# DECISION TREES

David Kauchak
CS 158 — Spring 2022

1

## Admin

Assignment 1 due tomorrow (Friday)

Assignment 2 out soon: start ASAP! (due next Sunday)

- Can (and are encouraged to) work in pairs

Slack

Office hours M-Th, 2:30-3:30pm, starting today (zoom link in sakai)

2

## Admin

Lecture notes posted (webpage)

Lecture recordings uploaded (box — see sakai for link)

Keep up with the reading

Videos before class

Class ends at 2:30 ☺

3

## Representing examples

examples



What is an example?
How is it represented?

4

## Features

examples



features

$f_1, f_2, f_3, \ldots, f_n$

$f_1, f_2, f_3, \ldots, f_n$

$f_1, f_2, f_3, \ldots, f_n$

$f_1, f_2, f_3, \ldots, f_n$

How our algorithms actually "view" the data

Features are the questions we can ask about the examples

5

## Features

examples



features

red, round, leaf, 3oz, …

green, round, no leaf, 4oz, …

yellow, curved, no leaf, 8oz, …

green, curved, no leaf, 7oz, …

How our algorithms actually "view" the data

Features are the questions we can ask about the examples

6

## Classification revisited

examples

red, round, leaf, 3oz, …

green, round, no leaf, 4oz, …

yellow, curved, no leaf, 8oz, …

green, curved, no leaf, 7oz, …

label

apple

apple

banana

banana

learn

model/ classifier

During learning/training/induction, learn a model of what distinguishes apples and bananas *based on the features*

7

## Classification revisited

red, round, no leaf, 4oz, …

model/ classifier

predict

Apple or banana?

The model can then classify a new example *based on the features*

8

2

## Classification revisited

red, round, no leaf, 4oz, … → model/ classifier → *predict* → Apple

Why?

The model can then classify a new example *based on the features*

9

## Classification revisited

| Training data | | Test set |
|---|---|---|
| examples | label | |
| red, round, leaf, 3oz, … | apple | |
| green, round, no leaf, 4oz, … | apple | red, round, no leaf, 4oz, … ? |
| yellow, curved, no leaf, 4oz, … | banana | |
| green, curved, no leaf, 5oz, … | banana | |

10

## Classification revisited

| Training data | | Test set |
|---|---|---|
| examples | label | |
| red, round, leaf, 3oz, … | apple | |
| green, round, no leaf, 4oz, … | apple | red, round, no leaf, 4oz, … ? |
| yellow, curved, no leaf, 4oz, … | banana | |
| green, curved, no leaf, 5oz, … | banana | |

Learning is about *generalizing* from the training data

What does this assume about the training and test set?

11

## A sample data set

| Features | | | | Label |
|---|---|---|---|---|
| Hour | Weather | Accident | Stall | Commute |
| 8 AM | Sunny | No | No | Long |
| 8 AM | Cloudy | No | Yes | Long |
| 10 AM | Sunny | No | No | Short |
| 9 AM | Rainy | Yes | No | Long |
| 9 AM | Sunny | Yes | Yes | Long |
| 10 AM | Sunny | No | No | Short |
| 10 AM | Cloudy | No | No | Short |
| 9 AM | Sunny | Yes | No | Long |
| 10 AM | Cloudy | Yes | Yes | Long |
| 10 AM | Rainy | No | No | Short |
| 8 AM | Cloudy | Yes | No | Long |
| 9 AM | Rainy | No | No | Short |

8 AM, Rainy, Yes, No?
10 AM, Rainy, No, No?

Can you describe a "model" that could be used to make decisions in general?

12

3

## Decision trees

Leave At

10 AM     8 AM     9 AM

Stall?     Long     Accident?

No   Yes     No   Yes

Short   Long     Short   Long

Tree with internal nodes labeled by features

Branches are labeled by tests on that feature

Leaves labeled with classes

13

## Decision trees

Leave At

10 AM     8 AM     9 AM

Stall?     Long     Accident?

No   Yes     No   Yes

Short   Long     Short   Long

Leave = 8 AM     Accident = Yes
Weather = Rainy     Stall = No

Tree with internal nodes labeled by features

Branches are labeled by tests on that feature

Leaves labeled with classes

14

## Decision trees

Leave At

10 AM     8 AM     9 AM

Stall?     Long     Accident?

No   Yes     No   Yes

Short   Long     Short   Long

Leave = 8 AM     Accident = Yes
Weather = Rainy     Stall = No

Tree with internal nodes labeled by features

Branches are labeled by tests on that feature

Leaves labeled with classes

15

## Decision trees

Leave At

10 AM     8 AM     9 AM

Stall?     Long     Accident?

No   Yes     No   Yes

Short   Long     Short   Long

Leave = 10 AM     Accident = No
Weather = Rainy     Stall = No

Tree with internal nodes labeled by features

Branches are labeled by tests on that feature

Leaves labeled with classes

16

---

Done thinking; writing final.

---

# Decision trees

Tree with internal nodes labeled by features

Branches are labeled by tests on that feature
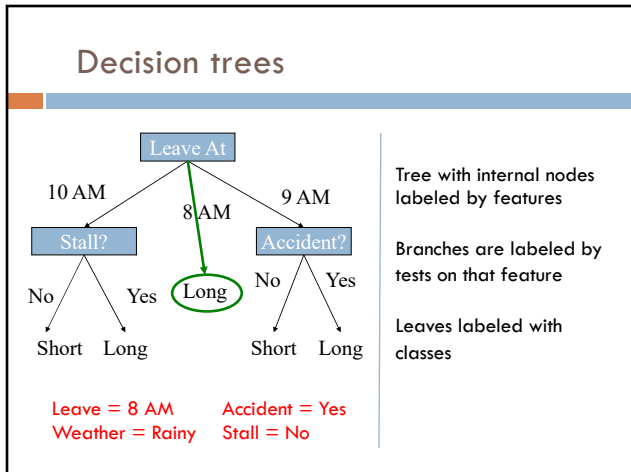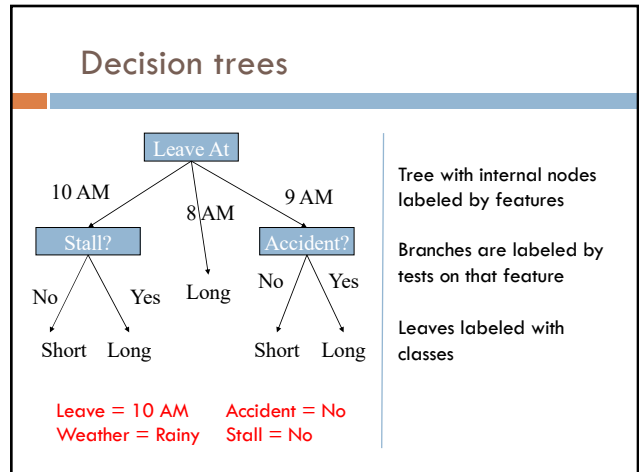
Leaves labeled with classes

Leave = 10 AM
Weather = Rainy
Accident = No
Stall = No

(Tree: Leave At → 10 AM / 8 AM / 9 AM; Stall? → No: Short, Yes: Long; 8 AM: Long; Accident? → No: Short, Yes: Long)

# To ride or not to ride, that is the question…

| Terrain | Unicycle-type | Weather | Go-For-Ride? |
|---|---|---|---|
| Trail | Normal | Rainy | NO |
| Road | Normal | Sunny | YES |
| Trail | Mountain | Sunny | YES |
| Road | Mountain | Rainy | YES |
| Trail | Normal | Snowy | NO |
| Road | Normal | Rainy | YES |
| Road | Mountain | Snowy | YES |
| Trail | Normal | Sunny | NO |
| Road | Normal | Snowy | NO |
| Trail | Mountain | Snowy | YES |

**Build a decision tree**

# Recursive approach

Base case: If all data belong to the same class, create a leaf node with that label

Otherwise:
- calculate the "score" for each feature if we used it to split the data
- pick the feature with the highest score, partition the data based on that data value and call recursively

# Partitioning the data

| Terrain | Unicycle-type | Weather | Go-For-Ride? |
|---|---|---|---|
| Trail | Normal | Rainy | NO |
| Road | Normal | Sunny | YES |
| Trail | Mountain | Sunny | YES |
| Road | Mountain | Rainy | YES |
| Trail | Normal | Snowy | NO |
| Road | Normal | Rainy | YES |
| Road | Mountain | Snowy | YES |
| Trail | Normal | Sunny | NO |
| Road | Normal | Snowy | NO |
| Trail | Mountain | Snowy | YES |

(Tree: Terrain → Road / Trail, ?)

17 18 19 20

## Partitioning the data (21)

| Terrain | Unicycle-type | Weather | Go-For-Ride? |
|---------|---------------|---------|--------------|
| Trail | Normal | Rainy | NO |
| Road | Normal | Sunny | YES |
| Trail | Mountain | Sunny | YES |
| Road | Mountain | Rainy | YES |
| Trail | Normal | Snowy | NO |
| Road | Normal | Rainy | YES |
| Road | Mountain | Snowy | YES |
| Trail | Normal | Sunny | NO |
| Road | Normal | Snowy | NO |
| Trail | Mountain | Snowy | YES |

Terrain
Road    Trail
?

21

## Partitioning the data (22)

| Terrain | Unicycle-type | Weather | Go-For-Ride? |
|---------|---------------|---------|--------------|
| Trail | Normal | Rainy | NO |
| Road | Normal | Sunny | YES |
| Trail | Mountain | Sunny | YES |
| Road | Mountain | Rainy | YES |
| Trail | Normal | Snowy | NO |
| Road | Normal | Rainy | YES |
| Road | Mountain | Snowy | YES |
| Trail | Normal | Sunny | NO |
| Road | Normal | Snowy | NO |
| Trail | Mountain | Snowy | YES |

Terrain
Road    Trail
YES: 4
NO: 1

22

## Partitioning the data (23)

| Terrain | Unicycle-type | Weather | Go-For-Ride? |
|---------|---------------|---------|--------------|
| Trail | Normal | Rainy | NO |
| Road | Normal | Sunny | YES |
| Trail | Mountain | Sunny | YES |
| Road | Mountain | Rainy | YES |
| Trail | Normal | Snowy | NO |
| Road | Normal | Rainy | YES |
| Road | Mountain | Snowy | YES |
| Trail | Normal | Sunny | NO |
| Road | Normal | Snowy | NO |
| Trail | Mountain | Snowy | YES |

Terrain
Road    Trail
YES: 4    ?
NO: 1

23

## Partitioning the data (24)

| Terrain | Unicycle-type | Weather | Go-For-Ride? |
|---------|---------------|---------|--------------|
| Trail | Normal | Rainy | NO |
| Road | Normal | Sunny | YES |
| Trail | Mountain | Sunny | YES |
| Road | Mountain | Rainy | YES |
| Trail | Normal | Snowy | NO |
| Road | Normal | Rainy | YES |
| Road | Mountain | Snowy | YES |
| Trail | Normal | Sunny | NO |
| Road | Normal | Snowy | NO |
| Trail | Mountain | Snowy | YES |

Terrain
Road    Trail
YES: 4    YES: 2
NO: 1    NO: 3

24

## Slide 25

### Partitioning the data

| Terrain | Unicycle-type | Weather | Go-For-Ride? |
|---------|---------------|---------|--------------|
| Trail | Normal | Rainy | NO |
| Road | Normal | Sunny | YES |
| Trail | Mountain | Sunny | YES |
| Road | Mountain | Rainy | YES |
| Trail | Normal | Snowy | NO |
| Road | Normal | Rainy | YES |
| Road | Mountain | Snowy | YES |
| Trail | Normal | Sunny | NO |
| Road | Normal | Snowy | NO |
| Trail | Mountain | Snowy | YES |

Terrain
Road → YES: 4 NO: 1
Trail → YES: 2 NO: 3

Unicycle
Mountain → ?
Normal → ?

## Slide 26

### Partitioning the data

| Terrain | Unicycle-type | Weather | Go-For-Ride? |
|---------|---------------|---------|--------------|
| Trail | Normal | Rainy | NO |
| Road | Normal | Sunny | YES |
| Trail | Mountain | Sunny | YES |
| Road | Mountain | Rainy | YES |
| Trail | Normal | Snowy | NO |
| Road | Normal | Rainy | YES |
| Road | Mountain | Snowy | YES |
| Trail | Normal | Sunny | NO |
| Road | Normal | Snowy | NO |
| Trail | Mountain | Snowy | YES |

Terrain
Road → YES: 4 NO: 1
Trail → YES: 2 NO: 3

Unicycle
Mountain → YES: 4 NO: 0
Normal → YES: 2 NO: 4

## Slide 27

### Partitioning the data

| Terrain | Unicycle-type | Weather | Go-For-Ride? |
|---------|---------------|---------|--------------|
| Trail | Normal | Rainy | NO |
| Road | Normal | Sunny | YES |
| Trail | Mountain | Sunny | YES |
| Road | Mountain | Rainy | YES |
| Trail | Normal | Snowy | NO |
| Road | Normal | Rainy | YES |
| Road | Mountain | Snowy | YES |
| Trail | Normal | Sunny | NO |
| Road | Normal | Snowy | NO |
| Trail | Mountain | Snowy | YES |

Terrain
Road → YES: 4 NO: 1
Trail → YES: 2 NO: 3

Unicycle
Mountain → YES: 4 NO: 0
Normal → YES: 2 NO: 4

Weather
Rainy → YES: 2 NO: 1
Snowy → YES: 2 NO: 2
Sunny → YES: 2 NO: 1

## Slide 28

### Partitioning the data

Terrain
Road → YES: 4 NO: 1
Trail → YES: 2 NO: 3

Unicycle
Mountain → YES: 4 NO: 0
Normal → YES: 2 NO: 4

Weather
Rainy → YES: 2 NO: 1
Snowy → YES: 2 NO: 2
Sunny → YES: 2 NO: 1

calculate the "score" for each feature
if we used it to split the data

What score should we use?
If we just stopped here, which tree would be best?
How could we make these into decision trees?

## Decision trees

| Terrain | | Unicycle | | Weather | | |
|---|---|---|---|---|---|---|
| Road | Trail | Mountain | Normal | Rainy | Snowy | Sunny |
| YES: 4 | YES: 2 | YES: 4 | YES: 2 | YES: 2 | YES: 2 | YES: 2 |
| NO: 1 | NO: 3 | NO: 0 | NO: 4 | NO: 1 | NO: 2 | NO: 1 |

How could we make these into decision trees?

29

## Decision trees

| Terrain | | Unicycle | | Weather | | |
|---|---|---|---|---|---|---|
| Road | Trail | Mountain | Normal | Rainy | Snowy | Sunny |
| **YES**: 4 | YES: 2 | **YES**: 4 | YES: 2 | **YES**: 2 | YES: 2 | **YES**: 2 |
| NO: 1 | **NO**: 3 | NO: 0 | **NO**: 4 | NO: 1 | **NO**: 2 | NO: 1 |

30

## Decision trees

| Terrain | | Unicycle | | Weather | | |
|---|---|---|---|---|---|---|
| Road | Trail | Mountain | Normal | Rainy | Snowy | Sunny |
| **YES**: 4 | YES: 2 | **YES**: 4 | YES: 2 | **YES**: 2 | YES: 2 | **YES**: 2 |
| NO: 1 | **NO**: 3 | NO: 0 | **NO**: 4 | NO: 1 | **NO**: 2 | NO: 1 |

Training error: the average error over the training set

For classification, the most common "error" is the number of mistakes

Training error for each of these?

31

## Decision trees

| Terrain | | Unicycle | | Weather | | |
|---|---|---|---|---|---|---|
| Road | Trail | Mountain | Normal | Rainy | Snowy | Sunny |
| **YES**: 4 | YES: 2 | **YES**: 4 | YES: 2 | **YES**: 2 | YES: 2 | **YES**: 2 |
| NO: 1 | **NO**: 3 | NO: 0 | **NO**: 4 | NO: 1 | **NO**: 2 | NO: 1 |

3/10          2/10          4/10

Training error: the average error over the training set

32

## Slide 33 — Training error vs. accuracy

Terrain
Road — Trail
YES: 4 / NO: 1    YES: 2 / NO: 3

Unicycle
Mountain — Normal
YES: 4 / NO: 0    YES: 2 / NO: 4

Weather
Rainy — Snowy — Sunny
YES: 2 / NO: 1    YES: 2 / NO: 2    YES: 2 / NO: 1

Training error: 3/10    2/10    4/10

Training accuracy: 7/10    8/10    6/10

training error = 1-accuracy    (and vice versa)

Training error: the average error over the training set

Training accuracy: the average proportion correct over the training set

33

## Slide 34 — Recurse

| Terrain | Unicycle-type | Weather | Go-For-Ride? |
|---|---|---|---|
| Trail | Normal | Rainy | NO |
| Road | Normal | Sunny | YES |
| Trail | Mountain | Sunny | YES |
| Road | Mountain | Rainy | YES |
| Trail | Normal | Snowy | NO |
| Road | Normal | Rainy | YES |
| Road | Mountain | Snowy | YES |
| Trail | Normal | Sunny | NO |
| Road | Normal | Snowy | NO |
| Trail | Mountain | Snowy | YES |

Unicycle
Mountain — Normal
YES: 4 / NO: 0    YES: 2 / NO: 4

34

## Slide 35 — Recurse

Unicycle
Mountain — Normal
YES: 4 / NO: 0    YES: 2 / NO: 4

| Terrain | Unicycle-type | Weather | Go-For-Ride? |
|---|---|---|---|
| Trail | Mountain | Sunny | YES |
| Road | Mountain | Rainy | YES |
| Road | Mountain | Snowy | YES |
| Trail | Mountain | Snowy | YES |

| Terrain | Unicycle-type | Weather | Go-For-Ride? |
|---|---|---|---|
| Trail | Normal | Rainy | NO |
| Road | Normal | Sunny | YES |
| Trail | Normal | Snowy | NO |
| Road | Normal | Rainy | YES |
| Trail | Normal | Sunny | NO |
| Road | Normal | Snowy | NO |

35

## Slide 36 — Recurse

Unicycle
Mountain — Normal
YES: 4 / NO: 0

| Terrain | Unicycle-type | Weather | Go-For-Ride? |
|---|---|---|---|
| Trail | Mountain | Sunny | YES |
| Road | Mountain | Rainy | YES |
| Road | Mountain | Snowy | YES |
| Trail | Mountain | Snowy | YES |

What should we do?

36

## Recurse

Unicycle — Mountain / Normal

YES: 4
NO: 0

| Terrain | Unicycle-type | Weather | Go-For-Ride? |
|---------|---------------|---------|--------------|
| Trail | Mountain | Sunny | YES |
| Road | Mountain | Rainy | YES |
| Road | Mountain | Snowy | YES |
| Trail | Mountain | Snowy | YES |

No need to examine other features since all examples have the same label.

37

## Recurse

Unicycle — Mountain / Normal

YES: 4          YES: 2
NO: 0           NO: 4

| Terrain | Unicycle-type | Weather | Go-For-Ride? |
|---------|---------------|---------|--------------|
| Trail | Normal | Rainy | NO |
| Road | Normal | Sunny | YES |
| Trail | Normal | Snowy | NO |
| Road | Normal | Rainy | YES |
| Trail | Normal | Sunny | NO |
| Road | Normal | Snowy | NO |

38

## Recurse

Unicycle — Mountain / Normal

YES: 4          YES: 2
NO: 0           NO: 4

Still two features left we can split on

| Terrain | Unicycle-type | Weather | Go-For-Ride? |
|---------|---------------|---------|--------------|
| Trail | Normal | Rainy | NO |
| Road | Normal | Sunny | YES |
| Trail | Normal | Snowy | NO |
| Road | Normal | Rainy | YES |
| Trail | Normal | Sunny | NO |
| Road | Normal | Snowy | NO |

39

## Recurse

Unicycle — Mountain / Normal

YES: 4          YES: 2
NO: 0           NO: 4

Terrain — Road / Trail

| Terrain | Unicycle-type | Weather | Go-For-Ride? |
|---------|---------------|---------|--------------|
| Trail | Normal | Rainy | NO |
| Road | Normal | Sunny | YES |
| Trail | Normal | Snowy | NO |
| Road | Normal | Rainy | YES |
| Trail | Normal | Sunny | NO |
| Road | Normal | Snowy | NO |

40

## Slide 41

# Recurse

Unicycle
Mountain — Normal
YES: 4    YES: 2
NO: 0     NO: 4

Terrain
Road — Trail
YES: 2    YES: 0
NO: 1     NO: 3

| Terrain | Unicycle-type | Weather | Go-For-Ride? |
|---------|---------------|---------|--------------|
| Trail | Normal | Rainy | NO |
| Road | Normal | Sunny | YES |
| Trail | Normal | Snowy | NO |
| Road | Normal | Rainy | YES |
| Trail | Normal | Sunny | NO |
| Road | Normal | Snowy | NO |

41

## Slide 42

# Recurse

Unicycle
Mountain — Normal
YES: 4    YES: 2
NO: 0     NO: 4

Terrain
Road — Trail
YES: 2    YES: 0
NO: 1     NO: 3

| Terrain | Unicycle-type | Weather | Go-For-Ride? |
|---------|---------------|---------|--------------|
| Trail | Normal | Rainy | NO |
| Road | Normal | Sunny | YES |
| Trail | Normal | Snowy | NO |
| Road | Normal | Rainy | YES |
| Trail | Normal | Sunny | NO |
| Road | Normal | Snowy | NO |

Weather
Rainy — Snowy — Sunny
YES: 1    YES: 0    YES: 1
NO: 1     NO: 2     NO: 1

42

## Slide 43

# Recurse

Unicycle
Mountain — Normal
YES: 4    YES: 2
NO: 0     NO: 4

Terrain
Road — Trail
YES: 2    YES: 0
NO: 1     NO: 3
1/6

| Terrain | Unicycle-type | Weather | Go-For-Ride? |
|---------|---------------|---------|--------------|
| Trail | Normal | Rainy | NO |
| Road | Normal | Sunny | YES |
| Trail | Normal | Snowy | NO |
| Road | Normal | Rainy | YES |
| Trail | Normal | Sunny | NO |
| Road | Normal | Snowy | NO |

Weather
Rainy — Snowy — Sunny
YES: 1    YES: 0    YES: 1
NO: 1     NO: 2     NO: 1
2/6

Which should we pick?

43

## Slide 44

# Recurse

Unicycle
Mountain — Normal
YES: 4    YES: 2
NO: 0     NO: 4

Terrain
Road — Trail
YES: 2    YES: 0
NO: 1     NO: 3

| Terrain | Unicycle-type | Weather | Go-For-Ride? |
|---------|---------------|---------|--------------|
| Road | Normal | Sunny | YES |
| Road | Normal | Rainy | YES |
| Road | Normal | Snowy | NO |

44

## Recurse



## Recurse



| Terrain | Unicycle-type | Weather | Go-For-Ride? |
|---|---|---|---|
| Trail | Normal | Rainy | NO |
| Road | Normal | Sunny | YES |
| Trail | Mountain | Sunny | YES |
| Road | Mountain | Rainy | YES |
| Trail | Normal | Snowy | NO |
| Road | Normal | Rainy | YES |
| Road | Mountain | Snowy | YES |
| Trail | Normal | Sunny | NO |
| Road | Normal | Snowy | NO |
| Trail | Mountain | Snowy | YES |

Training error?

Are we always guaranteed to get a training error of 0?

45

46

## Problematic data

| Terrain | Unicycle-type | Weather | Go-For-Ride? |
|---|---|---|---|
| Trail | Normal | Rainy | NO |
| Road | Normal | Sunny | YES |
| Trail | Mountain | Sunny | YES |
| Road | Mountain | Snowy | NO |
| Trail | Normal | Snowy | NO |
| Road | Normal | Rainy | YES |
| Road | Mountain | Snowy | YES |
| Trail | Normal | Sunny | NO |
| Road | Normal | Snowy | NO |
| Trail | Mountain | Snowy | YES |

When can this happen?

## Recursive approach

Base case: If all data belong to the same class, create a leaf node with that label **OR** all the data has the same feature values

Do we always want to go all the way to the bottom?

47

48

## What would the tree look like for…

| Terrain | Unicycle-type | Weather | Go-For-Ride? |
|---|---|---|---|
| Trail | Mountain | Rainy | YES |
| Trail | Mountain | Sunny | YES |
| Road | Mountain | Snowy | YES |
| Road | Mountain | Sunny | YES |
| Trail | Normal | Snowy | NO |
| Trail | Normal | Rainy | NO |
| Road | Normal | Snowy | YES |
| Road | Normal | Sunny | NO |
| Trail | Normal | Sunny | NO |

49

## What would the tree look like for…

| Terrain | Unicycle-type | Weather | Go-For-Ride? |
|---|---|---|---|
| Trail | Mountain | Rainy | YES |
| Trail | Mountain | Sunny | YES |
| Road | Mountain | Snowy | YES |
| Road | Mountain | Sunny | YES |
| Trail | Normal | Snowy | NO |
| Trail | Normal | Rainy | NO |
| Road | Normal | Snowy | YES |
| Road | Normal | Sunny | NO |
| Trail | Normal | Sunny | NO |



Is that what you would do?

50

## What would the tree look like for…

| Terrain | Unicycle-type | Weather | Go-For-Ride? |
|---|---|---|---|
| Trail | Mountain | Rainy | YES |
| Trail | Mountain | Sunny | YES |
| Road | Mountain | Snowy | YES |
| Road | Mountain | Sunny | YES |
| Trail | Normal | Snowy | NO |
| Trail | Normal | Rainy | NO |
| Road | Normal | Snowy | YES |
| Road | Normal | Sunny | NO |
| Trail | Normal | Sunny | NO |



Maybe…

51

## What would the tree look like for…

| Terrain | Unicycle-type | Weather | Go-For-Ride? |
|---|---|---|---|
| Trail | Mountain | Rainy | YES |
| Trail | Mountain | Sunny | YES |
| Road | Mountain | Snowy | YES |
| Road | Mountain | Sunny | YES |
| Trail | Normal | Snowy | NO |
| Trail | Normal | Rainy | NO |
| Road | Normal | Snowy | YES |
| Road | Normal | Sunny | NO |
| Trail | Normal | Sunny | NO |



An aside: how did we decide to pick the label for normal→road→rainy?

52

## Slide 53

### What would the tree look like for…

| Terrain | Unicycle-type | Weather | Jacket | ML grade | Go-For-Ride? |
|---------|---------------|---------|--------|----------|--------------|
| Trail | Mountain | Rainy | Heavy | D | YES |
| Trail | Mountain | Sunny | Light | C- | YES |
| Road | Mountain | Snowy | Light | B | YES |
| Road | Mountain | Sunny | Heavy | A | YES |
| … | Mountain | … | … | … | YES |
| Trail | Normal | Snowy | Light | D+ | NO |
| Trail | Normal | Rainy | Heavy | B- | NO |
| Road | Normal | Snowy | Heavy | C+ | YES |
| Road | Normal | Sunny | Light | A- | NO |
| Trail | Normal | Sunny | Heavy | B+ | NO |
| Trail | Normal | Snowy | Light | F | NO |
| … | Normal | … | … | … | NO |
| Trail | Normal | Rainy | Light | C | YES |

53

## Slide 54

### Overfitting

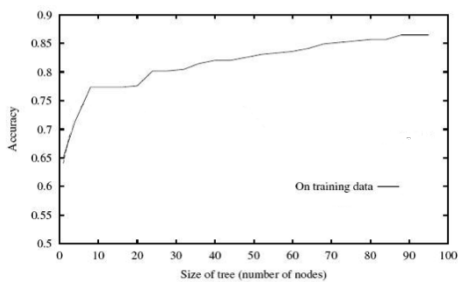| Terrain | Unicycle-type | Weather | Go-For-Ride? |
|---------|---------------|---------|--------------|
| Trail | Mountain | Rainy | YES |
| Trail | Mountain | Sunny | YES |
| Road | Mountain | Snowy | YES |
| Road | Mountain | Sunny | YES |
| Trail | Normal | Snowy | NO |
| Trail | Normal | Rainy | NO |
| Road | Normal | Snowy | YES |
| Road | Normal | Sunny | NO |
| Trail | Normal | Sunny | NO |

Unicycle
Mountain → YES
Normal → NO

*Overfitting* occurs when we bias our model too much towards the training data

Our goal is to learn a **general** model that will work on the training data as well as other data (i.e., test data)

54

## Slide 55

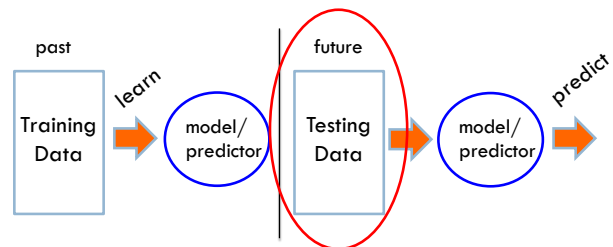### Overfitting

Our decision tree learning procedure always decreases training error
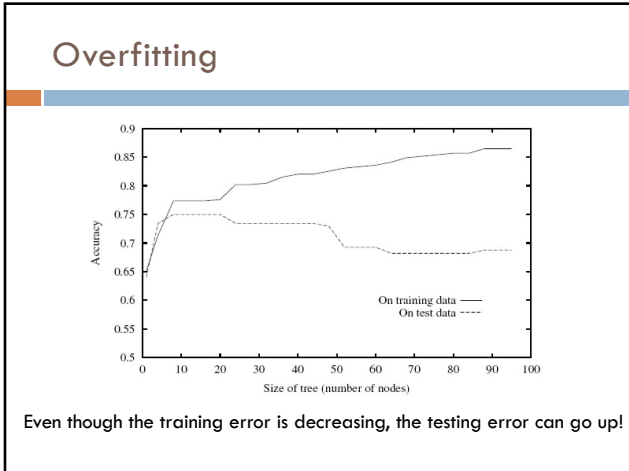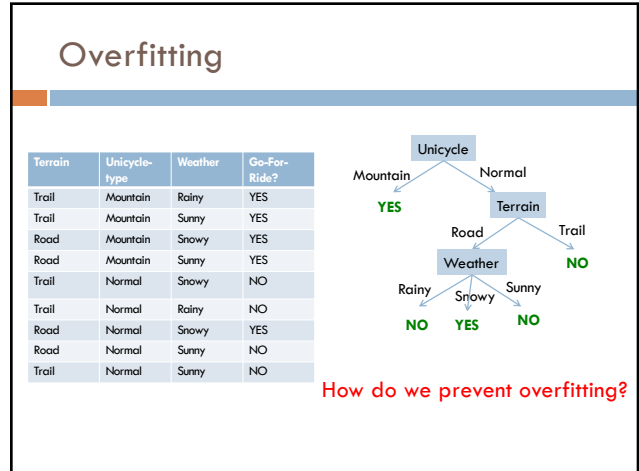
Is that what we want?

55

## Slide 56

### Test set error!

Machine learning is about predicting the future based on the past.
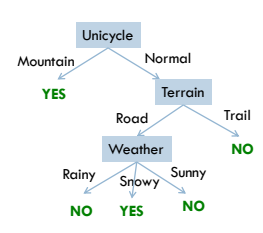-- Hal Daume III

past | future

Training Data → *learn* → model/predictor | Testing Data → model/predictor → *predict*

56

## Overfitting



Even though the training error is decreasing, the testing error can go up!

57

## Overfitting

| Terrain | Unicycle-type | Weather | Go-For-Ride? |
|---|---|---|---|
| Trail | Mountain | Rainy | YES |
| Trail | Mountain | Sunny | YES |
| Road | Mountain | Snowy | YES |
| Road | Mountain | Sunny | YES |
| Trail | Normal | Snowy | NO |
| Trail | Normal | Rainy | NO |
| Road | Normal | Snowy | YES |
| Road | Normal | Sunny | NO |
| Trail | Normal | Sunny | NO |



**How do we prevent overfitting?**

58

## Preventing overfitting

Base case:
- If all data belong to the same class, create a leaf node with that label
- **OR** all the data has the same feature values
- **OR** We've reached a particular depth in the tree
- ?

One idea: stop building the tree early

59

## Preventing overfitting

Base case:
- If all data belong to the same class, create a leaf node with that label
- **OR** all the data has the same feature values
- **OR** We've reached a particular depth in the tree
- We only have a certain number/fraction of examples remaining
- We've reached a particular training error
- Use development data (more on this later)
- …
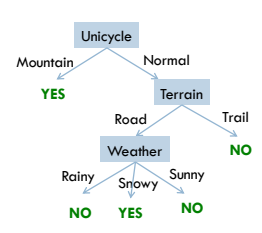
60

15

## Preventing overfitting: pruning



Pruning: after the tree is built, go back and "prune" the tree, i.e. remove some lower parts of the tree

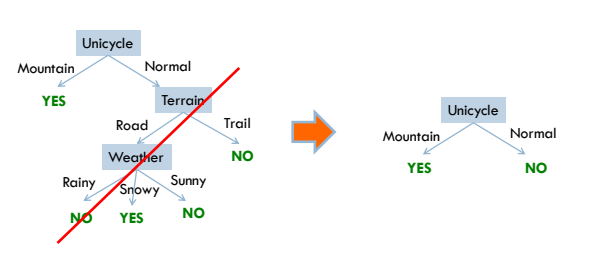Similar to stopping early, but done after the entire tree is built

61

## Preventing overfitting: pruning



Build the full tree
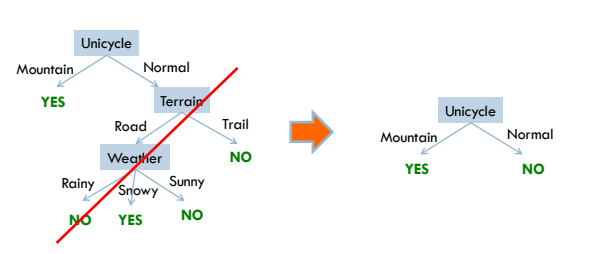
62

## Preventing overfitting: pruning



Build the full tree

Prune back leaves that are too specific

63

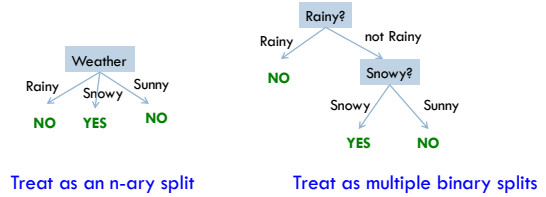## Preventing overfitting: pruning



Pruning criterion?

64

## Handling non-binary attributes

| PassengerId | Pclass | Sex | Age | SibSp | Parch | Ticket | Fare | Embarked | Survived |
|---|---|---|---|---|---|---|---|---|---|
| 804 | 3 | 0 | 0.42 | 0 | 1 | 2625 | 8.5167 | 0 | 1 |
| 756 | 2 | 0 | 0.67 | 1 | 1 | 250649 | 14.5 | 2 | 1 |
| 470 | 3 | 1 | 0.75 | 2 | 1 | 2666 | 19.2583 | 0 | 1 |
| 645 | 3 | 1 | 0.75 | 2 | 1 | 2666 | 19.2583 | 0 | 1 |
| 79 | 2 | 0 | 0.83 | 0 | 2 | 248738 | 29 | 2 | 1 |
| 832 | 2 | 0 | 0.83 | 1 | 1 | 29106 | 18.75 | 2 | 1 |
| 306 | 1 | 0 | 0.92 | 1 | 2 | 113781 | 151.55 | 2 | 1 |
| 165 | 3 | 0 | 1 | 4 | 1 | 3101295 | 39.6875 | 2 | 0 |
| 173 | 3 | 1 | 1 | 1 | 1 | 347742 | 11.1333 | 2 | 1 |
| 184 | 2 | 0 | 1 | 2 | 1 | 230136 | 39 | 2 | 1 |
| 382 | 3 | 1 | 1 | 0 | 2 | 2653 | 15.7417 | 0 | 1 |
| 387 | 3 | 0 | 1 | 5 | 2 | 2144 | 46.9 | 2 | 0 |
| 789 | 3 | 0 | 1 | 1 | 2 | 2315 | 20.575 | 2 | 1 |
| 828 | 2 | 0 | 1 | 0 | 2 | 2079 | 37.0042 | 0 | 1 |
| 8 | 3 | 0 | 2 | 3 | 1 | 349909 | 21.075 | 2 | 0 |
| 17 | 3 | 0 | 2 | 4 | 1 | 382652 | 29.125 | 1 | 0 |
| 120 | 3 | 1 | 2 | 4 | 2 | 347082 | 31.275 | 2 | 0 |
| 206 | 3 | 1 | 2 | 0 | 1 | 347054 | 10.4625 | 2 | 0 |
| 298 | 1 | 1 | 2 | 1 | 2 | 113781 | 151.55 | 2 | 0 |
| 341 | 2 | 0 | 2 | 1 | 1 | 230080 | 26 | 2 | 1 |
| 480 | 3 | 1 | 2 | 0 | 1 | 3101298 | 12.2875 | 2 | 1 |

**What do we do with features that have multiple values? Real-values?**

65

## Features with multiple values

Rainy?

Rainy — not Rainy

Weather

Rainy — Snowy — Sunny

**NO** — Snowy?

**NO** — **YES** — **NO**

Snowy — Sunny

**YES** — **NO**

**Treat as an n-ary split**     **Treat as multiple binary splits**

66

## Real-valued features

Use any comparison test ($>, <, \leq, \geq$) to split the data into two parts

Select a range filter, i.e. min < value < max

Fare < $20

Yes — No

Fare

0-10 — >50

10-20   20-50
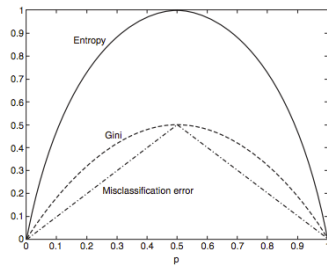
67

## Other splitting criterion

Otherwise:

- calculate the **"score"** for each feature if we used it to split the data
- pick the feature with the highest score, partition the data based on that data value and call recursively

**We used training error for the score. Any other ideas?**

68

## Other splitting criterion



- Entropy: how much uncertainty there is in the distribution over labels after the split
- Gini: sum of the square of the label proportions after split
- Training error = misclassification error

69

## Decision trees

Good?   Bad?



70

## Decision trees: the good

**Very intuitive and easy to interpret**

Fast to run and fairly easy to implement (Assignment 2 ☺)

Historically, perform fairly well (especially with a few more tricks we'll see later on)

No prior assumptions about the data

71

## Decision trees: the bad

Be careful with features with lots of values if you're not doing binary splits

| ID | Terrain | Unicycle-type | Weather | Go-For-Ride? |
|----|---------|---------------|---------|--------------|
| 1 | Trail | Normal | Rainy | NO |
| 2 | Road | Normal | Sunny | YES |
| 3 | Trail | Mountain | Sunny | YES |
| 4 | Road | Mountain | Rainy | YES |
| 5 | Trail | Normal | Snowy | NO |
| 6 | Road | Normal | Rainy | YES |
| 7 | Road | Mountain | Snowy | YES |
| 8 | Trail | Normal | Sunny | NO |
| 9 | Road | Normal | Snowy | NO |
| 10 | Trail | Mountain | Snowy | YES |

Which feature would be at the top here?

72

## Decision trees: the bad

Can be problematic (slow, bad performance) with large numbers of features

Can't learn some very simple data sets (e.g. some types of linearly separable data)

Pruning/tuning can be tricky to get right

73

## Final DT algorithm

DT_train(data):

Base cases:
1. If all data belong to the same class, pick that label
2. If all the data have the same feature values, pick majority label
3. If we're out of features to examine, pick majority label
4. If the we don't have any data left, pick majority label of *parent*
5. *If some other stopping criteria* exists to avoid overfitting, pick majority label

Otherwise (i.e. if none of the base cases apply):
- calculate the "score" for each feature if we used it to split the data
- pick the feature with the highest score, partition the data based on that data, e.g. data_left and data_right
- Recurse, i.e. DT_train(data_left) and DT_train(data_right)
- Make tree with feature as the splitting criterion with the decision trees returned from the recursive calls as the children

74