# CS60020: Foundations of Algorithm Design and Machine Learning

Sourangshu Bhattacharya

# ENSEMBLE METHODS

# Rationale for Ensemble Learning

- No Free Lunch thm: There is no algorithm that is always the most accurate

- Generate a group of <span style="color:red">base-learners</span> which when combined have higher accuracy

- Different learners use different
  - Algorithms
  - Parameters
  - Representations (Modalities)
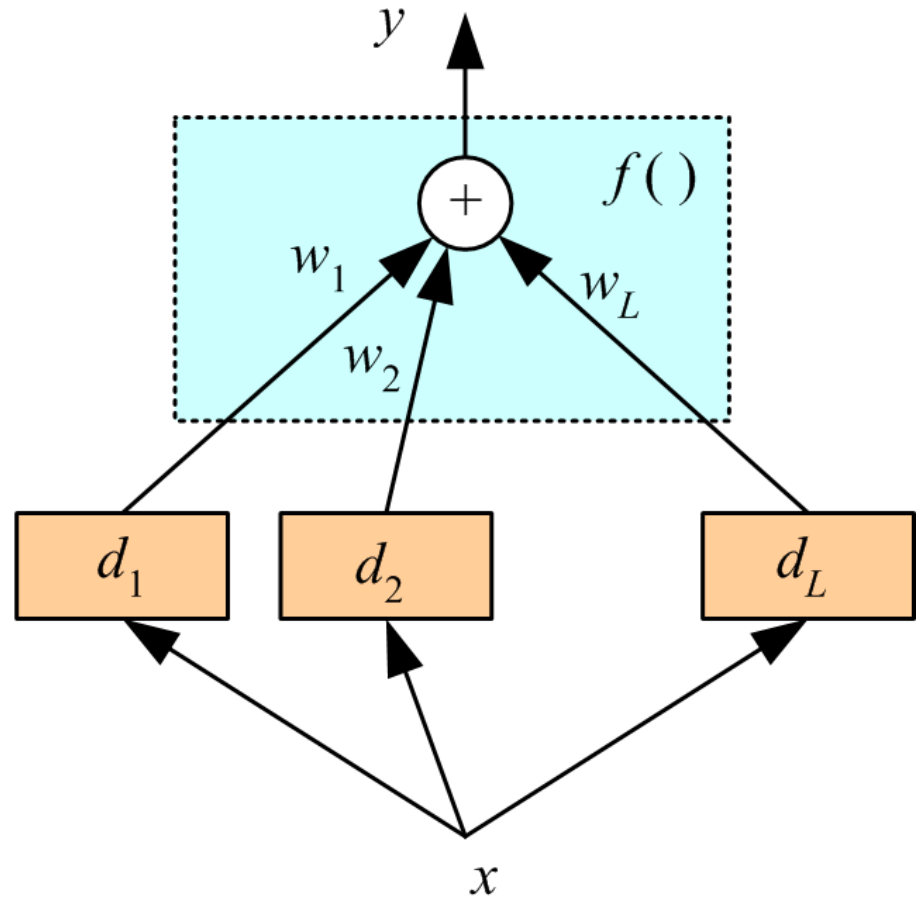  - Training sets
  - Subproblems

# Voting

- Linear combination

$$y = \sum_{j=1}^{L} w_j d_j$$

$$w_j \geq 0 \ \text{ and } \ \sum_{j=1}^{L} w_j = 1$$

- Classification

$$y_i = \sum_{j=1}^{L} w_j d_{ji}$$



4

# BAGGING

# Ensemble methods

- A single decision tree does not perform well

- But, it is super fast

- What if we learn multiple trees?

We need to make sure they do not all just learn the same

# Bagging

If we split the data in random different ways, decision trees give different results, **high variance.**

**Bagging: B**ootstrap **agg**regat**ing** is a method that result in low variance.

If we had multiple realizations of the data (or multiple samples) we could calculate the predictions multiple times and take the average of the fact that averaging  multiple onerous estimations produce less uncertain results

# Bagging

Say for each sample *b*, we calculate $f^b(x)$, then:

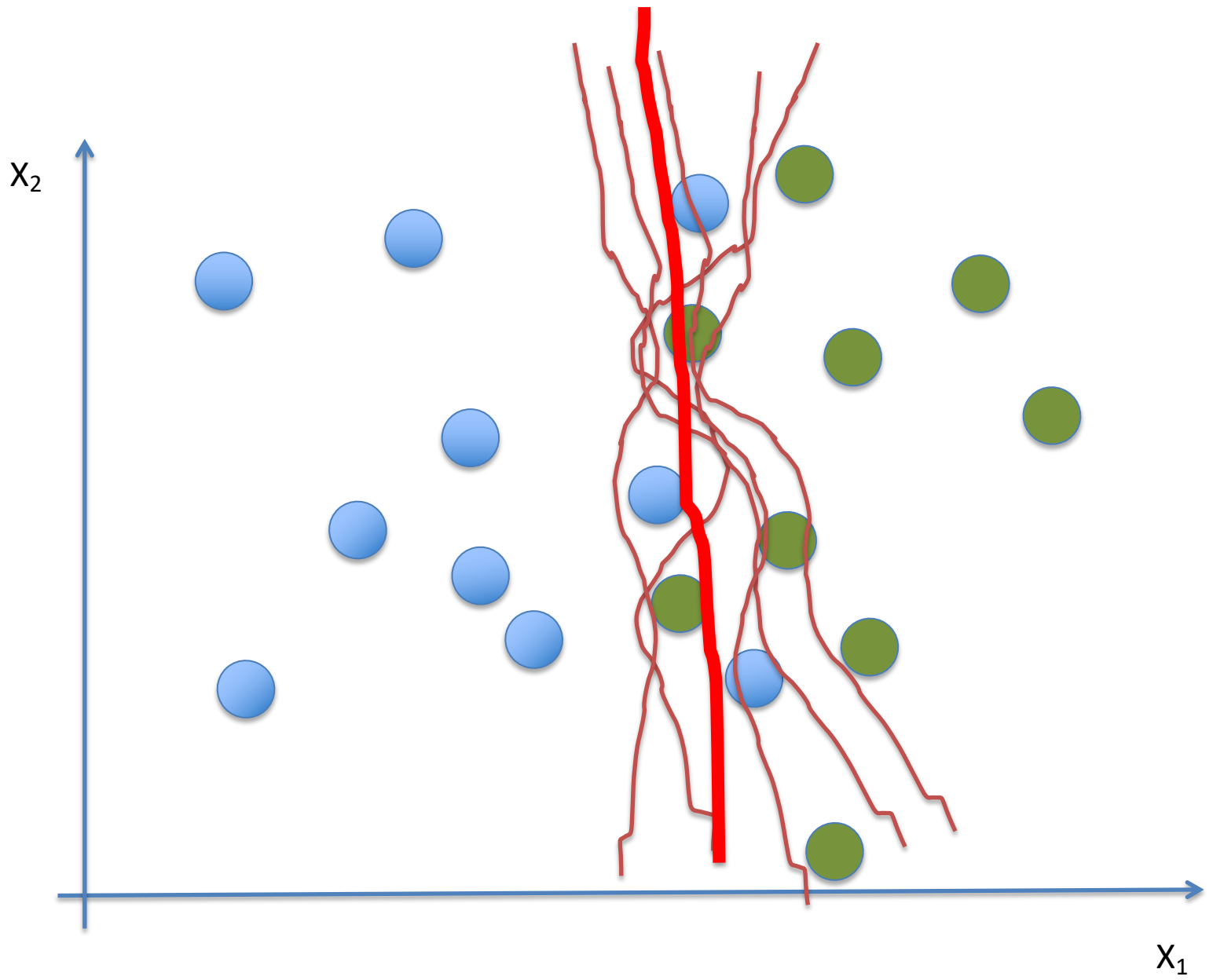$$\hat{f}_{avg}(x) = \frac{1}{B} \sum_{b=1}^{B} \hat{f}^b(x)$$

How?

---

**Bootstrap**

Construct B (hundreds of) trees (no pruning)

Learn a classifier for each bootstrap sample and average them
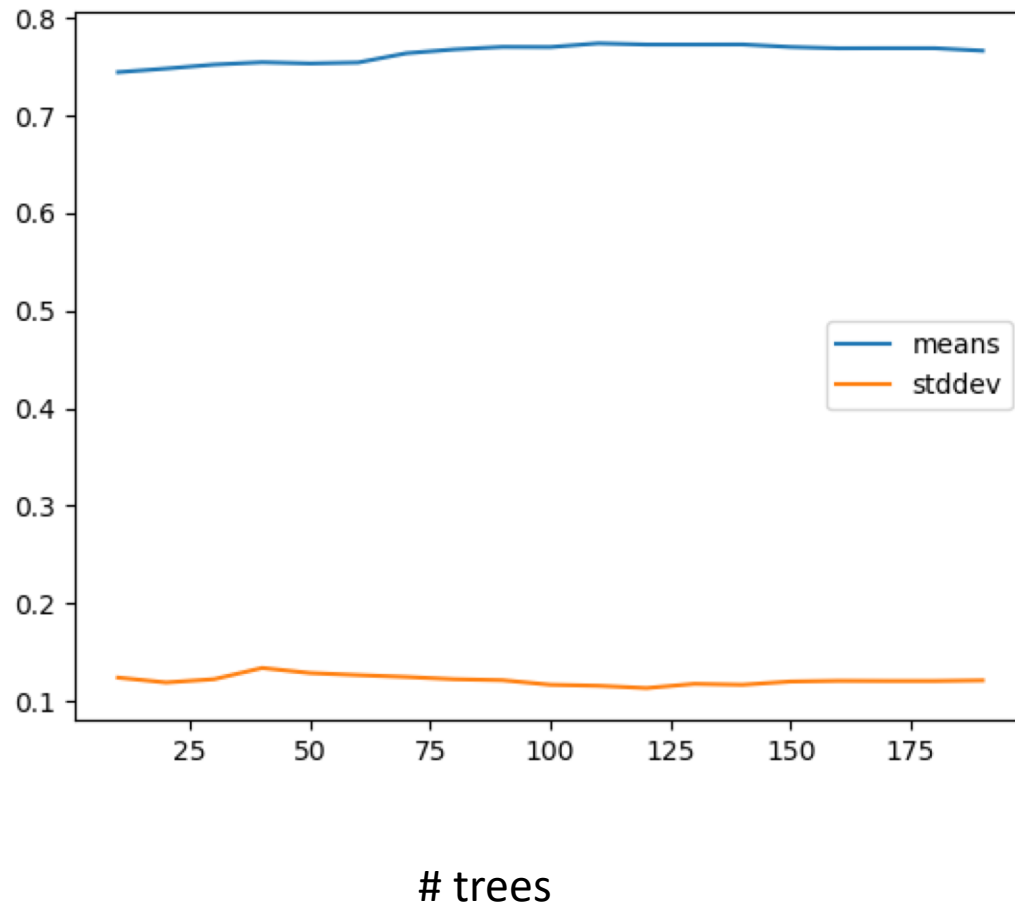
Very effective

# Example

- Pima Indians Diabetes data:
- 768 examples, 8 features
- Features: 'preg', 'plas', 'pres', 'skin', 'test', 'mass', 'pedi', 'age'

# Example

```python
for num_trees in range(start,stop,step):
    kfold = model_selection.KFold(n_splits=50,\
                                  random_state=int(time.time()))
    model = BaggingClassifier(base_estimator=cart,\
                              n_estimators=num_trees,\
                              random_state=seed)
    model = RandomForestClassifier(num_trees,\
                                   criterion="gini",\
                                   max_features=5,\
                                   random_state=int(time.time()))
    results = model_selection.cross_val_score(model, X, Y, cv=kfold)
    print(results.mean(), statistics.stdev(results))
    means.append(results.mean())
    stds.append(statistics.stdev(results))
```

# Example



# trees

# Out-of-Bag Error Estimation

- No cross validation?

- Remember, in bootstrapping we sample with replacement, and therefore **not all observations are used for each bootstrap sample**. On average 1/3 of them are not used!

- We call them out-of-bag samples (OOB)

-  We can predict the response for the *i-th* observation using each of the trees in which that observation was OOB and do this for *n* observations

-  Calculate overall OOB MSE or classification error

# Bagging

- Reduces overfitting (variance)

- Normally uses one type of classifier

- Decision trees are popular

- Easy to parallelize

# Bagging - issues

Each tree is identically distributed (i.d.)

➔ the expectation of the average of $B$ such trees is the same as the expectation of any one of them

➔the bias of bagged trees is the same as that of the individual trees

i.d. and not i.i.d

# Bagging - issues

An average of $B$ i.i.d. random variables, each with variance $\sigma^2$, has variance: $\sigma^2/B$

If i.d. (identical but not independent) and pair correlation $\rho$ is present, then the variance is:

$$\rho\,\sigma^2 + \frac{1-\rho}{B}\sigma^2$$

As $B$ increases the second term disappears but the first term remains

Why does bagging generate correlated trees?

# Bagging - issues

Suppose that there is one very strong predictor in the data set, along with a number of other moderately strong predictors.

Then all bagged trees will select the strong predictor at the top of the tree and therefore all trees will look similar.

How do we avoid this?

# RANDOM FORESTS

# Random Forests

As in bagging, we build a number of decision trees on bootstrapped training samples each time a split in a tree is considered, a random sample of $m$ predictors is chosen as split candidates from the full set of p predictors.

Note that if $m = p$, then this is bagging.

# Random Forests

Random forests are popular. Leo Breiman's and Adele Cutler maintains a random forest website where the software is freely available, and of course it is included in every ML/STAT package

http://www.stat.berkeley.edu/~breiman/RandomForests/

# Random Forests Algorithm

For b = 1 to B:

    (a) Draw a bootstrap sample Z∗ of size $N$ from the training data.

    (b) Grow a random-forest tree to the bootstrapped data, by recursively repeating the following steps for each terminal node of the tree, until the minimum node size $n_{min}$ is reached.

        i. Select $m$ variables at random from the $p$ variables.

        ii. Pick the best variable/split-point among the $m$.

        iii. Split the node into two daughter nodes.

Output the ensemble of trees.

To make a prediction at a new point $x$ we do:

    For regression: average the results

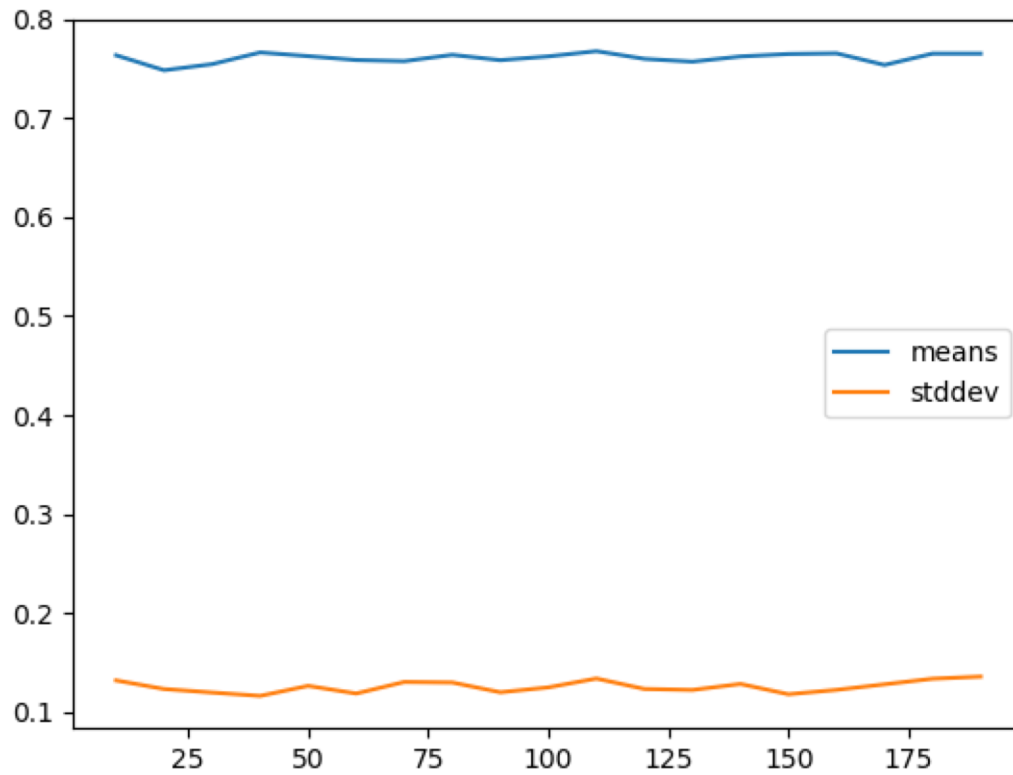    For classification: majority vote

# Random Forests Tuning

The inventors make the following recommendations:

- For classification, the default value for $m$ is $\sqrt{p}$ and the minimum node size is one.

- For regression, the default value for m is $p/3$ and the minimum node size is five.

In practice the best values for these parameters will depend on the problem, and they should be treated as tuning parameters.

Like with Bagging, we can use OOB and therefore RF can be fit in one sequence, with cross-validation being performed along the way. Once the OOB error stabilizes, the training can be terminated.
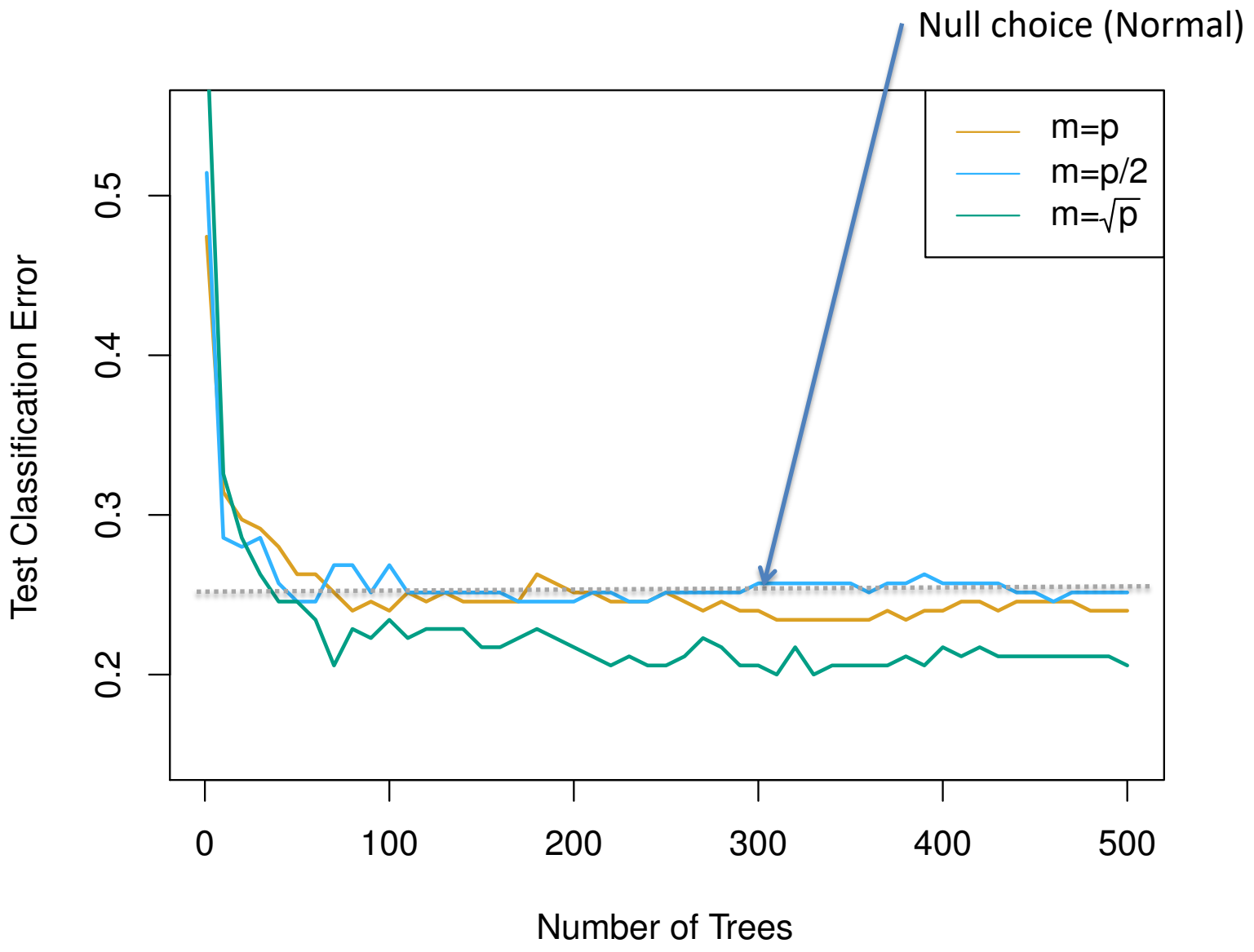
# Example

# Example

- 4,718 genes measured on tissue samples from 349 patients.

- Each gene has different expression

- Each of the patient samples has a qualitative label with 15 different levels: either normal or 1 of 14 different types of cancer.

Use random forests to predict cancer type based on the 500 genes that have the largest variance in the training set.

# Random Forests Issues

When the number of variables is large, but the fraction of relevant variables is small, random forests are likely to perform poorly when $m$ is small

Why?

Because:

At each split the chance can be small that the relevant variables will be selected

For example, with 3 relevant and 100 not so relevant variables the probability of any of the relevant variables being selected at any split is ~0.25

# Can RF overfit?

Random forests "cannot overfit" the data wrt to number of trees.

Why?

The number of trees, $B$ does not mean increase in the flexibility of the model