# 3050571 Practical Clin Data Sci

## Session 10: Machine learning framework
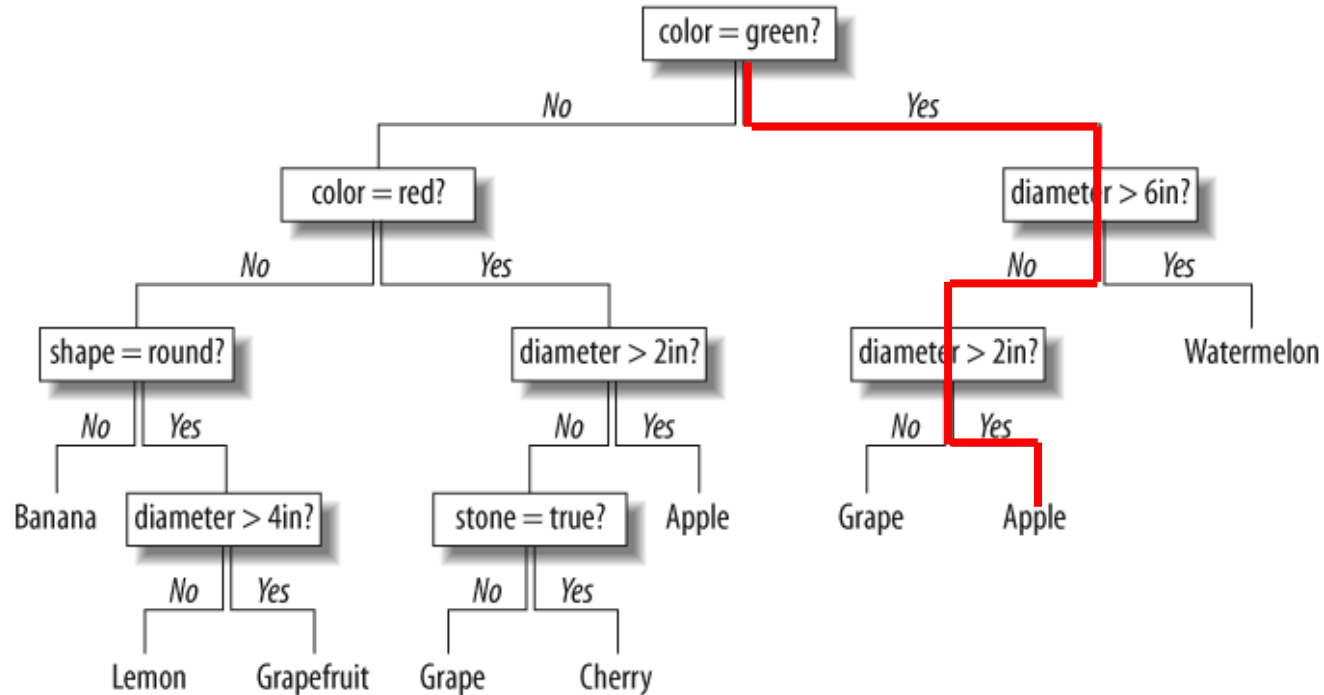
February 15, 2024

**Sira Sriswasdi, PhD**
- Research Affairs
- Center of Excellence in Computational Molecular Biology (CMB)
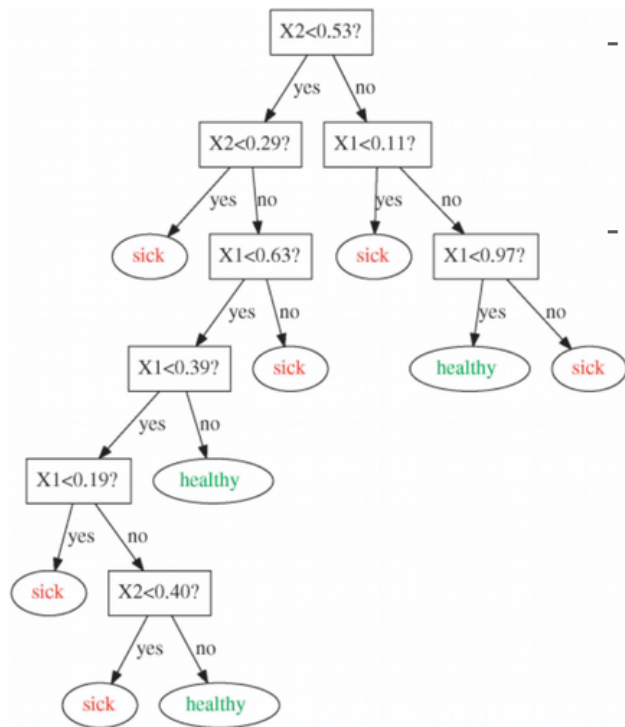- Center for Artificial Intelligence in Medicine (CU-AIM)
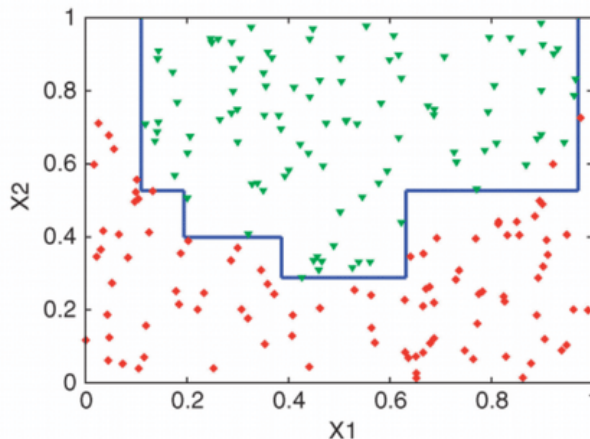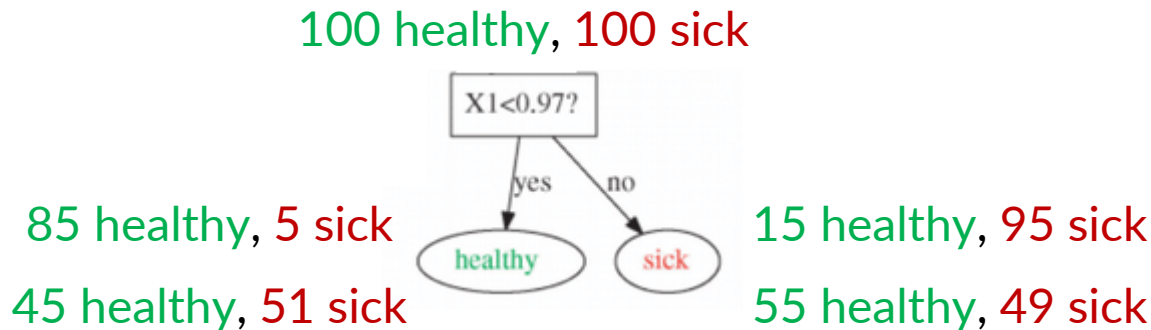
# Decision tree

# Decision tree



color = green?
- No → color = red?
  - No → shape = round?
    - No → Banana
    - Yes → diameter > 4in?
      - No → Lemon
      - Yes → Grapefruit
  - Yes → diameter > 2in?
    - No → stone = true?
      - No → Grape
      - Yes → Cherry
    - Yes → Apple
- Yes → diameter > 6in?
  - No → diameter > 2in?
    - No → Grape
    - Yes → Apple
  - Yes → Watermelon

Source: Programming collective intelligence by Toby Segaran

# Decision tree behaviors



- Each decision is a threshold on each feature
  - Piecewise linear
  - Parallel to an axis
- Good for criteria-based classification

Miller, C. "Screening meter data: Characterization of temporal energy data from large groups of non-residential buildings"

# Splitting quality

100 healthy, 100 sick



85 healthy, 5 sick

15 healthy, 95 sick

45 healthy, 51 sick

55 healthy, 49 sick

- Gini impurity: $\sum p(1-p)$
- Entropy: $-\sum p \ln(p)$
    - Minimal at p = 0 or 1 → Perfect split
    - Maximal at p = 0.5 → 50-50 split
- Search for feature and cutoff that yield lowest impurity or entropy

# Control mechanisms for tree building
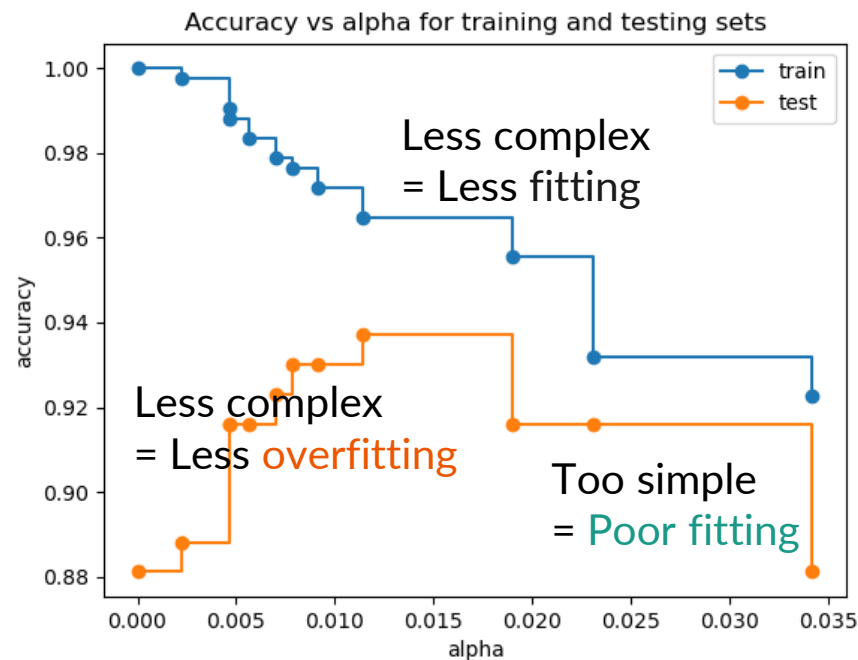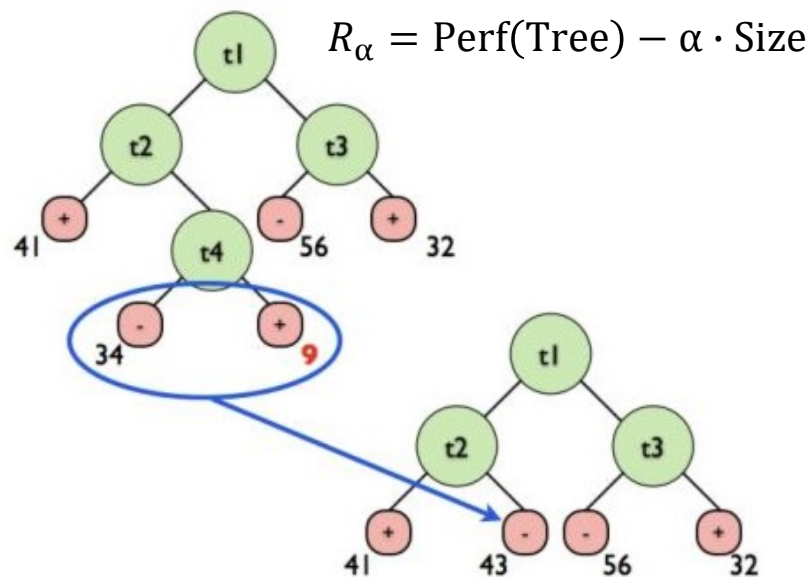
1. Too few samples to make a split

X1<0.97?

yes          no

healthy          sick

3. Impurity or entropy does not change much after the split

2. Too few samples on either branch

- Limit the tree size
- Limit the improvement in quality
- Limit the number of samples that support a split

# Tree pruning (post-processing)

$$R_\alpha = \text{Perf}(\text{Tree}) - \alpha \cdot \text{Size}$$



Less complex = Less fitting

Less complex = Less overfitting

Too simple = Poor fitting

Patel, N. and Upadhyay, S. "Study of Various Decision Tree Pruning Methods with their Empirical Comparison in WEKA" IJCA 2012

Image from scikit-learn,org

# Regularization on features

- **Linear model**: $\hat{y}_i = b_0 + b_1 x_{i,1} + \cdots + b_n x_{i,n}$
    - LASSO

- **Tree model**:
    - Repeatedly using the same feature
    - Early decision affects the rest

- Feature bagging
    - Look at only *N* features at each step
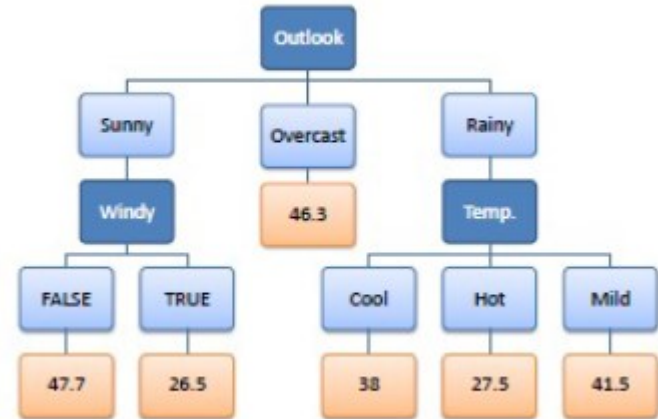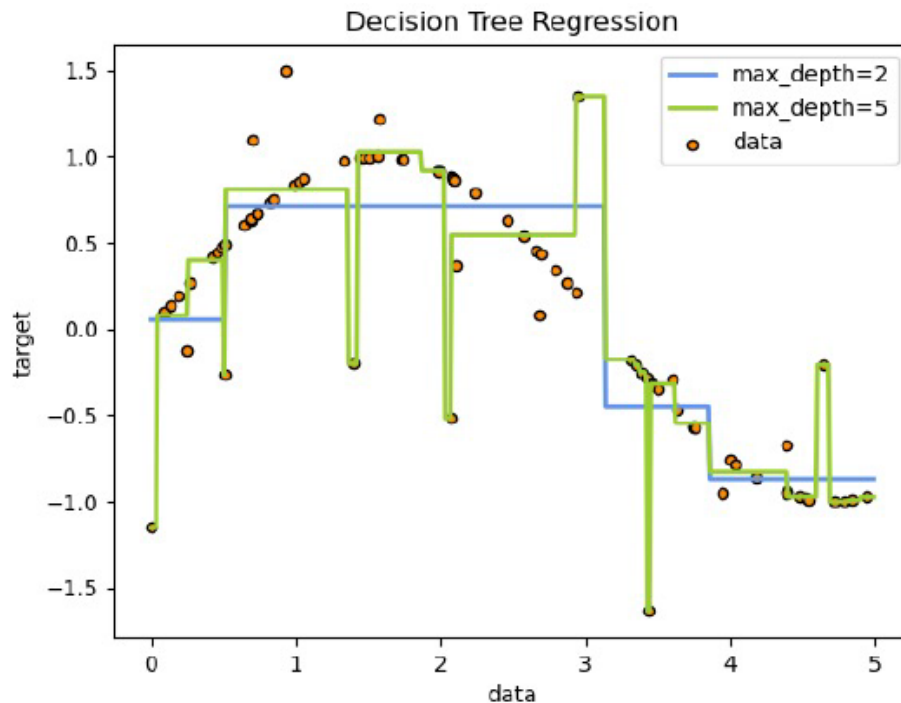    - Force model to use diverse features

# Decision tree for regression



Image from saedsayad.com

- Predict an average of samples in the same branch
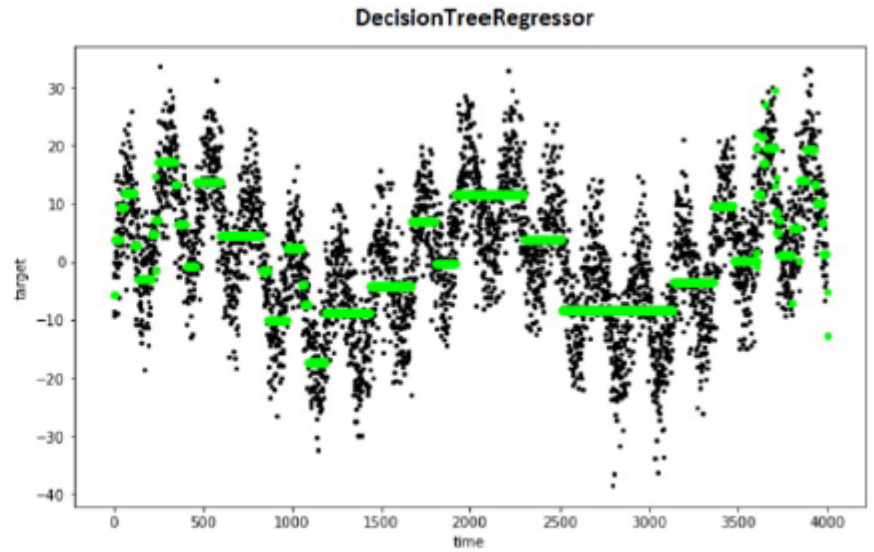
# Decision tree is a piecewise constant function
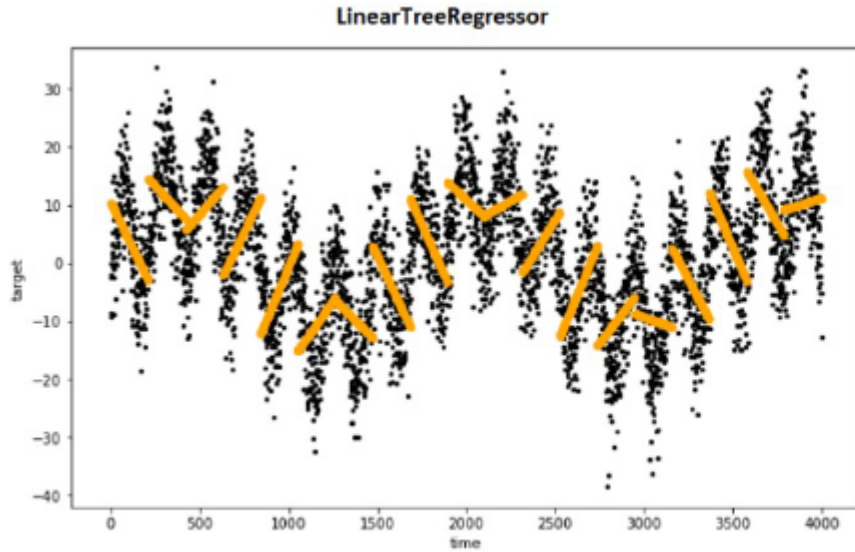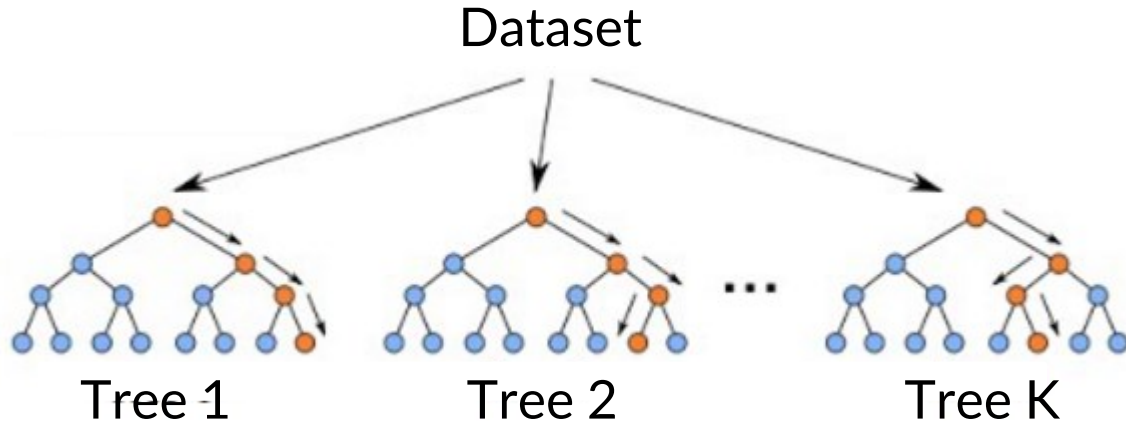


Image from scikit-learn.org

# Linear-Tree model



Image from towarddatascience.com

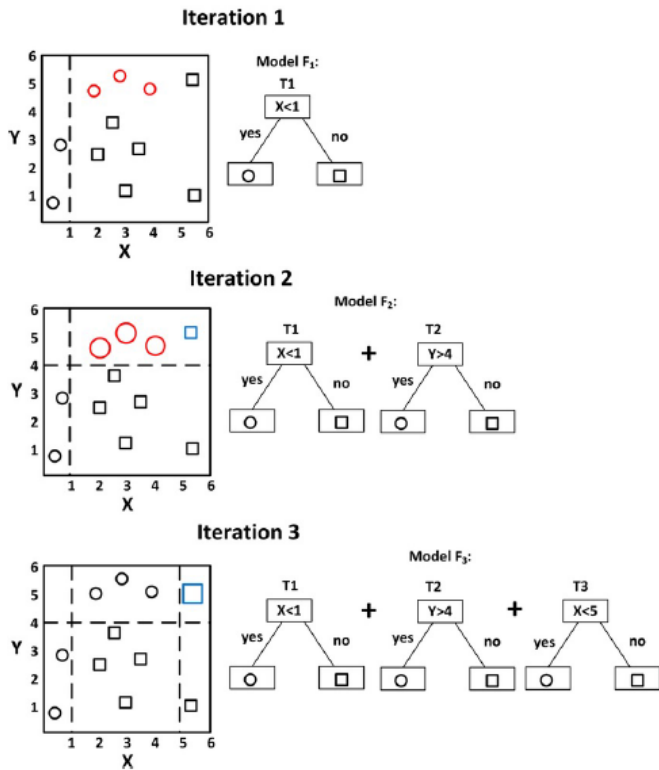- Fit a different linear model in each branch

# Ensemble approaches

# Bagging: Random forest

Dataset



Tree 1          Tree 2          Tree K

- Sample 80% of the dataset to train each decision tree
- Each tree may overfit to different part of the dataset
- But the consensus should be correct

# Boosting for classification



- The first model made some mistakes

- The mistakes were assigned higher weights for the subsequent models

- As more models are added, the ensemble should make less errors

- Ensemble = $w_1 f_1(x) + \cdots + w_n f_n(x)$
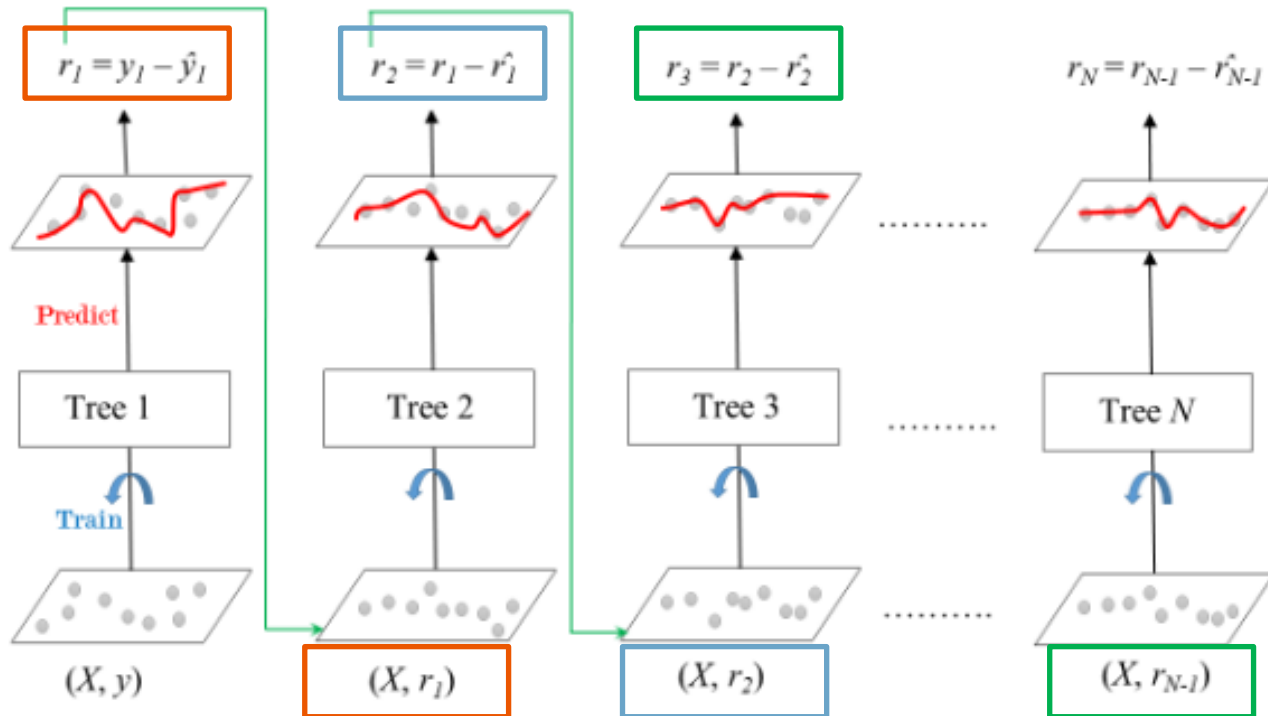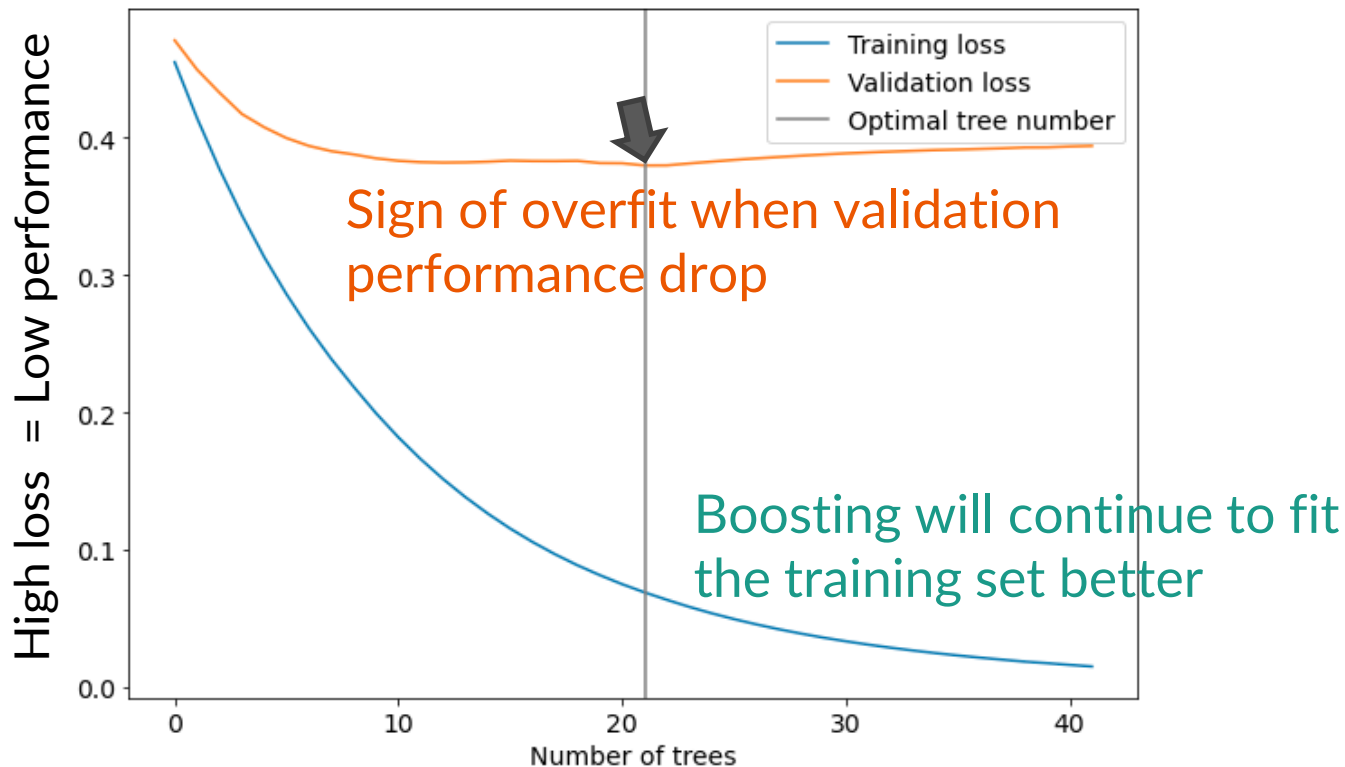
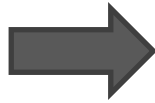# Boosting for regression



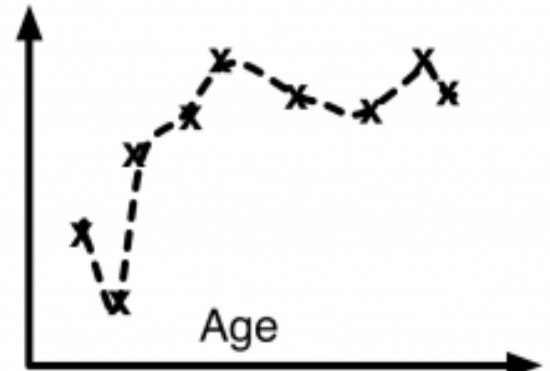Image from geeksforgeeks.org

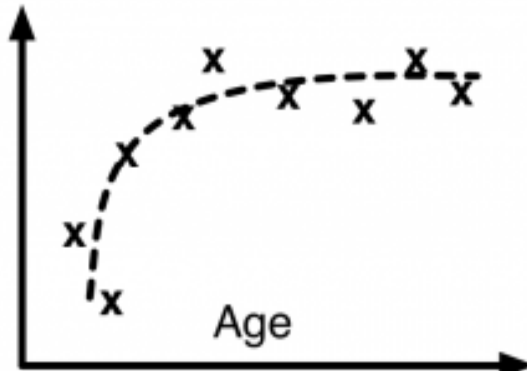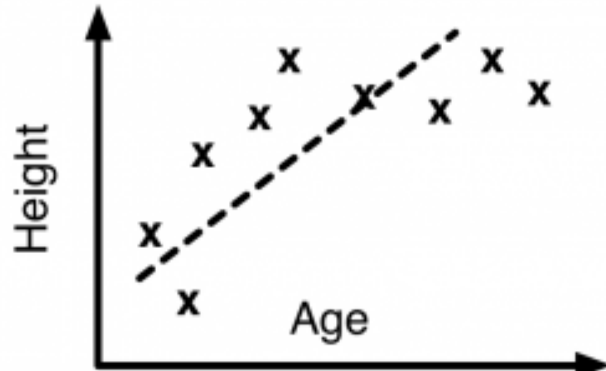# Controlling boosting with learning rate

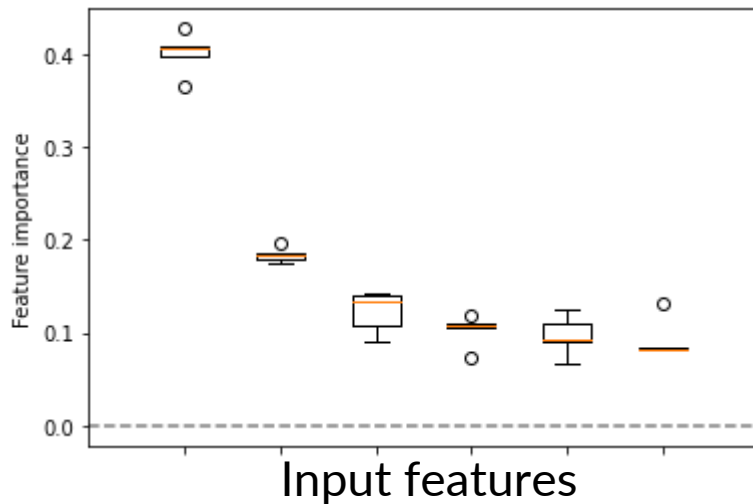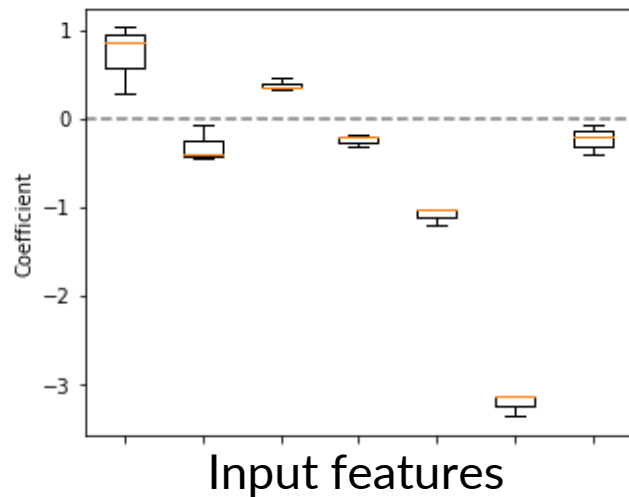# Impact of ensemble

Boosting solves underfitting
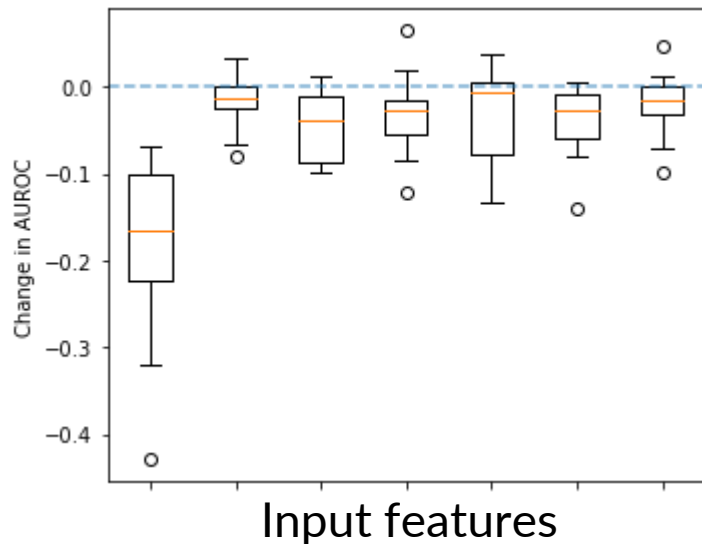
Bagging prevent overfitting

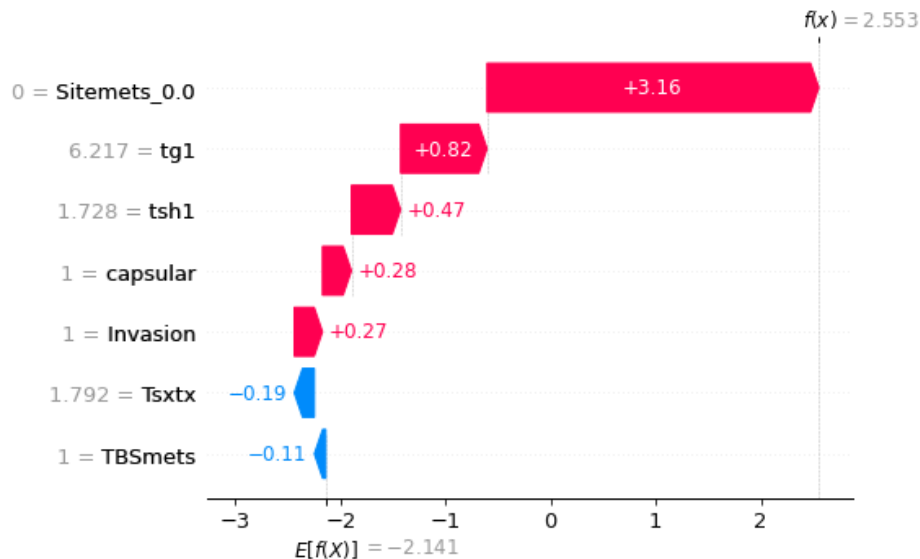# Explainability

# Feature importance



- Coefficients of linear, logistic, and SVM models
- Average improvement in impurity or entropy in tree models
- Model-level explanation

# Change in performance after dropping a feature



- Compare performance with and without each input feature
- Big drop = important

# Shapley value



$f(x) = 2.553$

0 = Sitemets_0.0    +3.16

6.217 = tg1    +0.82

1.728 = tsh1    +0.47

1 = capsular    +0.28

1 = Invasion    +0.27

1.792 = Tsxtx    −0.19

1 = TBSmets    −0.11

$E[f(X)] = -2.141$

- Change in predicted value due to the addition of a feature *i*

  - $\varphi_i(v) = \sum_{S \subseteq N \setminus \{i\}} \frac{|S|!(n-|S|-1)!}{n!} [v(S \cup \{i\}) - v(S)]$

- Sample-level explanation

# Any questions?

See you on February 23rd