# CSE/STAT 416

## Naïve Bayes and Decision Trees

**Tanmay Shah**
**Paul G. Allen School of Computer Science & Engineering**
**University of Washington**
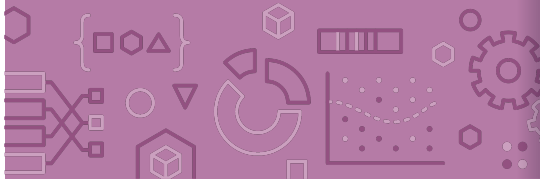
**April 24, 2024**

❓ **Questions?** Raise hand or **sli.do #cs416**
💬 **Before Class:** Pro-rain or anti-rain person?
🎵 **Listening to:** lecture

# Administrivia

- Midterm due tonight
- - Post questions on Edstem (Private post as needed)

- HW3 out Friday

# Probability Classifier

**Idea**: Estimate probabilities $\hat{P}(y|x)$ and use those for prediction

**Probability Classifier**

Input $x$: Sentence from review

Estimate class probability $\hat{P}(y = +1|x)$
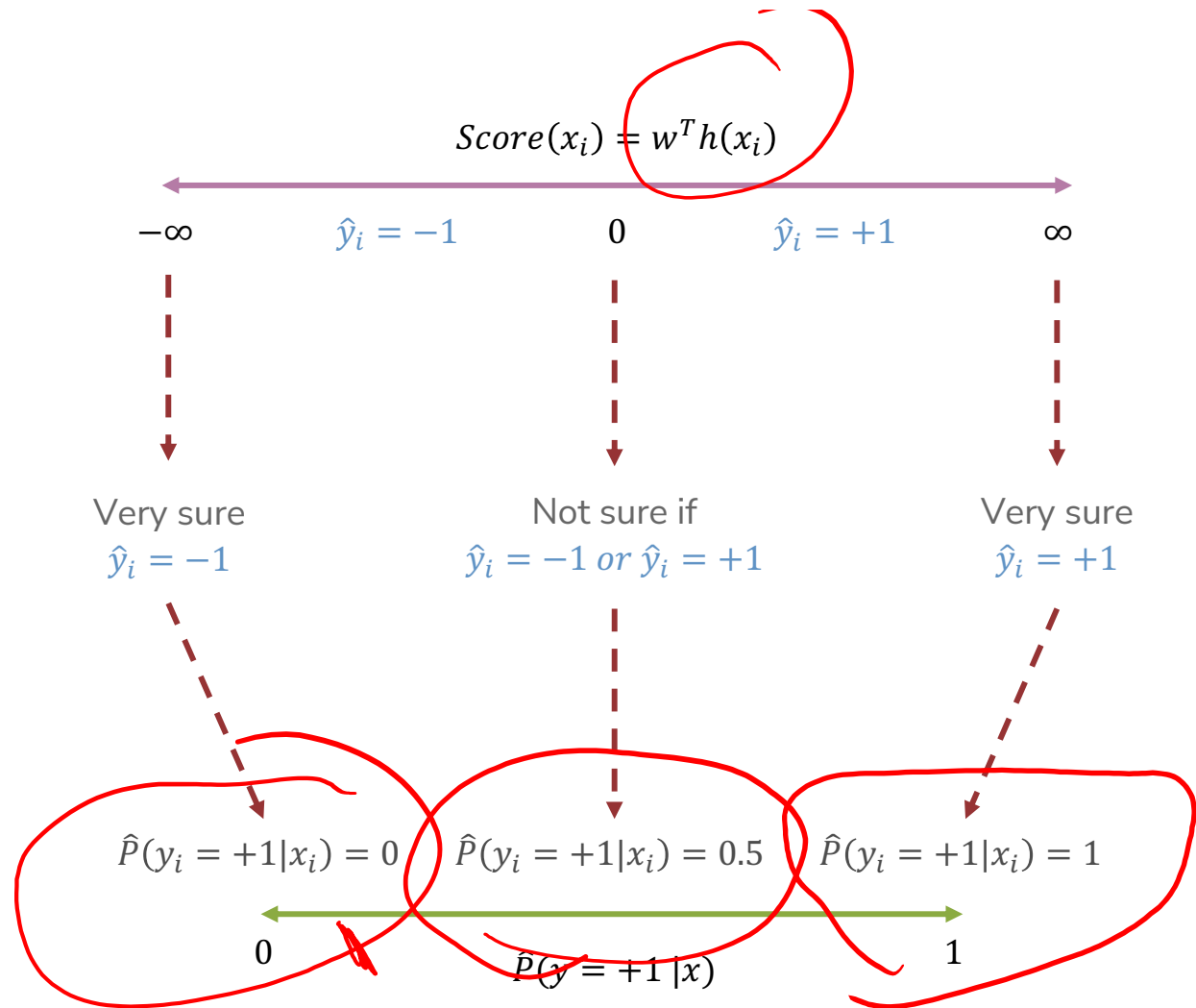
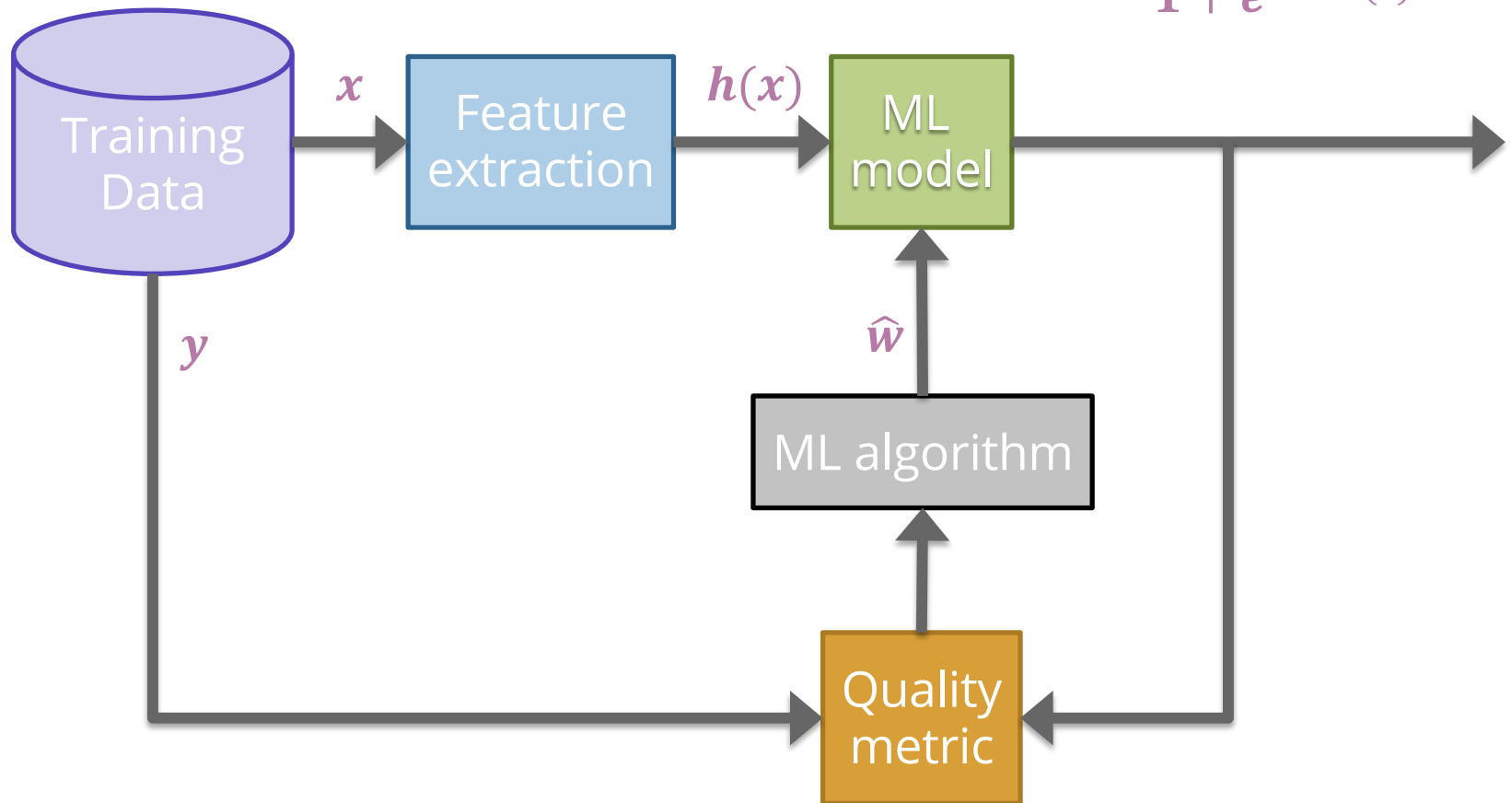If $\hat{P}(y = +1|x) > 0.5$:

- $\hat{y} = +1$

Else:

- $\hat{y} = -1$

**Notes**:

Estimating the probability improves **interpretability**

# Interpreting Score

$$Score(x_i) = w^T h(x_i)$$

$-\infty$  $\hat{y}_i = -1$  $0$  $\hat{y}_i = +1$  $\infty$

Very sure
$\hat{y}_i = -1$

Not sure if
$\hat{y}_i = -1 \; or \; \hat{y}_i = +1$

Very sure
$\hat{y}_i = +1$

$\hat{P}(y_i = +1|x_i) = 0$   $\hat{P}(y_i = +1|x_i) = 0.5$   $\hat{P}(y_i = +1|x_i) = 1$

$0$   $\hat{P}(y = +1|x)$   $1$

4

$$\widehat{P}(y = +1|x, \widehat{w}) = sigmoid\left(\widehat{w}^T h(x)\right) = \frac{1}{1 + e^{-\widehat{w}^T h(x)}}$$

# Naïve Bayes

# Idea: Naïve Bayes

$x = $ *"The sushi & everything else was awesome!"*

$P(y = +1 \mid x = $ *"The sushi & everything else was awesome!"*$)$?

$P(y = -1 \mid x = $ *"The sushi & everything else was awesome!"*$)$?

**Idea:** Select the class that is the most likely!

**Bayes Rule:**

$$P(y = +1|x) = \frac{P(x|y = +1)P(y = +1)}{P(x)}$$

Example

$$\frac{P(\text{"The sushi \&everything else was awesome!"} \mid y = +1)\, P(y = +1)}{P(\text{"The sushi \& everything else was awesome!"})}$$

Since we're just trying to find out which class has the greater probability, we can discard the divisor.

$$P(A \mid B) = \frac{P(B \mid A) \cdot P(B)}{P(A)}$$

7

# Naïve Assumption

**Idea**: Select the class with the highest probability!

**Problem**: We have not seen the sentence before.

**Assumption**: Words are independent from each other.

$x = $ *"The sushi & everything else was awesome!"*

$$\frac{P(\text{"The sushi \& everything else was awesome!"}|y = +1)\, P(y = +1)}{P(\text{"The sushi \& everything else was awesome!"})}$$

$P(\text{"The sushi \& everything else was awesome!"} | y = +1)$
$= P(\text{The} | y=+1) * P(sushi | y = +1) * P(\& | y = +1)$
$\quad * P(everything | y = +1) * P(else | y = +1) * P(was | y = +1)$
$\quad * P(awesome | y = +1)$

*Handwritten annotation:*

$P(A \cap B)$
$= P(A) \cdot P(B)$
$A \perp B$

## Compute Probabilities

How do we compute something like

$P(y = +1)?$

$$\frac{\#\ +ve\ reviews}{\#\ reviews}$$

How do we compute something like

$P(\text{``awesome''} \mid y = +1)?$

$$\frac{\#\ occurrences\ of\ \text{``awesome''}\ +ve\ reviews}{\#\ total\ words}$$

# Zeros

If a feature is missing in a class everything becomes zero.

$$P(\text{"The sushi \&everything else was awesome!"} \,|\, y = +1)$$
$$= P(\text{The} \,|\, y=+1) * P(sushi \,|\, y = +1) * P(\& | y = +1)$$
$$* P(everything | y = +1) * P(else | y = +1) * P(was | y = +1)$$
$$* P(awesome | y = +1)$$

Solutions?

    Take the log (product becomes a sum).

        -   Generally define $\log(0) = 0$ in these contexts

    Laplacian Smoothing (adding a constant to avoid multiplying by zero)

# Compare Models

**Logistic Regression:**

$$P(y = +1|x, w) = \frac{1}{1 + e^{-w^T h(x)}}$$

**Naïve Bayes:**

$$P(y|x_1, x_2, \ldots, x_d) = \prod_{j=1}^{d} P(x_j|y) \ P(y)$$

Based on counts of words/classes
- Laplace Smoothing

# Compare Models

**Generative:** defines a model for generating x (e.g. Naïve Bayes)

**Discriminative:** only cares about defining and optimizing a decision boundary (e.g. Logistic Regression)

**Recap**: What is the predicted class for this sentence assuming we have the following training set (no Laplace Smoothing). "he is not cool"

$$P(x \mid -1) = 0$$

$$P(\text{"he is not cool"} \mid +1) \Rightarrow +1$$

$$= P(he \mid +1) \cdot P(is \mid +1) \cdot P(not \mid +1)$$
$$\cdot P(cool \mid +1) \cdot P(y = +1) \to 213$$

$$= 2/11 \cdot 3/11 \cdot 1/11$$
$$\cdot \frac{1}{11} \cdot 2/3$$

$$\frac{12}{11^4 \cdot 3}$$

| Sentence | Label |
|---|---|
| this dog is cute | Positive |
| he does not like dogs | Negative |
| he is not bad he is cool | Positive |

13

# Decision Trees

Humans often make decisions based on **Flow Charts** or **Decision Trees**

# Parametric vs. Non-Parametric Methods

**Parametric Methods**: make assumptions about the data distribution

- Linear Regression ⇒ assume the data is linear
- Logistic Regression ⇒ assume probability has the shape of of a logistic curve and linear decision boundary
- Those assumptions result in a _parameterized_ function family. Our learning task is to learn the parameters.

**Non-Parametric Methods**: (mostly) don't make assumptions about the data distribution

- Decision Trees, k-NN (soon)
- We're still learning something, but not the parameters to a function family that we're assuming describes the data.
- Useful when you don't want to (or can't) make assumptions about the data distribution.

# XOR

A xor B

A line might not always support our decisions.

(0,1)

(1,1)

(1,0)

# Credit history explained

Did I pay previous loans on time?

Example: excellent, good, or fair

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**
★★★

# Intelligent application

**Loan Applications**

Intelligent loan application review system

Safe ✓

Risky ✗

Risky ✗

# Classifier review

Loan Application

Input: $\mathbf{x}_i$

Classifier MODEL

Output: $\hat{y}$
Predicted class

$\hat{y}_i = +1$

Safe

Risky

$\hat{y}_i = -1$

# Setup

Data (N observations, 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 | safe |
| fair | 5 yrs | low | safe |
| poor | 3 yrs | high | risky |
| poor | 5 yrs | low | safe |
| fair | 3 yrs | high | safe |

Evaluation: classification error

Many possible decisions: number of trees grows exponentially!

## Poll Everywhere

### Think  &

2 min

No poll

With our discussion of bias and fairness from last week, discuss the potential biases and fairness concerns that might be present in our dataset about loan safety.

# Decision Trees

node/ branch root

leaf



- **Branch/Internal node:** splits into possible values of a feature
- **Leaf node:** final decision (the class value)

## Brain Break

# Growing
Trees

# Visual Notation

Loan status:   Safe   Risky

Root
6   3

# of Risky loans

# of Safe loans

N = 9 examples

# Decision stump: 1 level

| 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 | safe |
| fair | 5 yrs | low | safe |
| poor | 3 yrs | high | risky |
| poor | 5 yrs | low | safe |
| fair | 3 yrs | high | safe |

Loan status: Safe Risky

Split on Credit

Root
6   3

Credit?

excellent
2   0

fair
3   1

poor
1   2

Subset of data with Credit = excellent

Subset of data with Credit = fair

Subset of data with Credit = poor

# Making predictions

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

Loan status:
Safe   Risky

Root
6   3

credit?

excellent
2   0

fair
3   1

poor
1   2

Safe

Safe

Risky

# How do we select the best feature?

- Select the split with lowest classification error

## Choice 1: Split on Credit

Loan status:
Safe  Risky

Root
6  3

Credit?

| excellent | fair | poor |
|-----------|------|------|
| 2   0     | 3  1 | 1  2 |

## Choice 2: Split on Term

Loan status:
Safe  Risky

Root
6  3

Term?

| 3 years | 5 years |
|---------|---------|
| 4  1    | 2  1    |

*Safe*          *Safe*

# Calculate the node values.

| 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 | safe |
| fair | 5 yrs | low | safe |
| poor | 3 yrs | high | risky |
| poor | 5 yrs | low | safe |
| fair | 3 yrs | high | safe |

## Choice 2: Split on Term

Loan status:
Safe  Risky

Root
6  3

Term?

3 years
4  1

5 years
2  2

# How do we select the best feature?

Select the split with lowest classification error



**Choice 1: Split on Credit**

Loan status:
Safe  Risky

Root
6  3

Credit?

excellent
2  0

fair
3  1

poor
1  2

**Choice 2: Split on Term**

Loan status:
Safe  Risky

Root
6  3

Term?

3 years
4  1

5 years
2  2

# How do we measure effectiveness of a split?

Loan status:
Safe  Risky

Root
6  3

Credit?

excellent
2  0

fair
3  1

poor
1  2

Safe

Safe

Risky

Idea: Calculate classification error of this decision stump

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
6    3

6 correct                    3 mistakes

Safe

$\hat{y}$ = majority class

Error = $\dfrac{3}{9}$ = $\dfrac{1}{3}$

=

| Tree | Classification error |
|------|----------------------|
| (root) | 0.33 |

37

# Choice 1: Split on Credit history?

Does a split on Credit reduce classification error below 0.33?

Choice 1: Split on Credit

Loan status:
Safe  Risky



Root
6  3

Credit?

excellent
2  0

fair
3  1

poor
1  2

# Split on Credit: Classification error

## Choice 1: Split on Credit

Loan status:
Safe  Risky

```
        Root
        6  3
          |
       Credit?
    /     |      \
excellent  fair   poor
  2   0    3  1   1  2
    |        |       |
  Safe     Safe    Risky
    |        |       |
0 mistakes 1 mistake 1 mistake
```

$$\text{Error} = \frac{2}{9}$$

$$=$$

| Tree | Classification error |
|------|----------------------|
| (root) | 0.33 |
| Split on credit | 0.22 |

39

# Choice 2: Split on Term?

## Choice 2: Split on Term

Loan status:
Safe Risky



40

# Evaluating the split on Term

## Choice 2: Split on Term

Loan status:
Safe  Risky

Root
6   3

Term?

3 years
4   1

5 years
2   2

Safe

Risky

1 mistake

2 mistakes

Error = 3/9

= 1/3

| Tree | Classification error |
|------|---------------------|
| (root) | 0.33 |
| Split on credit | 0.22 |
| Split on term | 0.33 |

Choice 1 vs Choice 2:
Comparing split on credit vs term

| Tree | Classification error |
|------|---------------------|
| (root) | 0.33 |
| split on credit | 0.22 |
| split on loan term | 0.33 |

## Choice 1: Split on Credit

Loan status:
Safe  Risky

Root
6  3

Credit?

excellent
2  0

poor
1  2

**WINNER**

## Choice 2: Split on Term

Loan status:
Safe  Risky

Root
6  3

Term?

3 years
4  1

5 years
2  2

# Split Selection

Split(node)

- o Given $M$, the subset of training data at a node
- o For each (remaining) feature $h_j(x)$ :
  - o Split data $M$ on feature $h_j(x)$
  - o Compute the classification error for the split
- o Chose feature $h_j^*(x)$ with the lowest classification error

# Greedy & Recursive Algorithm

*BuildTree(node)*

- If termination criterion is met:
  - Stop
- Else:
  - Split(node)
  - For child in node:
    - BuildTree(child)

Decision stump:
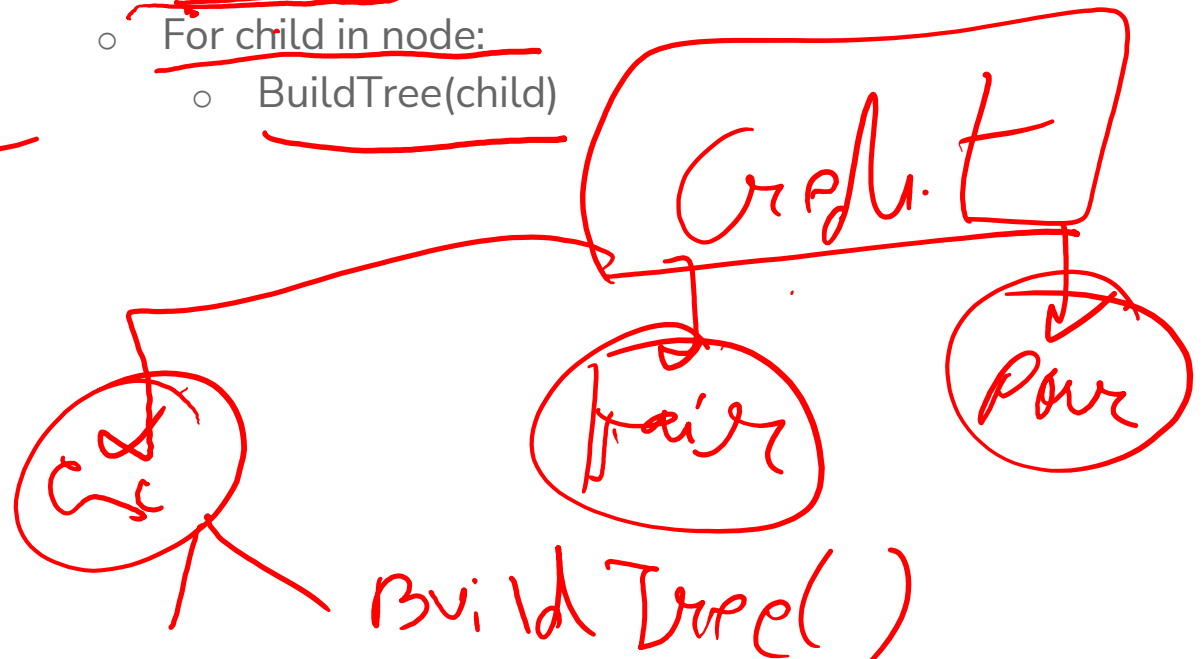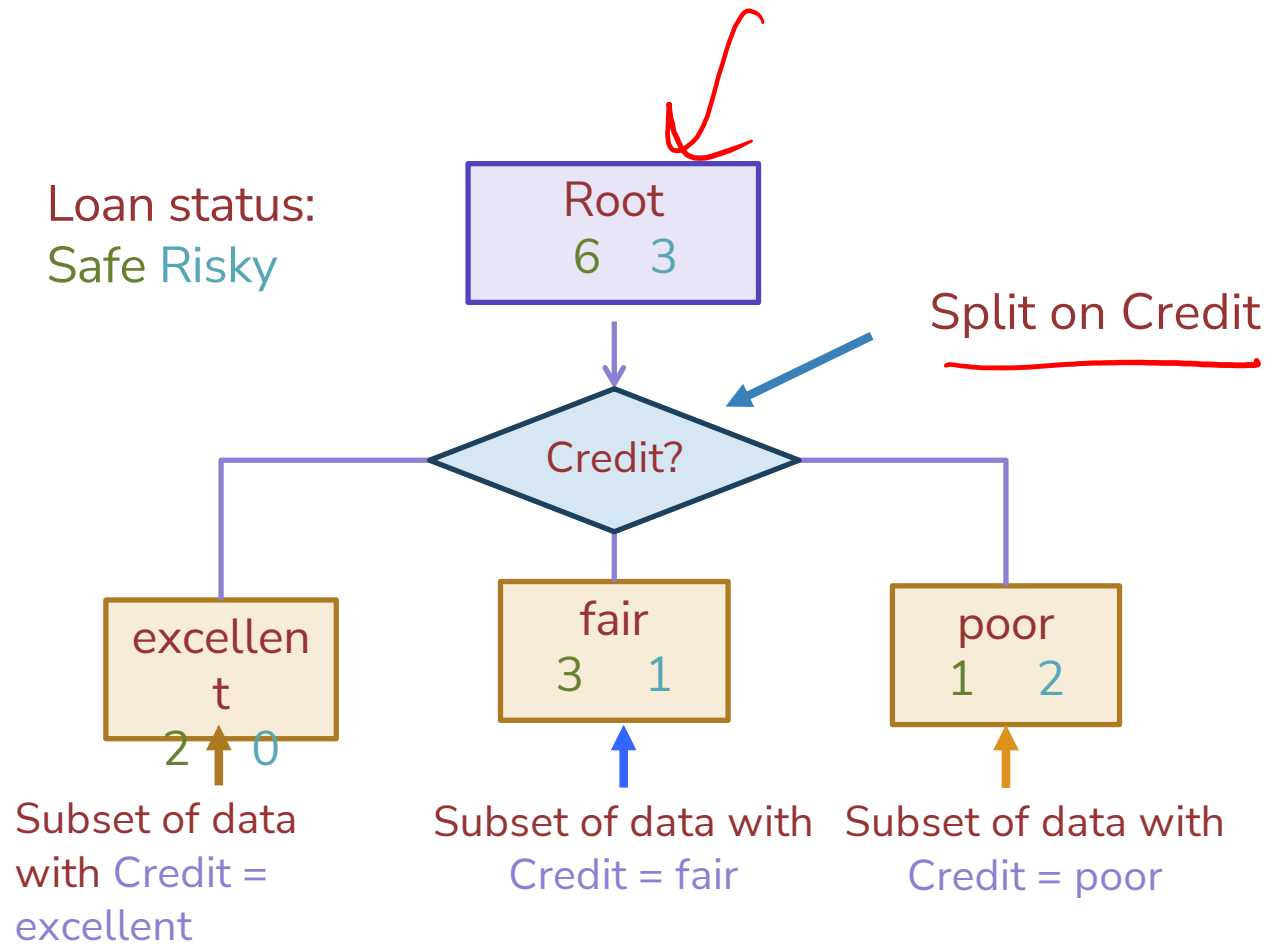1 level

Loan status:
Safe Risky

Root
6 3

Split on Credit

Credit?

excellent
2 0

fair
3 1

poor
1 2

Subset of data with Credit = excellent

Subset of data with Credit = fair

Subset of data with Credit = poor

# Stopping

*For now:* Stop when all points are in one class

**Loan status:**
Safe Risky



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
            6   3
              │
           ◇ Credit? ◇
      ┌───────┼───────┐
  excellent   fair    poor
    2   0     3   1    2   1
      │        │        │
    Safe   Build     Build
           decision  decision stump
           stump     with subset of data
           with      where Credit = poor
           subset
           of data
           where Credit
           = fair
```

# Second level

Loan status:
Safe Risky

slido

Think

1 min

What predictions **should** the below decision tree output for the following datapoints?

Loan status:
Safe Risky

| Credit | Term | Income |
|--------|-------|-----------|
| excellent | 5 yrs | high |
| fair | 3 yrs | low |
| poor | 5 yrs | (missing) |

Root
6   3

Credit?

excellent
2   0

fair
3   1

poor
1   2

Safe

Term?

Income?

3 years
1   0

5 years
2   1

high
0   2

low
1   0

Safe

...

Risky

Safe

*Imputation*

*= Feature*

49

What predictions **should** the below decision tree output for the following datapoints?

Loan status:
Safe Risky

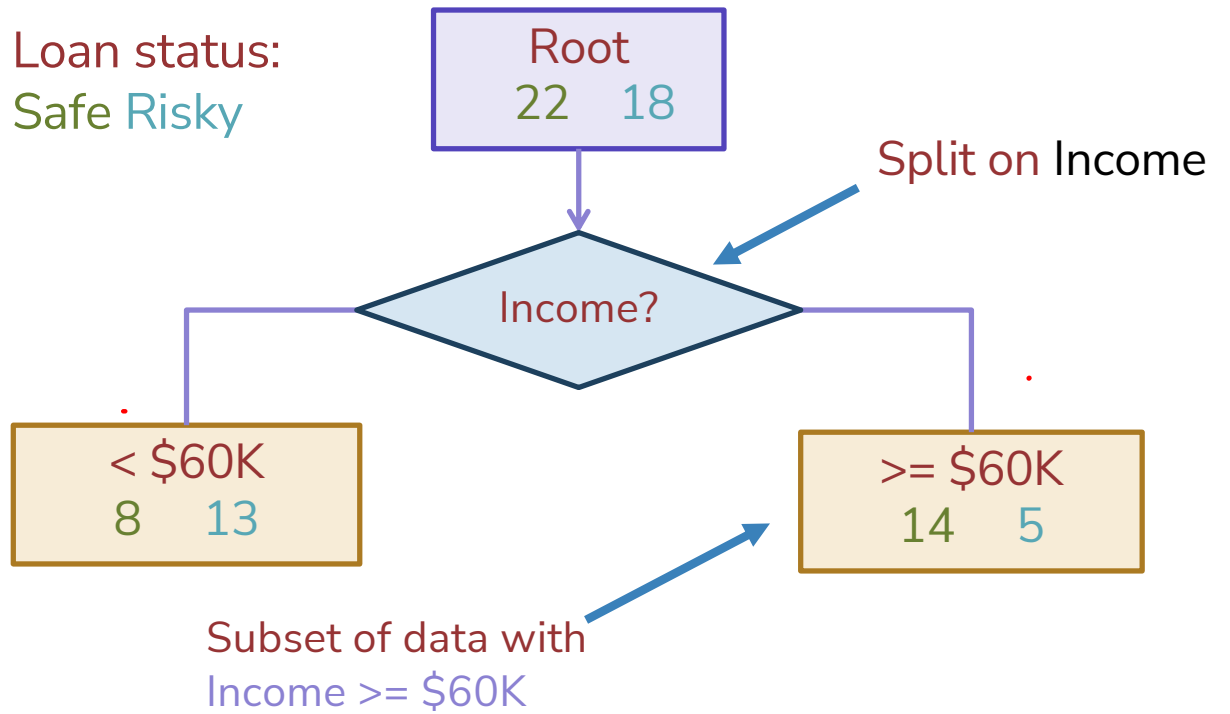| Credit | Term | Income |
|--------|-------|-----------|
| excellent | 5 yrs | high |
| fair | 3 yrs | low |
| poor | 5 yrs | (missing) |

Root
6   3

Credit?

excellent
2   0

fair
3   1

poor
1   2

Safe

Term?

Income?

3 years
1   0

5 years
2   1

high
0   2

low
1   0

Safe

…

Risky

Safe

slido

Group

2 min

**Brain Break**

*Real valued features*

| 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



Split on Income

Subset of data with
Income >= $60K

| Root | |
|------|------|
| 22 | 18 |

Income?

| < $60K | |
|--------|------|
| 8 | 13 |

| >= $60K | |
|---------|------|
| 14 | 5 |

# Best threshold?

*n-1 splits*

Similar to our simple, threshold model when discussing Fairness!
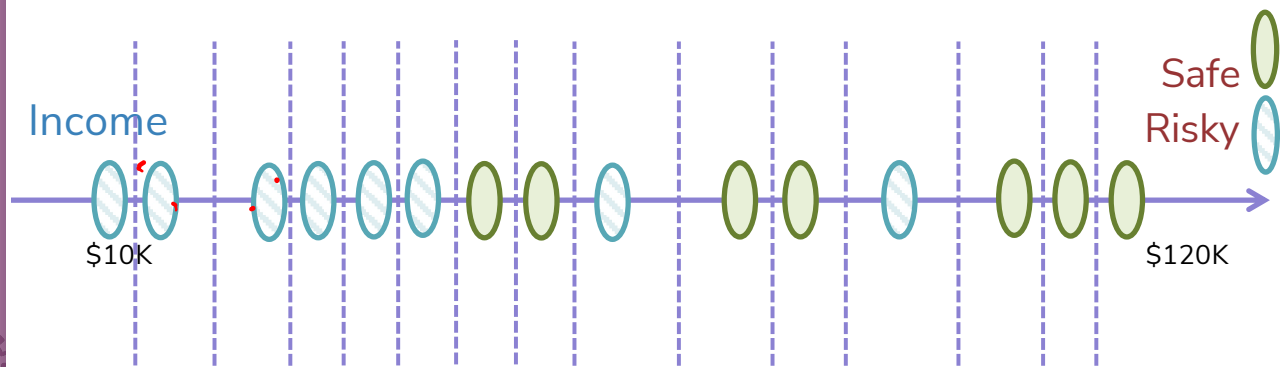
Infinite possible values of $t$

Income $= t^*$

Income $< t^*$

Income $\geq t^*$

Income

Safe

Risky

$10K

$120K

# Threshold between points

Same classification error for any threshold split between $v_a$ and $v_b$

# Only need to consider mid-points

$n - 1$

Finite number of splits to consider

Income

$10K

$120K

Safe

Risky

# Threshold split selection algorithm

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

Let $[v_1, v_2, \ldots, v_N]$ denote sorted values

**Step 2:**

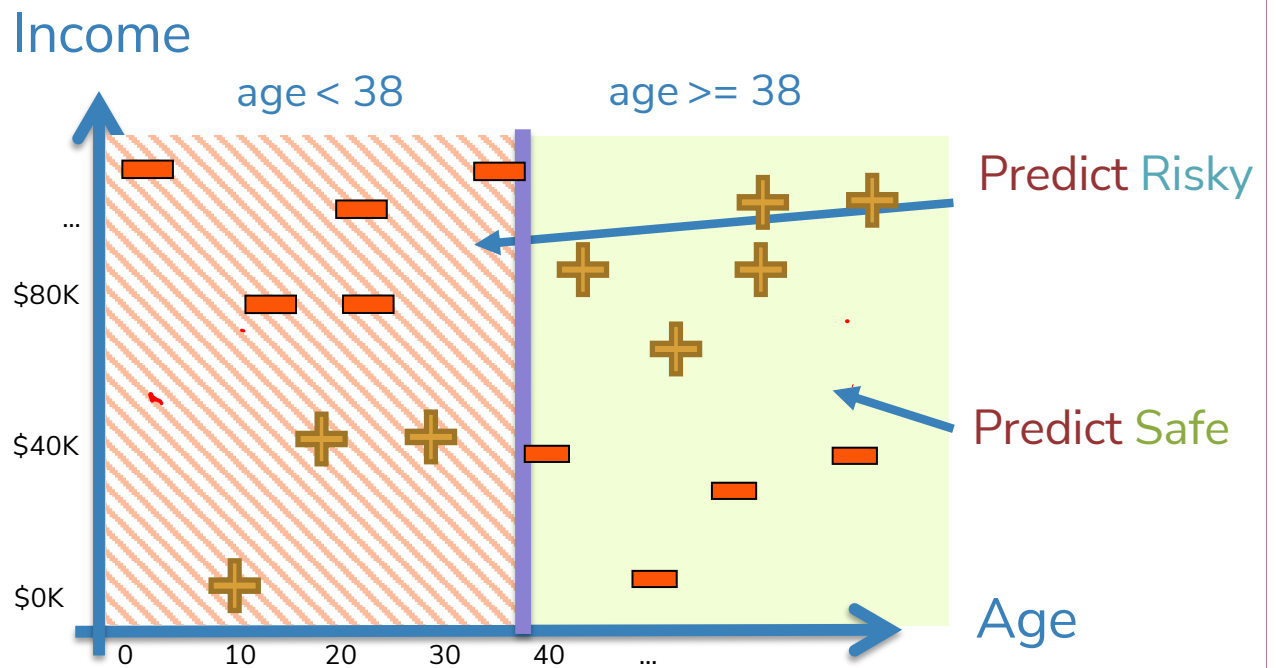- For $i = [1, \ldots, N-1]$

  - Consider split $t_i = \frac{v_i + v_{i+1}}{2}$

  - Compute classification error for threshold split $h_j(x) \geq t_i$
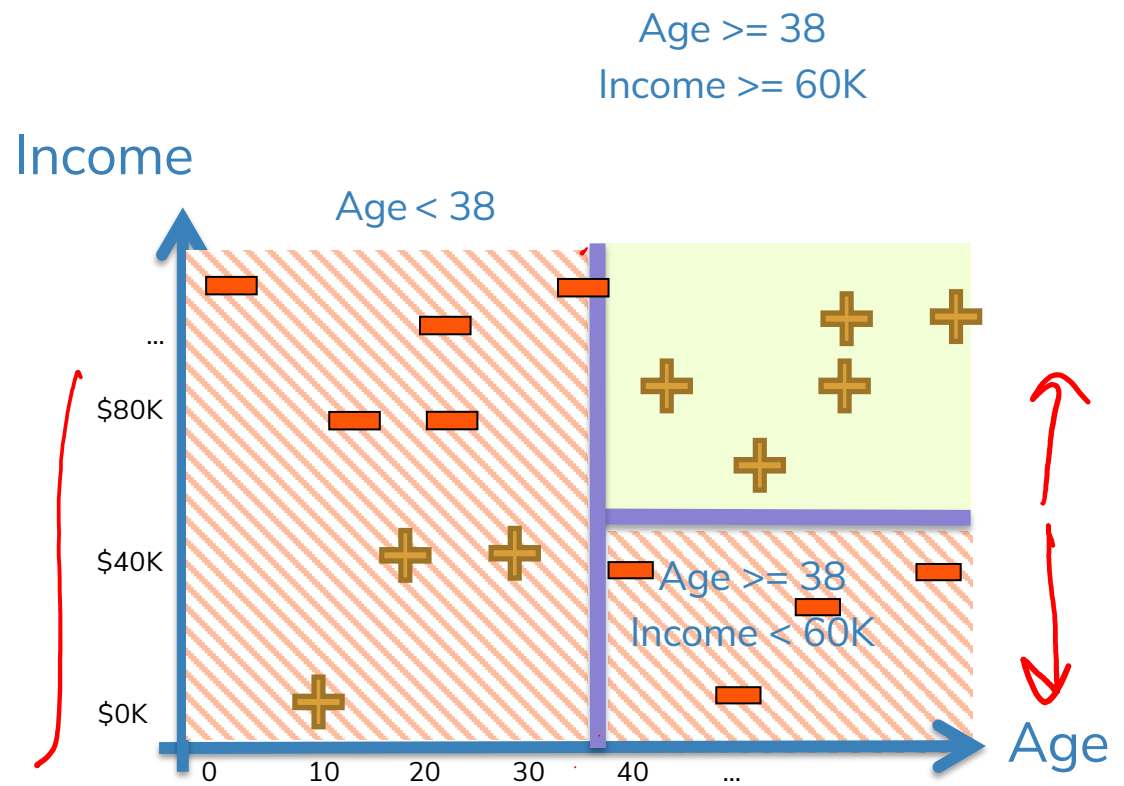
- Chose the $t^*$ with the lowest class. error
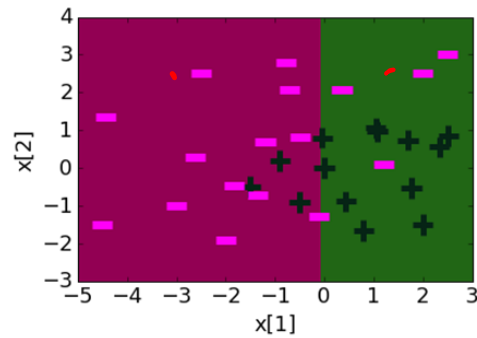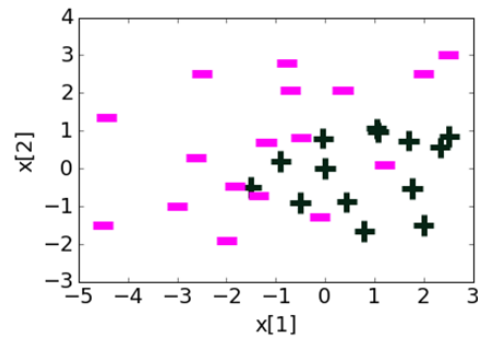
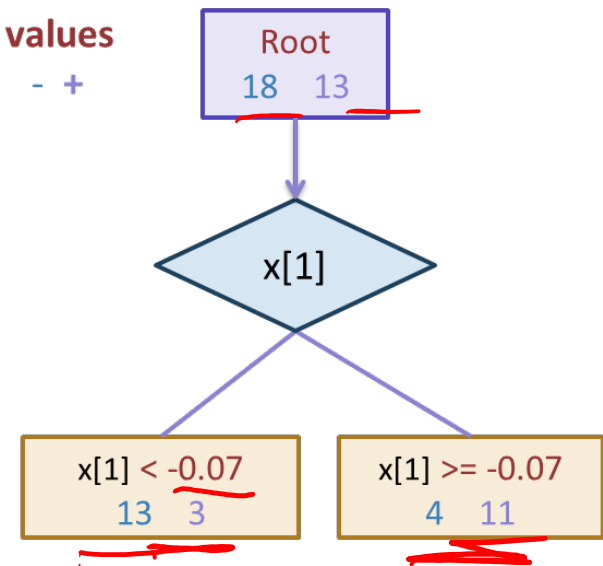Visualizing the threshold split

Each split partitions the 2-D space

# Depth 1:
# Split on x[1]
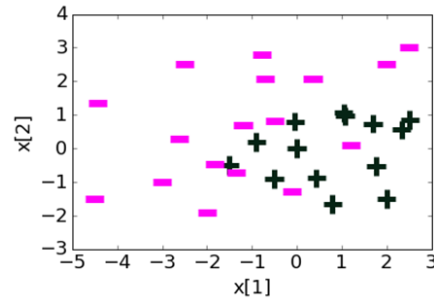


y values

  −  +

Root
18   13

x[1]

x[1] < -0.07
13   3

x[1] >= -0.07
4   11

# 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 |

# 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

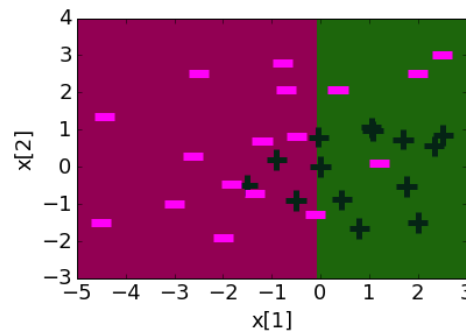x[1] >= -1.66
6   3

x[2] < 1.55
1   11

x[2] >= 1.55
3   0

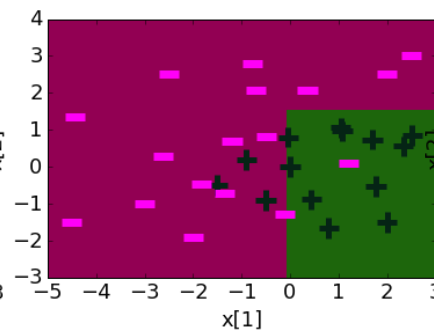# Decision boundaries

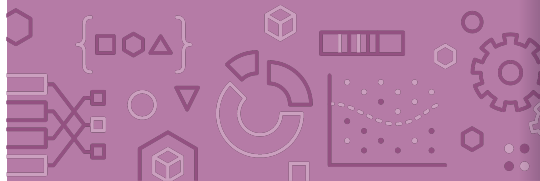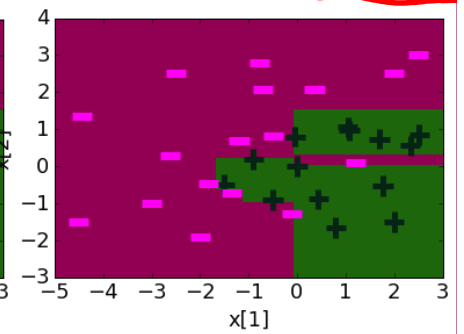Decision boundaries can be complex!



Depth 1          Depth 2          Depth

# Overfitting

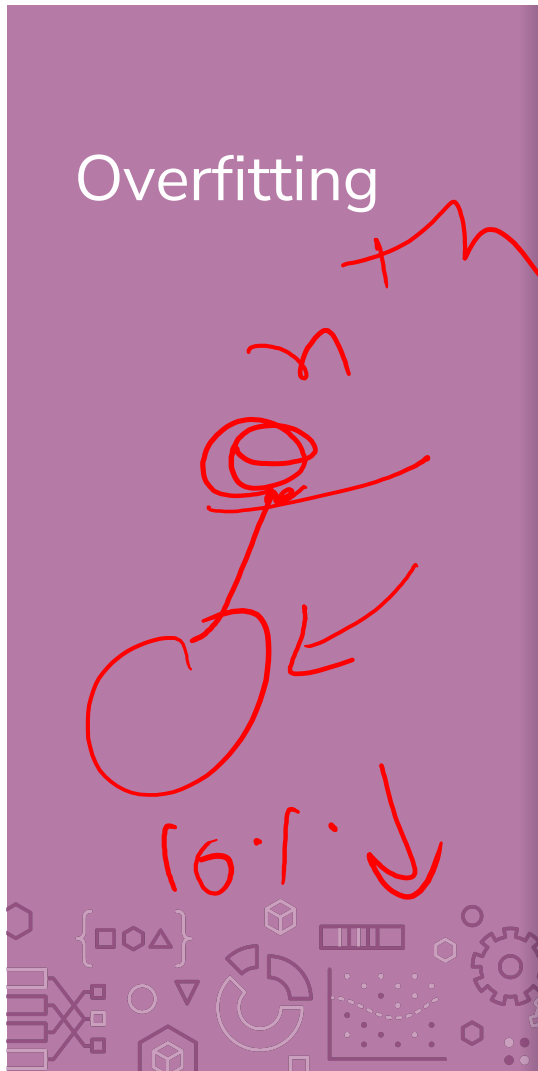Deep decision trees are prone to overfitting
- Decision boundaries are interpretable but not stable
- Small change in the dataset leads to big difference in the outcome

Overcoming Overfitting:
- Stop when tree reaches certain height (e.g., 4 levels)
- Stop when leaf has ≤ some num of points (e.g., 20 pts)
  - Will be the stopping condition for HW
- Stop if split won't significantly decrease error by more than some amount (e.g., 10%)

Other methods include growing full tree and pruning back

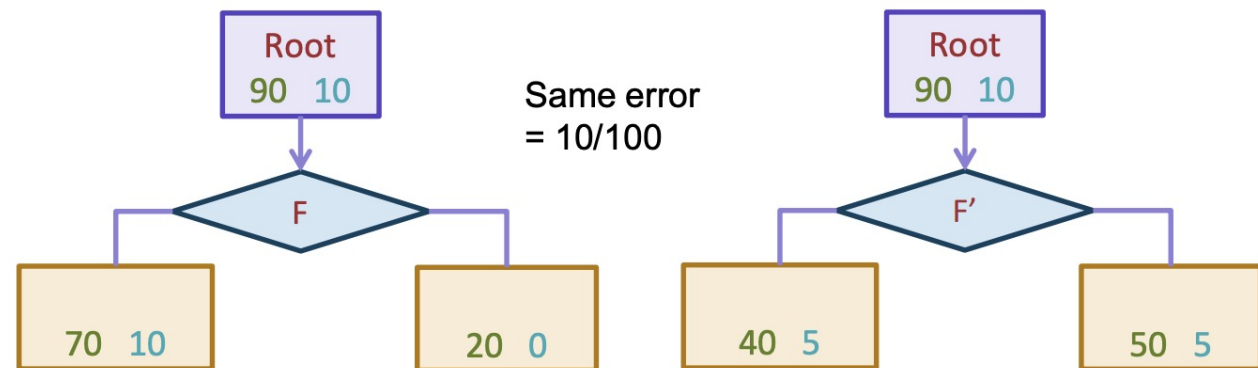Fine-tune hyperparameters with validation set or CV

# In Practice

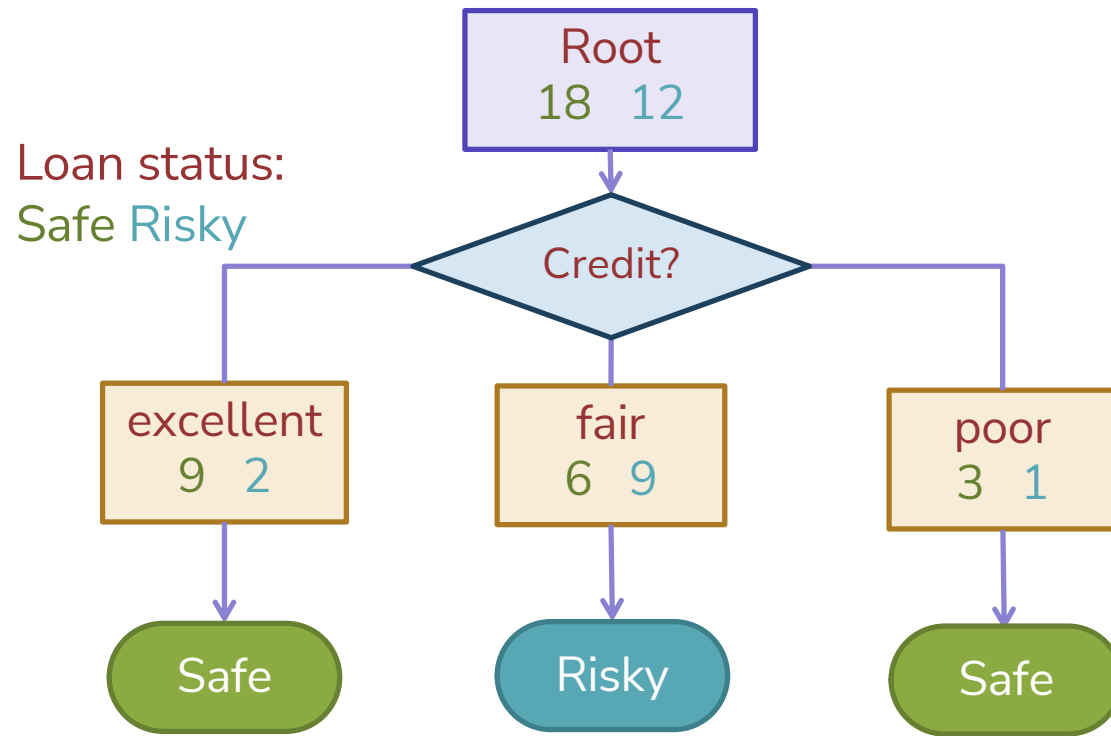Trees can be used for classification or regression (CART)
- Classification: Predict majority class for root node
- Regression: Predict average label for root node

In practice, we don't minimize classification error but instead some more complex metric to measure quality of split such as **Gini Impurity** or **Information Gain** (not covered in 416)



Same error = 10/100

Can also be used to predict probabilities

# Recap

What you can do now:

Define the assumptions and modeling for Naïve Bayes

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