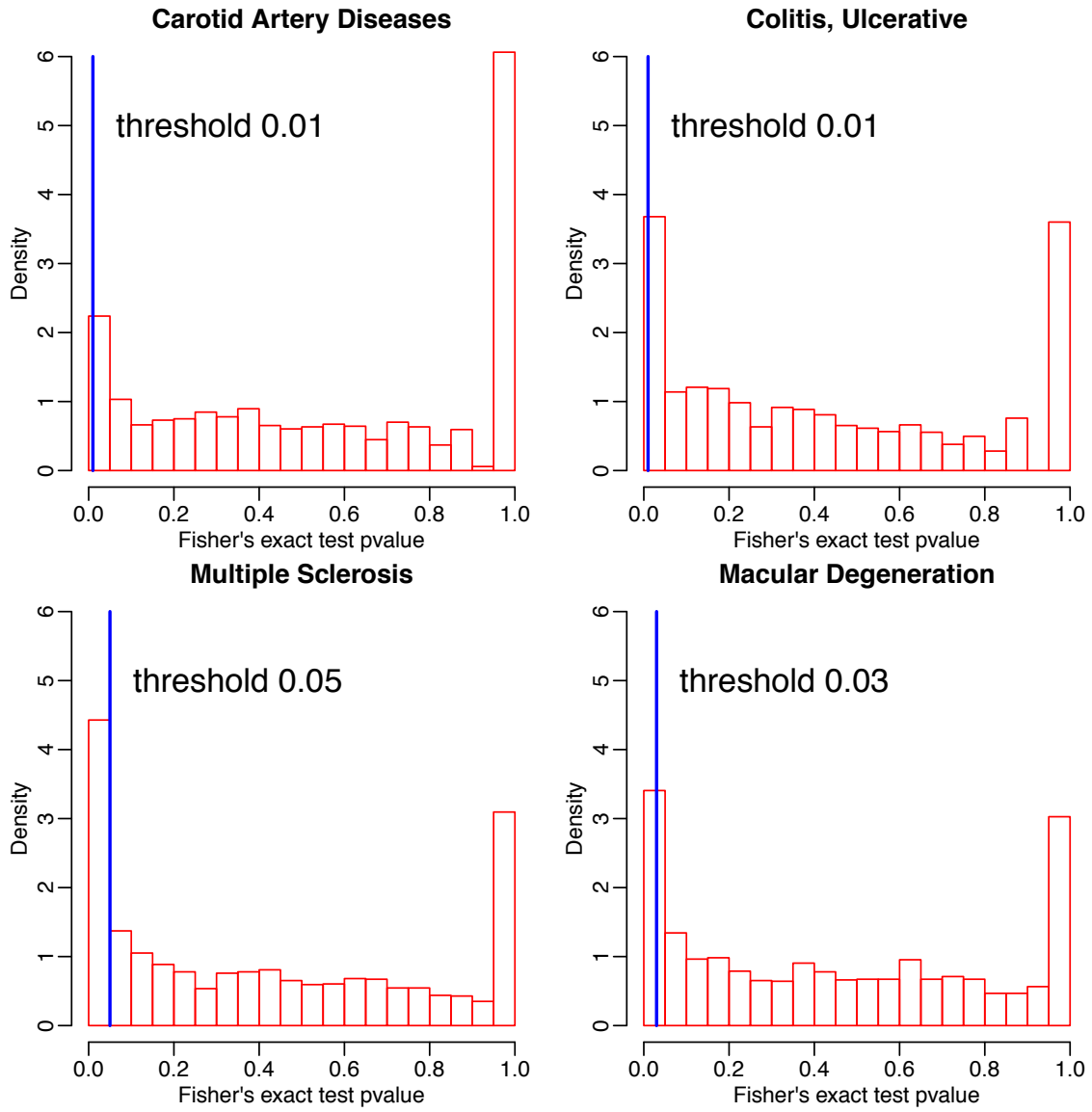
**Figure S10-A**



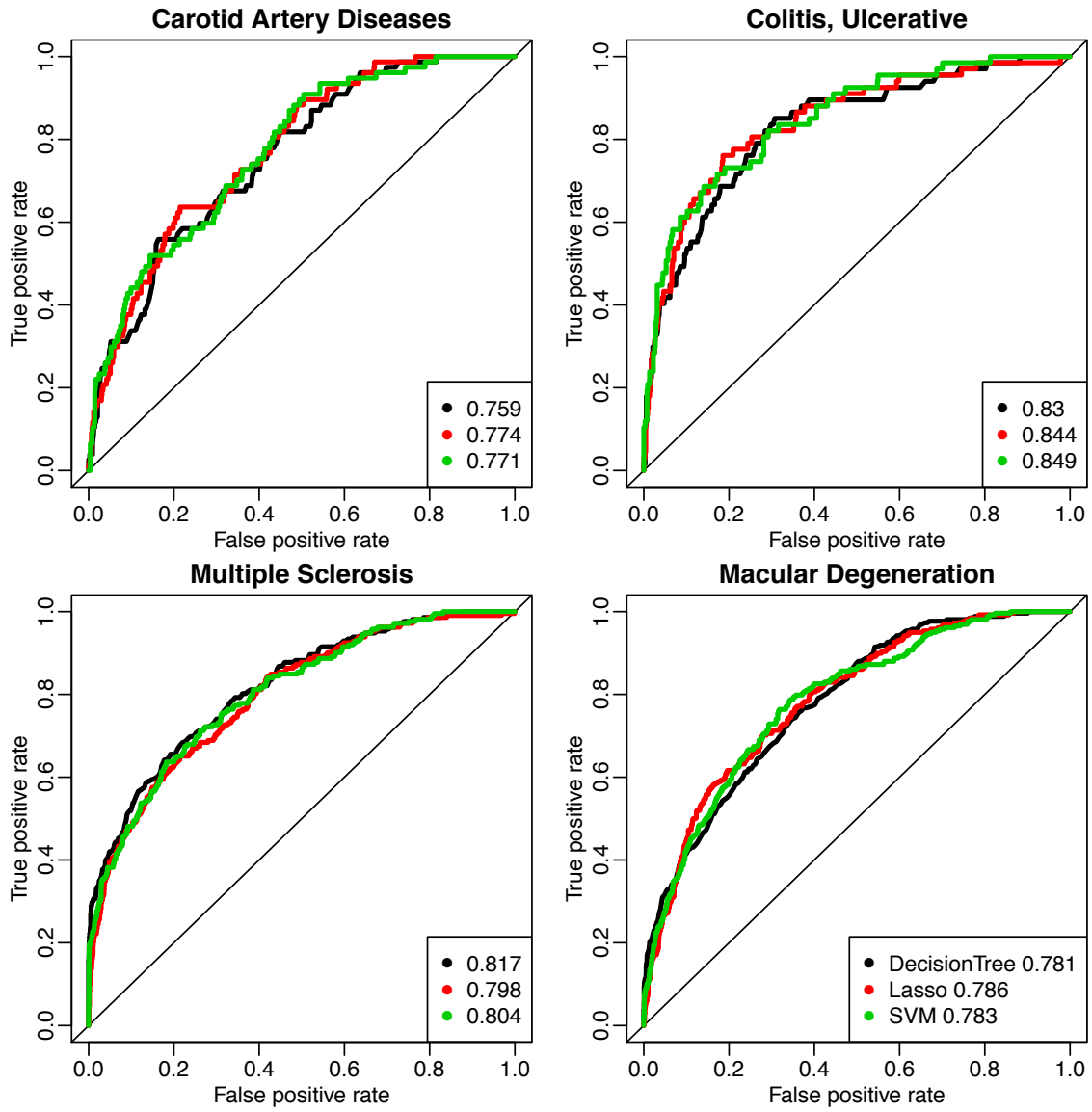
**Figure S10-B**



**Figure S10-C**



**Figure S11-A**



FigureS11-B

