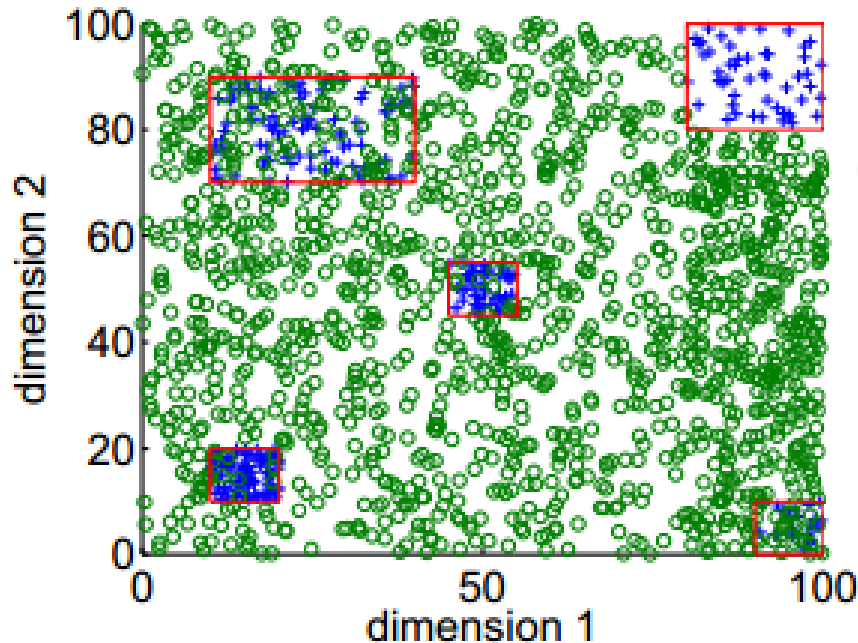# Box Drawings for Learning with Imbalanced Data

## Siong Thye Goh & Cynthia Rudin
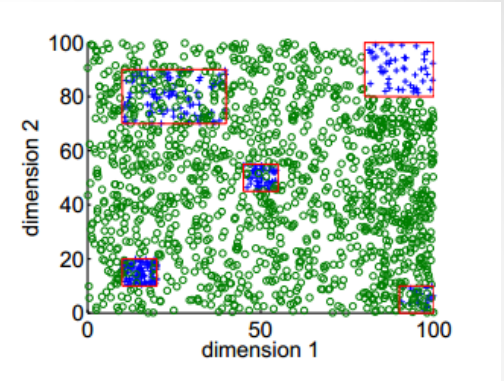## Massachusetts Institute of Technology

# Think of…

- A few positive examples in a sea of negative examples. For modeling rare events: Machine breakdown, etc.

A box drawing classifier is a union of axis-parallel rectangles.

# The usual way to do this...

Is by being greedy.



- Decision Tree (CART, C4.5 – Breiman 1984, Quinlan 1993)
  - Top down greedy: pick a features based on Gini Index or Information Gain, split data into two pieces, repeat. Prune afterwards.

- PRIM (Friedman, Fisher 1999)
  - Peel off subsets of data greedily, and if there is improvement, keep peeling off data. Occasionally put the data back.

This is too greedy for us.

# In fact, greed leads to stupid answers...

- It's often hard to get anything except a trivial model that always predicts the majority class.
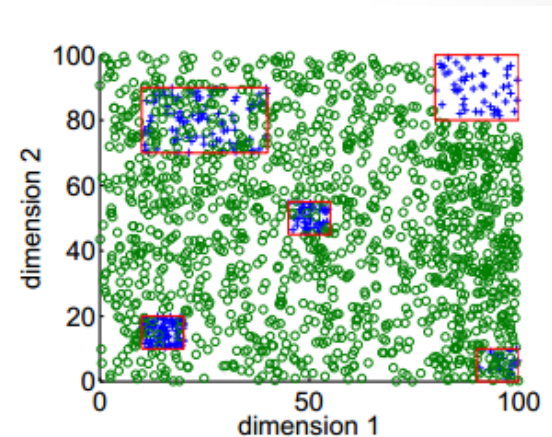
If our classifier says that all princesses have a thin waist, Princess Fiona from Shrek will feel unloved.

- This is true even if you fiddle like crazy with the algorithm's parameters. (We show some examples in the paper.)

- But since we are interested in the minority class, why don't we protect them?

# Our Approaches

- Approach 1: <span style="color:red">The Exact Boxes algorithm</span>
  - Optimize weighted accuracy, regularize by number of boxes
  - Mixed-Integer Programming formulation
  - Useful for not-huge problems, but solves exactly the problem we care about
  - Acts as a gold standard to compare with because it solves exactly the problem we want.

- Approach 2: <span style="color:red">The Fast Boxes algorithm</span>
  - Approximates the solution of Exact Boxes
  - Characterize (one class learning) before discriminating
  - Requires that features are continuous.

# Exact Boxes Algorithm

- Creates box drawing classifier
- A minority(positive) example is correctly classified if it resides within at least 1 box.
- A majority (negative) example is correctly classified if it does not reside in any box.
- Let *F* be the set of union of axis parallel rectangles.
- Here is Exact Boxes' objective:

$$\max_{f \in F} \sum_{i: y_i = 1} 1_{[f(x_i)=1]} + C_I \sum_{i: y_i = -1} 1_{[f(x_i)=-1]} - c_E(\text{\# of boxes})$$

Give less weight to negative examples

Simplicity/sparsity Favors fewer boxes.

# Exact Boxes Algorithm

$$\max_{l,\tilde{l},u,\tilde{u},w,z} \left[ -c_e K + \sum_{i \in S_+} z_i + c_I \sum_{i \in S_-} z_i \right] \text{ subject to}$$

$$x_{ij} - l_{jk} - v \le M\tilde{l}_{ijk}, \forall i \in S_+, \forall j, k$$

$$M(\tilde{l}_{ijk} - 1) + \epsilon \le x_{ij} - l_{jk} - v, \forall i \in S_+, \forall j, k$$

$$u_{jk} - v - x_{ij} \le M\tilde{u}_{ijk}, \forall i \in S_+, \forall j, k$$

$$M(\tilde{u}_{ijk} - 1) + \epsilon \le u_{jk} - x_{ij} - v, \forall i \in S_+, \forall j, k$$

$$\sum_{j=1}^n \tilde{l}_{ijk} + \sum_{j=1}^n \tilde{u}_{ijk} - 2n + 1 \le w_{ik}, \forall i \in S_+, \forall j, k$$

$$2nw_{ik} \le \sum_{j=1}^n \tilde{l}_{ijk} + \sum_{j=1}^n \tilde{u}_{ijk}, \forall i \in S_+, \forall j, k$$

$$\sum_{k=1}^K w_{ik} \le Kz_i, \forall i \in S_+, \forall k$$

$$z_i \le \sum_{k=1}^K w_{ik}, \forall i \in S_+, \forall k$$

$$l_{jk} - v - x_{ij} \le M\tilde{l}_{ijk}, \forall i \in S_-, \forall j, k$$

$$M(\tilde{l}_{ijk} - 1) + \epsilon \le l_{jk} - v - x_{ij}, \forall i \in S_-, \forall j, k$$

$$x_{ij} - u_{jk} - v \le M\tilde{u}_{ijk}, \forall i \in S_-, \forall j, k$$

$$M(\tilde{u}_{ijk} - 1) + \epsilon \le x_{ij} - u_{jk} - v, \forall i \in S_-, \forall j, k$$

$$\sum_{j=1}^n \tilde{l}_{ijk} + \sum_{j=1}^n \tilde{u}_{ijk} - 2n + 1 \le 2n(1 - w_{ik}),$$
$$\forall i \in S_-, \forall j, k$$

$$1 - w_{ik} \le \sum_{j=1}^n \tilde{l}_{ijk} + \sum_{j=1}^n \tilde{u}_{ijk}, \forall i \in S_-, \forall j, k$$

$$\sum_{k=1}^K w_{ik} \le K(1 - z_i), \forall i \in S_-, \forall k$$

$$1 - z_i \le \sum_{k=1}^K w_{ik}, \forall i \in S_-, \forall k$$

$$l_{jk} \le u_{jk}, \forall j, k.$$

- This is a mixed-integer *linear* programming formulation

- Works nicely for not-huge problems, acts as a gold standard

7

# Exact Boxes Algorithm

$$\max_{l,\tilde{l},u,\tilde{u},w,z} \left[ -c_e K + \sum_{i \in S_+} z_i + c_I \sum_{i \in S_-} z_i \right] \text{ subject to}$$

$$x_{ij} - l_{jk} - v \le M\tilde{l}_{ijk}, \forall i \in S_+, \forall j,k$$
$$M(\tilde{l}_{ijk} - 1) + \epsilon \le x_{ij} - l_{jk} - v, \forall i \in S_+, \forall j,k$$
$$u_{jk} - v - x_{ij} \le M\tilde{u}_{ijk}, \forall i \in S_+, \forall j,k$$
$$M(\tilde{u}_{ijk} - 1) + \epsilon \le u_{jk} - x_{ij} - v, \forall i \in S_+, \forall j,k$$
$$\sum_{j=1}^{n} \tilde{l}_{ijk} + \sum_{j=1}^{n} \tilde{u}_{ijk} - 2n + 1 \le w_{ik}, \forall i \in S_+, \forall j,k$$
$$2nw_{ik} \le \sum_{j=1}^{n} \tilde{l}_{ijk} + \sum_{j=1}^{n} \tilde{u}_{ijk}, \forall i \in S_-$$
$$\sum_{k=1}^{K} w_{ik} \le K z_i, \forall i \in S_-$$
$$z_i \le \sum_{k=1}^{K} w_{ik}$$
$$l_{jk} - v \le \ldots_{jk}, \forall i \in S_-, \forall j,k$$
$$M\tilde{\ldots} \le l_{jk} - v - x_{ij}, \forall i \in S_-, \forall j,k$$
$$u_{jk} - v \le M\tilde{u}_{ijk}, \forall i \in S_-, \forall j,k$$
$$(\tilde{u}_{ijk} - 1) + \epsilon \le x_{ij} - u_{jk} - v, \forall i \in S_-, \forall j,k$$
$$\sum_{j=1}^{n} \tilde{l}_{ijk} + \sum_{j=1}^{n} \tilde{u}_{ijk} - 2n + 1 \le 2n(1 - w_{ik}),$$
$$\forall i \in S_-, \forall j,k$$
$$1 - w_{ik} \le \sum_{j=1}^{n} \tilde{l}_{ijk} + \sum_{j=1}^{n} \tilde{u}_{ijk}, \forall i \in S_-, \forall j,k$$
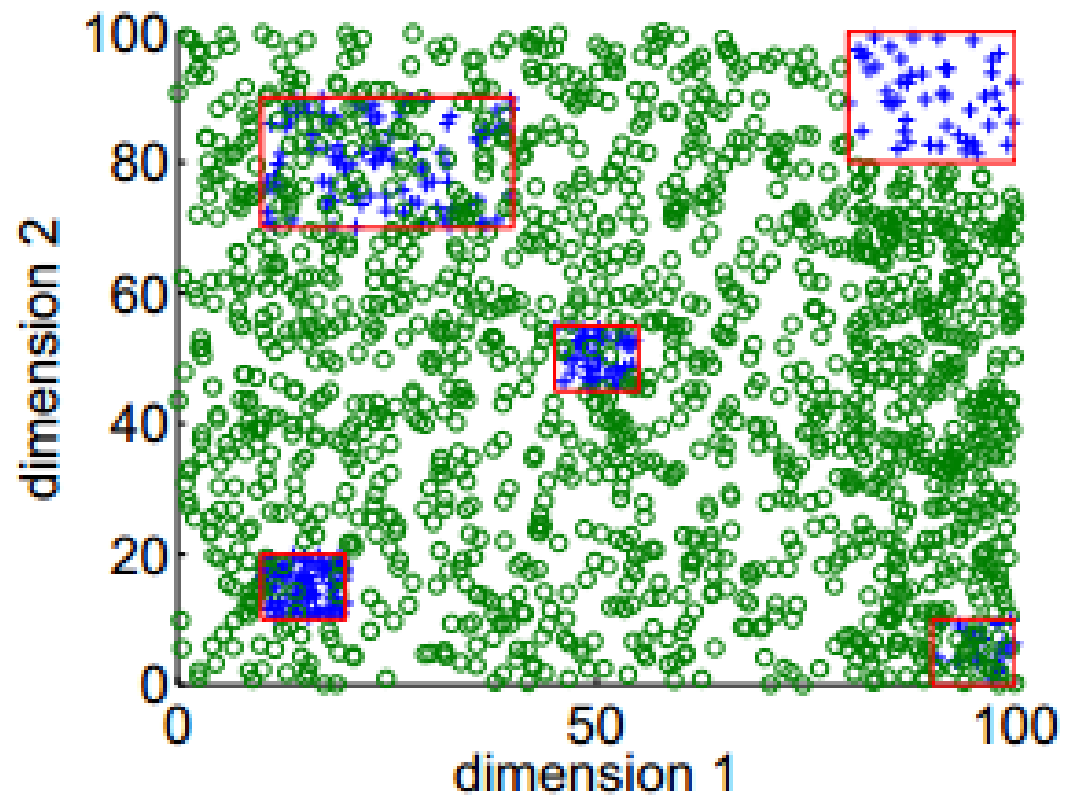$$\sum_{k=1}^{K} w_{ik} \le K(1 - z_i), \forall i \in S_-, \forall k$$
$$1 - z_i \le \sum_{k=1}^{K} w_{ik}, \forall i \in S_-, \forall k$$
$$l_{jk} \le u_{jk}, \forall j,k.$$
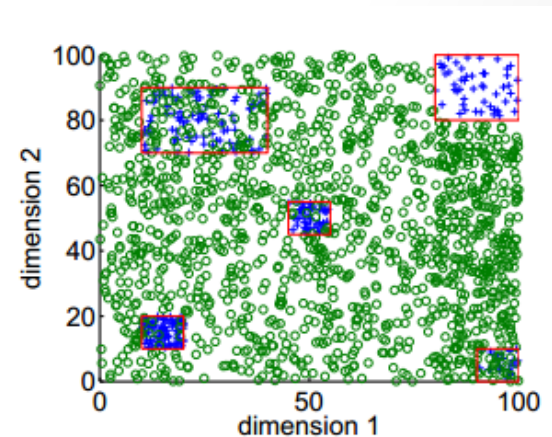
Life is too short for details

- This is a mixed-integer *linear* programming formulation

- Works nicely for not-huge problems, acts as a gold standard

8

# Our Approaches

- Approach 1: <span style="color:red">The Exact Boxes algorithm</span>
  - o Optimize accuracy, regularize by number of boxes
  - o Mixed-Integer Programming formulation
  - o Useful for not-huge problems, but solves exactly the problem we care about
  - o Acts as a gold standard to compare with because it solves exactly the problem we want.

- Approach 2: <span style="color:red">The Fast Boxes algorithm</span>
  - o Approximates the solution of Exact Boxes
  - o Characterize (one class learning) before discriminating
  - o Requires that features are continuous.

# Fast Boxes Algorithm

- Three main stages:

1. Clustering.
   o Characterization of positive data only.
2. Dividing space stage.
   o Decide which points will influence each boundary.
3. Boundary expansion stage.
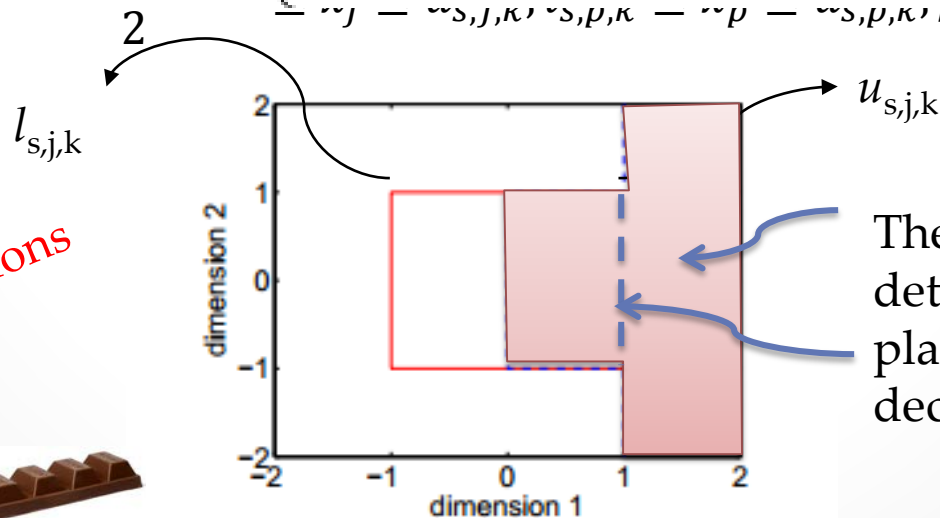   o Expand each box to give priority to the positive data.

# Fast Boxes Algorithm

- Three main stages:

1. Clustering.
   - Characterization of positive data only.
   - This is just K means, which assumes continuous features.
2. Dividing space stage.
   - Decide which points will influence each boundary.
3. Boundary expansion stage.
   - Expand each box to give priority to the positive data.

# Fast Boxes Dividing Space Stage

- Each boundary is determined by 1d supervised learning using data close to the boundary.
- We have analytical solution for each decision boundary.

$$X_{l,j,k} = \left\{ x\colon x_j \leq l_{s,j,k} \right\} \cup \left\{ x\colon l_{s,j,k} \leq x_j \leq \frac{l_{s,j,k} + u_{s,j,k}}{2} \,, l_{s,p,k} \leq x_p \leq u_{s,p,k}, p \neq j \right\}$$

$$X_{u,j,k} = \left\{ x\colon x_j \geq u_{s,j,k} \right\} \cup \left\{ x\colon \frac{l_{s,j,k} + u_{s,j,k}}{2} \leq x_j \leq u_{s,j,k}, l_{s,p,k} \leq x_p \leq u_{s,p,k}, p \neq j \right\}$$

$l_{s,j,k}$

$u_{s,j,k}$

*Analytical solutions are better than chocolate.*

The points in here determine the placement of this decision boundary.

# Boundary Expansion Stage

- We push decision boundaries for each box outwards as a form of regularization.
- This is a final optional push to almost the next negative point.



Analytical solutions are even better than cupcakes.

# Summary of Fast Boxes Algorithm

1. Normalize the features to be between -1 and 1.
2. Cluster the minority data into K clusters (take standard deviation into consideration)
3. Construct the minimal enclosing box for each cluster.
4. Pick which data participate in each decision boundary calculation.
5. Compute decision boundaries analytically and separately
6. Perform expansion (optional)
7. Un-normalize by rescaling the features back.

# Experiments

# Baseline Algorithms

- Non-interpretable Models:
  - o Logistic Regression.
  - o SVM with RBF kernels
  - o Random Forests
  - o AdaBoost (boosted trees)
- Interpretable Models:
  - o CART
  - o C4.5
  - o C5.0
  - o Hellinger Distance Decision Tree

déjà vu, didn't you see this in the previous talk?

# Performance analysis

- Compare Fast Boxes with baseline algorithms.
- It is often (but not always) one of the best performers.
- However the solution is often more interpretable.

# Experimental Comparison



Excessive amount of experiments that show:

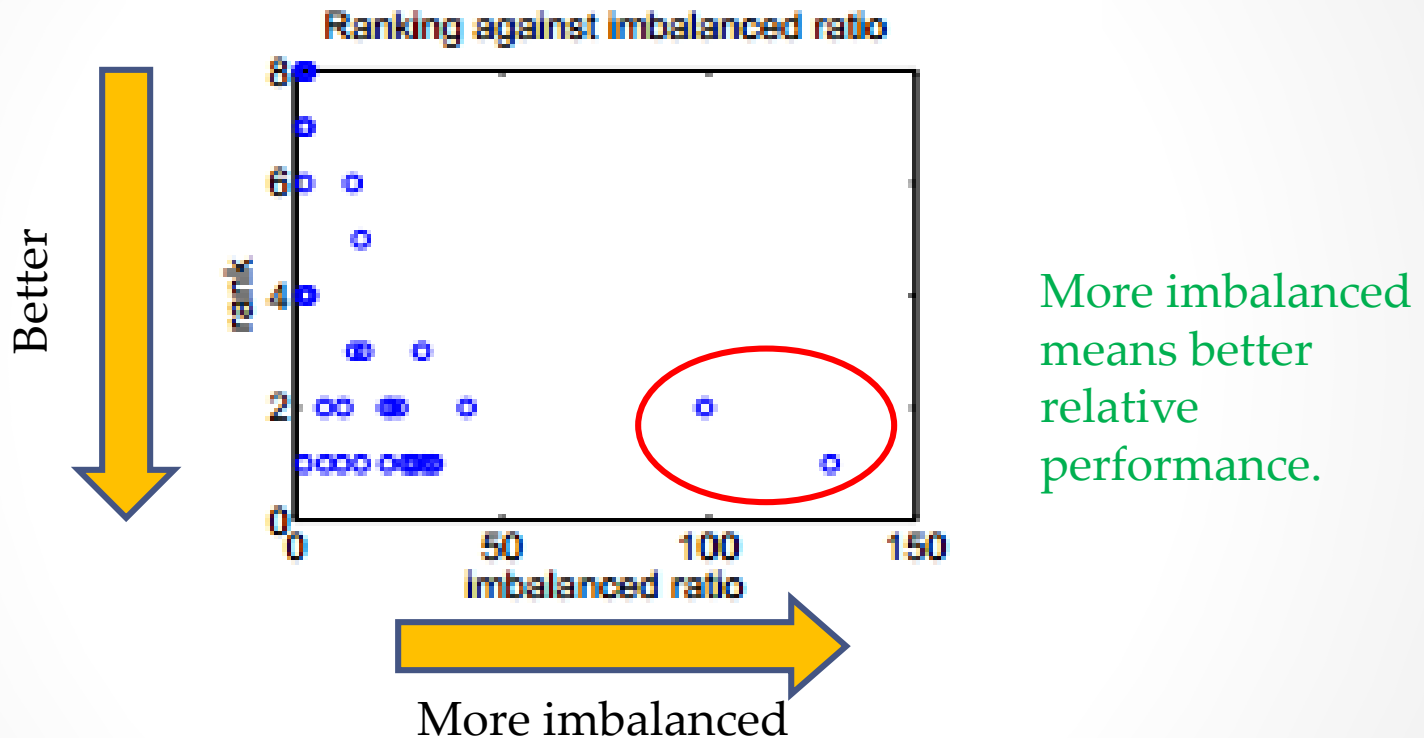1) Fast Boxes has the necessary complexity to produce high accuracy results (like SVM)

2) Interpretability doesn't hurt a bit.

# When does Fast Boxes perform well? When data are more imbalanced.



- Each dot represents a dataset we tried, and the rank of Fast Boxes as compared with the other methods.
- Bottom line: if the data are more imbalanced, Fast Boxes does better.

20

# When does Fast Boxes perform well?
# When data are more imbalanced.



**Better** (vertical arrow pointing down)

Ranking against imbalanced ratio

**More imbalanced** (horizontal arrow pointing right)

More imbalanced means better relative performance.

- Each dot represents a dataset we tried, and the rank of Fast Boxes as compared with the other methods.
- Bottom line: if the data are more imbalanced, Fast Boxes
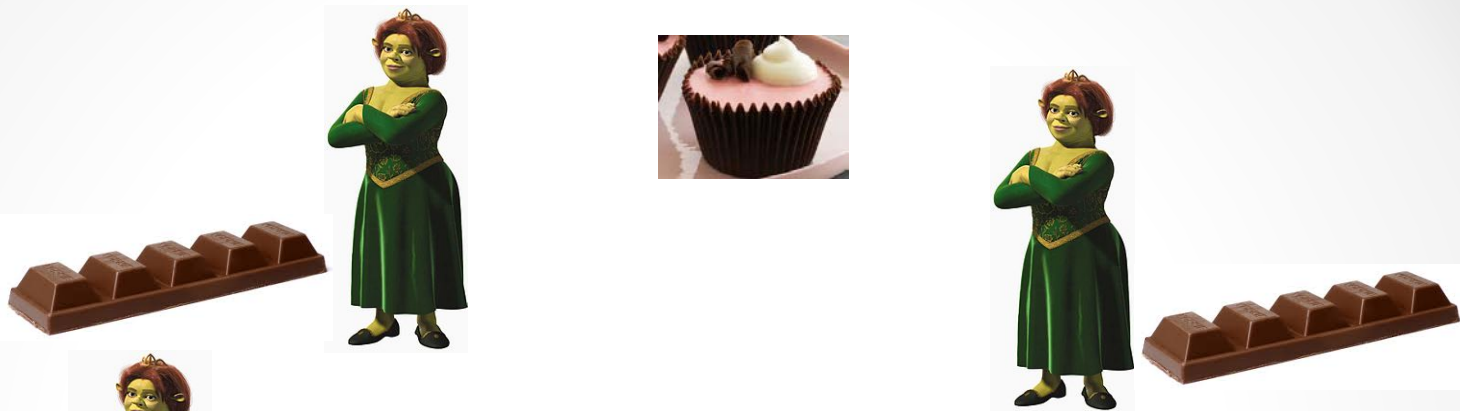- does better.

# Exact Boxes is frequently one of the Best Performers

| Data | Best Performance | Fast Boxes | Exact Boxes | Exact Boxes ranking |
|------|------------------|------------|-------------|---------------------|
| vehicle2 | 0.9496 (0.015) | 0.9191 (0.0242) | 0.9496 (0.015) | 1 |
| haberman | 0.6699 (0.0276) | 0.5290 (0.0265) | **0.6632** (0.0303) | 2 |
| yeast1 | 0.7641 (0.0133) | 0.5903 (0.0286) | 0.7392 (0.0172) | 2 |
| glass0 | 0.8312 (0.0345) | 0.7937 (0.0212) | 0.7977 (0.0421) | 2 |
| iris0 | 1 (0) | 1 (0) | **1** (0) | 1 |
| wisconsin | 0.9741 (0.0075) | 0.8054 (0.1393) | **0.9726** (0.0079) | 2 |
| ecoli01 | 0.9840 (0.0105) | 0.9433 (0.0300) | **0.9839** (0.0109) | 2 |
| glass1 | 0.7922 (0.0377) | 0.6654 (0.0356) | **0.7922** (0.0337) | 1 |

- Restricting the algorithm to produce a box drawing classifier does not generally seem to hinder performance.
- Selecting optimal K is difficult for Exact Boxes. Optimal K for Fast Boxes is used instead.

# Summary

- Exact Boxes
  - Mixed integer Programming

- Fast Boxes
  - Characterize-then-discriminate approach.

- Take away: Aim for both interpretability and accuracy, because you can often get both.

Thank you!