



NYU

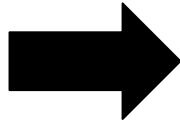
Exploring the Relationship Between Feature Attribution Methods and Model Performance

Priscylla Silva
Claudio Silva
Luis Gustavo Nonato

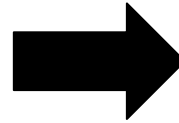


Motivation

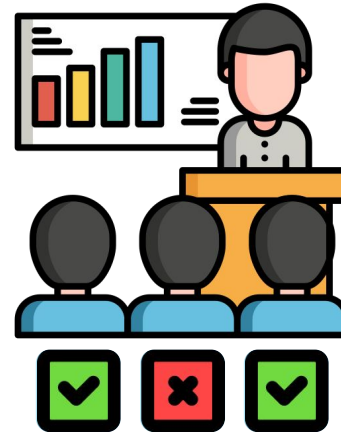
Educational Context



ML models



Predictions



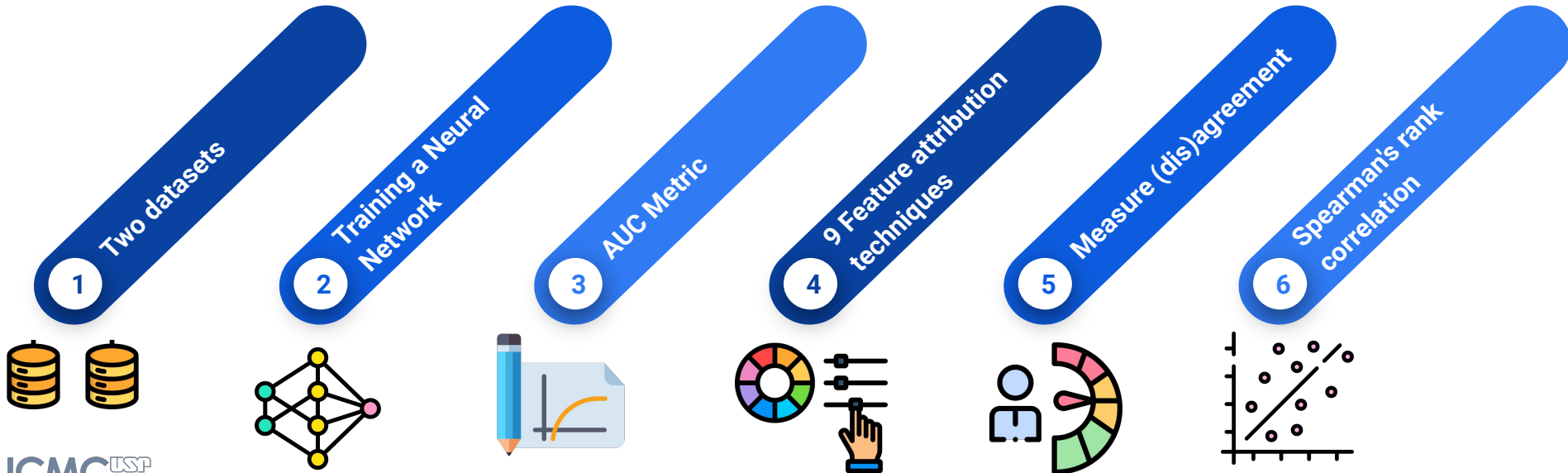
Why?



Research Question

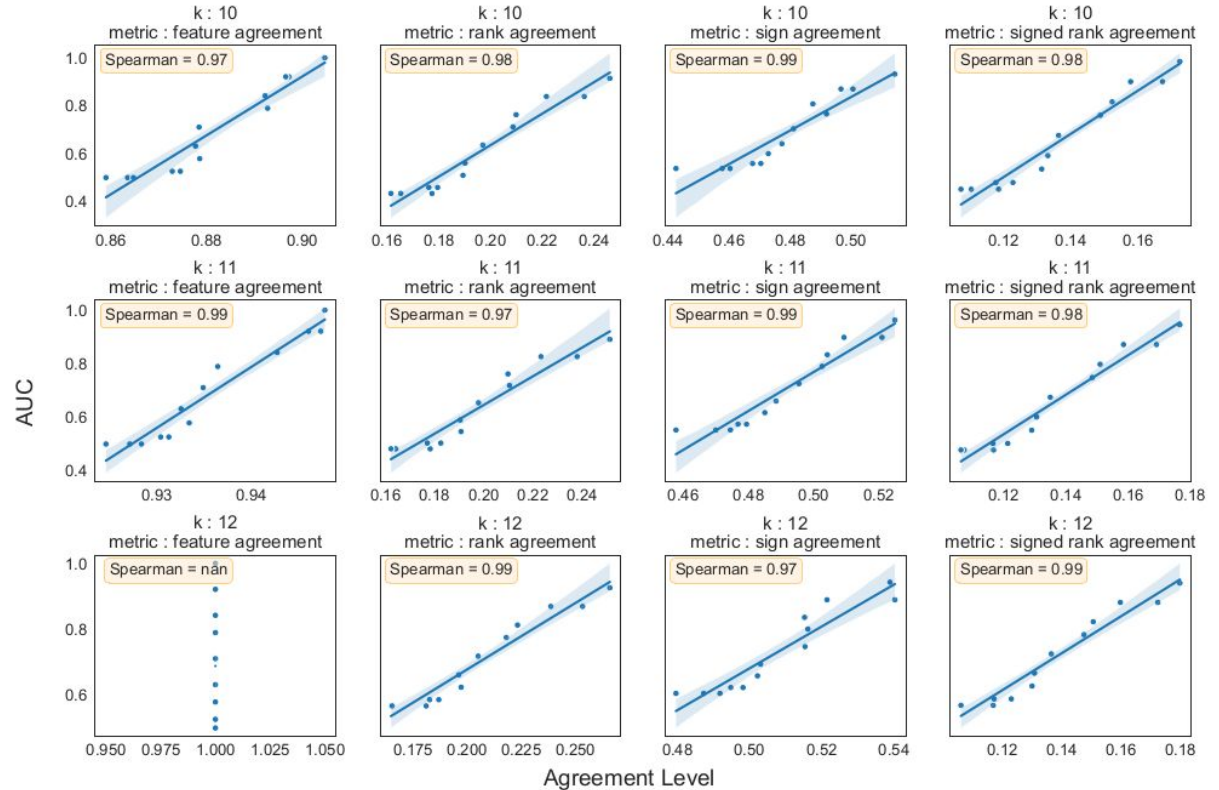
Is there a correlation between the model's performance and the disagreement level observed among explanation methods?

Study methodology



Results

Our results show that models with an AUC greater than or equal to 0.8 consistently exhibit the highest levels of agreement among explanation methods.





NYU

Exploring the Relationship Between Feature Attribution Methods and Model Performance

Priscylla Silva
Claudio Silva
Luis Gustavo Nonato