

DEBIASING ALGORITHMS AND GENERALIZATION IN PROTEIN-LIGAND BINDING



VIKRAM SUNDAR AND LUCY COLWELL
DEPARTMENT OF CHEMISTRY, UNIVERSITY OF CAMBRIDGE

INTRODUCTION

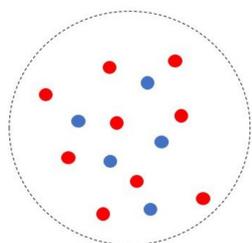
- Many machine learning (ML) approaches to protein/ligand binding have achieved outstanding success on benchmark datasets. [1]
- Unclear whether performance indicates generalizability or overfitting to training data. [2, 3]
- Chemical space is naturally clustered and our datasets are further clustered due to biases from experimentalists. [2, 3]

DATASET

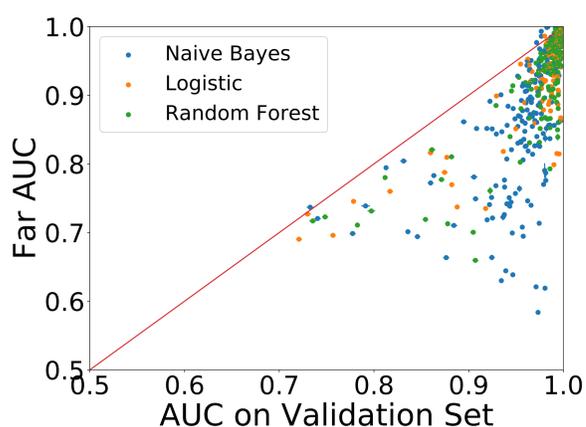
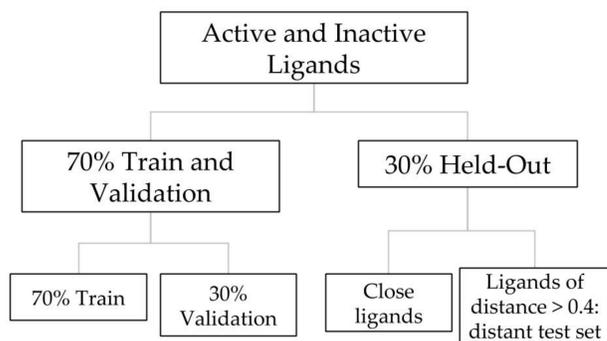
- Active ligands acquired from ChEMBL. Activity threshold of 1 μ M.
- Inactive ligands acquired from PubChem. Further inactives randomly drawn until an even split of actives and inactives was achieved.
- ECFP6 fingerprints with 2048 bits used as feature set. Tanimoto similarity as distance metric.
- Tested 189 targets with at least 500 active ligands.

FAR AUC

- To measure generalizability, hold out a distant test set to measure algorithm performance.



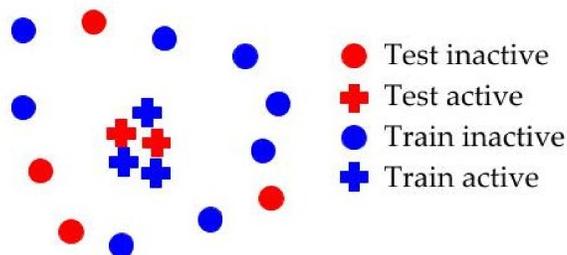
Training/Validation Sets Held-Out Test Set



- All models show poor generalization.
- Validation set AUC may not be representative of true performance of model.

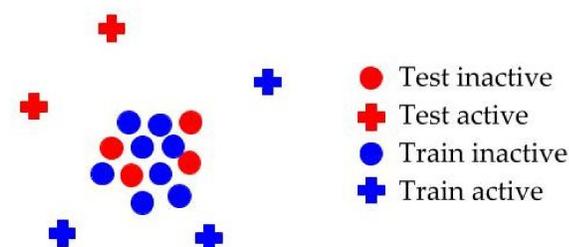
DEBIASING ALGORITHMS

- Goal of data bias and debiasing algorithm is to prevent models from memorizing and improve their generalization ability.
- Maximum Unbiased Validation (MUV) requires active ligands be uniformly embedded within inactive ligands. [2] This set is MUV-biased:



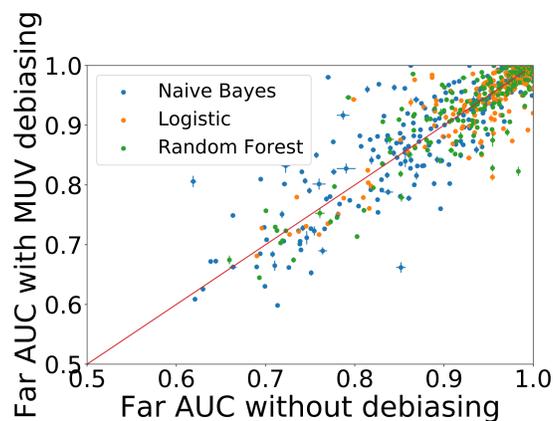
- **Key Question: Do these debiasing algorithms improve generalization ability as measured by the far AUC?**

- Asymmetric Validation Embedding (AVE) prevents inactives from being clustered. [3] This set is AVE-biased:

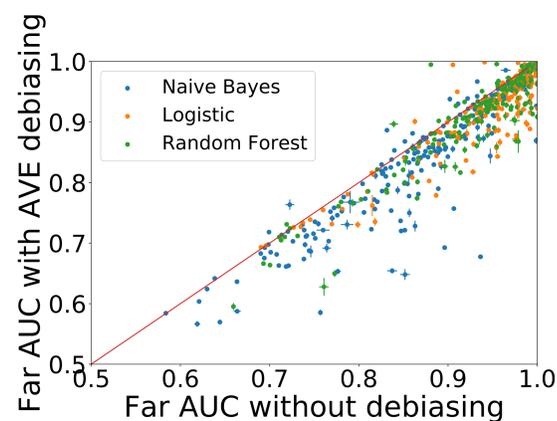


- We use genetic algorithms to debias a given train/validation split. [2, 3]

RESULTS



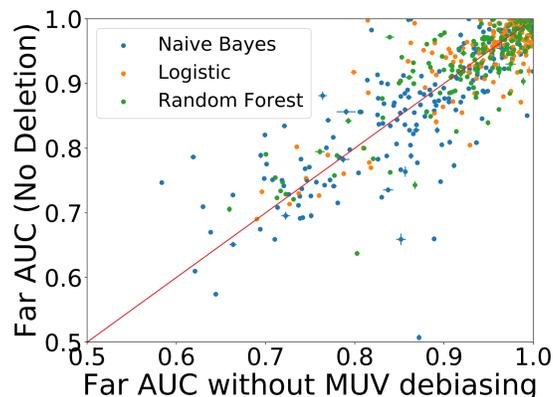
- Debiasing algorithms do not consistently improve, and often worsen, the ability of models to generalize.



- Independent of number of active ligands, number of inactive ligands (or random decoys added), and degree of debiasing.

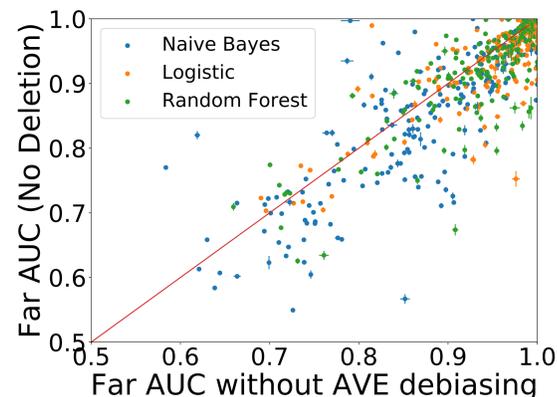
AVOIDING DATA DELETION

- Debiasing algorithms allow for deletion of data from both training and validation set.



- Slight improvement over the results while allowing deletion.

- Still possible to partially debias datasets without deleting data.



- Still no consistent improvement over performance without debiasing.

CONCLUSIONS

- Debiasing algorithms cannot distinguish between signal and bias and remove relevant information from the training data.
- Some clustering among actives expected in protein/ligand binding. Need to distinguish between this and artificial clustering.
- Generalization in ML is hard. Better approach: understand domains where our models are applicable.

ACKNOWLEDGMENTS

I acknowledge the support of the Winston Churchill Foundation of the USA.

REFERENCES

- [1] Lucy J. Colwell. Statistical and machine learning approaches to predicting protein-ligand interactions. *Current Opinion in Structural Biology*, 49:123–128, 2018.
- [2] Sebastian G. Rohrer and Knut Baumann. Maximum unbiased validation (MUV) data sets for virtual screening based on PubChem bioactivity data. *Journal of Chemical Information and Modeling*, 2009.
- [3] Izhar Wallach and Abraham Heifets. Most Ligand-Based Classification Benchmarks Reward Memorization Rather than Generalization. *Journal of Chemical Information and Modeling*, 2018.