# IDEA: ITERATIVE DICHOTOMIZATION ON EVERY ATTRIBUTE

## 1. TODO LIST

- Competency
    - Compare performance for logged vs. unlogged coordinates.
    - Compare performance for pruned vs. unpruned with lives=3 system. Performance should increase once we prune.
- Cosmetics
    - ~~Thinner lines for block outlines.~~
    - Only display clustered quadrants.
    - Try only showing the top 80% densest clusters.
- Performance Scores
    - Eliminate self testing by setting aside train/test sets.
- Clean up and profile code.
- Clustering
    - ~~Use a color scale to represent cluster scores.~~
    - ~~Graph legend showing colors with scores.~~
    - Mark (gray out) clusters which have a large delta between their performance scores (neighboring clusters which have high/low scores.)
    - Instead of using a 10% similarity rule in the GRIDCLUS algorithm [2], find the largest drop in block similarity and use that as a basis for similar.
      $\sum \frac{\sqrt{23}}{3} i_t$
- Takeaways from the Teak Experiment [1]
    - Prune subtrees using a decreasing performance rule with lives=3.
    - Re-cluster on remaining blocks. (Grid 2)
    - Test on Grid 2.
- Feature Subset Selection
    - What's best among the best clusters?
- Contrast Sets
    - What makes cluster X different from unclustered or neighboring clusters?
- Interactive Clustering
    - Import data.
    - Assign attributes.
    - Show statistics for active cluster / all clusters.

## References

[1] Jacky W. Keung Ekrem Kocaguneli, Tim Menzies. Teak: Learning better case selection strategies for analogy based software cost estimation. IEEE Transactions on Software Engineering, 6, 2007.

[2] E. Schikuta. Grid-clustering: An efficient hierarchical clustering method for very large data sets. Pattern Recognition, International Conference on, 2:101, 1996.