# 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?



I want a to buy a new house!

Loan Application

Credit History ★★★★

Income ★★★

Term ★★★★★

Personal Info ★★★

# Credit history explained

Did I pay previous loans on time?

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

Credit History
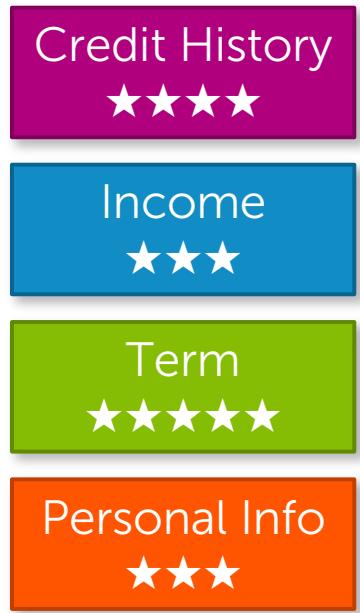★★★★

Income
★★★

Term
★★★★★

Personal Info
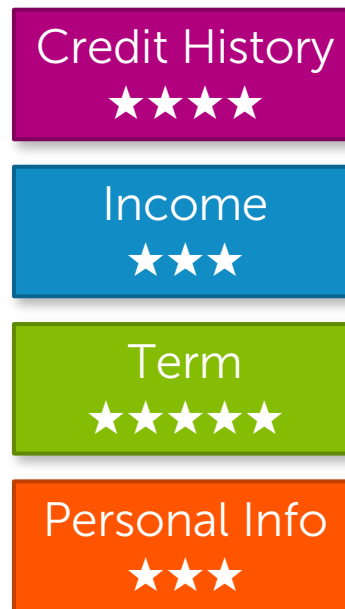★★★

# Income

What's my income?

**Example:**
$80K per year

Credit History
★★★★

Income
★★★

Term
★★★★★

Personal Info
★★★

# Loan terms

How soon do I need to pay the loan?

**Example:** 3 years, 5 years,...

Credit History
★★★★

Income
★★★

Term
★★★★★

Personal Info
★★★
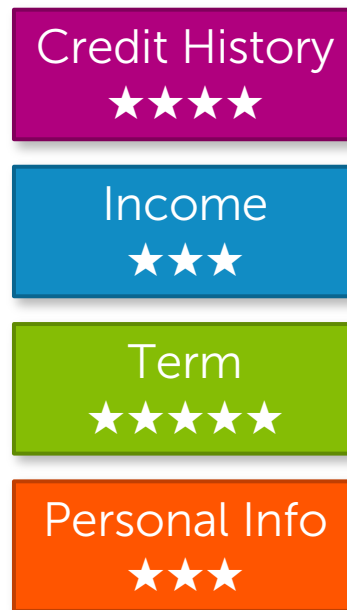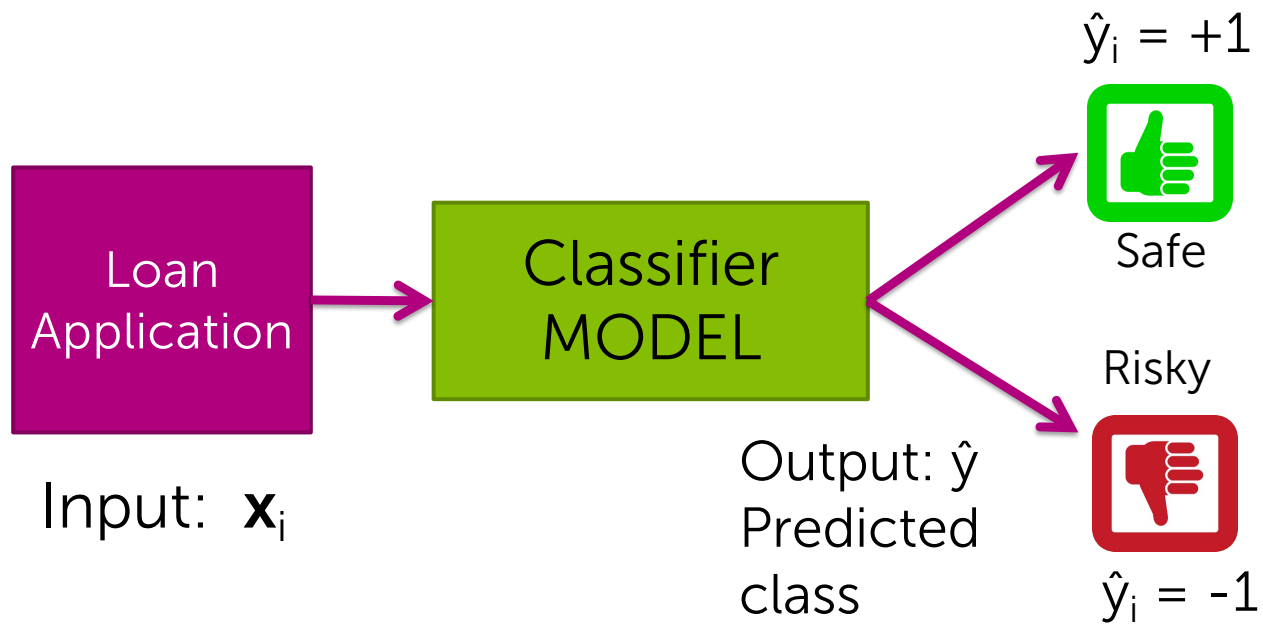
# Personal information

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

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

Credit History ★★★★

Income ★★★

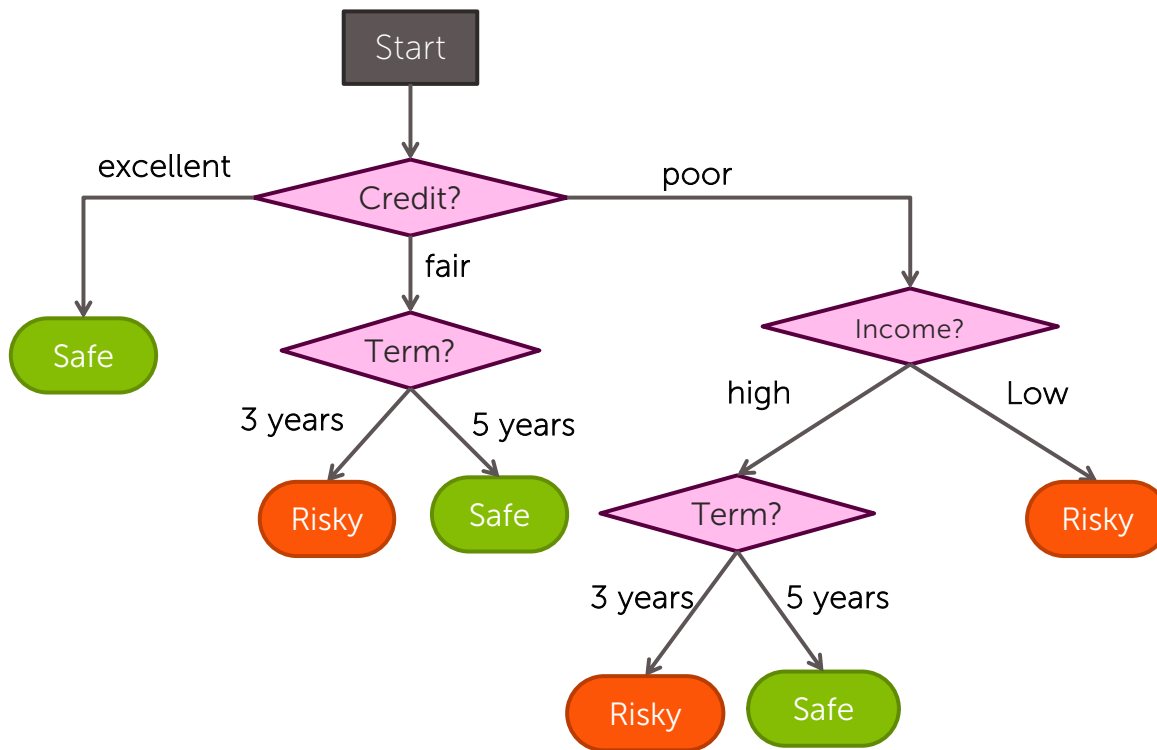Term ★★★★★

Personal Info ★★★

CS229: Machine Learning

# Classifier review



Input: $\mathbf{x}_i$

Classifier MODEL

Output: ŷ
Predicted class

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

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

# This module ... decision trees

©2022 Carlos Guestrin
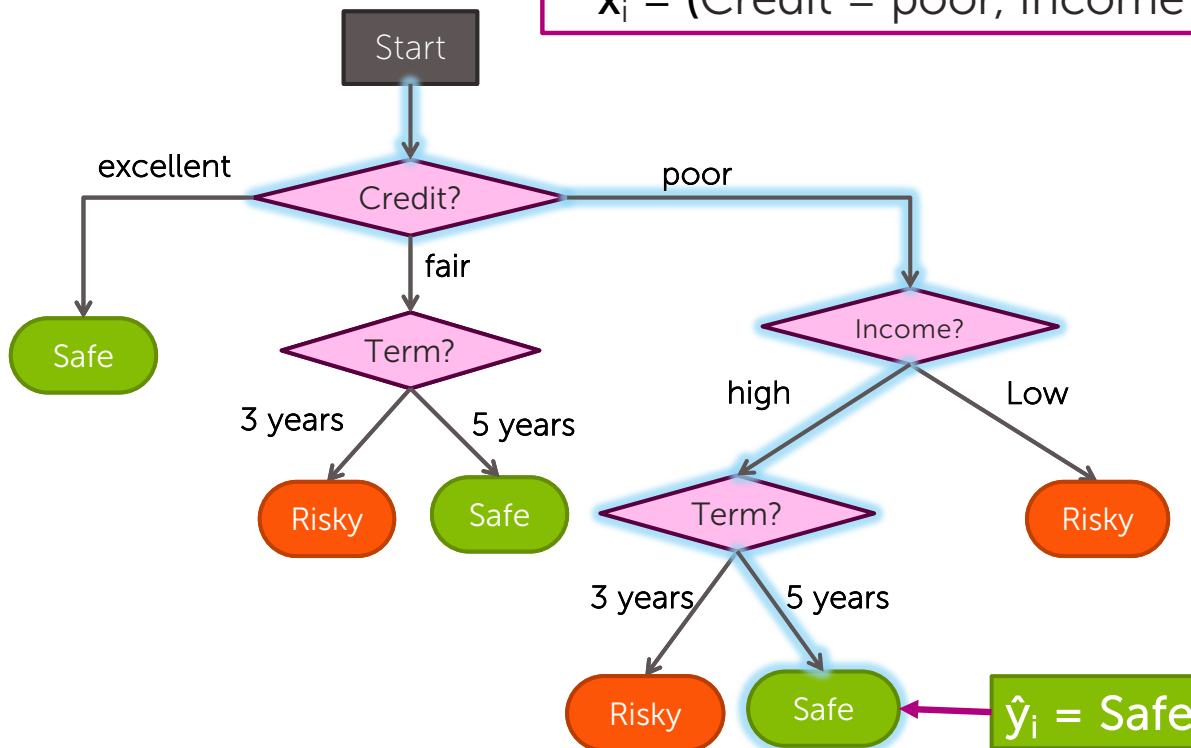
# Scoring a loan application



$\mathbf{x}_i = (\text{Credit} = \text{poor}, \text{Income} = \text{high}, \text{Term} = 5 \text{ years})$

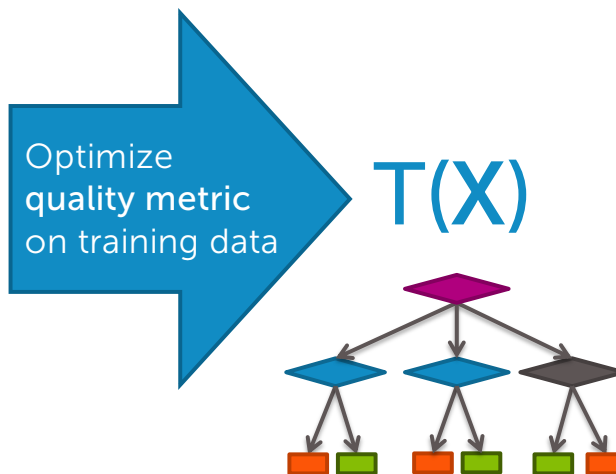$\hat{y}_i = \text{Safe}$

©2022 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

$T(X)$

# Quality metric: Classification error
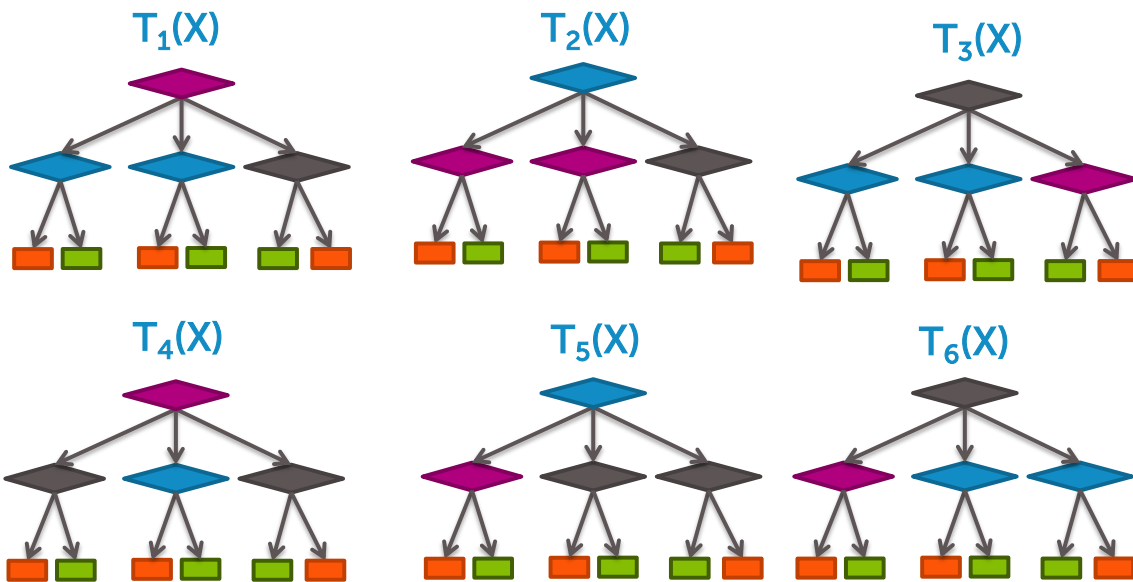
- Error measures fraction of mistakes

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

  – 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]



$T_1(X)$  $T_2(X)$  $T_3(X)$

$T_4(X)$  $T_5(X)$  $T_6(X)$

# 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 |

# 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**



Root
22    18

# of **Risky** loans

# of **Safe** loans

N = 40 examples

# Decision stump: Single level tree



Loan status:
Safe Risky

Root
22    18

Split on **Credit**

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**

# 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?

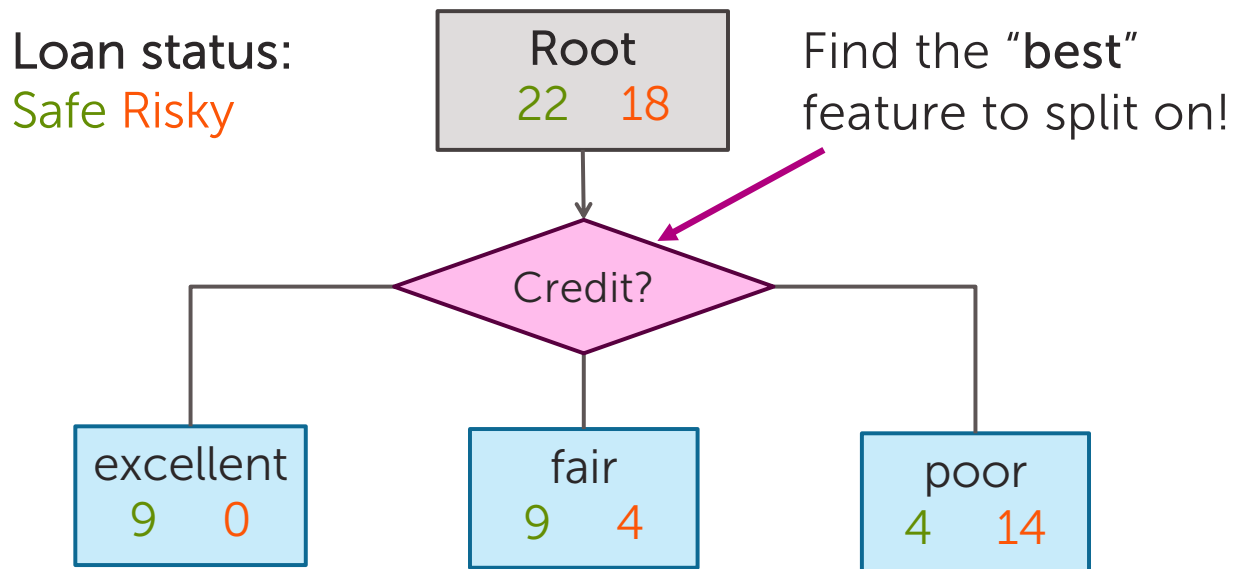| excellent | fair | poor |
| 9  0 | 9  4 | 4  14 |

Safe    Safe    Risky

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

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

©2022 Carlos Guestrin

CS229: Machine Learning
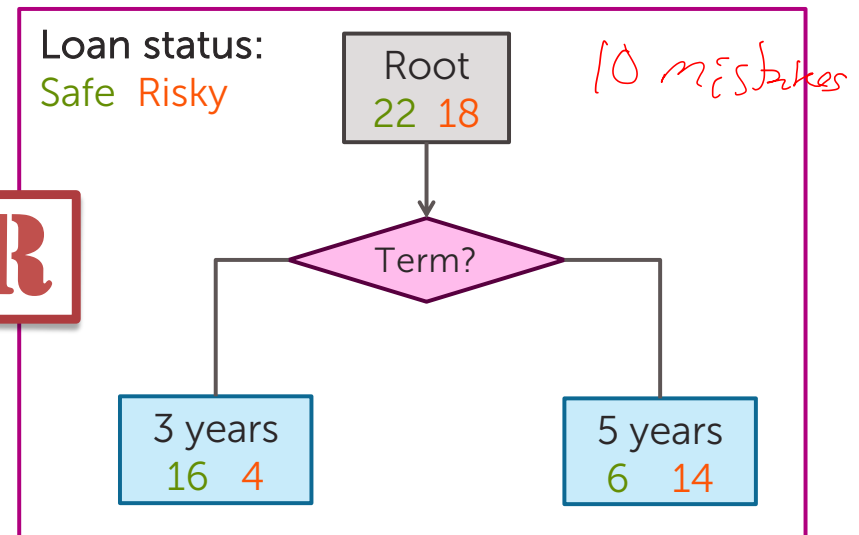
# How do we select the best feature?

*Which one makes fewer mistakes?*

## Choice 1: Split on **Credit**

Loan status:
Safe  Risky

*8 mistakes*

```
        Root
        22  18
          |
       Credit?
     /     |     \
excellent  fair   poor
  9   0    9  4   4  14
```

**OR**

## Choice 2: Split on **Term**

Loan status:
Safe  Risky

*10 mistakes*

```
        Root
        22  18
          |
        Term?
       /      \
  3 years    5 years
  16   4      6  14
```

©2022 Carlos Guestrin

CS229: Machine Learning

# How do we measure effectiveness of a split?

Loan status:
Safe  Risky

Root
22  18

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

Credit?

excellent
9   0

fair
9   4

poor
4   14

Error =  # mistakes
# data points

# Calculating classification error

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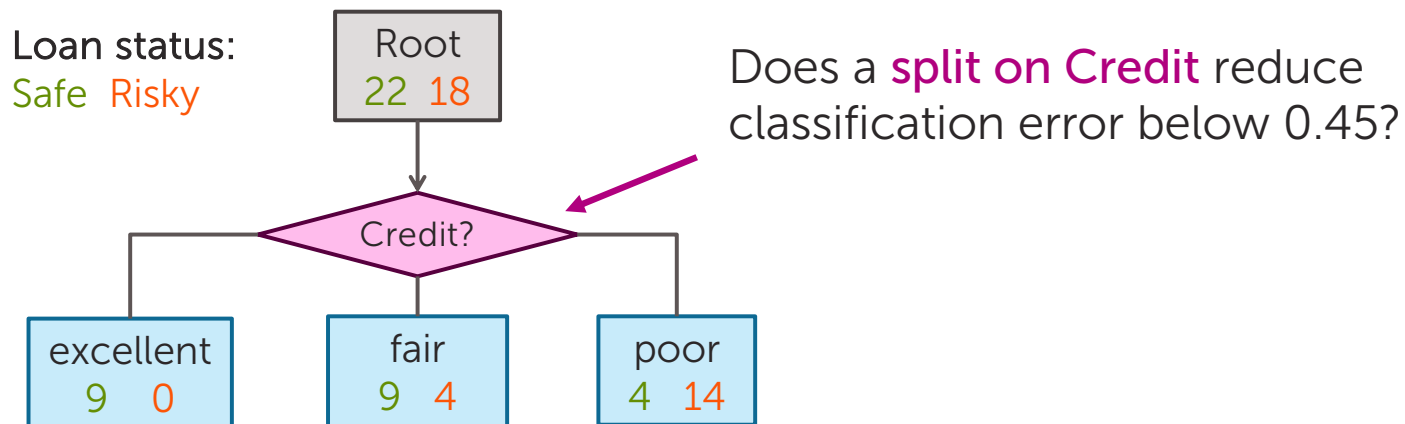- **Step 2:** Calculate classification error of predicting ŷ for this data

Loan status:
Safe  Risky

Root
22   18

22 correct          18 mistakes

Safe

ŷ = majority class

Error = ___*18*___.

$= \frac{18}{40}$

$= 0.45$

| 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?

©2022 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
```

$$\text{Error} = \frac{8}{40}$$

$$= 0.2$$

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

©2022 Carlos Guestrin

# Choice 2: 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

# 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 = $\frac{10}{40}$

= 0.25

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

# 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**

Loan status:
Safe  Risky

Root
22  18

Credit?

excellent
9   0

poor
4   14

WINNER

**OR**

**Choice 2:** Split on **Term**

Loan status:
Safe  Risky

Root
22  18

Term?

3 years
16   4

5 years
6   14

# 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

# Recursion & Stopping conditions
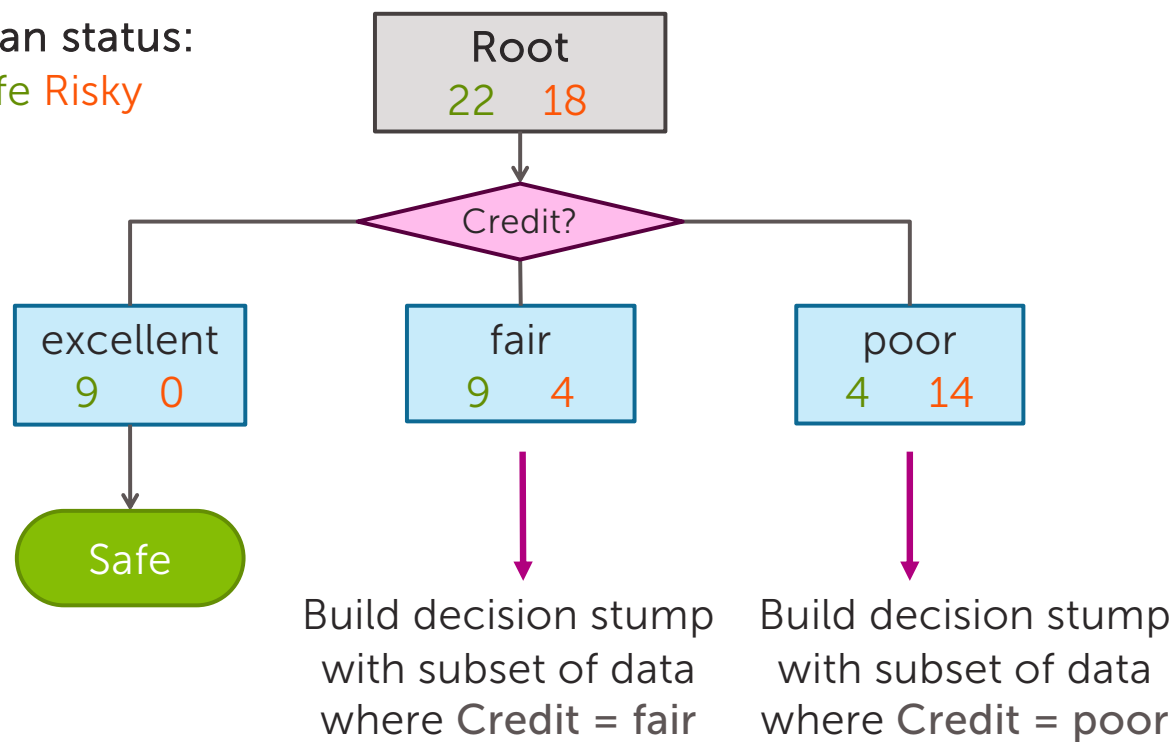
# We've learned a decision stump, what next?

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

# Tree learning = Recursive stump learning

Loan status:
Safe Risky



Root
22   18

Credit?

excellent
9    0

fair
9    4

poor
4    14

Safe

Build decision stump
with subset of data
where **Credit = fair**

Build decision stump
with subset of data
where **Credit = poor**

# Second level

Loan status:
Safe Risky

Root
22  18

Credit?

excellent
9  0

fair
9  4

poor
4  14

Safe

Term?

Income?

3 years
0  4

5 years
9  0

high
4  5

Low
0  9

Risky

Safe

Risky

**Build another stump these data points**

# Final decision tree

Loan status:
Safe Risky

©2022 Carlos Guestrin
CS229: 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???

©2022 Carlos Guestrin

CS229: Machine Learning

# 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

©2022 Carlos Guestrin

# Stopping condition 2: Already split on all features



**Already split on all possible features ➔ Nothing to do**

41

# 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

Recursion

©2022 Carlos Guestrin

CS229: Machine Learning

# Is this a good idea?

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

©2022 Carlos Guestrin

CS229: Machine Learning

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

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

| $\mathbf{x}[1]$ | $\mathbf{x}[2]$ | $y$ |
|------|------|------|
| False | False | False |
| False | True | True |
| True | False | True |
| True | True | False |

y values
True  False

Root
2    2

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

= 0.5

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

# 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
1    1

False
1    1

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

= 0.5

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

©2022 Carlos Guestrin

# 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

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

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 |

# 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

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

| 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

©2022 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 — turn continuous var into binary

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

# Finding the best threshold split

**Infinite possible values of t**

Income = t*

Income < t*          Income >= t*

Safe ○
Risky ○

**Income**

$10K                    $120K

# Consider a threshold between points

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



Income

$v_A$   $v_B$

Safe ○
Risky ○

$10K                                      $120K

# Only need to consider mid-points

*Sort data:*

Finite number of
splits to consider



Safe ◯
Risky ◯

**Income**

$10K

$120K

©2022 Carlos Guestrin

CS229: Machine Learning

# Threshold split selection algorithm

- Step 1: Sort the values of a feature $h_j(\mathbf{x})$ :

    Let $\{v_1, v_2, v_3, \ldots v_N\}$ 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 $h_j(\mathbf{x}) >= t_i$
  - Chose the $\mathbf{t}^*$ with the lowest classification error

# Visualizing the threshold split



Threshold split is the line Age = 38

Income

...

$80K

$40K

$0K

Age

0    10    20    30    40    ...

©2022 Carlos Guestrin

CS229: Machine Learning

# Split on Age >= 38



Income

age < 38    age >= 38

... 

$80K

$40K

$0K

0    10    20    30    40    ...

Age

Predict **Risky**

Predict **Safe**

©2022 Carlos Guestrin

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

Threshold split is the line **Income = 60K**

# Each split partitions the 2-D space

©2022 Carlos Guestrin

Decision trees vs logistic regression:
*Example*

# 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 |

# Depth 1: Split on x[1]



y values

&minus;  +

Root
18    13

x[1]

x[1] < -0.07
13    3

x[1] >= -0.07
4    11

# Depth 2



y values
   –  +

```
                    Root
                   18    13
                     |
                    x[1]
                   /      \
        x[1] < -0.07      x[1] >= -0.07
          13    3            4    11
            |                   |
          x[1]                x[2]
         /      \            /      \
x[1] < -1.66  x[1] >= -1.66  x[2] < 1.55  x[2] >= 1.55
   7    0       6    3         1    11       3    0
```

©2022 Carlos Guestrin

# Threshold split caveat

y values
-  +

Root
18   13

x[1]

For threshold splits, same feature can be used multiple times

x[1] < -0.07
13   3

x[1] >= -0.07
4   11

x[1]

x[2]

x[1] < -1.66
7   0
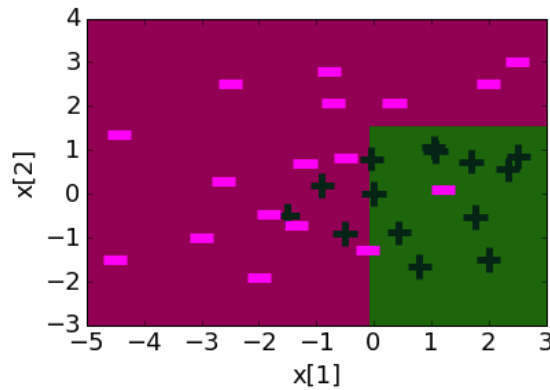
x[1] >= -1.66
6   3

x[2] < 1.55
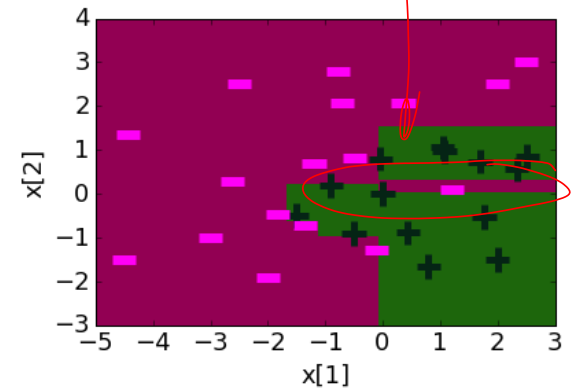1   11

x[2] >=  1.55
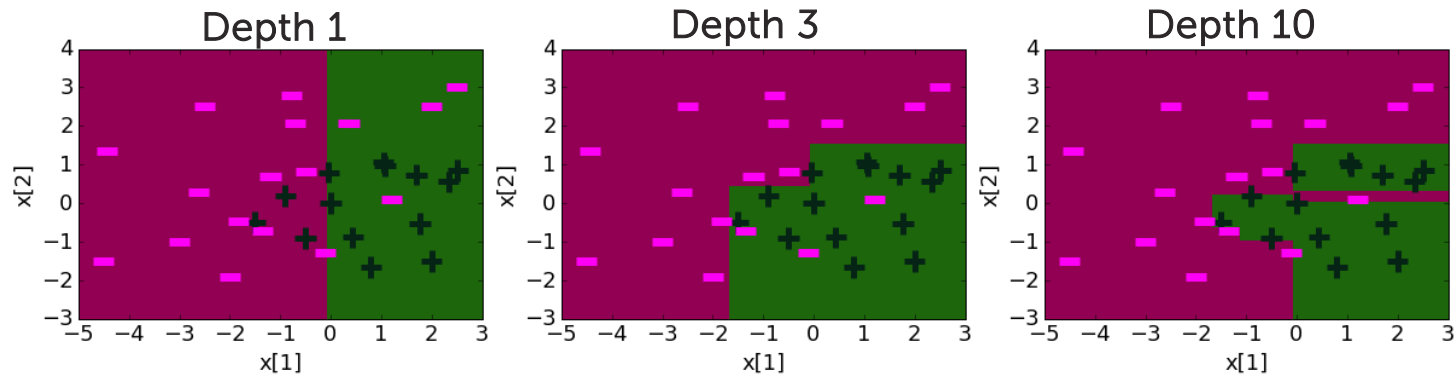3   0

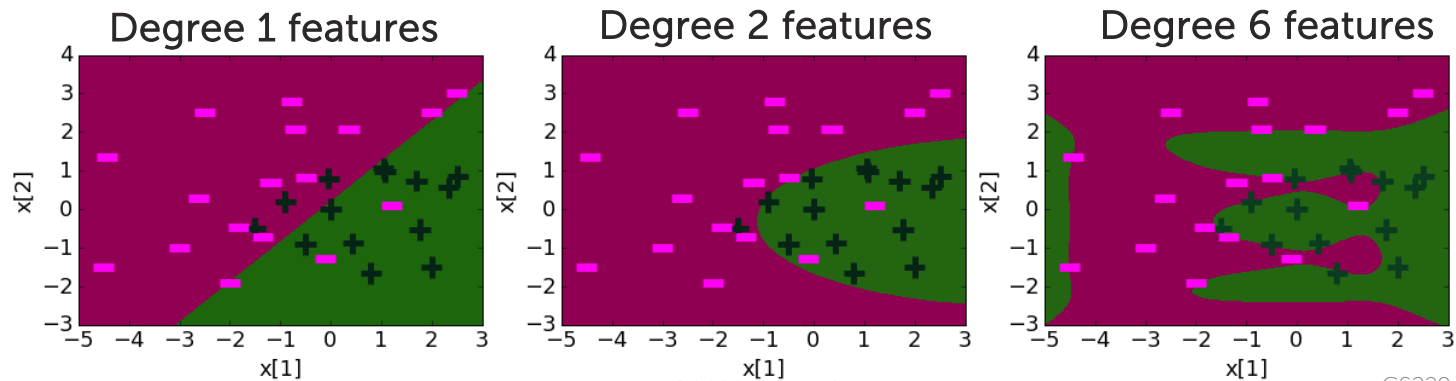# Decision boundaries



Depth 1

Depth 2

Depth 10

# Comparing decision boundaries

**Decision Tree**



**Logistic Regression**

# 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