

# Dataset Structure and Active Learning

---

A Comparative  
Regression Study

# Presentation Structure

- Motivation
- Datasets
- Methods
- Hypotheses
- Results
- Conclusions



# Presentation Structure

- **Motivation**
- Datasets
- Methods
- Hypotheses
- Results
- Conclusions



# Motivation

- The efficacy of an active learning method depends on the characteristics of the data to which it is applied
- We wanted to study this relationship in the context of regression



# Synthetic vs. real data

- Synthetic data allows for better isolation of dataset characteristics
- Real data can give us domain-specific insights and – at least to us – is more interesting to work with





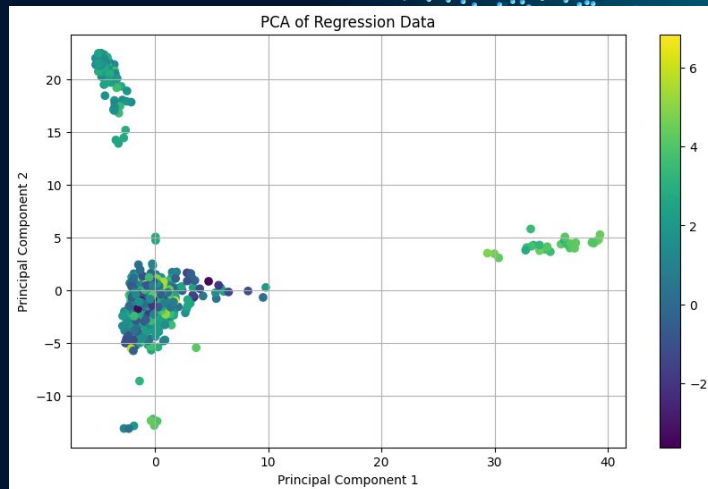
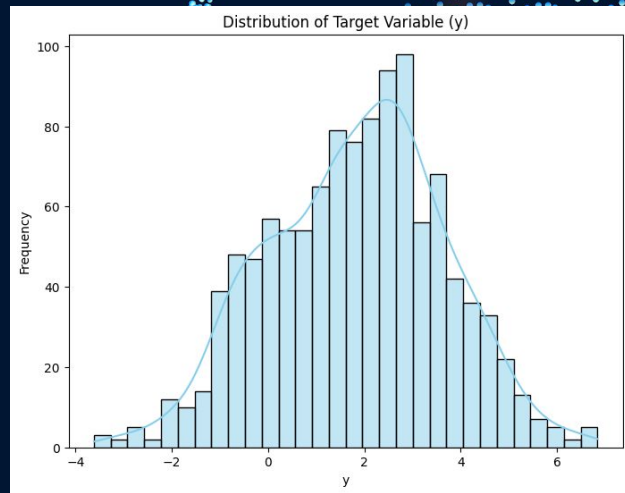
# Presentation Structure

- Motivation
- **Datasets**
- Methods
- Hypotheses
- Results
- Conclusions



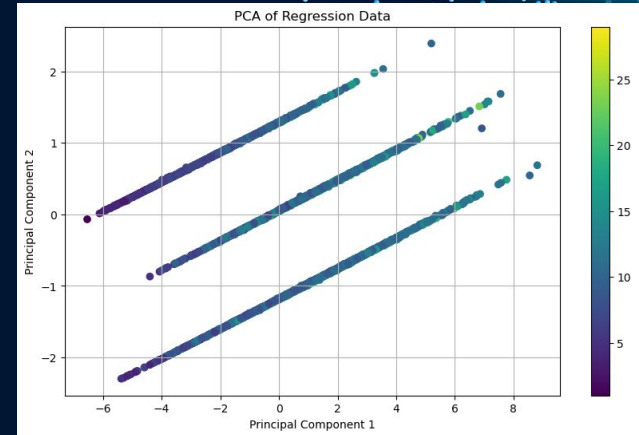
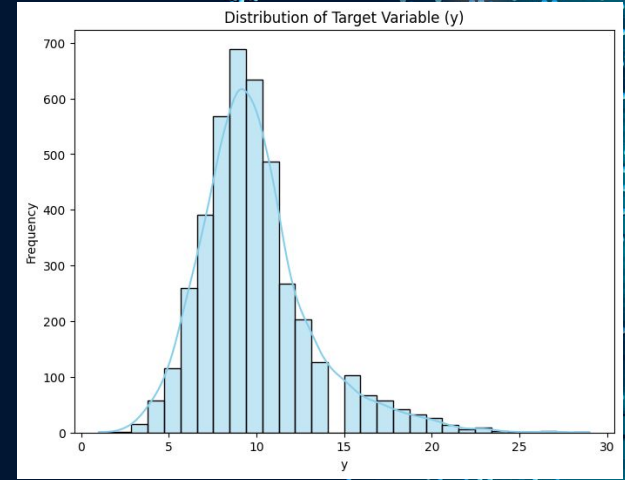
# Logd74

- A dataset of 1130 organic compounds and their lipophilicity
- Featurized with Morgan fingerprinting
- Lipophilicity is relevant to drug discovery, as some targets may be within lipid membranes
- (1130 samples, 2048 features)



# Abalone Age

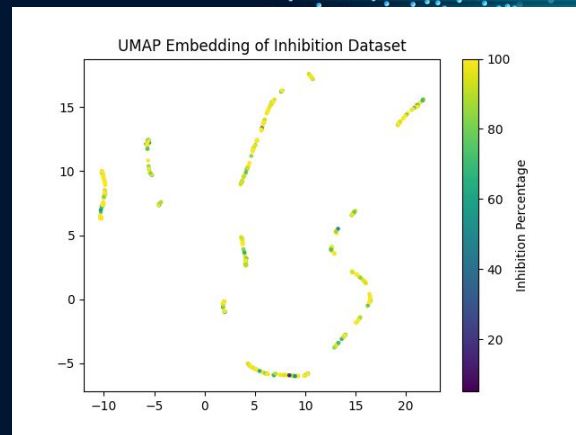
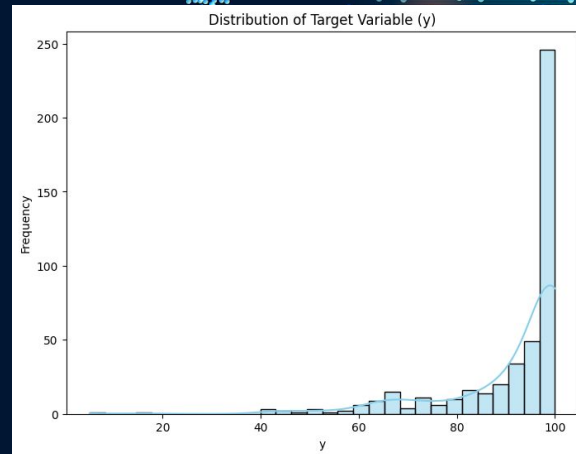
- UCI ML repository classic dataset
- Regression Target: Abalone Age
- Features: Sex, Length, Diameter, Height, Weights x6
- Use? Abalone age-identification process is inconvenient
- (4177 samples, 8 features)





# Kinase Inhibition Data

- A dataset of 569 proteins and their inhibition levels by a small molecule drug
- Featurized with protein physicochemical properties using biopython
- The goal of using this dataset is to capture which chemical properties cause a certain small molecule inhibitor to be most effective
- (569 samples, 28 features)



# Presentation Structure

- Motivation
- Datasets
- **Methods**
- Hypotheses
- Results
- Conclusions



# Random forest batch selection

- Random forest base learner
- Batch query methods:
  - Random
  - Uncertainty
  - K-means diversity
  - Hierarchical diversity



## ALICE

1. Data Access:  
Membership  
Query  
Synthesis
2. DOE Criterion:  
A-optimality
3. Assumes:  
 $P_{\text{test}}(X)$  known  
or accurately  
approximated

- Aggressive +  
IW Least Squares
- DOE Inspired
- PAC Guarantees
- Linear combo basis  
regression models
- Robust against  
model  
misspecification

## P-ALICE

1. Data Access:  
Pool-Based
2. DOE Criterion:  
Q-optimality  
[Ref 1,2]
3. Assumes:  
 $p_{\text{te}}(x) > 0$ , for  
all samples in  
pool

# P-ALICE/ALICE Base Learner

$$\hat{f}(x) = \sum_{\ell=1}^t \theta_{\ell} \varphi_{\ell}(x).$$

Eqn 3

## P-ALICE

**Input:** Test input points  $\{\mathbf{x}_j^{\text{te}}\}_{j=1}^{n_{\text{te}}}$  and basis functions  $\{\varphi_{\ell}(\mathbf{x})\}_{\ell=1}^t$ .

**Output:** Learned parameter  $\hat{\theta}_{\text{w}}$

Compute the  $t \times t$  matrix  $\hat{U}$  with  $\hat{U}_{\ell, \ell'} = \frac{1}{n_{\text{te}}} \sum_{j=1}^{n_{\text{te}}} \varphi_{\ell}(\mathbf{x}_j^{\text{te}}) \varphi_{\ell'}(\mathbf{x}_j^{\text{te}})$ ;

**For** several different values of  $\lambda$  (possibly around  $\lambda = 1/2$ )

    Compute  $\{b_{\lambda}(\mathbf{x}_j^{\text{te}})\}_{j=1}^{n_{\text{te}}}$  with  $b_{\lambda}(\mathbf{x}) = (\sum_{\ell, \ell'=1}^t [\hat{U}^{-1}]_{\ell, \ell'} \varphi_{\ell}(\mathbf{x}) \varphi_{\ell'}(\mathbf{x}))^{\lambda}$ ;

    Choose  $\mathcal{X}_{\lambda}^{\text{tr}} = \{\mathbf{x}_i^{\text{tr}}\}_{i=1}^{n_{\text{tr}}}$  from  $\{\mathbf{x}_j^{\text{te}}\}_{j=1}^{n_{\text{te}}}$  with probability proportional to  $\{b_{\lambda}(\mathbf{x}_j^{\text{te}})\}_{j=1}^{n_{\text{te}}}$ ;

    Compute the  $n_{\text{tr}} \times t$  matrix  $X_{\lambda}$  with  $[X_{\lambda}]_{i, \ell} = \varphi_{\ell}(\mathbf{x}_i^{\text{tr}})$ ;

    Compute the  $n_{\text{tr}} \times n_{\text{tr}}$  diagonal matrix  $W_{\lambda}$  with  $[W_{\lambda}]_{i, i} = (b_{\lambda}(\mathbf{x}_i^{\text{tr}}))^{-1}$ ;

    Compute  $L_{\lambda} = (X_{\lambda}^{\top} W_{\lambda} X_{\lambda})^{-1} X_{\lambda}^{\top} W_{\lambda}$ ;

    Compute  $\text{P-ALICE}(\lambda) = \text{tr}(\hat{U} L_{\lambda} L_{\lambda}^{\top})$ ;

**End**

Compute  $\hat{\lambda} = \text{argmin}_{\lambda} \text{P-ALICE}(\lambda)$ ;

Gather training output values  $\mathbf{y}^{\text{tr}} = (y_1^{\text{tr}}, y_2^{\text{tr}}, \dots, y_{n_{\text{tr}}}^{\text{tr}})^{\top}$  at  $\mathcal{X}_{\hat{\lambda}}^{\text{tr}}$ ;

Compute  $\hat{\theta}_{\text{w}} = L_{\hat{\lambda}} \mathbf{y}^{\text{tr}}$ ;

**Fig. 3** Pseudo code of proposed pool-based active learning algorithm



# P-ALICE Base Learners

1. Linear Regression
2. EDA Spearman Significant Features for Squared, Cubic, and First-Order Interactions
3. -1/+1 map for 0-1 Binary Features

Used for Abalone Age & Kinase Inhibition

Used for Logd74

$$\varphi_\ell(\vec{x}) : \{ \underline{x_l}, \underline{x_l^2}, \underline{x_l^3}, \sin(x_l), \text{ReLu}(x_l), \underline{\mathbb{1}(x_l)}, \dots \}$$
$$x_l : \{ x_i, \underline{x_i x_j}, x_i x_j x_k, \dots \}$$

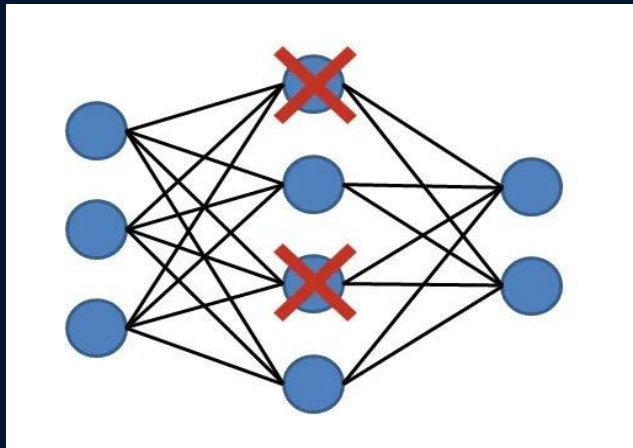
# Neural network batch selection

- Base learner: a multi-layer perceptron (MLP)
- Batch query methods:
  - Random
  - Uncertainty
  - COVDROP



## Neural network uncertainty estimation

- To estimate uncertainty, perform forward passes with dropout repeatedly and calculate variance of predictions



# COVDROP

- Premise: the  $k$  most uncertain instances may be highly correlated
- Instead, choose points that maximize joint entropy
  - Equivalent to choosing points maximizing determinant of covariance matrix





# Determinant Maximization Sampling

**Input:**  $S = \{x_0, \dots, x_N\}$ , sample pool

**Input:**  $\mathcal{M}$ , trained model

**Input:**  $K$  batch size;  $M$  # of starts

**Data:**  $\text{Cov} \in \mathbb{R}^{N \times N}$ , covariance for  $\mathcal{M}$  on  $S$ , (computation given elsewhere)

2 **begin**

3   Extract variances  $\text{Var}$  from  $\text{Cov}$  ;

4    $\mu$  probability measure on  $S$ , proportional to  $\text{Quantile}(\text{Var})$ ;

5   Choose  $B_m^0 = (x_{m1}^0, \dots, x_{mK}^0)$ ,  $0 \leq m < M$  random starting batches ;

6    $\text{Cov}(B_m^0)$ , principal submatrices of  $\text{Cov}$ ;

7   Compute scores  $S_m^0 := \log \det \text{Cov}(B_m^0)$  via Cholesky decomposition;

8   **Optional** Keep only  $M' < M$  highest scoring starts;

9    $i \leftarrow 0$ ,  $i \in \mathbb{Z}/K\mathbb{Z}$ ;

10   **while not converged do**

11     **for**  $m = 1$  **to**  $M$  **do**

12       **for**  $j = 1$  **to**  $N$  **do**

13          Let  $B_{m,j}^i$  be  $B_m^i$  with  $x_j$  substituted at the  $i$ th position;

14          Compute score  $S_{m,j}^i := \log \det \text{Cov}(B_{m,j}^i)$  via rank-1 Cholesky update;

15           $B_m^{i+1} \leftarrow B_{m,J}^i$  where  $J = \arg \max_j S_{m,j}^i$ ;

16        $i \leftarrow i + 1$ ;

**Result:**  $B = B_b^i$  where  $b = \arg \max_m B_m^i$



# Presentation Structure

- Motivation
- Datasets
- Methods
- **Hypotheses**
- Results
- Conclusions



## Hypotheses (Batch sampling)

- We expected batch methods that consider the makeup of a batch to outperform methods that do not
- Thus, both methods below should outperform selecting the  $k$  most uncertain instances:
  - Random forest cluster-based batch sampling
  - Neural network COVDROP



## Hypotheses (Dataset structure)

- We expected batch methods that consider the makeup a batch to have a greater performance advantage when the dataset is highly heterogeneous
  - E.g. logd74
- We expected P-ALICE active sampling to be less effective on noisy and high-dim. datasets
  - E.g. logd74

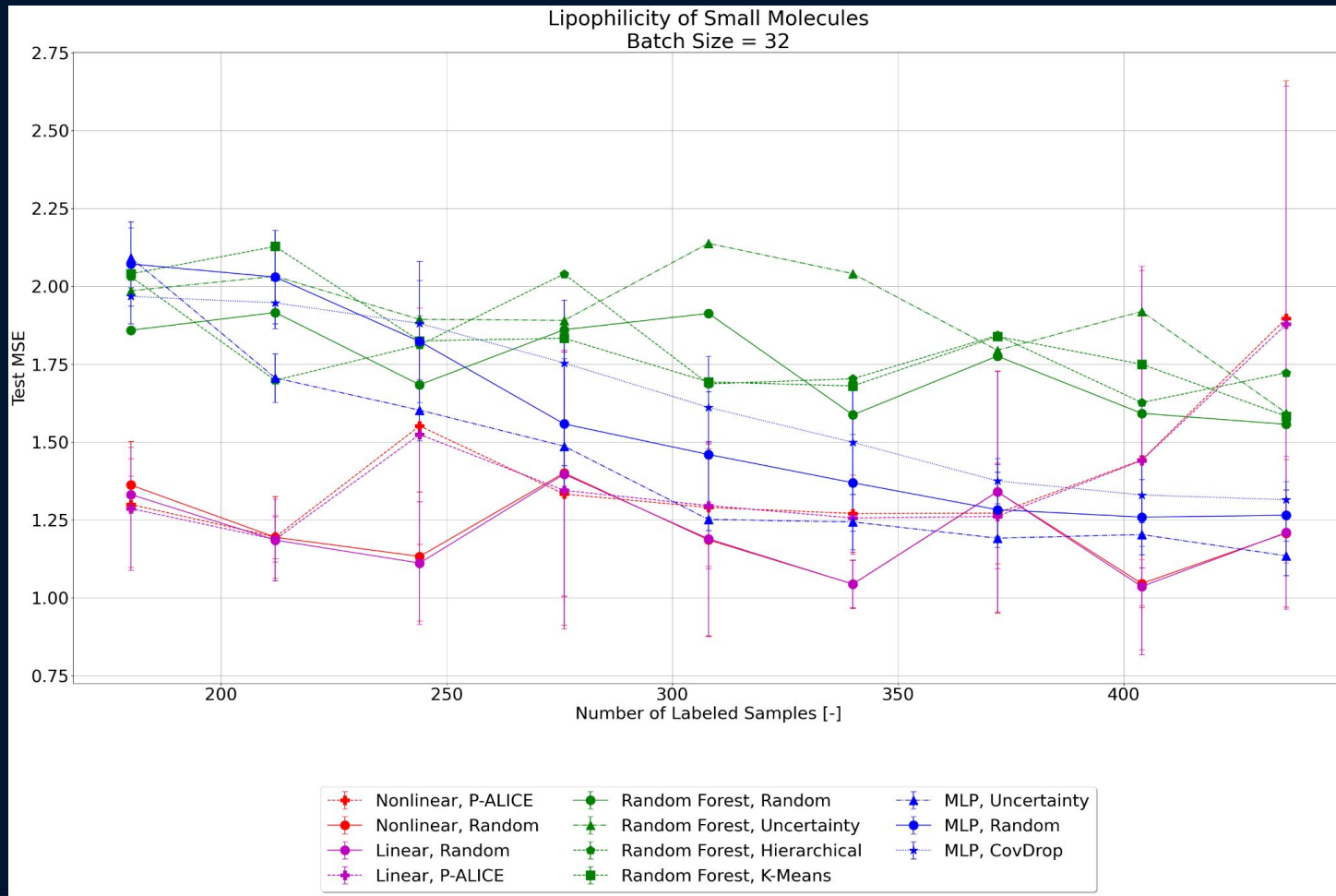


# Presentation Structure

- Motivation
- Datasets
- Methods
- Hypotheses
- **Results**
- Conclusions

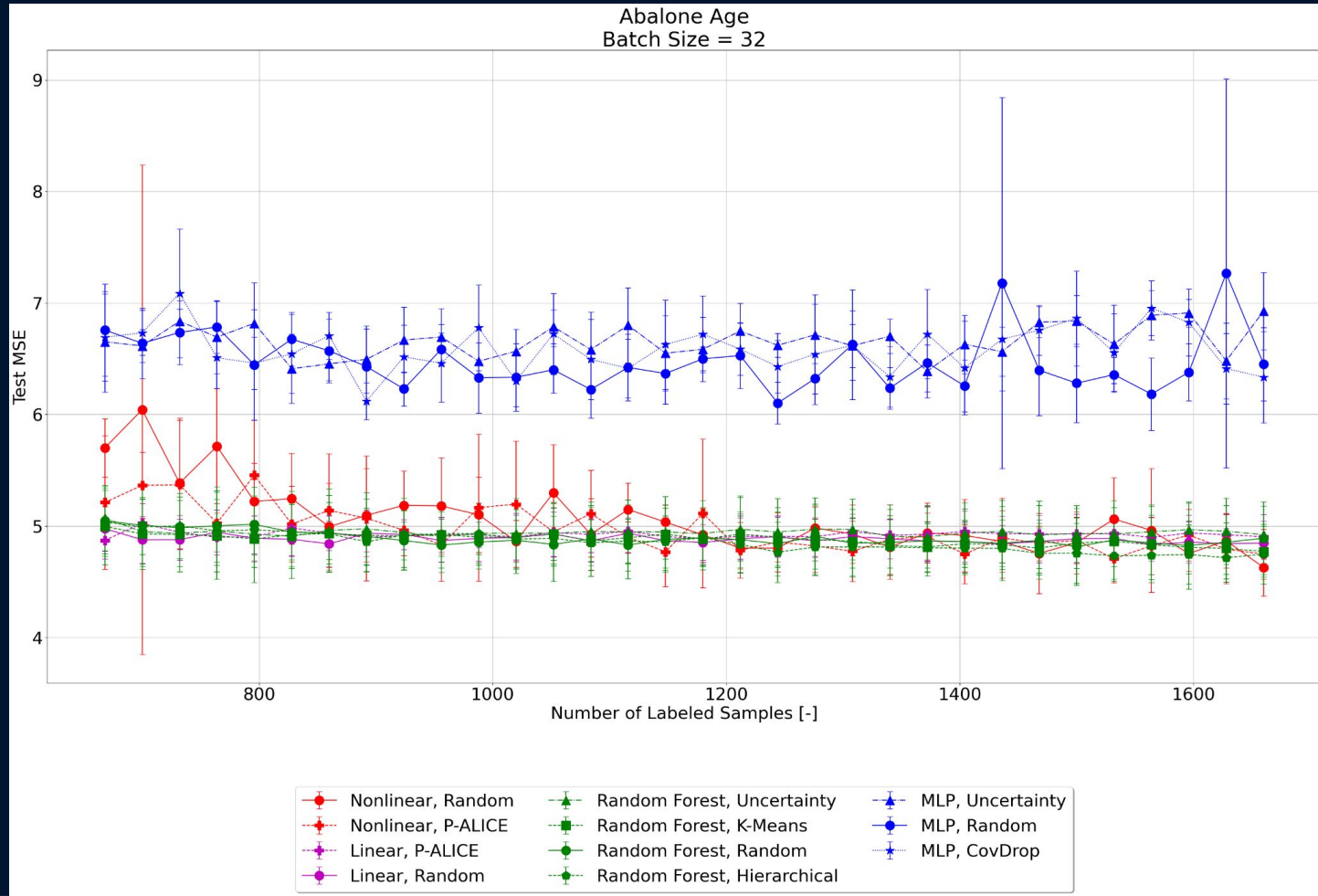


# Results Logd74

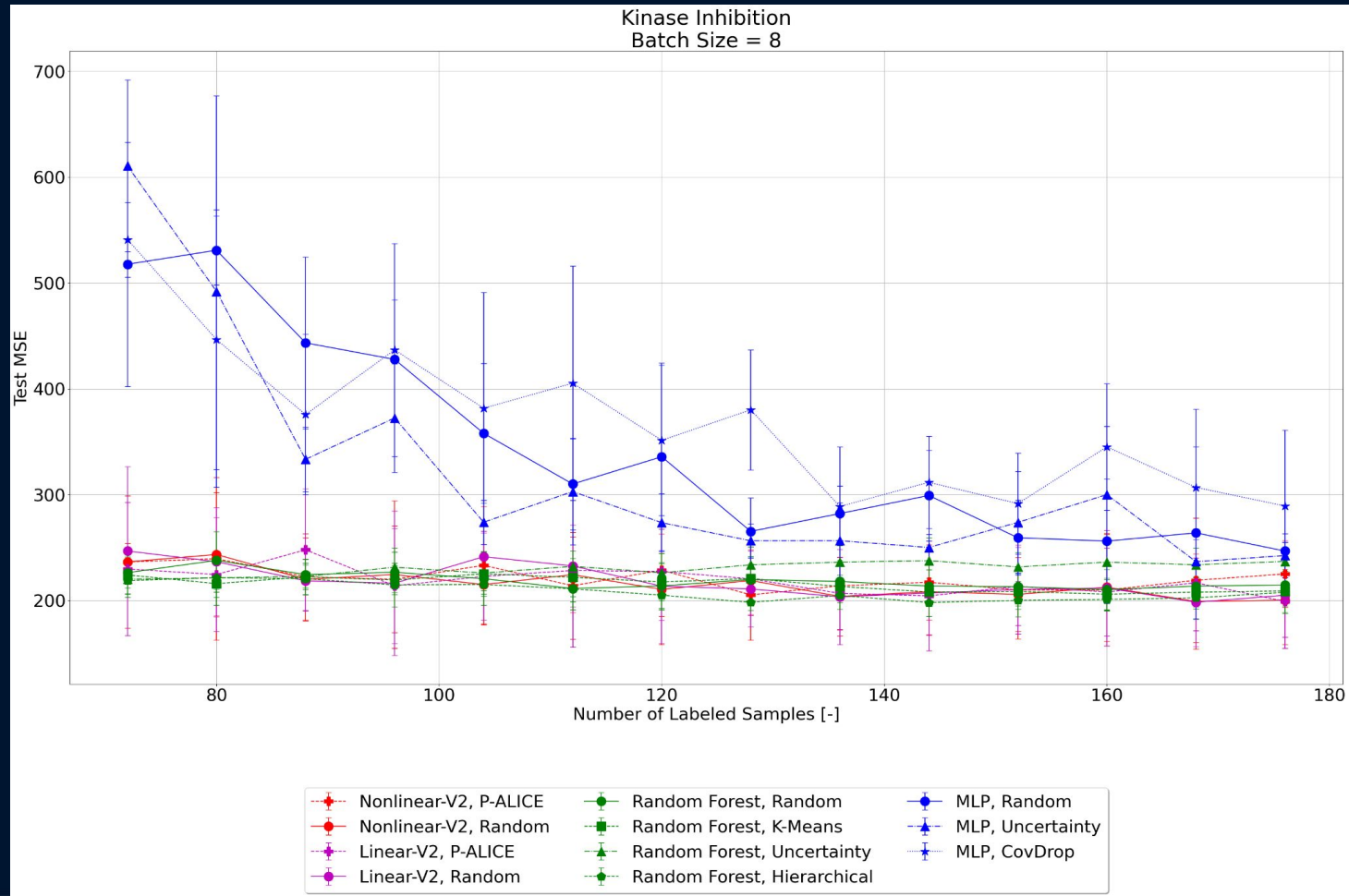




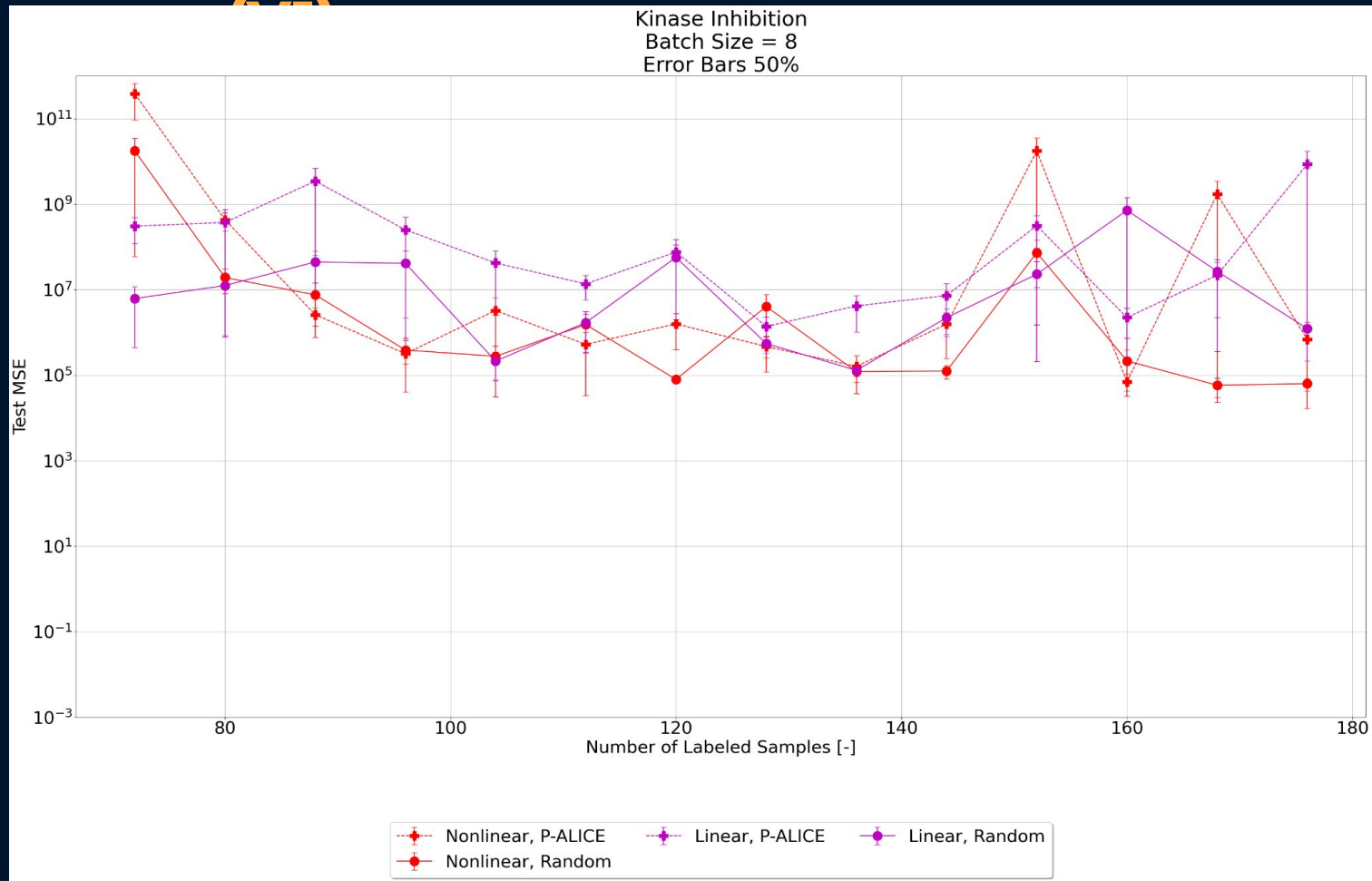
# Results Abalone



# Results Kinase Inhibition (V2)



# Results Kinase Inhibition



# Presentation Structure

- Motivation
- Datasets
- Methods
- Hypotheses
- Results
- **Conclusions**



# Conclusions

---

- Base learners can be a confounding variable when evaluating active learning methods across different base learners
- Real datasets can be rewarding to work with, but summarizing their differences can be challenging
- Let us know if you have ideas about how we can improve our analysis!





# References

- [1] <https://link.springer.com/article/10.1007/s10994-009-5100-3>
- [2] <https://statweb.rutgers.edu/buyske/591/lect11.pdf>
- [3] [https://www.researchgate.net/publication/254196943\\_Using\\_Random\\_Forest\\_to\\_Learn\\_Imbalanced\\_Data#:~:text=Lastly%2C%20ensemble%20methods%2C%20such%20as,class%20%5B7%2C%2018%5D.](https://www.researchgate.net/publication/254196943_Using_Random_Forest_to_Learn_Imbalanced_Data#:~:text=Lastly%2C%20ensemble%20methods%2C%20such%20as,class%20%5B7%2C%2018%5D.)
- [4] <https://elifesciences.org/reviewed-preprints/89679>

