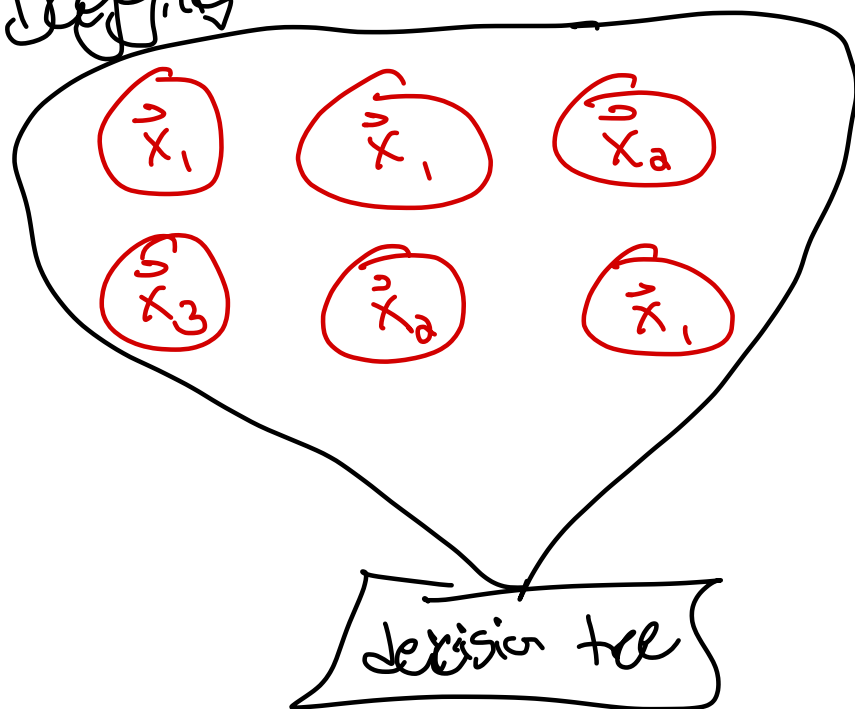
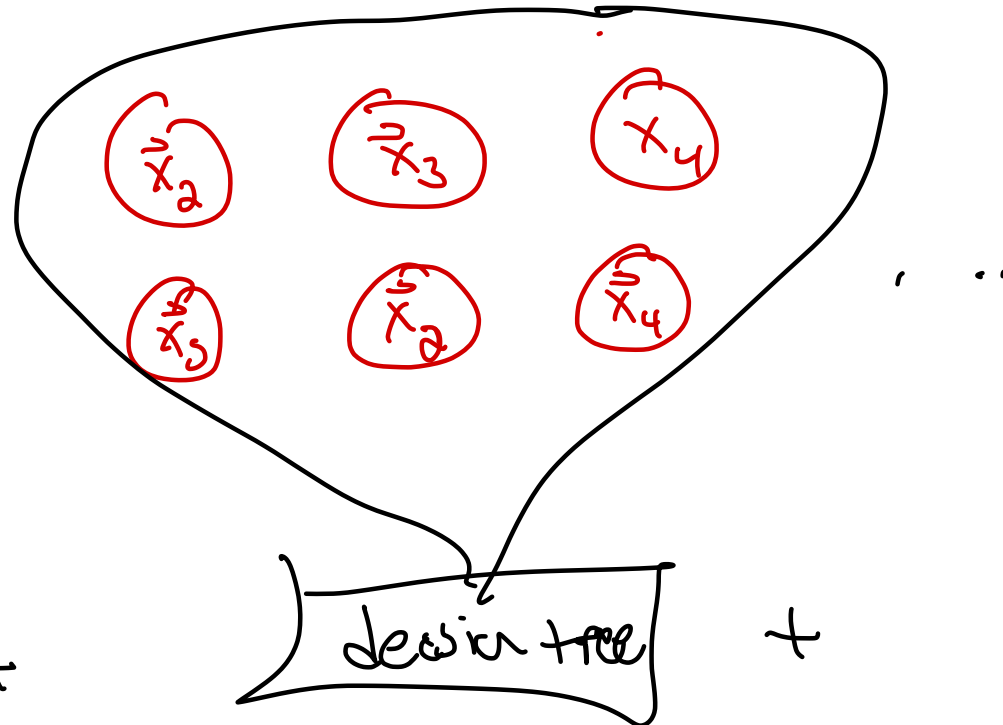


random forest:
a bunch of these
that all vote

Bagging



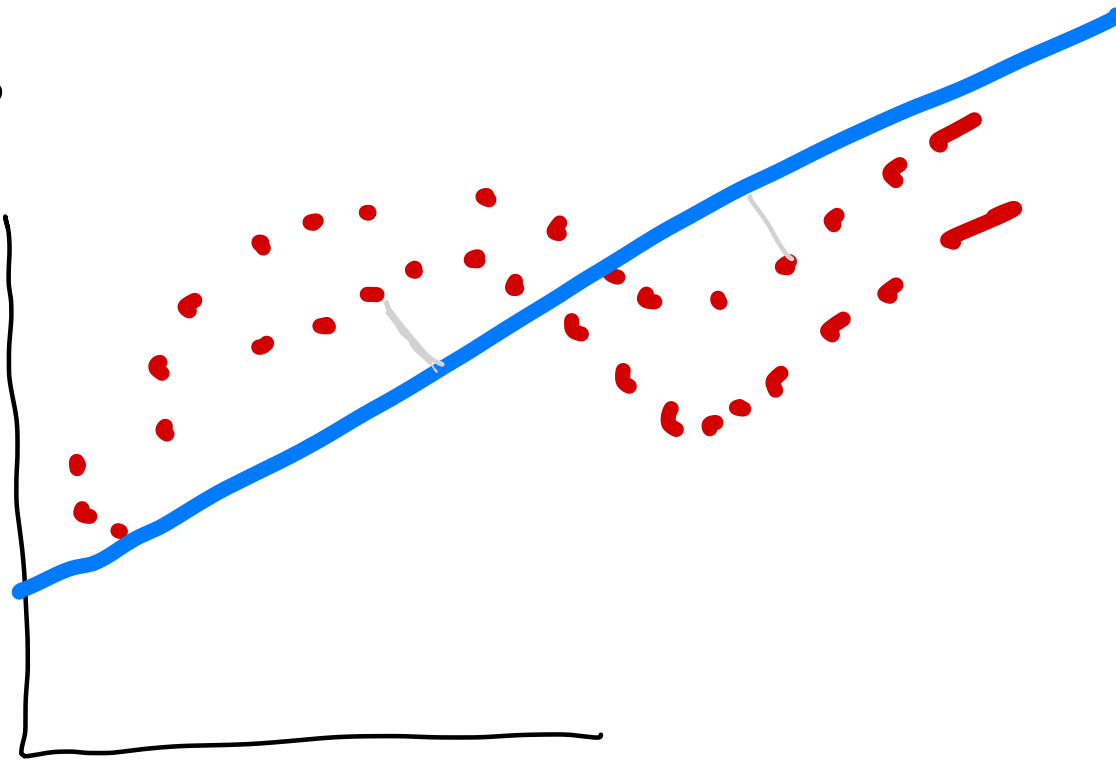
+



Variance

bagging \rightarrow designed to combat high
variance

bias



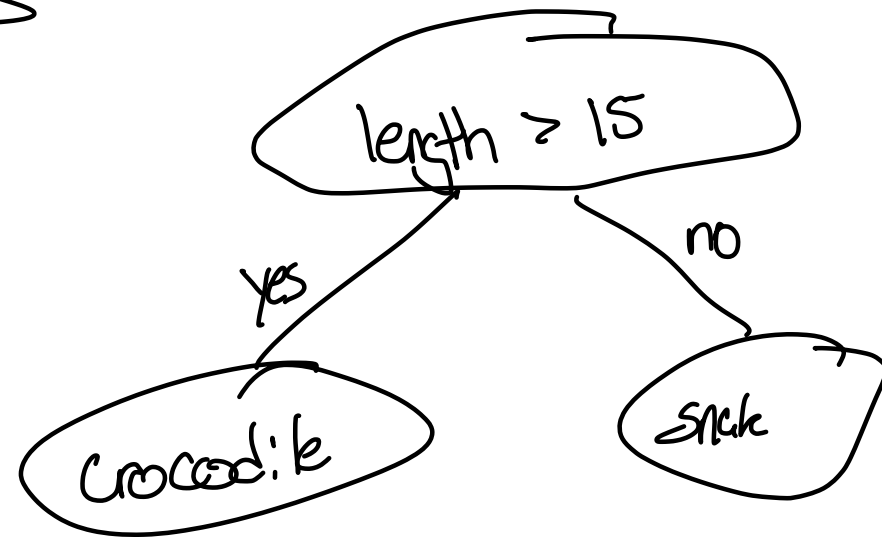
high bias

Boosting

↳ combat high bias

↳ "weak learners"

↳ trained sequentially



- AdaBoost ("adaptive boosting")

- adjust weights based on what it got

wrong

- Gradient boosting

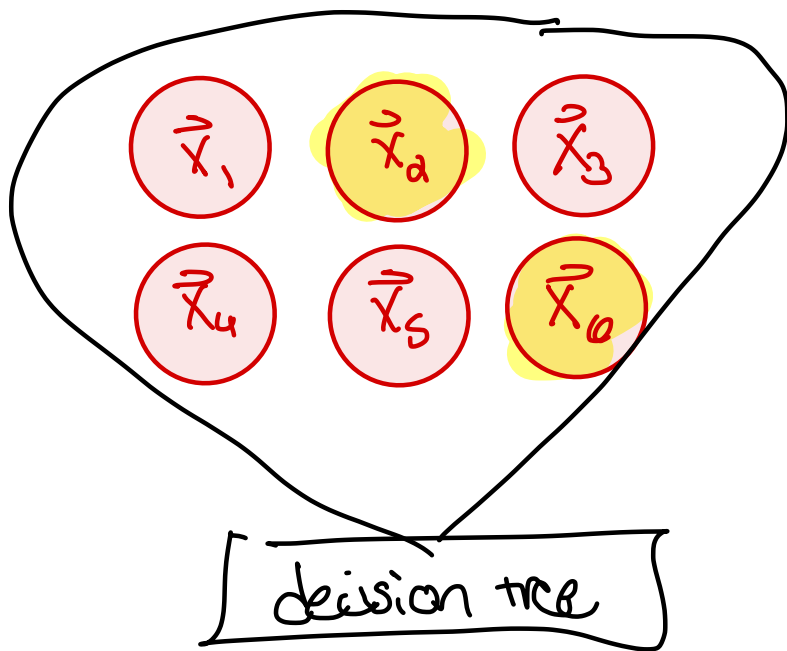
- ↳ next predictor is fit to errors made by the previous predictor

- XGBoost

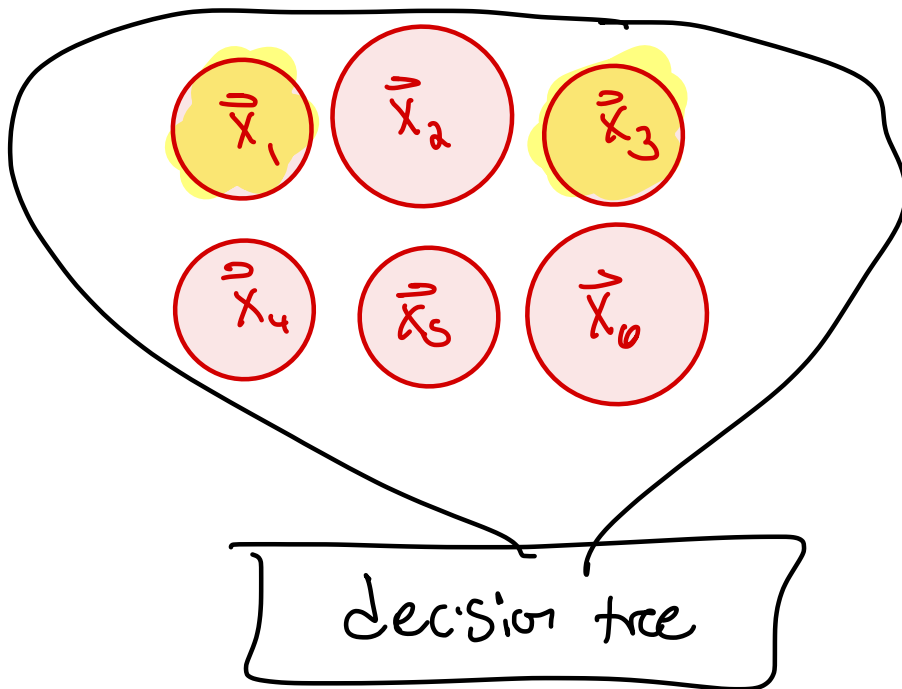
↳ uses more regularized model

- Light GBM

↳ different gradient boost



$h_1(x)$



$h_2(x)$

$$\alpha_1 h_1(x) + \alpha_2 h_2(x) + \dots + \alpha_m h_m(x)$$

weight	length	class	sample weight
		S	1/8
		C	1/8
		C	1/8
		C	1/8
		S	1/8
		S	1/8
		C	1/8
		S	1/8

+1 or -1

$$h_1(x)$$

$$\alpha = \frac{1}{2} \ln \left(\frac{1 - \text{error rate}}{\text{error rate}} \right)$$

$$\text{error rate} = \frac{\# \text{ wrong}}{\# \text{ total}}$$

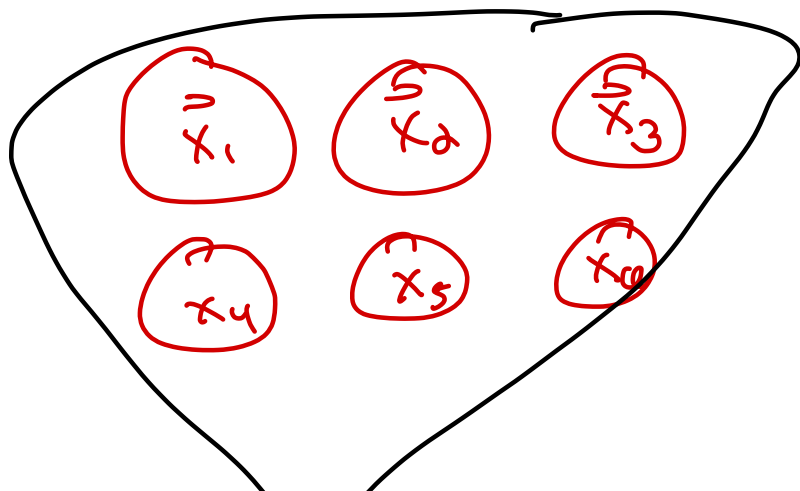
say
0

$$\sum_i w_i \delta(h(x_i) \neq y(x_i))$$

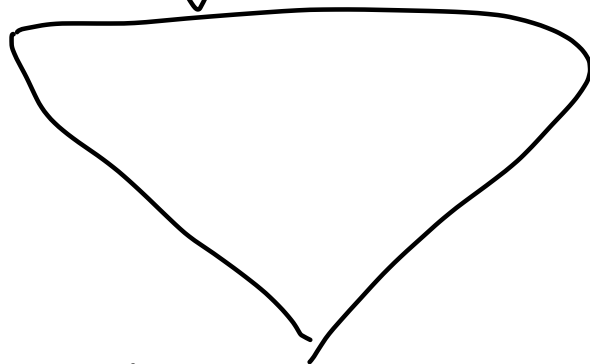
1 if not equal
0 otherwise

$$w_i e^{-\alpha y h(x)}$$

error rate



decision tree



decision tree

Other consideration with trees/forests

1) imbalanced data

a) be aware \rightarrow compare to baseline accuracy

b) oversample minority classes

\rightarrow down sample majority class

} more balanced dataset to train on

c) penalize minority mistakes more

\rightarrow Balanced Random Forest

\rightarrow draw bootstrap samples from minority class

\rightarrow draw same # of samples, with replacement, from majority class

\rightarrow repeat: train decision tree incorporating randomness

\rightarrow Weighted Random Forests +

a) assigns weight to each

b) weight GINI when choosing split

c) uses weighted majority vote when determining prediction

- different costs for features
↳ replace information gain with $\frac{\text{gain}}{\text{cost}}$