



Linear classifiers:

Handling overfitting, categorical inputs, & multiple classes

STAT/CSE 416: Machine Learning
Emily Fox
University of Washington
April 24, 2018

©2018 Emily Fox

Encoding categorical inputs

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Categorical inputs

- Numeric inputs:
 - #awesome, age, salary,...
 - Intuitive when multiplied by coefficient
 - e.g., 1.5 #awesome
- Categorical inputs:



Gender
(Male, Female,...)



Country of birth
(Argentina, Brazil, USA,...)



Zipcode
(10005, 98195,...)

Numeric value, but should be interpreted as category
(98195 not about 9x larger than 10005)

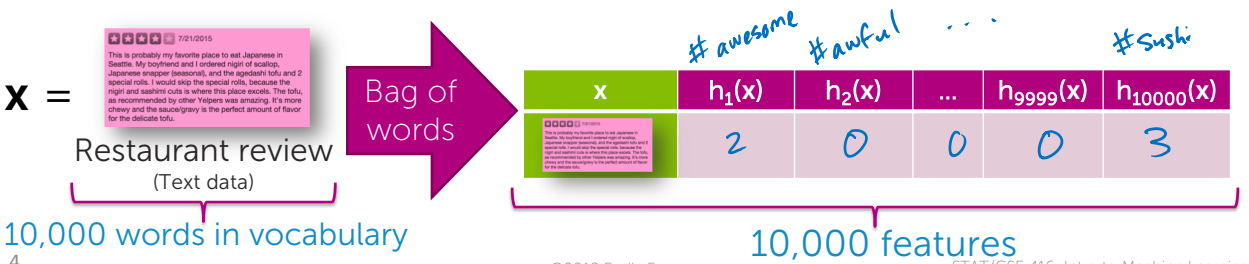
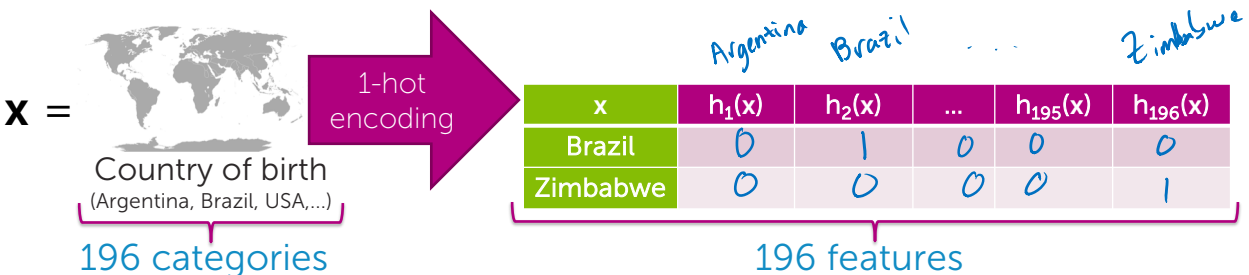
How do we multiply category by coefficient???
Must convert categorical inputs into numeric features

3

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Encoding categories as numeric features



4

©2018 Emily Fox

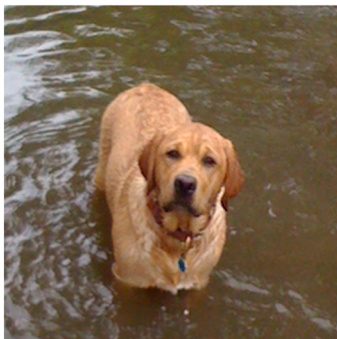
STAT/CSE 416: Intro to Machine Learning

Multiclass classification using 1 versus all

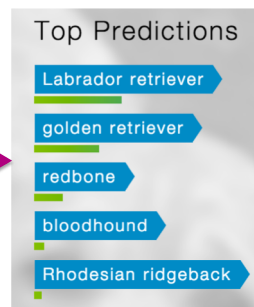
©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Multiclass classification



Input: x
Image pixels



Output: y
Object in image

6

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Multiclass classification formulation

- C possible classes:
 - y can be 1, 2, ..., C
- N datapoints:

before $y \in \{-1, 1\}$

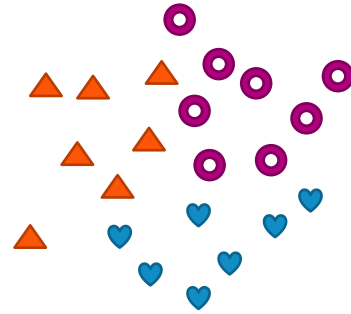
Data point	x[1]	x[2]	y
\mathbf{x}_1, y_1	2	1	▲
\mathbf{x}_2, y_2	0	2	♥
\mathbf{x}_3, y_3	3	3	◯
\mathbf{x}_4, y_4	4	1	◯

Learn:

$$\hat{P}(y = \blacktriangle | \mathbf{x})$$

$$\hat{P}(y = \heartsuit | \mathbf{x})$$

$$\hat{P}(y = \bigcirc | \mathbf{x})$$



7

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

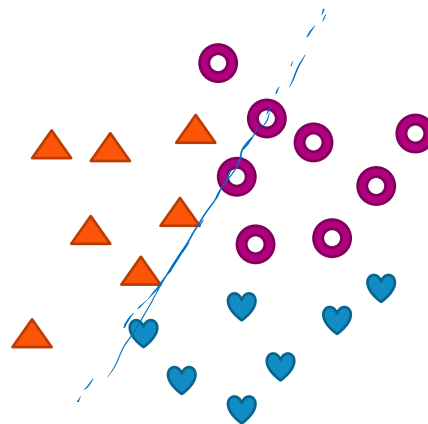
1 versus all:

Estimate $\hat{P}(y = \blacktriangle | \mathbf{x})$ using 2-class model

+1 class: points with $y_i = \blacktriangle$
 -1 class: points with $y_i = \heartsuit$ OR \bigcirc

Train classifier: $\hat{P}_{\blacktriangle}(y = +1 | \mathbf{x})$

Predict: $\hat{P}(y = \blacktriangle | \mathbf{x}_i) = \hat{P}_{\blacktriangle}(y = +1 | \mathbf{x}_i)$



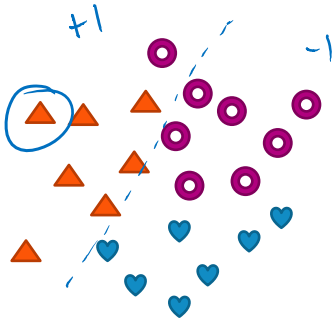
8

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

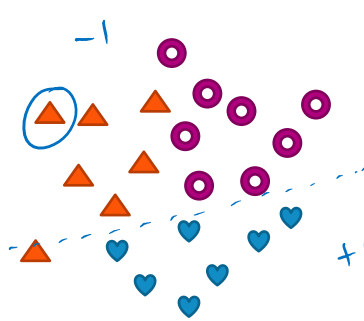
1 versus all: simple multiclass classification using C 2-class models

$$\hat{P}(y=\triangle | \mathbf{x}_i) = \hat{P}_{\triangle}(y=+1 | \mathbf{x}_i, \mathbf{w}_{\triangle})$$



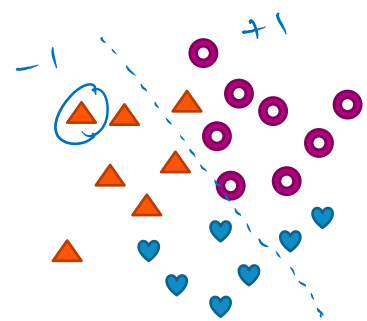
9

$$\hat{P}(y=\heartsuit | \mathbf{x}_i) = \hat{P}_{\heartsuit}(y=+1 | \mathbf{x}_i, \mathbf{w}_{\heartsuit})$$



©2018 Emily Fox

$$\hat{P}(y=\circ | \mathbf{x}_i) = \hat{P}_{\circ}(y=+1 | \mathbf{x}_i, \mathbf{w}_{\circ})$$



STAT/CSE 416: Intro to Machine Learning

Multiclass training

$\hat{P}_c(y=+1 | \mathbf{x})$ = estimate of
1 vs all model for each class

Predict most likely class

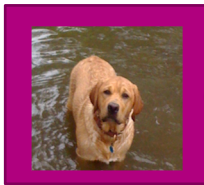
max_prob = 0; $\hat{y} = 0$

For $c = 1, \dots, C$:

If $\hat{P}_c(y=+1 | \mathbf{x}_i) > \text{max_prob}$:

$\hat{y} = c$

max_prob = $\hat{P}_c(y=+1 | \mathbf{x}_i)$



Input: \mathbf{x}_i

10

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Summary of overfitting in logistic regression, categorical inputs, and multiclass classification

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

What you can do now...

- Describe symptoms and effects of overfitting in classification
 - Identify when overfitting is happening
 - Relate large learned coefficients to overfitting
 - Describe the impact of overfitting on decision boundaries and predicted probabilities of linear classifiers
- Use regularization to mitigate overfitting
 - Motivate the form of L2 regularized logistic regression quality metric
 - Describe the use of L1 regularization to obtain sparse logistic regression solutions
 - Describe what happens to estimated coefficients as tuning parameter λ is varied
 - Interpret coefficient path plot
- Use 1-hot encoding to represent categorical inputs
- Perform multiclass classification using the 1-versus-all approach

12

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning



Decision Trees

STAT/CSE 416: Machine Learning
Emily Fox
University of Washington
April 24, 2018

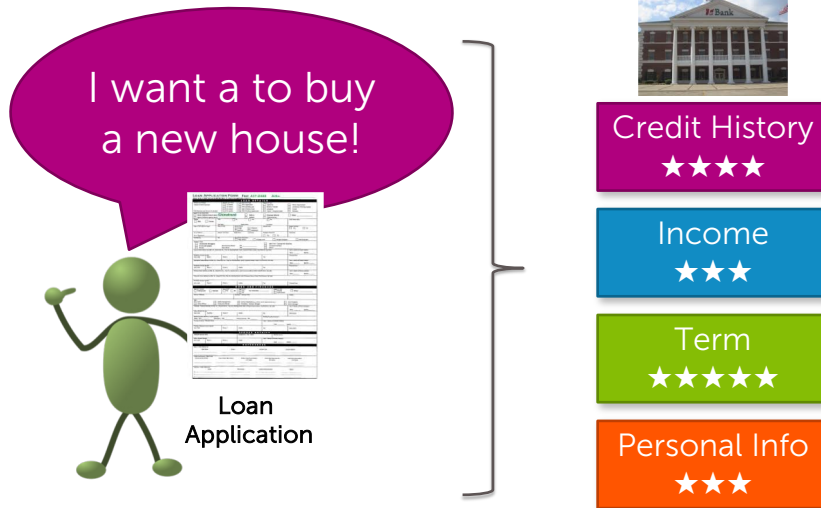
©2018 Emily Fox

Predicting potential loan defaults

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

What makes a loan risky?



15

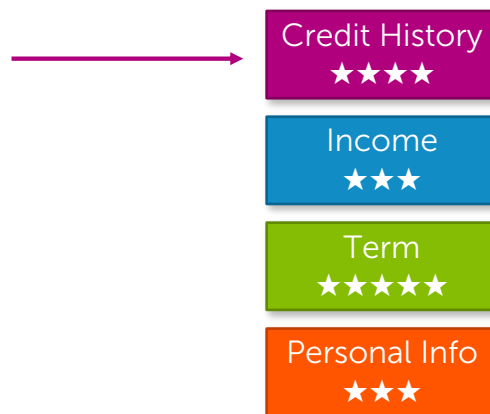
©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Credit history explained

Did I pay previous loans on time?

Example:
excellent, good, or fair



16

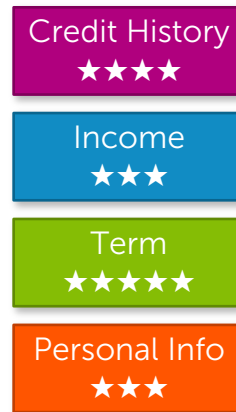
©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Income

What's my income?

Example: \$80K per year



17

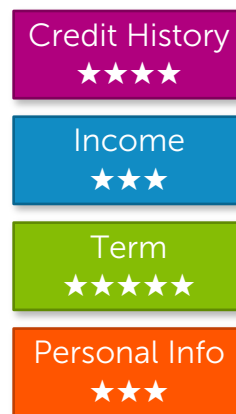
©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Loan terms

How soon do I need to pay the loan?

Example: 3 years, 5 years,...



18

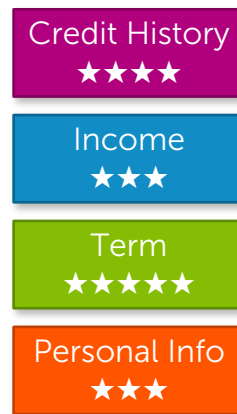
©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Personal information

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

Example: Home loan for a married couple



19

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Intelligent application

Loan Applications

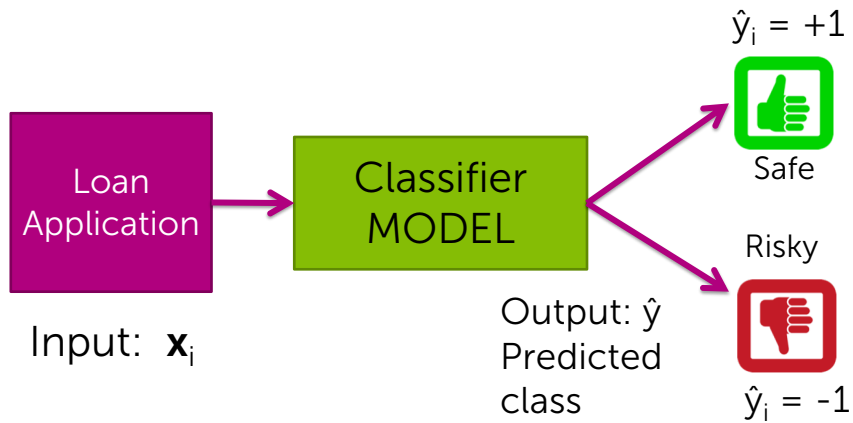


20

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Classifier review

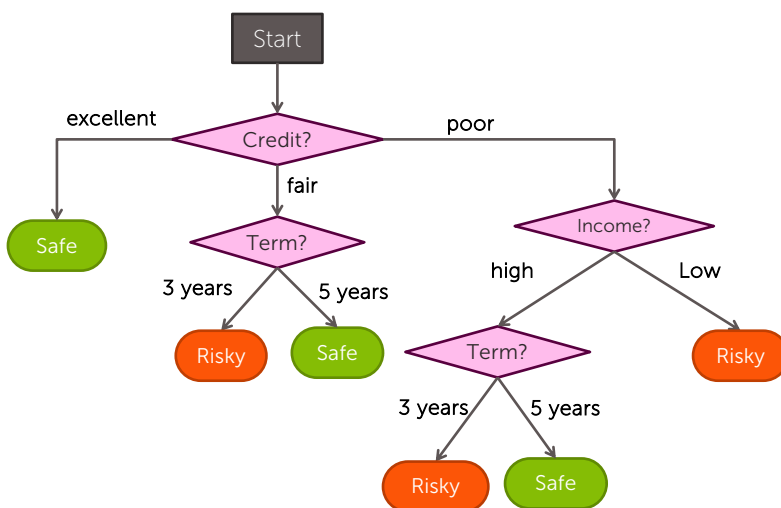


21

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

This module ... decision trees

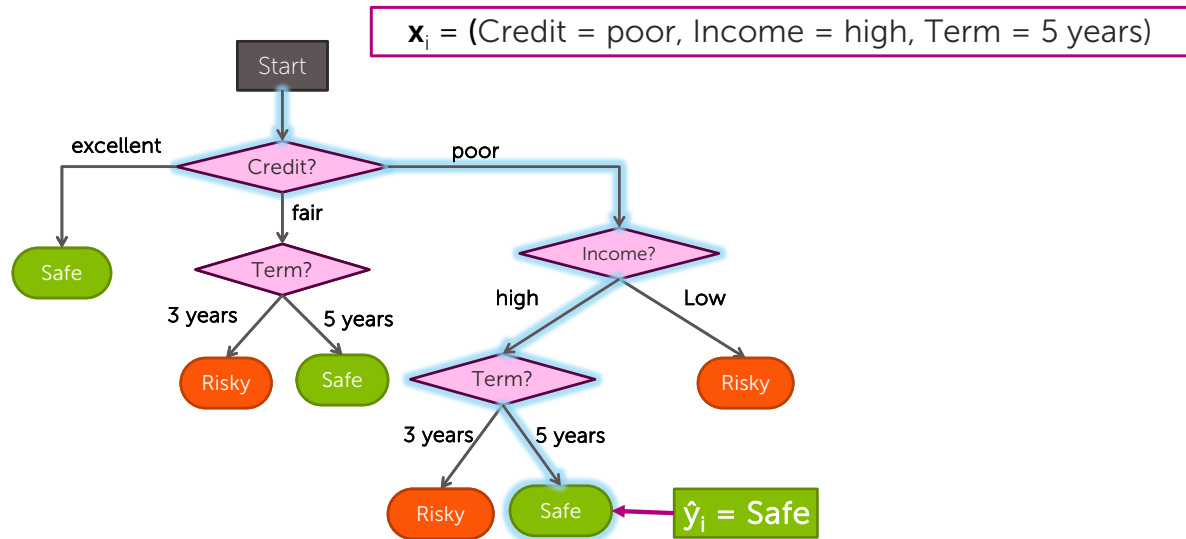


22

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Scoring a loan application



23

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Decision tree learning task

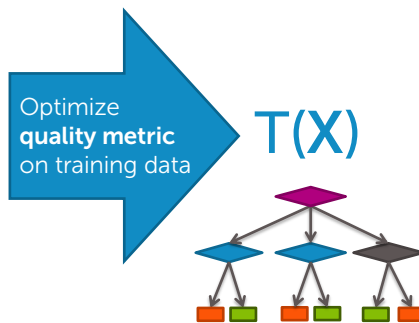
©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

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



25

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Quality metric: 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

26

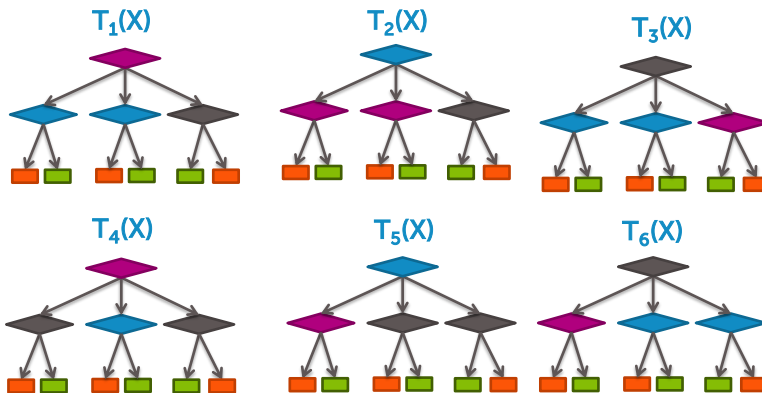
©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

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]



27

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Greedy decision tree learning

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

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

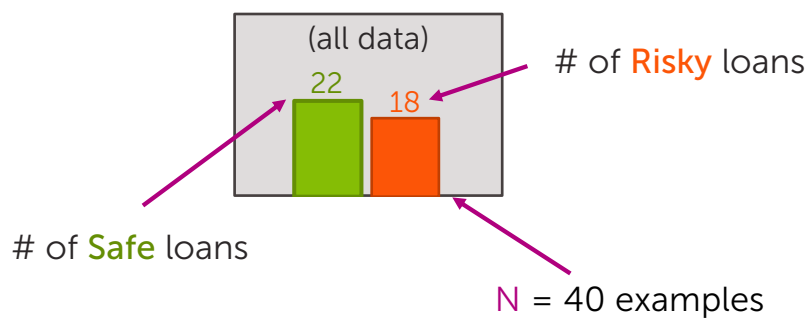
29

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Start with all the data

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



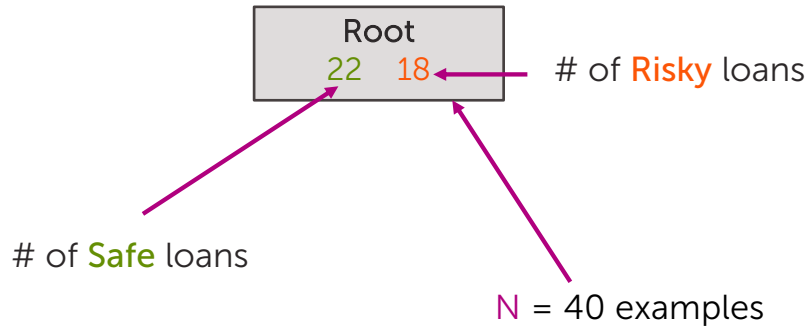
30

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Compact visual notation: Root node

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



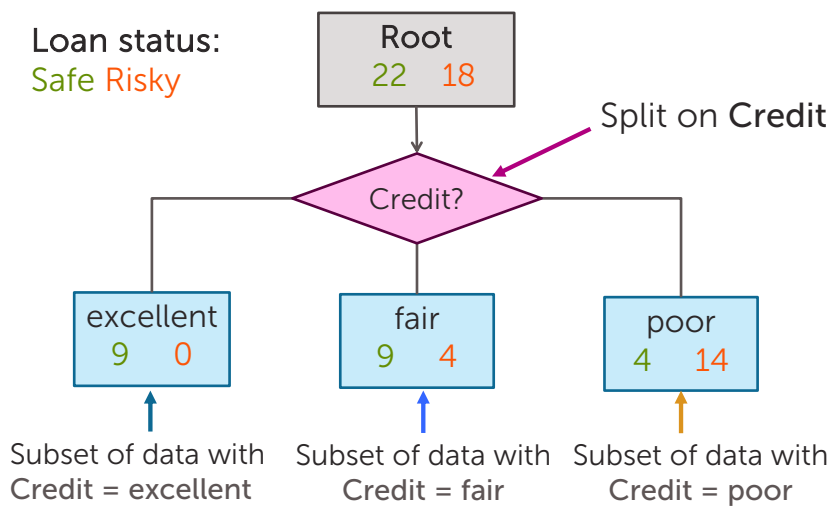
31

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Decision stump: Single level tree

Loan status:
Safe **Risky**



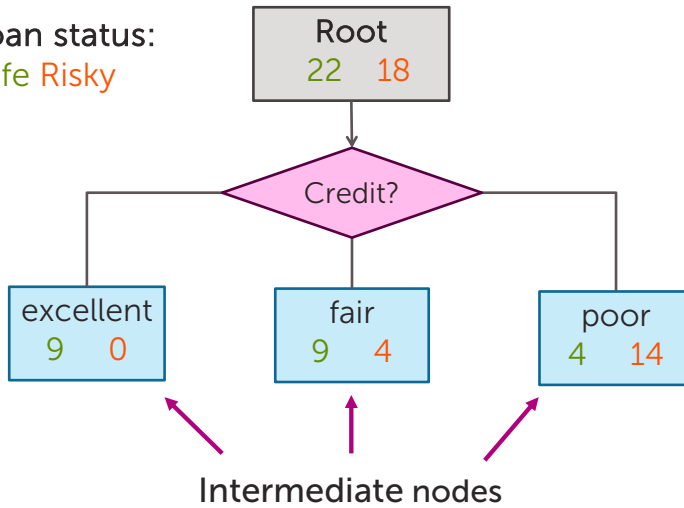
32

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Visual notation: Intermediate nodes

Loan status:
Safe Risky



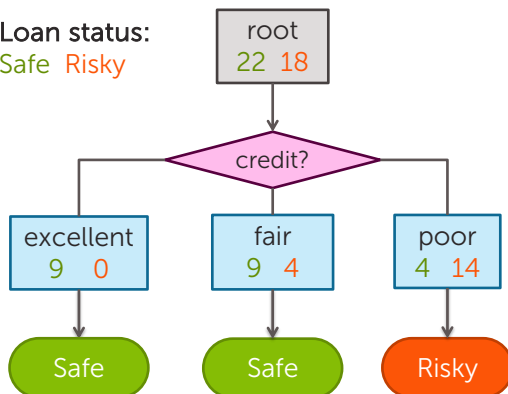
33

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Making predictions with a decision stump

Loan status:
Safe Risky



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

34

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

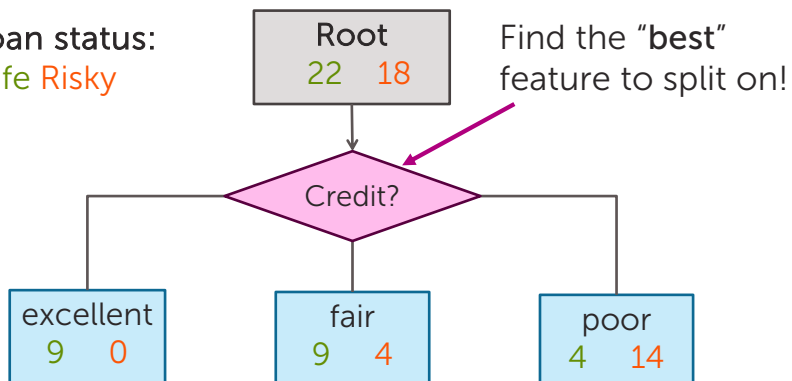
Selecting best feature to split on

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

How do we learn a decision stump?

Loan status:
Safe Risky



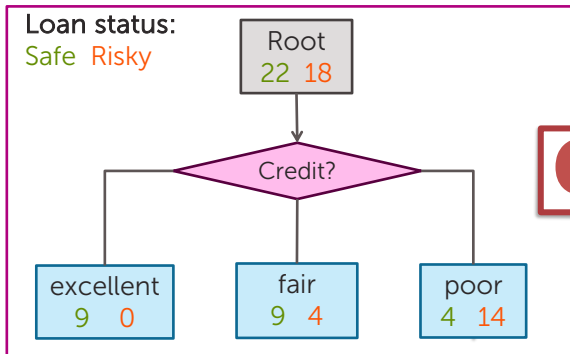
36

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

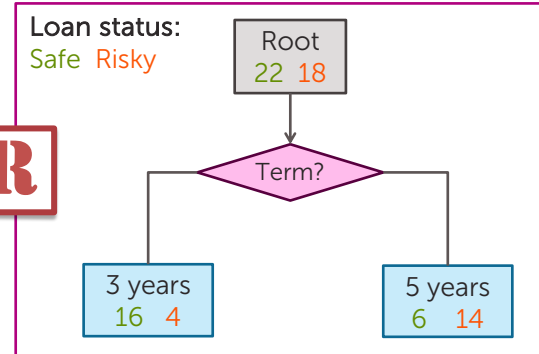
How do we select the best feature?

Choice 1: Split on Credit



OR

Choice 2: Split on Term



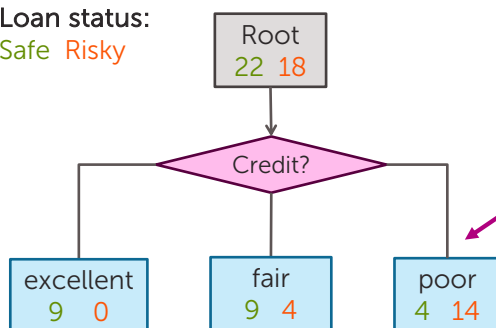
37

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

How do we measure effectiveness of a split?

Loan status:
Safe Risky



Idea: Calculate classification error of this decision stump

$$\text{Error} = \frac{\# \text{ mistakes}}{\# \text{ data points}}$$

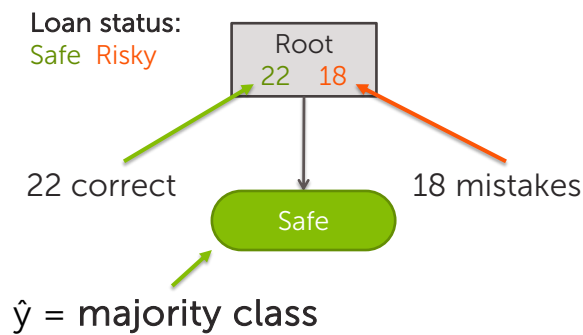
38

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

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



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

$$= 0.45$$

Tree	Classification error
(root)	0.45

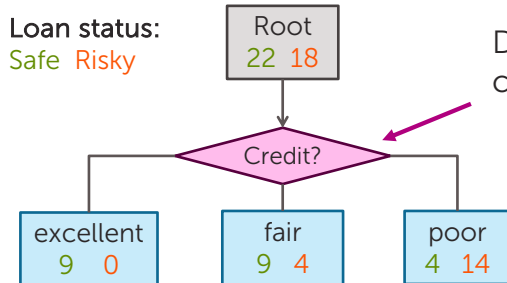
39

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Choice 1: Split on Credit history?

Choice 1: Split on Credit



Does a **split on Credit** reduce classification error below 0.45?

40

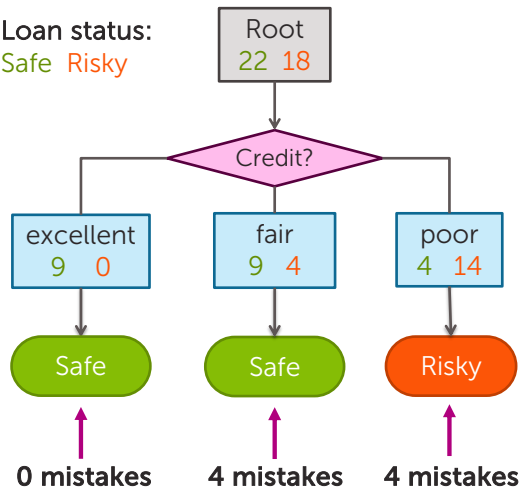
©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Split on Credit: Classification error

Choice 1: Split on Credit

Loan status:
Safe Risky



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

Tree	Classification error
(root)	0.45
Split on credit	0.2

41

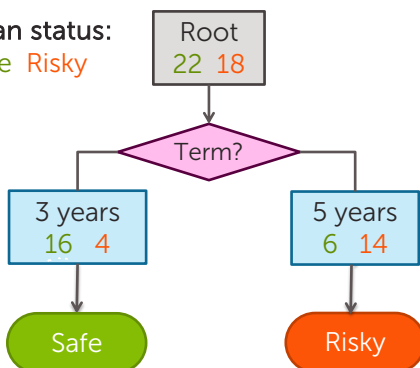
©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Choice 2: Split on Term?

Choice 2: Split on Term

Loan status:
Safe Risky



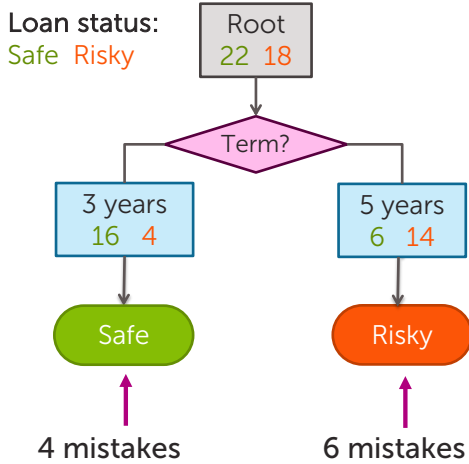
42

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

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

43

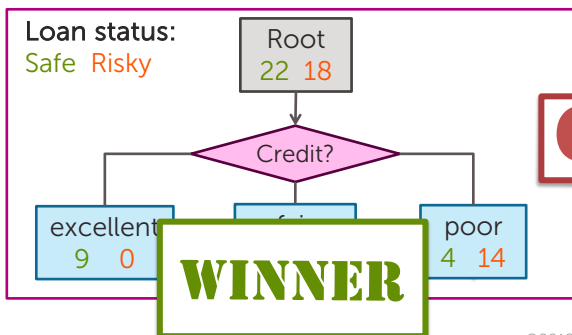
©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

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

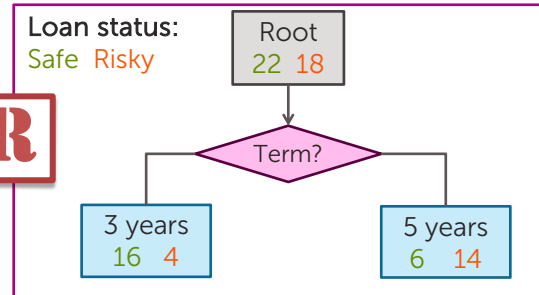
Tree	Classification error
(root)	0.45
split on credit	0.2
split on loan term	0.25

Choice 1: Split on Credit



OR

Choice 2: Split on Term



44

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Feature split selection algorithm

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

45

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

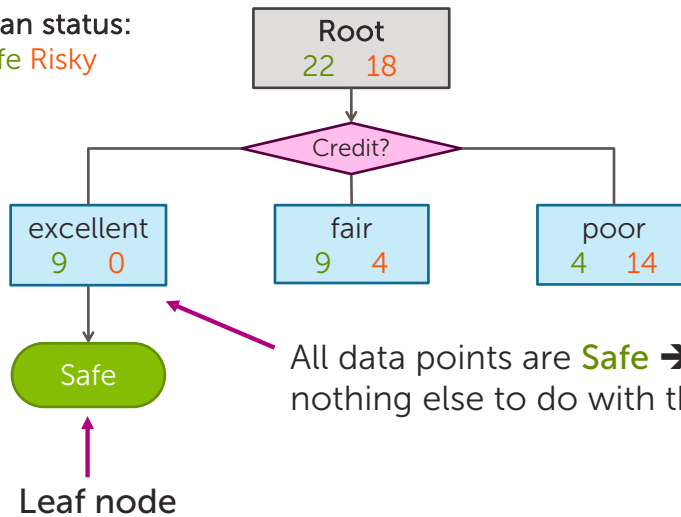
Recursion & Stopping conditions

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

We've learned a decision stump, what next?

Loan status:
Safe Risky



All data points are **Safe** →
nothing else to do with this subset of data

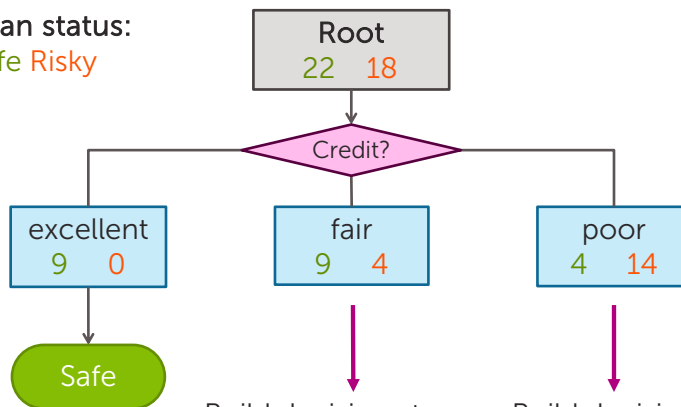
47

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Tree learning = Recursive stump learning

Loan status:
Safe Risky



Build decision stump
with subset of data
where Credit = fair

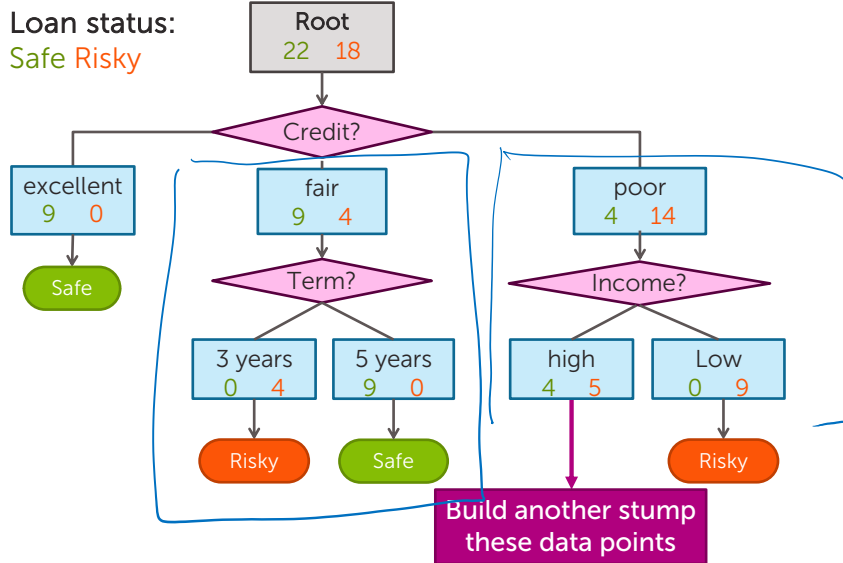
Build decision stump
with subset of data
where Credit = poor

48

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Second level

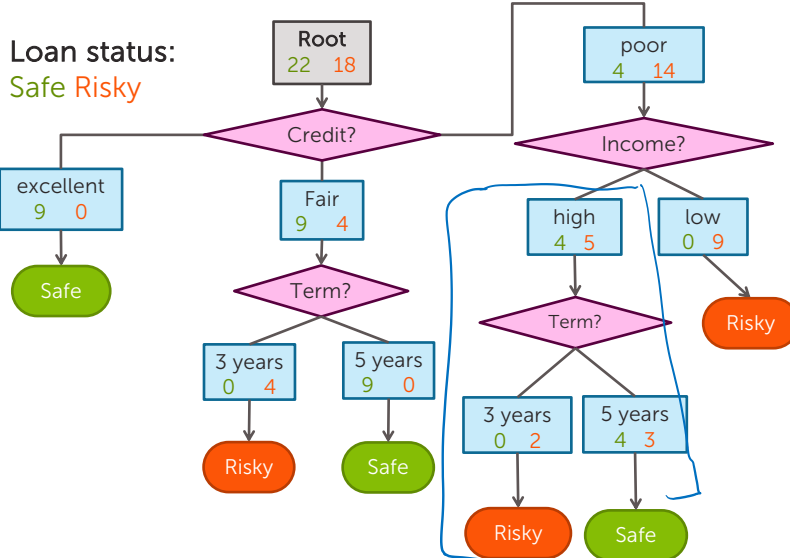


49

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Final decision tree



50

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

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???

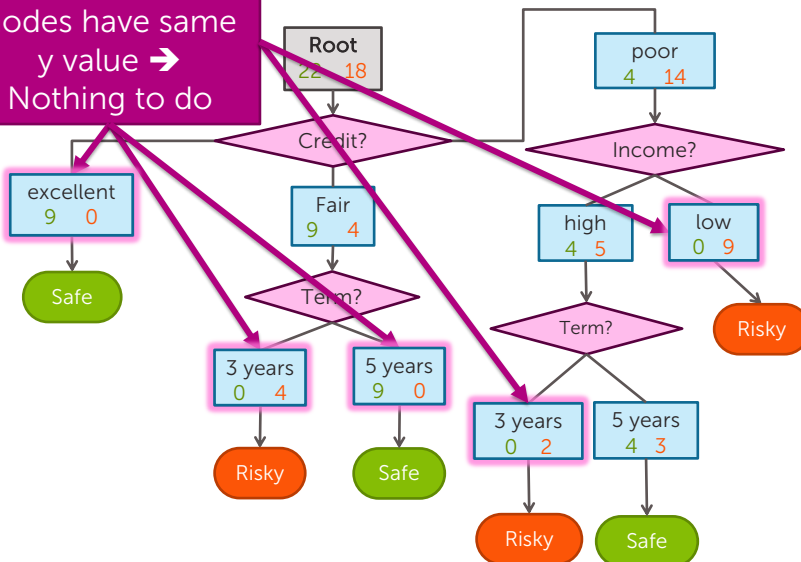
51

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Stopping condition 1: All data agrees on y

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



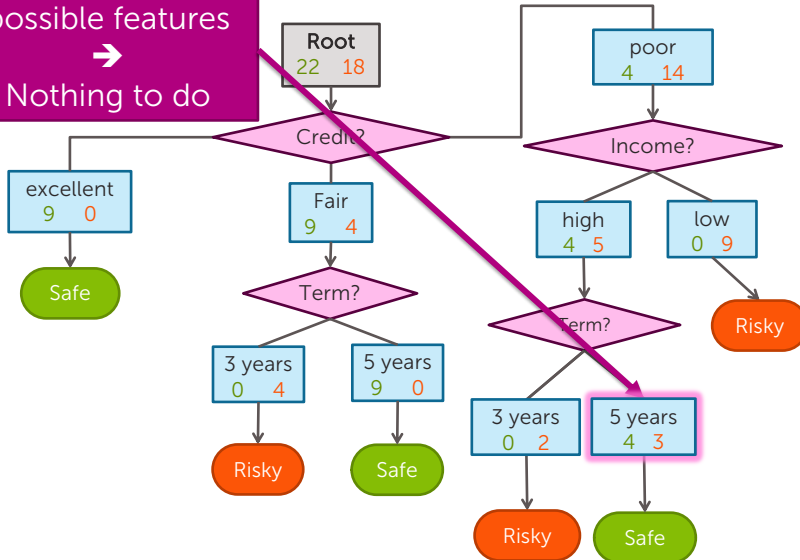
52

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Stopping condition 2: Already split on all features

Already split on all possible features
 →
 Nothing to do



53

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

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

54

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Is this a good idea?

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

55

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Stopping condition 3:

Don't 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

y values
True False

Root
2 2

↓
↑
y = safe

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

$$= 0.5$$

Tree	Classification error
(root)	0.5

56

©2018 Emily Fox

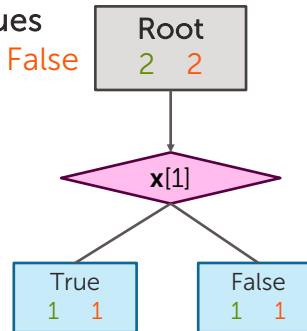
STAT/CSE 416: Intro to Machine Learning

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

57

©2018 Emily Fox

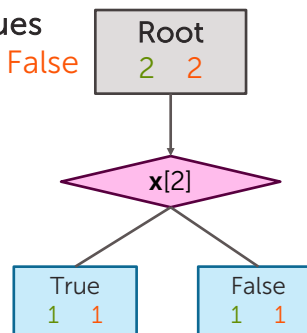
STAT/CSE 416: Intro to Machine Learning

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

58

©2018 Emily Fox

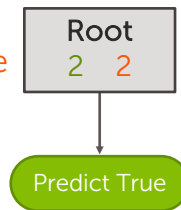
STAT/CSE 416: Intro to Machine Learning

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

59

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

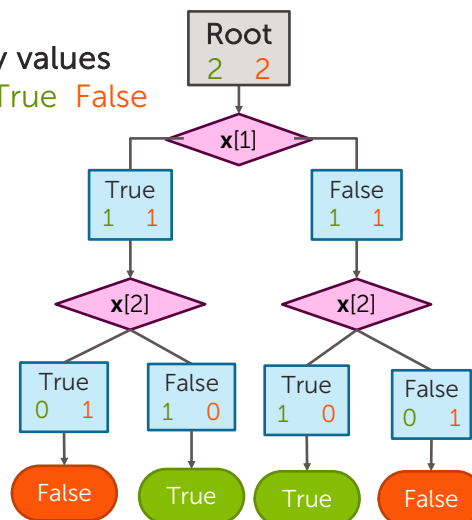
Without stopping condition 3

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

y values
True False



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

60

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Decision tree learning: *Real valued features*

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

How do we use real values inputs?

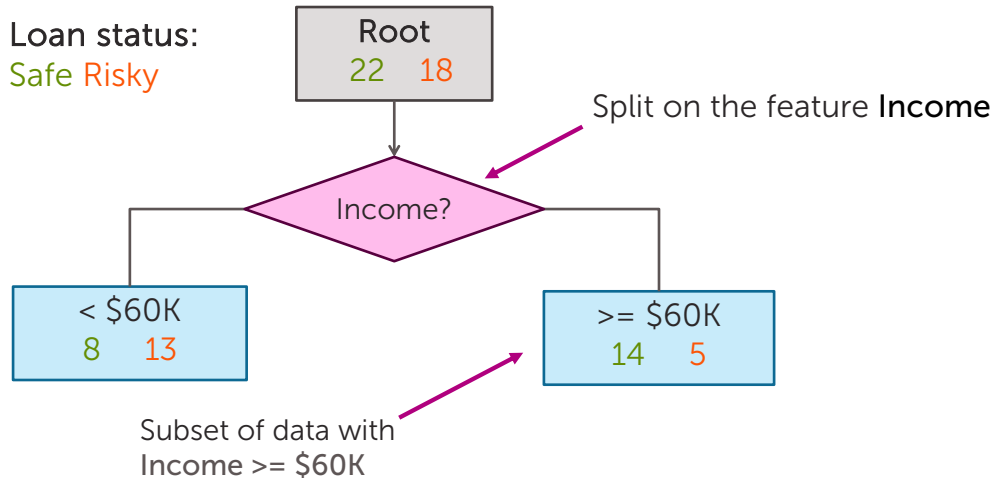
Income	Credit	Term	y
\$105 K	excellent	3 yrs	Safe
\$112 K	good	5 yrs	Risky
\$73 K	fair	3 yrs	Safe
\$69 K	excellent	5 yrs	Safe
\$217 K	excellent	3 yrs	Risky
\$120 K	good	5 yrs	Safe
\$64 K	fair	3 yrs	Risky
\$340 K	excellent	5 yrs	Safe
\$60 K	good	3 yrs	Risky

62

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Threshold split



63

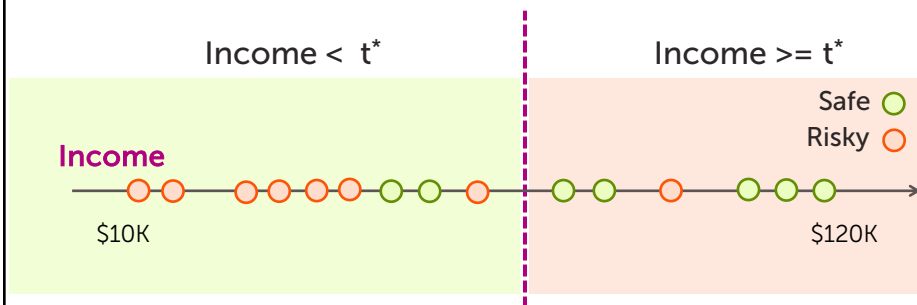
©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Finding the best threshold split

Infinite possible
values of t

Income = t^* *threshold to choose*



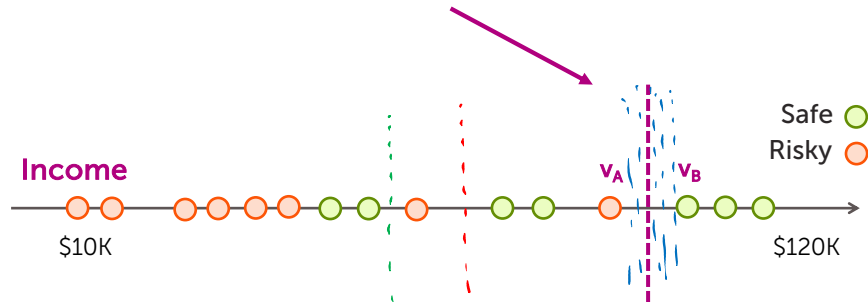
64

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Consider a threshold between points

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



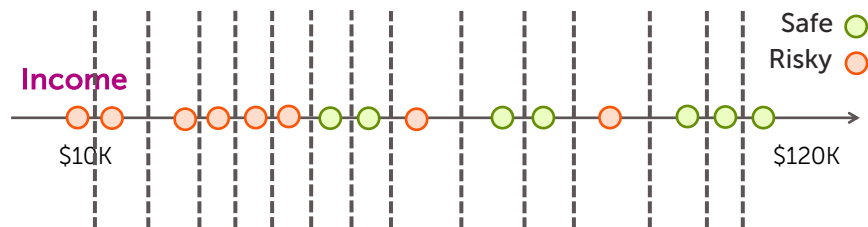
65

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Only need to consider mid-points

Finite number of splits to consider



66

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

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 $t_i = (v_i + v_{i+1}) / 2$
 - Compute classification error for threshold split $h_j(\mathbf{x}) \geq t_i$
 - Chose the t^* with the lowest classification error

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

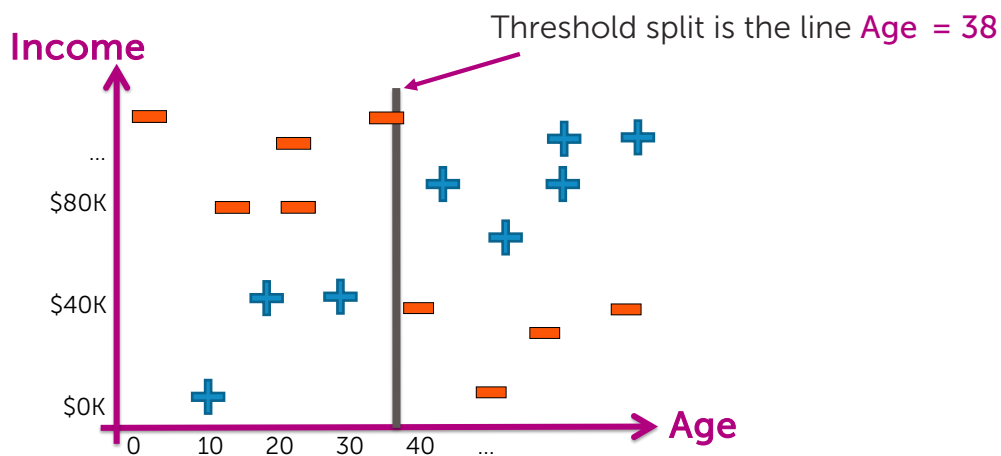
decision
stump

67

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Visualizing the threshold split

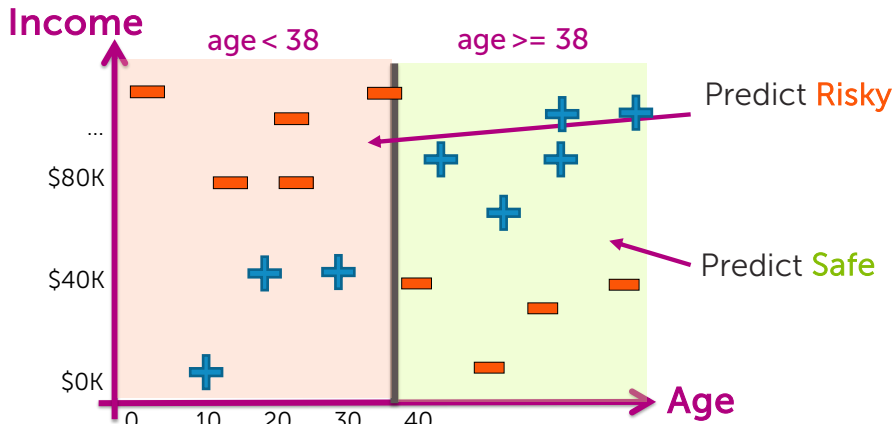


68

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Split on Age ≥ 38

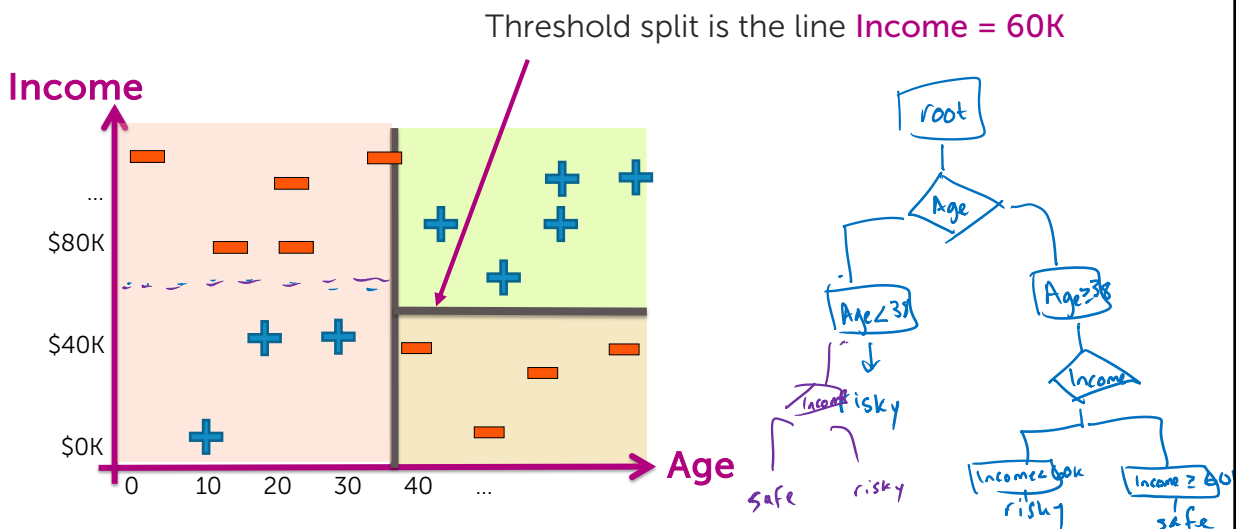


69

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

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

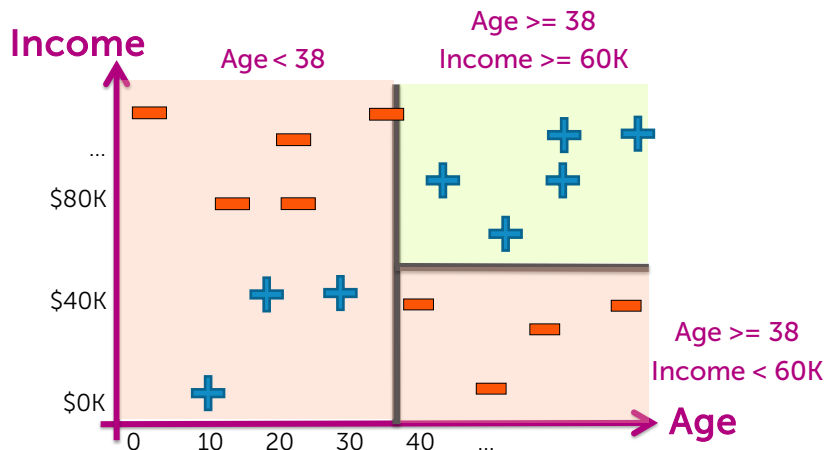


70

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Each split partitions the 2-D space



71

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

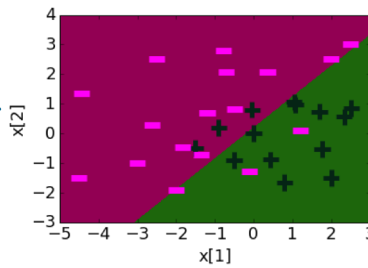
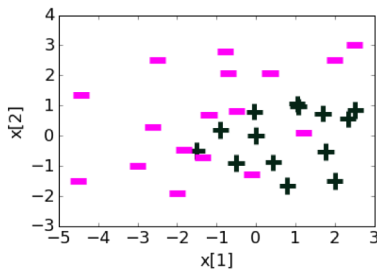
Decision trees vs logistic regression: Example

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Logistic regression

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

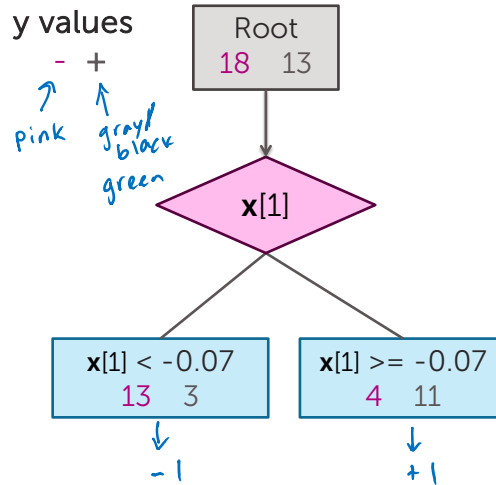
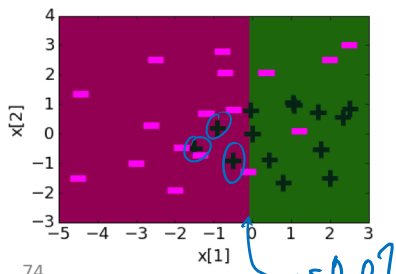
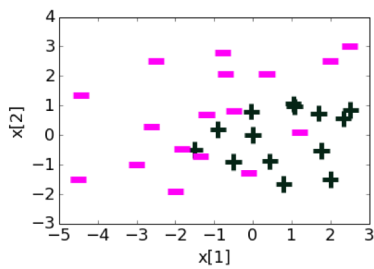


73

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Depth 1: Split on $x[1]$

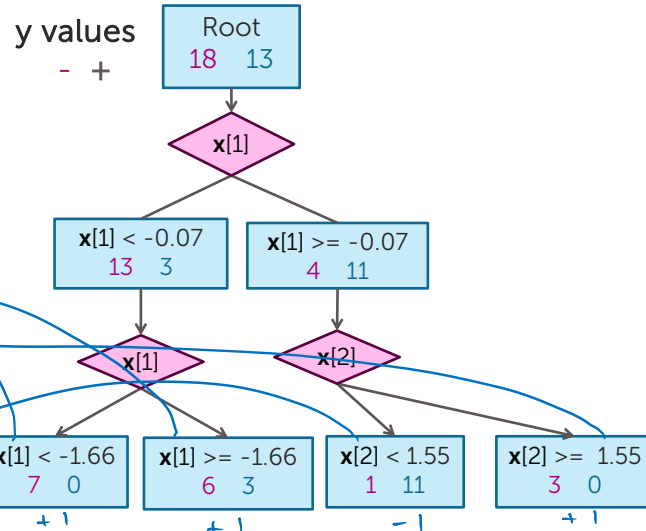
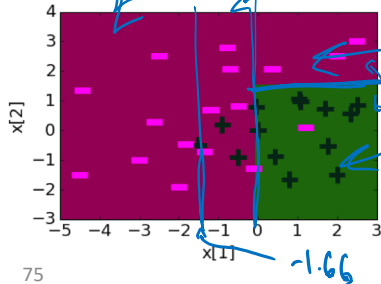
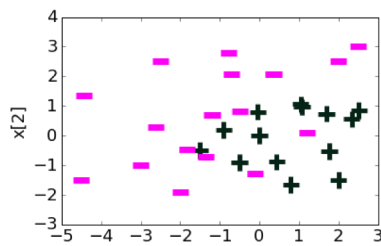


74

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Depth 2



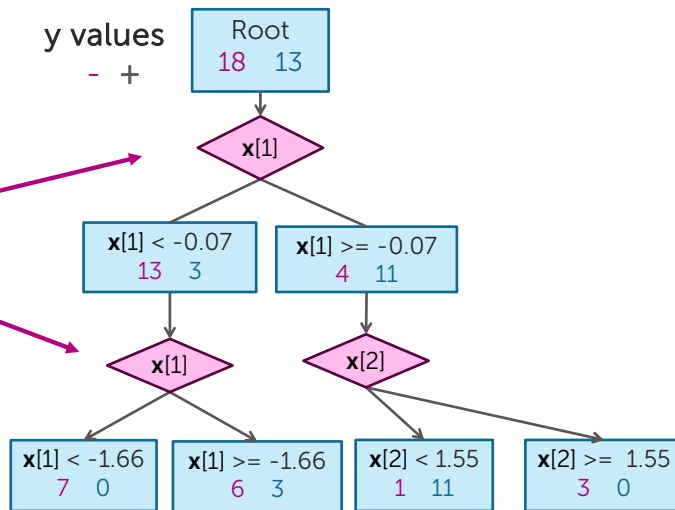
75

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Threshold split caveat

For threshold splits, same feature can be used multiple times

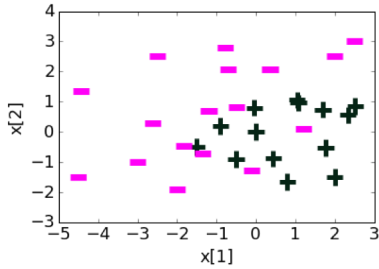


76

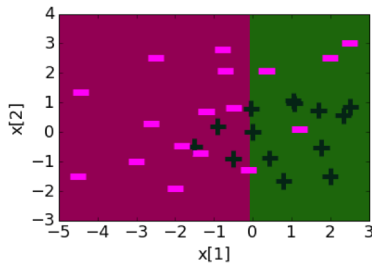
©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

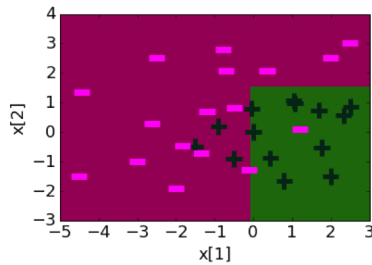
Decision boundaries



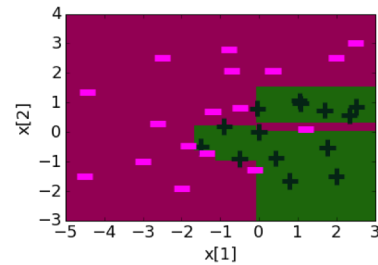
Depth 1



Depth 2



Depth 10



77

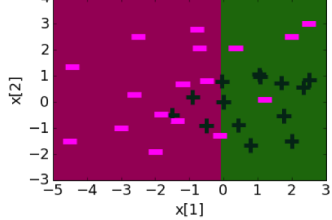
©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

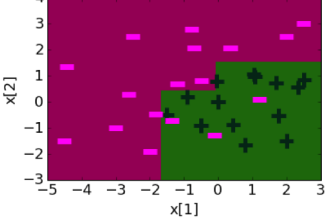
Comparing decision boundaries

Decision Tree

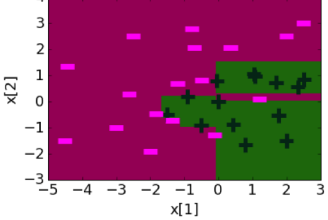
Depth 1



Depth 3

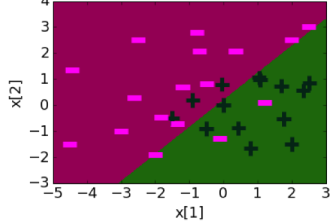


Depth 10

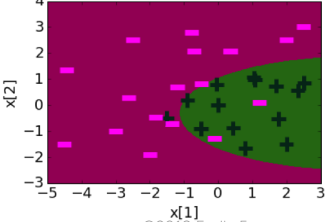


Logistic Regression

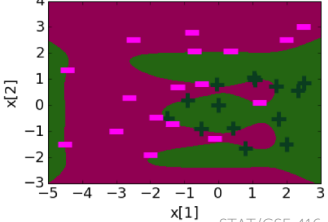
Degree 1 features



Degree 2 features



Degree 6 features



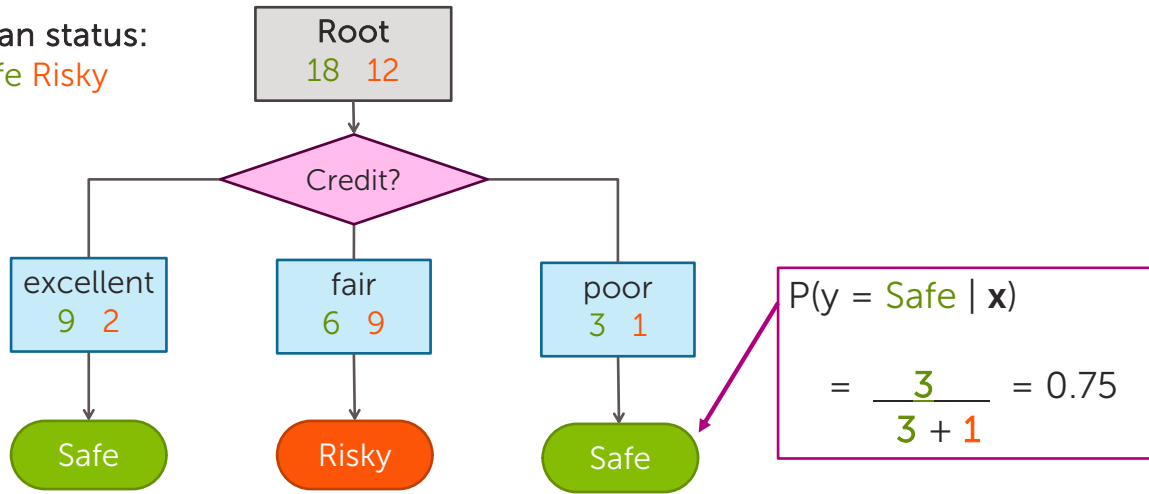
78

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Predicting probabilities with decision trees

Loan status:
Safe Risky

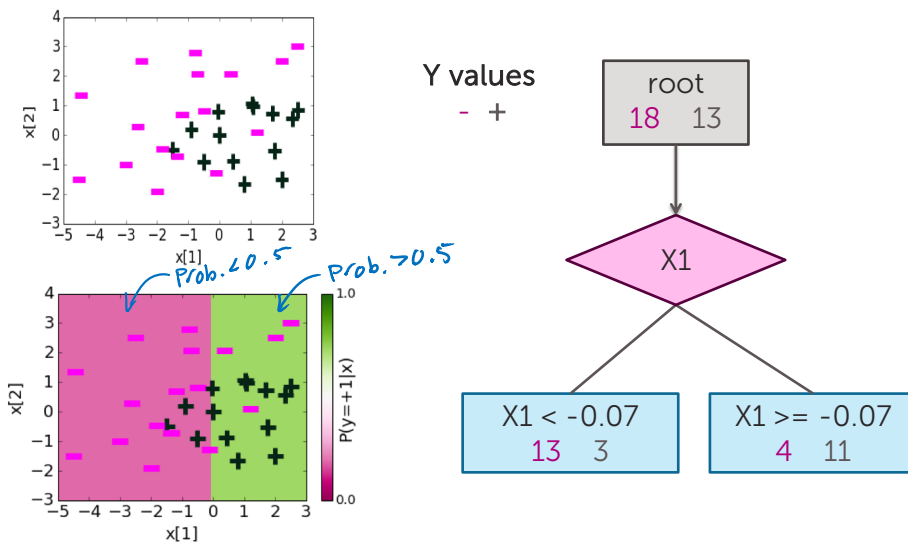


79

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Depth 1 probabilities

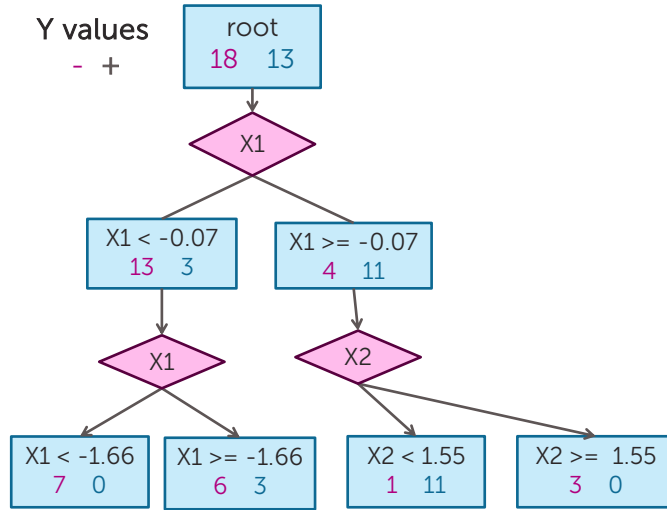
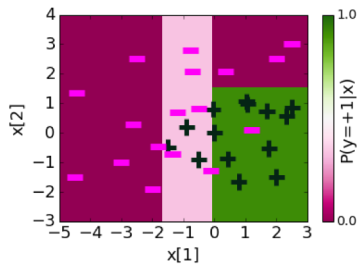
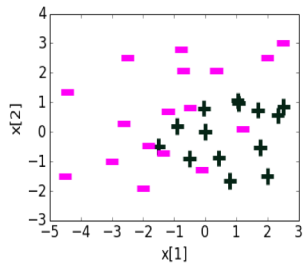


80

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Depth 2 probabilities



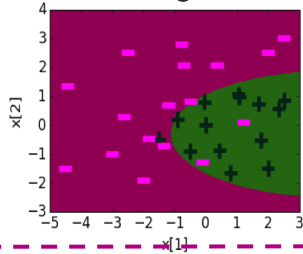
81

©2018 Emily Fox

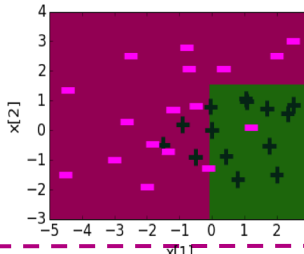
STAT/CSE 416: Intro to Machine Learning

Comparison with logistic regression

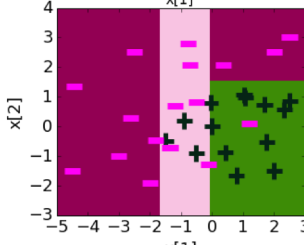
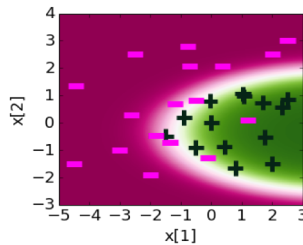
Logistic Regression
(Degree 2)



Decision Trees
(Depth 2)



Class



Probability

82

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

Summary of decision trees

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning

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

84

©2018 Emily Fox

STAT/CSE 416: Intro to Machine Learning