

# Decision trees

Fraida Fund

## Contents

In this lecture . . . . .	2
Recap . . . . .	3
Models for regression . . . . .	3
Models for classification . . . . .	3
Evaluating models (review) . . . . .	3
Demo: digit classification . . . . .	3
Flexible decisions with cheap prediction? . . . . .	3
Decision tree . . . . .	4
Tree terminology . . . . .	4
Note on notation . . . . .	4
Stratification of feature space (1) . . . . .	4
Stratification of feature space (2) . . . . .	4
Tree representation . . . . .	4
Stratification of feature space - illustration . . . . .	4
Training a decision tree . . . . .	5
Basic idea . . . . .	5
Recursive binary splitting . . . . .	5
Recursive binary splitting steps . . . . .	5
Loss function for regression tree . . . . .	6
Loss function for classification tree . . . . .	6
Classification error rate . . . . .	6
GINI index . . . . .	6
Entropy (1) . . . . .	6
Entropy (2) . . . . .	6
Information gain . . . . .	7
Example: should I play tennis? (1) . . . . .	7
Example: should I play tennis? (2) . . . . .	7
Example: should I play tennis? (3) . . . . .	7
Example: should I play tennis? (4) . . . . .	8
Example: should I play tennis? (5) . . . . .	8
Feature importance . . . . .	8
Illustration of measures of impurity . . . . .	8
Demo notebooks . . . . .	8
Bias and variance . . . . .	9
Managing tree depth . . . . .	9
Bias in decision tree . . . . .	10
Variance in decision tree . . . . .	10
Stopping criteria . . . . .	11
Pruning . . . . .	11
Pruning classification trees . . . . .	11
Weakest link pruning (1) . . . . .	11
Weakest link pruning (2) . . . . .	11

Cost complexity pruning . . . . .	11
Summary - so far . . . . .	12
The good and the bad (1) . . . . .	12
The good and the bad (2) . . . . .	12

## In this lecture

- Evaluating model cost, interpretability
- Decision trees
- Training decision trees
- Bias and variance of decision trees

## Recap

### Models for regression

Model	Function shape	Loss fn.	Training	Prediction	↓ complexity
Linear regression	Linear (or transformed)	$(\hat{y} - y)^2$	$\hat{\beta} = (A^T A)^{-1} A^T y$	$\hat{y} = [1, x^T] \hat{\beta}$	Regularization
KNN	Arbitrarily complicated	NA	Non-parametric, store training data	$\hat{y} = \frac{1}{K} \sum_{K_x} y_i$	Increase K

### Models for classification

Model	Function shape	Loss fn.	Training	$P(y = m x) =$	↓ complexity
Logistic regression	Linear (or transformed)	$-\ln P(y X)$	No closed form soln., use solver	$\frac{e^{z_m}}{\sum_{\ell=1}^M e^{z_\ell}}$	Regularization
KNN	Arbitrarily complicated	NA	Non-parametric, store training data	$\frac{1}{K} \sum_{K_x} I(y_i = m)$	Increase K

### Evaluating models (review)

- Performance (on different types of data)
- Cost/time for training and prediction
- Interpretability

### Demo: digit classification

- Shows how to use “magic commands” to time a line of code
- Interpretability of logistic regression, KNN

### Flexible decisions with cheap prediction?

KNN was very flexible, but prediction is **slow**.

Next up: flexible decisions, non-parametric approach, fast prediction

## Decision tree

### Tree terminology

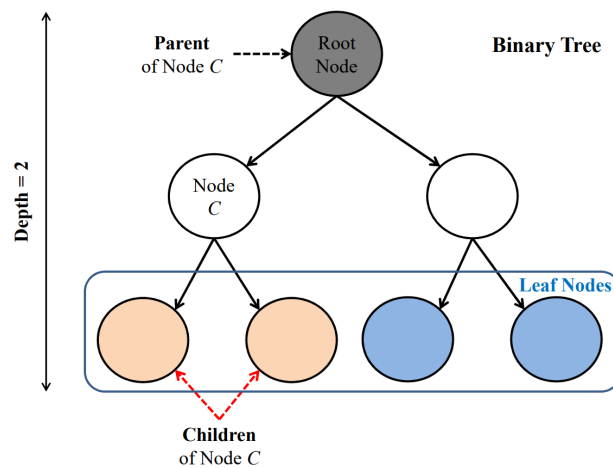


Figure 1: A binary tree.

### Note on notation

Following notation of ISLR, Chapter 8:

- $X_j$  is feature  $j$
- $x_i$  is sample  $i$

### Stratification of feature space (1)

- Given set of possible predictors,  $X_1, \dots, X_p$
- Training: Divide predictor space (set of possible values of  $X$ ) into  $J$  non-overlapping regions:  $R_1, \dots, R_J$ , by splitting sequentially on one feature at a time.

### Stratification of feature space (2)

- Prediction: For each observation that falls in region  $R_j$ , predict
  - mean of labels of training points in  $R_j$  (regression)
  - mode of labels of training points in  $R_j$  (classification)

### Tree representation

- Each node of tree that is not a leaf node: test one feature  $X_i$
- Each branch from a node: selects one value for  $X_i$
- Each leaf node: predict  $\hat{y}_{R_m}$

Tree is characterized by \* size of tree  $|T|$  (number of leaf nodes) \* depth (max length from root node to a leaf node)

### Stratification of feature space - illustration

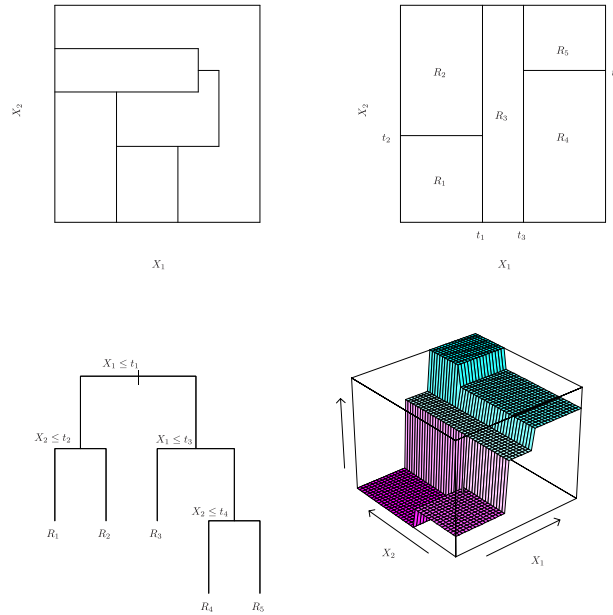


Figure 2: ISLR, Fig. 8.3. The stratification on the top left cannot be produced by a decision tree using recursive binary splitting. The other three subfigures represent a single stratification.

## Training a decision tree

### Basic idea

- Goal: find the high-dimensional rectangles that minimize error
- Computationally expensive to consider every possible partition
- Instead: top-down, greedy approach: recursive binary splitting
- Greedy: at each step, make the best decision at that step, without looking ahead and making a decision that might yield better results at future steps

### Recursive binary splitting

For any feature  $j$  and *cutpoint*  $s$ , define the regions

$$R_1(j, s) = \{X | X_j < s\}, \quad R_2(j, s) = \{X | X_j \geq s\}$$

where  $\{X | X_j < s\}$  is the region of predictor space in which  $X_j$  takes on a value less than  $s$ .

### Recursive binary splitting steps

Start at root of the tree, considering all training samples.

1. At the current node,
2. Find feature  $X_j$  and cutpoint  $s$  that minimizes some loss function (?)
3. Split training samples at that node into two leaf nodes
4. Stop when no training error (?)
5. Otherwise, repeat at leaf nodes

### Loss function for regression tree

For regression: look for feature  $j$  and cutpoint  $s$  that leads to the greatest possible reduction in RSS:

$$\sum_{i: x_i \in R_1(j,s)} (y_i - \hat{y}_{R_1})^2 + \sum_{i: x_i \in R_2(j,s)} (y_i - \hat{y}_{R_2})^2$$

(where  $\hat{y}_{R_j}$  is the prediction for the samples in  $R_j$ .)

### Loss function for classification tree

For classification, the basic idea is to find a split that leads to the greatest reduction in some measure of node *impurity*:

- A node whose samples all belong to the same class - most *pure*
- A node whose samples are evenly distributed among all classes - highly *impure*

### Classification error rate

For classification: one possible way is to split on 0-1 loss or *misclassification rate*:

$$\sum_{x_i \in R_m} 1(y_i \neq \hat{y}_{R_m})$$

Not a good metric for impurity!

### GINI index

GINI index:

$$\sum_{k=1}^K \hat{p}_{mk} (1 - \hat{p}_{mk})$$

### Entropy (1)

Entropy of a random variable  $X$  (from information theory):

$$H(X) = - \sum_{i=1}^N P(X = i) \log_2 P(X = i)$$

### Entropy (2)

Entropy as a measure of impurity on subset of samples:

$$- \sum_{k=1}^K \hat{p}_{mk} \log_2 \hat{p}_{mk}$$

where  $\hat{p}_{mk}$  is the proportion of training samples in  $R_m$  belonging to class  $k$ .

### Information gain

- Splitting on  $X$  creates subsets  $S_1$  and  $S_2$  with different entropies
- Conditional entropy:

$$\text{Entropy}(S|X) = \sum_v \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

- Choose feature to split so as to maximize information gain, the expected reduction in entropy due to splitting on  $X$ :

$$\text{Gain}(S, X) := \text{Entropy}(S) - \text{Entropy}(S|X)$$

### Example: should I play tennis? (1)

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Figure 3: Via Tom Mitchell.

### Example: should I play tennis? (2)

For top node:  $S = \{9+, 5-\}$ ,  $|S| = 14$

$$\text{Entropy}(S) = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} = 0.94$$

### Example: should I play tennis? (3)

If we split on Wind:

Considering the Weak branch:

- $S_{\text{weak}} = \{6+, 2-\}$ ,  $|S_{\text{weak}}| = 8$
- $\text{Entropy}(S_{\text{weak}}) = -\frac{6}{8} \log_2 \left(\frac{6}{8}\right) - \frac{2}{8} \log_2 \left(\frac{2}{8}\right) = 0.81$

Considering the Strong branch:

- $S_{\text{strong}} = \{3+, 3-\}, |S_{\text{strong}}| = 6$
- $\text{Entropy}(S_{\text{strong}}) = 1$

**Example: should I play tennis? (4)**

- $\text{Entropy}(S) = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} = 0.94$
- $\text{Entropy}(S|\text{Wind}) = \frac{8}{14} \text{Entropy}(S_{\text{weak}}) + \frac{6}{14} \text{Entropy}(S_{\text{strong}}) = 0.89$
- $\text{Gain}(S, \text{Wind}) = 0.94 - 0.89 = 0.05$

**Example: should I play tennis? (5)**

- $\text{Gain}(S, \text{Outlook}) = 0.246$
- $\text{Gain}(S, \text{Humidity}) = 0.151$
- $\text{Gain}(S, \text{Wind}) = 0.048$
- $\text{Gain}(S, \text{Temperature}) = 0.029$

→ Split on Outlook!

### Feature importance

- For each feature  $X_j$ , find all nodes where the feature was used as the split variable
- Add up information gain due to split (or for GINI index, different in loss weighted by number of samples.)
- This sum reflects feature importance

### Illustration of measures of impurity

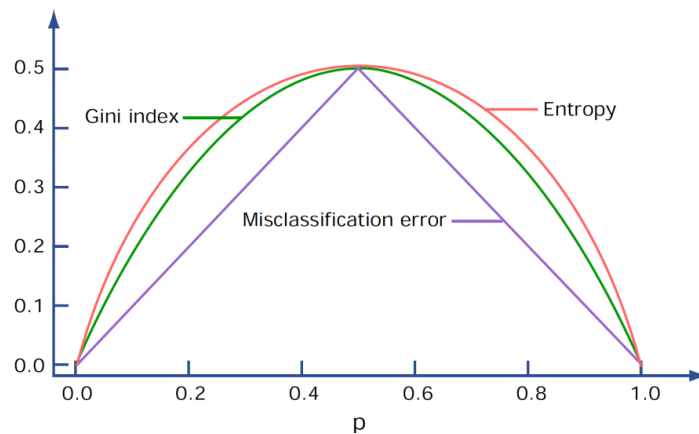


Figure 4: Misclassification error, GINI index, scaled entropy.

### Demo notebooks

- Demo: classification on mammal vs. non-mammal dataset
- Demo: digits dataset



## **Bias and variance**

### **Managing tree depth**

- If tree is too deep - likely to overfit (high variance)
- If tree is not deep enough - likely to have high bias

## Bias in decision tree

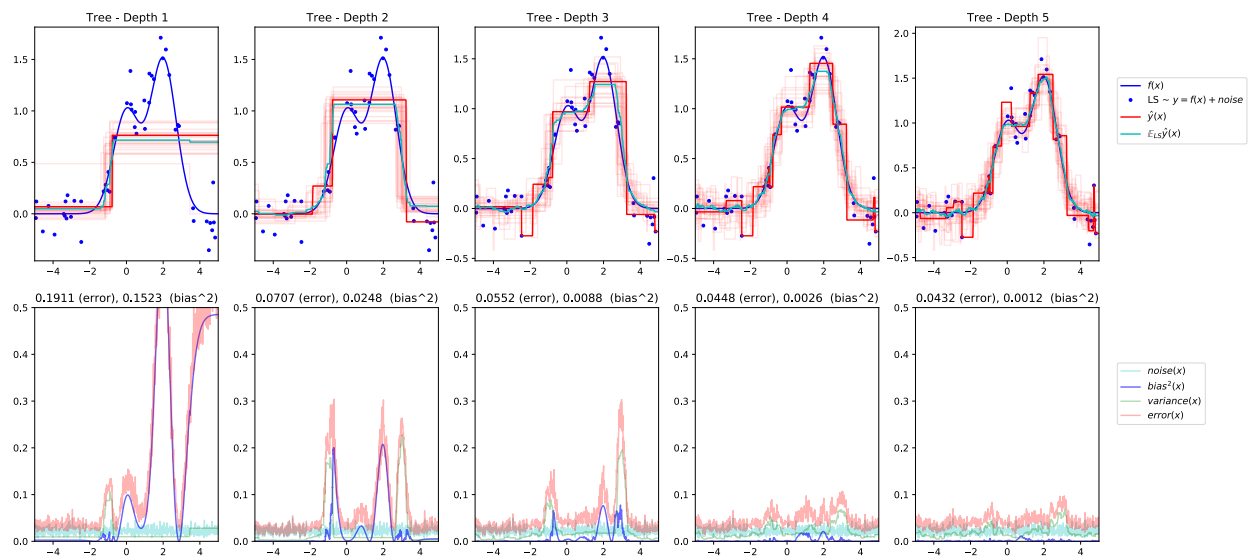


Figure 5: A tree that is too shallow has high bias.

## Variance in decision tree

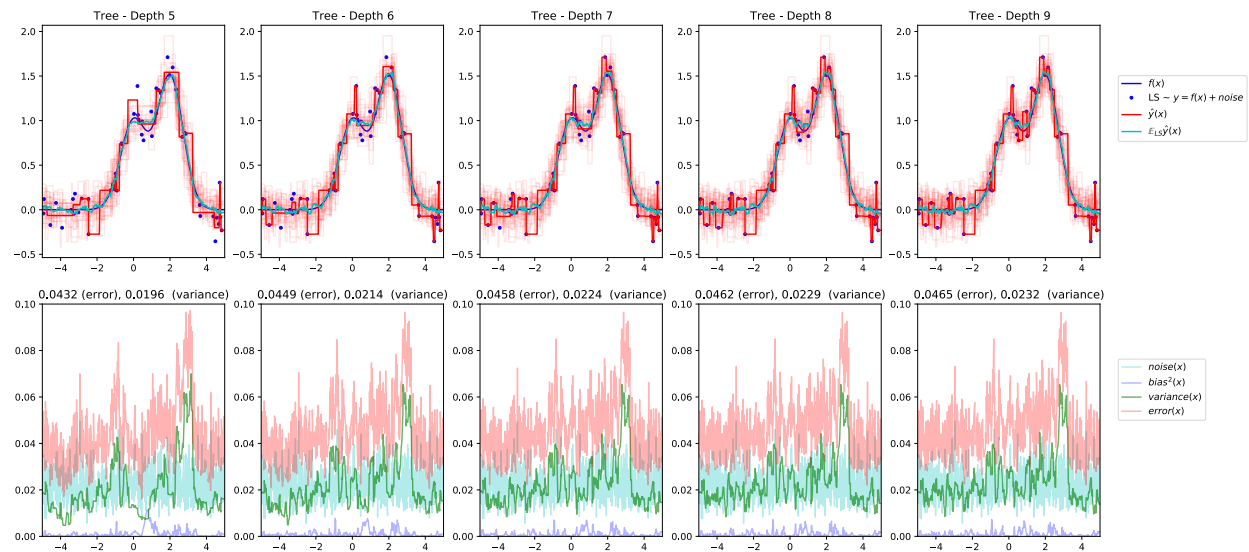


Figure 6: A tree that is too deep has high variance.

## Stopping criteria

If we build tree until there is zero error on training set, we have “memorized” training data.

Other stopping criteria:

- Max depth
- Max size (number of leaf nodes)
- Min number of samples to split
- Min number of samples in leaf node
- Min decrease in loss function due to split

(Can select depth, etc. by CV)

## Pruning

- Alternative to stopping criteria: build entire tree, then *prune*
- With greedy algorithm - a very good split may descend from a less-good split

## Pruning classification trees

We usually prune classification trees using classification error rate as loss function, even if tree was built using GINI or entropy.

In the next slides - loss function shown is for regression tree.

## Weakest link pruning (1)

Prune a large tree from leaves to root:

- Start with full tree  $T_0$
- Merge two adjacent leaf nodes into their parent to obtain  $T_1$  by minimizing:

$$\frac{RSS(T_1) - RSS(T_0)}{|T_0| - |T_1|}$$

## Weakest link pruning (2)

- Iterate to produce a sequence of trees  $T_0, T_1, \dots, T_m$  where  $T_m$  is a tree of minimum size.
- Select optimal tree by CV

## Cost complexity pruning

Minimize

$$\sum_{m=1}^{|T|} \sum_{x_i}^{R_m} (y_i - \hat{y}_{R_m})^2 + \alpha |T|$$

Choose  $\alpha$  by CV, 1-SE rule ( $\uparrow \alpha, \downarrow |T|$ ).

## **Summary - so far**

### **The good and the bad (1)**

Good:

- Easy to interpret, close to human decision-making
- Can derive feature importance
- Easily handles mixed types of features and different ranges
- Can find interactions that linear classifiers can't

### **The good and the bad (2)**

Bad:

- Need deep tree to overcome bias
- Deep trees have large variance
- Non-robust: Small change in data can cause large change in estimated tree