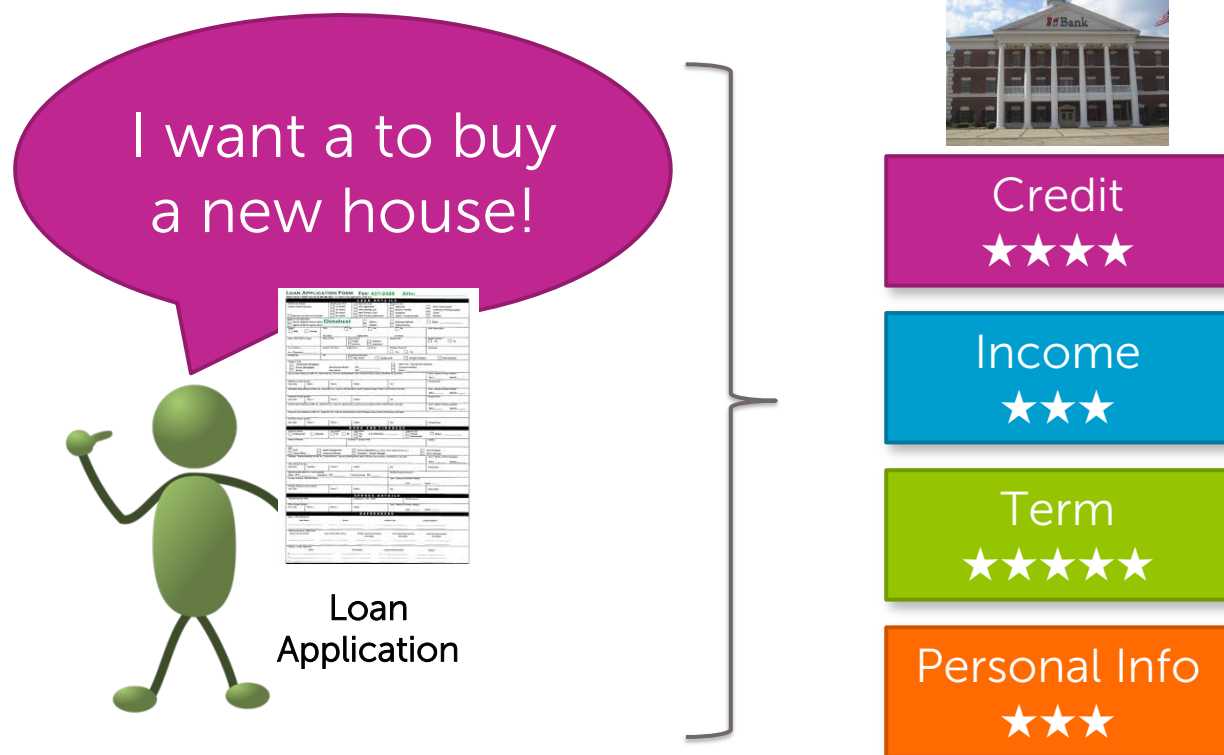# Decision Trees

Emily Fox & Carlos Guestrin

Machine Learning Specialization

University of Washington

# Predicting potential loan defaults

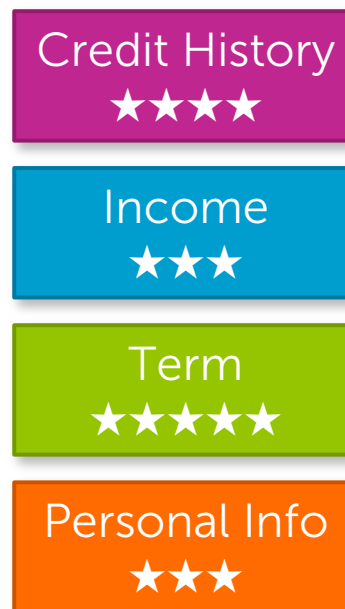# What makes a loan risky?



I want a to buy a new house!

Loan Application

Credit ★★★★

Income ★★★

Term ★★★★★

Personal Info ★★★

Machine Learning Specialization

# Credit history explained

Did I pay previous
loans on time?

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

Credit History
★★★★

Income
★★★

Term
★★★★★

Personal Info
★★★

Machine Learning Specialization

# Income

What's my income?

**Example:**
$80K per year



Credit History
★★★★

Income
★★★

Term
★★★★★

Personal Info
★★★

Machine Learning Specialization

# Loan terms

How soon do I need to pay the loan?

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

Credit History ★★★★

Income ★★★

Term ★★★★★

Personal Info ★★★

Machine Learning Specialization

# Personal information

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

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

Credit History
★★★★

Income
★★★

Term
★★★★★

Personal Info
★★★

Machine Learning Specialization

# Intelligent application



Loan Applications → Intelligent loan application review system → Safe ✓ / Risky ✗ / Risky ✗

Machine Learning Specialization

# Classifier review



Input: $\mathbf{x}_i$

Output: $\hat{y}$
Predicted
class

$\hat{y}_i = +1$

Safe

Risky

$\hat{y}_i = -1$

Machine Learning Specialization

# This module ... decision trees

```
                    ┌─────────┐
                    │  Start  │
                    └─────────┘
                         │
    excellent        ◇ Credit? ◇        poor
  ┌──────────────────                ──────────────────┐
  │                     │                              │
  ▼                   fair                             ▼
┌──────┐               │                          ◇ Income? ◇
│ Safe │           ◇ Term? ◇                   high          Low
└──────┘        3 years   5 years          ◇ Term? ◇      ┌───────┐
              ┌────────   ────────┐      3 years  5 years │ Risky │
              ▼                   ▼                        └───────┘
          ┌───────┐          ┌──────┐    ┌───────┐  ┌──────┐
          │ Risky │          │ Safe │    │ Risky │  │ Safe │
          └───────┘          └──────┘    └───────┘  └──────┘
```
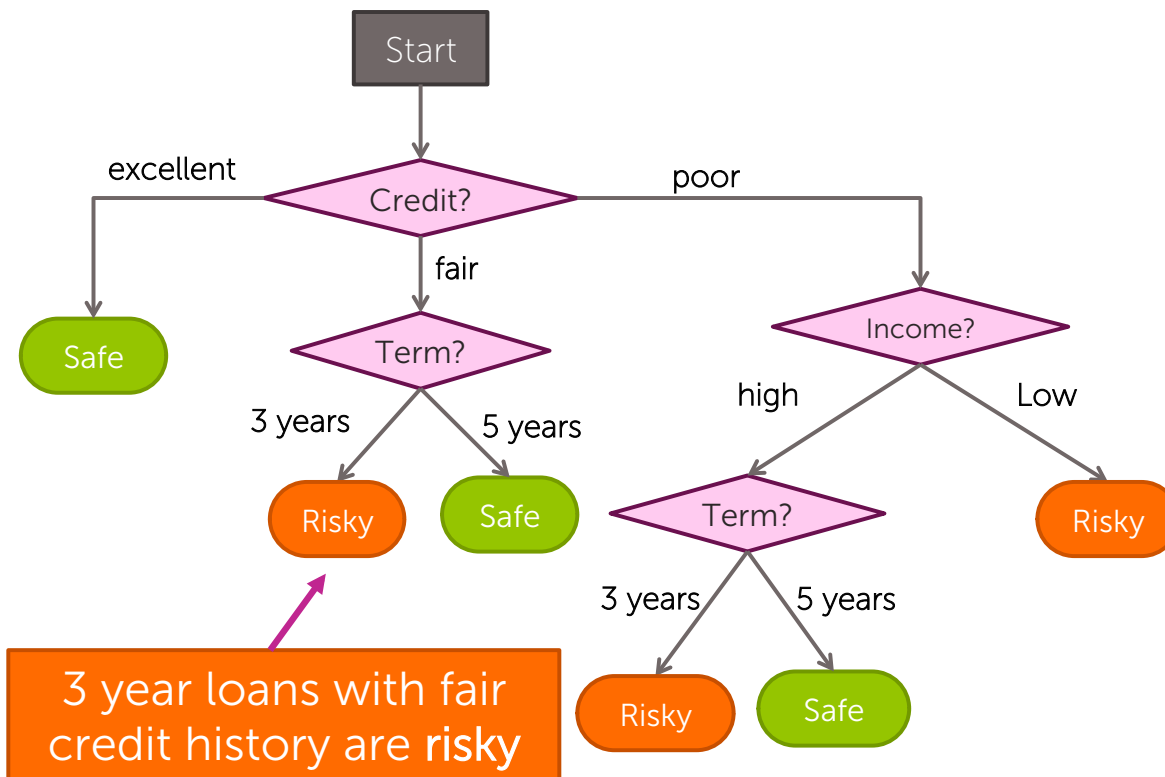
# Decision trees: *Intuition*

# What does a decision tree represent?

# What does a decision tree represent?



3 year loans with high income & poor credit history are **risky**

©2015-2016 Emily Fox & Carlos Guestrin

Machine Learning Specialization

# Scoring a loan application

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

©2015-2016 Emily Fox & Carlos Guestrin

Machine Learning Specialization

# Decision tree model



$T(\mathbf{x}_i)$ = Traverse decision tree

Loan Application

Input: $\mathbf{x}_i$

start

excellent — Credit? — poor

fair

Safe

Term?

3 years — Risky    5 years — Safe

Income?

high — Term?    Low — Risky

3 years — Risky    5 years — Safe

$\hat{y}_i$ = Safe

# Decision tree learning task

©2015-2016 Emily Fox & Carlos Guestrin

Machine Learning Specialization

# Learn decision tree from data?

$h_1(x)$   $h_2(x)$   $h_3(x)$   Loan Status

| 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 |

# 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 **quality metric** on training data

T(X)

# Quality metric: Classification error

- Error measures fraction of mistakes

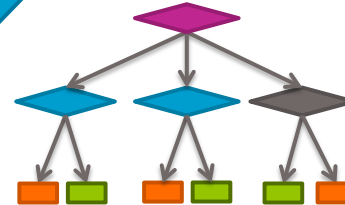$$Error = \frac{\# \text{ incorrect predictions}}{\# \text{ examples}}$$

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

Machine Learning Specialization

# Find the tree with lowest classification error

| 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 |

Minimize **classification error** on training data

T(X)

Machine Learning Specialization

# How do we find the best tree?

Exponentially large number of possible trees makes decision tree learning hard! *(NP-hard problem)*

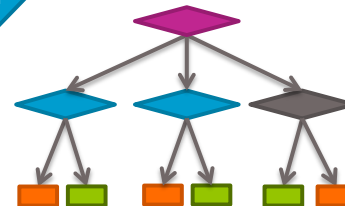$T_1(X)$    $T_2(X)$    $T_3(X)$

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

Machine Learning Specialization

# Simple (greedy) algorithm finds "good" tree

| 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 |

Approximately minimize **classification error** on training data
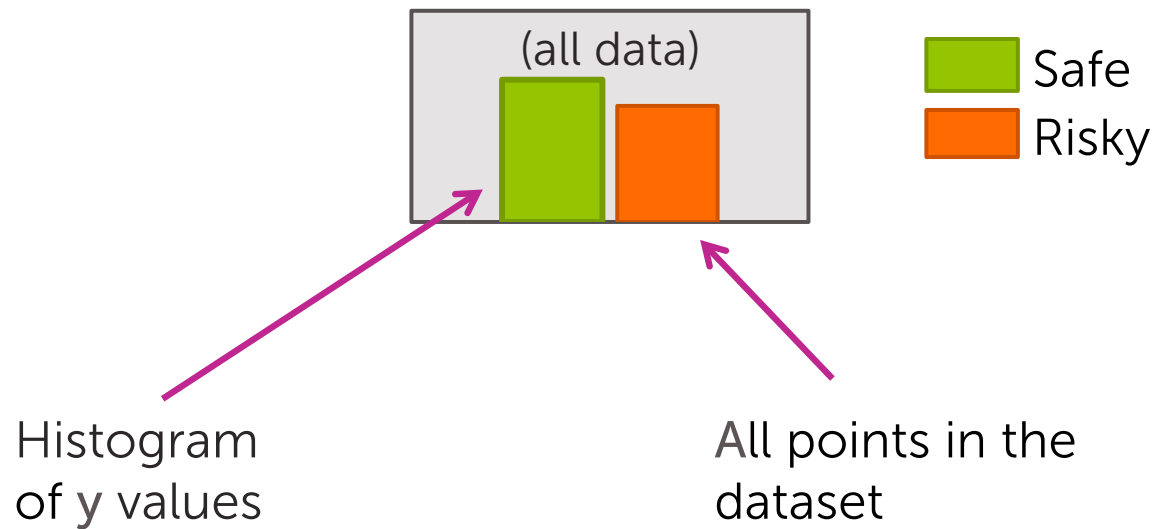
T(X)
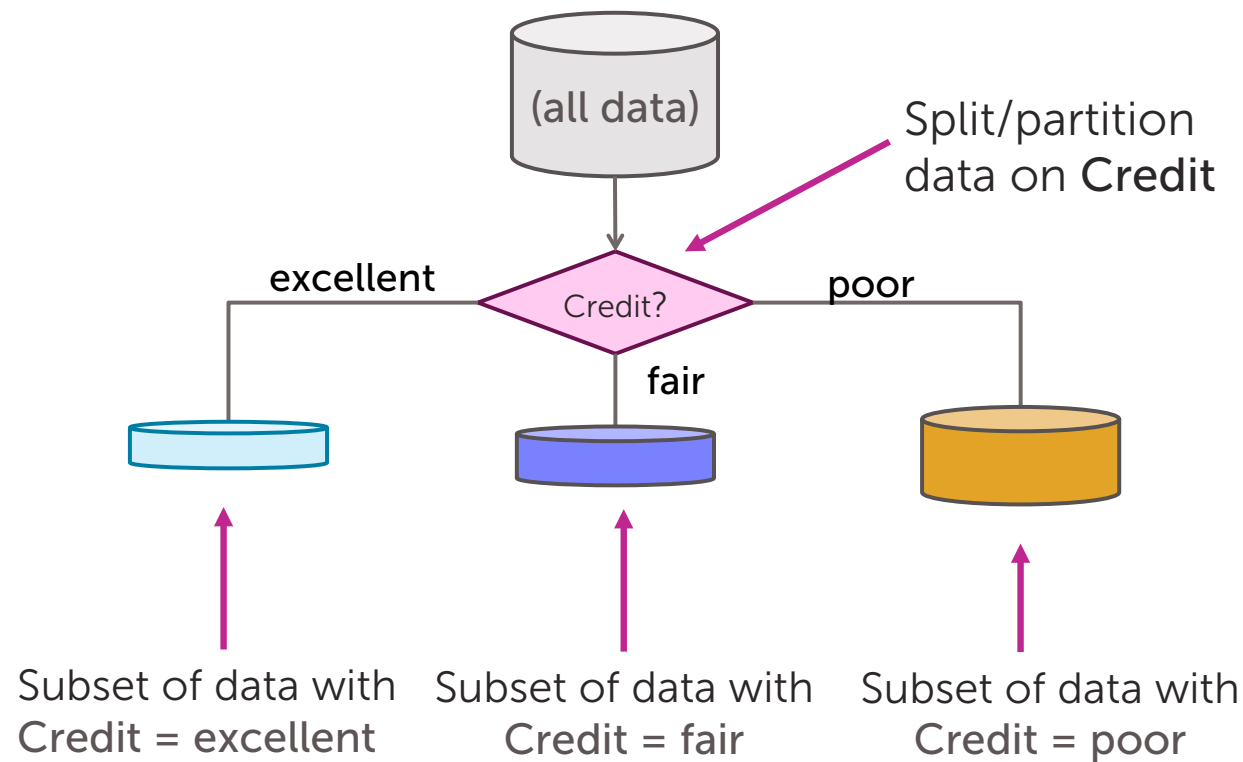
# Greedy decision tree learning:
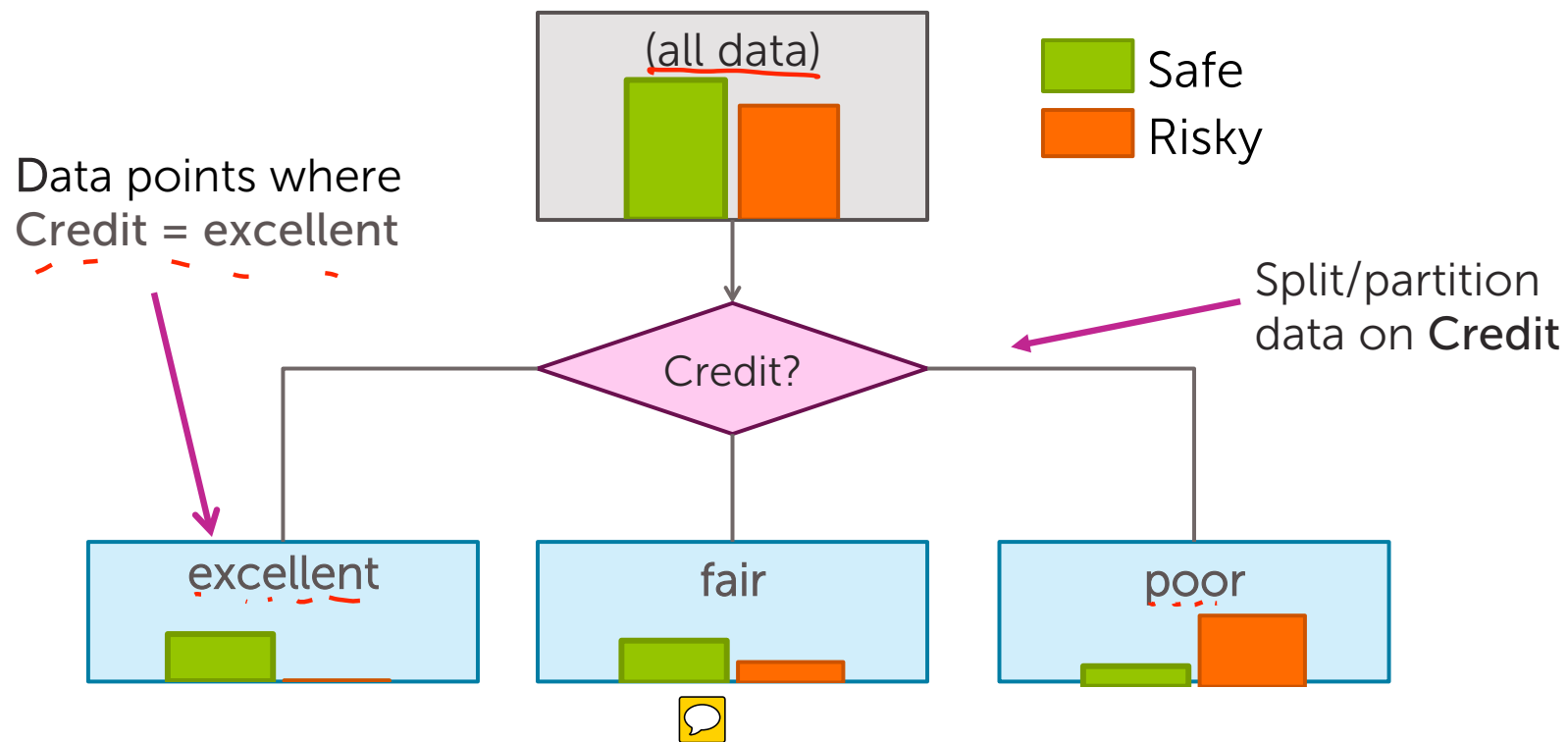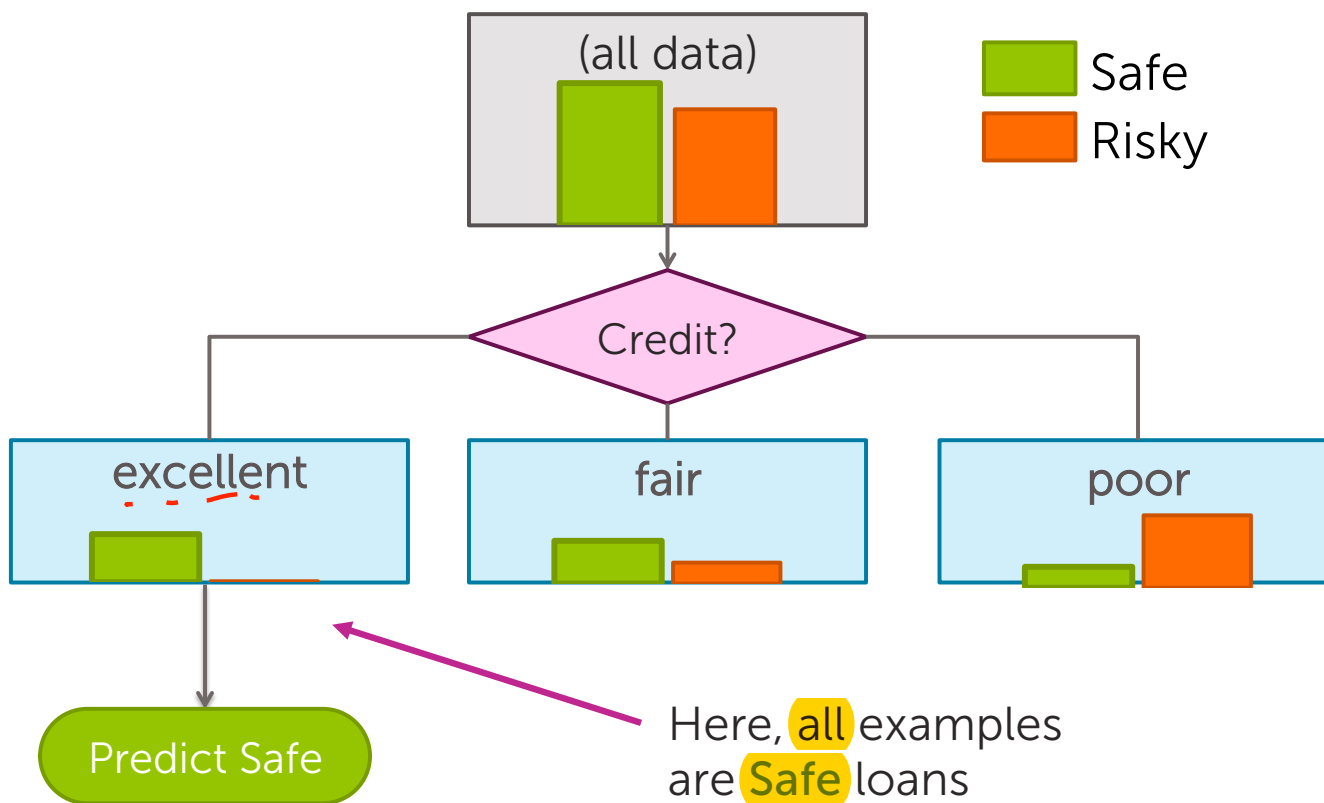## *Algorithm outline*

# Step 1: Start with an empty tree



(all data)

Safe
Risky

Histogram
of y values

All points in the
dataset

Machine Learning Specialization

# Step 2: Split on a feature



(all data)

Split/partition data on **Credit**

excellent    Credit?    poor

fair

Subset of data with **Credit = excellent**

Subset of data with **Credit = fair**

Subset of data with **Credit = poor**

Machine Learning Specialization

# Feature split explained

# Step 3: Making predictions



Here, **all** examples are **Safe** loans

©2015-2016 Emily Fox & Carlos Guestrin

Machine Learning Specialization

# Step 4: Recursion

Machine Learning Specialization

# 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

**Problem 1:** Feature split selection

**Problem 2:** Stopping condition

Recursion

Machine Learning Specialization

# Feature split learning

## =

# Decision stump learning

# Start with the data

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 |

Machine Learning Specialization

# Start with all the data

Loan status:    Safe    Risky



(all data)

22

18

Number of **Risky** loans

Number of **Safe** loans

N = 40 examples

Machine Learning Specialization

# Compact visual notation: Root node

Loan status:    Safe    Risky



Number of **risky** loans

Number of **safe** loans

N = 40 examples

Machine Learning Specialization

# Decision stump: Single level tree

Loan status:
Safe Risky



(all data)

Split on **Credit**

Credit?

excellent | fair | poor

excellent
9 | 0

fair
9 | 4

4 | 14

Subset of data with
Credit = excellent

Subset of data with
Credit = fair

Subset of data with
Credit = poor

Machine Learning Specialization

# Visual Notation: Intermediate nodes

Loan status:
Safe Risky

Root
22  18

Credit?

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

Intermediate nodes

Machine Learning Specialization

# Making predictions with a decision stump

Loan status:
Safe  Risky



```
                    root
                   22  18
                     |
                     v
                 < credit? >
          /          |          \
   excellent        fair         poor
     9   0          9   4        4   14
       |             |             |
       v             v             v
     Safe          Safe          Risky
```

For each intermediate node,
set ŷ = **majority value**

Machine Learning Specialization

# 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

Machine Learning Specialization

# How do we select the best feature?

**Better?**

**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

# How do we measure effectiveness of a split?

Loan status:
Safe  Risky

Root
22  18

Credit?

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

excellent
9    0

fair
9    4

poor
4    14

Error =  # mistakes
          # data points

# Calculating classification error

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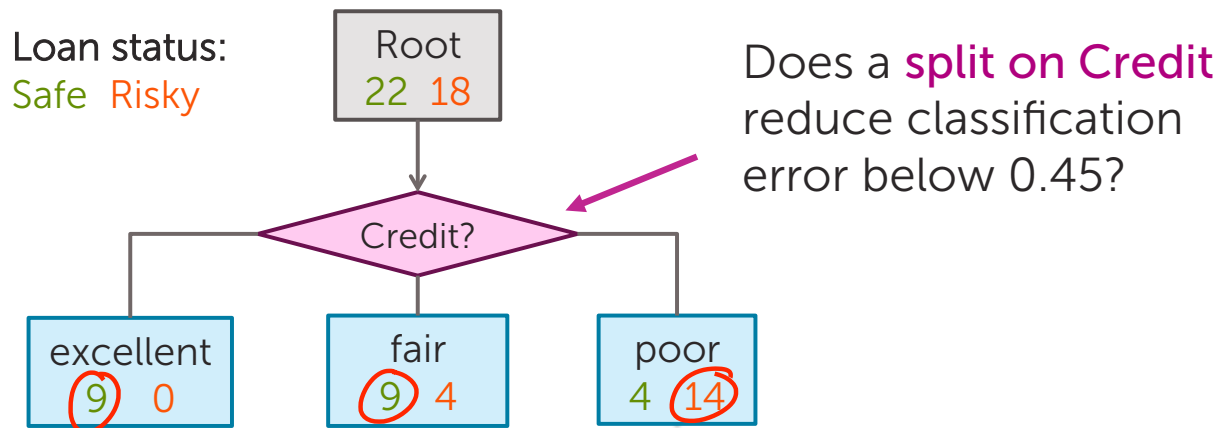- **Step 2:** Calculate classification error of predicting $\hat{y}$ for this data

Loan status:
Safe  Risky

Root
22    18

22 correct

18 mistakes

Safe

$\hat{y}$ = majority class

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

$$= 0.45$$

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

Machine Learning Specialization

# Choice 1: Split on **credit history?**

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

Loan status:
Safe  Risky

Root
22  18

Credit?

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

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

# How good is the split on **Credit**?

**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

**Step 1:** For each intermediate node, set ŷ = **majority value**

Machine Learning Specialization

# 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
```

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

$$= 0.2b$$

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

Machine Learning Specialization

# 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

Machine Learning Specialization

# Evaluating the split on **Term**

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

Loan status:
Safe  Risky

```
         ┌──────────┐
         │   Root   │
         │  22  18  │
         └──────────┘
               │
               ▼
           ◇ Term? ◇
          ╱          ╲
    ┌──────────┐   ┌──────────┐
    │ 3 years  │   │ 5 years  │
    │  16  4   │   │  6   14  │
    └──────────┘   └──────────┘
          │              │
          ▼              ▼
     (  Safe  )     (  Risky  )
          ▲              ▲
          │              │
    4 mistakes      6 mistakes
```

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

$$= 0.25$$

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

# Choice 1 vs Choice 2

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

← First split!

## 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 $h_i(x)$: ← *credit, term, income*

  1. Split data of $M$ according to feature $h_i(x)$

  2. Compute classification error split

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

# 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**

Machine Learning Specialization

# Decision Tree Learning:
*Recursion & Stopping conditions*

# Learn decision tree from data?

| 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 |

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

Loan status:
Safe Risky

Root
22    18

Credit?

excellent
9    0

fair
9    4

poor
4    14

Safe

Leaf node

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

# 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

©2015-2016 Emily Fox & Carlos Guestrin

Machine Learning Specialization

# Final decision tree



Loan status:
Safe Risky

©2015-2016 Emily Fox & Carlos Guestrin

Machine Learning Specialization

# 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???

# 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

Machine Learning Specialization

# 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

# Predictions with decision trees

Machine Learning Specialization

# Decision tree model



T($\mathbf{x}_i$) = Traverse decision tree

**Input:** $\mathbf{x}_i$

$\hat{y}_i$

Machine Learning Specialization

# Traversing a decision tree

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

©2015-2016 Emily Fox & Carlos Guestrin

Machine Learning Specialization

# Decision tree prediction algorithm

predict(tree_node, input)

- **If** current **tree_node** is a leaf:
  - o **return** majority class of data points in leaf

- **else:**
  - o next_note = child node of tree_node whose feature value agrees with **input**
  - o **return** predict(next_note, input)

# Multiclass classification & predicting probabilities

# Multiclass prediction

Input: $\mathbf{x}_i$

**Loan Application**

**Classifier MODEL**

Output: $\hat{y}_i$
Predicted class

Safe

Risky

Danger

# Multiclass decision stump

N = 40,
1 feature,
3 classes

| Credit | y |
|---|---|
| excellent | safe |
| fair | risky |
| fair | safe |
| poor | danger |
| excellent | risky |
| fair | safe |
| poor | danger |
| poor | safe |
| fair | safe |
| ... | ... |

Loan status:
Safe Risky Danger

Root
18  12  10

Credit?

excellent
9  2  1

fair
6  9  2

poor
3  1  7

Safe

Risky

Danger

Machine Learning Specialization

# Predicting probabilities with decision trees

Loan status:
Safe Risky Danger

Root
18  12  10

Credit?

excellent
9  2  1

fair
6  9  2

poor
3  1  7

Safe

Risky

Danger

$y =$

$$P(y = \text{danger} \mid \mathbf{x})$$

$$= \frac{7}{3 + 1 + 7} = 0.64$$

Machine Learning Specialization

# 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 |

Machine Learning Specialization

# Split on each numeric value?

Danger: May only contain one data point per node

Loan status: Safe Risky

Root
22   18

Income?

| $30K | $31.4K | $39.5K | | $61.1K | $91.3K |
|------|--------|--------|--|--------|--------|
| 0  1 | 1  0 | 0  1 | | 0  1 | 0  1 |

Can't trust prediction (overfitting)

# Alternative: 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

Many data points ➔
lower chance of overfitting

# Threshold splits in 1-D

Threshold split is the line
Income = $60K

Income <  $60K     Income >= $60K

Safe ○
Risky ○

Income

$10K                              $120K

Machine Learning Specialization

# Visualizing the threshold split

Threshold split is
the line **Age = 38**

# Split on Age >= 38



Income

age < 38    age >= 38

Predict **Risky**

... 

$80K

Predict **Safe**

$40K

$0K

0   10   20   30   40   ...   Age

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



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

©2015-2016 Emily Fox & Carlos Guestrin

Machine Learning Specialization

# Each split partitions the 2-D space

©2015-2016 Emily Fox & Carlos Guestrin

Machine Learning Specialization

# Finding the best threshold split

OPTIONAL

# Finding the best threshold split

**Infinite possible values of t**

Income = t*

Income < t*          Income >= t*

Safe ◯
Risky ◯

Income

$10K          $120K

Machine Learning Specialization

# 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

# Only need to consider mid-points

Finite number of
splits to consider

Income

Safe ○
Risky ○

$10K                                    $120K

Machine Learning Specialization

# Threshold split selection algorithm

*Income*

- 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 … 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 $t^*$ with the lowest classification error

Machine Learning Specialization

# Decision trees vs logistic regression:
## *Example*

# Logistic regression

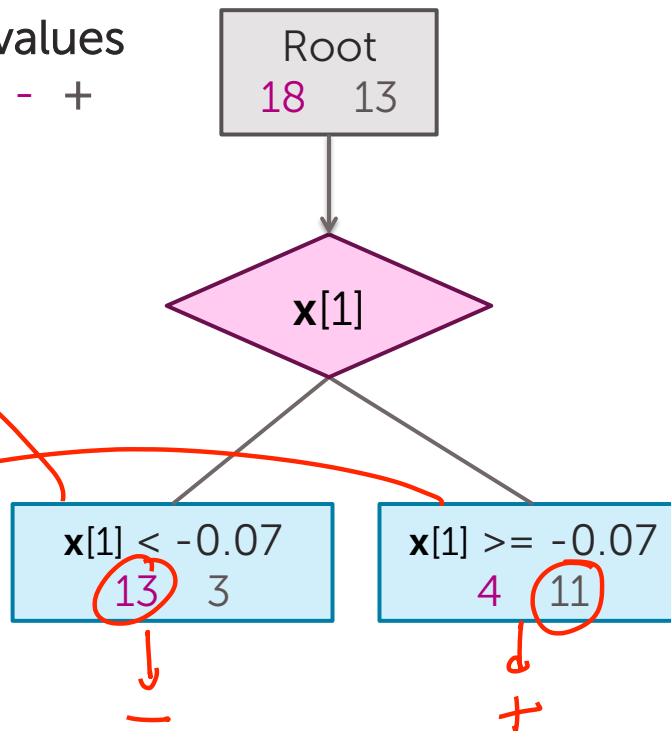| Feature | Value | Weight Learned |
|---------|-------|----------------|
| $h_0(\mathbf{x})$ | 1 | 0.22 |
| $h_1(\mathbf{x})$ | $\mathbf{x}[1]$ | 1.12 |
| $h_2(\mathbf{x})$ | $\mathbf{x}[2]$ | -1.07 |



data

Machine Learning Specialization

# 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

−0.07

Machine Learning Specialization

# Depth 2



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
1   0

**x**[1] >= -1.66
6   3

**x**[2] < 1.55
1   11

**x**[2] >= 1.55
3   0

©2015-2016 Emily Fox & Carlos Guestrin

Machine Learning Specialization

# Threshold split caveat

y values
−  +

Root
18   13

For threshold splits, same feature can be used multiple times

**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

Machine Learning Specialization

# Decision boundaries



Depth 1          Depth 2          Depth 10

©2015-2016 Emily Fox & Carlos Guestrin                    Machine Learning Specialization

# Comparing decision boundaries

## Decision Tree



Depth 1          Depth 3          Depth 10

## Logistic Regression



Degree 1 features      Degree 2 features      Degree 6 features

Machine Learning Specialization

# 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
  - Probability predictions
  - Multiclass classification

Machine Learning Specialization

# Thank you to Dr. Krishna Sridhar



Dr. Krishna Sridhar

Staff Data Scientist, Dato, Inc.

Machine Learning Specialization