



Decision Trees

CS 229: Machine Learning

Emily Fox

Stanford University

February 12, 2024

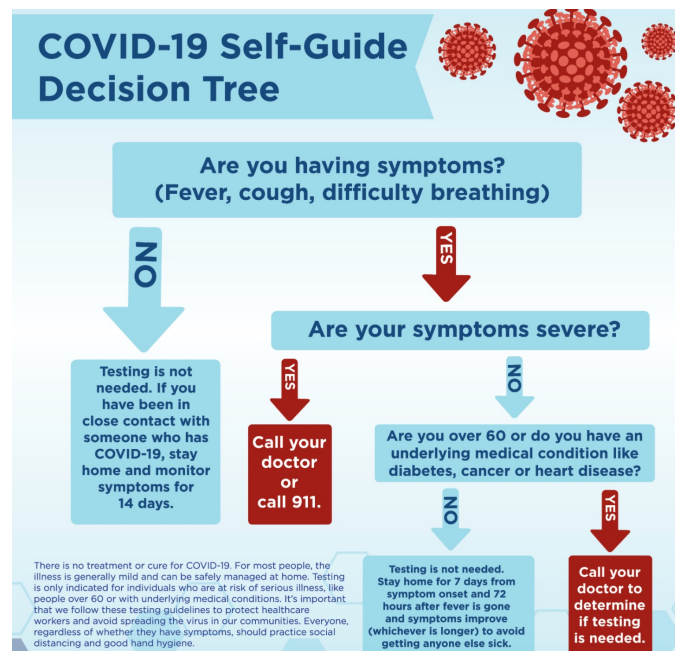
Slides include content developed by and co-developed with Carlos Guestrin

©2024 Emily Fox

1

How do
we make
decisions?

COVID-19 Self-Guide Decision Tree



<https://www.holzer.org/coronavirus-covid-19-updates/>

©2024 Emily Fox

CS 229: Machine Learning

2

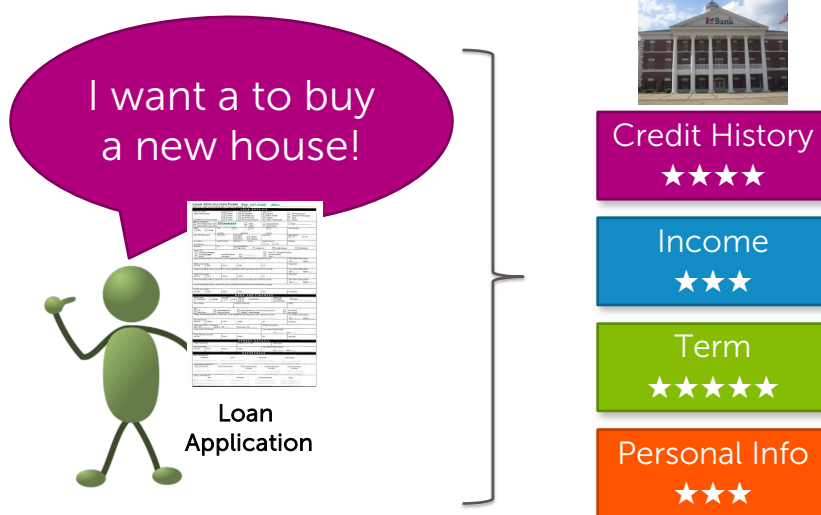
Predicting potential loan defaults

©2024 Emily Fox

CS 229: Machine Learning

3

What makes a loan risky?



©2024 Emily Fox

CS 229: Machine Learning

4

Credit history explained

Did I pay previous loans on time?

Example:
excellent, good, or fair



Credit History
★★★★

Income
★★★

Term
★★★★★

Personal Info
★★★

©2024 Emily Fox

CS 229: Machine Learning

5

Income

What's my income?

Example:
\$80K per year



Credit History
★★★★

Income
★★★

Term
★★★★★

Personal Info
★★★

©2024 Emily Fox

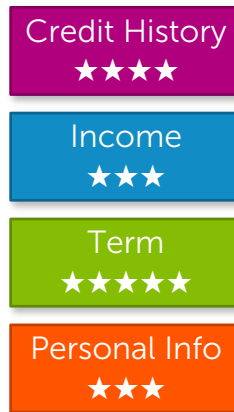
CS 229: Machine Learning

6

Loan terms

How soon do I need to pay the loan?

Example: 3 years,
5 years,...



©2024 Emily Fox

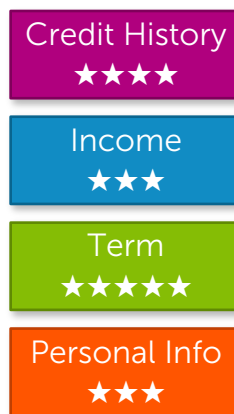
CS 229: Machine Learning

7

Personal information

Age, reason for the loan, marital status,...

Example: Home loan
for a married couple



©2024 Emily Fox

CS 229: Machine Learning

8

Intelligent application

Loan Applications

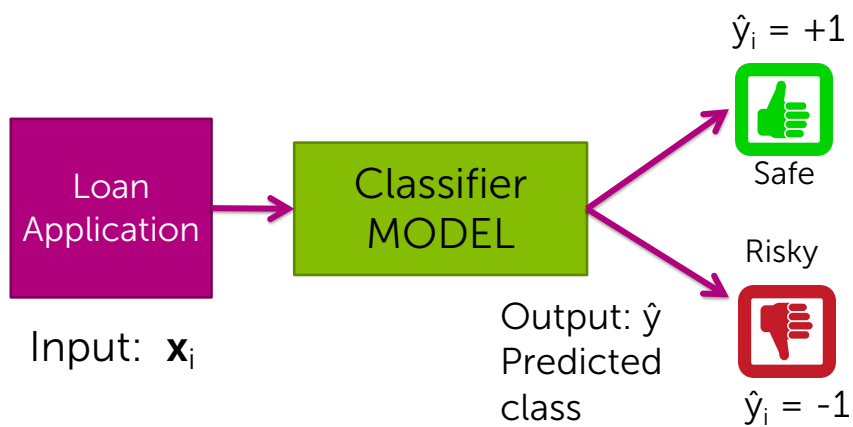


©2024 Emily Fox

CS 229: Machine Learning

9

Classifier review

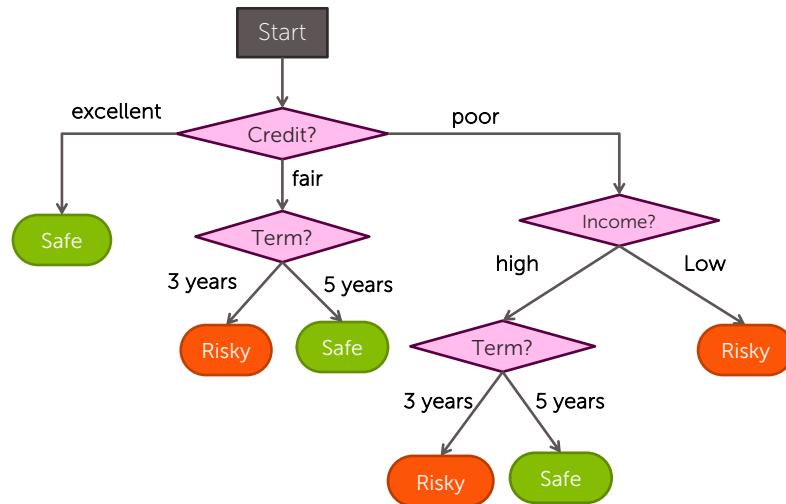


©2024 Emily Fox

CS 229: Machine Learning

10

This module ... decision trees

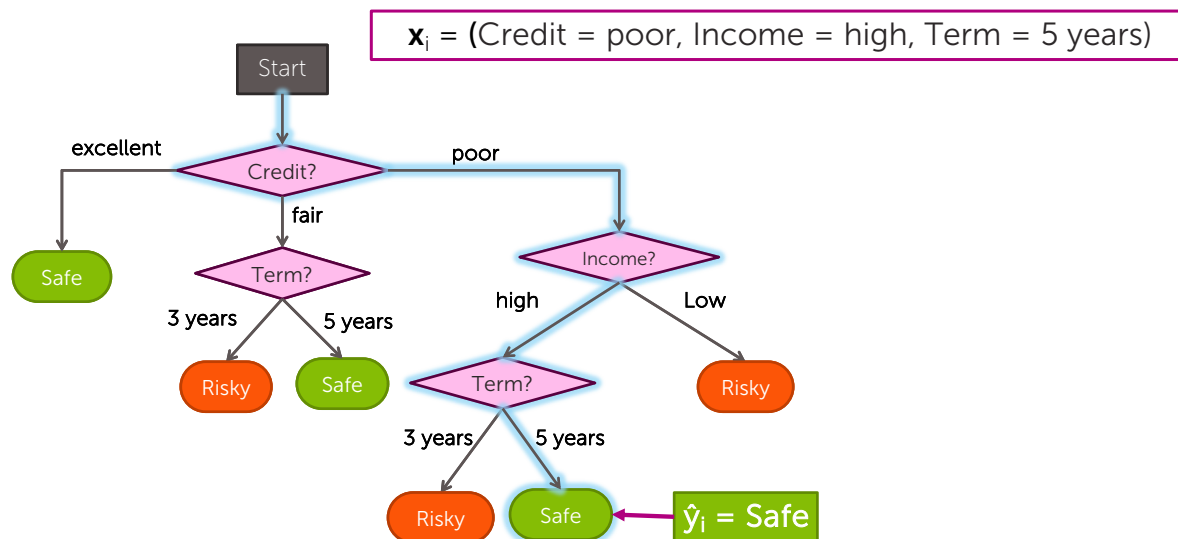


©2024 Emily Fox

CS 229: Machine Learning

11

Scoring a loan application



©2024 Emily Fox

CS 229: Machine Learning

12

Decision tree learning task

©2024 Emily Fox

CS 229: Machine Learning

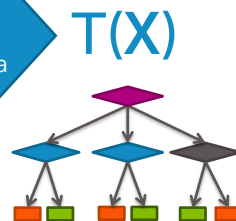
13

Decision tree learning problem

Training data: N observations (\mathbf{x}_i, y_i)

Credit	Term	Income	y
excellent	3 yrs	high	safe
fair	5 yrs	low	risky
fair	3 yrs	high	safe
poor	5 yrs	high	risky
excellent	3 yrs	low	risky
fair	5 yrs	low	safe
poor	3 yrs	high	risky
poor	5 yrs	low	safe
fair	3 yrs	high	safe

Optimize
cost function
on training data



©2024 Emily Fox

CS 229: Machine Learning

14

Cost function: Classification error

- Error measures fraction of mistakes

$$\text{Error} = \frac{\text{\# incorrect predictions}}{\text{\# examples}}$$

- Best possible value : 0.0
- Worst possible value: 1.0

©2024 Emily Fox

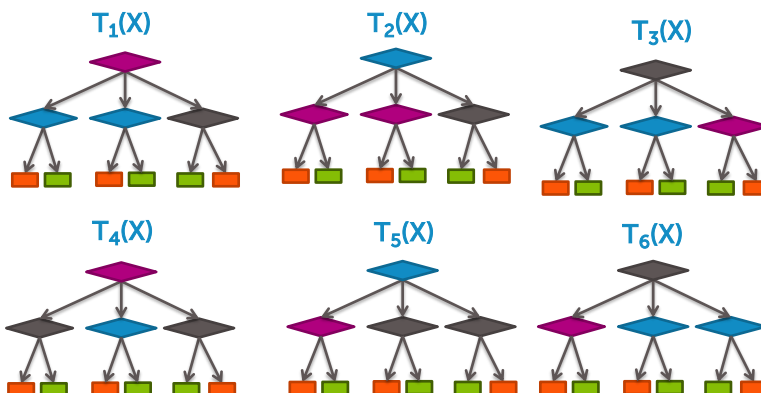
CS 229: Machine Learning

15

How do we find the best tree?

Exponentially large number of possible trees makes decision tree learning **hard**!

Learning the smallest decision tree is an *NP-hard* problem
[Hyafil & Rivest '76]



©2024 Emily Fox

CS 229: Machine Learning

16

Greedy decision tree learning

©2024 Emily Fox

CS 229: Machine Learning

17

Our training data table

Assume $N = 40$, 3 features

Credit	Term	Income	y
excellent	3 yrs	high	safe
fair	5 yrs	low	risky
fair	3 yrs	high	safe
poor	5 yrs	high	risky
excellent	3 yrs	low	risky
fair	5 yrs	low	safe
poor	3 yrs	high	risky
poor	5 yrs	low	safe
fair	3 yrs	high	safe

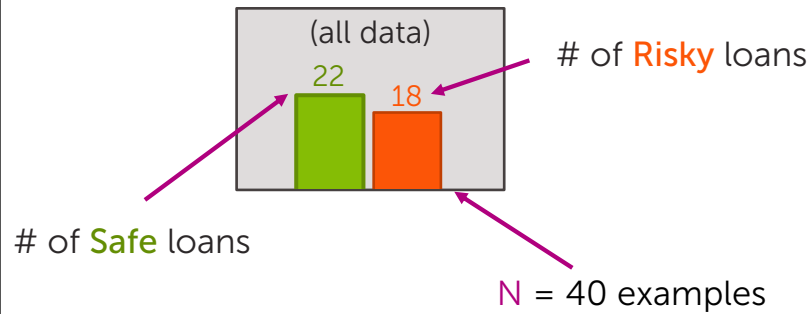
©2024 Emily Fox

CS 229: Machine Learning

18

Start with all the data

Loan status: **Safe** **Risky**



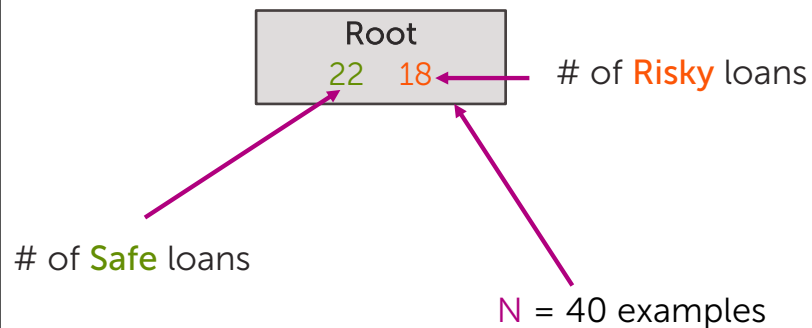
©2024 Emily Fox

CS 229: Machine Learning

19

Compact visual notation: Root node

Loan status: **Safe** **Risky**

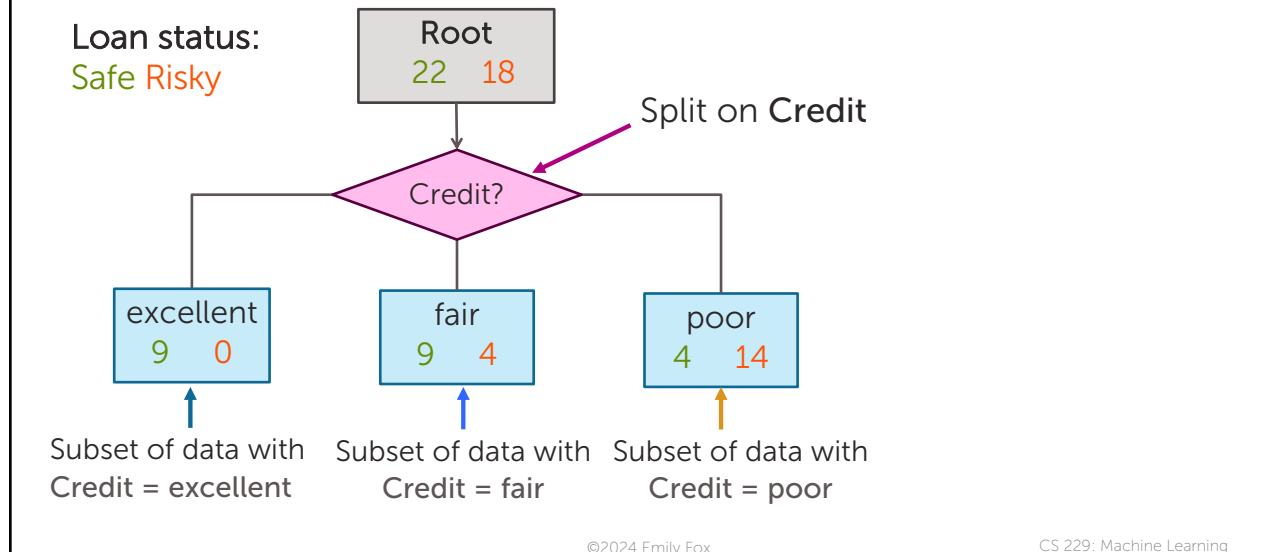


©2024 Emily Fox

CS 229: Machine Learning

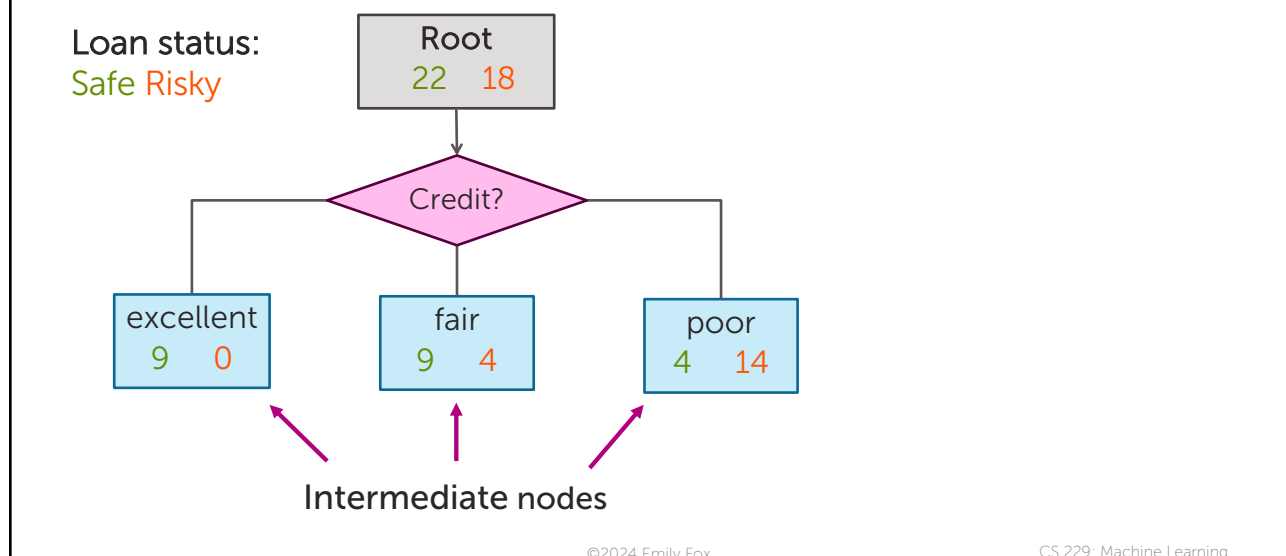
20

Decision stump: Single level tree



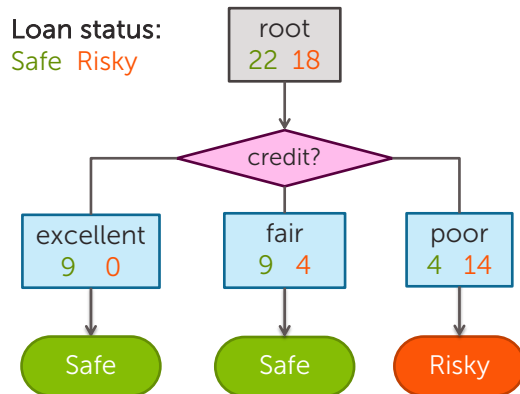
21

Visual notation: Intermediate nodes



22

Making predictions with a decision stump



For each intermediate node,
set \hat{y} = majority value

©2024 Emily Fox

CS 229: Machine Learning

23

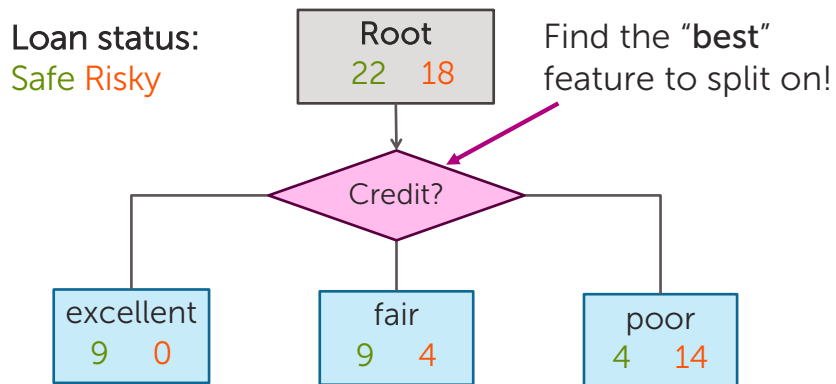
Selecting best feature to split on

©2024 Emily Fox

CS 229: Machine Learning

24

How do we learn a decision stump?

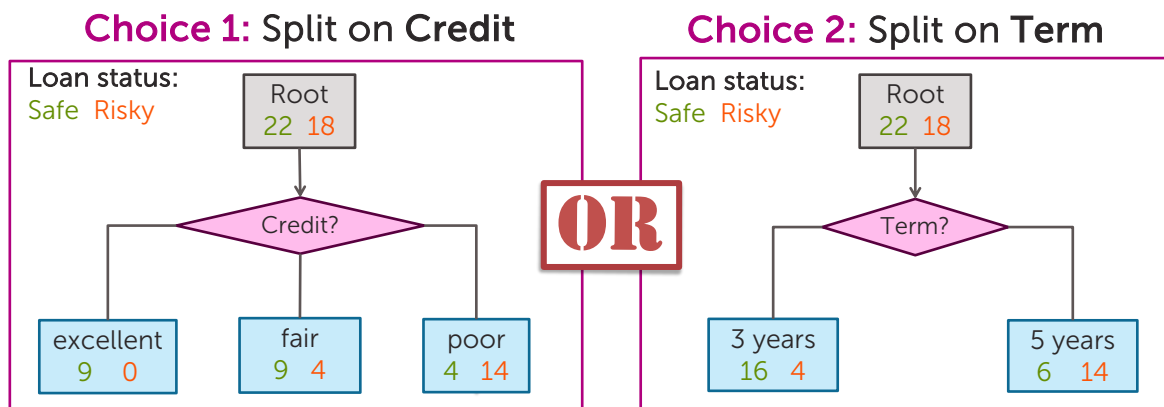


©2024 Emily Fox

CS 229: Machine Learning

25

How do we select the best feature?

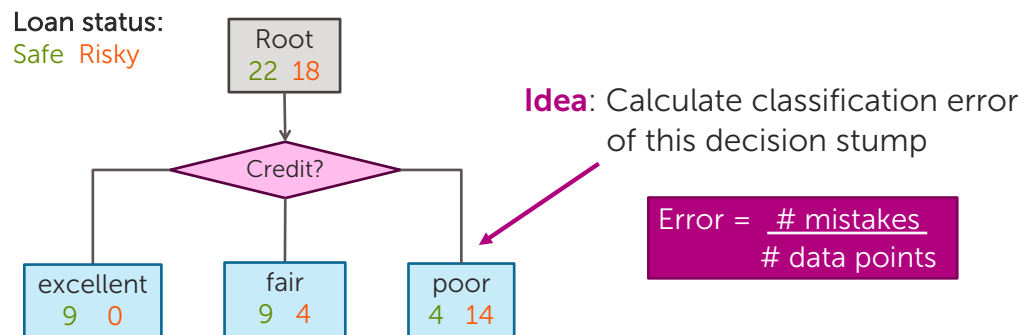


©2024 Emily Fox

CS 229: Machine Learning

26

How do we measure effectiveness of a split?



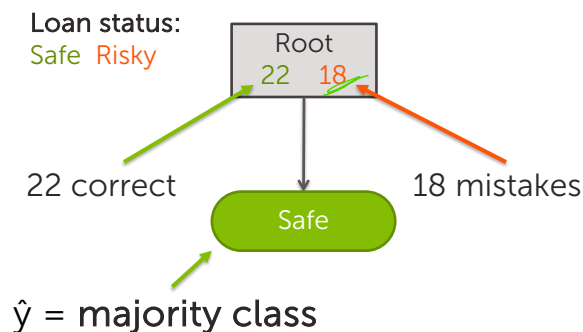
©2024 Emily Fox

CS 229: Machine Learning

27

Calculating classification error

- **Step 1:** \hat{y} = class of majority of data in node
- **Step 2:** Calculate classification error of predicting \hat{y} for this data



$$\text{Error} = \frac{18}{18 + 22 = 40} = 0.45$$

Tree	Classification error
(root)	0.45

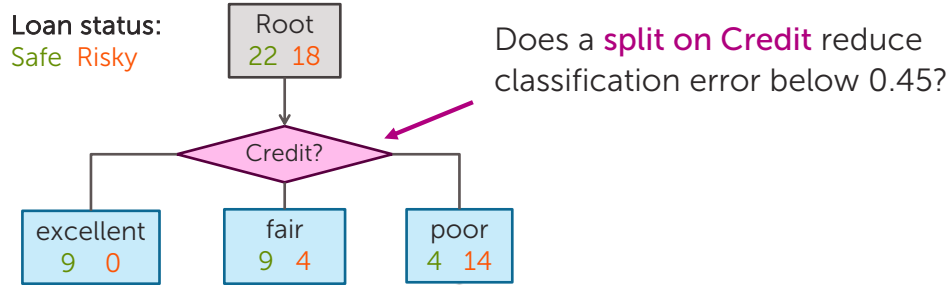
©2024 Emily Fox

CS 229: Machine Learning

28

Choice 1: Split on **Credit** history?

Choice 1: Split on Credit



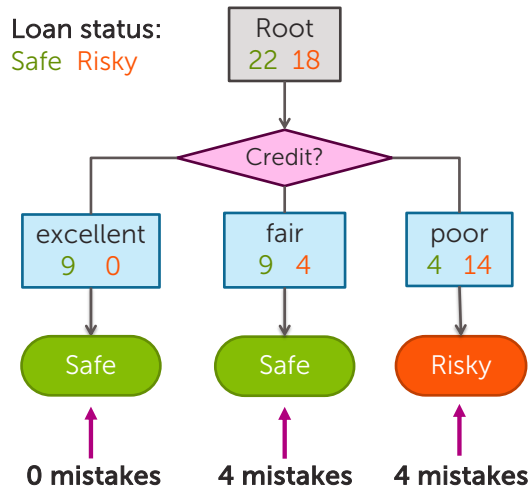
©2024 Emily Fox

CS 229: Machine Learning

29

Split on **Credit**: Classification error

Choice 1: Split on Credit



$$\text{Error} = \frac{0 + 4 + 4}{40} = 0.2$$

Tree	Classification error
(root)	0.45
Split on credit	0.2

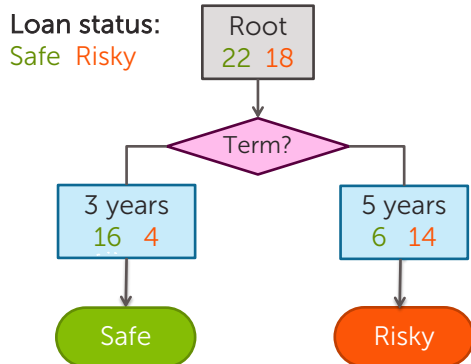
©2024 Emily Fox

CS 229: Machine Learning

30

Choice 2: Split on Term?

Choice 2: Split on Term



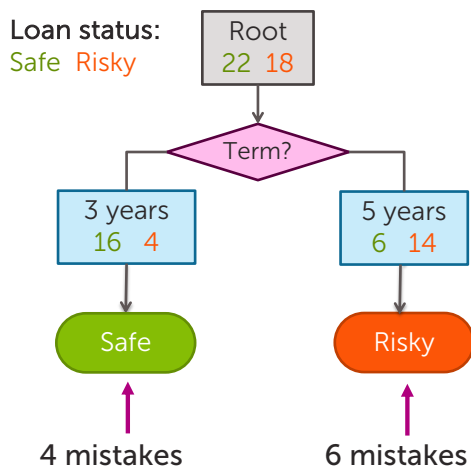
©2024 Emily Fox

CS 229: Machine Learning

31

Evaluating the split on Term

Choice 2: Split on Term



$$\text{Error} = \frac{4+6}{40} = 0.25$$

Tree	Classification error
(root)	0.45
Split on credit	0.2
Split on term	0.25

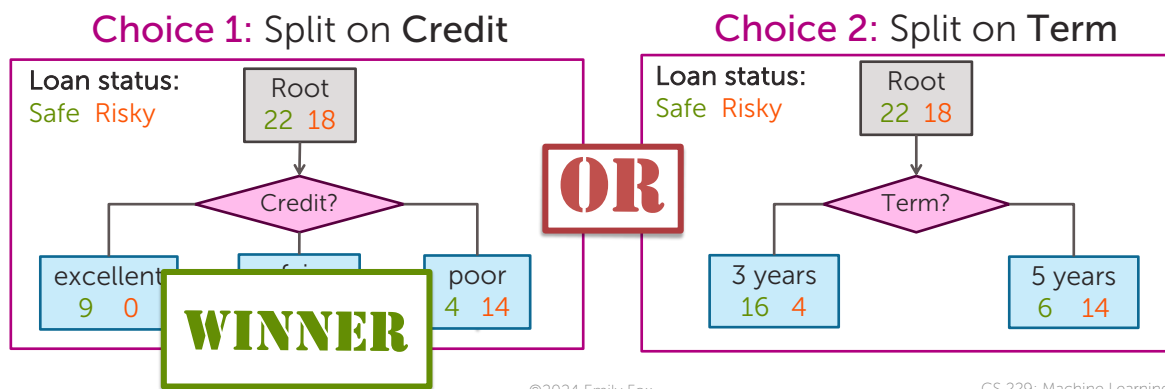
©2024 Emily Fox

CS 229: Machine Learning

32

Choice 1 vs Choice 2: Comparing split on **Credit** vs **Term**

Tree	Classification error
(root)	0.45
split on credit	0.2
split on loan term	0.25



33

Feature split selection algorithm

- Given a subset of data M (a node in a tree)
- For each feature $h_i(x)$:
 - Split data of M according to feature $h_i(x)$
 - Compute classification error of split
- Chose feature $h^*(x)$ with lowest classification error

©2024 Emily Fox

CS 229: Machine Learning

34

Recursion & Stopping conditions

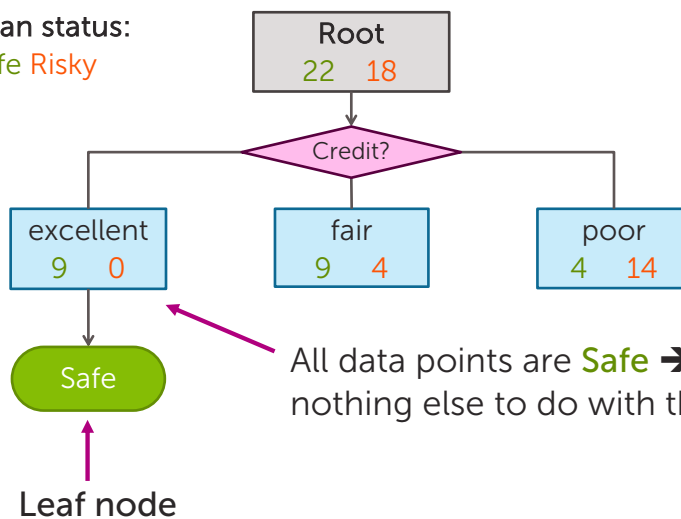
©2024 Emily Fox

CS 229: Machine Learning

35

We've learned a decision stump, what next?

Loan status:
Safe Risky



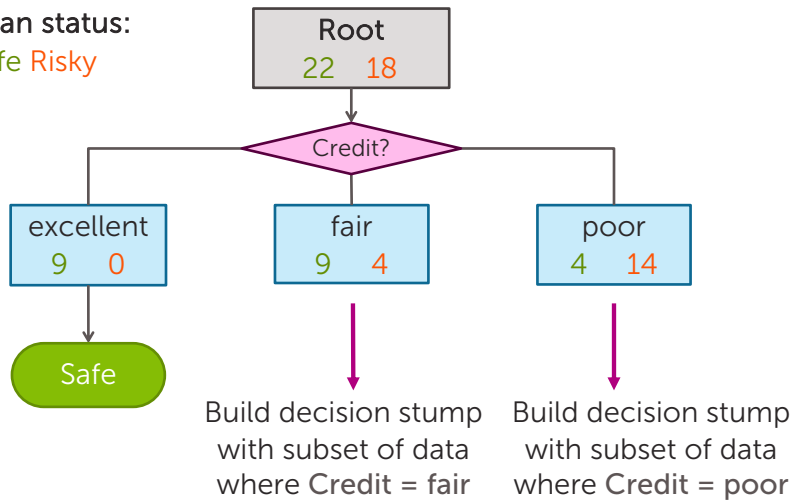
©2024 Emily Fox

CS 229: Machine Learning

36

Tree learning = Recursive stump learning

Loan status:
Safe Risky



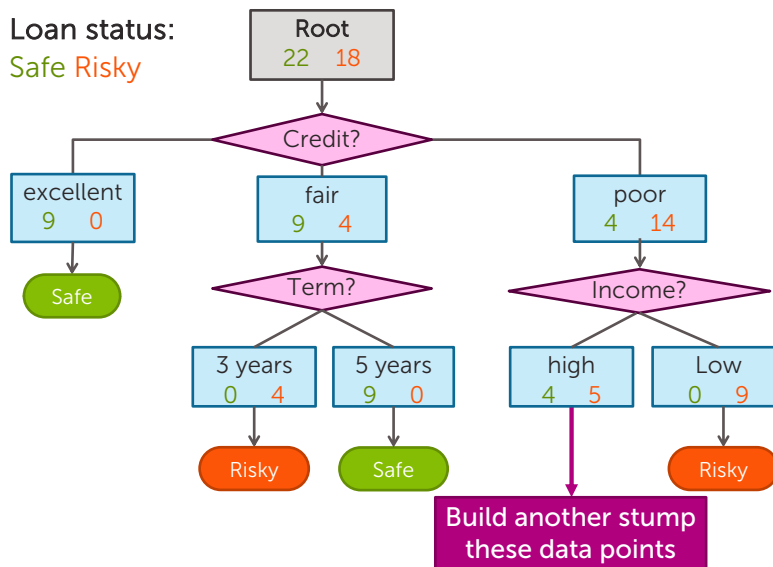
©2024 Emily Fox

CS 229: Machine Learning

37

Second level

Loan status:
Safe Risky

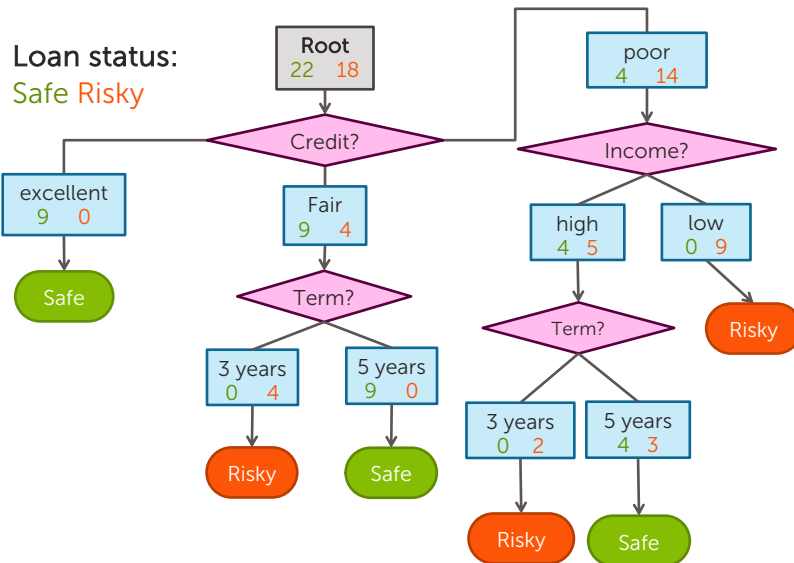


©2024 Emily Fox

CS 229: Machine Learning

38

Final decision tree



©2024 Emily Fox

CS 229: Machine Learning

39

Simple greedy decision tree learning

Pick best feature to split on

Learn decision stump with this split

For each leaf of decision stump,
recurse

When do we stop???

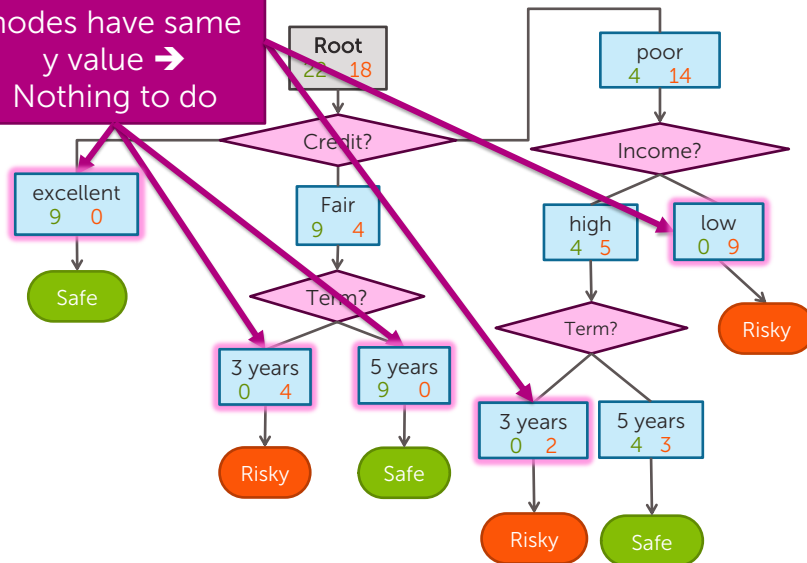
©2024 Emily Fox

CS 229: Machine Learning

40

Stopping condition 1: All data agrees on y

All data in these nodes have same y value → Nothing to do



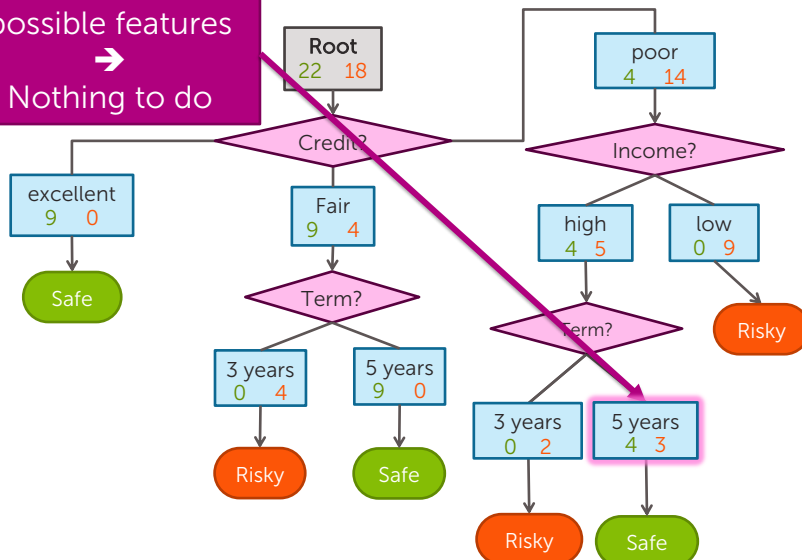
©2024 Emily Fox

CS 229: Machine Learning

41

Stopping condition 2: Already split on all features

Already split on all possible features → Nothing to do



©2024 Emily Fox

CS 229: Machine Learning

42

Greedy decision tree learning

- **Step 1:** Start with an empty tree

- **Step 2:** Select a feature to split data

- For each split of the tree:

- **Step 3:** If nothing more to, make predictions

- **Step 4:** Otherwise, go to **Step 2** & continue (recurse) on this split

Pick feature split leading to lowest classification error

Stopping conditions 1 & 2

Recursion

©2024 Emily Fox

CS 229: Machine Learning

43

Is this a good idea?

Proposed stopping condition 3:
Stop if no split reduces the classification error

©2024 Emily Fox

CS 229: Machine Learning

44

Stopping condition 3:

Stop if error doesn't decrease???

$$y = x[1] \text{ xor } x[2]$$

x[1]	x[2]	y
False	False	False
False	True	True
True	False	True
True	True	False

X

y values
True False

Root
2 2

$$\text{Error} = \frac{2}{4} = 0.5$$

Tree	Classification error
(root)	0.5

©2024 Emily Fox

CS 229: Machine Learning

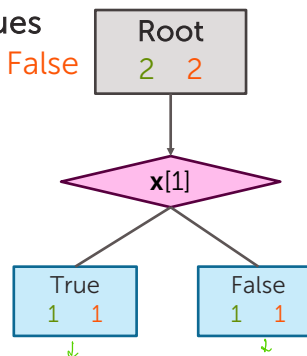
45

Consider split on x[1]

$$y = x[1] \text{ xor } x[2]$$

x[1]	x[2]	y
False	False	False
False	True	True
True	False	True
True	True	False

y values
True False



$$\text{Error} = \frac{1+1}{4} = 0.5$$

Tree	Classification error
(root)	0.5
Split on x[1]	0.5

©2024 Emily Fox

CS 229: Machine Learning

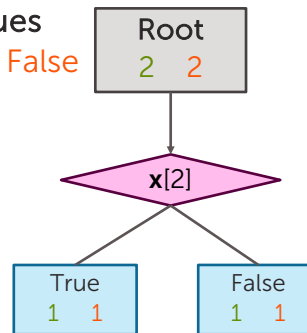
46

Consider split on $x[2]$

$$y = x[1] \text{ xor } x[2]$$

$x[1]$	$x[2]$	y
False	False	False
False	True	True
True	False	True
True	True	False

y values
True False



$$\text{Error} = \frac{1+1}{2+2} = 0.5$$

Neither features
improve training error...
Stop now???

Tree	Classification error
(root)	0.5
Split on $x[1]$	0.5
Split on $x[2]$	0.5

©2024 Emily Fox

CS 229: Machine Learning

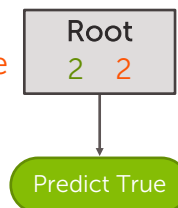
47

Final tree with stopping condition 3

$$y = x[1] \text{ xor } x[2]$$

$x[1]$	$x[2]$	y
False	False	False
False	True	True
True	False	True
True	True	False

y values
True False



Tree	Classification error
with stopping condition 3	0.5

©2024 Emily Fox

CS 229: Machine Learning

48

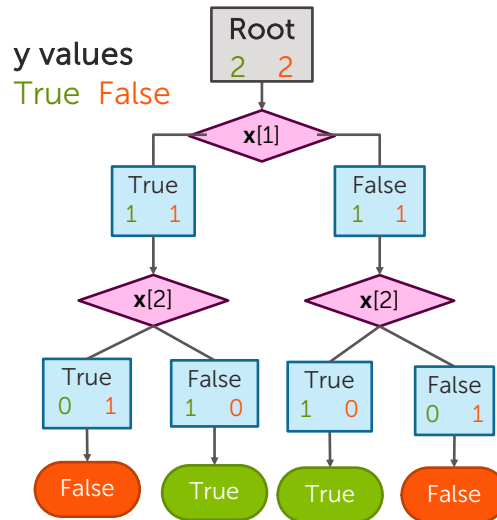
Without stopping condition 3

Condition 3 (stopping when training error doesn't improve) is **not** recommended!

$$y = x[1] \text{ xor } x[2]$$

x[1]	x[2]	y
False	False	False
False	True	True
True	False	True
True	True	False

Tree	Classification error
with stopping condition 3	0.5
without stopping condition 3	0



©2024 Emily Fox

CS 229: Machine Learning

49

Decision tree learning:
Real valued features

©2024 Emily Fox

CS 229: Machine Learning

50

How do we use real values inputs?

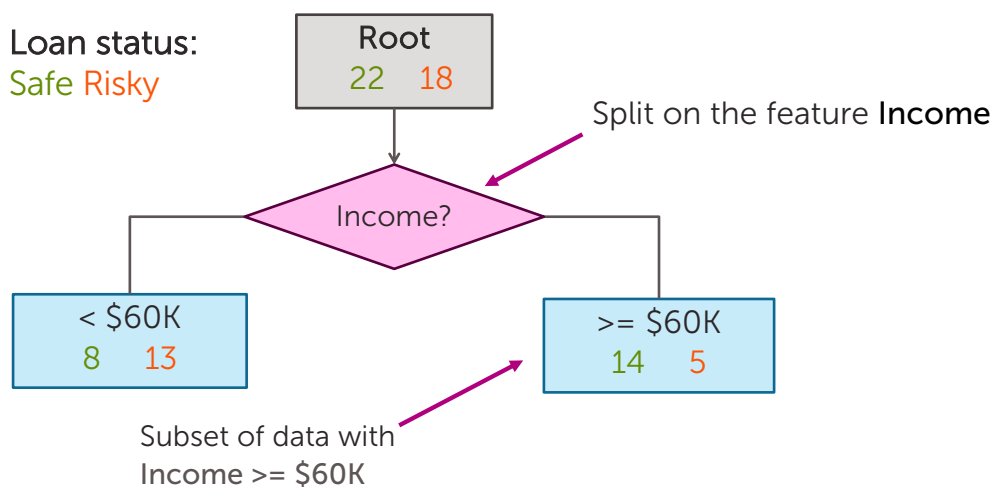
Income	Credit	Term	y
\$105 K	excellent	3 yrs	Safe
\$112 K	good	5 yrs	Risky
\$73 K	fair	3 yrs	Safe
\$69 K	excellent	5 yrs	Safe
\$217 K	excellent	3 yrs	Risky
\$120 K	good	5 yrs	Safe
\$64 K	fair	3 yrs	Risky
\$340 K	excellent	5 yrs	Safe
\$60 K	good	3 yrs	Risky

©2024 Emily Fox

CS 229: Machine Learning

51

Threshold split



©2024 Emily Fox

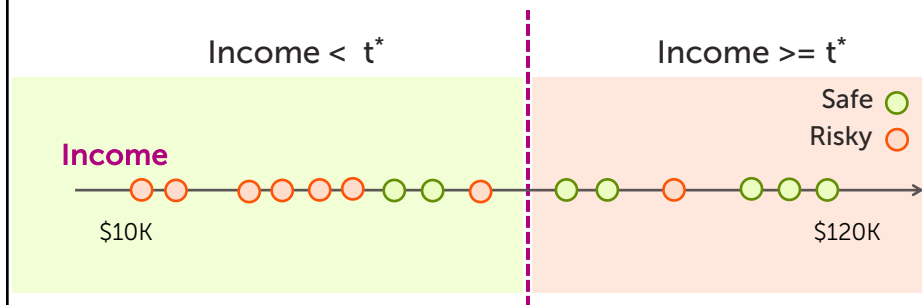
CS 229: Machine Learning

52

Finding the best threshold split

Infinite possible
values of t

Income = t^* *threshold to choose*



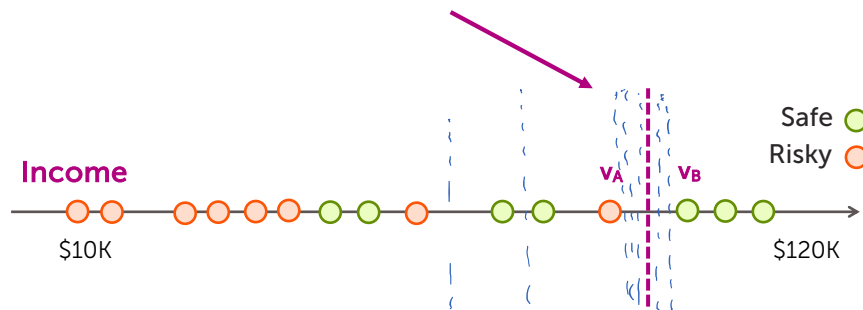
©2024 Emily Fox

CS 229: Machine Learning

53

Consider a threshold between points

Same **classification error** for any
threshold split between v_A and v_B



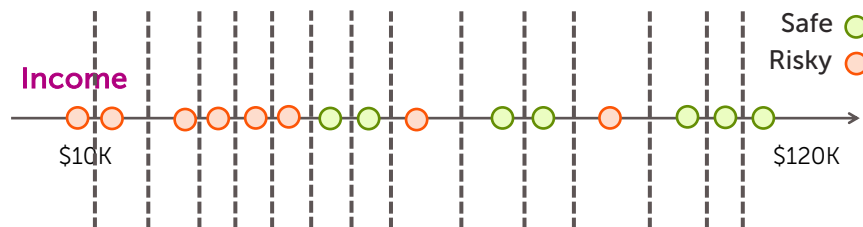
©2024 Emily Fox

CS 229: Machine Learning

54

Only need to consider mid-points

Finite number of splits to consider



©2024 Emily Fox

CS 229: Machine Learning

55

Threshold split selection algorithm

- **Step 1:** Sort the values of a feature $h_j(\mathbf{x})$:
Let $\{v_1, v_2, v_3, \dots, v_N\}$ denote sorted values
- **Step 2:**
 - For $i = 1 \dots N-1$
 - Consider split $t_i = (v_i + v_{i+1}) / 2$ ← midpoint
 - Compute classification error for threshold split $h_j(\mathbf{x}) \geq t_i$
 - Chose the t^* with the lowest classification error

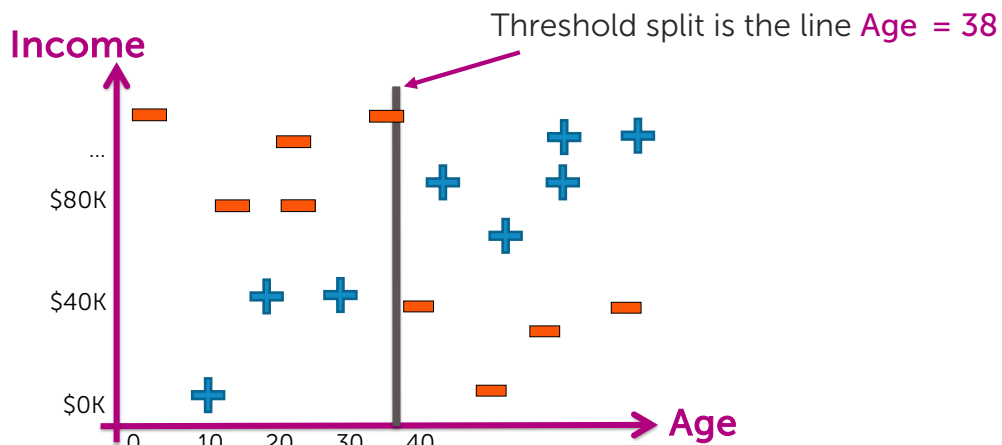
	$h_1(x)$	$h_2(x)$	$h_3(x)$...	$h_{10}(x)$
t^*	39 yrs	\$60k			
error	0.1	0.4			

©2024 Emily Fox

CS 229: Machine Learning

56

Visualizing the threshold split

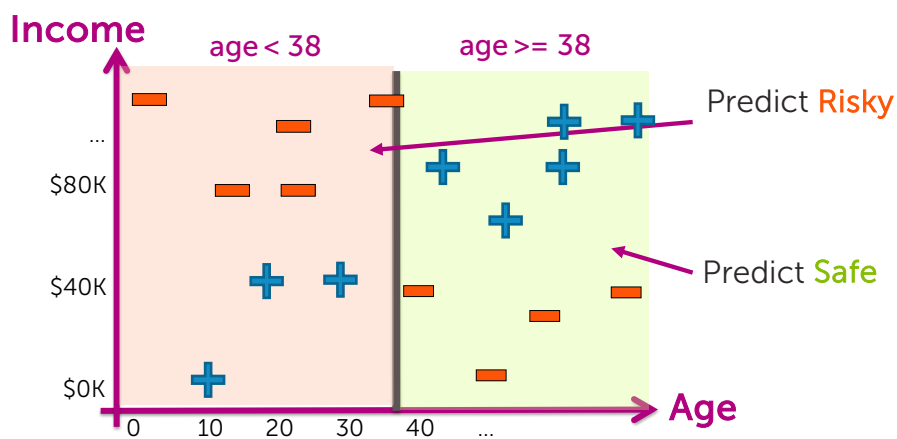


©2024 Emily Fox

CS 229: Machine Learning

57

Split on Age ≥ 38

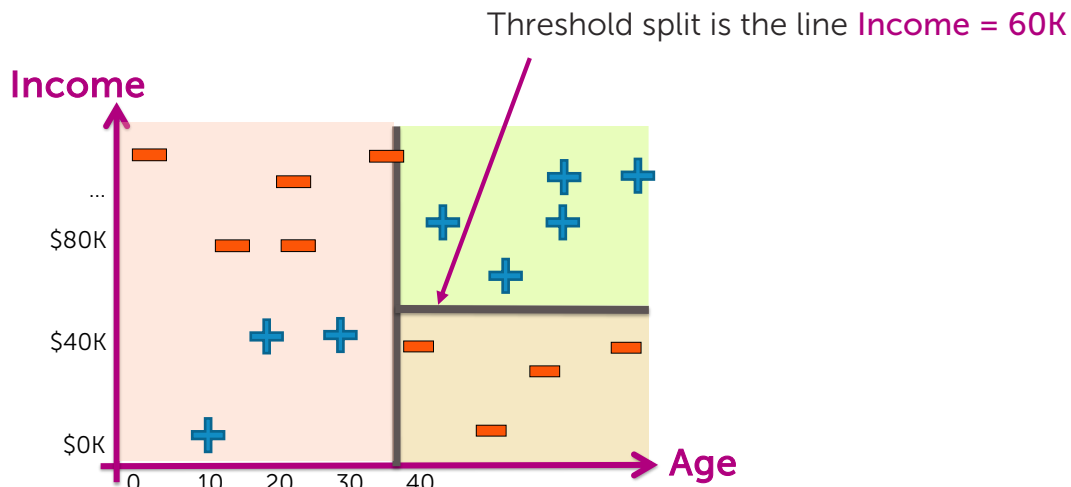


©2024 Emily Fox

CS 229: Machine Learning

58

Depth 2: Split on Income \geq \$60K

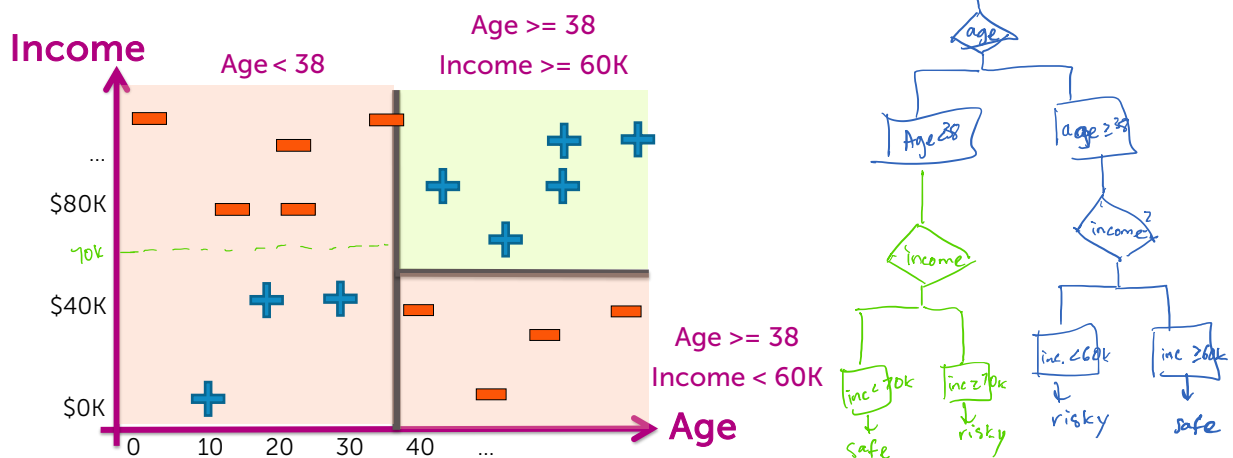


©2024 Emily Fox

CS 229: Machine Learning

59

Each split partitions the 2-D space



©2024 Emily Fox

CS 229: Machine Learning

60

Decision trees vs logistic regression: *Example*

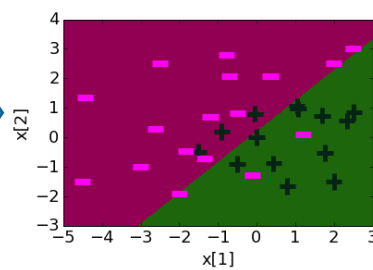
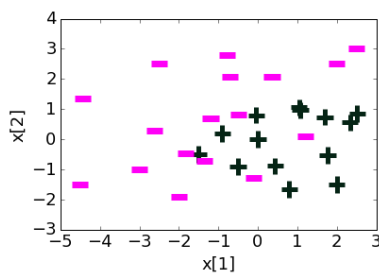
©2024 Emily Fox

CS 229: Machine Learning

61

Logistic regression

Feature	Value	Weight Learned
$h_0(\mathbf{x})$	1	0.22
$h_1(\mathbf{x})$	$\mathbf{x}[1]$	1.12
$h_2(\mathbf{x})$	$\mathbf{x}[2]$	-1.07

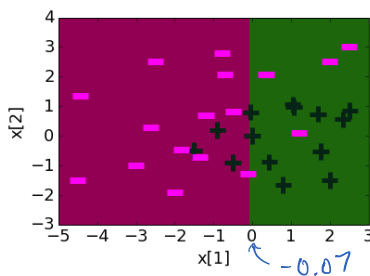
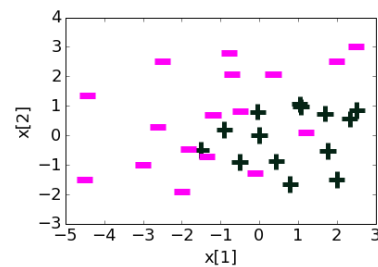


©2024 Emily Fox

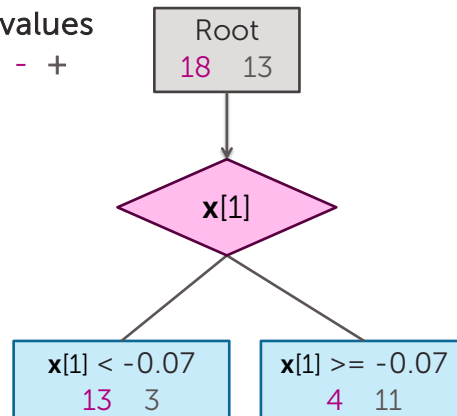
CS 229: Machine Learning

62

Depth 1: Split on $x[1]$



y values
- +

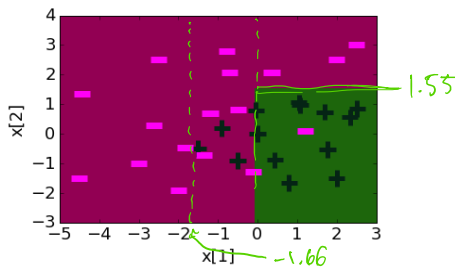
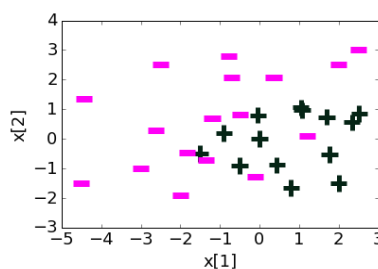


©2024 Emily Fox

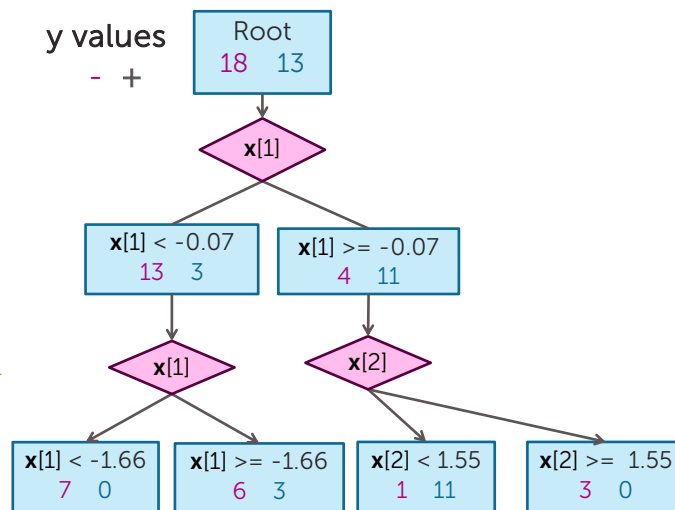
CS 229: Machine Learning

63

Depth 2



y values
- +

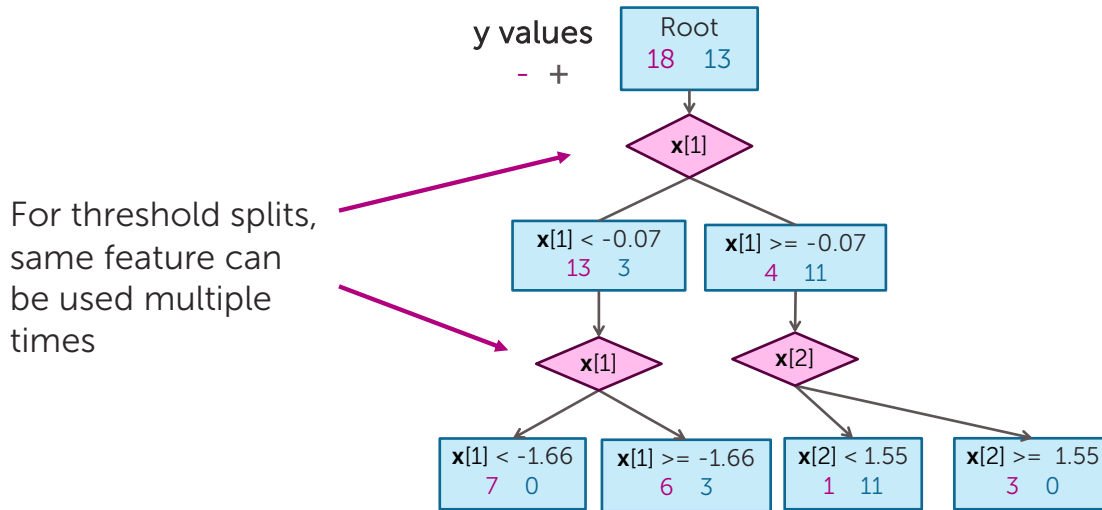


©2024 Emily Fox

CS 229: Machine Learning

64

Threshold split caveat

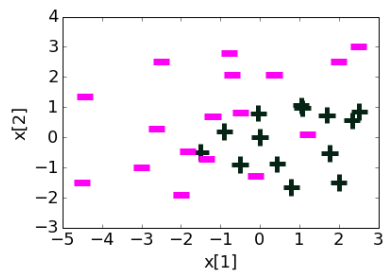


©2024 Emily Fox

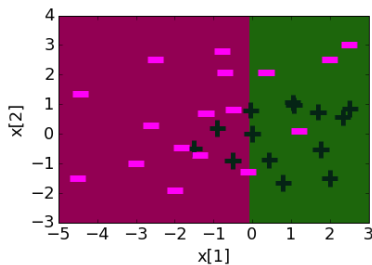
CS 229: Machine Learning

65

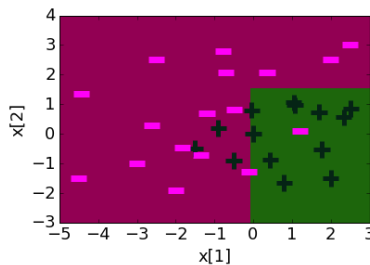
Decision boundaries



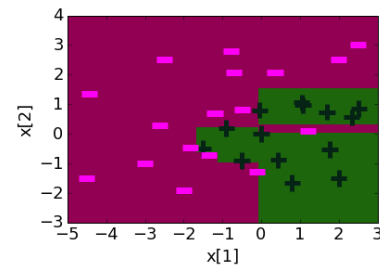
Depth 1



Depth 2



Depth 10

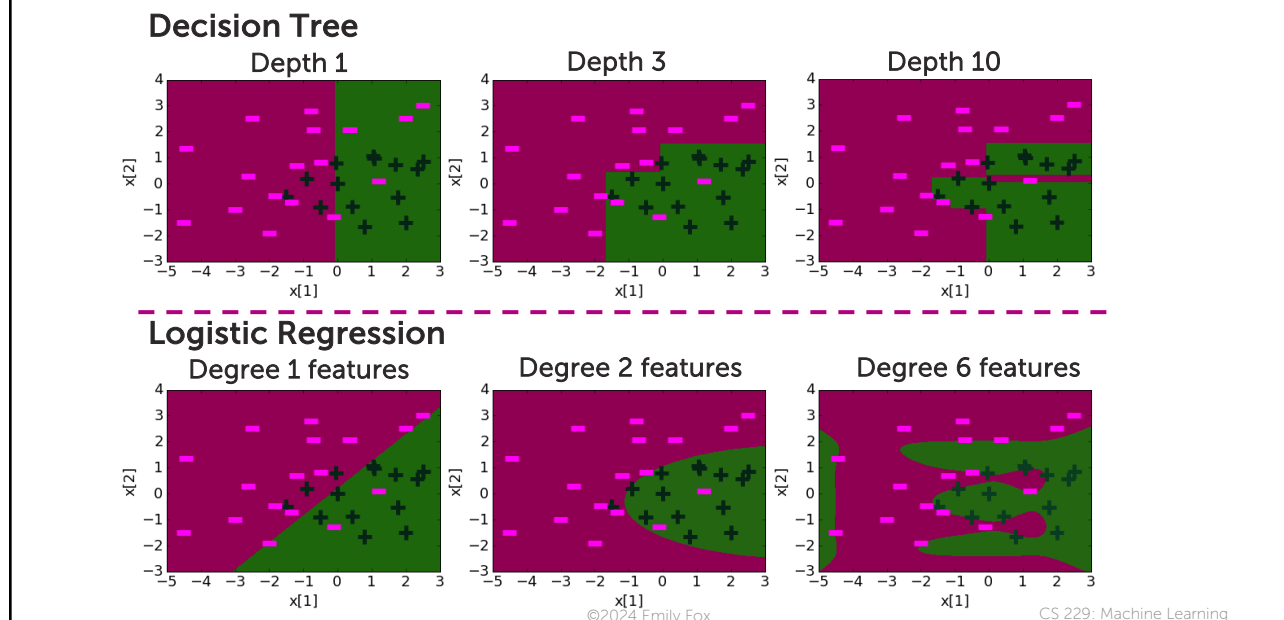


©2024 Emily Fox

CS 229: Machine Learning

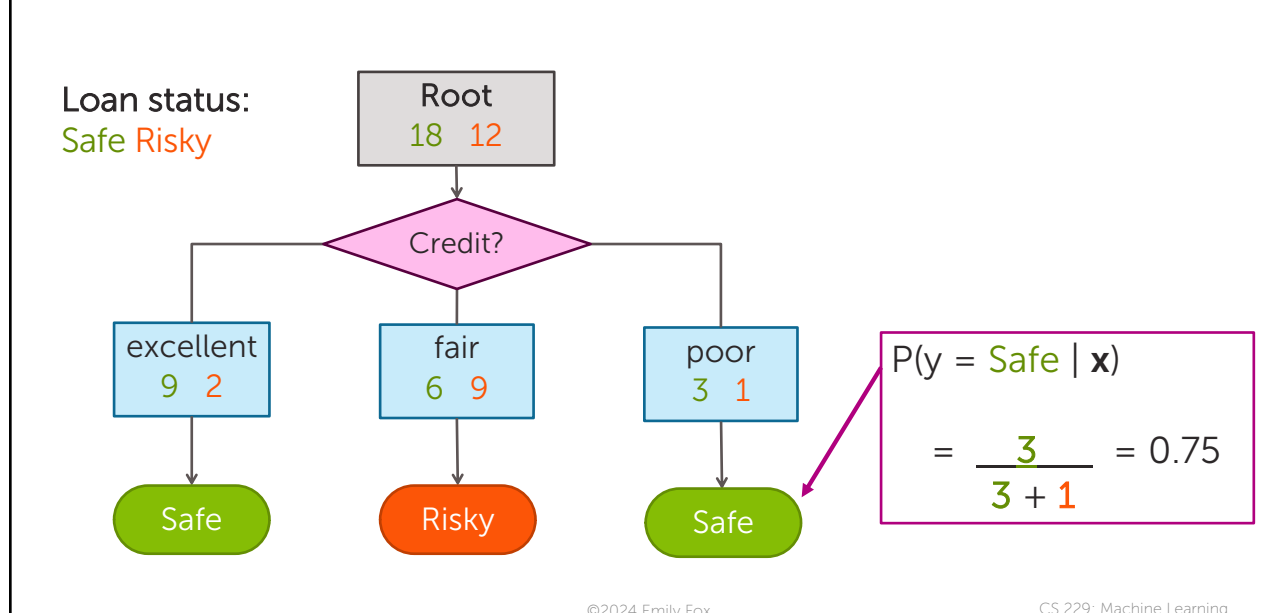
66

Comparing decision boundaries



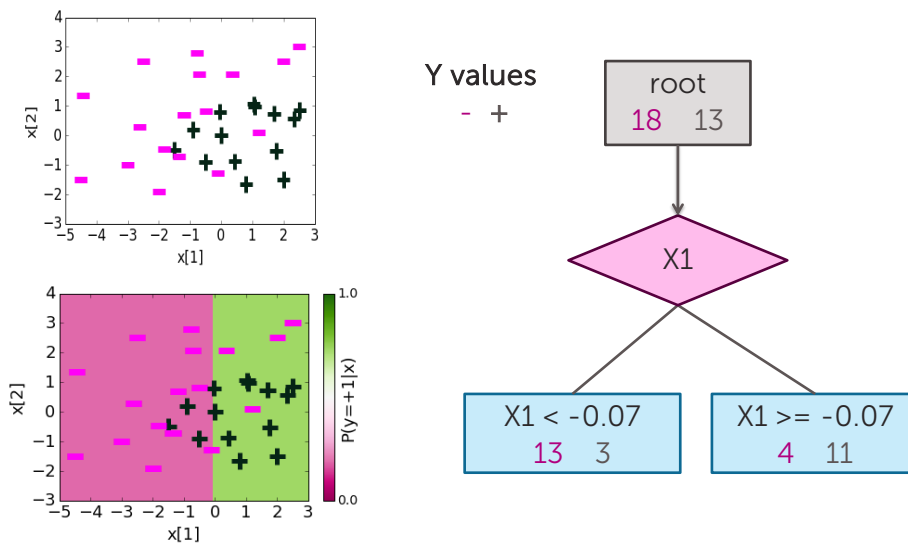
67

Predicting probabilities with decision trees



68

Depth 1 probabilities

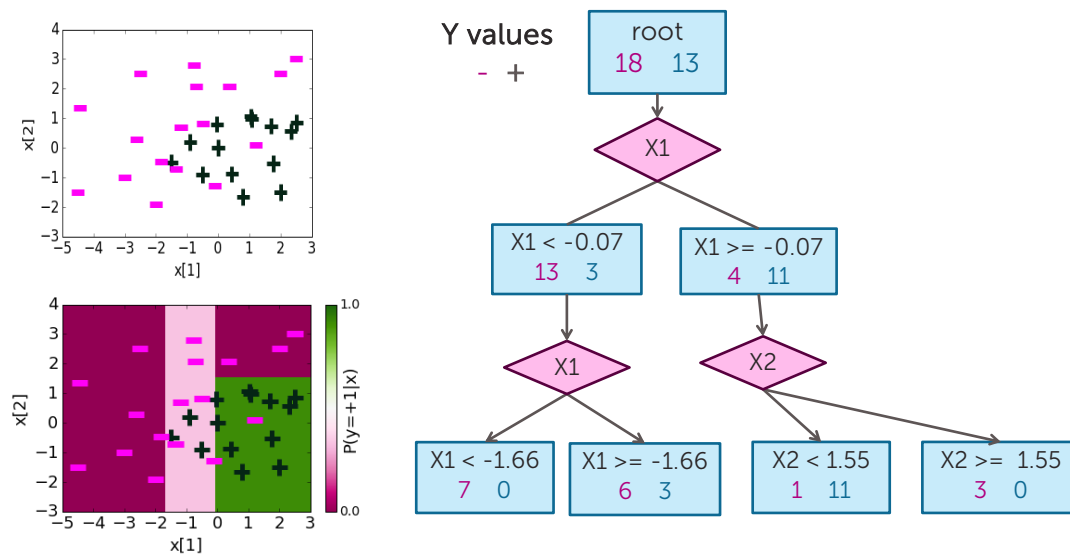


©2024 Emily Fox

CS 229: Machine Learning

69

Depth 2 probabilities

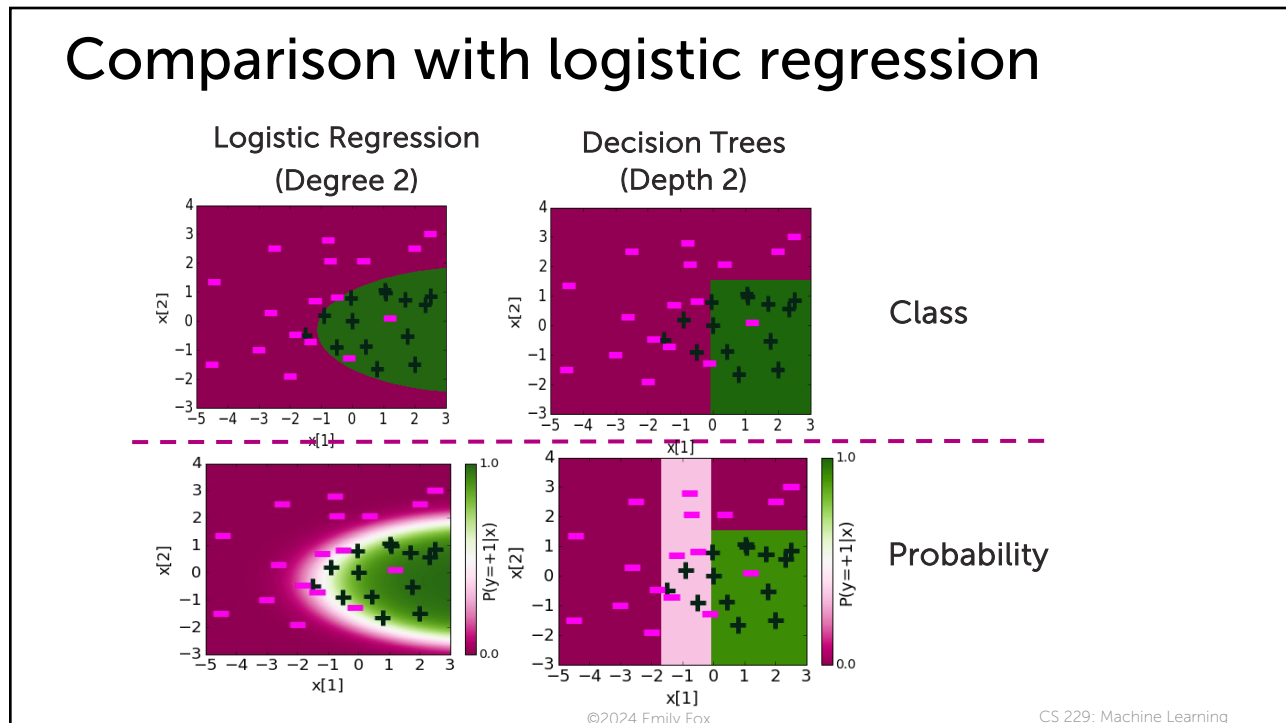


©2024 Emily Fox

CS 229: Machine Learning

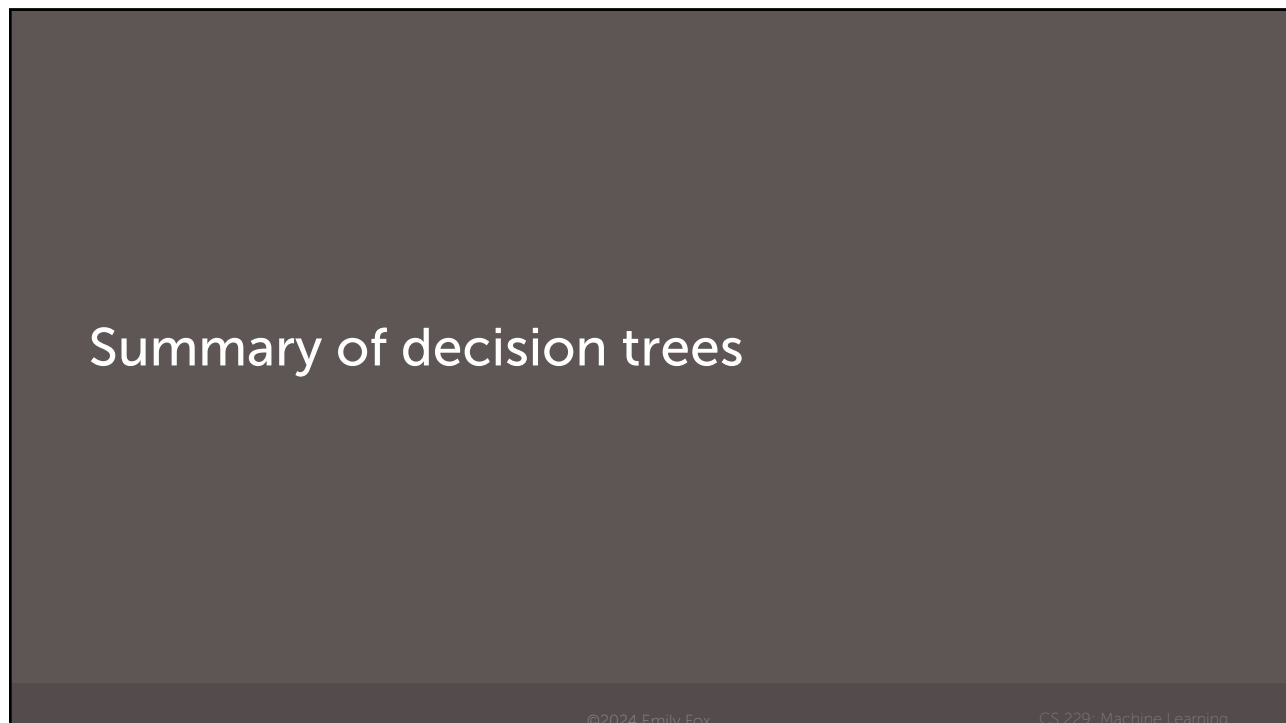
70

Comparison with logistic regression



71

Summary of decision trees



72

What you can do now

- Define a decision tree classifier
- Interpret the output of a decision trees
- Learn a decision tree classifier using greedy algorithm
- Traverse a decision tree to make predictions
 - Majority class predictions

©2024 Emily Fox

CS 229: Machine Learning

73



Decision Trees: Overfitting

CS 229: Machine Learning
Emily Fox
Stanford University
February 12, 2024

Slides include content developed by and co-developed with Carlos Guestrin

©2024 Emily Fox

74

Overfitting in decision trees

©2024 Emily Fox

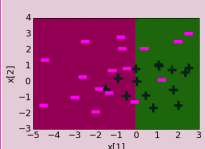
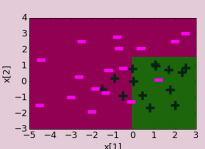
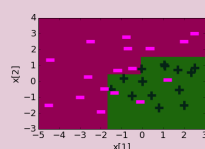
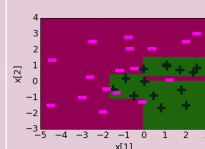
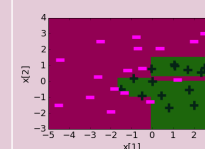
CS 229: Machine Learning

75

What happens when we increase depth?

Training error reduces with depth



Tree depth	depth = 1	depth = 2	depth = 3	depth = 5	depth = 10
Training error	0.22	0.13	0.10	0.03	0.00
Decision boundary					

big warning!

©2024 Emily Fox

CS 229: Machine Learning

76

Two approaches to picking simpler trees

1. **Early Stopping:**
Stop the learning algorithm **before** tree becomes too complex
2. **Pruning:**
Simplify the tree **after** the learning algorithm terminates

©2024 Emily Fox

CS 229: Machine Learning

77

Technique 1: Early stopping

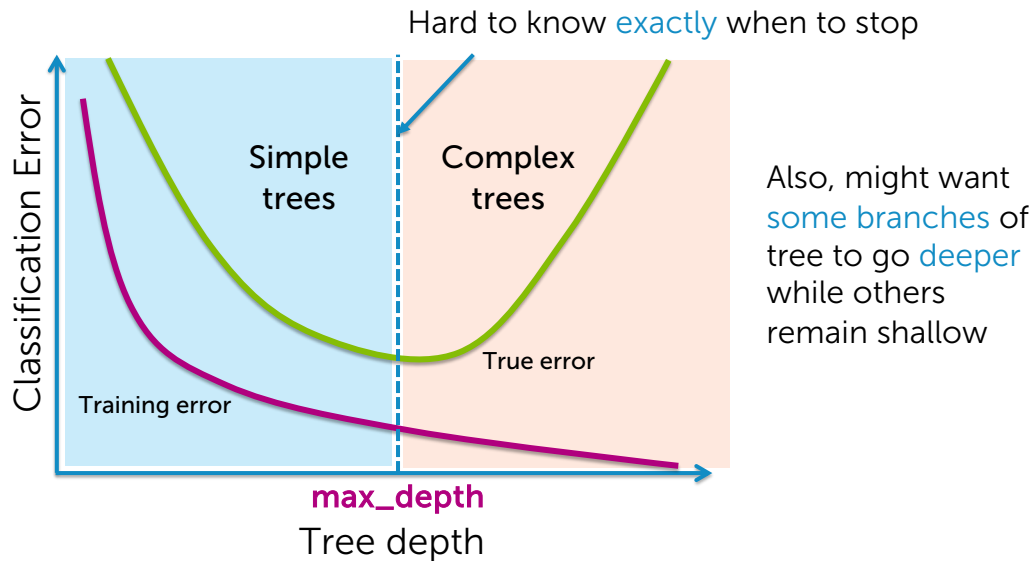
- **Stopping conditions (recap):**
 1. All examples have the same target value
 2. No more features to split on
- **Early stopping conditions:**
 1. Limit tree depth (choose *max_depth* using validation set)
 2. Do not consider splits that do not cause a sufficient decrease in classification error
 3. Do not split an intermediate node which contains too few data points

©2024 Emily Fox

CS 229: Machine Learning

78

Challenge with early stopping condition 1



©2024 Emily Fox

CS 229: Machine Learning

79

Early stopping condition 2: Pros and Cons

- **Pros:**
 - A reasonable heuristic for early stopping to avoid useless splits
- **Cons:**
 - **Too short sighted:** We may miss out on “good” splits may occur right after “useless” splits
 - Saw this with “xor” example

©2024 Emily Fox

CS 229: Machine Learning

80

Two approaches to picking simpler trees

1. Early Stopping:

Stop the learning algorithm **before** tree becomes too complex

2. Pruning:

Simplify the tree **after** the learning algorithm terminates

Complements early stopping

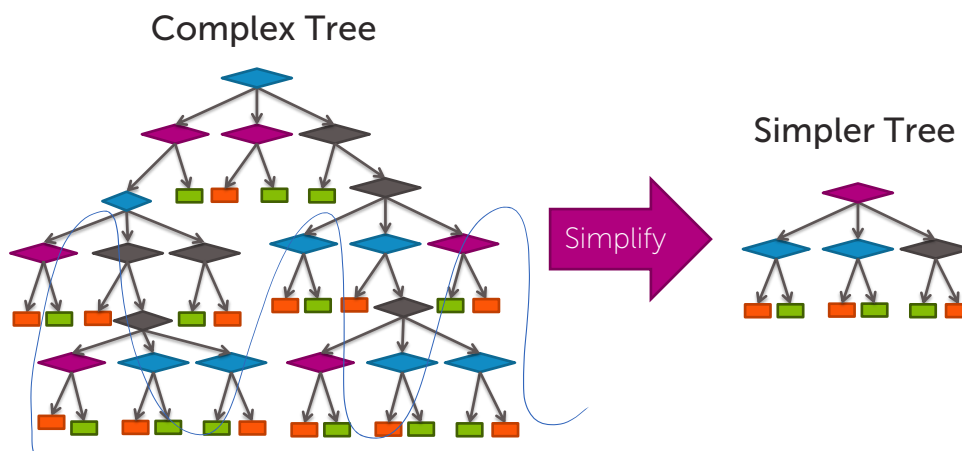
©2024 Emily Fox

CS 229: Machine Learning

81

Pruning: *Intuition*

Train a complex tree, simplify later

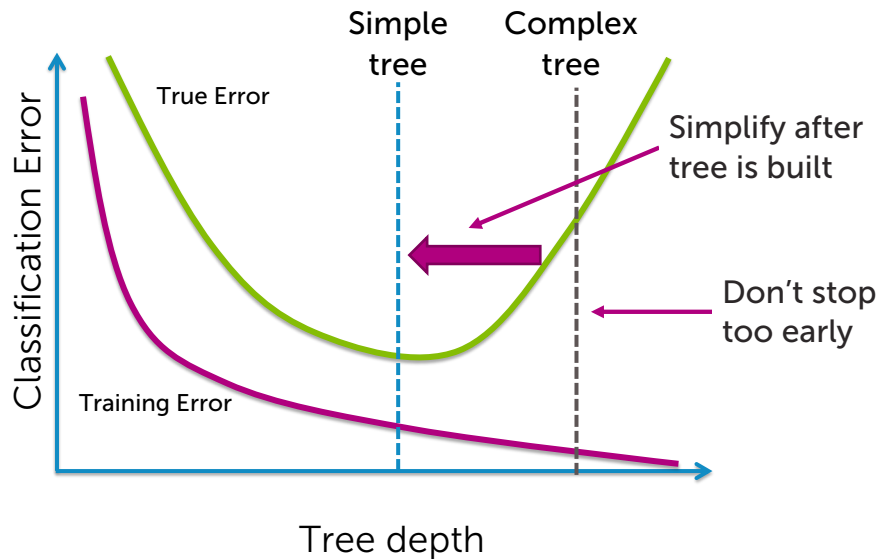


©2024 Emily Fox

CS 229: Machine Learning

82

Pruning motivation



©2024 Emily Fox

CS 229: Machine Learning

83

Scoring trees: Desired total quality format

Want to balance:

- i. How well tree fits data
- ii. Complexity of tree

Total cost =

$$\text{measure of fit} + \text{measure of complexity}$$

want to balance

↑ (classification error)
Large # = bad fit to training data

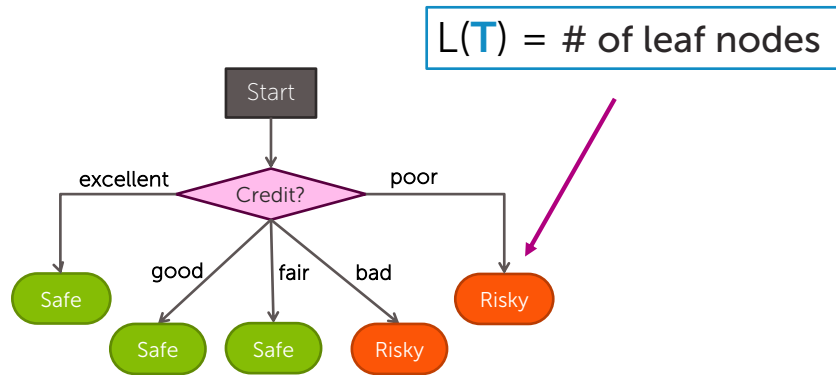
↑ Large # = likely to overfit

©2024 Emily Fox

CS 229: Machine Learning

84

Simple measure of complexity of tree



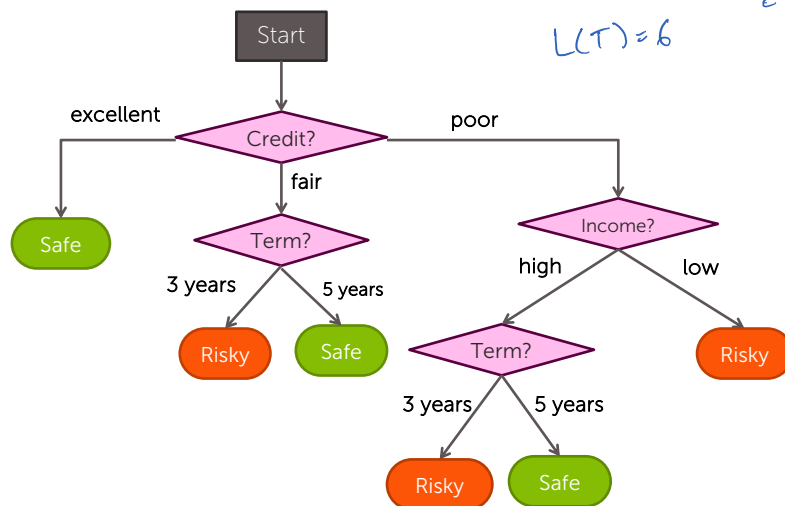
©2024 Emily Fox

CS 229: Machine Learning

85

Balance simplicity & predictive power

Too complex, risk of overfitting



©2024 Emily Fox

CS 229: Machine Learning

86

Too simple, high classification error




solution in between

Balancing fit and complexity

$$\text{Total cost } C(\mathbf{T}) = \text{Error}(\mathbf{T}) + \lambda L(\mathbf{T})$$

 tuning parameter

If $\lambda=0$: standard decision tree learning

If $\lambda=\infty$: no penalty \rightarrow  \hat{y} = majority vote
(of all training data)

If λ in between: balance of fit + complexity

©2024 Emily Fox

CS 229: Machine Learning

87

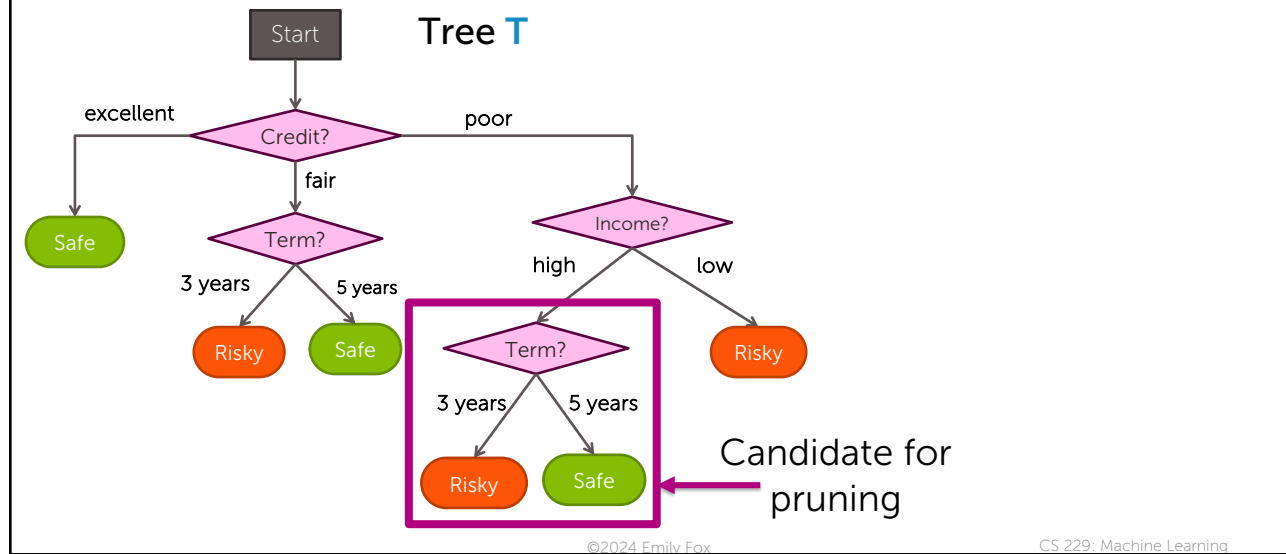
Tree pruning algorithm

©2024 Emily Fox

CS 229: Machine Learning

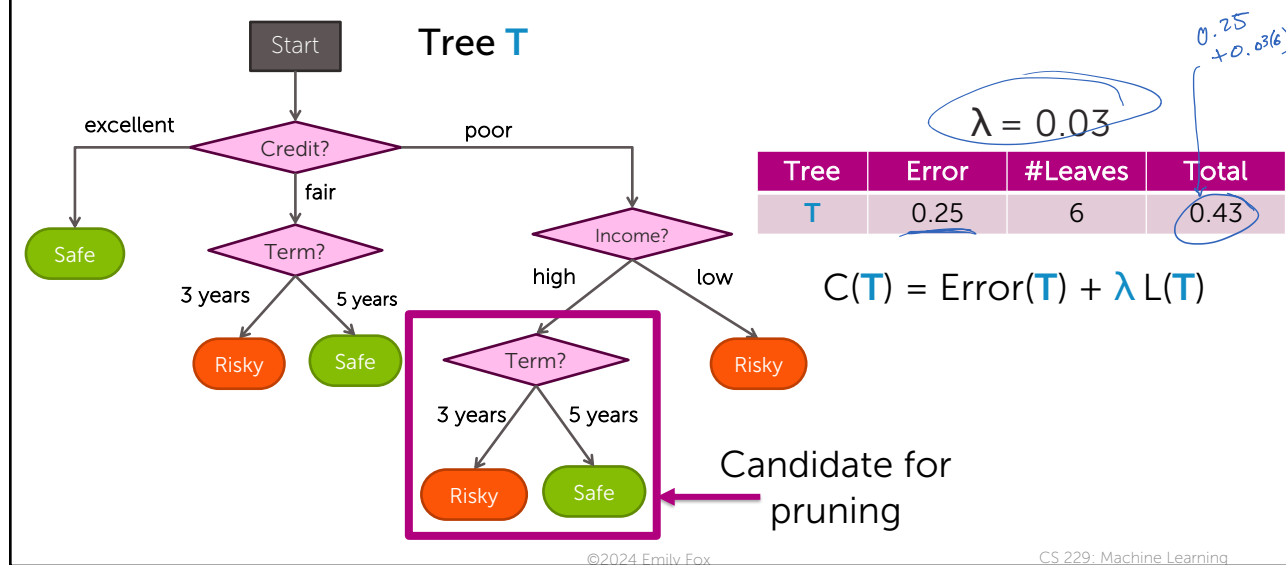
88

Step 1: Consider a split



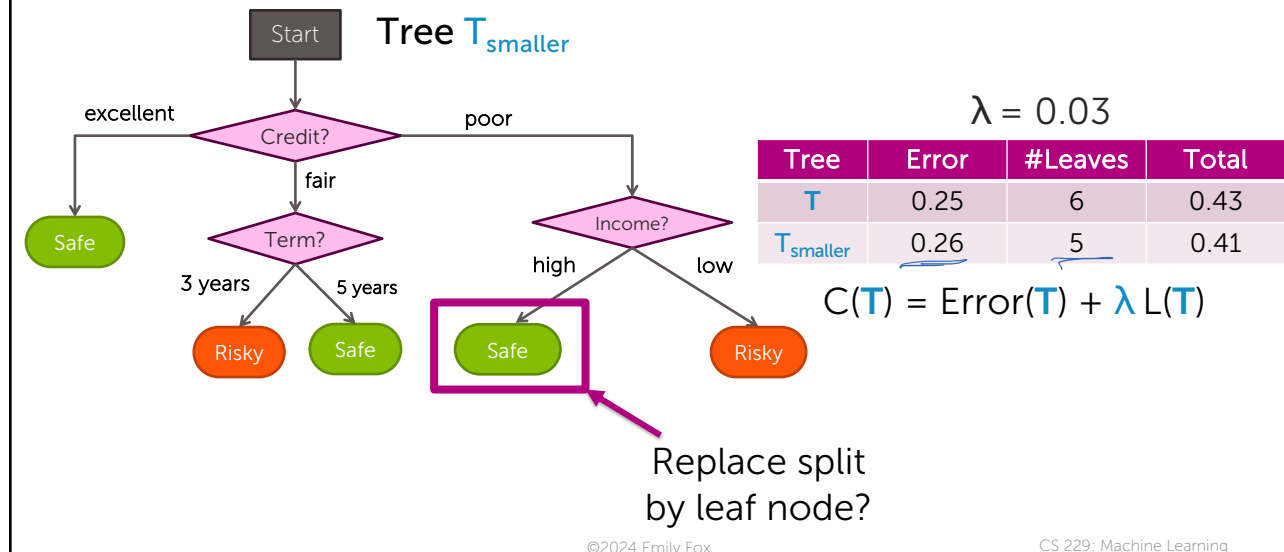
89

Step 2: Compute total cost $C(T)$ of split



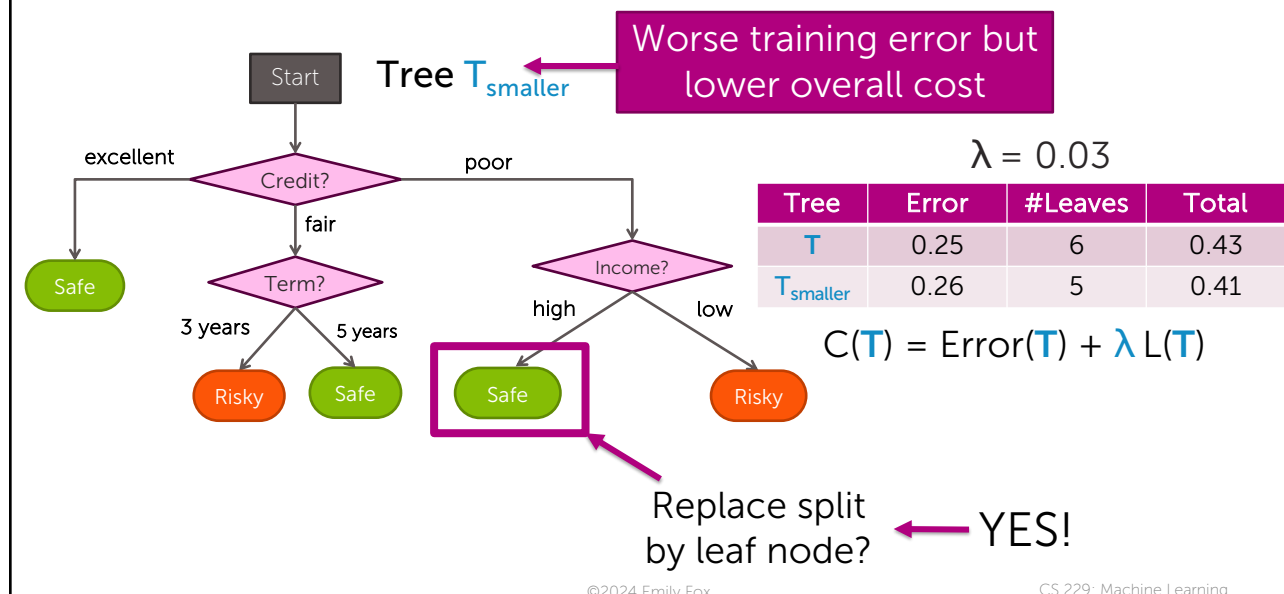
90

Step 2: "Undo" the splits on T_{smaller}



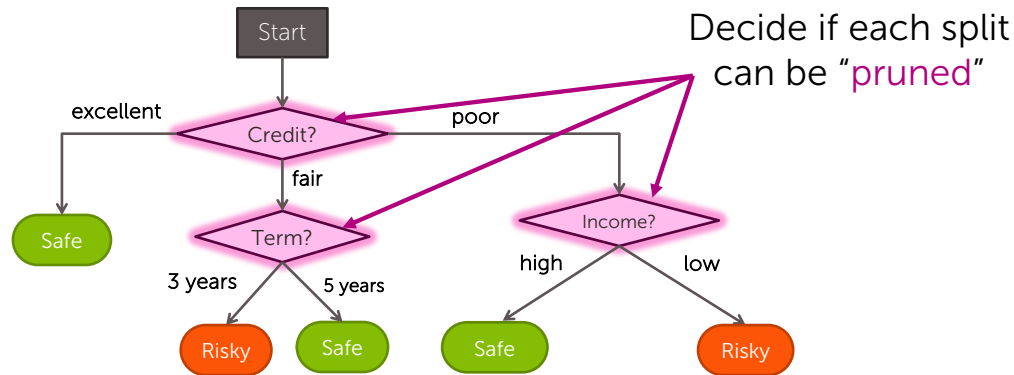
91

Prune if total cost is lower: $C(T_{\text{smaller}}) \leq C(T)$



92

Step 5: Repeat Steps 1-4 for every split



©2024 Emily Fox

CS 229: Machine Learning

93

How to choose hyperparameters?

(e.g., λ or max_depth) *← tuning params*

*validation error on valid. set
or
Cross validations*

©2024 Emily Fox

CS 229: Machine Learning

94

Summary of overfitting in decision trees

©2024 Emily Fox

CS 229: Machine Learning

95

What you can do now...

- Identify when overfitting in decision trees
- Prevent overfitting with early stopping
 - Limit tree depth
 - Do not consider splits that do not reduce classification error
 - Do not split intermediate nodes with only few points
- Prevent overfitting by pruning complex trees
 - Use a total cost formula that balances classification error and tree complexity
 - Use total cost to merge potentially complex trees into simpler ones

©2024 Emily Fox

CS 229: Machine Learning

96