# Decision Trees
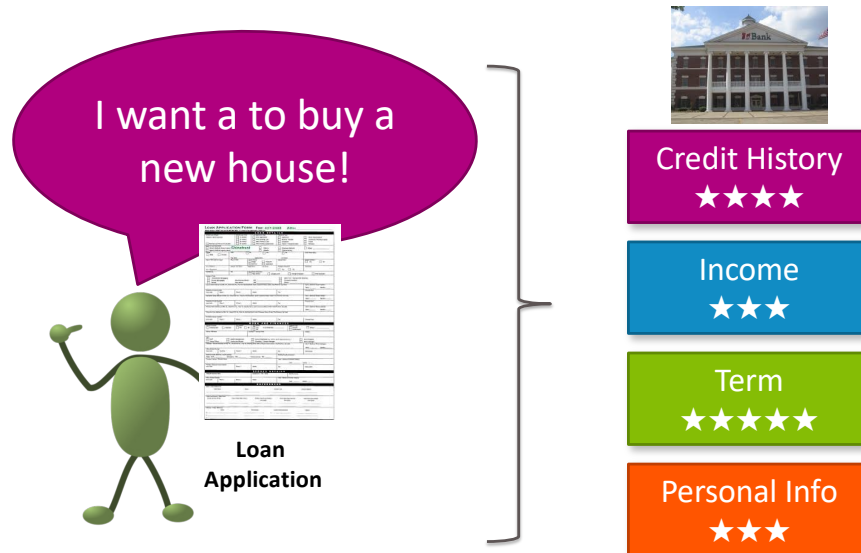
CS229: Machine Learning
Carlos Guestrin
Stanford University

**Slides include content developed by and co-developed with Emily Fox**
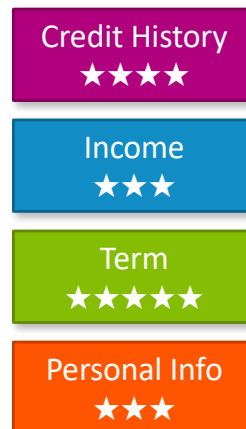
# Predicting potential loan defaults

# What makes a loan risky?

©2021 Carlos Guestrin

CS229: Machine Learning

# Credit history explained

Did I pay previous
loans on time?

**Example:** excellent,
good, or fair

| | |
|---|---|
| Credit History ★★★★ | |
| Income ★★★ | |
| Term ★★★★★ | |
| Personal Info ★★★ | |

©2021 Carlos Guestrin

CS229: Machine Learning

# Income

What's my income?

**Example:**
$80K per year

Credit History
★★★★

Income
★★★

Term
★★★★★

Personal Info
★★★

©2021 Carlos Guestrin

CS229: Machine Learning

# Loan terms

How soon do I need to
pay the loan?

**Example:** 3 years,
5 years,…

Credit History
★★★★

Income
★★★

Term
★★★★★

Personal Info
★★★

©2021 Carlos Guestrin

CS229: Machine Learning
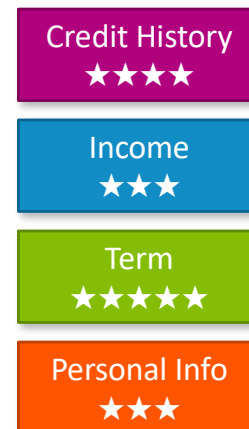
# Personal information

Age, reason for the loan,
marital status,…

**Example:** Home loan for a
married couple

| | |
|---|---|
| Credit History ★★★★ | |
| Income ★★★ | |
| Term ★★★★★ | |
| Personal Info ★★★ | |

©2021 Carlos Guestrin

CS229: Machine Learning

# Classifier review

Input: $\mathbf{x}_i$

Loan Application → Classifier MODEL

Output: $\hat{y}$ Predicted class

$\hat{y}_i = +1$ Safe

Risky $\hat{y}_i = -1$

# This module ... decision trees

©2021 Carlos Guestrin

CS229: Machine Learning

# Scoring a loan application

$x_i$ = (Credit = poor, Income = high, Term = 5 years)

Start

**excellent**  Credit?  **poor**

**fair**

Safe

Term?  Income?

**3 years**  **5 years**  **high**  **Low**

Risky  Safe  Term?  Risky

**3 years**  **5 years**

Risky  Safe  ◄— **ŷ_i = Safe**

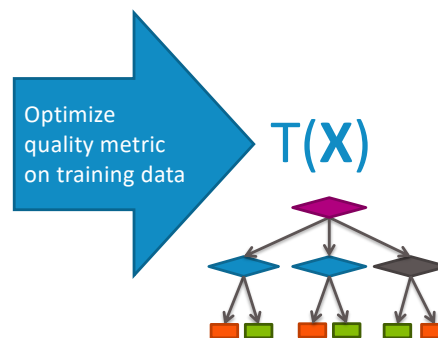©2021 Carlos Guestrin

CS229: Machine Learning

# Decision tree learning task

# Decision tree learning problem

Training data: N observations ($x_i, y_i$)

| 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 quality metric on training data

T(**X**)

CS229: Machine Learning

# Quality metric: Classification error

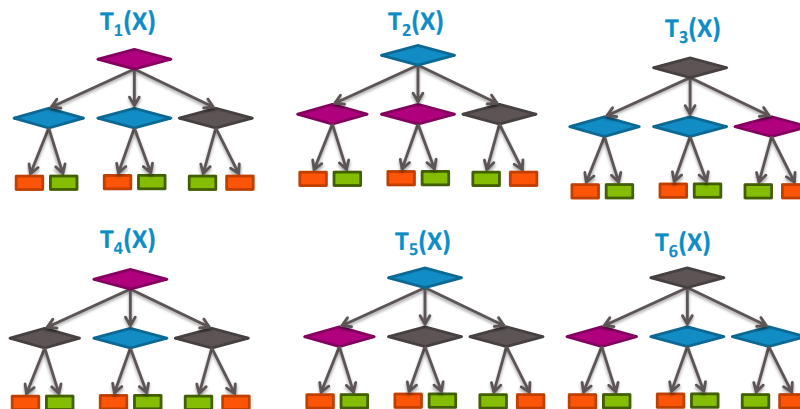• Error measures fraction of mistakes

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

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

# 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]
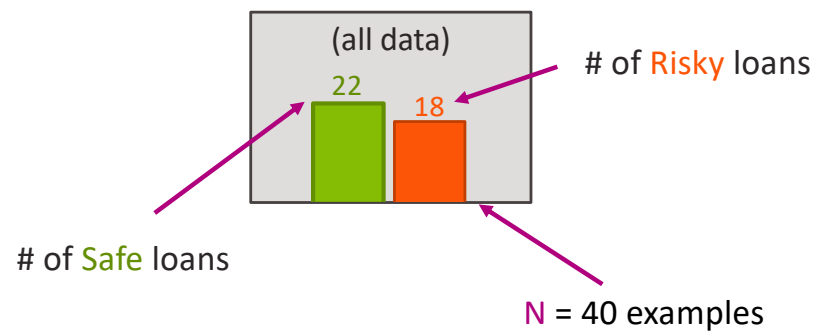
Greedy decision tree learning
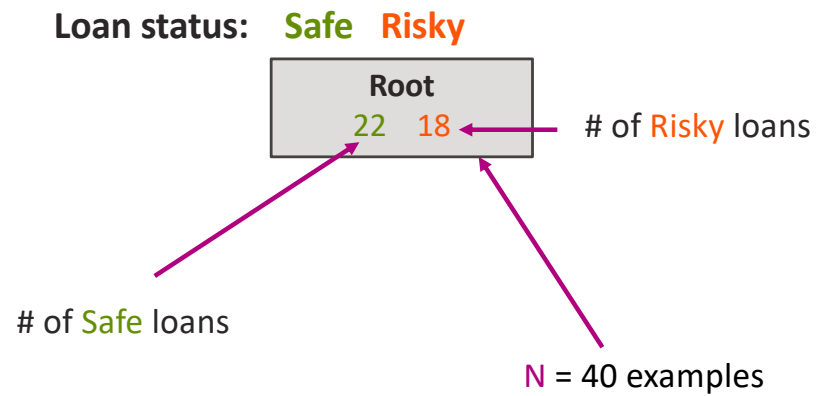
# Our training data table

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 |

CS229: Machine Learning

# Start with all the data

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



(all data)
22    18

# of Risky loans

# of Safe loans

N = 40 examples

# Compact visual notation: Root node

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



# of Risky loans

# of Safe loans

N = 40 examples

©2021 Carlos Guestrin                                    CS229: Machine Learning

# Decision stump: Single level tree

**Loan status:**
Safe Risky



Split on Credit

Root
22   18

Credit?

excellent
9   0

fair
9   4

poor
4   14

Subset of data with
Credit = excellent

Subset of data with
Credit = fair

Subset of data with
Credit = poor

CS229: Machine Learning

# Visual notation: Intermediate nodes

**Loan status:**
Safe Risky

Root
22   18

Credit?

excellent
9   0

fair
9   4

poor
4   14

Intermediate nodes

# Making predictions with a decision stump

**Loan status:**
Safe  Risky

root
22  18

credit?

| excellent | fair | poor |
|-----------|------|------|
| 9  0 | 9  4 | 4  14 |

Safe    Safe    Risky

For each intermediate node,
set ŷ = majority value

CS229: Machine Learning

Selecting best feature to split on

# How do we learn a decision stump?



**Loan status:**
Safe Risky

Root
22  18

Find the "best"
feature to split on!

Credit?

excellent
9  0

fair
9  4

poor
4  14

©2021 Carlos Guestrin

CS229: Machine Learning

# How do we select the best feature?

**Choice 1:** Split on Credit

Loan status:
Safe  Risky

Root
22  18

Credit?

excellent
9   0

fair
9   4

poor
4   14

**OR**

**Choice 2:** Split on Term

Loan status:
Safe  Risky

Root
22  18

Term?

3 years
16   4

5 years
6   14

CS229: Machine Learning

# How do we measure effectiveness of a split?

**Loan status:**
Safe  Risky

Root
22  18

Credit?

excellent
9  0

fair
9  4

poor
4  14

**Idea:** Calculate classification error of this decision stump

Error = # mistakes
# data points

# Calculating classification error

- Step 1: ŷ = class of majority of data in node
- Step 2: Calculate classification error of predicting ŷ for this data

**Loan status:**
Safe  Risky

Root
22  18
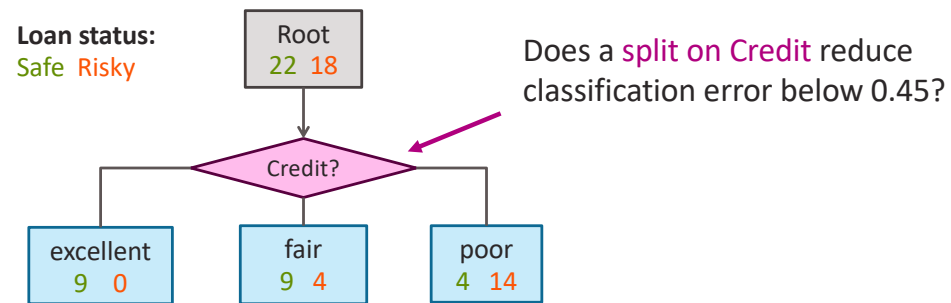
22 correct

18 mistakes

Safe

ŷ = **majority class**

Error = _____.

=

| Tree | Classification error |
|------|---------------------|
| (root) | 0.45 |

# Choice 1: Split on Credit history?

Choice 1: Split on Credit

**Loan status:**
Safe  Risky



Root
22  18

Credit?

excellent
9  0

fair
9  4

poor
4  14

Does a split on Credit reduce classification error below 0.45?

©2021 Carlos Guestrin

CS229: Machine Learning

# Split on Credit: Classification error

Choice 1: Split on Credit

**Loan status:**
Safe  Risky



Root
22  18

Credit?

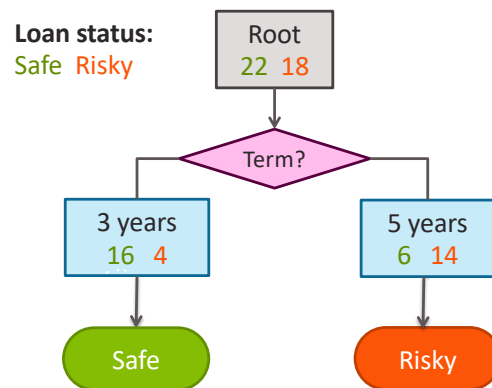| excellent | fair | poor |
| 9  0 | 9  4 | 4  14 |

Safe    Safe    Risky

0 mistakes    4 mistakes    4 mistakes

Error = _____.

=

| Tree | Classification error |
|------|---------------------|
| (root) | 0.45 |
| Split on credit | 0.2 |

29

©2021 Carlos Guestrin
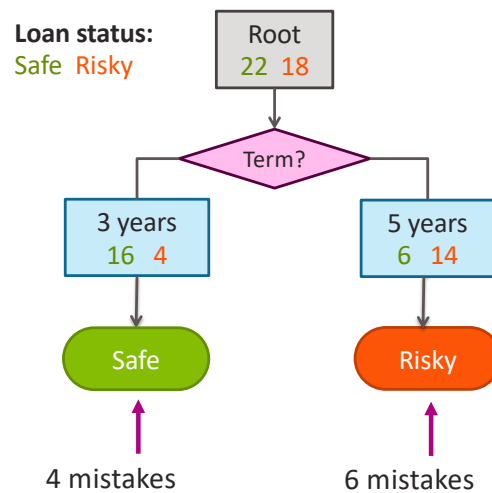
CS229: Machine Learning

# Choice 2: Split on Term?

Choice 2: Split on Term

**Loan status:**
Safe  Risky

©2021 Carlos Guestrin

CS229: Machine Learning

# Evaluating the split on Term

### Choice 2: Split on Term

**Loan status:**
Safe  Risky

Root
22  18

Term?

3 years
16  4

5 years
6  14

Safe

Risky

4 mistakes

6 mistakes

Error = _____

=

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

©2021 Carlos Guestrin

CS229: Machine Learning

# Choice 1 vs Choice 2: Comparing split on Credit vs Term
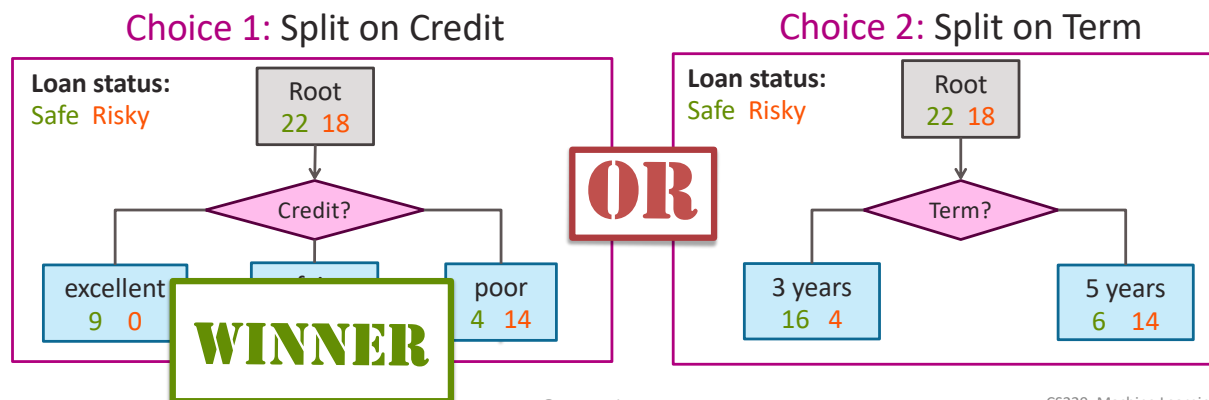
| Tree | Classification error |
|------|---------------------|
| (root) | 0.45 |
| split on credit | 0.2 |
| split on loan term | 0.25 |

### Choice 1: Split on Credit



### Choice 2: Split on Term



OR

WINNER

©2021 Carlos Guestrin

CS229: Machine Learning

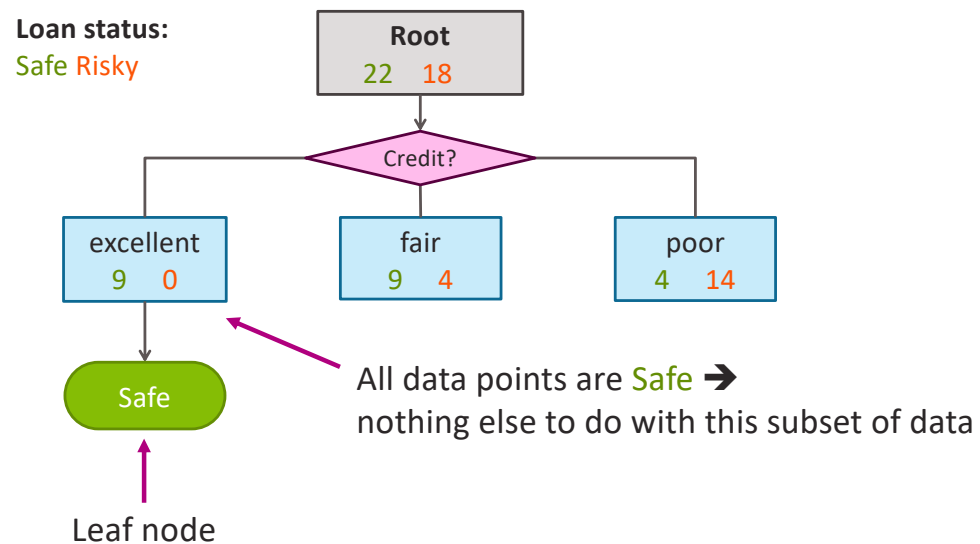# Feature split selection algorithm

- Given a subset of data $M$ (a node in a tree)

- For each feature $h_i(x)$:
    1. Split data of $M$ according to feature $h_i(x)$
    2. Compute classification error of split

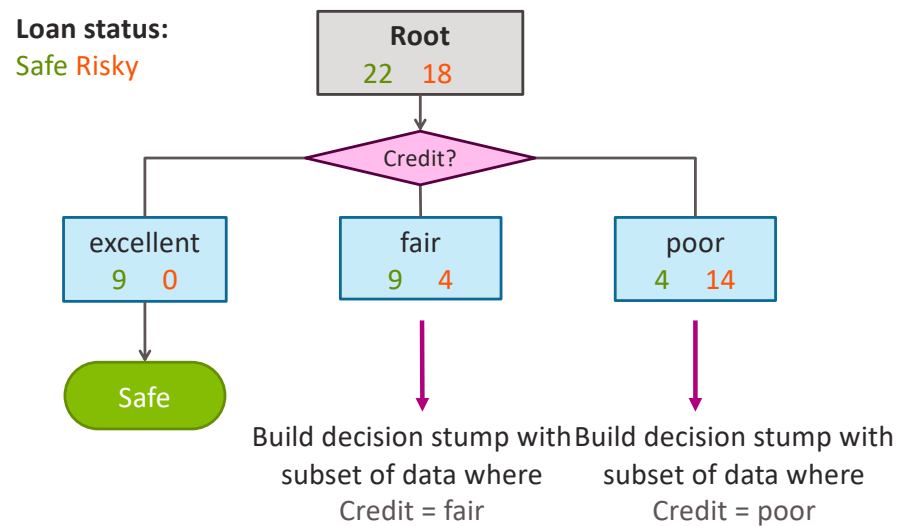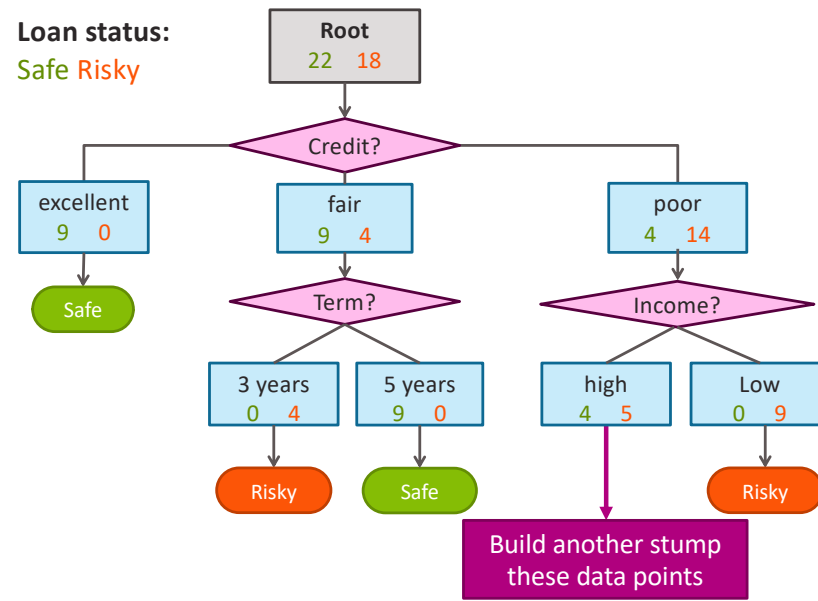- Chose feature $h^*(x)$ with lowest classification error

©2021 Carlos Guestrin CS229: Machine Learning

# Recursion & Stopping conditions

# We've learned a decision stump, what next?

**Loan status:**
Safe Risky

Root
22  18

Credit?

| excellent | fair | poor |
|-----------|------|------|
| 9    0    | 9    4 | 4    14 |

Safe

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

Leaf node

©2021 Carlos Guestrin

CS229: Machine Learning

# Tree learning = Recursive stump learning

**Loan status:**
Safe Risky

```
        Root
        22  18
         │
      ◇ Credit? ◇
   ┌──────┼───────┐
excellent  fair    poor
 9   0    9   4    4   14
   │        │        │
 [Safe]     ↓        ↓
```

Build decision stump with subset of data where Credit = fair

Build decision stump with subset of data where Credit = poor

©2021 Carlos Guestrin

CS229: Machine Learning

# Second level

**Loan status:**
Safe Risky

```
                    ┌──────────┐
                    │   Root   │
                    │  22  18  │
                    └────┬─────┘
                         ▼
                    ◇ Credit? ◇
         ┌───────────────┼───────────────┐
    ┌─────────┐     ┌─────────┐      ┌─────────┐
    │excellent│     │  fair   │      │  poor   │
    │  9   0  │     │  9   4  │      │  4  14  │
    └────┬────┘     └────┬────┘      └────┬────┘
         ▼               ▼                ▼
      ( Safe )       ◇ Term? ◇       ◇ Income? ◇
                  ┌──────┴──────┐  ┌──────┴──────┐
              ┌───────┐   ┌───────┐┌───────┐ ┌───────┐
              │3 years│   │5 years││ high  │ │  Low  │
              │ 0   4 │   │ 9   0 ││ 4   5 │ │ 0   9 │
              └───┬───┘   └───┬───┘└───┬───┘ └───┬───┘
                  ▼           ▼        ▼         ▼
              ( Risky )   ( Safe )             ( Risky )
```

**Build another stump these data points**

37

# Final decision tree

**Loan status:**
Safe Risky

©2021 Carlos Guestrin

# 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???

©2021 Carlos Guestrin

# Stopping condition 1: All data agrees on y



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

Root
22  18

poor
4  14

Credit?

Income?

excellent
9  0

Fair
9  4

high
4  5

low
0  9

Safe

Term?

Term?

Risky

3 years
0  4

5 years
9  0

3 years
0  2

5 years
4  3

Risky

Safe

Risky

Safe

40

# Stopping condition 2: Already split on all features



Already split on all possible features ➜ Nothing to do

Root
22  18

Credit?

excellent
9   0

Safe

Fair
9   4

Term?

3 years
0   4

Risky

5 years
9   0

Safe

poor
4   14

Income?

high
4   5

low
0   9

Risky

Term?

3 years
0   2

Risky

5 years
4   3

Safe

41

©2021 Carlos Guestrin

CS229: 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 do, 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

©2021 Carlos Guestrin

CS229: Machine Learning

Is this a good idea?

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

# Stopping condition 3:
## Don't stop if error doesn't decrease???

**y** = **x**[1] **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 |

Error = _____.

    =

| Tree | Classification error |
|------|----------------------|
| (root) | 0.5 |

©2021 Carlos Guestrin

CS229: Machine Learning

# Consider split on x[1]

**y = x[1] 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  │
        └────┬─────┘
             │
             ▼
          ◇ x[1] ◇
         ╱        ╲
   ┌────────┐  ┌────────┐
   │ True   │  │ False  │
   │ 1   1  │  │ 1   1  │
   └────────┘  └────────┘
```

Error = _____.

=

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

# Consider split on x[2]

**y** = **x**[1] **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

x[2]

True
1   1

False
1   1

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 |

46

# Final tree with stopping condition 3

**y** = **x**[1] **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 |

Predict True

| Tree | Classification error |
|------|---------------------|
| with stopping condition 3 | 0.5 |

# Without stopping condition 3

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

## y = x[1] 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 | |



y values
True  False

©2021 Carlos Guestrin

CS229: Machine Learning

Decision tree learning:
*Real valued features*

# 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 |

# Threshold split

**Loan status:**
Safe Risky



Root
22  18

Split on the feature Income

Income?

< $60K
8   13

>= $60K
14   5

Subset of data with
Income >= $60K

©2021 Carlos Guestrin

CS229: Machine Learning

# Finding the best threshold split



**Infinite possible values of t**

Income = t$^*$

Income < t$^*$    Income >= t$^*$

Safe ○
Risky ○

**Income**

$10K    $120K

©2021 Carlos Guestrin    CS229: Machine Learning

# Consider a threshold between points

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



Safe ⬤
Risky ⬤

Income

$v_A$  $v_B$

$10K                                    $120K

©2021 Carlos Guestrin

CS229: Machine Learning

# Only need to consider mid-points

Finite number of splits
to consider

Safe ◯
Risky ◯

**Income**

$10K

$120K

©2021 Carlos Guestrin

CS229: Machine Learning

# Threshold split selection algorithm

- Step 1: Sort the values of a feature $h_j(x)$ :

    Let $\{v_1, v_2, v_3, ... v_N\}$ denote sorted values

- Step 2:
    - For i = 1 ... N-1
        - Consider split $t_i = (v_i + v_{i+1}) / 2$
        - Compute classification error for treshold split $h_j(x) >= t_i$
    - Chose the $t^*$ with the lowest classification error

# Visualizing the threshold split



Threshold split is the line Age = 38

Income

Age

0   10   20   30   40   ...

$0K

$40K

$80K

...

©2021 Carlos Guestrin

CS229: Machine Learning

# Split on Age >= 38

©2021 Carlos Guestrin

CS229: Machine Learning

# Depth 2: Split on Income >= $60K

Threshold split is the line Income = 60K
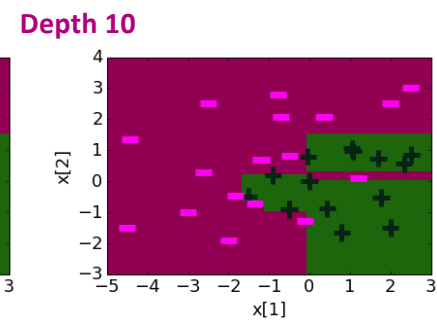
©2021 Carlos Guestrin

CS229: Machine Learning

# Each split partitions the 2-D space

©2021 Carlos Guestrin

CS229: Machine Learning

Decision trees vs logistic regression:
*Example*

# Logistic regression

| Feature | Value | Weight Learned |
|---------|-------|----------------|
| $h_0(x)$ | 1 | 0.22 |
| $h_1(x)$ | x[1] | 1.12 |
| $h_2(x)$ | x[2] | -1.07 |

©2021 Carlos Guestrin

CS229: Machine Learning

# Depth 1: Split on x[1]



**y values**

    −  +

Root
18   13

x[1]

x[1] < -0.07
13   3
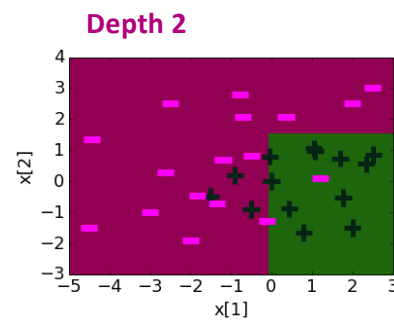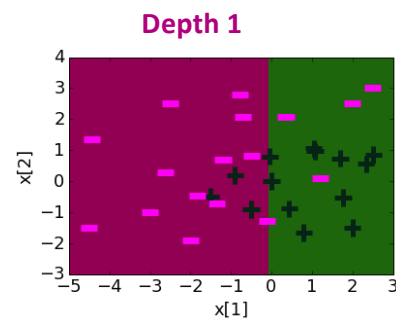
x[1] >= -0.07
4   11

CS229: Machine Learning

# Depth 2

©2021 Carlos Guestrin
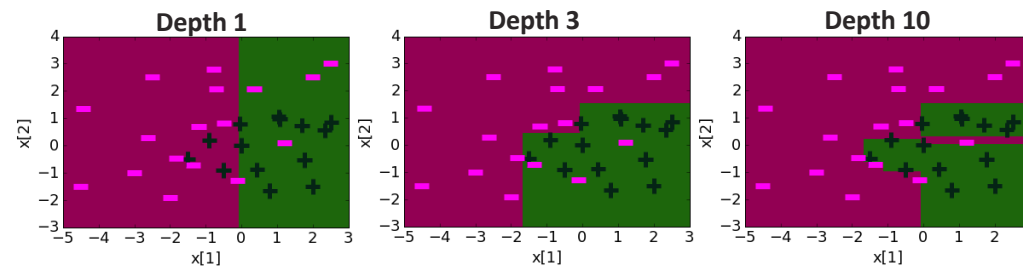
CS229: Machine Learning

# Threshold split caveat



For threshold splits, same feature can be used multiple times

y values
- +

Root
18   13

x[1]

x[1] < -0.07
13   3

x[1] >= -0.07
4   11

x[1]

x[2]

x[1] < -1.66
7   0

x[1] >= -1.66
6   3

x[2] < 1.55
1   11

x[2] >=  1.55
3   0

©2021 Carlos Guestrin

CS229: Machine Learning

# Decision boundaries



**Depth 1**　　　　　　　　**Depth 2**　　　　　　　　**Depth 10**

©2021 Carlos Guestrin　　　　　　　　　　　　CS229: Machine Learning

# Comparing decision boundaries

**Decision Tree**



**Logistic Regression**

CS229: Machine Learning

# Summary of decision trees

# 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
- Tackle continuous and discrete features