

Supplementary Material Based on Illustrations

Paper: Explanations Based on Item Response Theory (*eXirt*):
A Model-Specific Method to Explain Tree-Ensemble
Model in Trust Perspective

41 binary classification problem datasets provided by the *OpenML*.

15 properties of 41 datasets organized in table format.

Use of 6 XAI methods present in the current literature and 1 proposed new method.

4 different algorithms



LightGBM



CatBoost

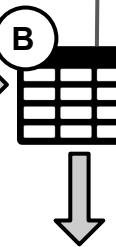
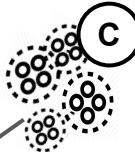


Random Forest



Gradient Boosting

Application of the *K-means* clustering and *MCA* algorithm in dataset properties. Four datasets clusters were identified and your profiles.



The following procedures were performed for each of the datasets in each cluster.

4 different models were created for each of the 41 datasets.

Pre-processing of features and application of min-max normalization in all 41 datasets.

4 different algorithms



Ciu



Dalex



Eli5



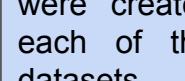
eXirt



Lofo

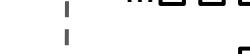


Shap

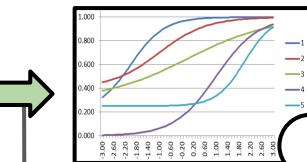
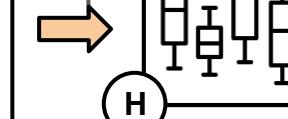


Skater

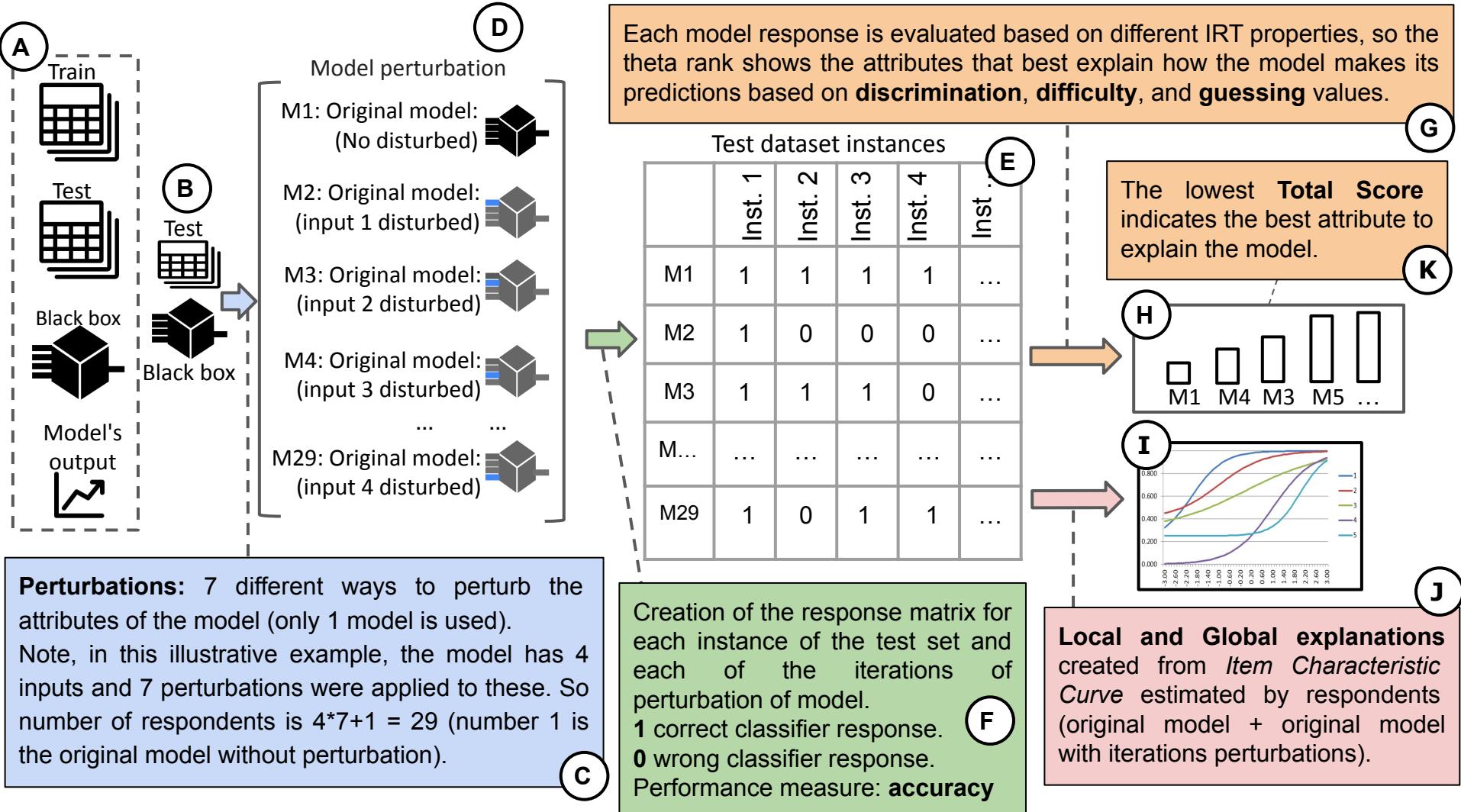
Global Feature Relevance Rank



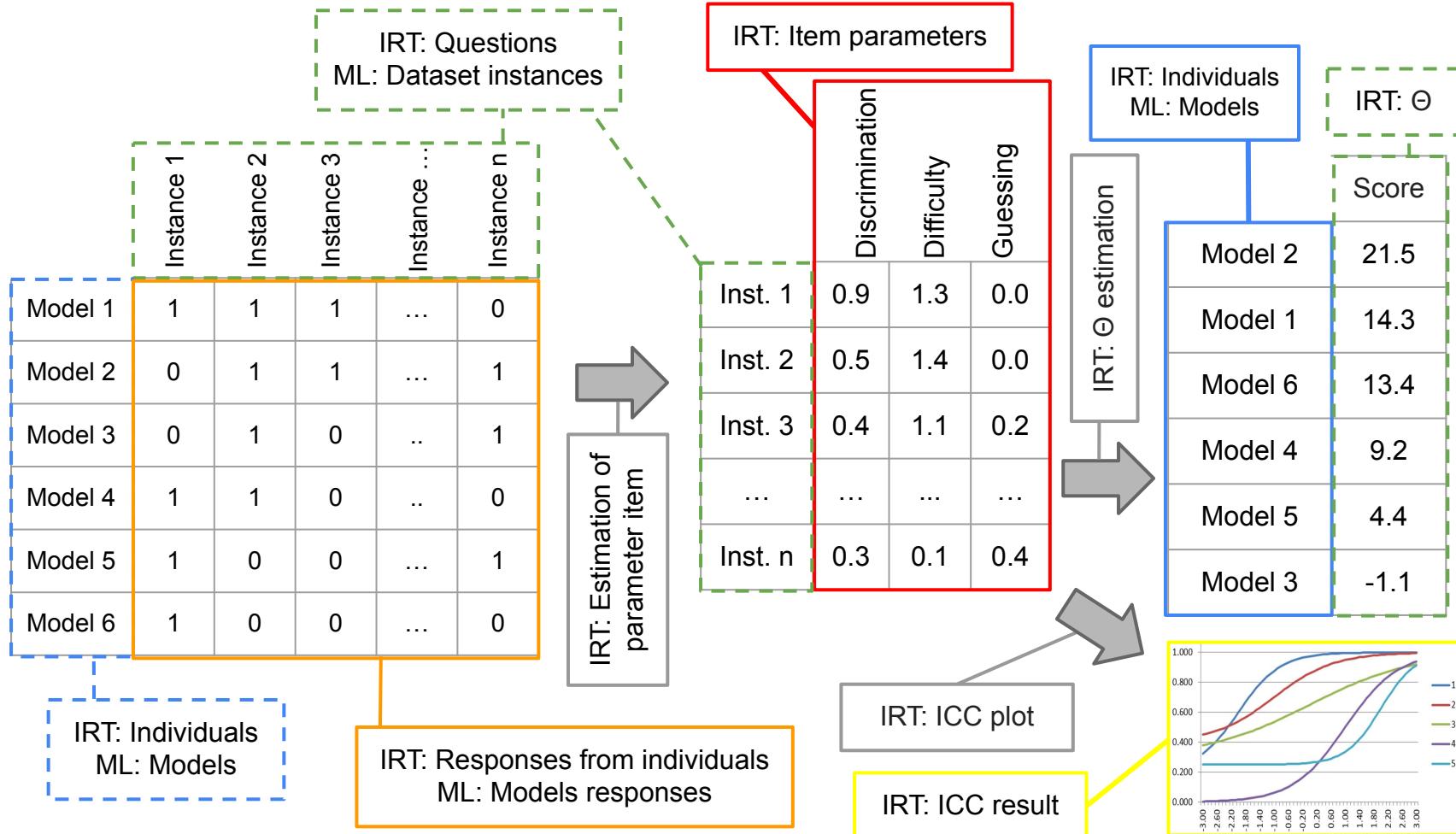
Considering different dataset clusters, the creation of boxplots containing comparisons of the *eXirt* ranks with the ranks of the other XAI methods.



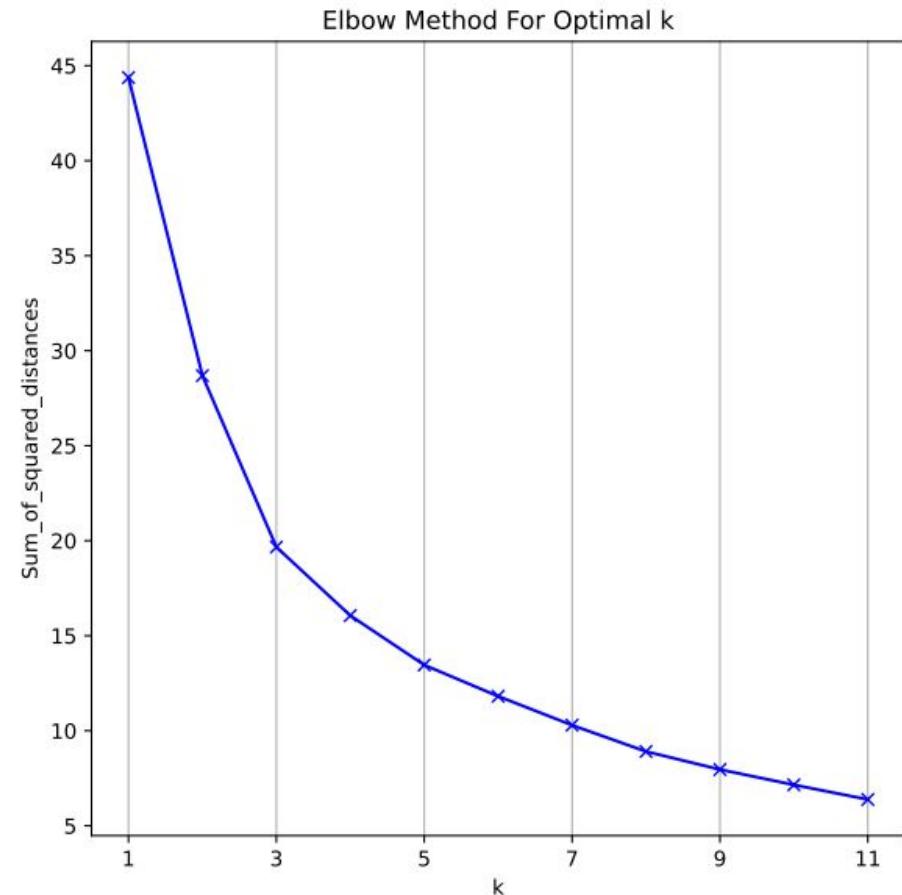
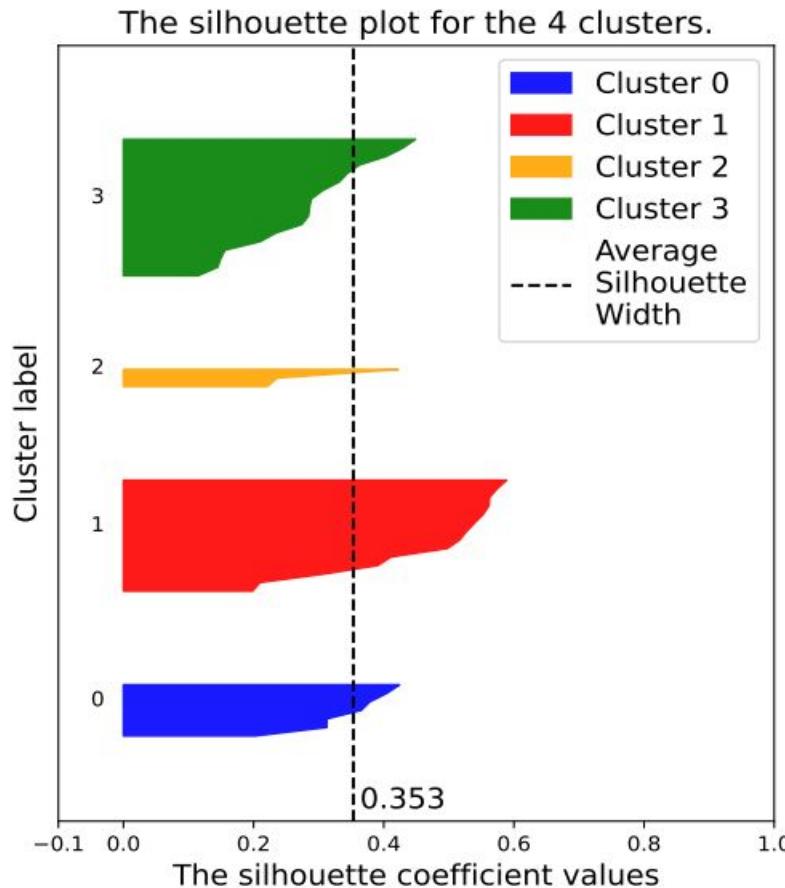
Local and Global explanations based on *Item Characteristic Curve*. These explanations are differentiators of *eXirt*, as they generate insights into confidence in the model.



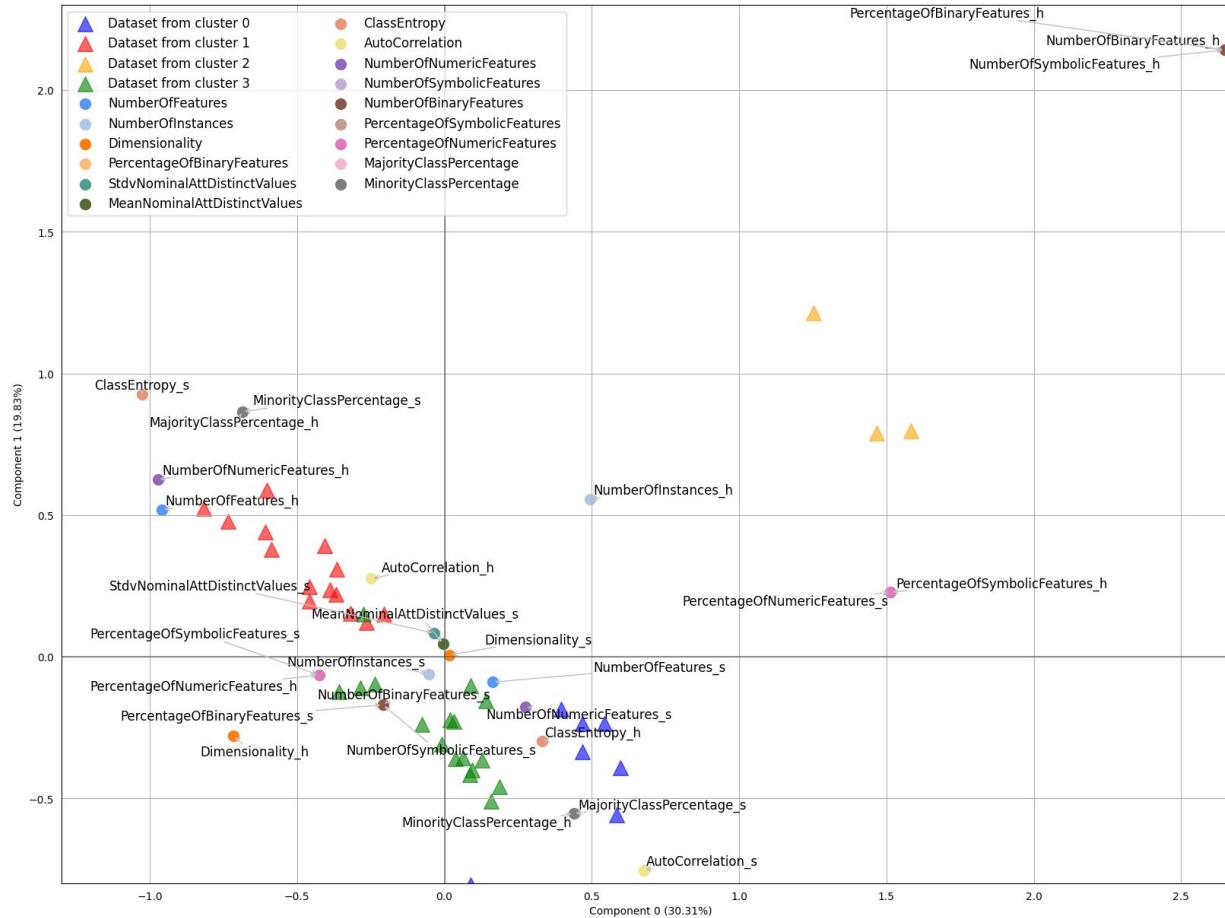
Relation of terms IRT and ML



K-size identification of clustering using Silhouette and Elbow



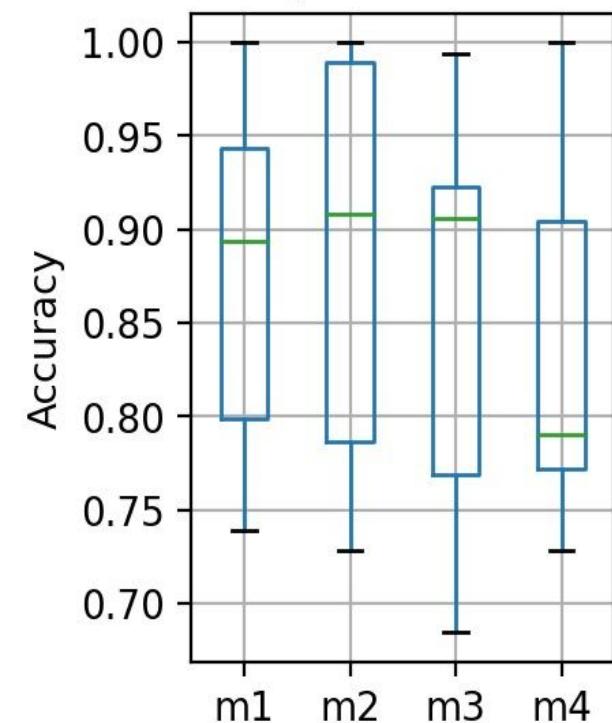
Profiles of datasets identified through Multiple Correspondence Analysis



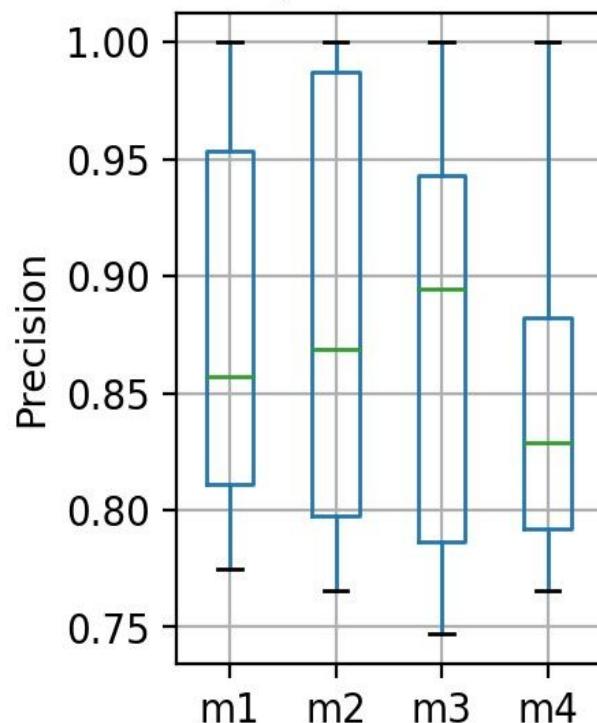
Cluster 0

Data regarding the performance of models M1 to M4 based on cluster 0

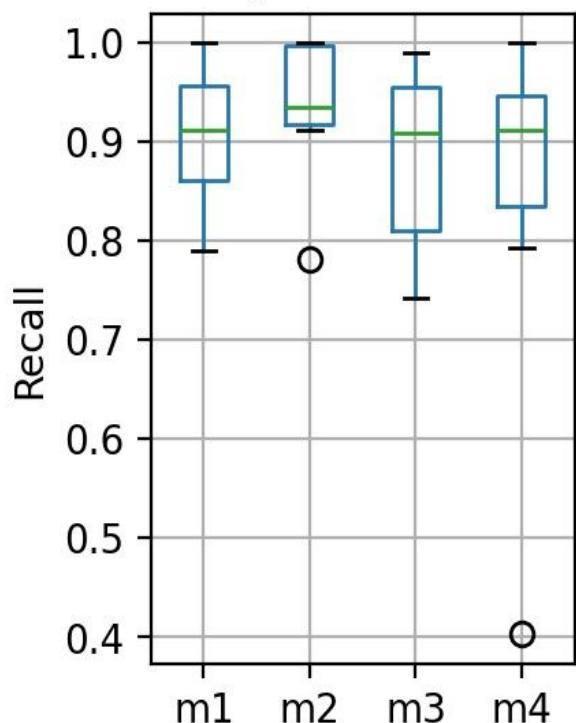
Accuracy of models
by cluster: 0



Precision of models
by cluster: 0

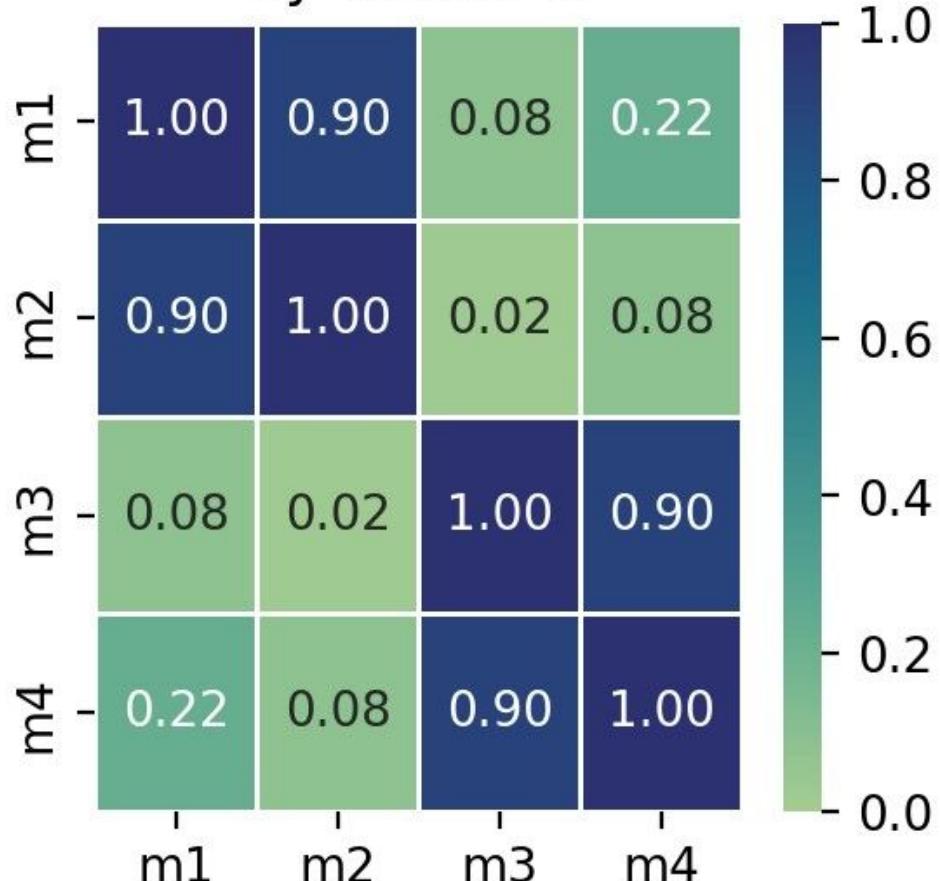


Recall of models
by cluster: 0

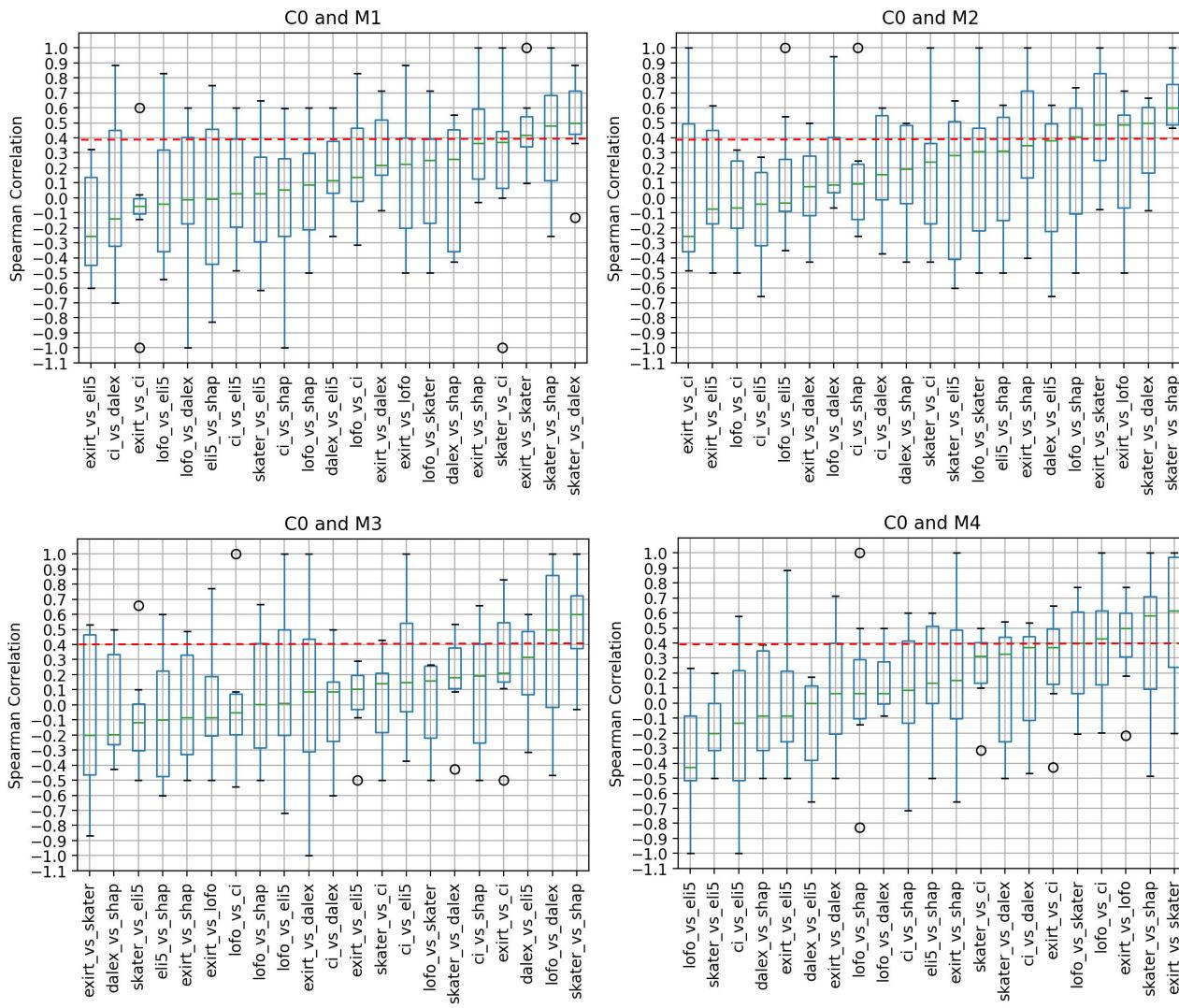


Friedman test result - Models by Cluster 0

Friedman test of models by cluster 0

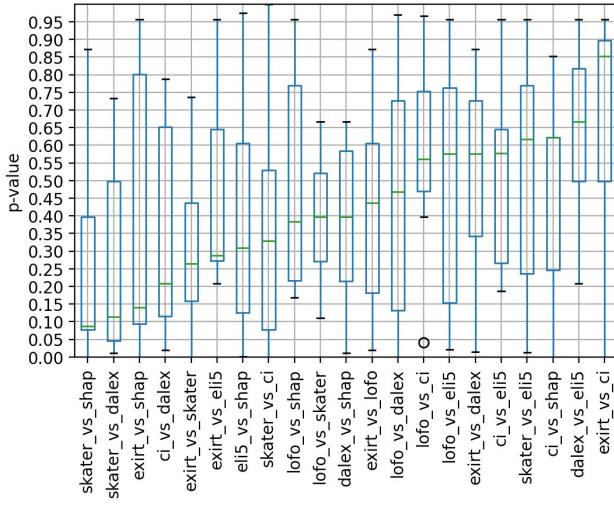


All pair spearman comparisons of feature relevance ranks for models (M1 to M4) based on cluster 0

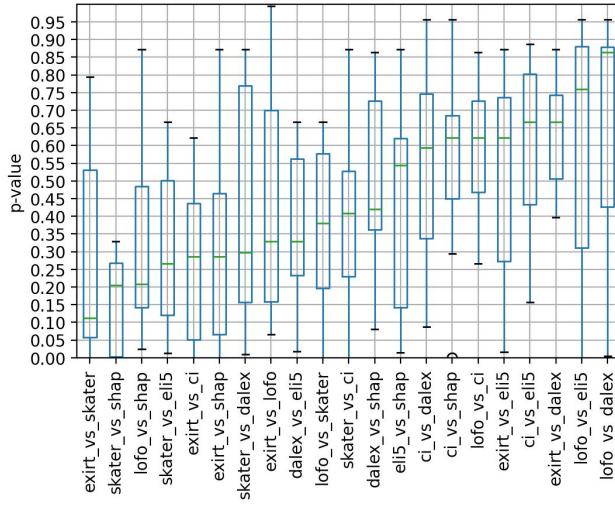


All pair p -values comparisons of feature relevance ranks for models (M1 to M4) based on cluster 0

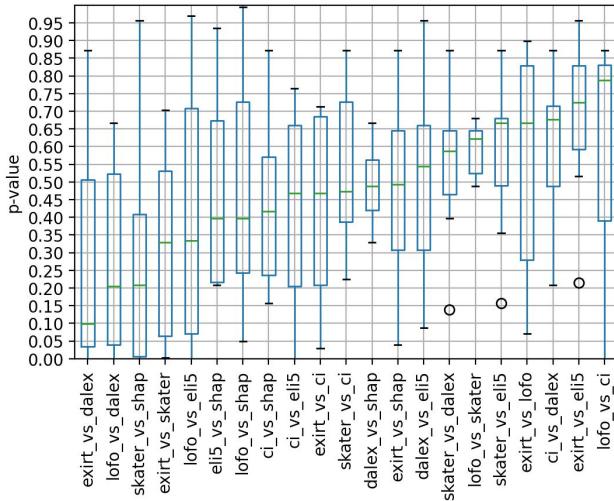
C0 and M1



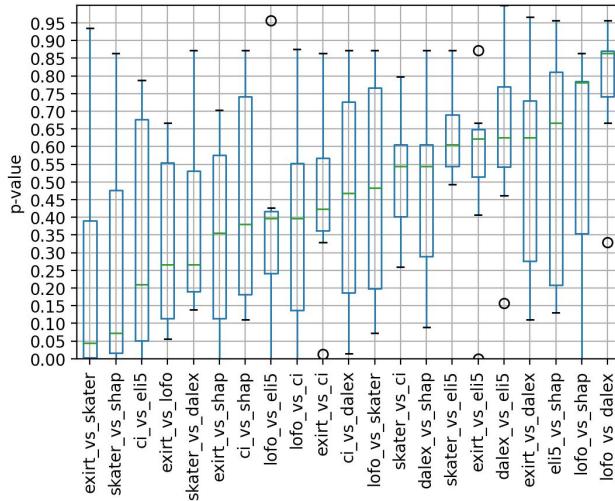
C0 and M2



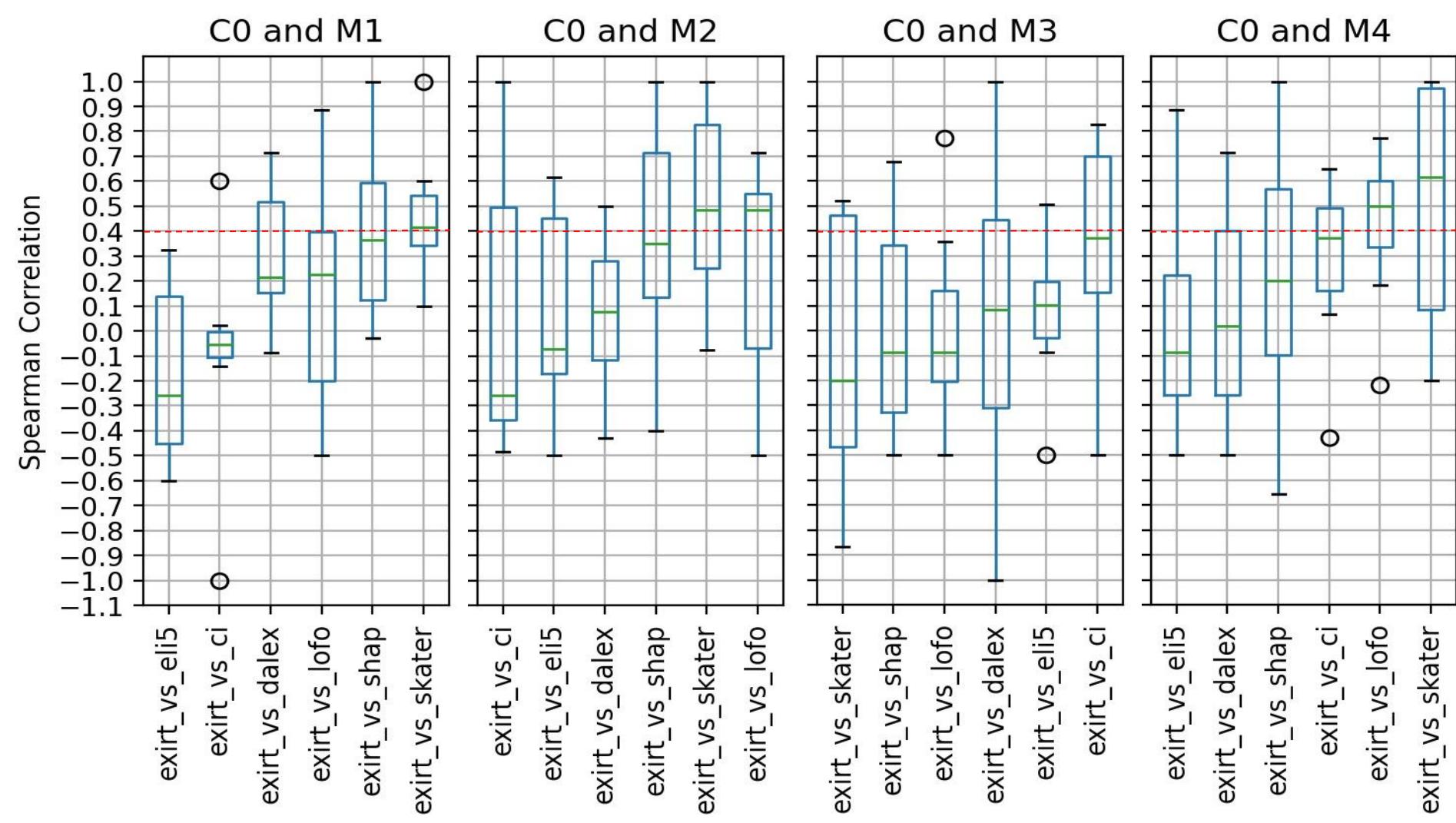
C0 and M3



C0 and M4

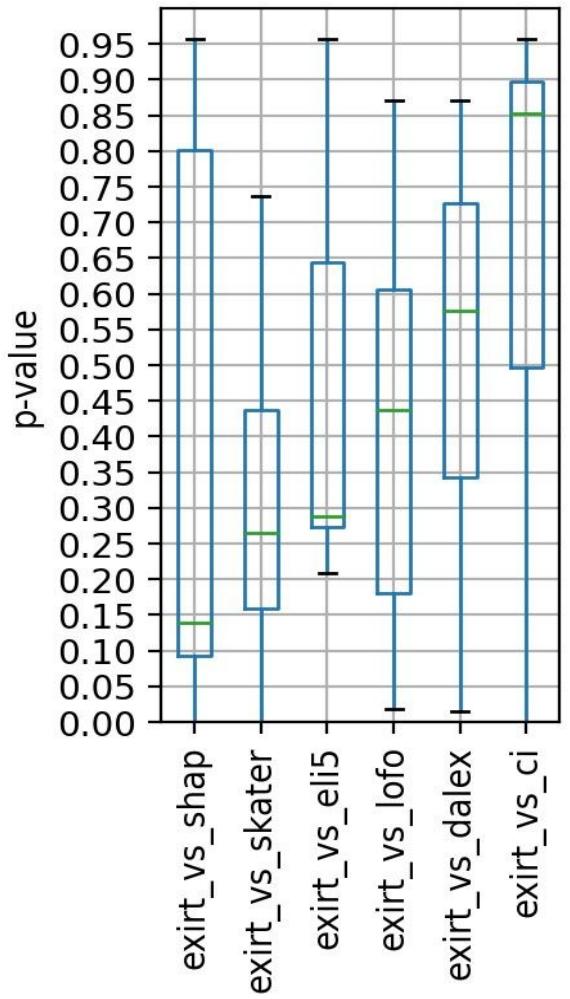


Only eXirt spearman comparisons of feature relevance ranks for models (M1 to M4) based on cluster 0

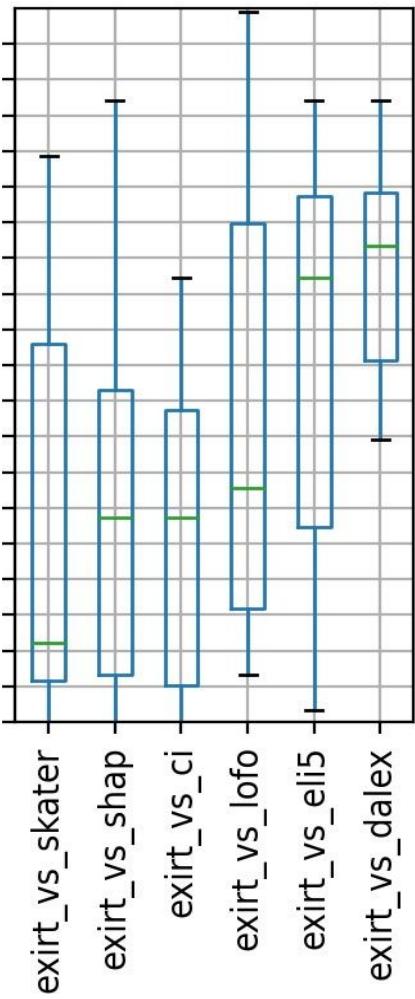


Only *eXirt p-values* comparisons of feature relevance ranks for models (M1 to M4) based on cluster 0

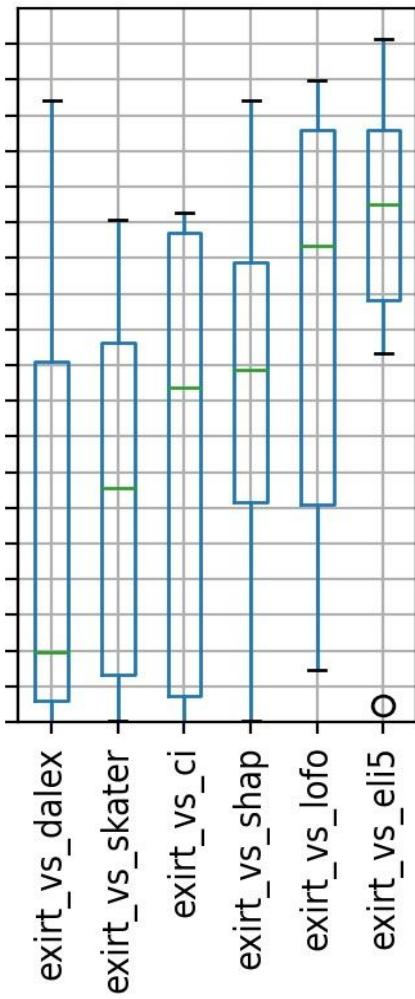
C0 and M1



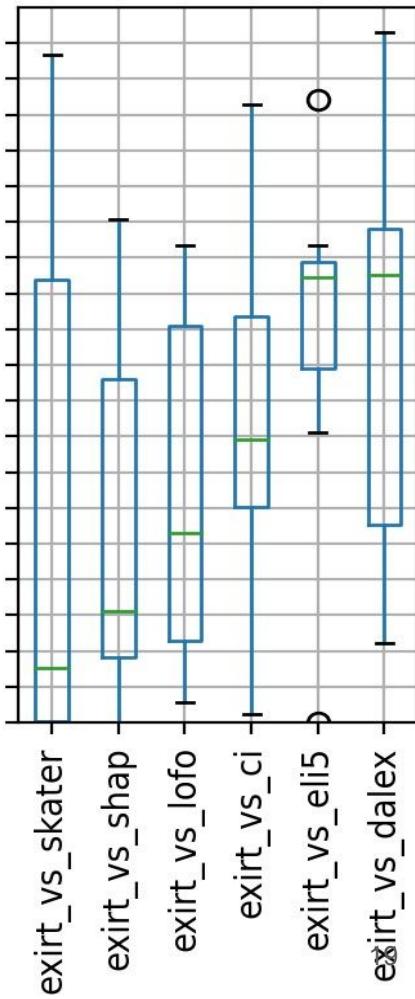
C0 and M2



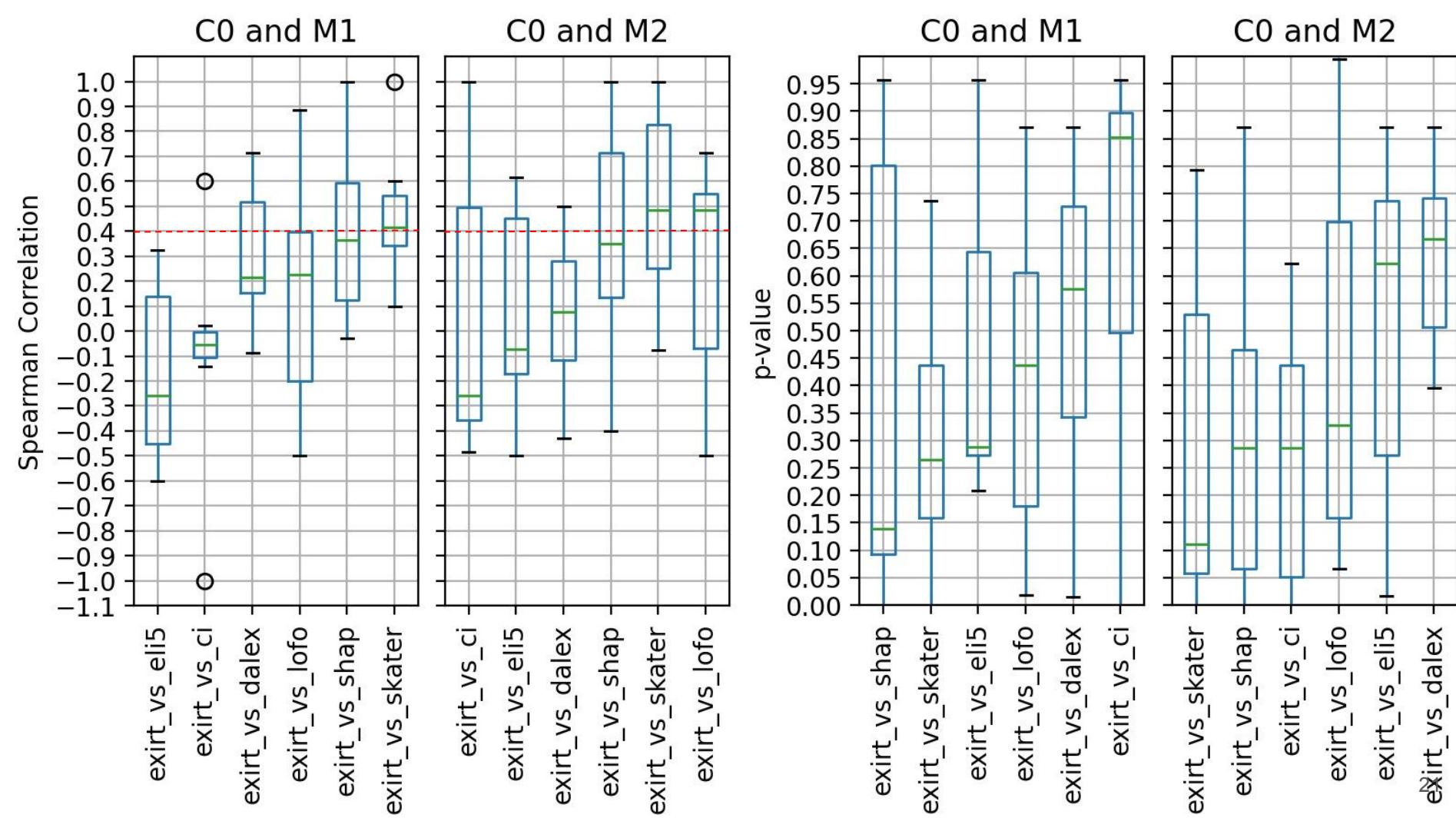
C0 and M3



C0 and M4



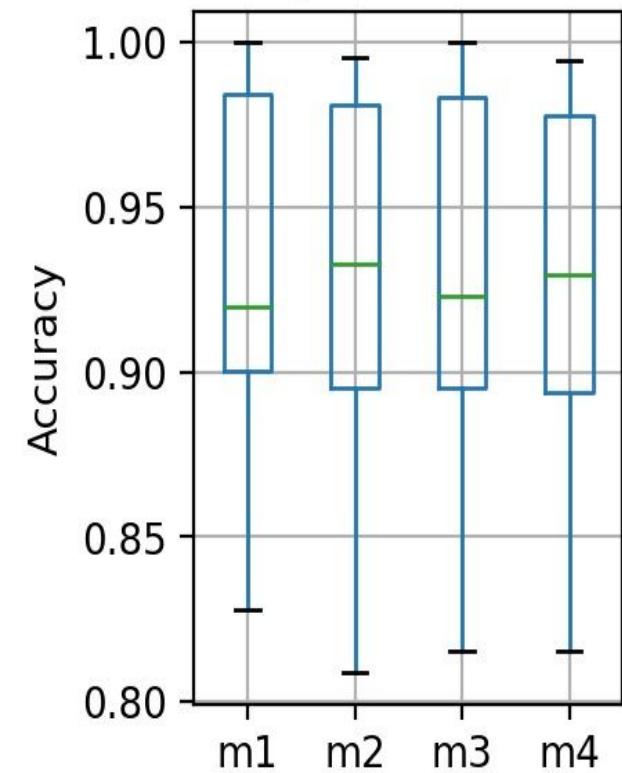
**Only *eXirt* comparisons of feature relevance ranks for models (M1 and M2)
based on cluster 0**



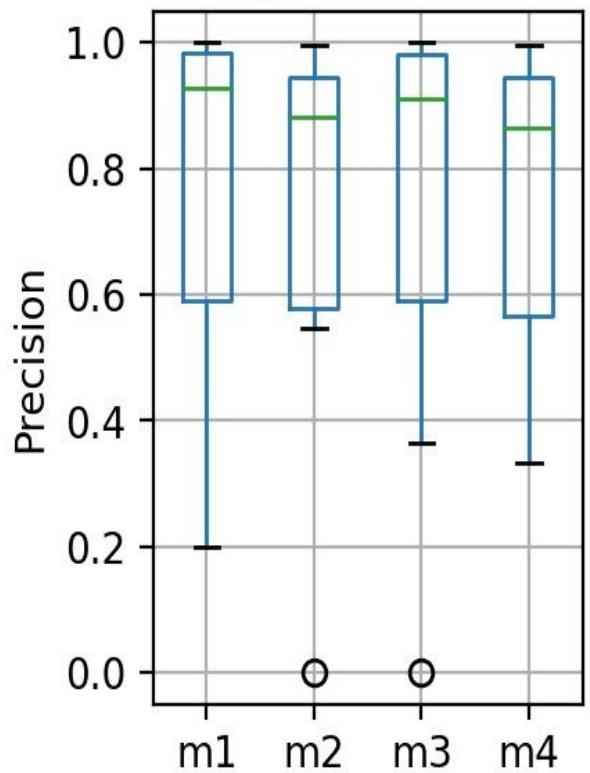
Cluster 1

Data regarding the performance of models M1 to M4 based on cluster 1

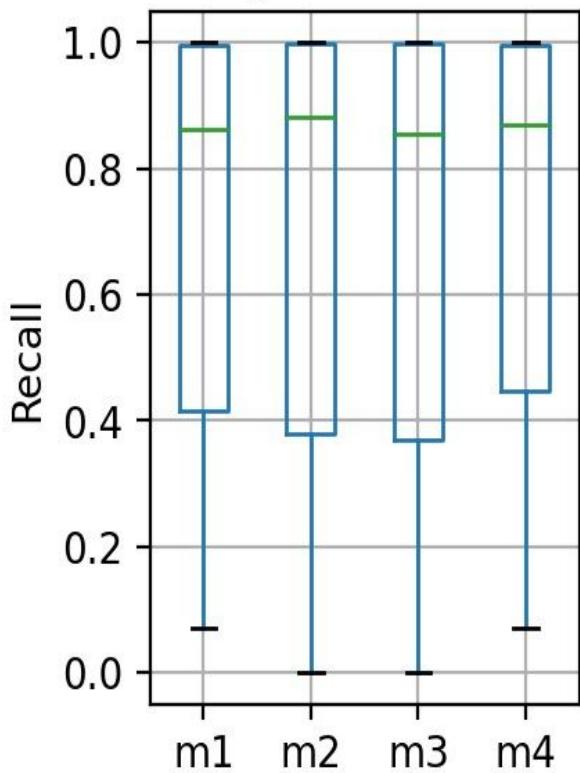
Accuracy of models
by cluster: 1



Precision of models
by cluster: 1

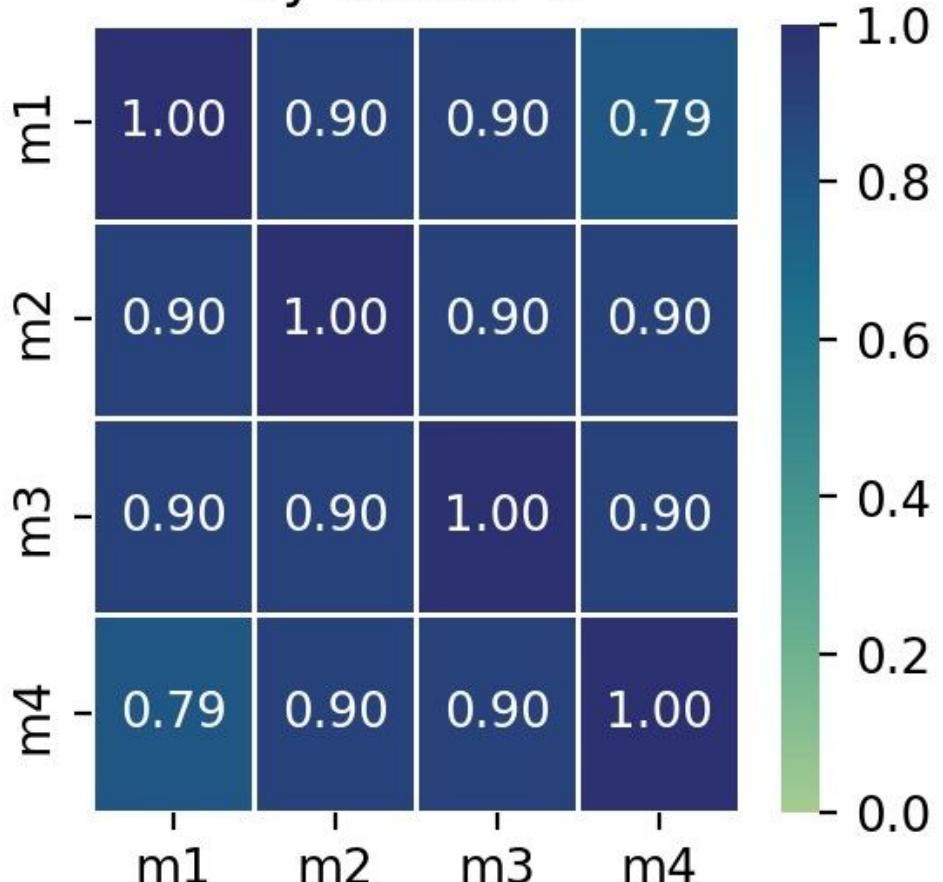


Recall of models
by cluster: 1

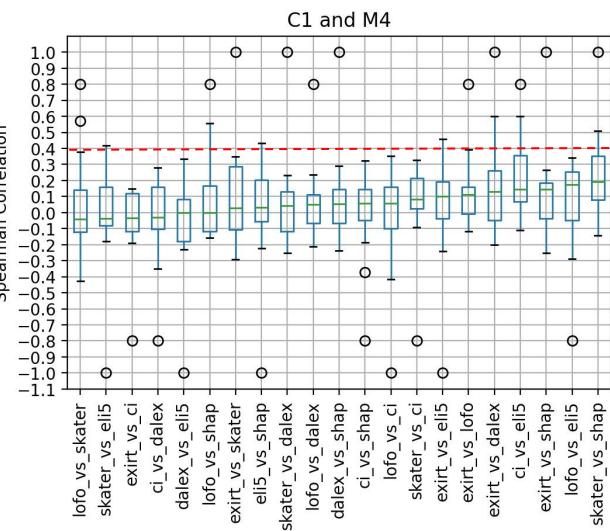
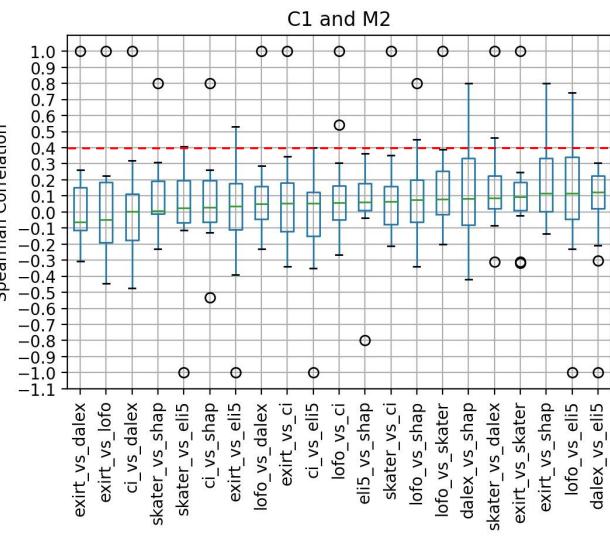
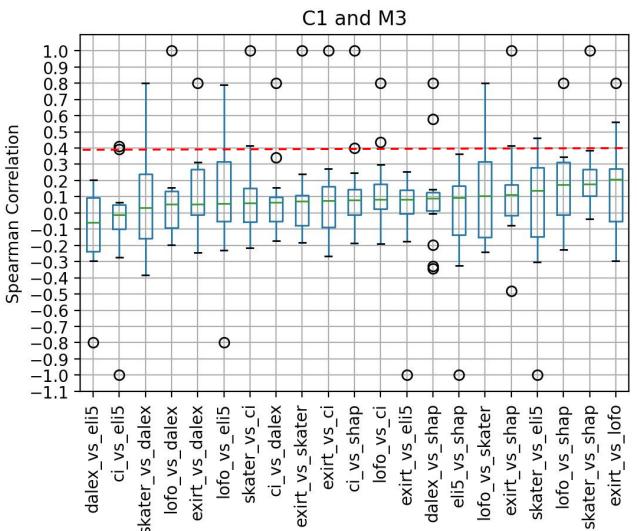
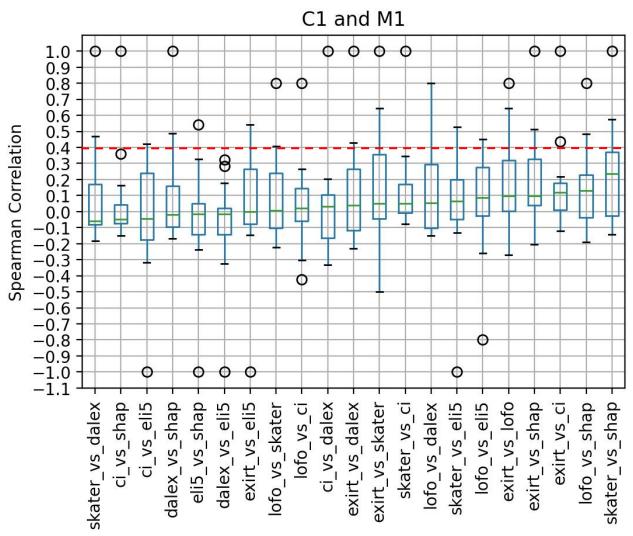


Friedman test result - Models by Cluster 1

Friedman test of models by cluster 1

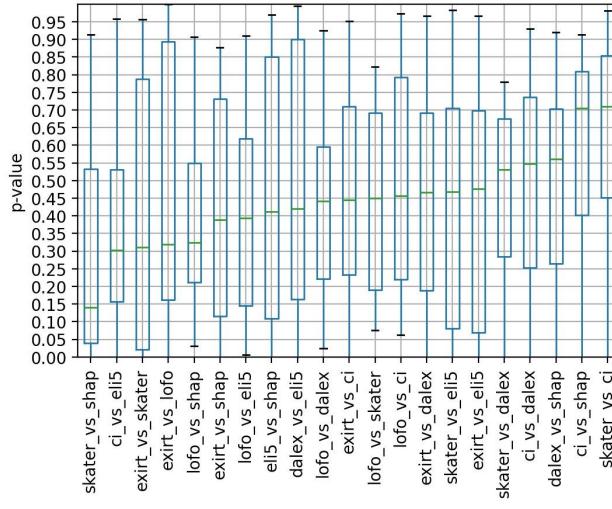


All pair spearman comparisons of feature relevance ranks for models (M1 to M4) based on cluster 1

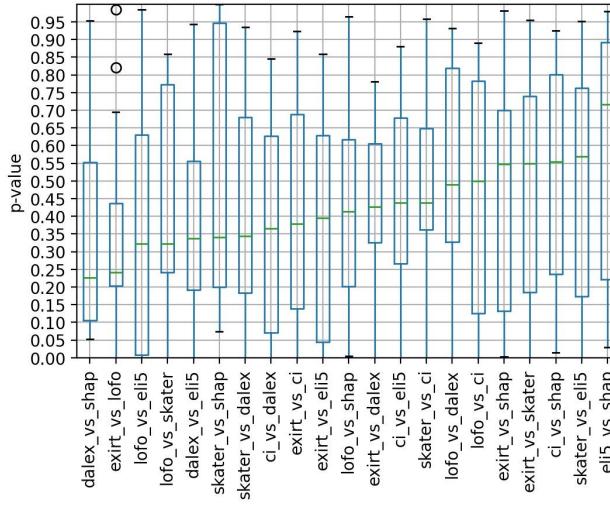


All pair p -values comparisons of feature relevance ranks for models (M1 to M4) based on cluster 1

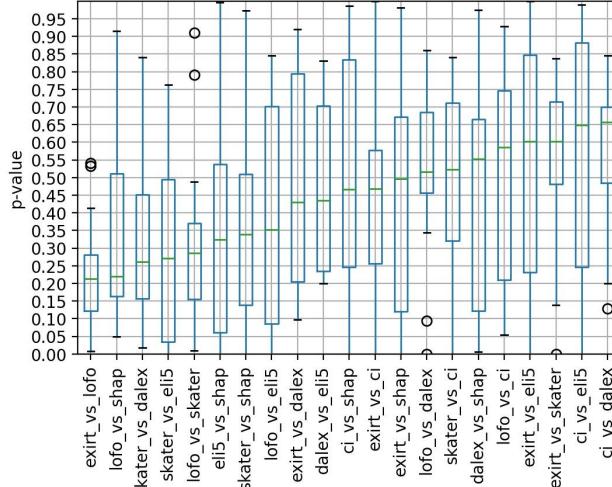
C1 and M1



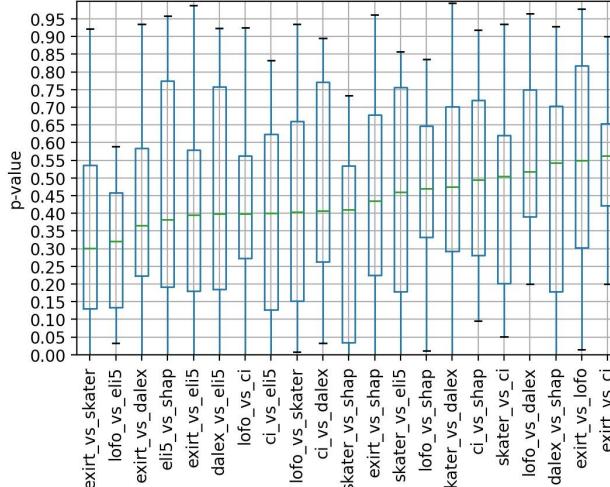
C1 and M2



C1 and M3

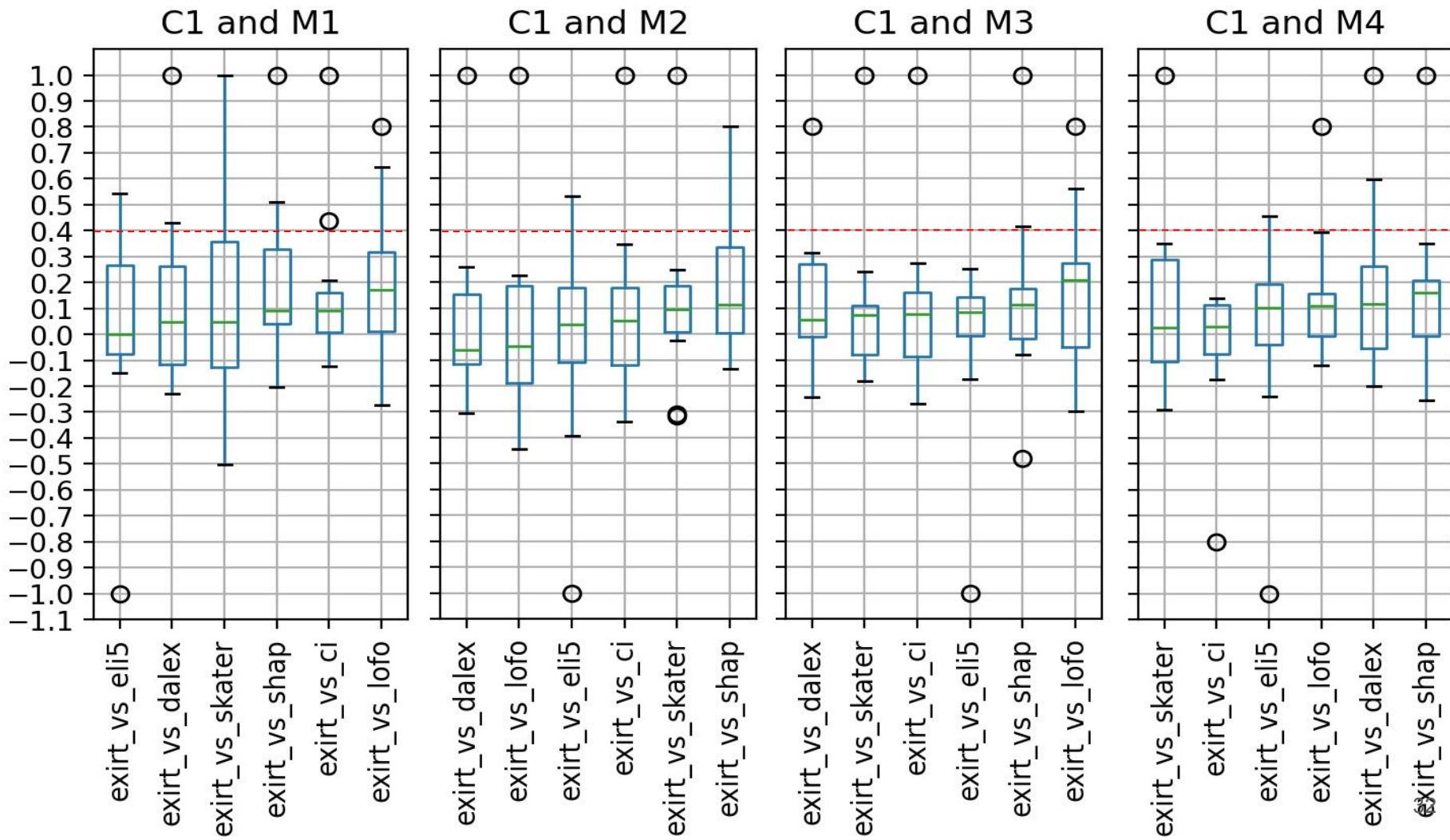


C1 and M4



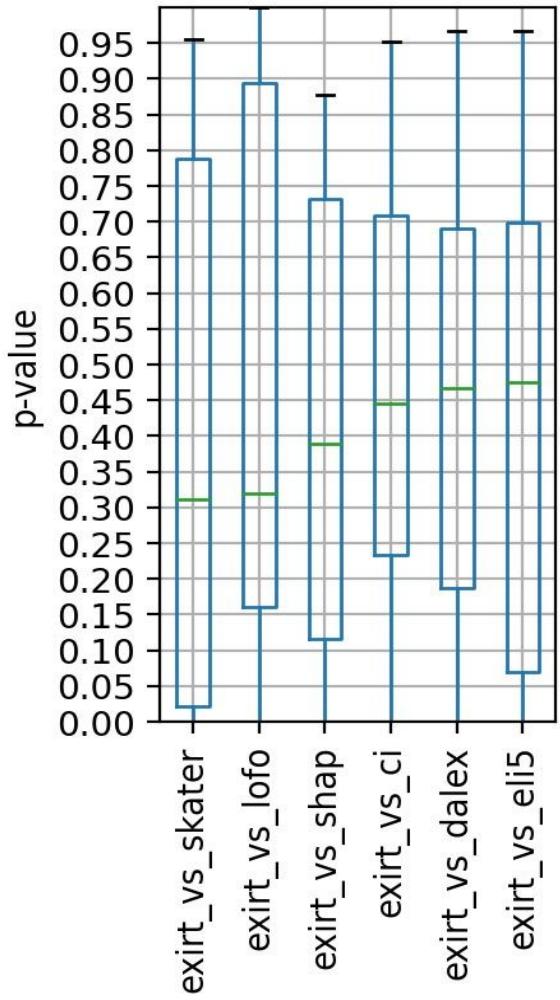
Only eXirt spearman comparisons of feature relevance ranks for models (M1 to M4) based on cluster 1

Spearman Correlation

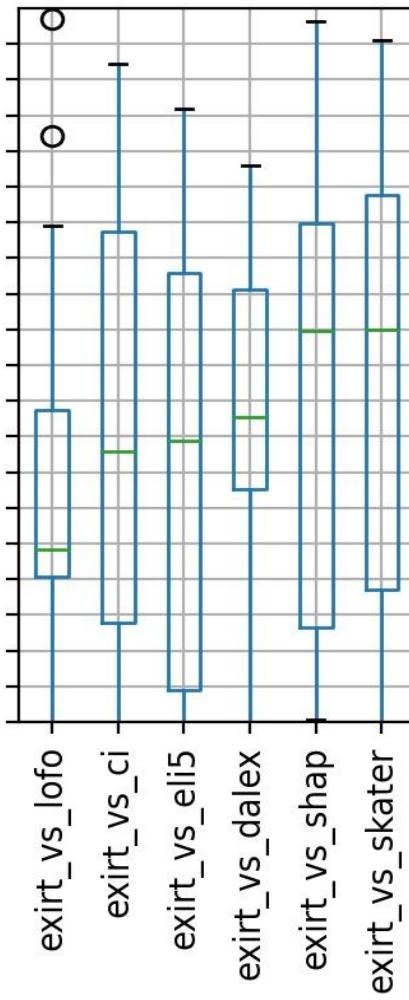


Only eXirt *p*-values comparisons of feature relevance ranks for models (M1 to M4) based on cluster 1

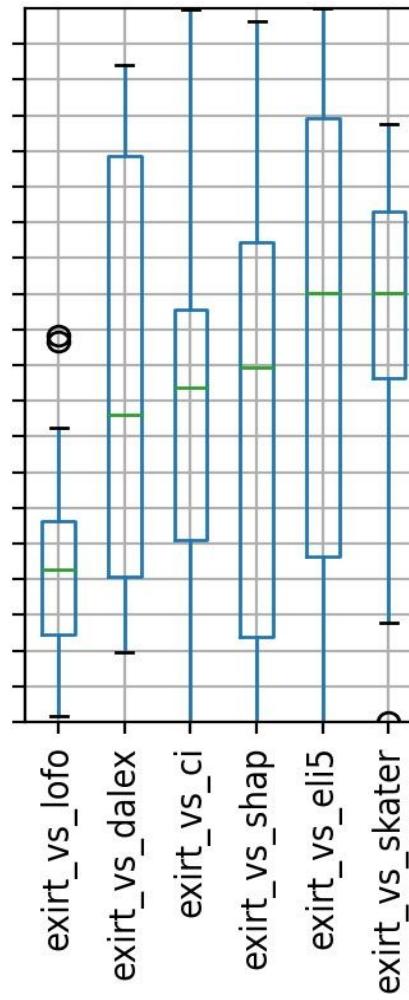
C1 and M1



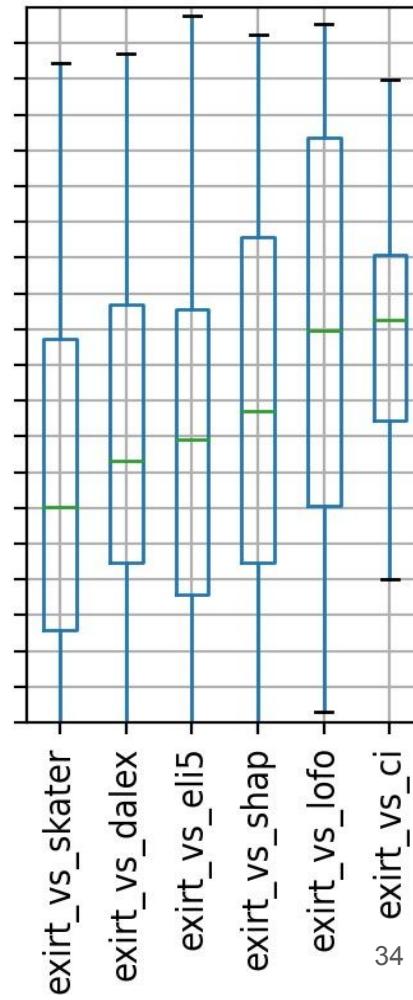
C1 and M2



C1 and M3



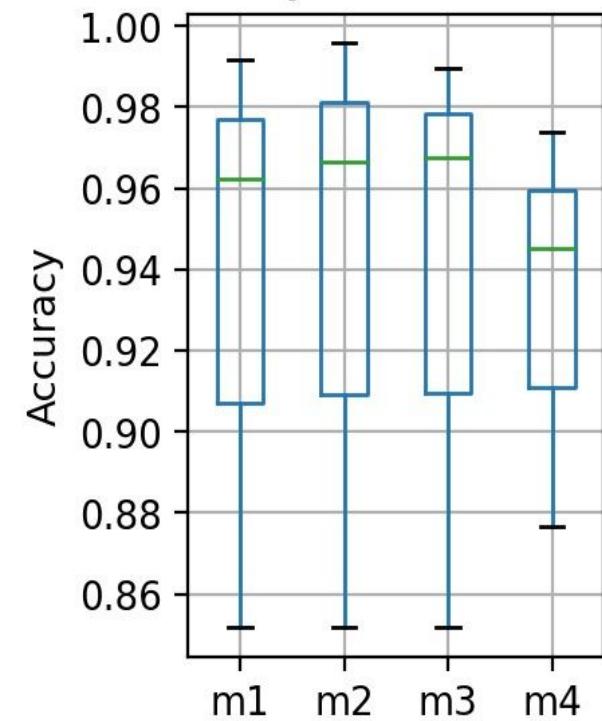
C1 and M4



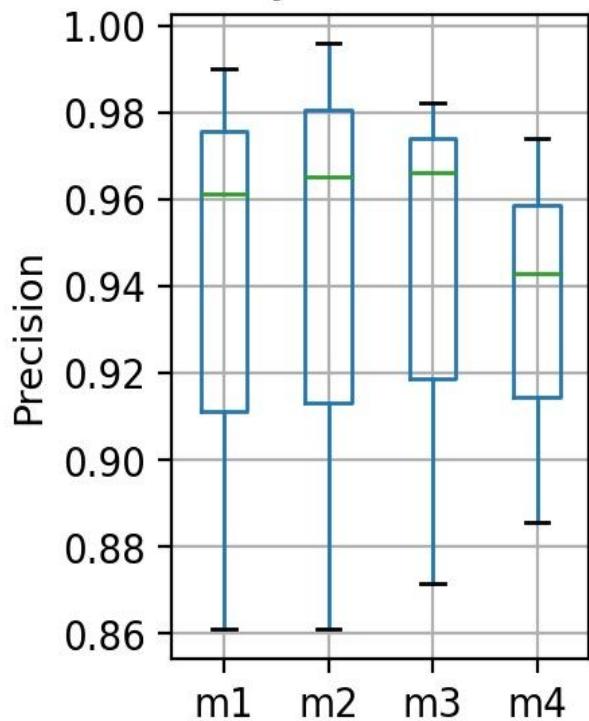
Cluster 2

Data regarding the performance of models M1 to M4 based on cluster 2

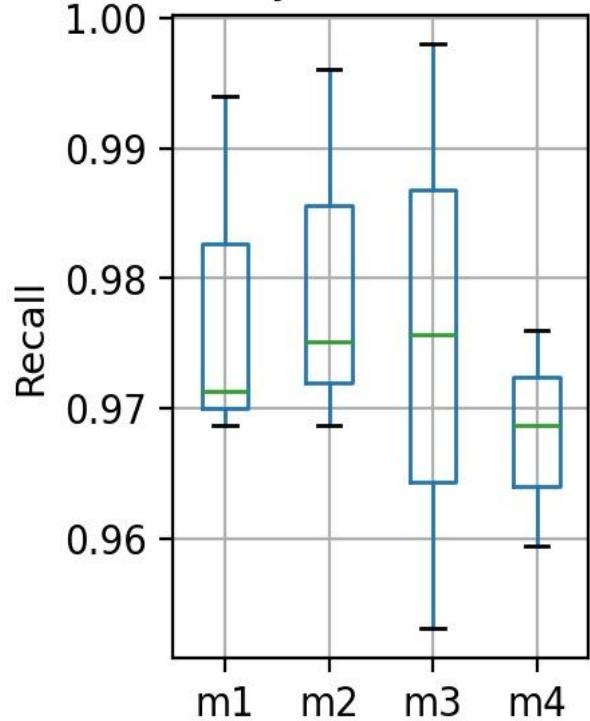
Accuracy of models
by cluster: 2



Precision of models
by cluster: 2

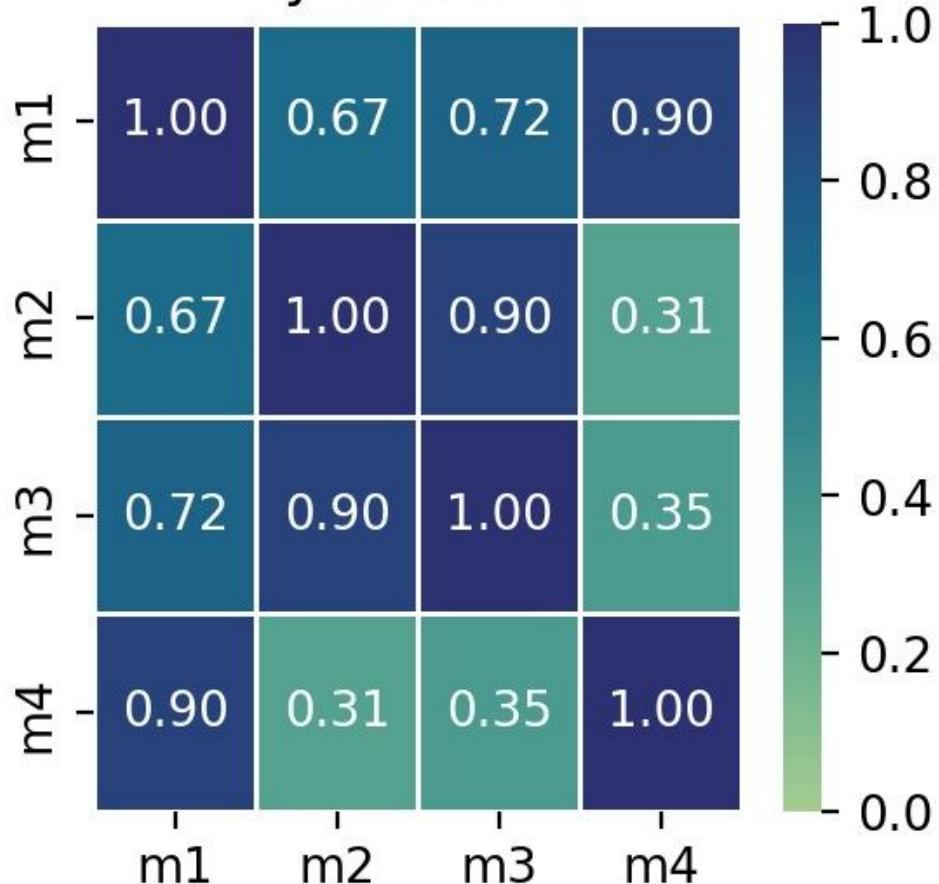


Recall of models
by cluster: 2



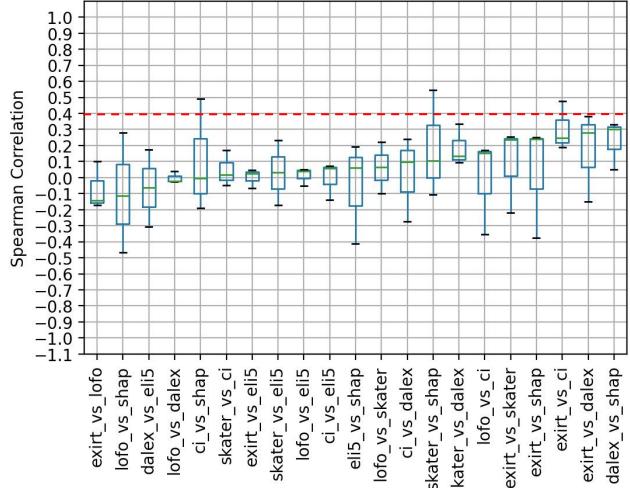
Friedman test result - Models by Cluster 2

Friedman test of models by cluster 2

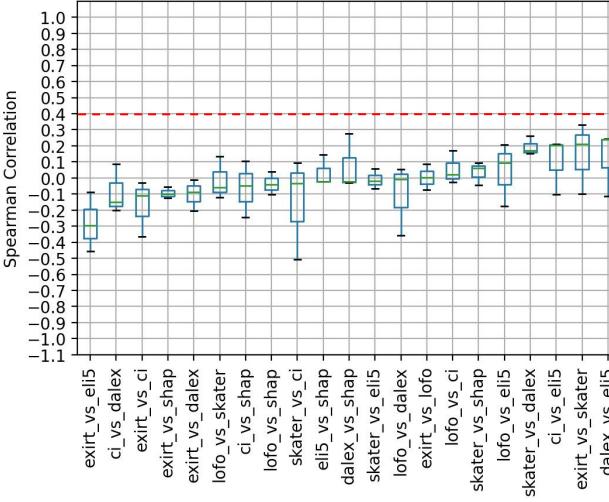


All pair spearman comparisons of feature relevance ranks for models (M1 to M4) based on cluster 2

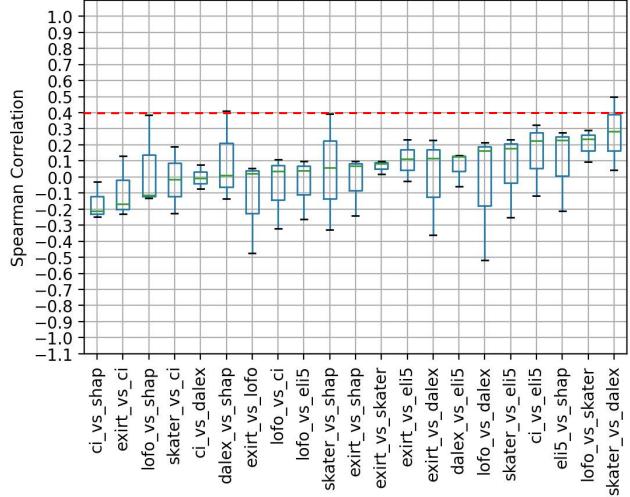
C2 and M1



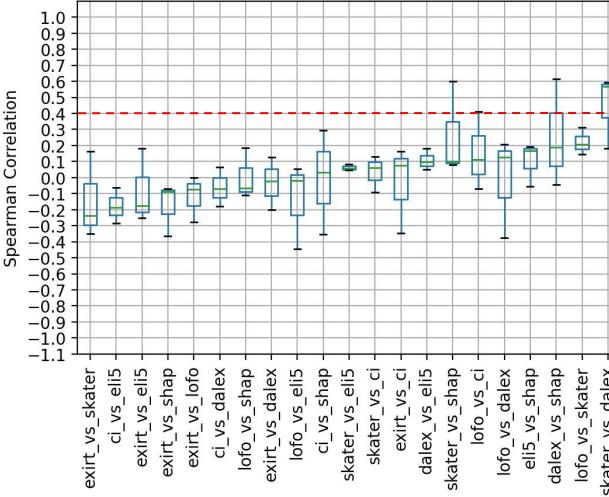
C2 and M2



C2 and M3

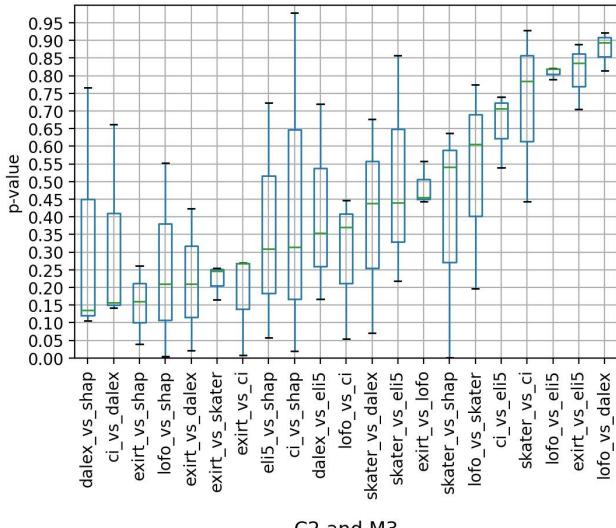


C2 and M4

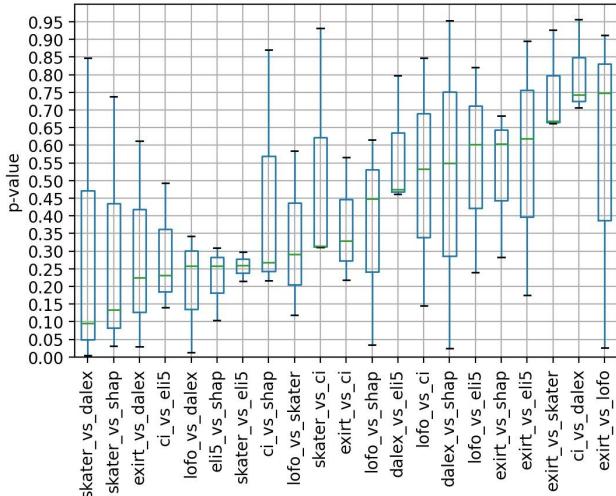


All pair p -values comparisons of feature relevance ranks for models (M1 to M4) based on cluster 2

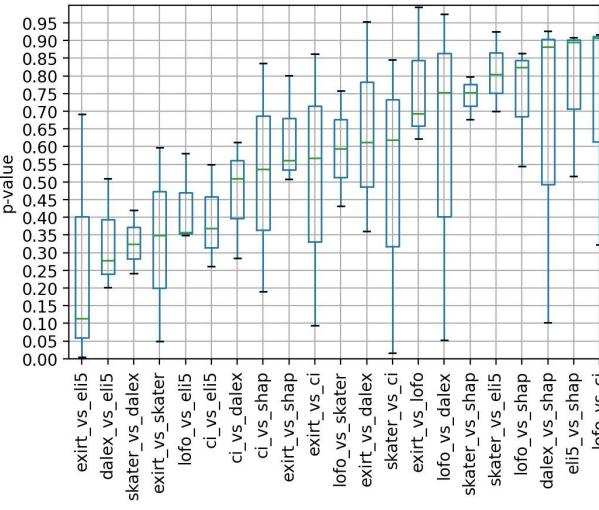
C2 and M1



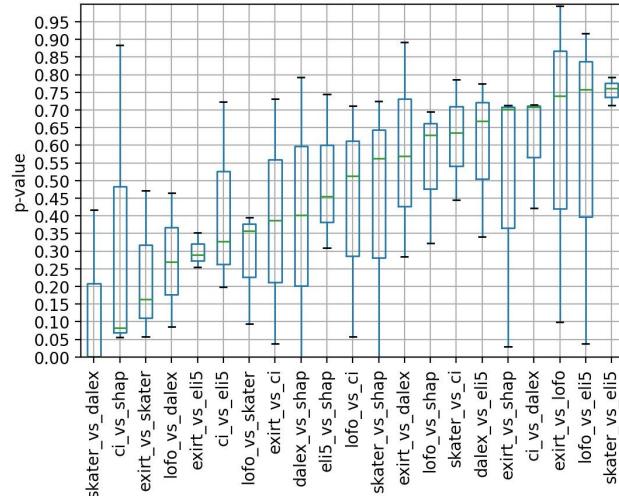
C2 and M3



C2 and M2

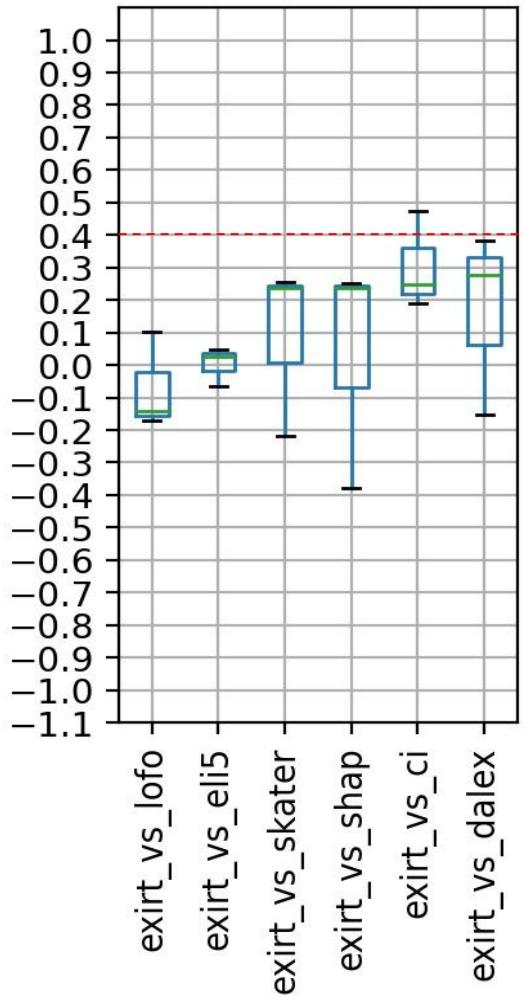


C2 and M4

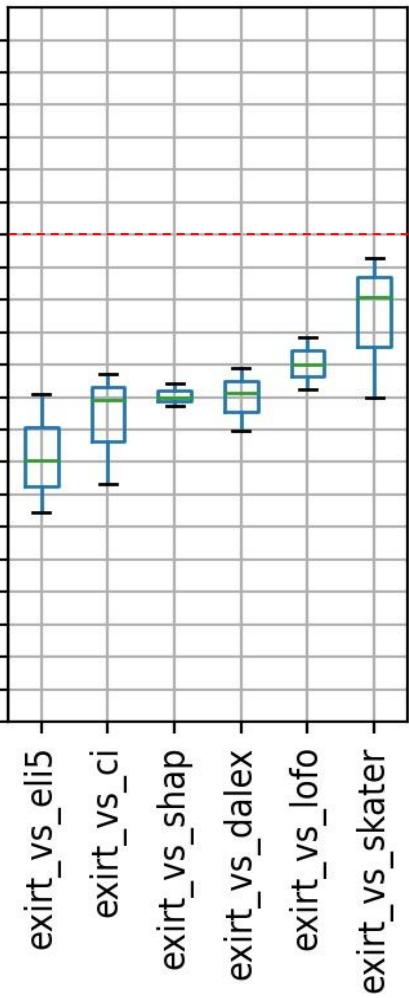


Only eXirt spearman comparisons of feature relevance ranks for models (M1 to M4) based on cluster 2

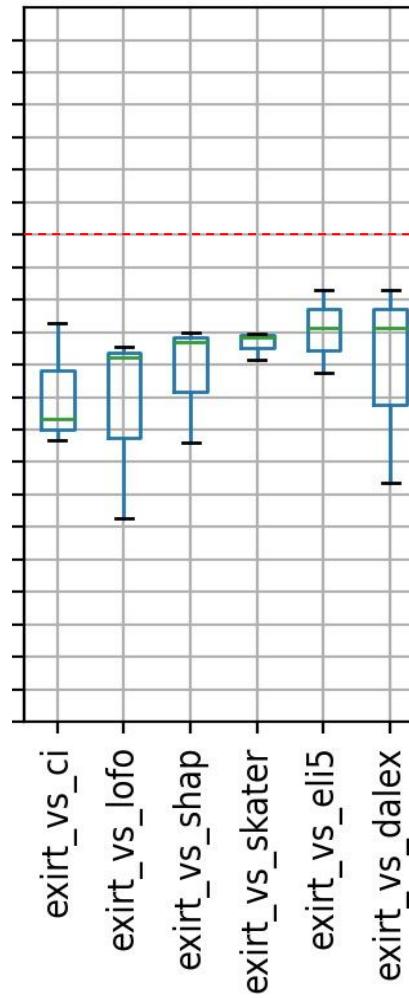
C2 and M1



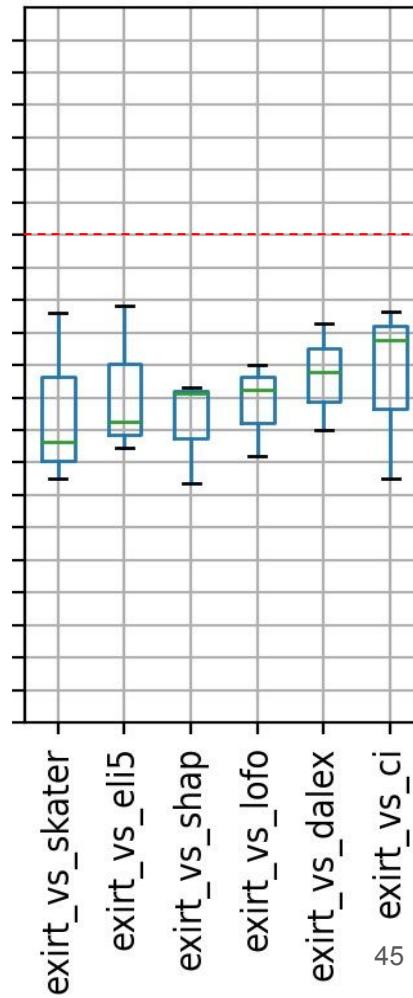
C2 and M2



C2 and M3

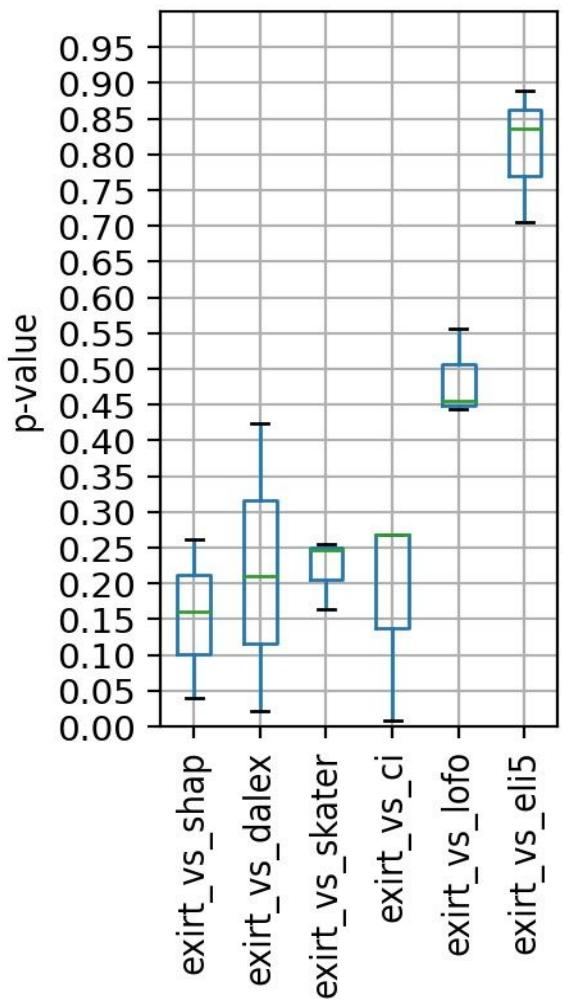


C2 and M4

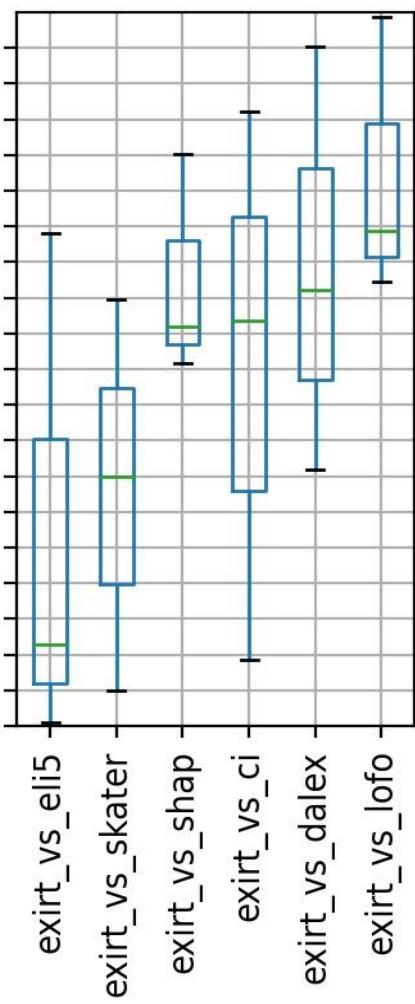


Only eXirt *p*-values comparisons of feature relevance ranks for models (M1 to M4) based on cluster 2

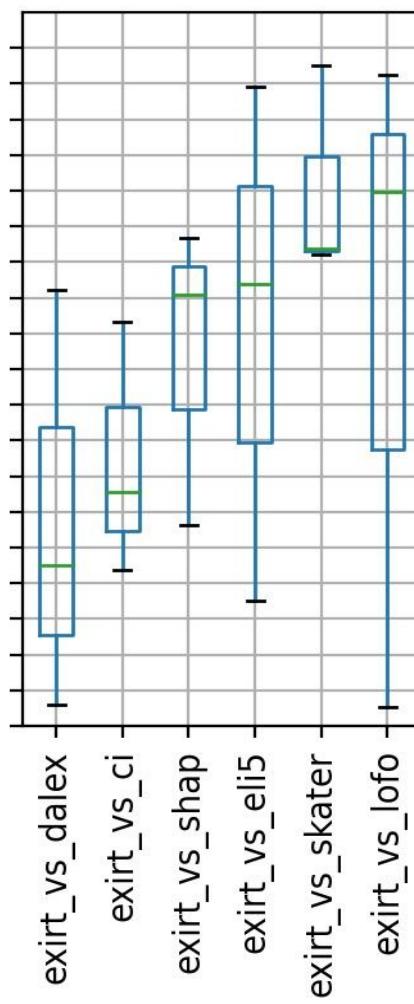
C2 and M1



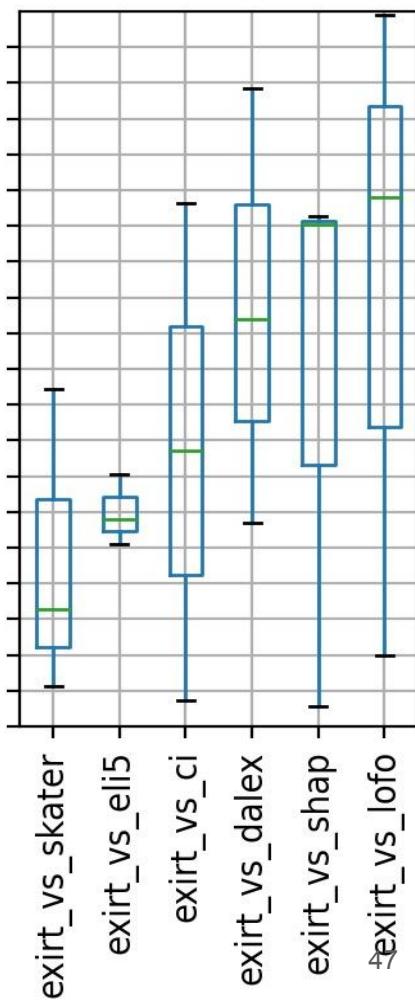
C2 and M2



C2 and M3



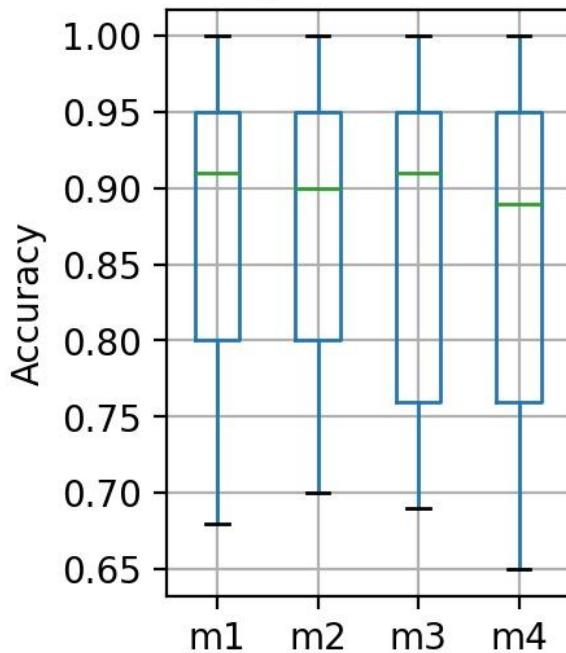
C2 and M4



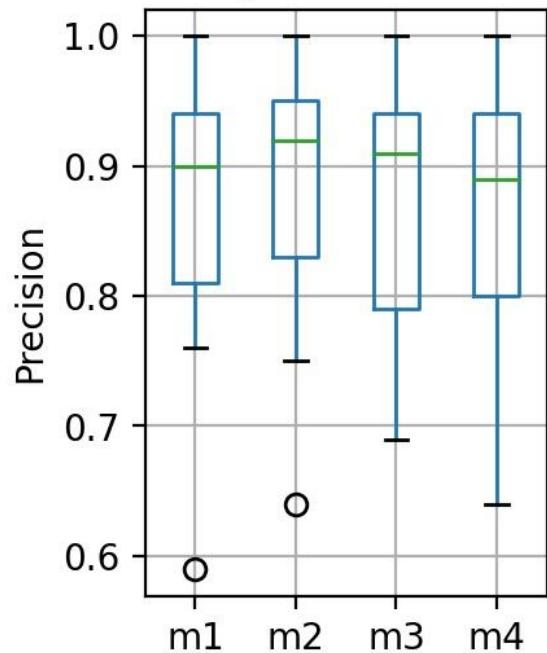
Cluster 3

Data regarding the performance of models M1 to M4 based on cluster 3

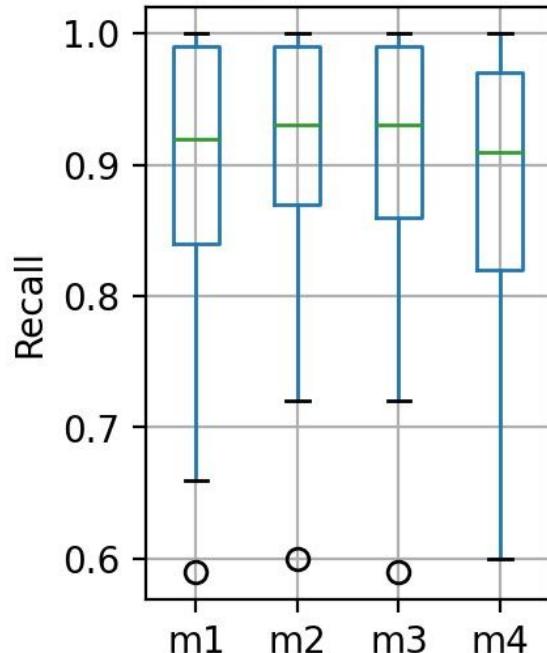
Accuracy of models
by cluster: 3



Precision of models
by cluster: 3

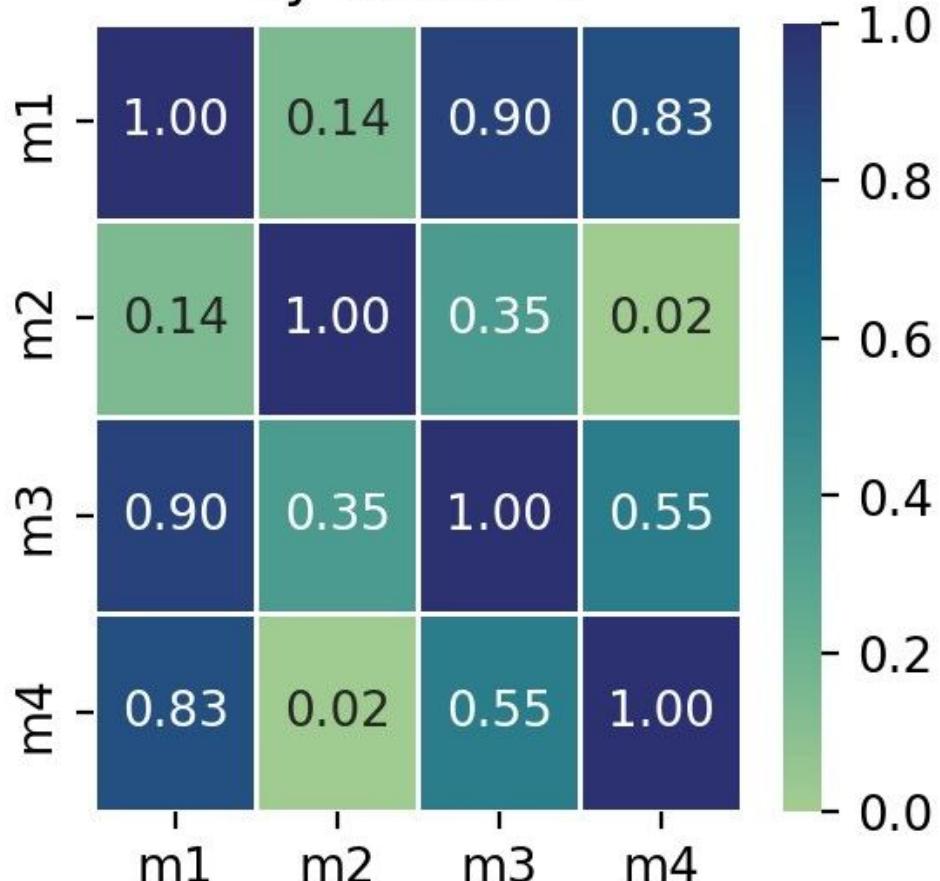


Recall of models
by cluster: 3



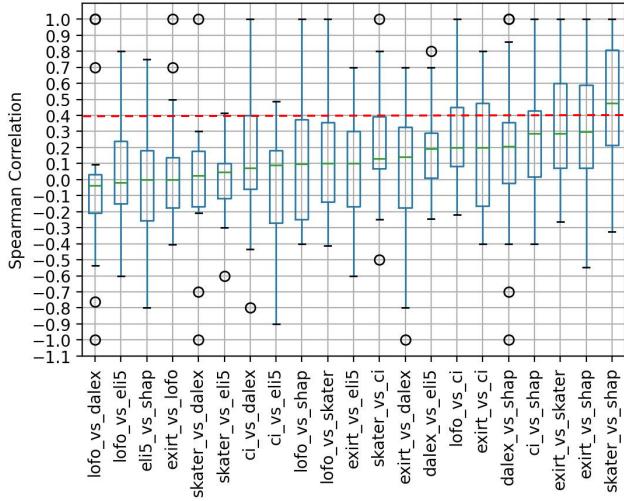
Friedman test result - Models by Cluster 3

Friedman test of models by cluster 3

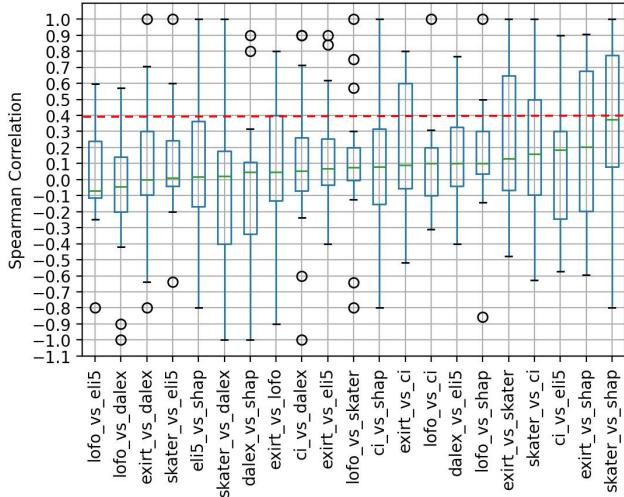


All pair spearman comparisons of feature relevance ranks for models (M1 to M4) based on cluster 3

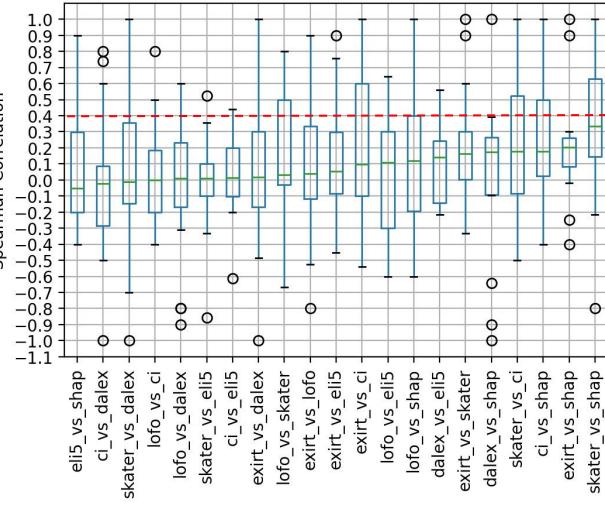
C3 and M1



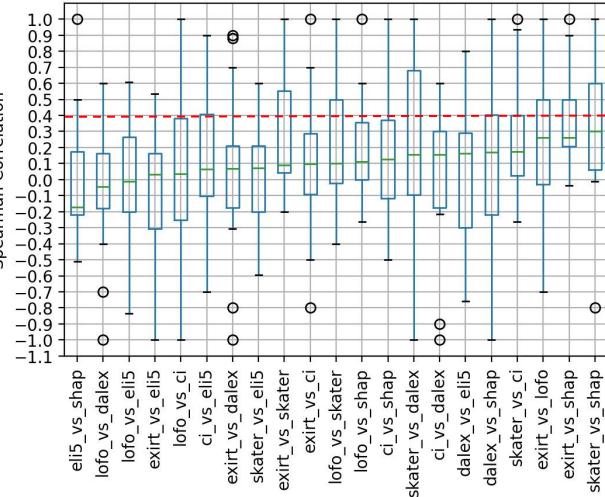
C3 and M3



C3 and M2

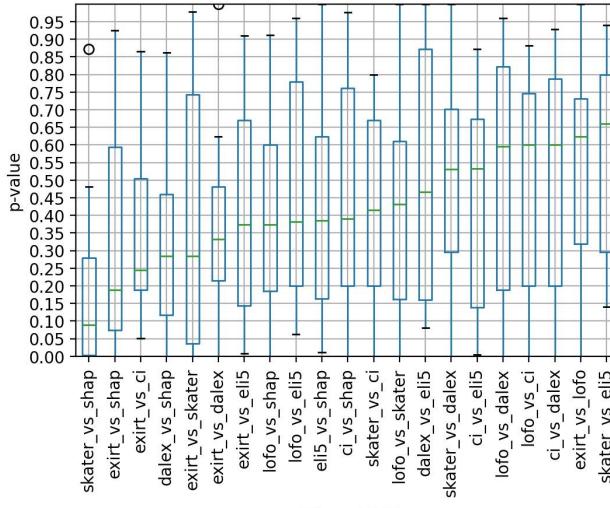


C3 and M4

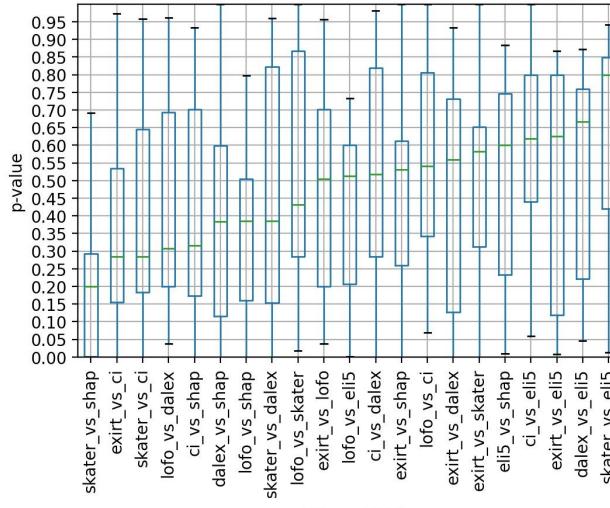


All pair p -values comparisons of feature relevance ranks for models (M1 to M4) based on cluster 3

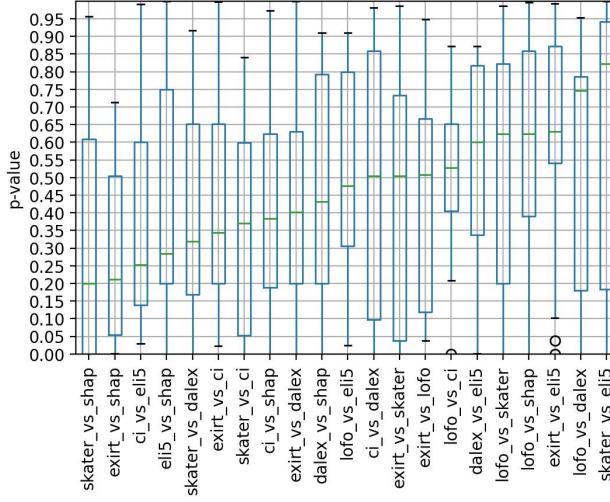
C3 and M1



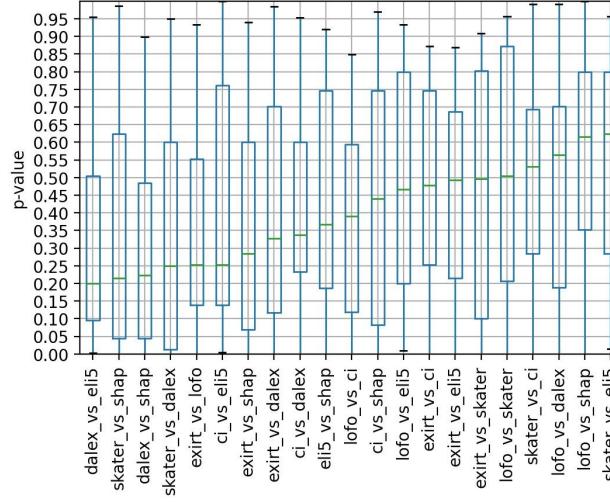
C3 and M2



C3 and M3

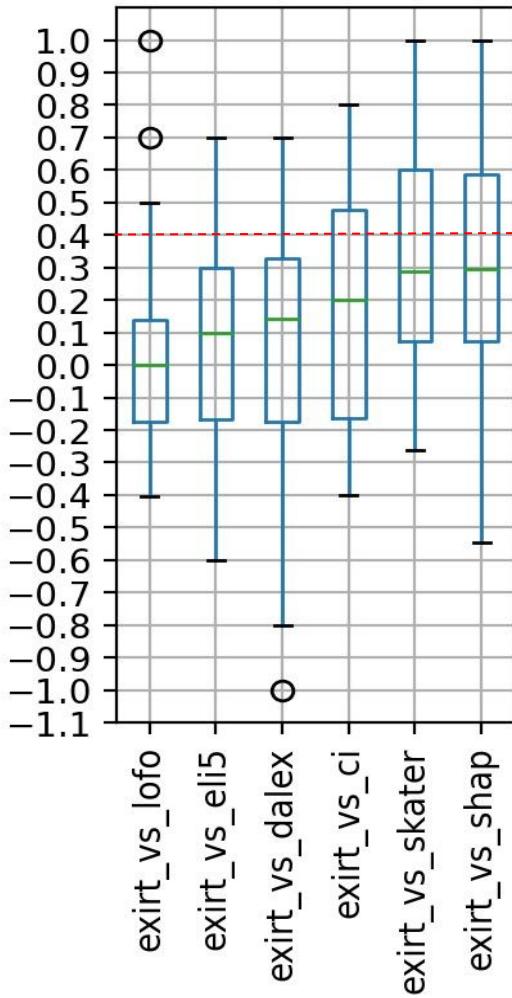


C3 and M4

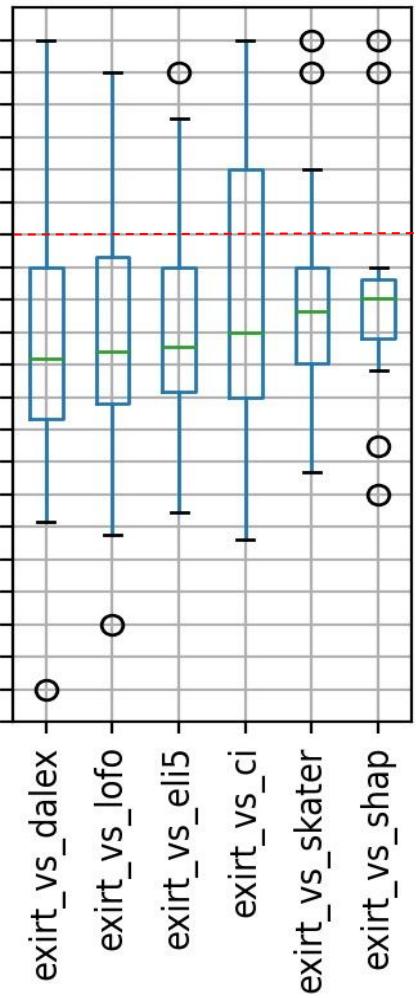


Only eXirt spearman comparisons of feature relevance ranks for models (M1 to M4) based on cluster 3

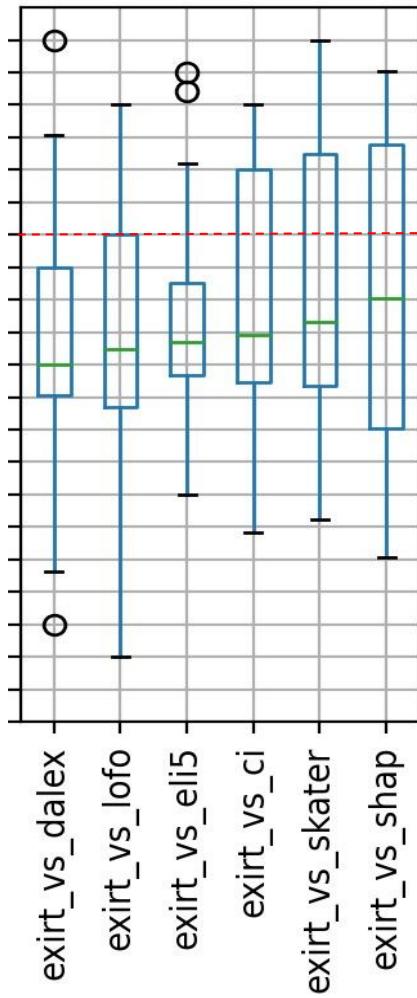
C3 and M1



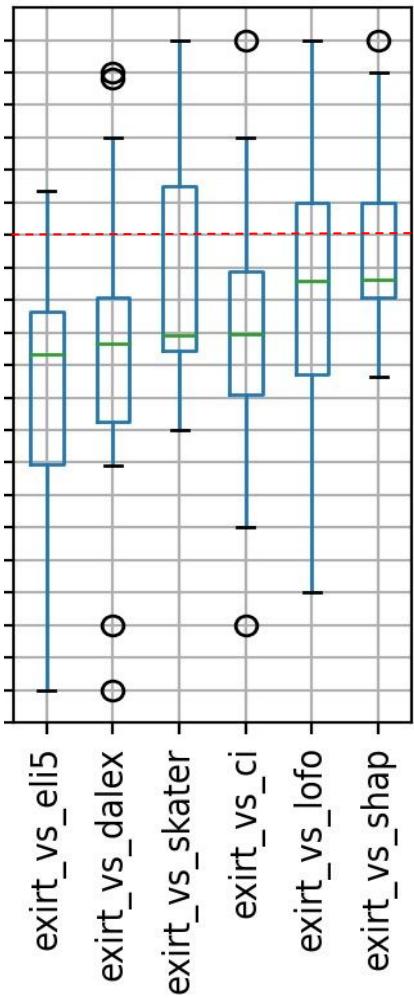
C3 and M2



C3 and M3

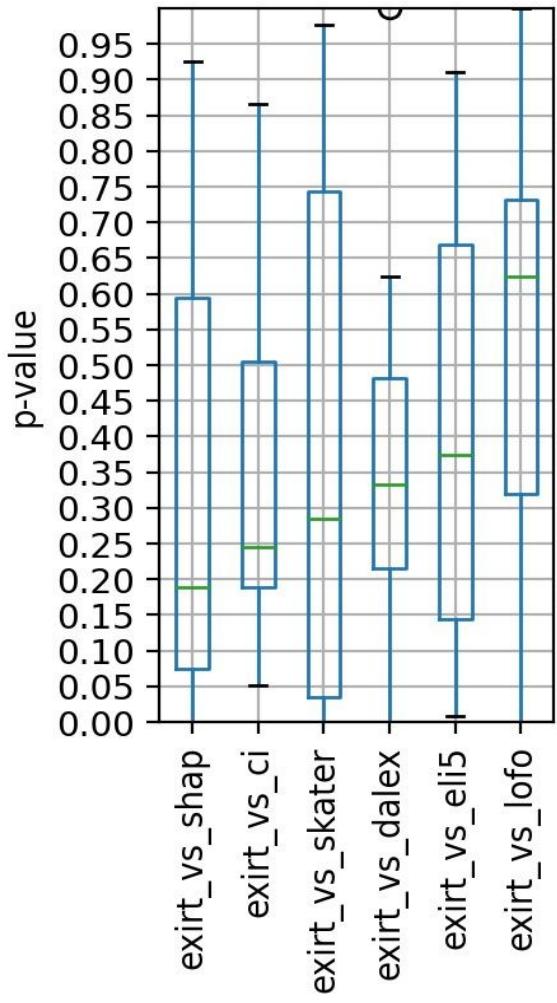


C3 and M4

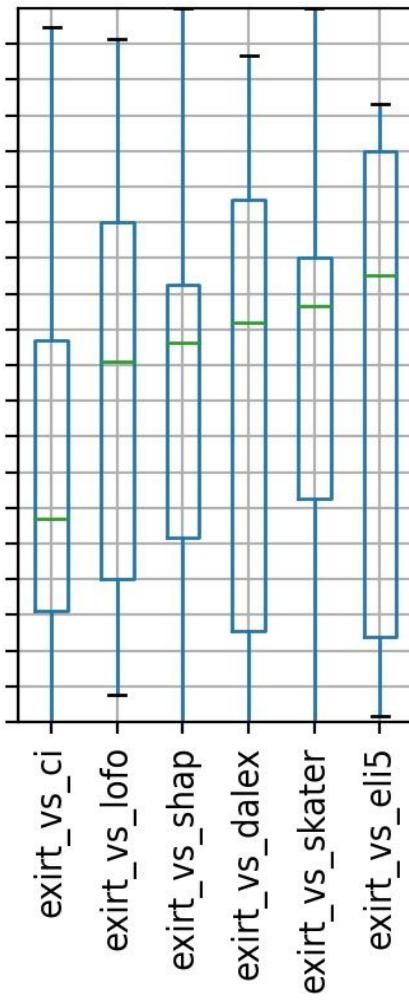


Only eXirt *p*-values comparisons of feature relevance ranks for models (M1 to M4) based on cluster 3

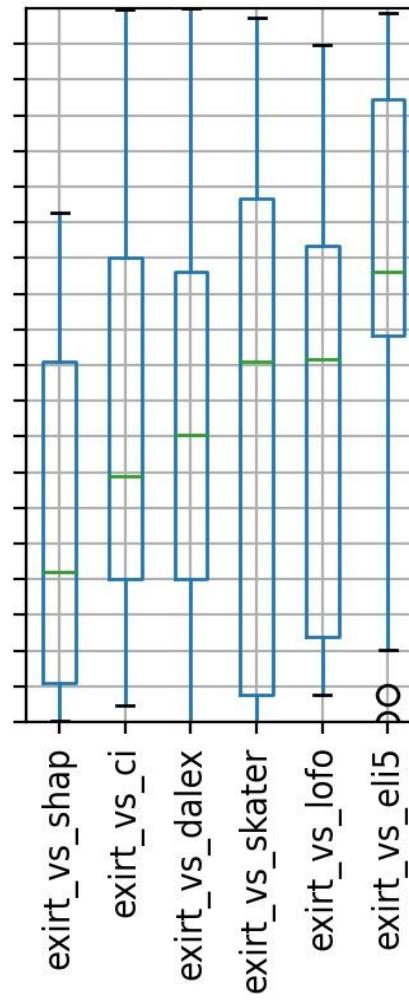
C3 and M1



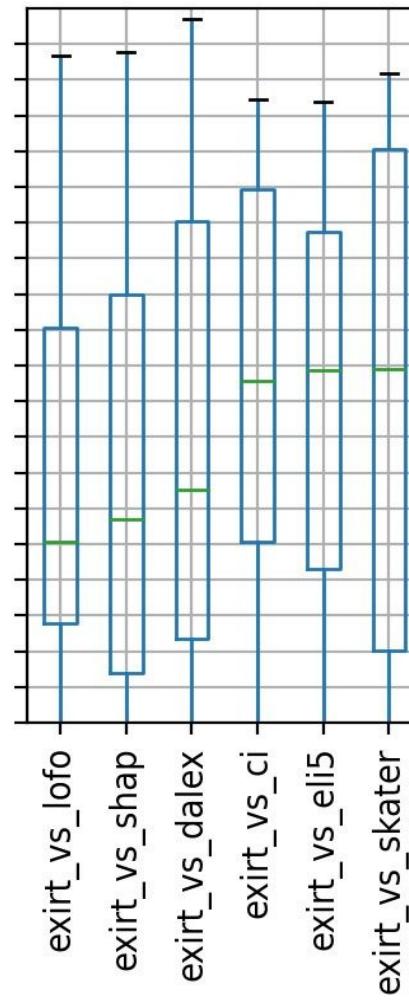
C3 and M2



C3 and M3



C3 and M4



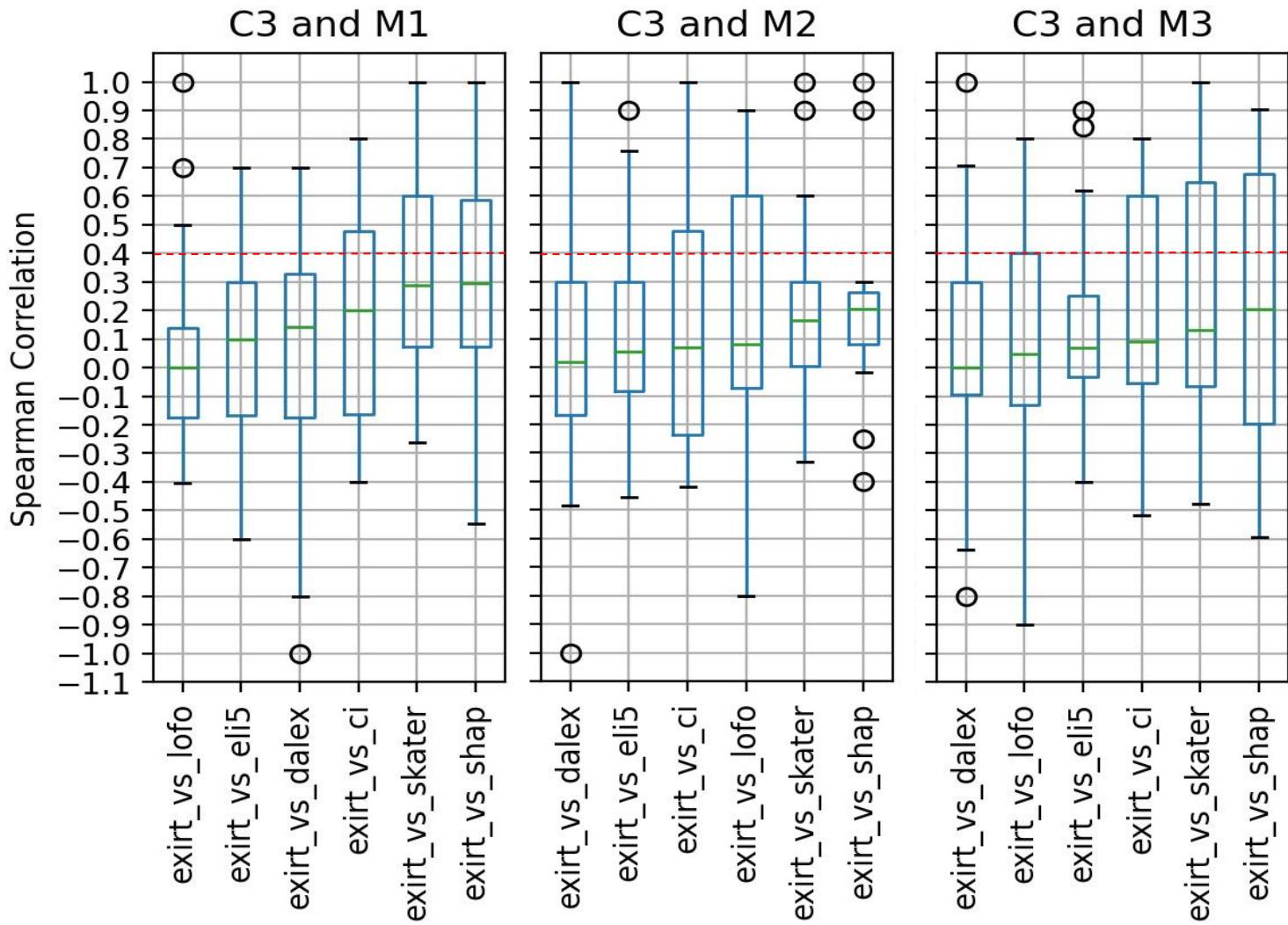
exirt_vs_shap
exirt_vs_ci
exirt_vs_skater
exirt_vs_dalex
exirt_vs_el5
exirt_vs_lofo

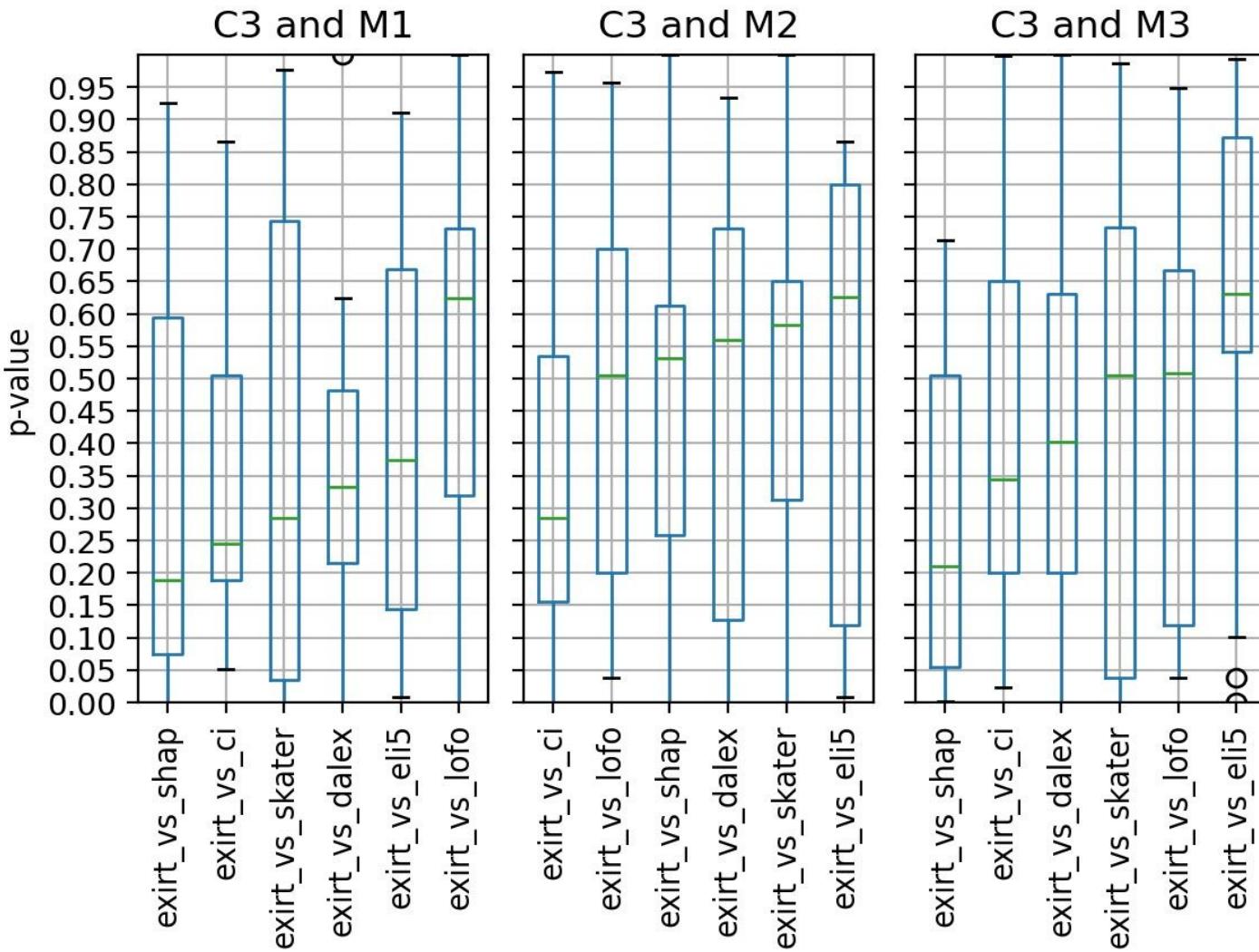
exirt_vs_ci
exirt_vs_lofo
exirt_vs_shap
exirt_vs_dalex
exirt_vs_skater
exirt_vs_el5

exirt_vs_shap
exirt_vs_ci
exirt_vs_dalex
exirt_vs_skater
exirt_vs_lofo
exirt_vs_el5

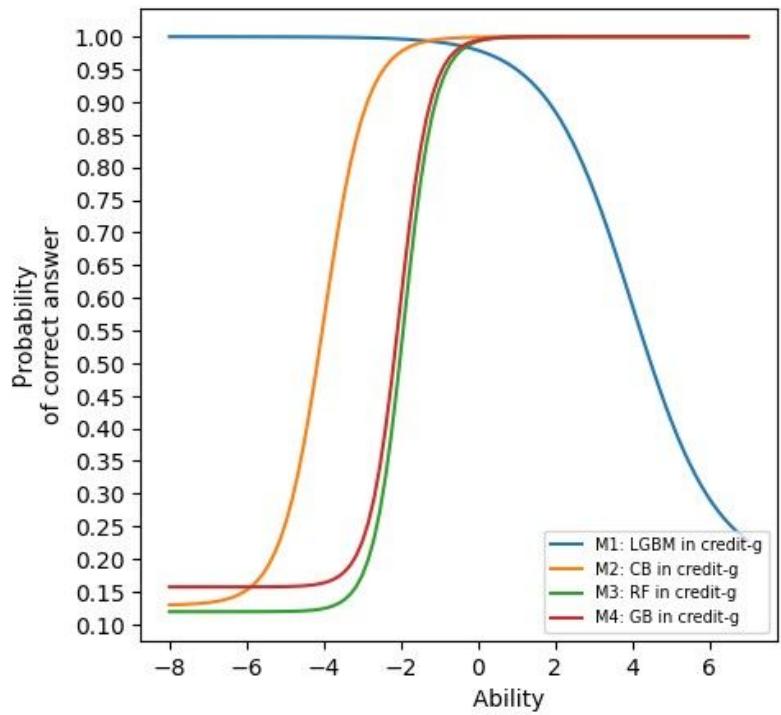
exirt_vs_lofo
exirt_vs_shap
exirt_vs_dalex
exirt_vs_ci
exirt_vs_el5
exirt_vs_skater

Only *eXirt* comparisons of feature relevance ranks for models (M1, M2, and M3) based on cluster 3

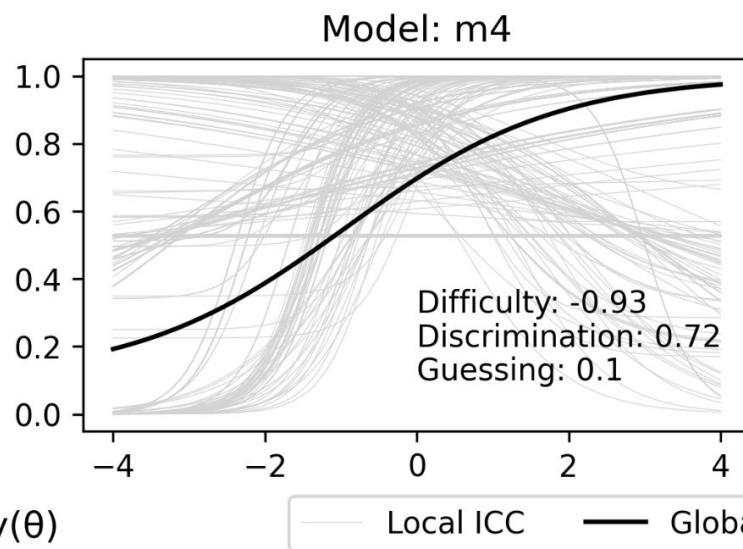
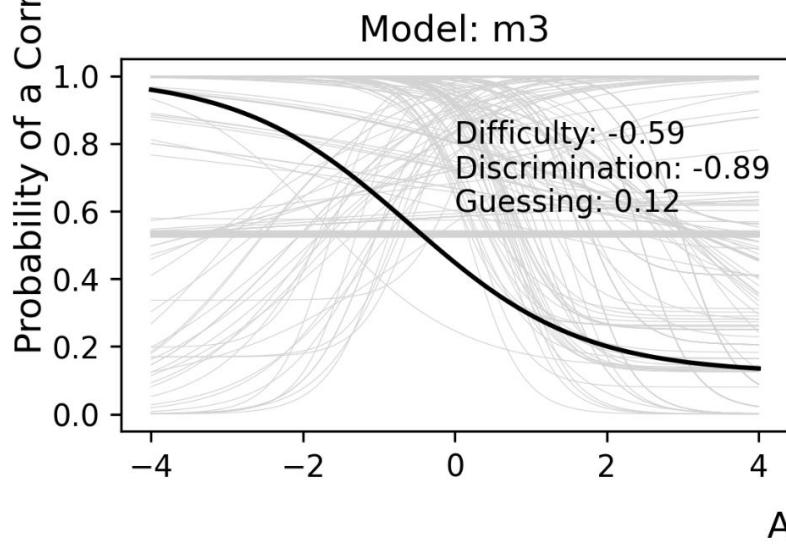
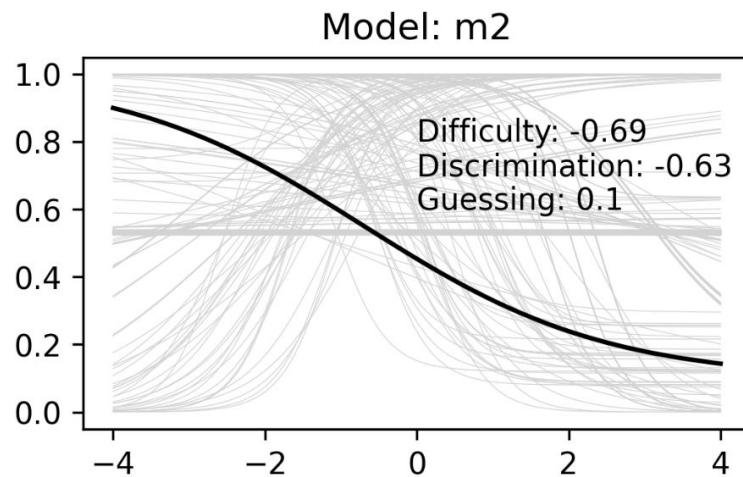
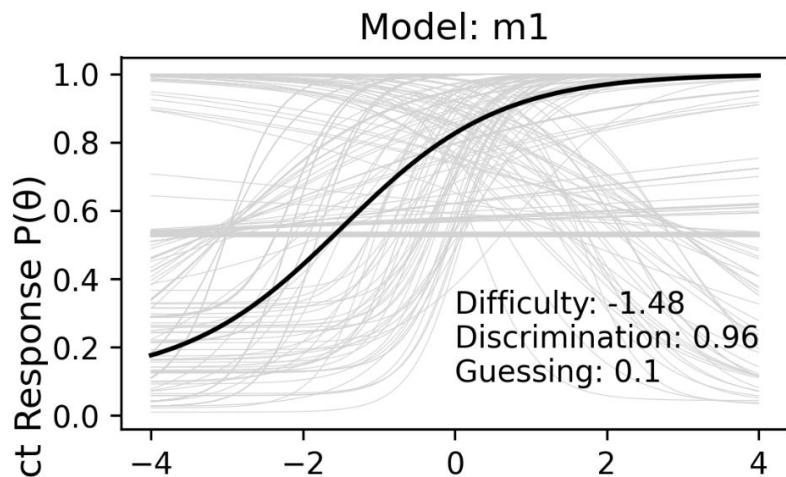




Item Characteristic Curve of the credit-g dataset (test) for models M1 to M4



Mean Item Characteristic Curve of the credit-g dataset (test) for models M1 to M4



Local ICC — Global ICC