# BELLWETHER 2.0

By:

Aayushi Agrawal

Akshay Ravichandran

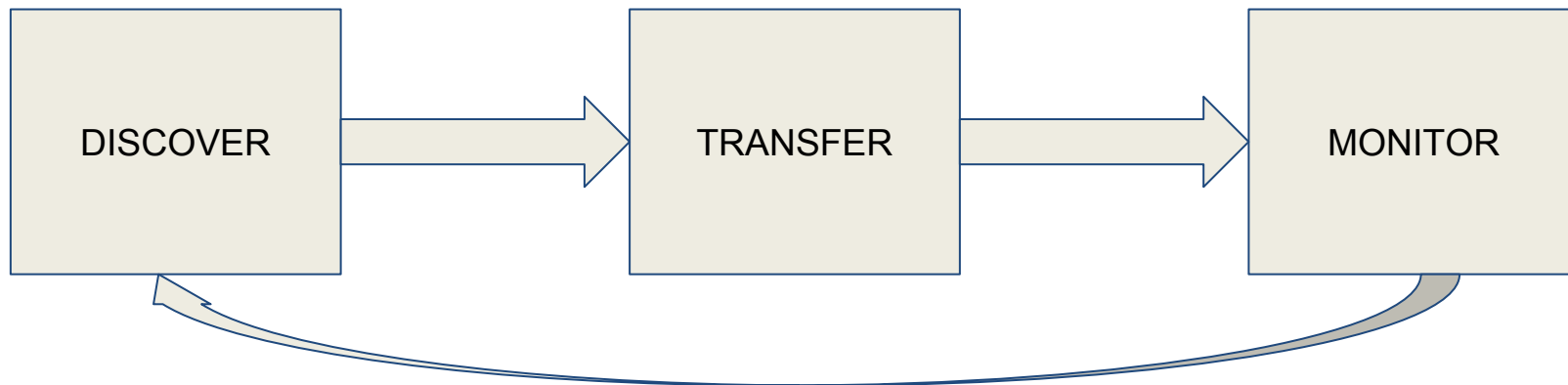Yudong RAO

# INTRODUCTION

- What is a '*bellwether*'?

  **A:** Given N projects from a community of projects, the '*bellwether*' is that project whose data yields the best predictions on all others

# ORIGINAL STUDY

- Conducted by Rahul Krishna and Prof. Menzies[5]
- Introduced the ***bellwether method:***

# MOTIVATION

- Bellwethers appear prevalent in SE data sets - *defect prediction, code smell, effort estimation*

- Useful baseline method for transfer learning

- Bellwethers mitigate conclusion instability

# ISSUE ADDRESSED

- The *Discover* phase is slow

```
def discover(datasets):
"Identify Bellwether Datasets"
  for data_1, data_2 in datasets:
    def train(data_1):
      "Construct quality predictor"
      return predictor
    def predict(data_1):
      "Predict for quality"
      return predictions
    def score(data_1, data_2):
      "Return accuracy of Prediction"
      return accuracy(train(data_1),\
      test(data_2))

"Return data with best prediction score"
```

# RESEARCH QUESTION

**RQ:** Can the time taken to find the bellwether dataset be reduced?

# INSPIRATION - RFE[6]

---

**Algorithm 1:** Recursive feature elimination

---

1.1 Tune/train the model on the training set using all predictors

1.2 Calculate model performance

1.3 Calculate variable importance or rankings

1.4 **for** *Each subset size $S_i$, $i = 1 \ldots S$* **do**

1.5     Keep the $S_i$ most important variables

1.6     [Optional] Pre–process the data

1.7     Tune/train the model on the training set using $S_i$ predictors

1.8     Calculate model performance

1.9     [Optional] Recalculate the rankings for each predictor

1.10 **end**

1.11 Calculate the performance profile over the $S_i$

1.12 Determine the appropriate number of predictors

1.13 Use the model corresponding to the optimal $S_i$

---

# INSPIRATION - BEETLE PAPER[4]

```
1   def FindBellwether(sources, step_size, budget, thres, ←
            lives):
2     while lives or cost > budget:
3       "Sample configurations"
4       sampled = list()
5       for source in sources:
6           "Sample step_size number of configurations"
7           sampled += source.sample(step_size)
8       "Get cost"
9       cost = get_cost(sampled)
10      "Evaluate pair−wise performances"
11      perf = get_perf(sampled)
12      "Remove non−bellwether environments"
13      sources=remove_non_bellwethers(sources, perf, ←
                thres)
14      "Loose life if no sources are removed"
15      if prev == len(sources): lives −= 1
16    "Return a bellwether"
17    return sources[argmin(perf)]
```

# OUR APPROACH

```
1  discoverBellwether(datasets, n_rep):
2      repeat n_rep times:
3              Initialize lives, x, step_size
4              while lives > 0:
5                      for dataset_1, dataset_2 in datasets:
6                              data1 = sample x% of dataset_1
7                              data2 = sample x% of dataset_2
8                              predictor = train(data1)
9                              G_score[dataset_1, dataset_2] = score(test(predictor, data2))
10                     rank and eliminate datasets based on median g-score
11                     if no dataset to eliminate:
12                             lives  = lives - 1
13                     x = x + step_size
14             for each dataset not eliminated:
15                     increment its wins
16     return dataset with most wins
```

9

# TWO VERSIONS FOR ELIMINATION

- Throw away bottom ⅓ of datasets sorted with median G-Score in each iteration

- Throw away datasets with median G-Score lower than threshold in each iteration
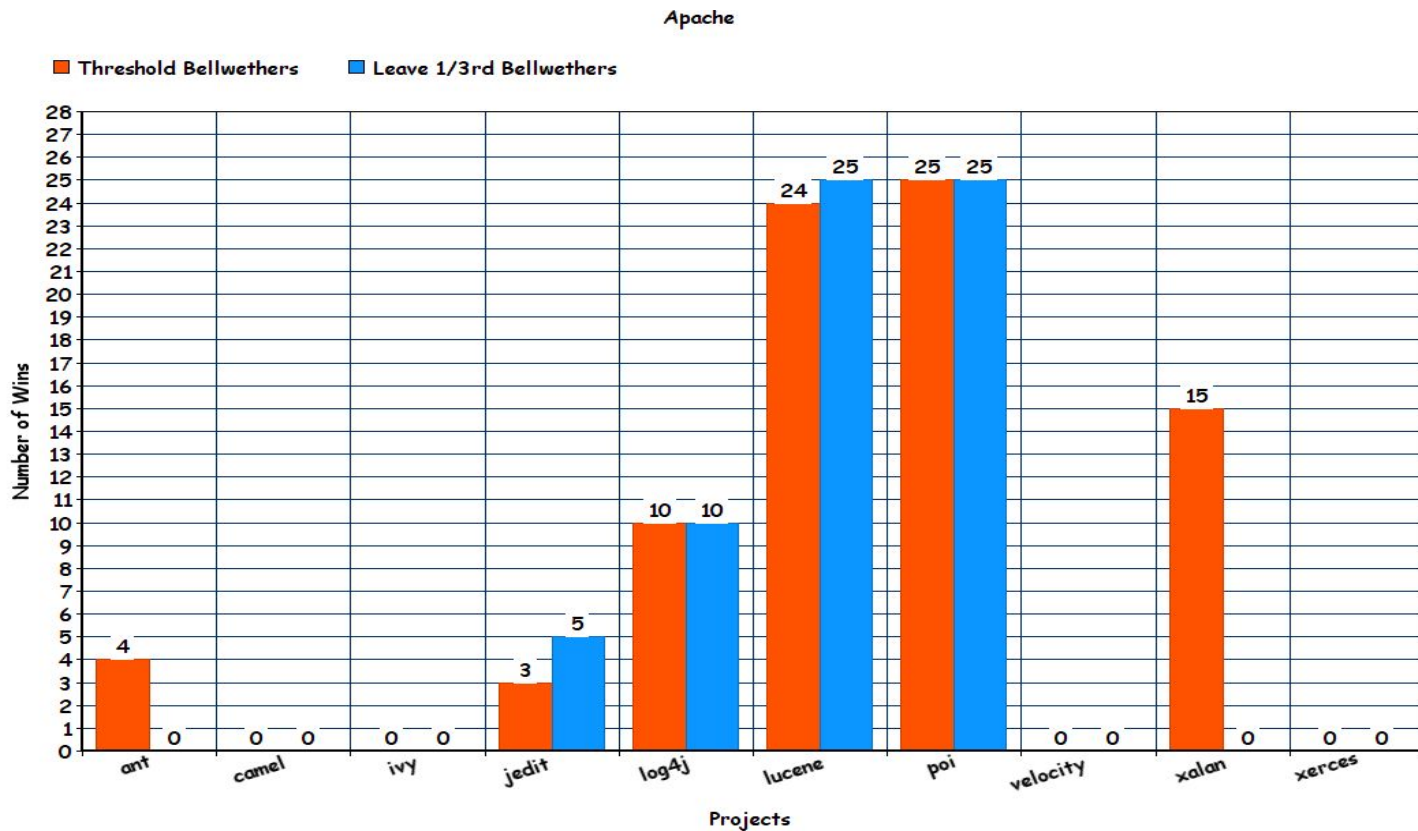
# DATA SET DESCRIPTION

- The AEEEM dataset was gathered by D'Amborse et al. [1], it contains 61 metrics, 5371 instances, 893 defective  instances.

- The RELINK community data was obtained from work by Wu et al. [2] , it contains 26 metrics, 649 total instances, 238 defective instances

- The Apache community data was gathered by Jureczko et al. [3]. This dataset contains records of the number of known defects for each class using a post-release bug tracking system. The classes are described in terms of 20 OO metrics, including CK metrics and McCabes complexity metrics.
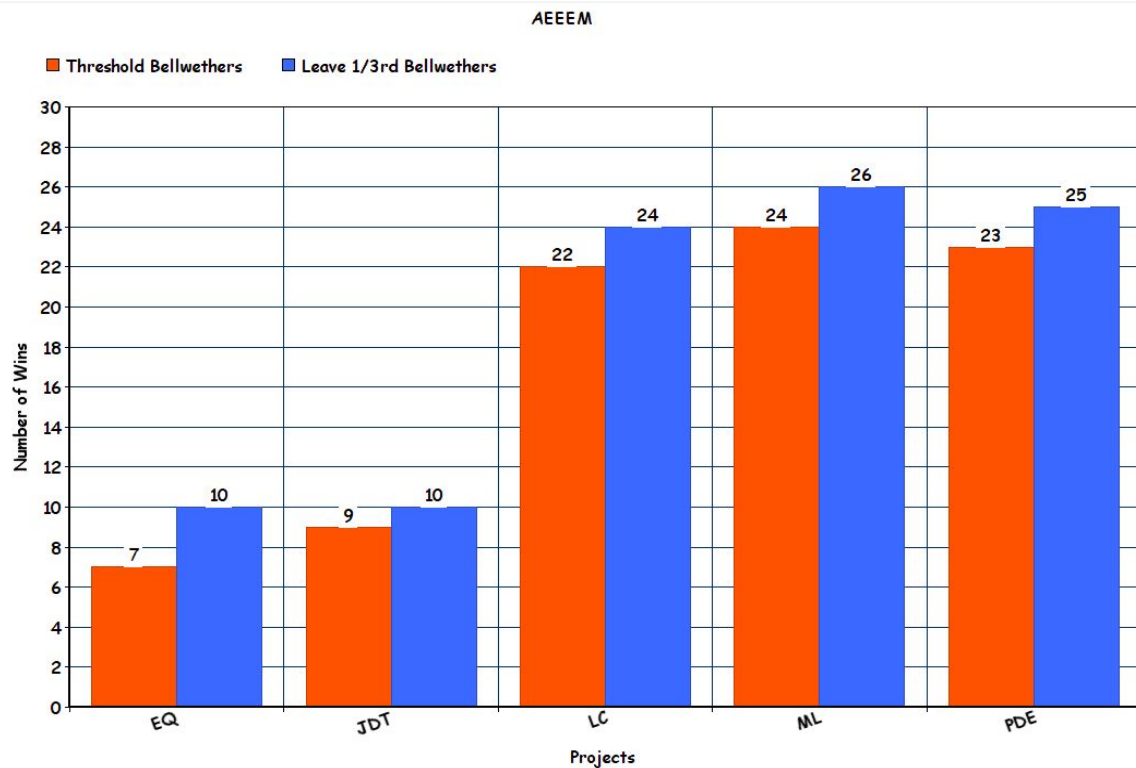
# EXPERIMENTAL SETUP

- Threshold value = 52
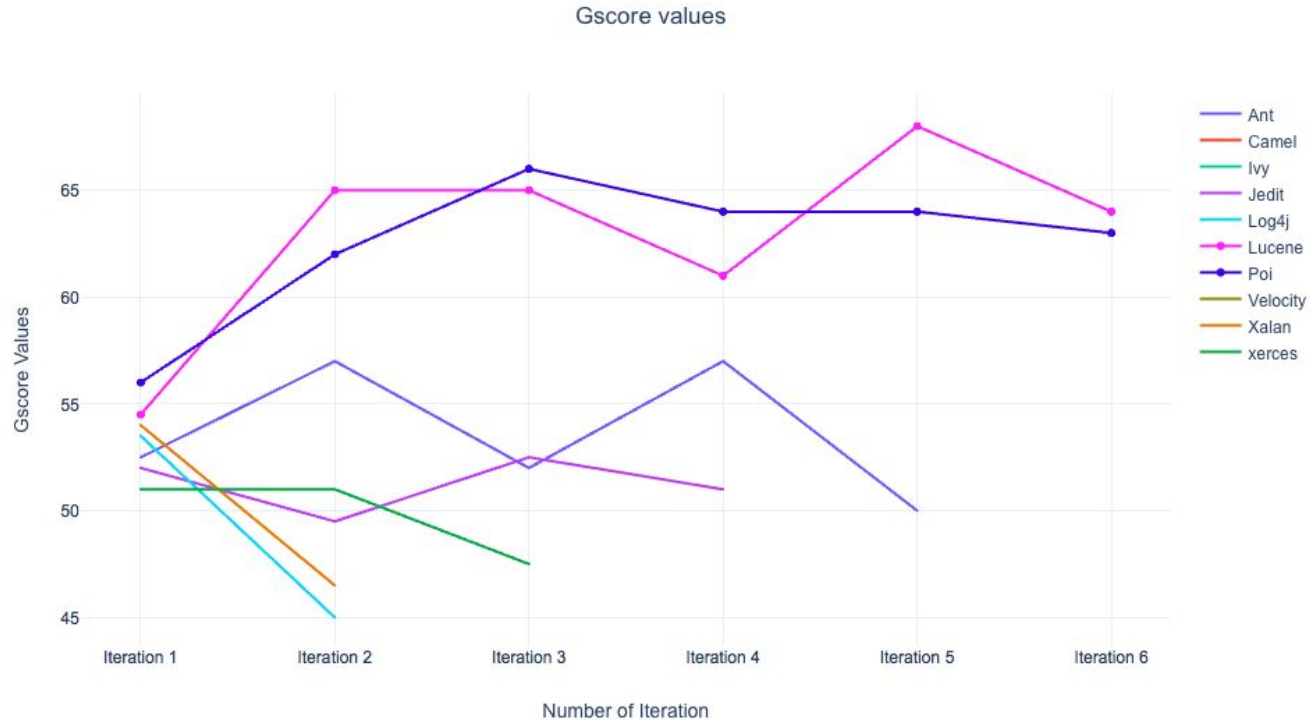- sample size = 0.25, incremented by 0.05
- Lives = 10

# RESULTS

# RESULTS

# G-SCORES

# FUTURE WORK

- Predicting the threshold values for optimal bellwethers

- Analysis of trade-off between different parameters

- Combining our method with Hoeffding Racing

# References

[1] D'Ambros, Marco, Michele Lanza, and Romain Robbes. "Evaluating defect prediction approaches: a benchmark and an extensive comparison." *Empirical Software Engineering 17*. 4-5 (2012): 531-577.

[2] Bavota, Gabriele. "Mining unstructured data in software repositories: current and future trends." *Software Analysis, Evolution, and Reengineering (SANER), 2016 IEEE 23rd International Conference on*. Vol. 5. IEEE, 2016.

[3] Jureczko, Marian, and Lech Madeyski. "Towards identifying software project clusters with regard to defect prediction." *Proceedings of the 6th International Conference on Predictive Models in Software Engineering*. ACM, 2010.

[4] Nair, Vivek, et al. "Transfer Learning with Bellwethers to find Good Configurations." *arXiv preprint arXiv:1803.03900* (2018).

[5] Krishna, Rahul, and Tim Menzies. "Bellwethers: A Baseline Method For Transfer Learning." *IEEE Transactions on Software Engineering* (2018).

[6] Max, "Recursive Feature Elimation" https://topepo.github.io/caret/recursive-feature-elimination.html#search

# THANK YOU!