# Supplementary Material

Breno Caetano da Silva
Jorge Luiz Franco
Hermes Manoel Galvão Castelo Branco
Eric K. Tokuda
Alexandre Delbem

July 16, 2024

## Contents

# 1 Introduction

In Active Learning (AL), selecting an effective initial set of labeled samples is crucial for developing robust models with limited data. Traditional AL methods often rely on an oracle to label samples, but without a well-chosen initial set, model performance can be suboptimal.

This supplementary material presents PhyIL, our novel approach using phylogram analysis and community detection to enhance cold-start selection in AL. By choosing samples that represent data diversity and complexity, PhyIL ensures a robust initial model.
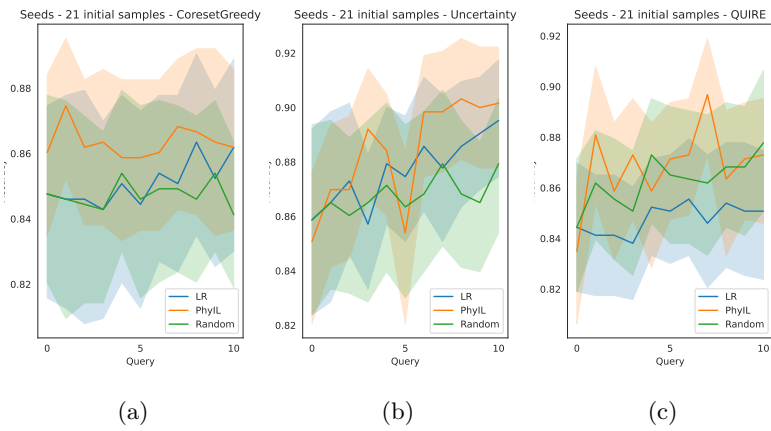
We evaluated PhyIL on nine diverse real-world datasets: seeds, sonar, parkinson, wine, haberman, glass, heart, and musk. These datasets were chosen for their variety in features and challenges, offering a comprehensive test of PhyIL's effectiveness. For example, the seeds dataset involves grain classification, while the sonar dataset deals with the differentiation between mines and rocks based on sonar signals. The parkinson dataset contains speech signals to detect Parkinson's disease, and the wine dataset is used for wine classification based on chemical properties.

We compared PhyIL against established methods: Query By Committee (QBC), which uses multiple models to select informative samples; QUIRE, which selects samples based on label ambiguity; Coreset-Greedy, which optimizes sample selection for model accuracy; and Uncertainty Sampling, which picks samples where the model is least certain.
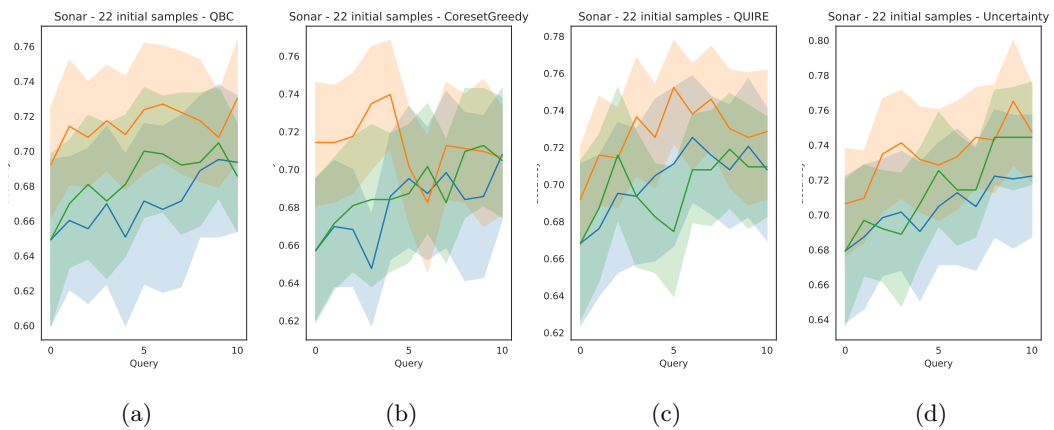
Our results show that PhyIL outperforms these methods, leading to better initial sample selection and improved model performance. This document provides detailed results for various classifiers, including ADABOOST, Linear Regression, Gaussian Process, Support Vector Machine (SVM), Decision Tree, Multi-Layer Perceptron (MLP), Random Forest, and K-Nearest Neighbors (KNN), demonstrating PhyIL's advantages across these datasets.
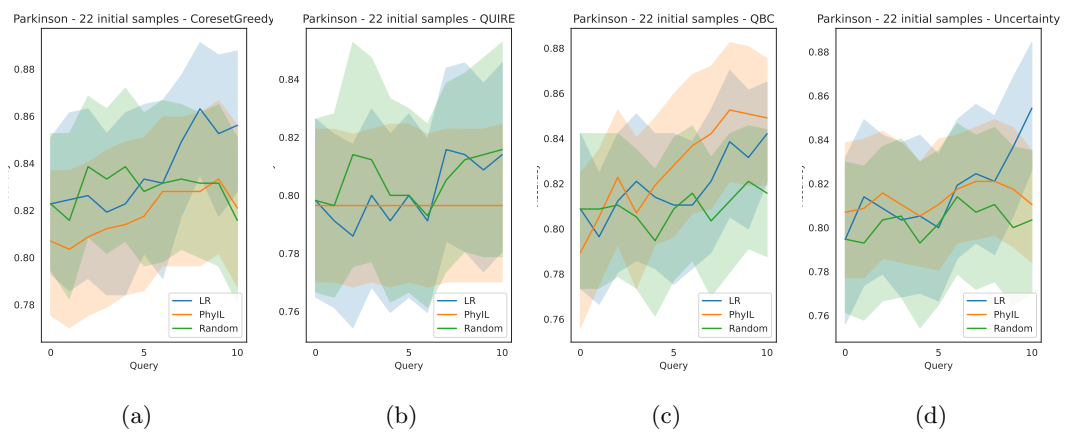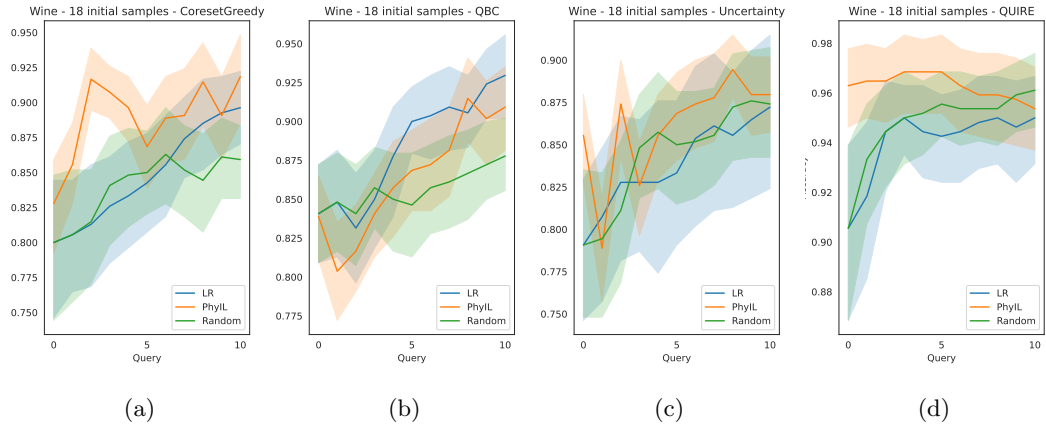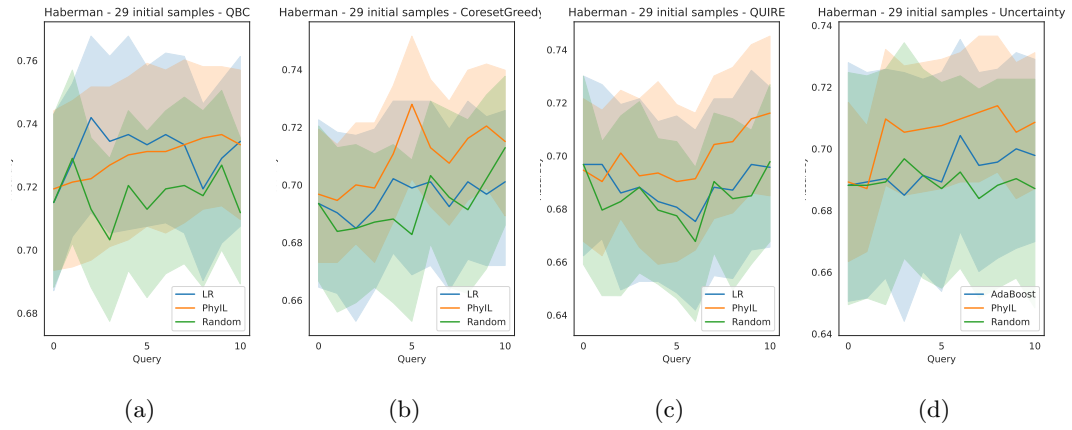
# 2 ADABOOST

## 2.1 seeds



(a)                    (b)                    (c)

## 2.2 sonar



(a)                (b)                (c)                (d)

## 2.3 parkinson



(a)                (b)                (c)                (d)

## 2.4 wine



Wine - 18 initial samples - CoresetGreedy

Wine - 18 initial samples - QBC

Wine - 18 initial samples - Uncertainty

Wine - 18 initial samples - QUIRE

(a)　　　　　　　(b)　　　　　　　(c)　　　　　　　(d)

## 2.5 haberman



Haberman - 29 initial samples - QBC

Haberman - 29 initial samples - CoresetGreedy

Haberman - 29 initial samples - QUIRE

Haberman - 29 initial samples - Uncertainty

(a)　　　　　　　(b)　　　　　　　(c)　　　　　　　(d)

## 2.6 glass



Glass - 22 initial samples - Uncertainty

Glass - 22 initial samples - QBC

Glass - 22 initial samples - QUIRE

Glass - 22 initial samples - CoresetGreedy

(a)　　　　　　　(b)　　　　　　　(c)　　　　　　　(d)

## 2.7 heart

(a)     (b)     (c)     (d)

## 2.8 musk



(a)     (b)     (c)     (d)

# 3 Linear Regression

## 3.1 seeds



(a)     (b)     (c)     (d)

## 3.2 sonar

(a)       (b)       (c)       (d)

## 3.3   parkinson



(a)       (b)       (c)       (d)

## 3.4   wine



(a)       (b)       (c)       (d)

## 3.5   glass

Glass - 22 initial samples - CoresetGreedy (a)
Glass - 22 initial samples - Uncertainty (b)
Glass - 22 initial samples - QBC (c)
Glass - 22 initial samples - QUIRE (d)

## 3.6 heart



Heart - 26 initial samples - QUIRE (a)
Heart - 26 initial samples - QBC (b)
Heart - 26 initial samples - CoresetGreedy (c)
Heart - 26 initial samples - Uncertainty (d)

## 3.7 musk



Musk_v1 - 34 initial samples - CoresetGreedy (a)
Musk_v1 - 34 initial samples - QUIRE (b)
Musk_v1 - 34 initial samples - Uncertainty (c)
Musk_v1 - 34 initial samples - QBC (d)

## 3.8 haberman

(a)        (b)        (c)        (d)

# 4 Gaussian Process

## 4.1 seeds



(a)        (b)        (c)        (d)

## 4.2 sonar



(a)        (b)        (c)        (d)

## 4.3 parkinson

(a)                (b)                (c)

## 4.4   wine



(a)           (b)           (c)           (d)

## 4.5   haberman



(a)                (b)                (c)

## 4.6   glass

(a)        (b)        (c)        (d)

## 4.7   heart



(a)        (b)        (c)        (d)

## 4.8   musk



(a)        (b)        (c)        (d)

# 5   SVM

## 5.1   seeds

(a)       (b)       (c)       (d)

## 5.2    sonar



(a)       (b)       (c)       (d)

## 5.3    parkinson



(a)       (b)       (c)

## 5.4    haberman

(a)  (b)  (c)  (d)

## 5.5 glass



(a)  (b)  (c)  (d)

## 5.6 heart



(a)  (b)  (c)  (d)

## 5.7 wine

(a)       (b)       (c)       (d)

## 5.8   musk



(a)       (b)       (c)       (d)

# 6   Decision Tree

## 6.1   heart



(a)       (b)       (c)       (d)

## 6.2   musk

(a)                    (b)                    (c)                    (d)

# 7 MLP

## 7.1 seeds



(a)                    (b)                    (c)                    (d)

## 7.2 sonar



(a)                    (b)                    (c)                    (d)

## 7.3 parkinson

Parkinson - 22 initial samples - CoresetGreedy    Parkinson - 22 initial samples - QUIRE    Parkinson - 22 initial samples - QBC    Parkinson - 22 initial samples - Uncertainty

(a)      (b)      (c)      (d)

## 7.4   wine



Wine - 18 initial samples - CoresetGreedy    Wine - 18 initial samples - QBC    Wine - 18 initial samples - Uncertainty    Wine - 18 initial samples - QUIRE

(a)      (b)      (c)      (d)

## 7.5   haberman



Haberman - 29 initial samples - QBC    Haberman - 29 initial samples - CoresetGreedy    Haberman - 29 initial samples - QUIRE    Haberman - 29 initial samples - Uncertainty

(a)      (b)      (c)      (d)

## 7.6   glass

(a)  (b)  (c)  (d)

## 7.7 heart



(a)  (b)  (c)  (d)

## 7.8 musk



(a)  (b)  (c)  (d)

# 8 Random Forest

## 8.1 seeds

(a)　　　　　　　(b)　　　　　　　(c)　　　　　　　(d)

## 8.2    parkinson



(a)　　　　　　　(b)　　　　　　　(c)　　　　　　　(d)

## 8.3    wine



(a)　　　　　　　(b)　　　　　　　(c)

## 8.4    glass

(a)                     (b)                     (c)                     (d)

## 8.5   heart



(a)                     (b)                     (c)                     (d)

## 8.6   musk



(a)                     (b)                     (c)                     (d)

# 9   KNN

## 9.1   seeds

(a)          (b)          (c)          (d)

## 9.2 sonar



(a)          (b)          (c)          (d)

## 9.3 parkinson



(a)          (b)          (c)          (d)

## 9.4 wine

(a)           (b)           (c)           (d)

## 9.5    haberman



(a)           (b)           (c)           (d)

## 9.6    glass



(a)           (b)           (c)           (d)

## 9.7    heart

Heart - 26 initial samples - QUIRE (a)
Heart - 26 initial samples - QBC (b)
Heart - 26 initial samples - CoresetGreedy (c)
Heart - 26 initial samples - Uncertainty (d)

## 9.8 musk



Musk_v1 - 34 initial samples - CoresetGreedy (a)
Musk_v1 - 34 initial samples - QUIRE (b)
Musk_v1 - 34 initial samples - Uncertainty (c)
Musk_v1 - 34 initial samples - QBC (d)