

Machine Learning I

Mohamed Hussien

Decision Trees & Random Forest

Lecture Overview

- Simple Decision Tree
- Mobile Apps Recommendation Example
- Split Using Accuracy
- Split Using Gini Impurity
- Split Using Entropy
- Tree Hyperparameters
- One-hot Encoded Features
- Continues Features
- Decision Tree For Regression
- Random Forest



Wear a jacket?

Wear a jacket?

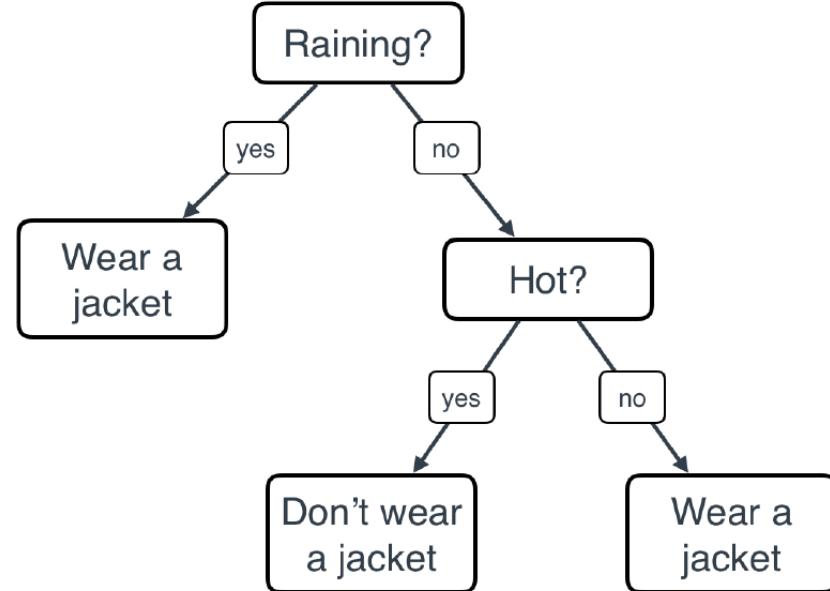
You want to decide if you should wear a jacket today. What does the decision process look like?

- You may look outside and check if it's raining.
If it's raining, then you definitely wear a jacket.
- If it's not, then it could be that it's hot outside, or cold. So then you check the temperature, and if it is hot, then you don't wear a jacket, but if it is cold, then you wear a jacket.

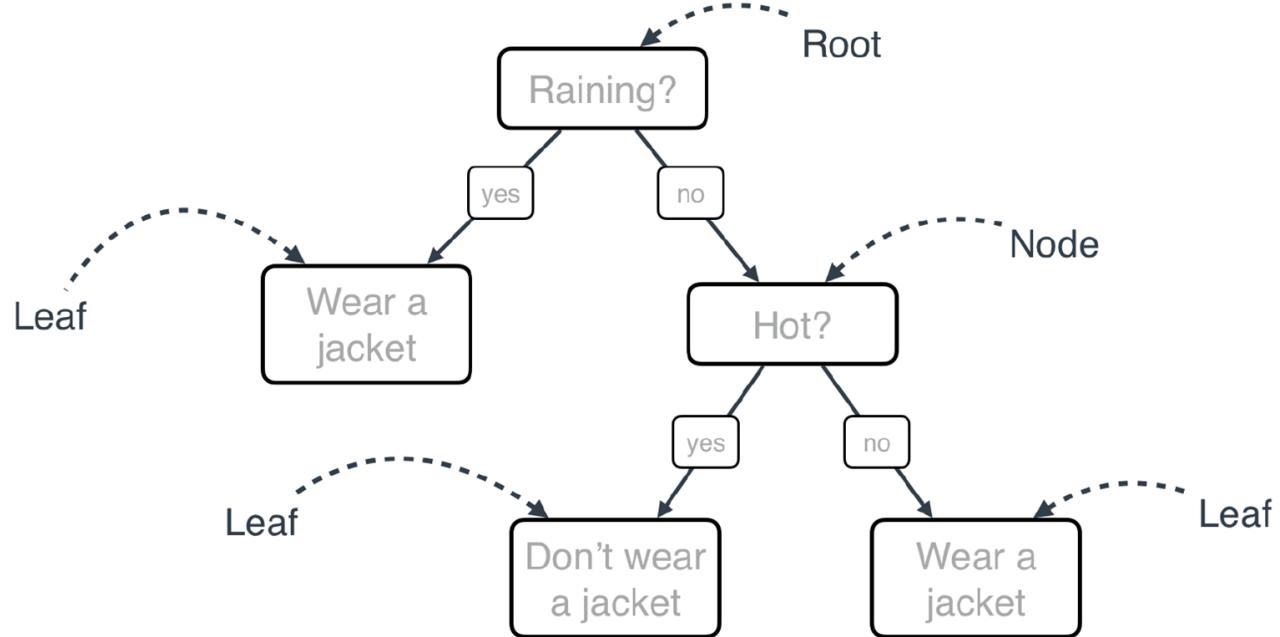
Wear a jacket?

You want to decide if you should wear a jacket today. What does the decision process look like?

- You may look outside and check if it's raining. If it's raining, then you definitely wear a jacket.
- If it's not, then it could be that it's hot outside, or cold. So then you check the temperature, and if it is hot, then you don't wear a jacket, but if it is cold, then you wear a jacket.



Simple Decision Tree



Simple Decision Tree

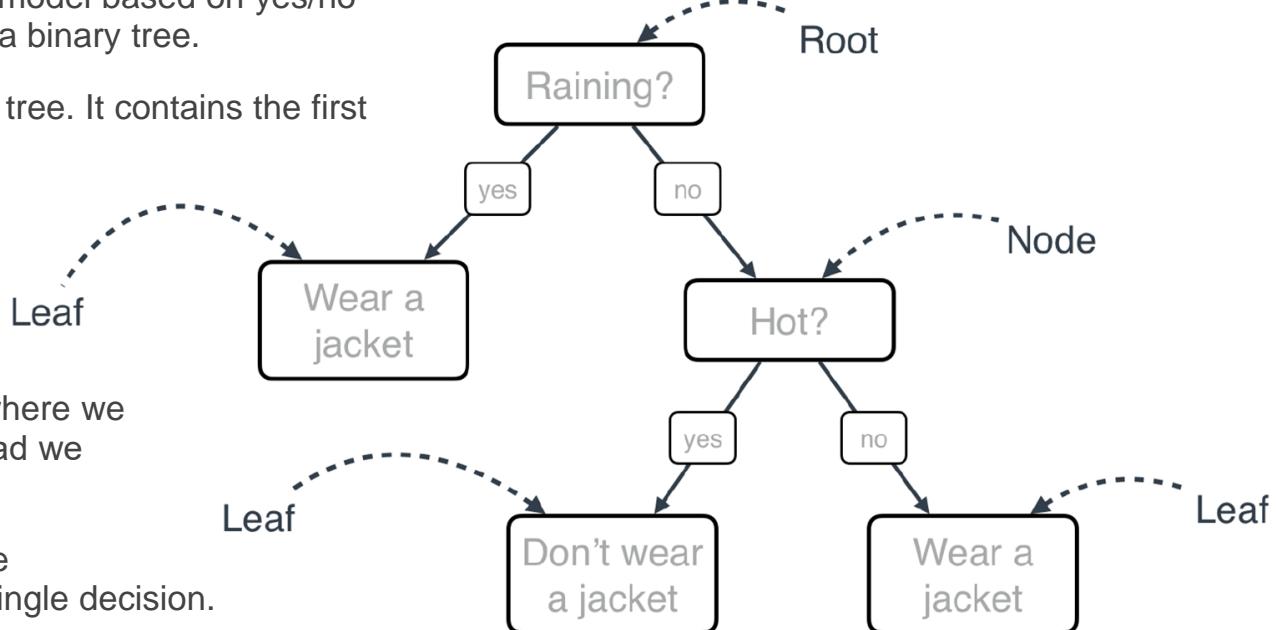
Decision tree A classification model based on yes/no questions and represented by a binary tree.

Root The topmost node of the tree. It contains the first yes/no question.

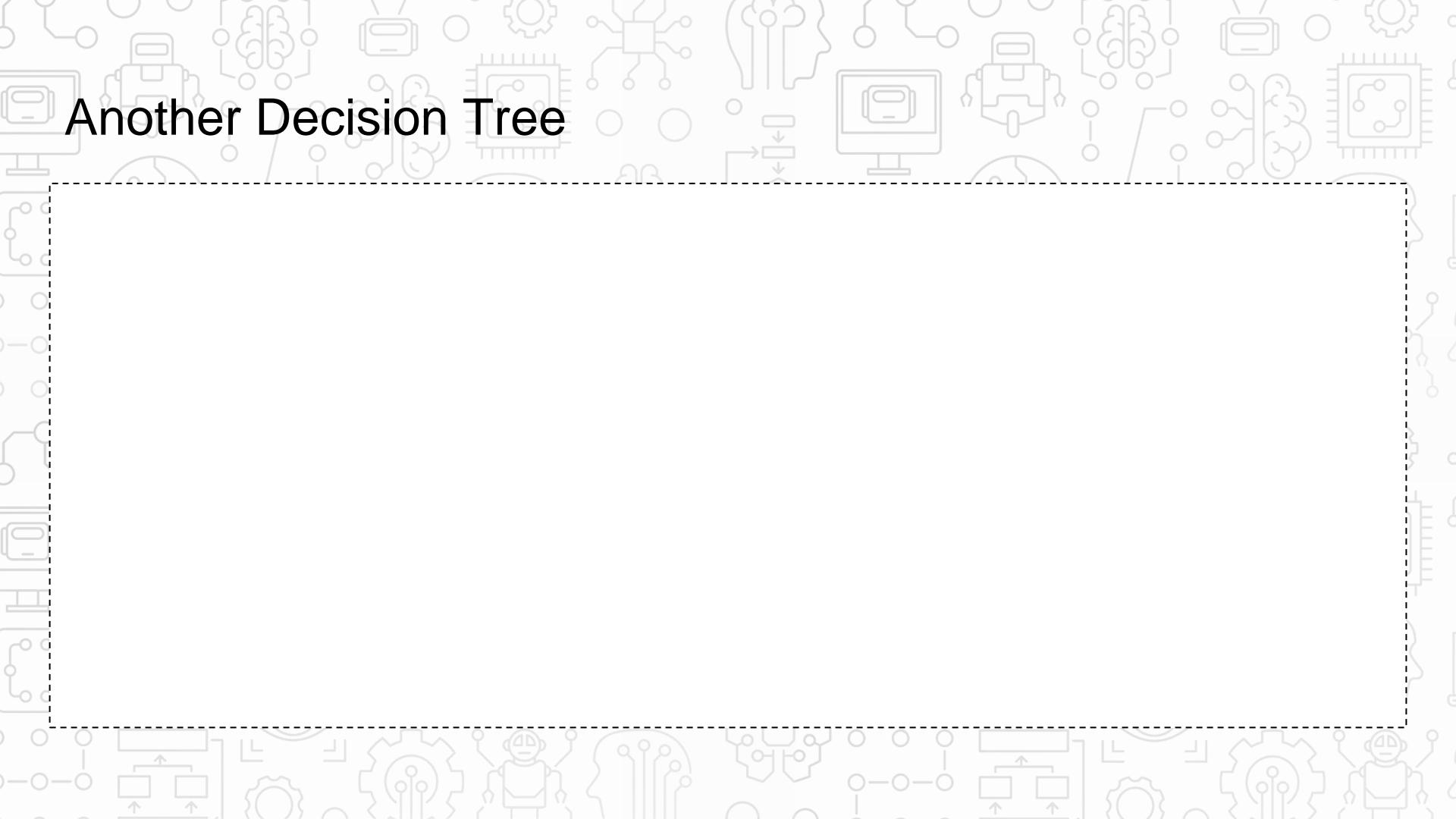
Node Each yes/no question is represented by a node, or decision stump, with two branches emanating from it.

Leaf When we reach a point where we don't ask a question and instead we make a decision

Stub A tree with only one node and 2 leaves. It represents a single decision.



Another Decision Tree



Another Decision Tree

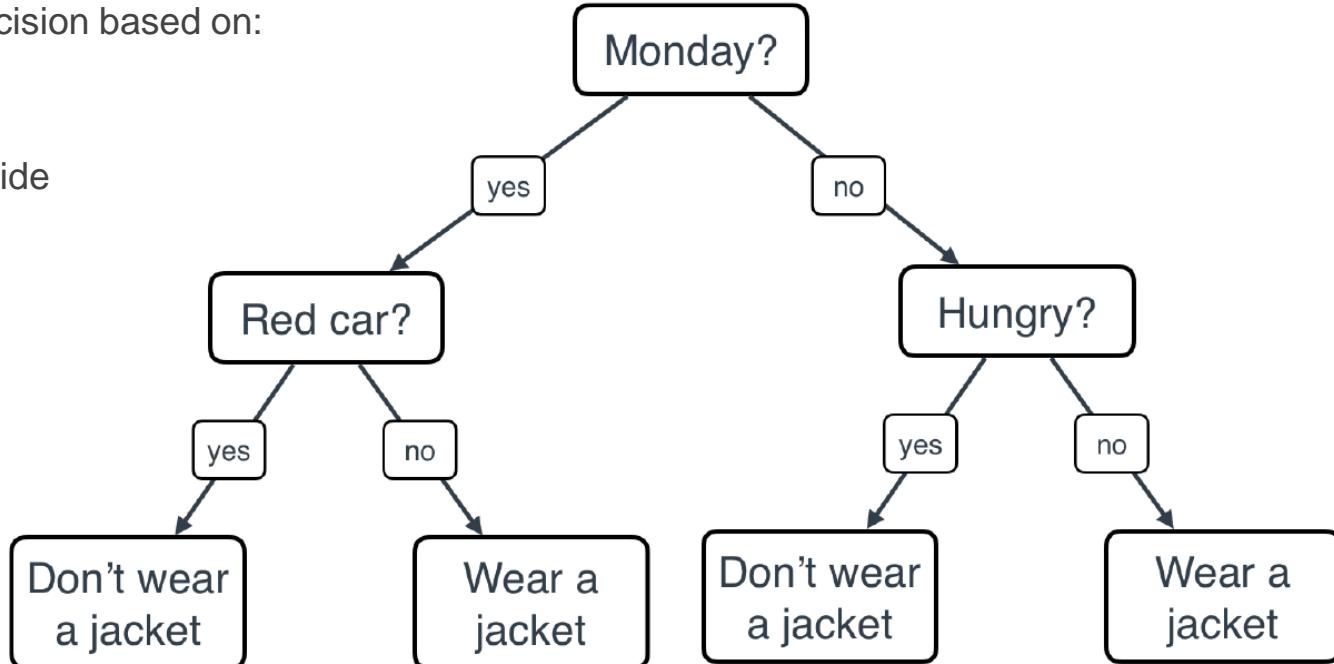
What if we takes our decision based on:

- Hungry or not
- Monday or not
- There is a red car outside

Another Decision Tree

What if we takes our decision based on:

- Hungry or not
- Monday or not
- There is a red car outside

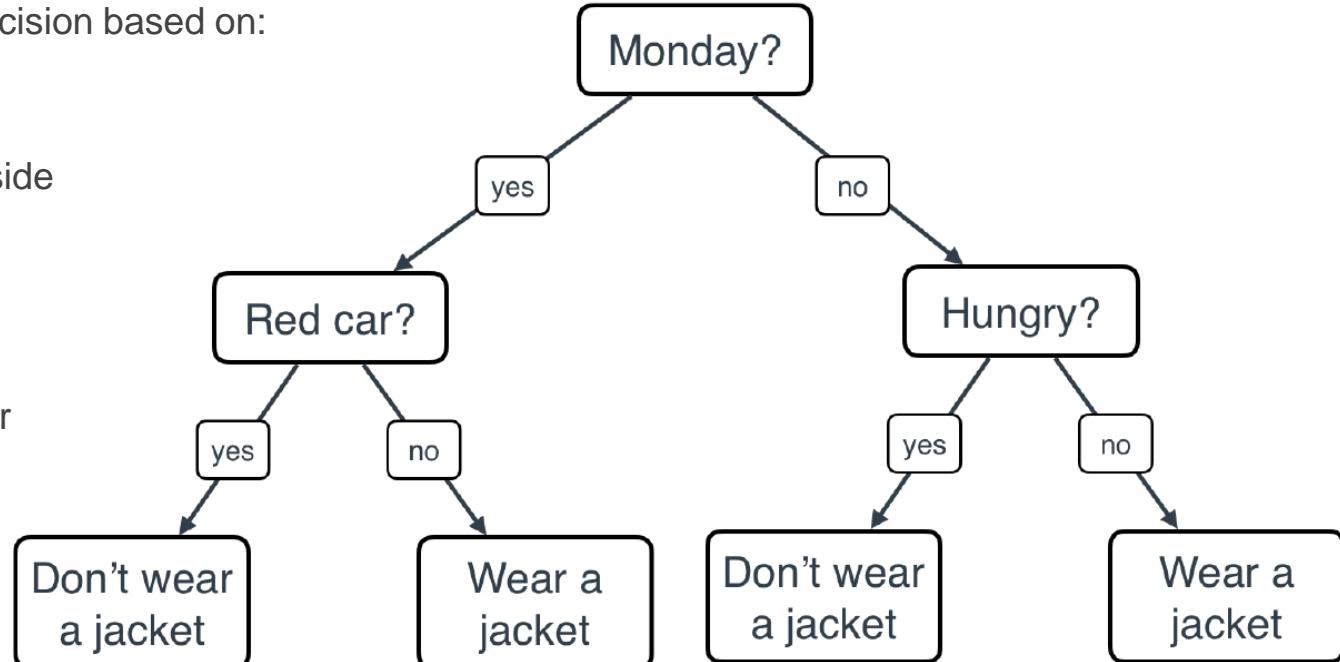


Another Decision Tree

What if we takes our decision based on:

- Hungry or not
- Monday or not
- There is a red car outside

Which mode is better?
And how does computer
know that?



Another Decision Tree

Picking our first question:

1. Is it raining?
2. Is it hot outside?
3. Am I hungry?
4. Is there a red car outside?
5. Is it Monday?

Another Decision Tree

Picking our first question:

1. Is it raining?
2. Is it hot outside?
3. Am I hungry?
4. Is there a red car outside?
5. Is it Monday?

Neglecting the last three questions, data says that wearing a jacket was true:

- 250 times based on raining or not
- 200 times based on hot or not

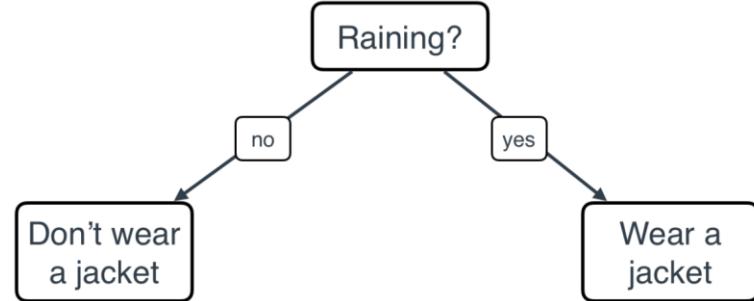
Another Decision Tree

Picking our first question:

1. Is it raining?
2. Is it hot outside?
3. Am I hungry?
4. Is there a red car outside?
5. Is it Monday?

Neglecting the last three questions, data says that wearing a jacket was true:

- 250 times based on raining or not
- 200 times based on hot or not



Another Decision Tree

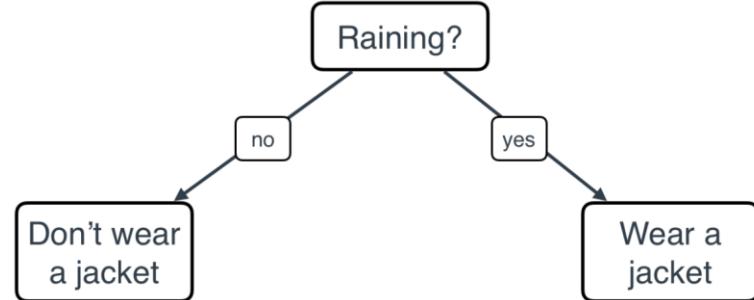
Picking our first question:

1. Is it raining?
2. Is it hot outside?
3. Am I hungry?
4. Is there a red car outside?
5. Is it Monday?

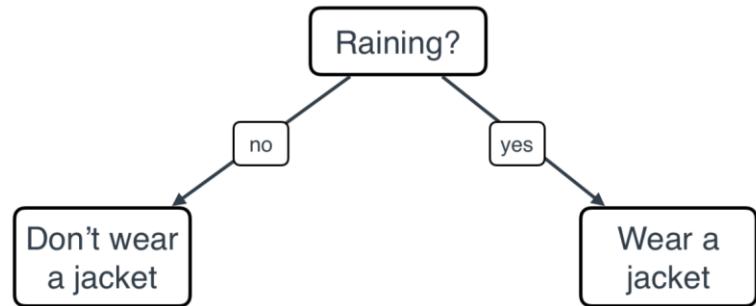
Neglecting the last three questions, data says that wearing a jacket was true:

- 250 times based on raining or not
- 200 times based on hot or not

Can we enhance this model?



Another Decision Tree

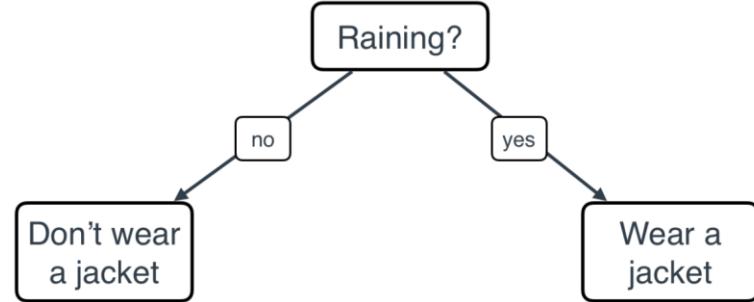


Another Decision Tree

Let's say we noticed that the right part of the tree is very accurate.

But the left part of the tree is not so correct.

So we can add a new question at the left to enhance the model.

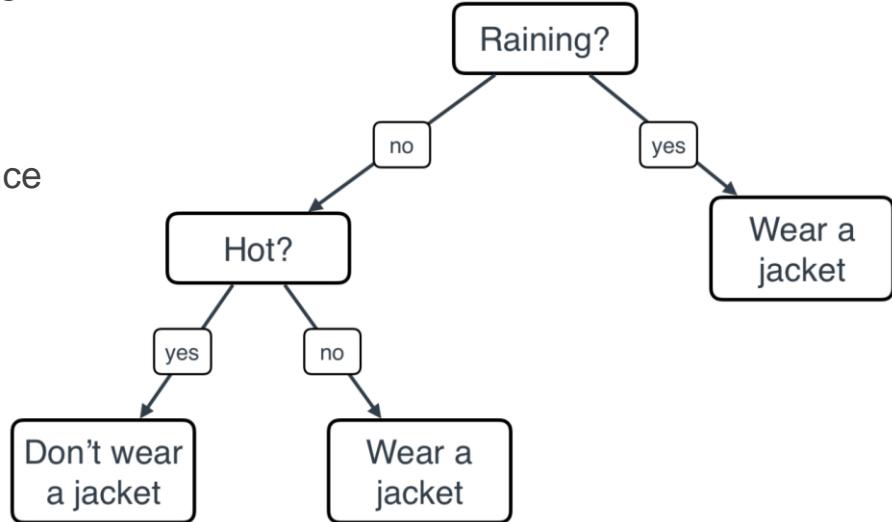


Another Decision Tree

Let's say we noticed that the right part of the tree is very accurate.

But the left part of the tree is not so correct.

So we can add a new question at the left to enhance the model.

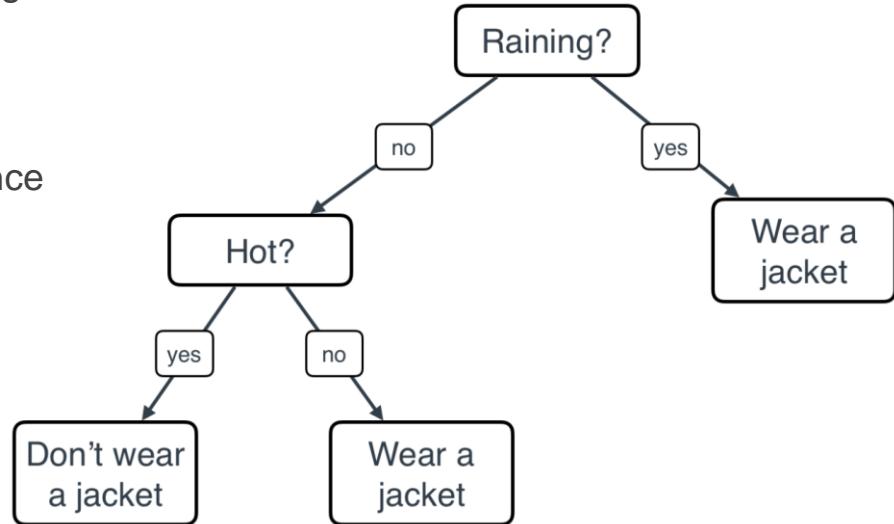


Another Decision Tree

Let's say we noticed that the right part of the tree is very accurate.

But the left part of the tree is not so correct.

So we can add a new question at the left to enhance the model.



How we choose the best question each time?

Lecture Overview

Simple Decision Tree

Mobile Apps Recommendation Example

Split Using Accuracy

Split Using Gini Impurity

Split Using Entropy

Tree Hyperparameters

One-hot Encoded Features

Continues Features

Decision Tree For Regression

Random Forest

Mobile Apps Recommendation System Example

Gender	Age	App
F	15	 Snapchat
F	25	 SHEIN
M	32	 Fantasy Premier League
F	35	 SHEIN
M	12	 Snapchat
M	14	 Snapchat

Mobile Apps Recommendation System Example

Guess which app to recommend to these customers:

- A girl aged 13
- A woman aged 28
- A man aged 34

Gender	Age	App
F	15	 Snapchat
F	25	 SHEIN
M	32	 Fantasy Premier League
F	35	 SHEIN
M	12	 Snapchat
M	14	 Snapchat

Mobile Apps Recommendation System Example

Guess which app to recommend to these customers:

- A girl aged 13
- A woman aged 28
- A man aged 34

How we know that?

Gender	Age	App
F	15	 Snapchat
F	25	 SHEIN
M	32	 Fantasy Premier League
F	35	 SHEIN
M	12	 Snapchat
M	14	 Snapchat

Mobile Apps Recommendation System Example

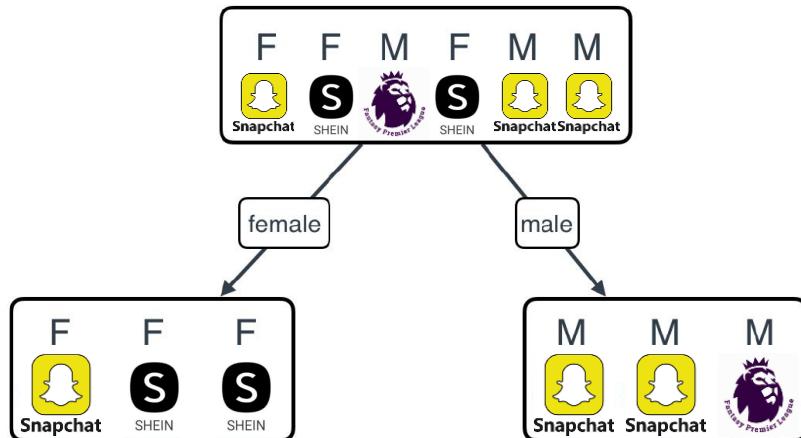
Which question would be better to start with?

Gender	Age	App
F	young	 Snapchat
F	adult	 SHEIN
M	adult	 Fantasy Premier League
F	adult	 SHEIN
M	young	 Snapchat
M	young	 Snapchat

Mobile Apps Recommendation System Example

Which question would be better to start with?

Gender?



Gender	Age	App
F	young	Snapchat
F	adult	S SHEIN
M	adult	Fantasy Premier League
F	adult	S SHEIN
M	young	Snapchat
M	young	Snapchat

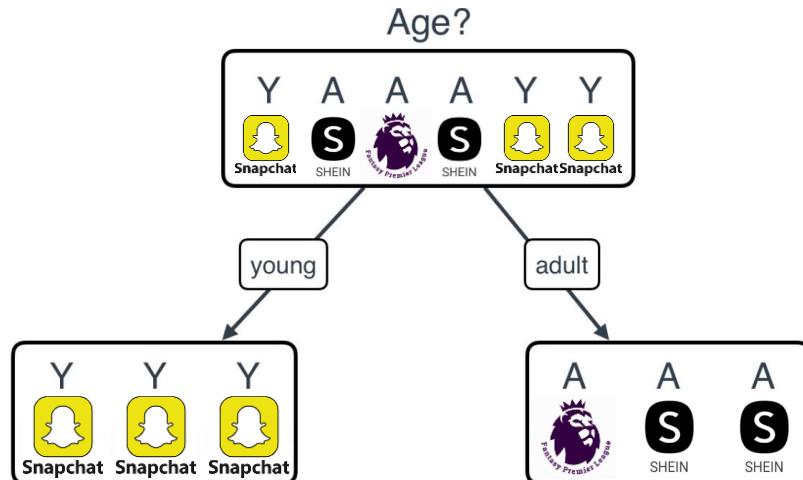
Mobile Apps Recommendation System Example

Which question would be better to start with?

Gender	Age	App
F	young	 Snapchat
F	adult	 SHEIN
M	adult	 Fantasy Premier League
F	adult	 SHEIN
M	young	 Snapchat
M	young	 Snapchat

Mobile Apps Recommendation System Example

Which question would be better to start with?

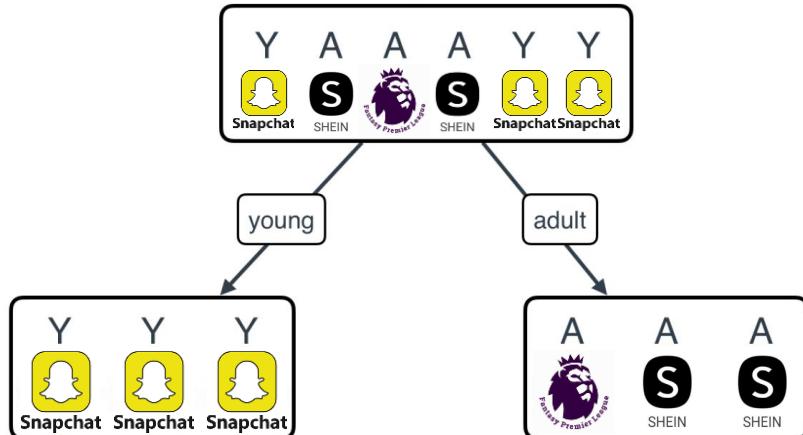


Gender	Age	App
F	young	Snapchat
F	adult	S SHEIN
M	adult	Fantasy Premier League
F	adult	S SHEIN
M	young	Snapchat
M	young	Snapchat

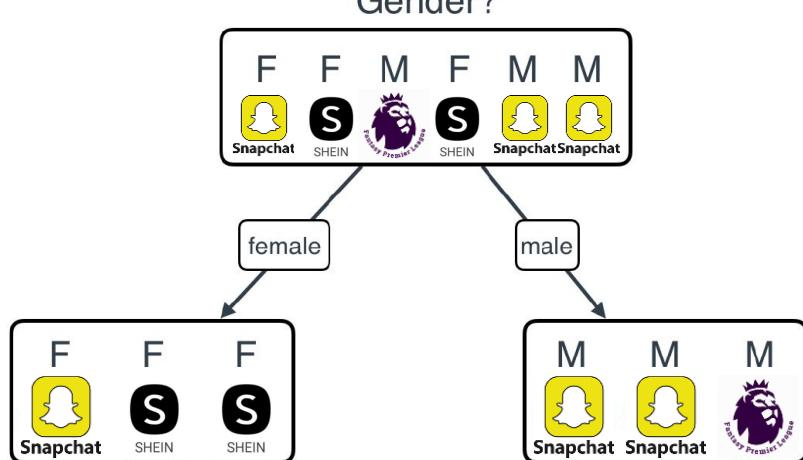
Mobile Apps Recommendation System Example

Which question would be better to start with?

Age?



Gender?



Mobile Apps Recommendation System Example

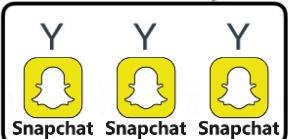
Which question would be better to start with?

Age?



young

adult



Gender?



female

male



How does computer know that asking about age is better?

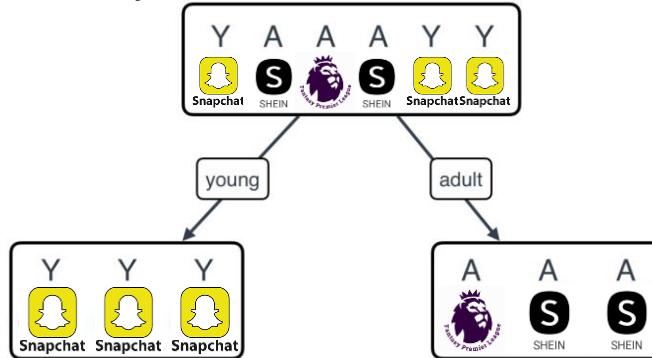
Mobile Apps Recommendation System Example

Using Accuracy:

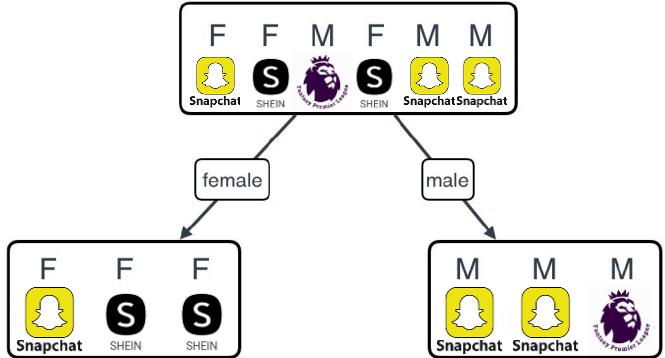
Mobile Apps Recommendation System Example

Using Accuracy:

Age?



Gender?



Mobile Apps Recommendation System Example

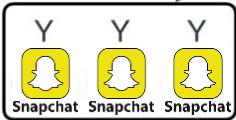
Using Accuracy:

Age?



young

adult



Gender?



female

male



Mobile Apps Recommendation System Example

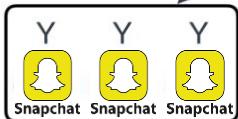
Using Accuracy:

Age?



young

adult



Snapchat



SHEIN

Gender?



female

male



SHEIN



Snapchat

Correct 5 out of 6 times

Mobile Apps Recommendation System Example

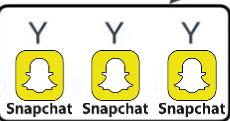
Using Accuracy:

Age?



young

adult



Snapchat



S
SHEIN

Correct 5 out of 6 times

Gender?



female

male



Snapchat
SHEIN
SHEIN

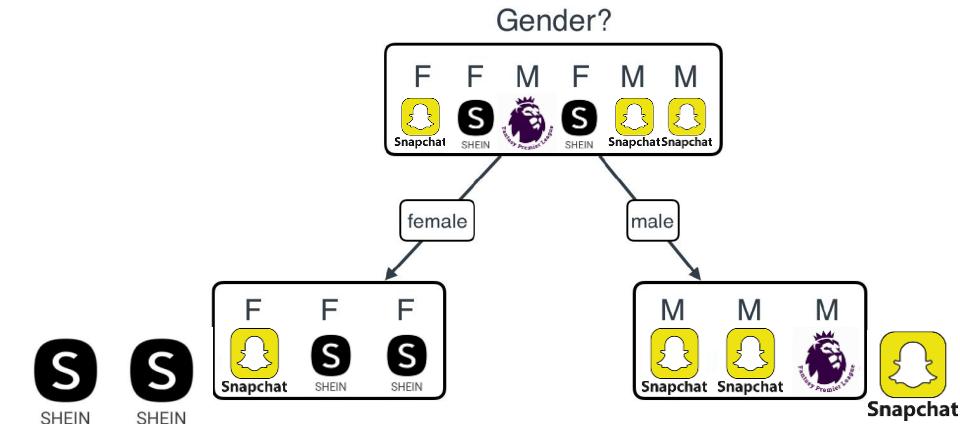
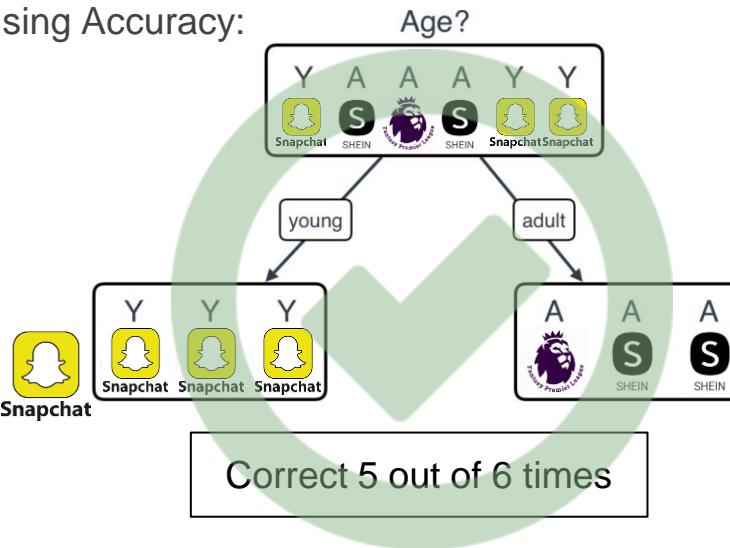


Snapchat

Correct 4 out of 6 times

Mobile Apps Recommendation System Example

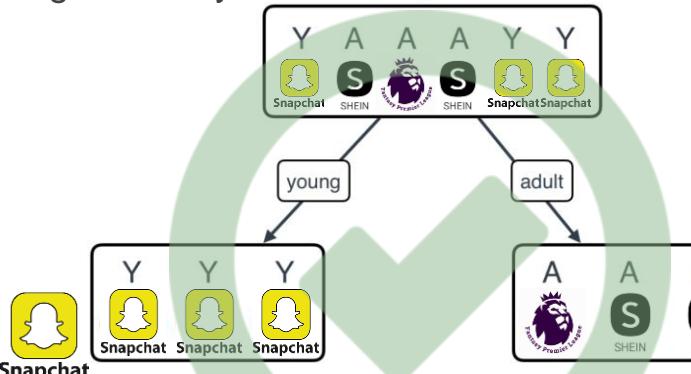
Using Accuracy:



Mobile Apps Recommendation System Example

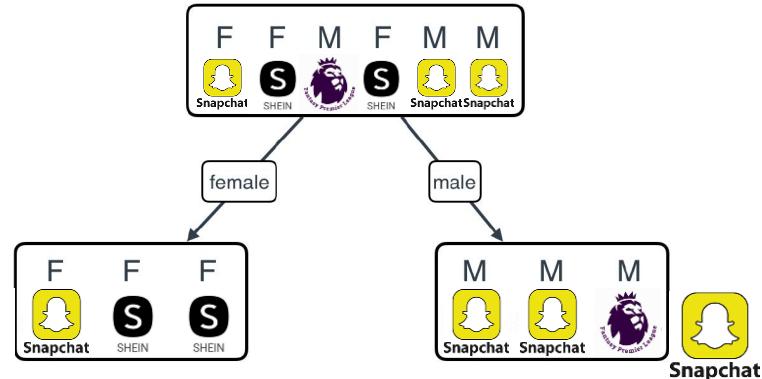
Using Accuracy:

Age?



Correct 5 out of 6 times

Gender?

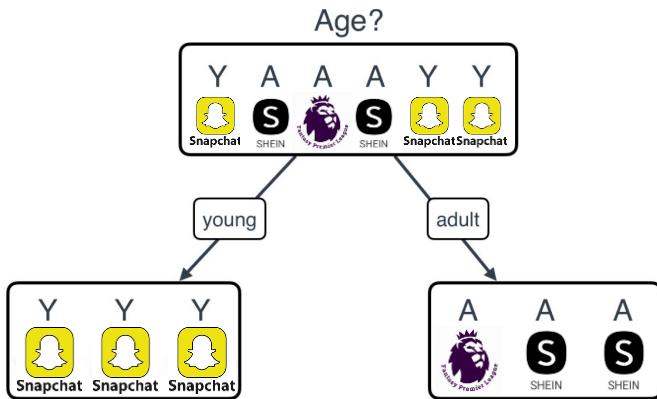


Correct 4 out of 6 times

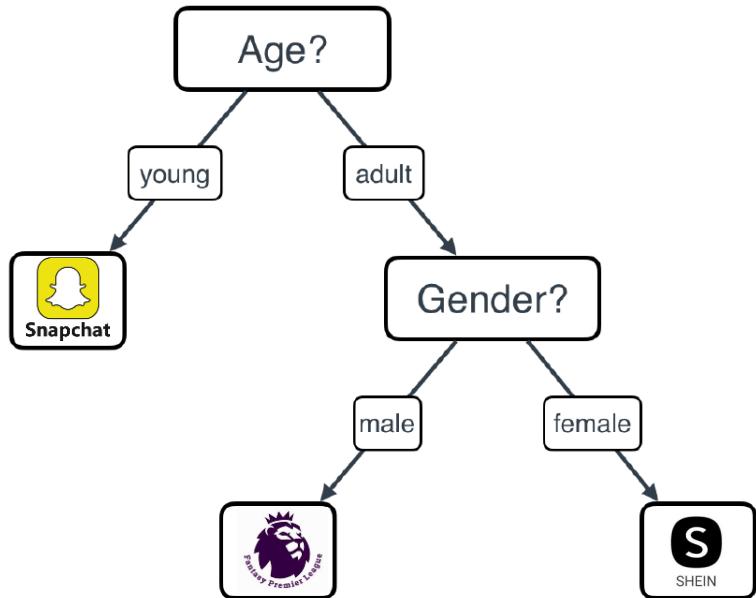
What are next steps?

Iterate by asking the best question every time

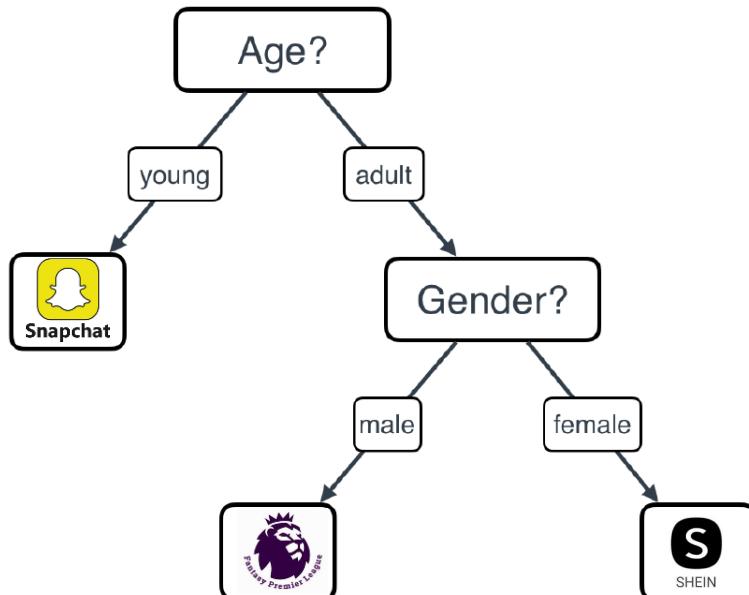
Mobile Apps Recommendation System Example



Mobile Apps Recommendation System Example

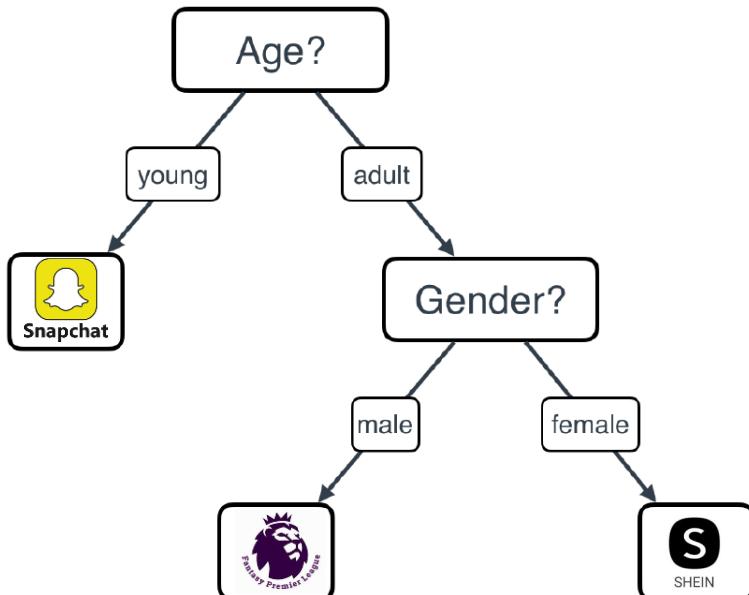


Mobile Apps Recommendation System Example



Gender	Age	App
F	adult	 SHEIN
M	adult	 Fantasy Premier League
F	adult	 SHEIN

Mobile Apps Recommendation System Example



Gender	Age	App
F	adult	SHEIN
M	adult	Fantasy Premier League
F	adult	SHEIN

Predict these examples: 1) A girl aged 13

2) A woman aged 28

3) A man aged 34

Lecture Overview

Simple Decision Tree

Mobile Apps Recommendation Example

Split Using Accuracy

Split Using Gini Impurity

Split Using Entropy

Tree Hyperparameters

One-hot Encoded Features

Continues Features

Decision Tree For Regression

Random Forest

How to pick the right feature to split

Gender	Age	Location	App
Female	Young	A	1
Male	Young	A	1
Female	Young	A	1
Male	Adult	A	1
Female	Young	B	2
Male	Adult	B	2
Female	Adult	B	2
Male	Adult	B	2

How to pick the right feature to split

We can use three methods:

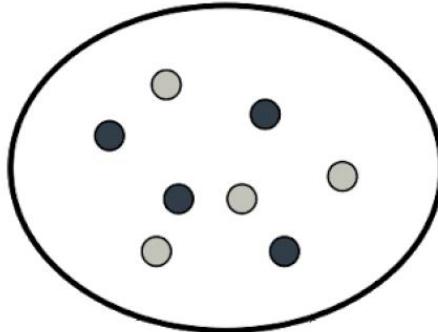
- 1- Accuracy
- 2- Gini impurity
- 3- Entropy and information gain

Gender	Age	Location	App
Female	Young	A	1
Male	Young	A	1
Female	Young	A	1
Male	Adult	A	1
Female	Young	B	2
Male	Adult	B	2
Female	Adult	B	2
Male	Adult	B	2

How to pick the right feature to split

We can use three methods:

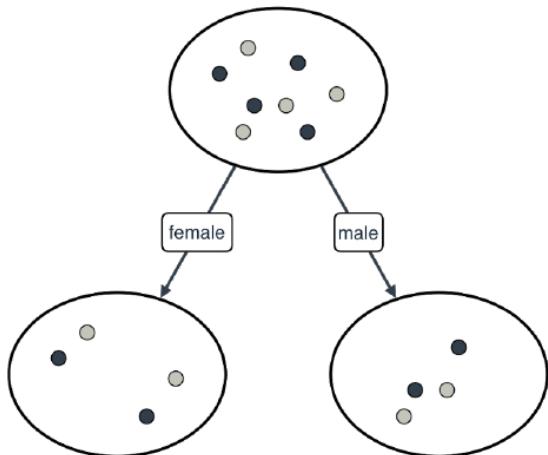
- 1- Accuracy
- 2- Gini impurity
- 3- Entropy and information gain



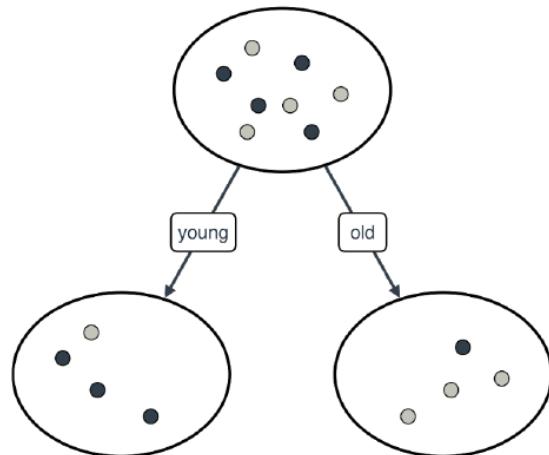
Gender	Age	Location	App
Female	Young	A	1
Male	Young	A	1
Female	Young	A	1
Male	Adult	A	1
Female	Young	B	2
Male	Adult	B	2
Female	Adult	B	2
Male	Adult	B	2

How to pick the right feature to split

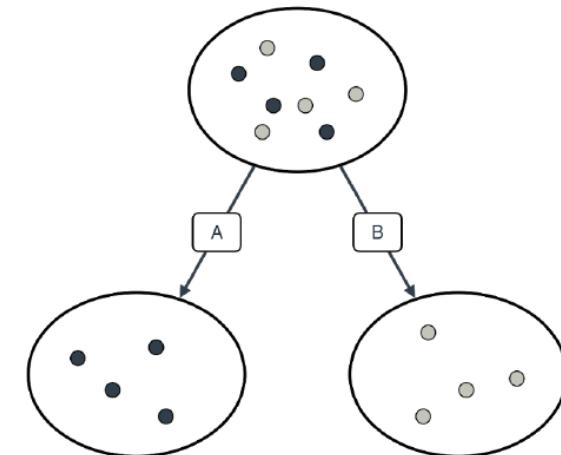
Split on gender



Split on age

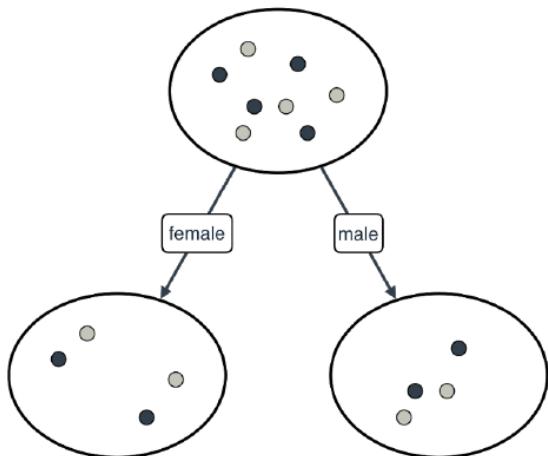


Split on location

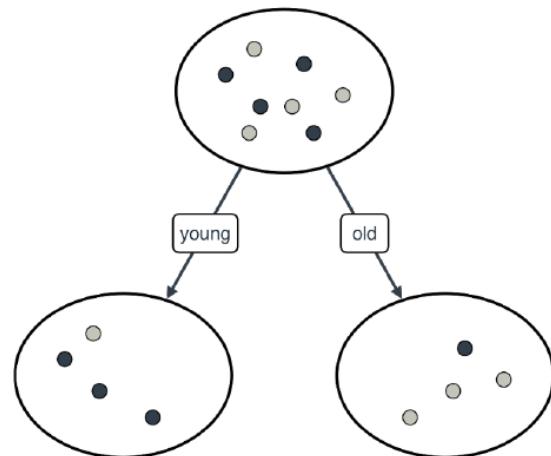


How to pick the right feature to split

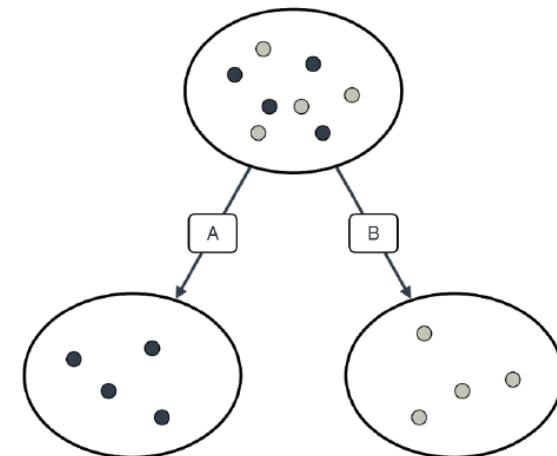
Split on gender



Split on age



Split on location



What is the right feature to split from your point of view?

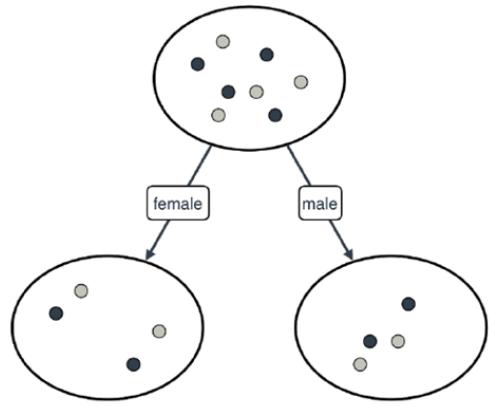
How to pick the right feature to split

Using **Accuracy**:

How to pick the right feature to split

Using **Accuracy**:

Split on gender



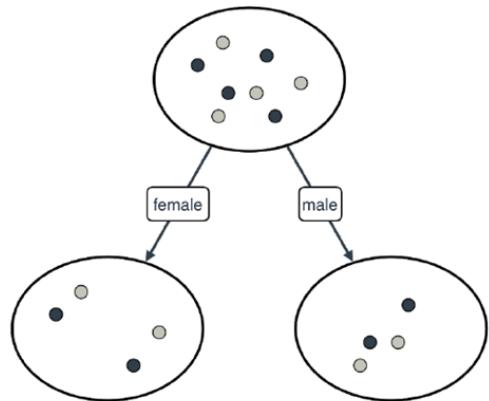
Predict App 1

Predict App 2

How to pick the right feature to split

Using **Accuracy**:

Split on gender



Predict App 1

Predict App 2

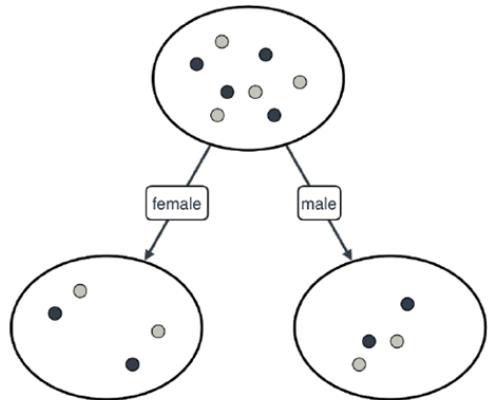
Classifier 1

Accuracy: 50%

How to pick the right feature to split

Using **Accuracy**:

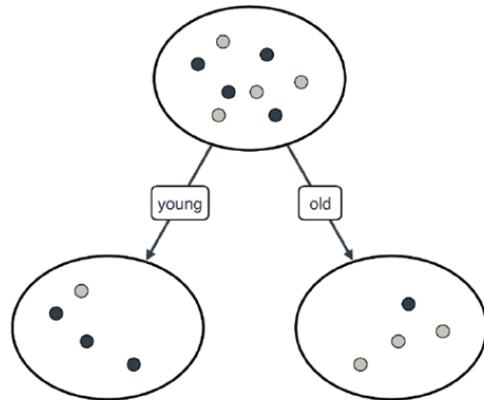
Split on gender



Predict App 1

Predict App 2

Split on age



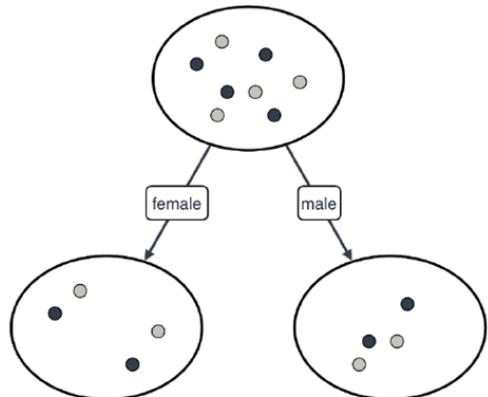
Predict App 1

Classifier 1
Accuracy: 50%

How to pick the right feature to split

Using **Accuracy**:

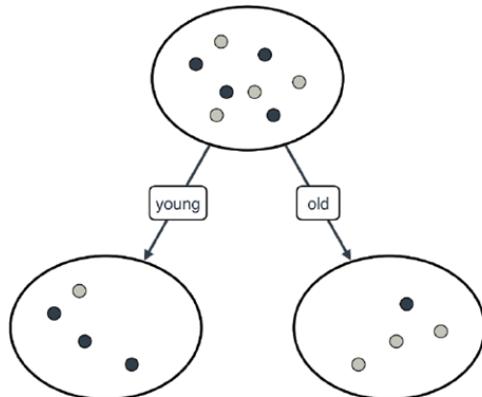
Split on gender



Predict App 1

Classifier 1
Accuracy: 50%

Split on age



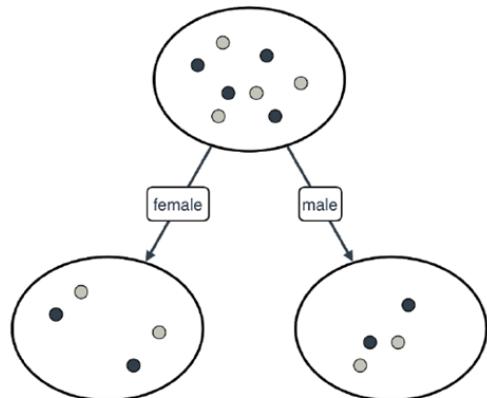
Predict App 1

Classifier 2
Accuracy: 75%

How to pick the right feature to split

Using **Accuracy**:

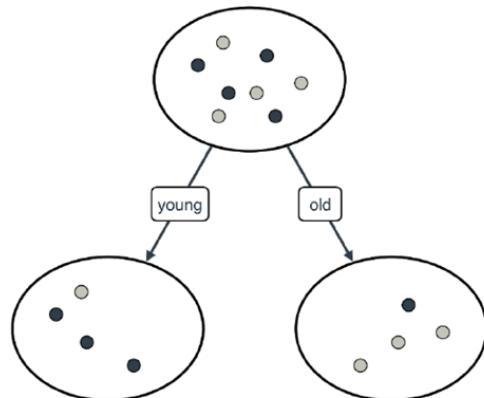
Split on gender



Predict App 1

Classifier 1
Accuracy: 50%

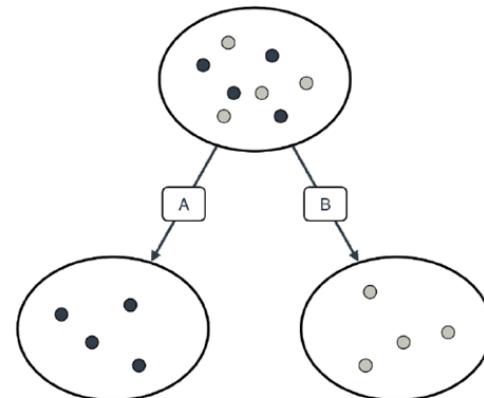
Split on age



Predict App 1

Classifier 2
Accuracy: 75%

Split on location



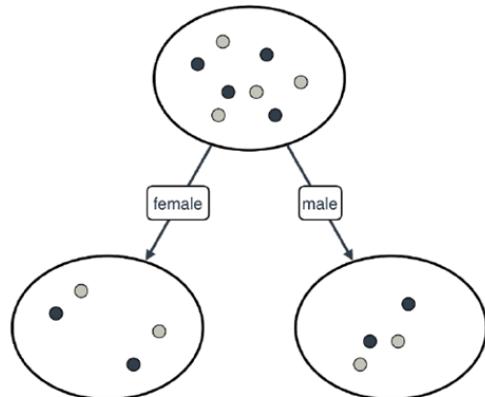
Predict App 1

Predict App 2

How to pick the right feature to split

Using **Accuracy**:

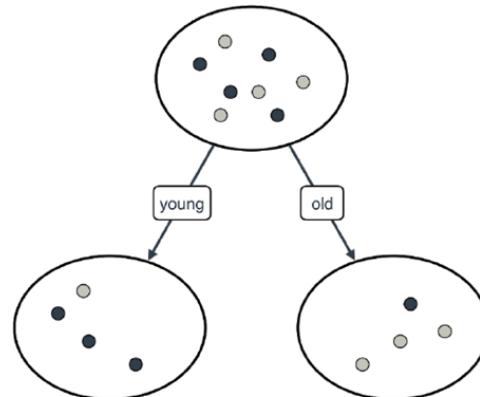
Split on gender



Predict App 1

Classifier 1
Accuracy: 50%

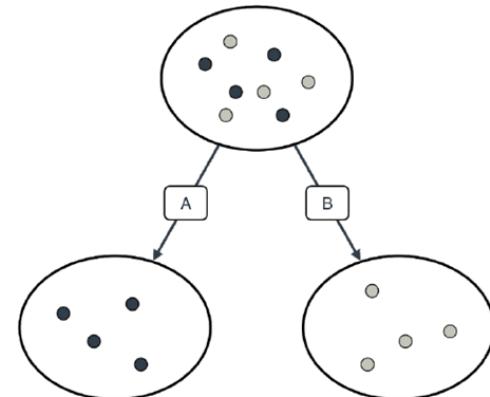
Split on age



Predict App 1

Classifier 2
Accuracy: 75%

Split on location



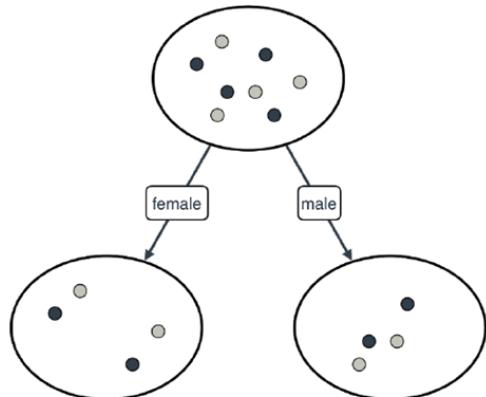
Predict App 1

Classifier 3
Accuracy: 100%

How to pick the right feature to split

Using **Accuracy**:

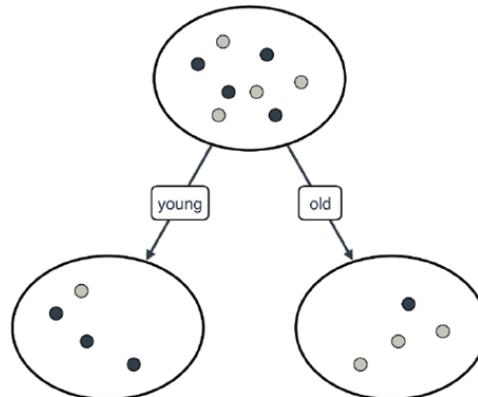
Split on gender



Predict App 1

Classifier 1
Accuracy: 50%

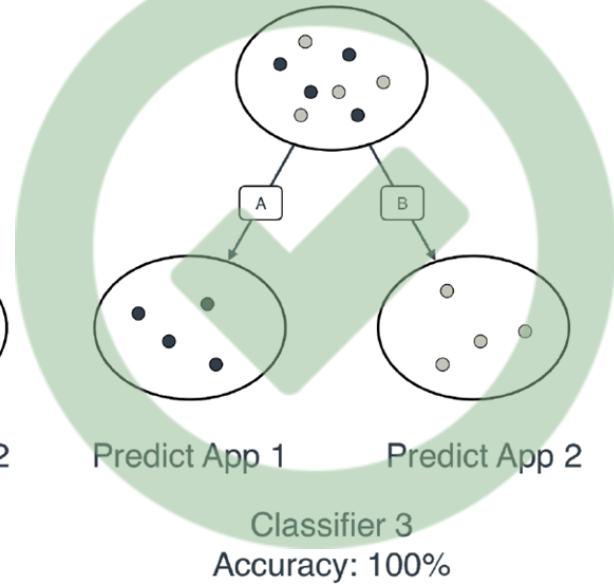
Split on age



Predict App 1

Classifier 2
Accuracy: 75%

Split on location



Predict App 1

Classifier 3
Accuracy: 100%

Lecture Overview

Simple Decision Tree

Mobile Apps Recommendation Example

Split Using Accuracy

Split Using Gini Impurity

Split Using Entropy

Tree Hyperparameters

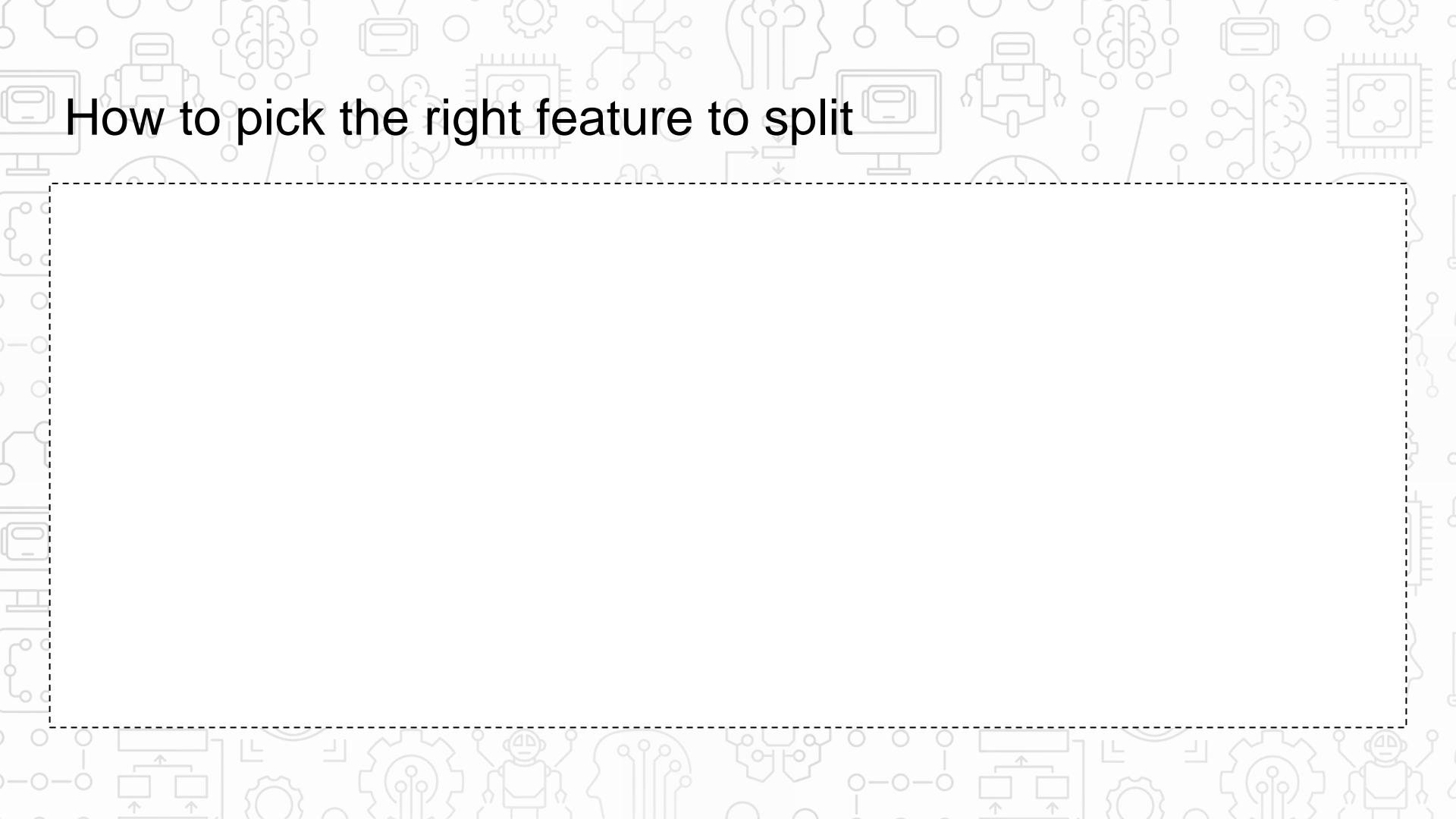
One-hot Encoded Features

Continues Features

Decision Tree For Regression

Random Forest

How to pick the right feature to split



How to pick the right feature to split

Using **Gini Impurity**:

Gini impurity is a measure of
Non-homogeneity in a set.

How to pick the right feature to split

Using **Gini Impurity**:

Gini impurity is a measure of
Non-homogeneity in a set.



Low Gini
impurity index



High Gini
impurity index

How to pick the right feature to split

Using **Gini Impurity**:

Gini impurity is a measure of Non-homogeneity in a set.

Possible sets:

- Set 1: App 1, App 1, App 1, App 1.
- Set 2: App 1, App 1, App 1, App 2.
- Set 3: App 1, App 1, App 2, App 2.
- Set 4: App 1, App 2, App 2, App 2.
- Set 5: App 2, App 2, App 2, App 2.



Low Gini
impurity index



High Gini
impurity index

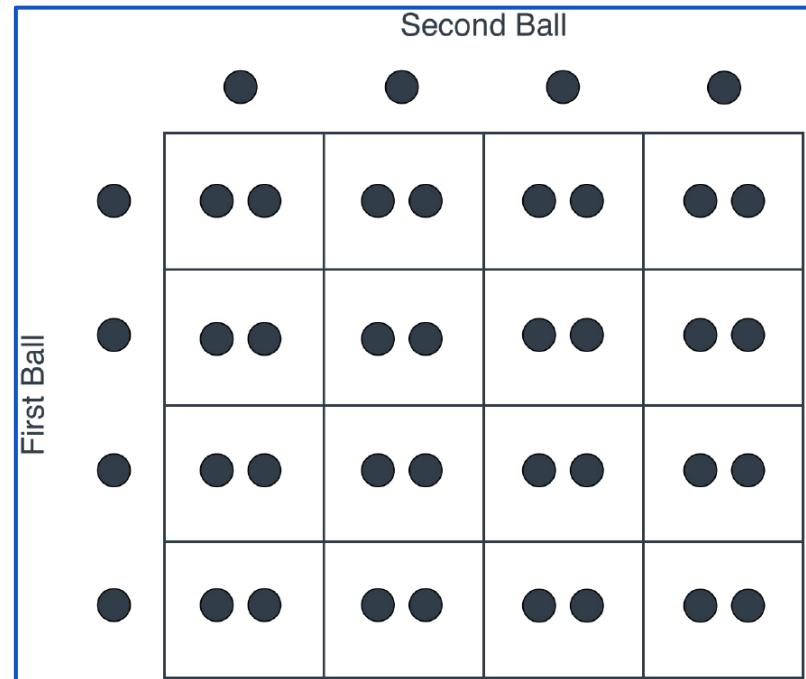
How to pick the right feature to split

Using **Gini Impurity**:

Gini impurity is a measure of Non-homogeneity in a set.

Possible sets:

- Set 1: App 1, App 1, App 1, App 1.
- Set 2: App 1, App 1, App 1, App 2.
- Set 3: App 1, App 1, App 2, App 2.
- Set 4: App 1, App 2, App 2, App 2.
- Set 5: App 2, App 2, App 2, App 2.



How to pick the right feature to split

Using **Gini Impurity**:

Gini impurity is a measure of Non-homogeneity in a set.

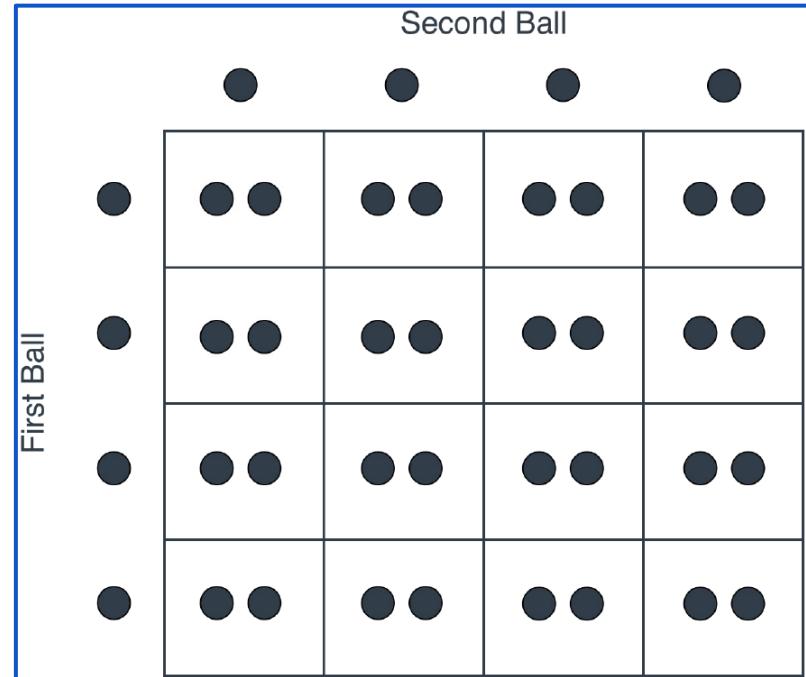
Possible sets:

- Set 1: App 1, App 1, App 1, App 1.
- Set 2: App 1, App 1, App 1, App 2.
- Set 3: App 1, App 1, App 2, App 2.
- Set 4: App 1, App 2, App 2, App 2.
- Set 5: App 2, App 2, App 2, App 2.

Gini Index

0

0



How to pick the right feature to split

Using **Gini Impurity**:

Gini impurity is a measure of Non-homogeneity in a set.

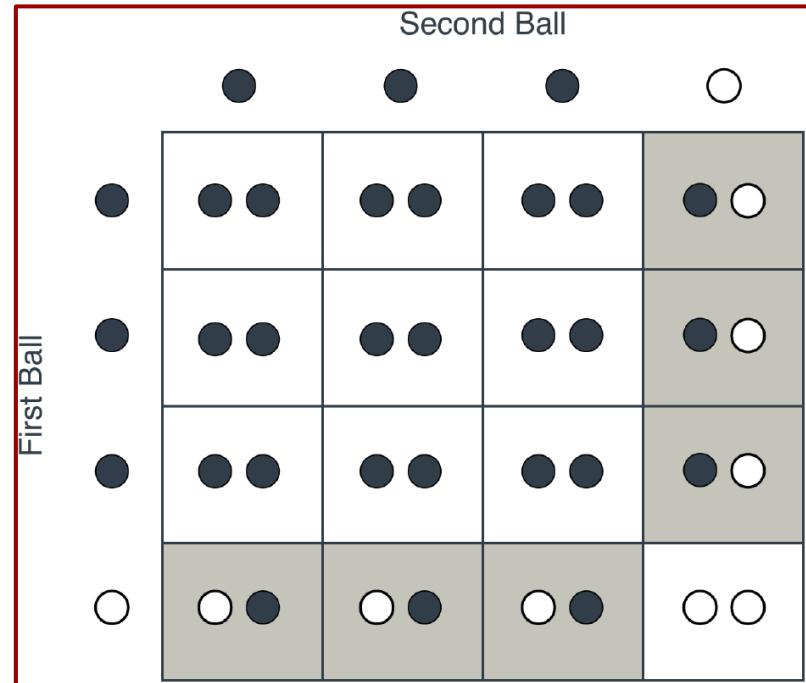
Possible sets:

- Set 1: App 1, App 1, App 1, App 1.
- Set 2: App 1, App 1, App 1, App 2.
- Set 3: App 1, App 1, App 2, App 2.
- Set 4: App 1, App 2, App 2, App 2.
- Set 5: App 2, App 2, App 2, App 2.

Gini Index

0

0



How to pick the right feature to split

Using **Gini Impurity**:

Gini impurity is a measure of Non-homogeneity in a set.

Possible sets:

- Set 1: App 1, App 1, App 1, App 1.
- Set 2: App 1, App 1, App 1, App 2.
- Set 3: App 1, App 1, App 2, App 2.
- Set 4: App 1, App 2, App 2, App 2.
- Set 5: App 2, App 2, App 2, App 2.

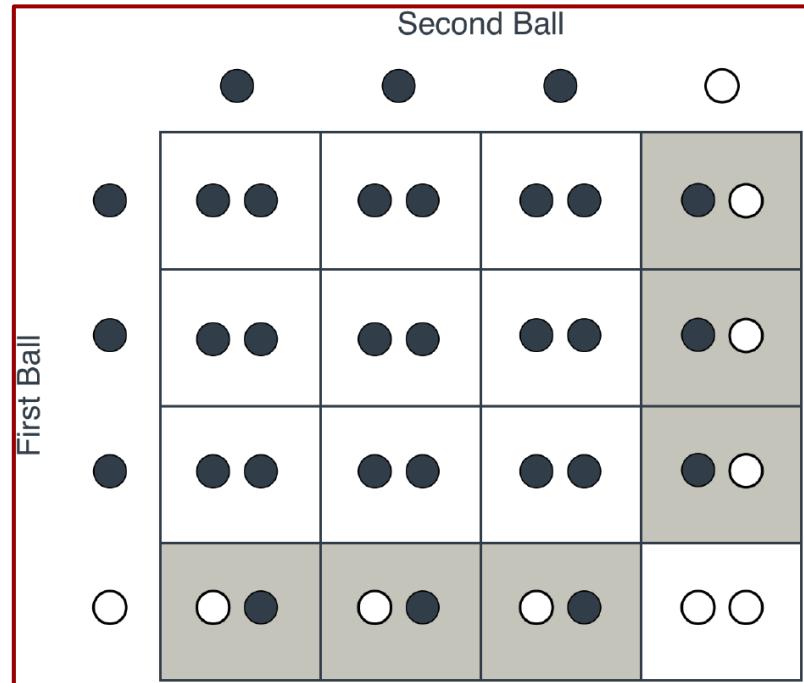
Gini Index

0

0.375

0.375

0



How to pick the right feature to split

Using **Gini Impurity**:

Gini impurity is a measure of Non-homogeneity in a set.

Possible sets:

- Set 1: App 1, App 1, App 1, App 1.
- Set 2: App 1, App 1, App 1, App 2.
- Set 3: App 1, App 1, App 2, App 2.
- Set 4: App 1, App 2, App 2, App 2.
- Set 5: App 2, App 2, App 2, App 2.

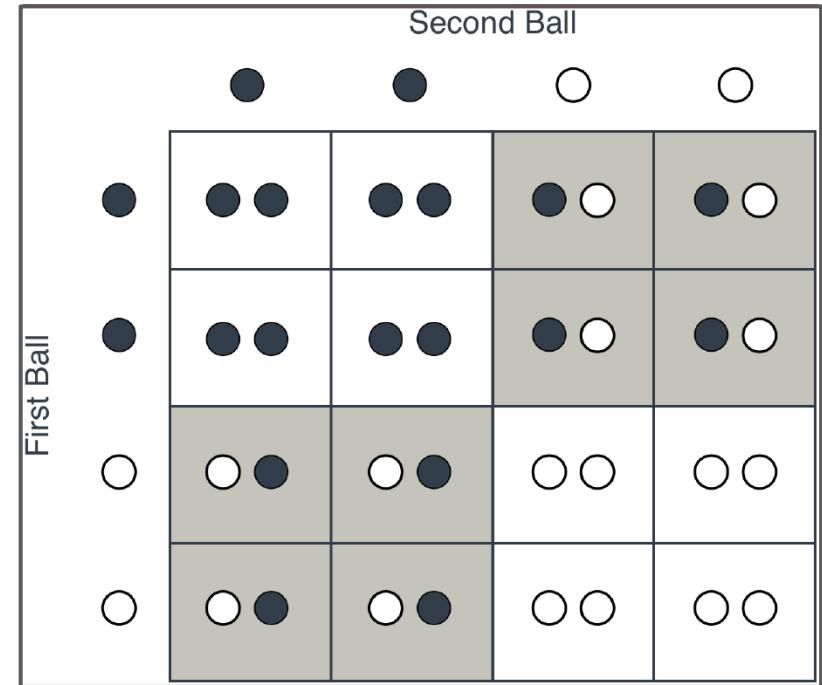
Gini Index

0

0.375

0.375

0



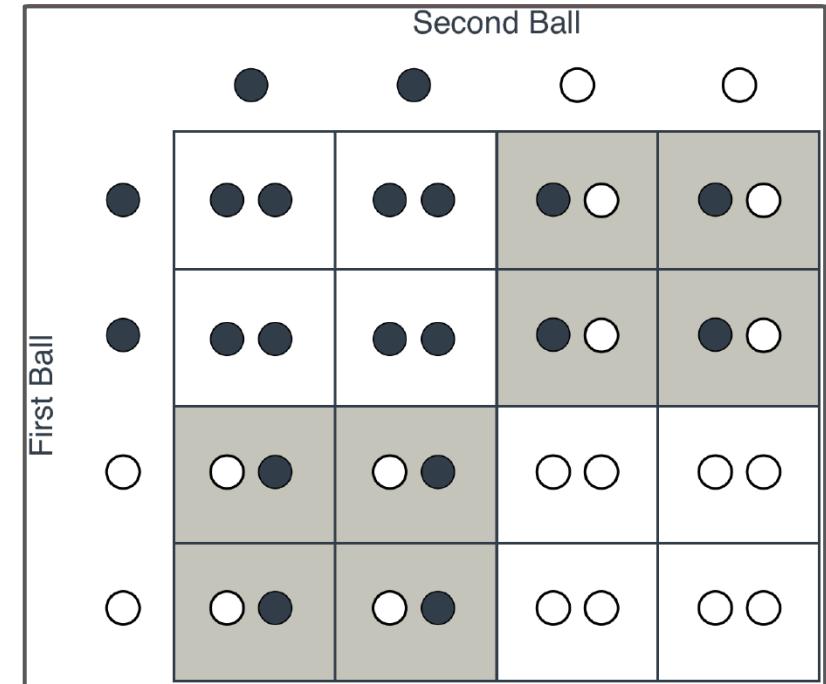
How to pick the right feature to split

Using **Gini Impurity**:

Gini impurity is a measure of Non-homogeneity in a set.

Possible sets:

- | | Gini Index |
|--------------------------------------|-------------------|
| • Set 1: App 1, App 1, App 1, App 1. | 0 |
| • Set 2: App 1, App 1, App 1, App 2. | 0.375 |
| • Set 3: App 1, App 1, App 2, App 2. | 0.5 |
| • Set 4: App 1, App 2, App 2, App 2. | 0.375 |
| • Set 5: App 2, App 2, App 2, App 2. | 0 |



How to pick the right feature to split

Using **Gini Impurity**:

Set 1



Gini impurity index = 0

Set 2



Gini impurity index = 0.375

Set 3



Gini impurity index = 0.5

Set 4



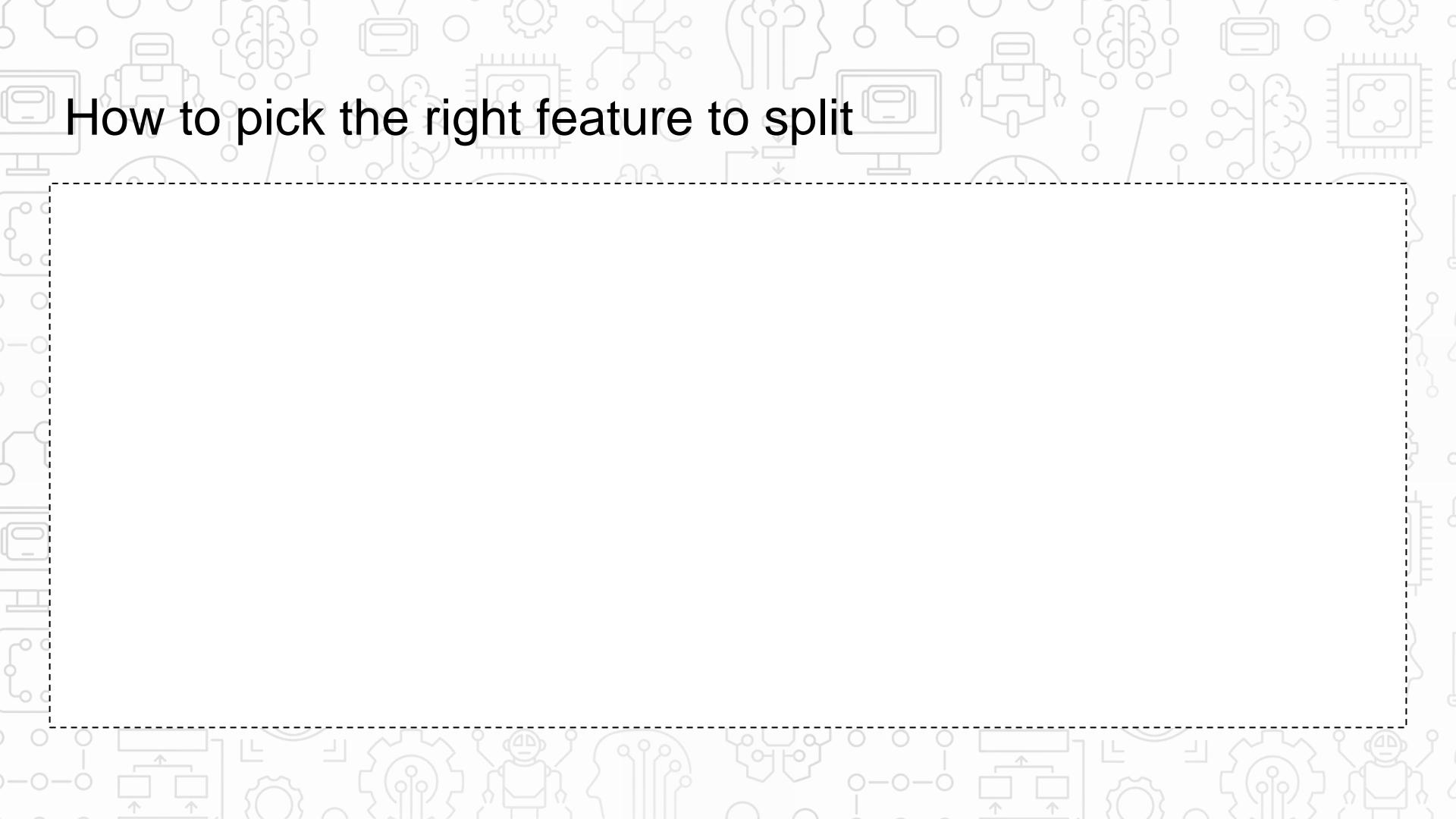
Gini impurity index = 0.375

Set 5



Gini impurity index = 0

How to pick the right feature to split



How to pick the right feature to split

Using **Gini Impurity**:

Formula:

What is the probability of picking the same ball color twice with replacement if probability of this color is P?

Probability of picking the color twice: $P_1^2 + P_2^2 + P_3^2 + P_4^2 + \dots$

Gini Impurity index: $1 - P_1^2 - P_2^2 - P_3^2 - P_4^2 - \dots$

How to pick the right feature to split

Using **Gini Impurity**:

Formula:

What is the probability of picking the same ball color twice with replacement if probability of this color is P ?

Probability of picking the color twice: $P_1^2 + P_2^2 + P_3^2 + P_4^2 + \dots$

Gini Impurity index: $1 - P_1^2 - P_2^2 - P_3^2 - P_4^2 - \dots$



Gini Impurity Index = $P(\text{picking two balls of different color})$

$$= 1 - \left(\frac{3}{6}\right)^2 - \left(\frac{2}{6}\right)^2 - \left(\frac{1}{6}\right)^2$$

$P(\text{Both balls are black})$

$P(\text{Both balls are grey})$

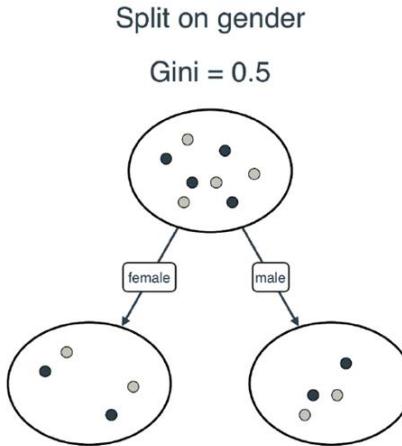
$P(\text{Both balls are white})$

How to pick the right feature to split

Using **Gini Impurity**:

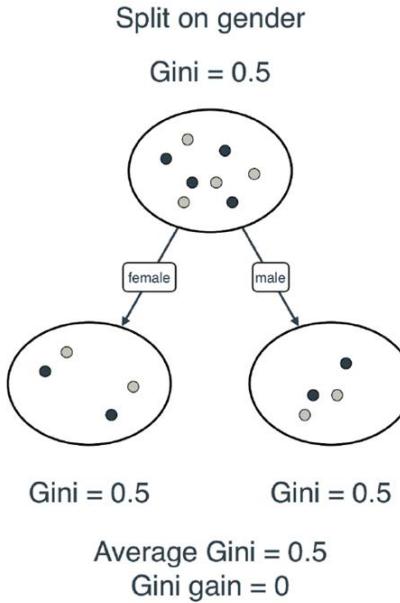
How to pick the right feature to split

Using **Gini Impurity**:



How to pick the right feature to split

Using **Gini Impurity**:

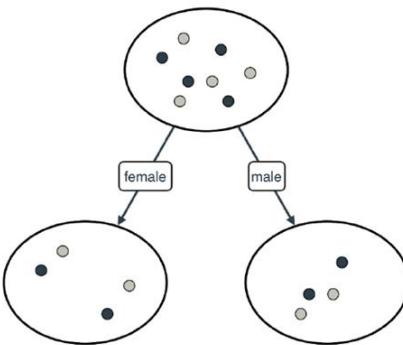


How to pick the right feature to split

Using **Gini Impurity**:

Split on gender

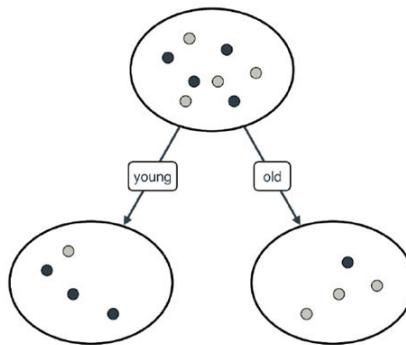
Gini = 0.5



Gini = 0.5

Split on age

Gini = 0.5



Gini = 0.5

Average Gini = 0.5

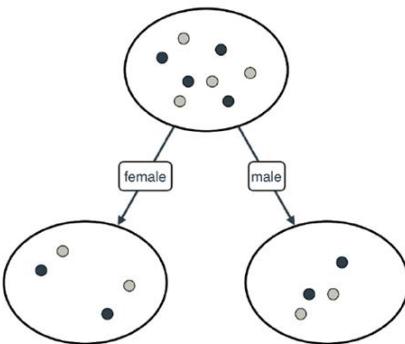
Gini gain = 0

How to pick the right feature to split

Using **Gini Impurity**:

Split on gender

Gini = 0.5



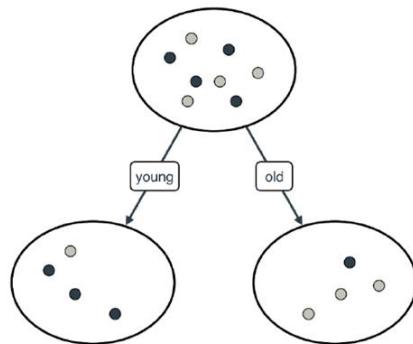
Gini = 0.5

Average Gini = 0.5

Gini gain = 0

Split on age

Gini = 0.5



Gini = 0.375

Average Gini = 0.375

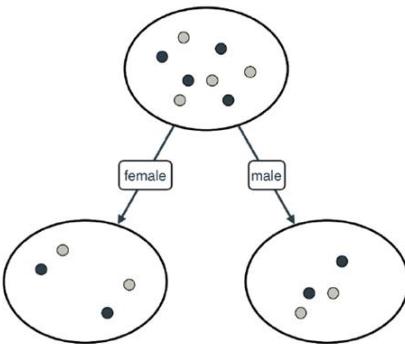
Gini gain = 0.125

How to pick the right feature to split

Using **Gini Impurity**:

Split on gender

Gini = 0.5

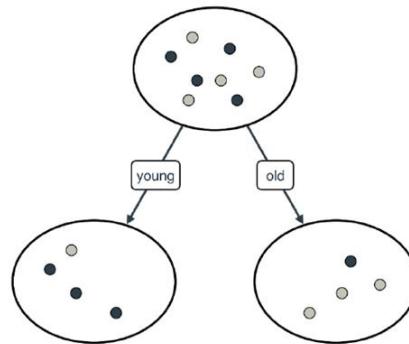


Gini = 0.5

Average Gini = 0.5
Gini gain = 0

Split on age

Gini = 0.5

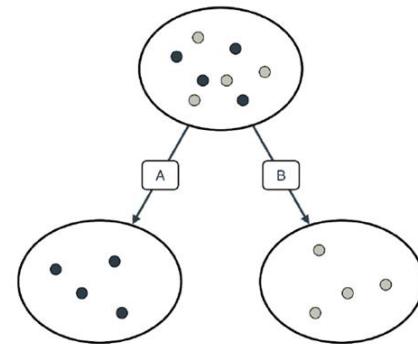


Gini = 0.375

Average Gini = 0.375
Gini gain = 0.125

Split on location

Gini = 0.5

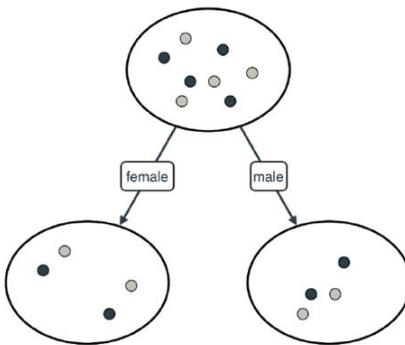


How to pick the right feature to split

Using **Gini Impurity**:

Split on gender

Gini = 0.5

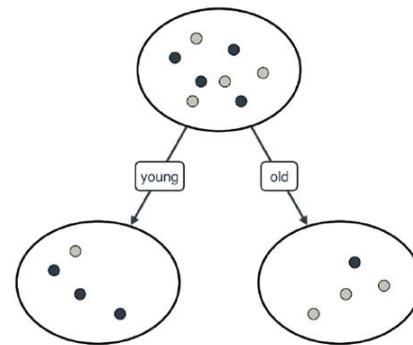


Gini = 0.5

Average Gini = 0.5
Gini gain = 0

Split on age

Gini = 0.5

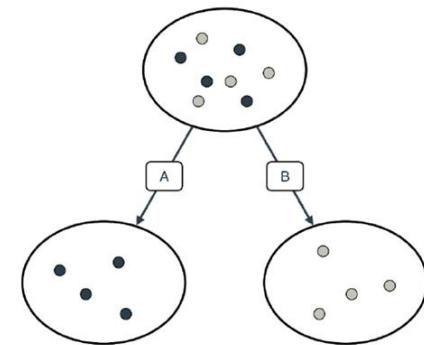


Gini = 0.375

Average Gini = 0.375
Gini gain = 0.125

Split on location

Gini = 0.5



Gini = 0

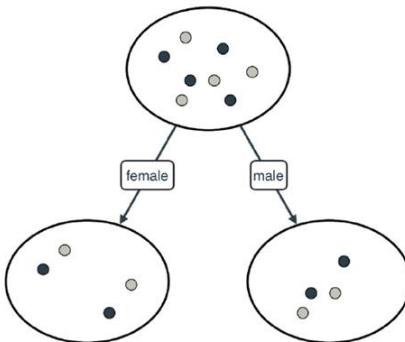
Average Gini = 0
Gini gain = 0.5

How to pick the right feature to split

Using **Gini Impurity**:

Split on gender

Gini = 0.5

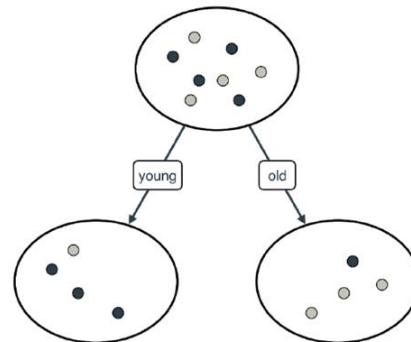


Gini = 0.5

Average Gini = 0.5
Gini gain = 0

Split on age

Gini = 0.5

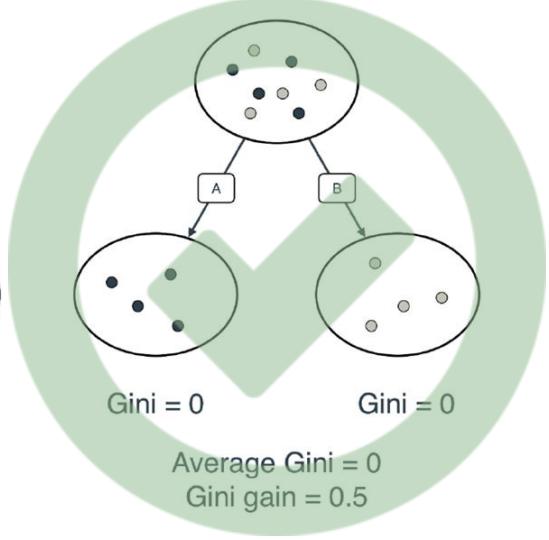


Gini = 0.375

Average Gini = 0.375
Gini gain = 0.125

Split on location

Gini = 0.5



Gini = 0

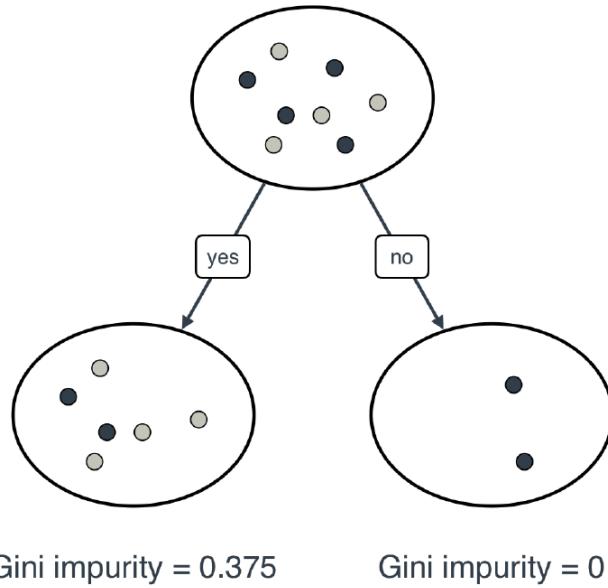
Gini = 0

Average Gini = 0
Gini gain = 0.5

How to pick the right feature to split

Using **Gini Impurity**:

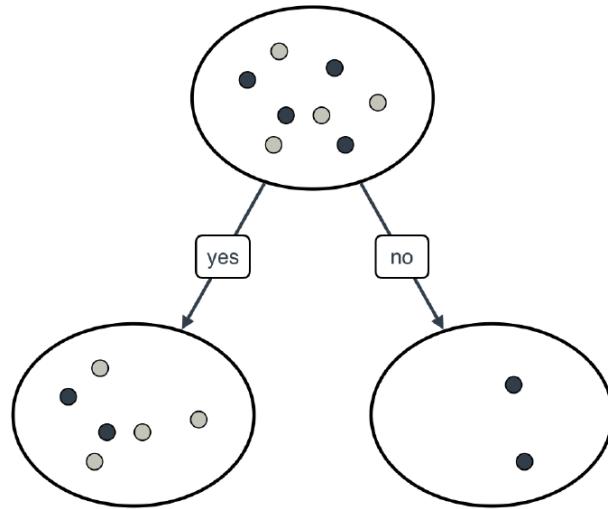
Weighted Gini Impurity:



How to pick the right feature to split

Using **Gini Impurity**:

Weighted Gini Impurity:



Gini impurity = 0.375

Gini impurity = 0

$$\text{Weighted average} = 0.375 \left(\frac{6}{8} \right) + 0 \left(\frac{2}{8} \right) = 0.25$$

Lecture Overview

Simple Decision Tree

Mobile Apps Recommendation Example

Split Using Accuracy

Split Using Gini Impurity

Split Using Entropy

Tree Hyperparameters

One-hot Encoded Features

Continues Features

Decision Tree For Regression

Random Forest

How to pick the right feature to split

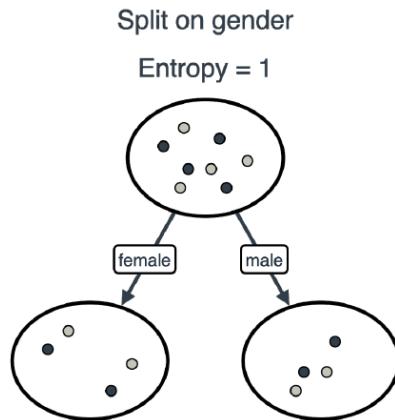
Using **Entropy**:

Formula: $E = -p_1 \log_2(p_1) - p_2 \log_2(p_2) - \dots - p_k \log_2(p_k)$.

How to pick the right feature to split

Using **Entropy**:

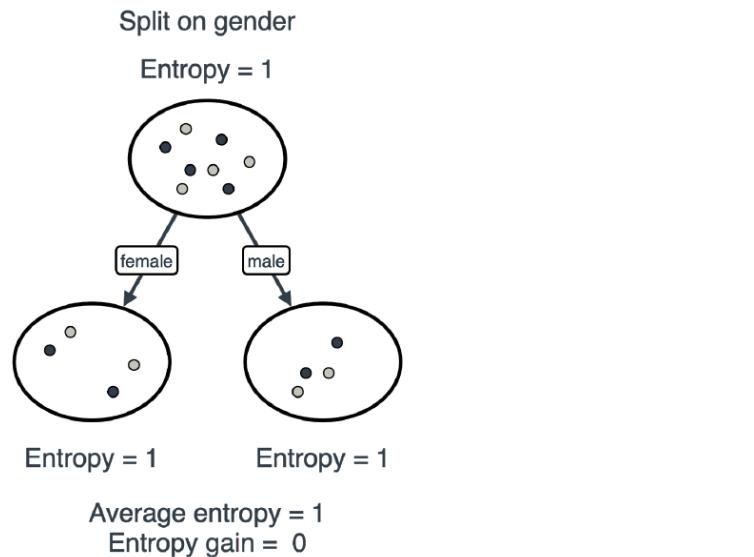
Formula: $E = -p_1 \log_2(p_1) - p_2 \log_2(p_2) - \dots - p_k \log_2(p_k)$.



How to pick the right feature to split

Using **Entropy**:

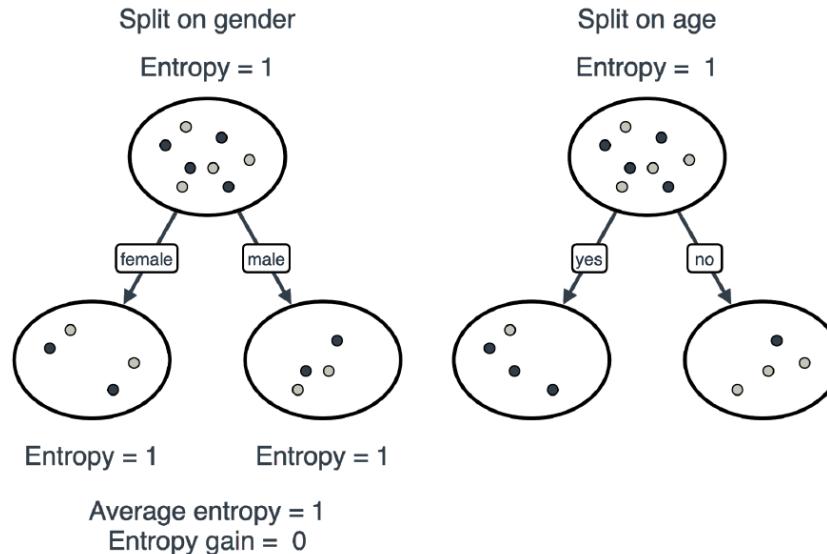
Formula: $E = -p_1 \log_2(p_1) - p_2 \log_2(p_2) - \dots - p_k \log_2(p_k)$.



How to pick the right feature to split

Using **Entropy**:

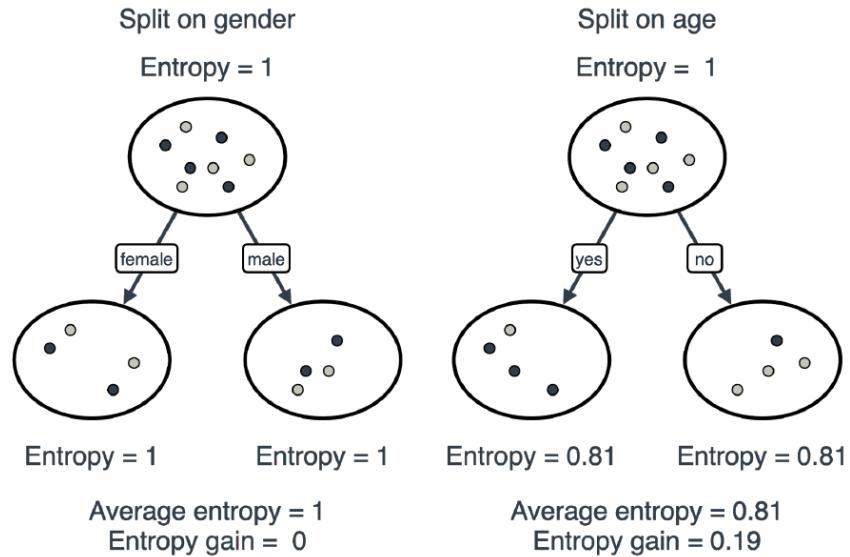
Formula: $E = -p_1 \log_2(p_1) - p_2 \log_2(p_2) - \dots - p_k \log_2(p_k)$.



How to pick the right feature to split

Using **Entropy**:

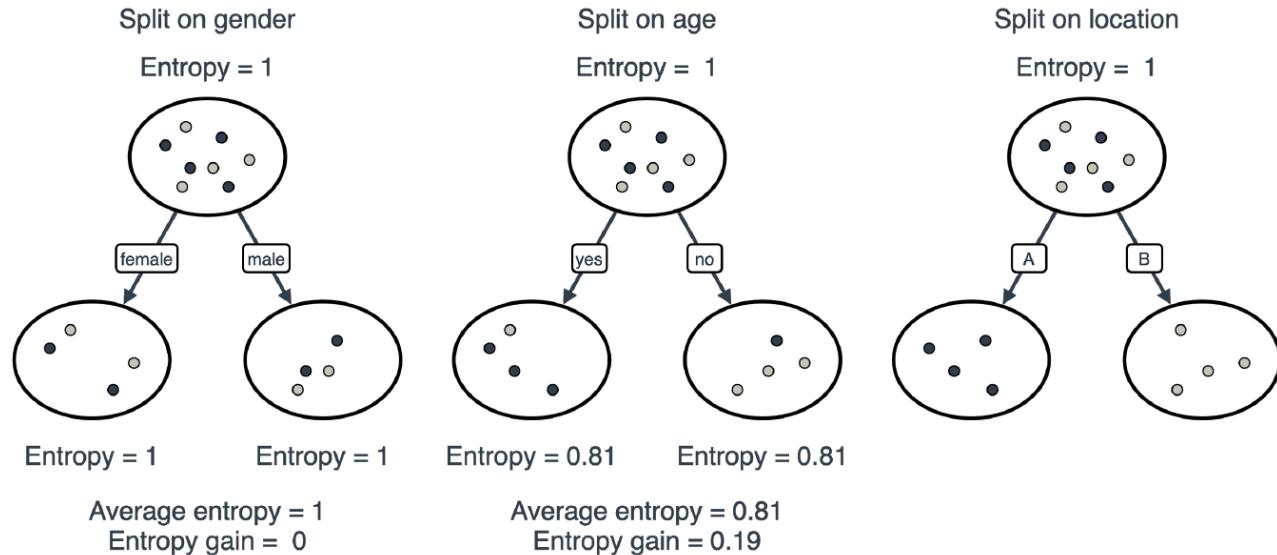
Formula: $E = -p_1 \log_2(p_1) - p_2 \log_2(p_2) - \dots - p_k \log_2(p_k)$.



How to pick the right feature to split

Using **Entropy**:

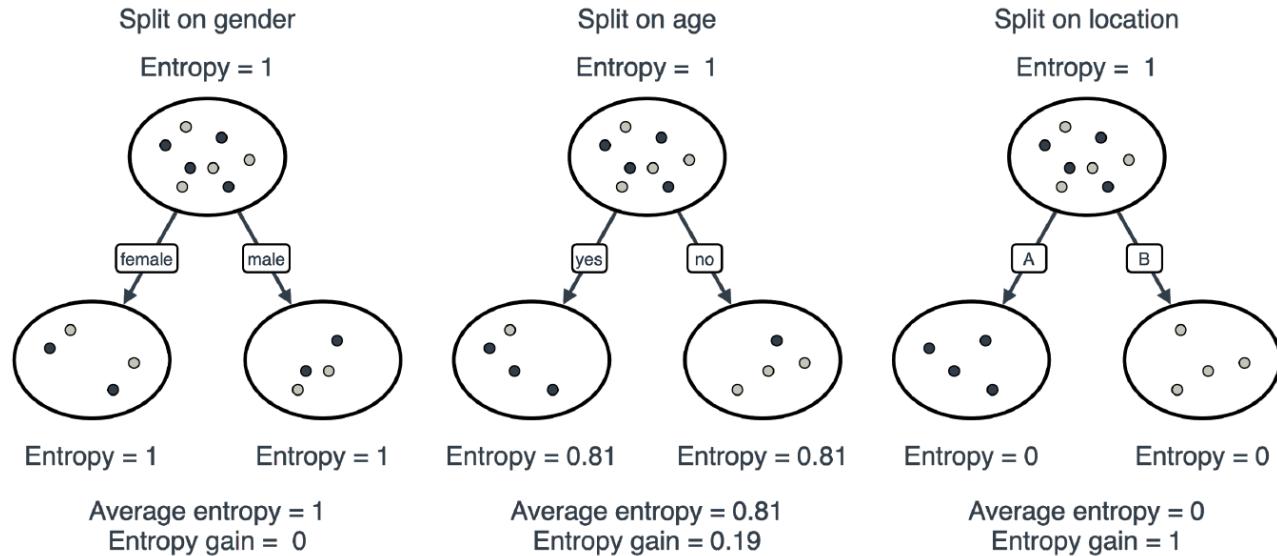
Formula: $E = -p_1 \log_2(p_1) - p_2 \log_2(p_2) - \dots - p_k \log_2(p_k)$.



How to pick the right feature to split

Using **Entropy**:

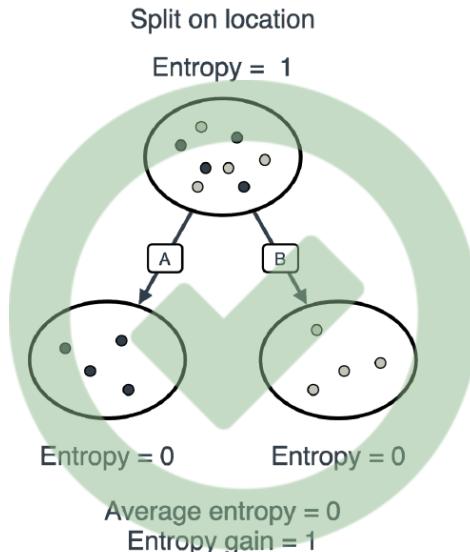
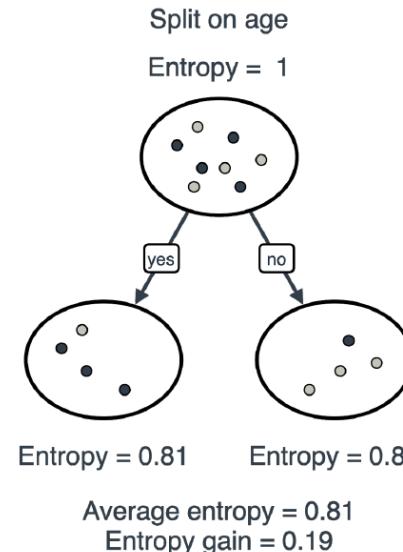
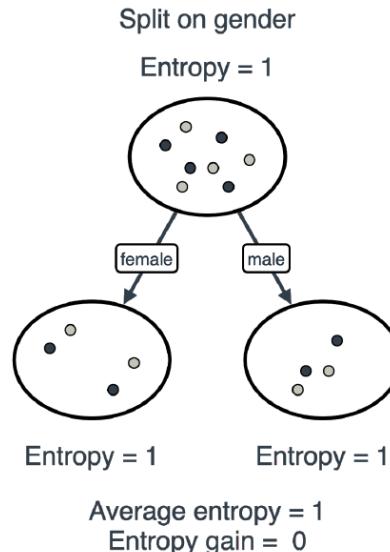
Formula: $E = -p_1 \log_2(p_1) - p_2 \log_2(p_2) - \dots - p_k \log_2(p_k)$.



How to pick the right feature to split

Using **Entropy**:

Formula: $E = -p_1 \log_2(p_1) - p_2 \log_2(p_2) - \dots - p_k \log_2(p_k)$.

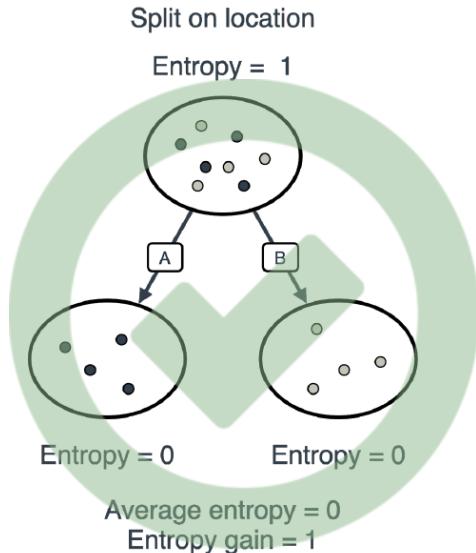
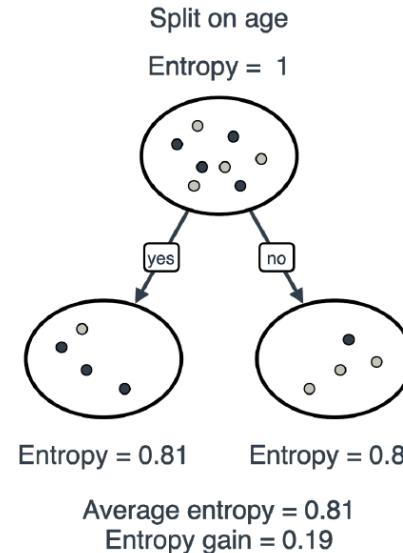
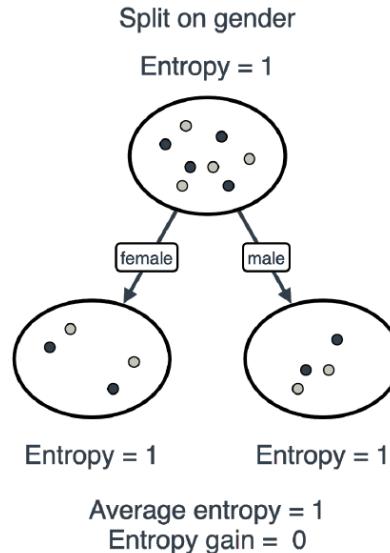


How to pick the right feature to split

Using **Entropy**:

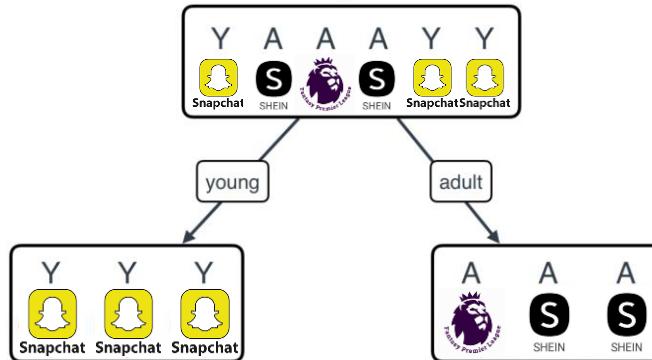
Formula: $E = -p_1 \log_2(p_1) - p_2 \log_2(p_2) - \dots - p_k \log_2(p_k)$.

As we did with the Gini impurity, we can also find the weighted information gain when the leaves don't have the same size.

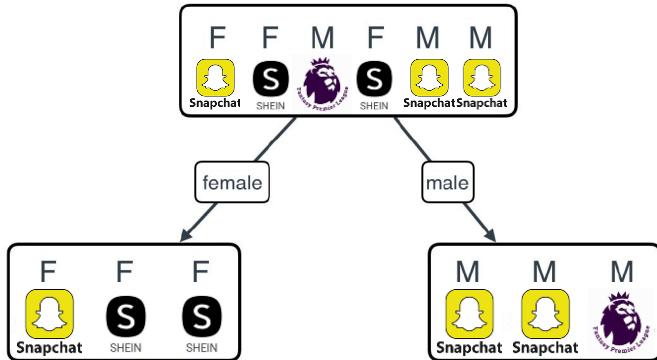


Back To Our Recommendation System

Age?

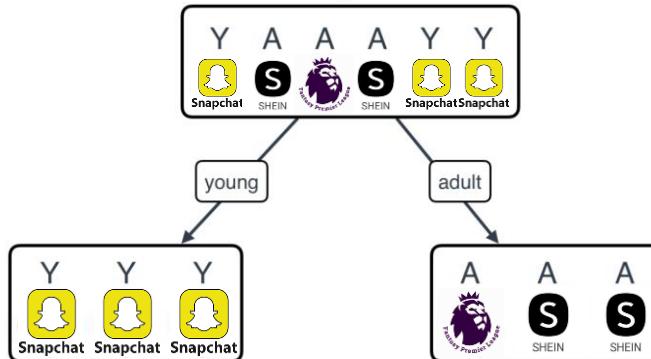


Gender?

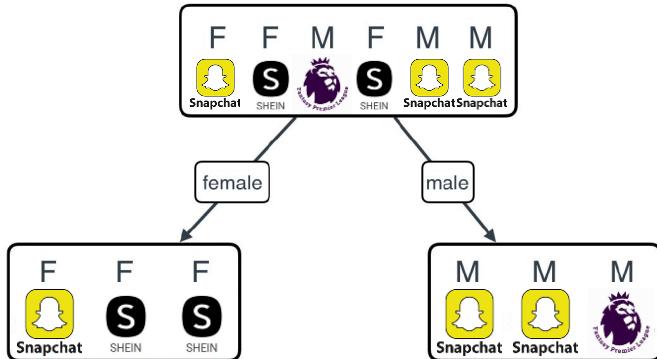


Back To Our Recommendation System

Age?



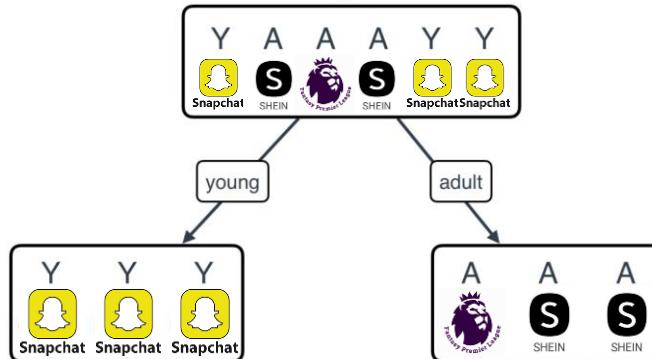
Gender?



Correct 5 out of 6 times

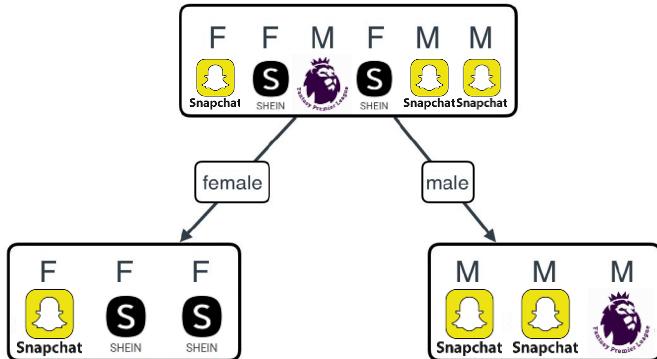
Back To Our Recommendation System

Age?



Correct 5 out of 6 times

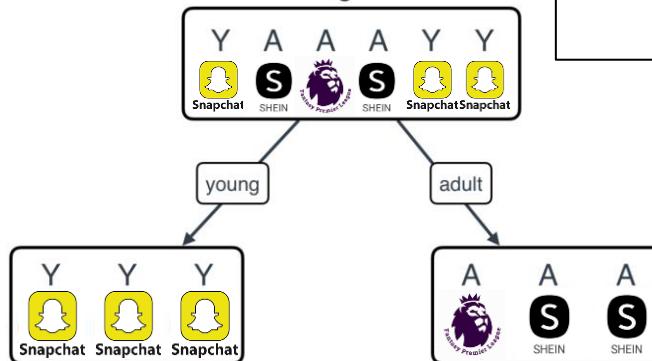
Gender?



Correct 4 out of 6 times

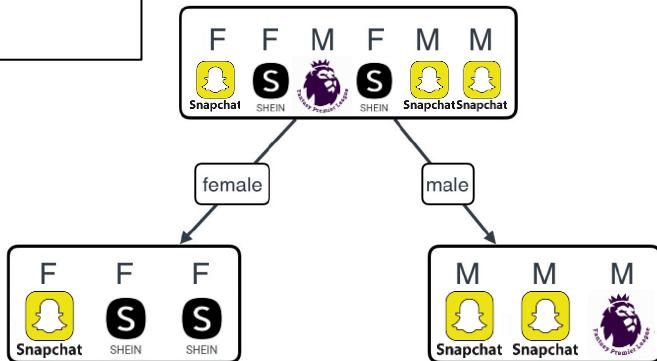
Back To Our Recommendation System

Age?



Gini: 0.611

Gender?



Correct 5 out of 6 times

Correct 4 out of 6 times

Back To Our Recommendation System



Gini: 0.611



Correct 5 out of 6 times

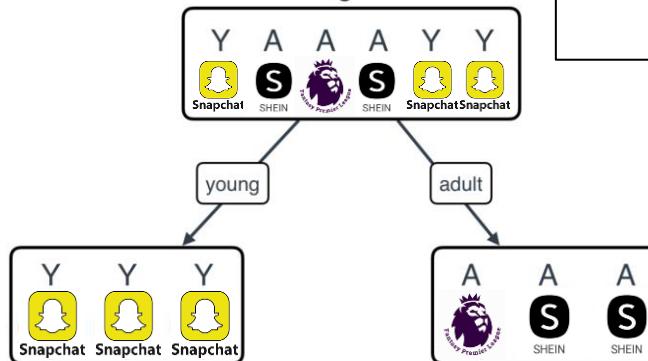


Correct 4 out of 6 times

Gini Average: 0.222 Gain= 0.389

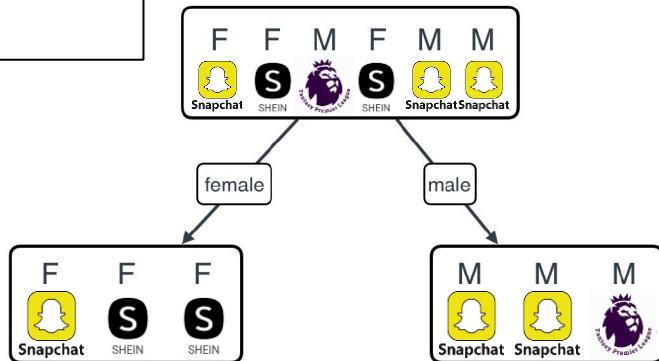
Back To Our Recommendation System

Age?



Gini: 0.611

Gender?



Correct 5 out of 6 times

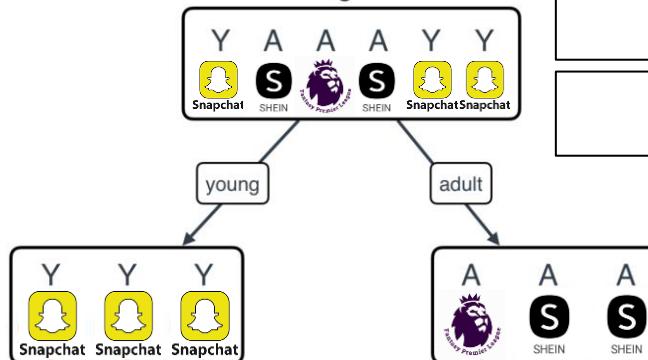
Gini Average: 0.222 Gain= 0.389

Correct 4 out of 6 times

Gini Average: 0.444 Gain= 0.167

Back To Our Recommendation System

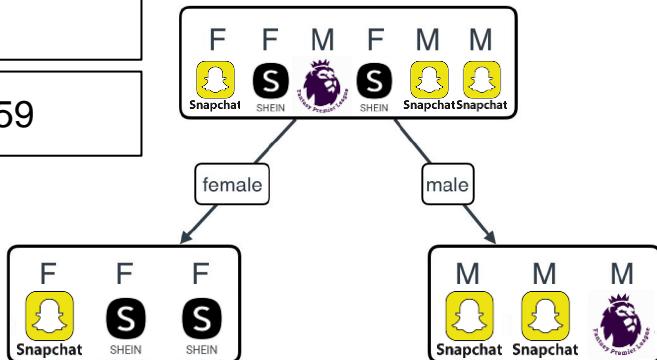
Age?



Gini: 0.611

Entropy: 1.459

Gender?



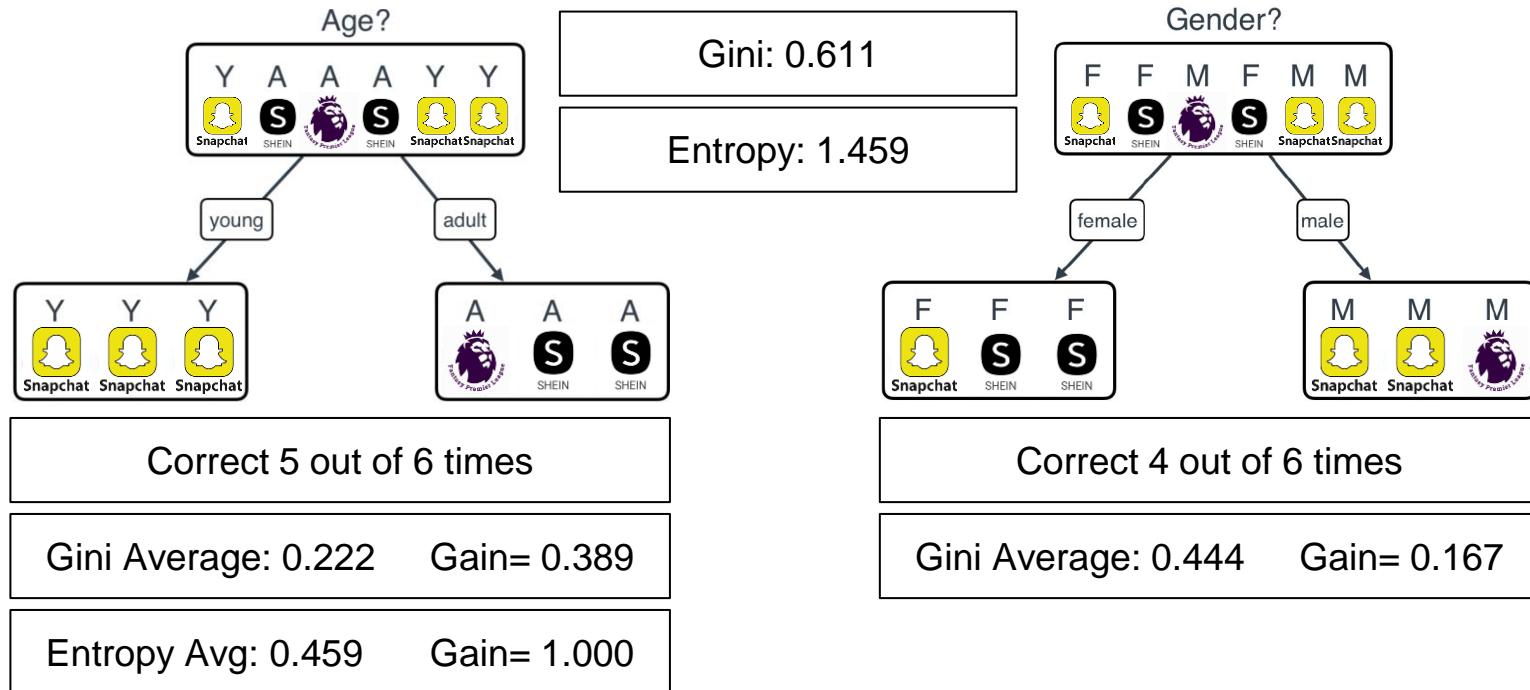
Correct 5 out of 6 times

Gini Average: 0.222 Gain= 0.389

Correct 4 out of 6 times

Gini Average: 0.444 Gain= 0.167

Back To Our Recommendation System



Back To Our Recommendation System



Gini: 0.611

Entropy: 1.459



female

male



Correct 5 out of 6 times

Gini Average: 0.222 Gain= 0.389

Entropy Avg: 0.459 Gain= 1.000

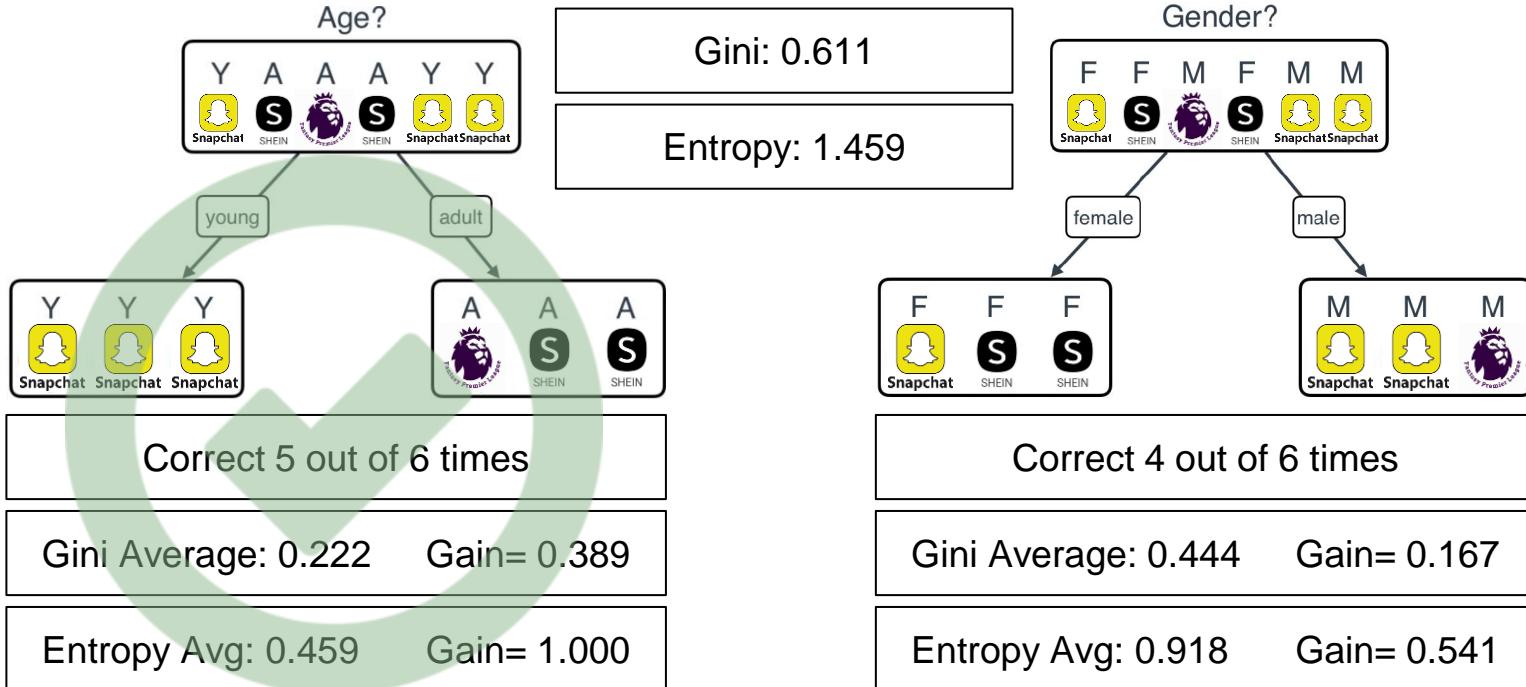


Correct 4 out of 6 times

Gini Average: 0.444 Gain= 0.167

Entropy Avg: 0.918 Gain= 0.541

Back To Our Recommendation System



Lecture Overview

Simple Decision Tree

Mobile Apps Recommendation Example

Split Using Accuracy

Split Using Gini Impurity

Split Using Entropy

Tree Hyperparameters

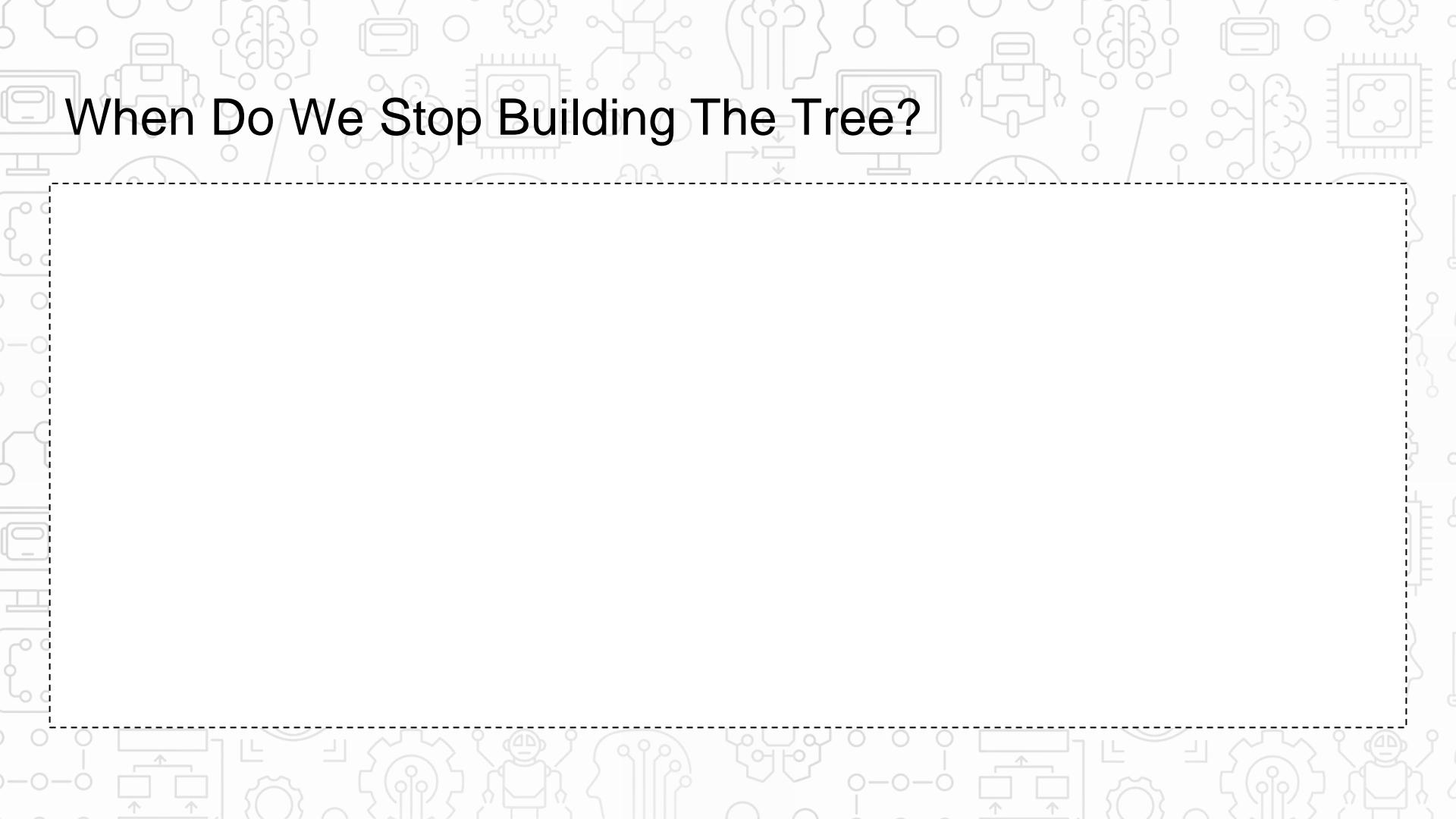
One-hot Encoded Features

Continues Features

Decision Tree For Regression

Random Forest

When Do We Stop Building The Tree?



When Do We Stop Building The Tree?

These are hyperparameters:

- The gain of Gini index or Entropy is below some threshold.
- If node has less than a certain number of elements.
- Reach a certain number of levels.

Lecture Overview

Simple Decision Tree

Mobile Apps Recommendation Example

Split Using Accuracy

Split Using Gini Impurity

Split Using Entropy

Tree Hyperparameters

One-hot Encoded Features

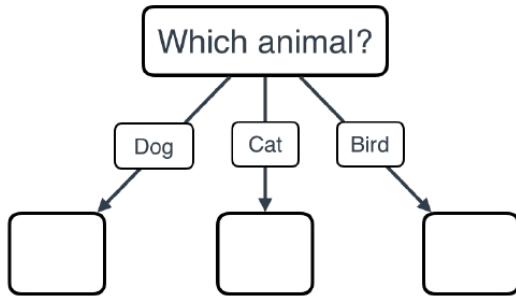
Continues Features

Decision Tree For Regression

Random Forest

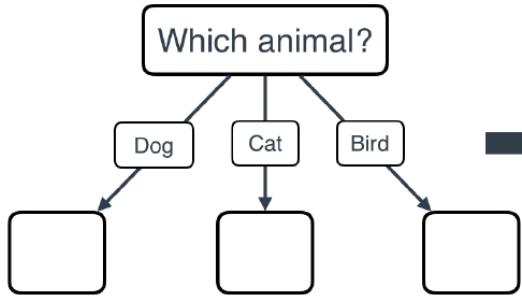
Features With More Than Two Categories

Non-binary feature

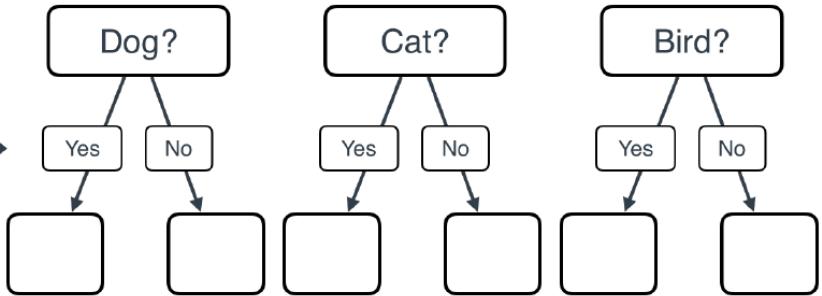


Features With More Than Two Categories

Non-binary feature

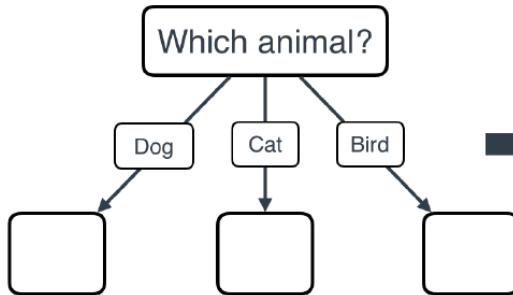


More binary features

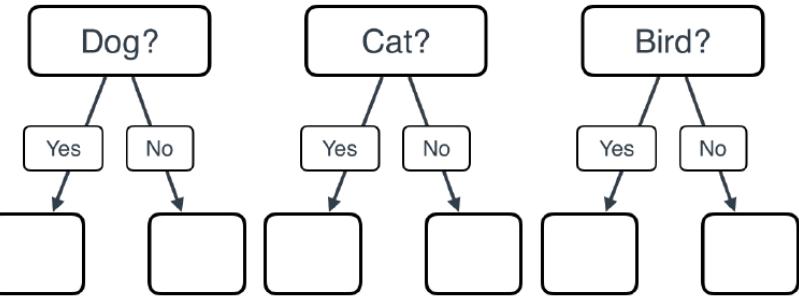


Features With More Than Two Categories

Non-binary feature



More binary features



Animal
Dog
Cat
Bird
Dog
Bird

One-hot
encoding

Dog?	Cat?	Bird?
1	0	0
0	1	0
0	0	1
1	0	0
0	0	1

Lecture Overview

Simple Decision Tree

Mobile Apps Recommendation Example

Split Using Accuracy

Split Using Gini Impurity

Split Using Entropy

Tree Hyperparameters

One-hot Encoded Features

Continues Features

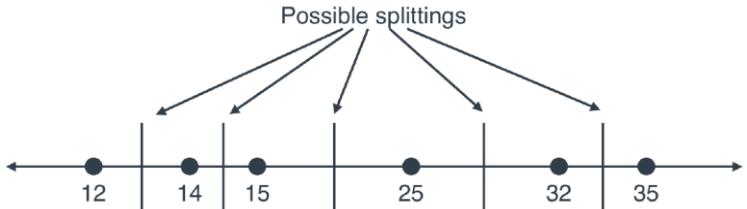
Decision Tree For Regression

Random Forest

Continuous Features

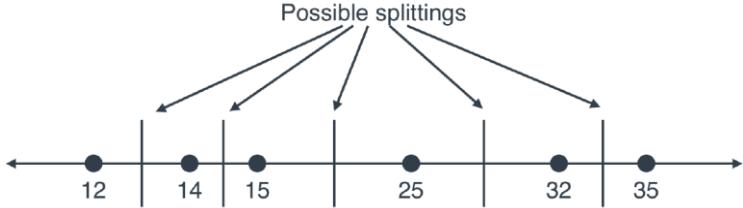
Gender	Age	App
F	15	 Snapchat
F	25	 SHEIN
M	32	 Fantasy Premier League
F	35	 SHEIN
M	12	 Snapchat
M	14	 Snapchat

Continuous Features



Gender	Age	App
F	15	Snapchat
F	25	SHEIN
M	32	Fantasy Premier League
F	35	SHEIN
M	12	Snapchat
M	14	Snapchat

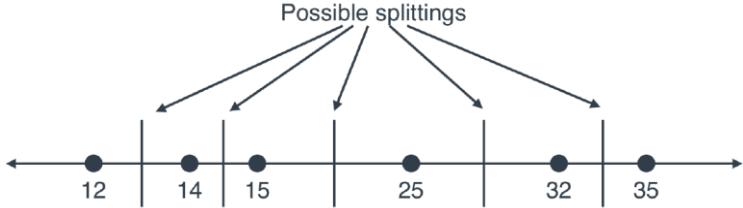
Continuous Features



Question	First set	Second set
Is the user younger than 13?	12	14, 15, 25, 32, 35
Is the user younger than 14.5?	12, 14	15, 25, 32, 35
Is the user younger than 20?	12, 14, 15	25, 32, 35
Is the user younger than 27?	12, 14, 15, 25	32, 35
Is the user younger than 33?	12, 14, 15, 25, 32	35

Gender	Age	App
F	15	Snapchat
F	25	SHEIN
M	32	Fantasy Premier League
F	35	SHEIN
M	12	Snapchat
M	14	Snapchat

Continuous Features



Question	First set	Second set
Is the user younger than 13?	12	14, 15, 25, 32, 35
Is the user younger than 14.5?	12, 14	15, 25, 32, 35
Is the user younger than 20?	12, 14, 15	25, 32, 35
Is the user younger than 27?	12, 14, 15, 25	32, 35
Is the user younger than 33?	12, 14, 15, 25, 32	35

Add all these 5 questions to the age question and use accuracy, Gini-index or entropy to choose the question to split

Gender	Age	App
F	15	Snapchat
F	25	SHEIN
M	32	Fantasy Premier League
F	35	SHEIN
M	12	Snapchat
M	14	Snapchat

Lecture Overview

Simple Decision Tree

Mobile Apps Recommendation Example

Split Using Accuracy

Split Using Gini Impurity

Split Using Entropy

Tree Hyperparameters

One-hot Encoded Features

Continues Features

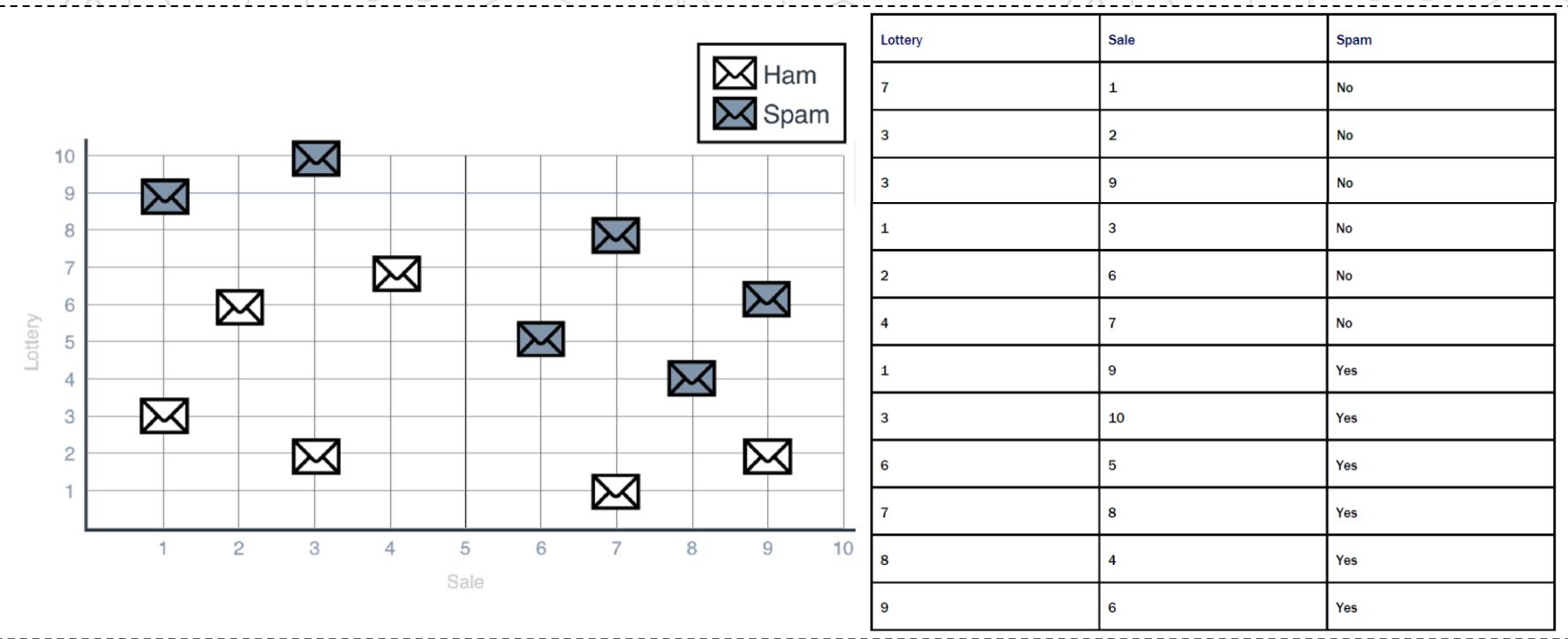
Decision Tree For Regression

Random Forest

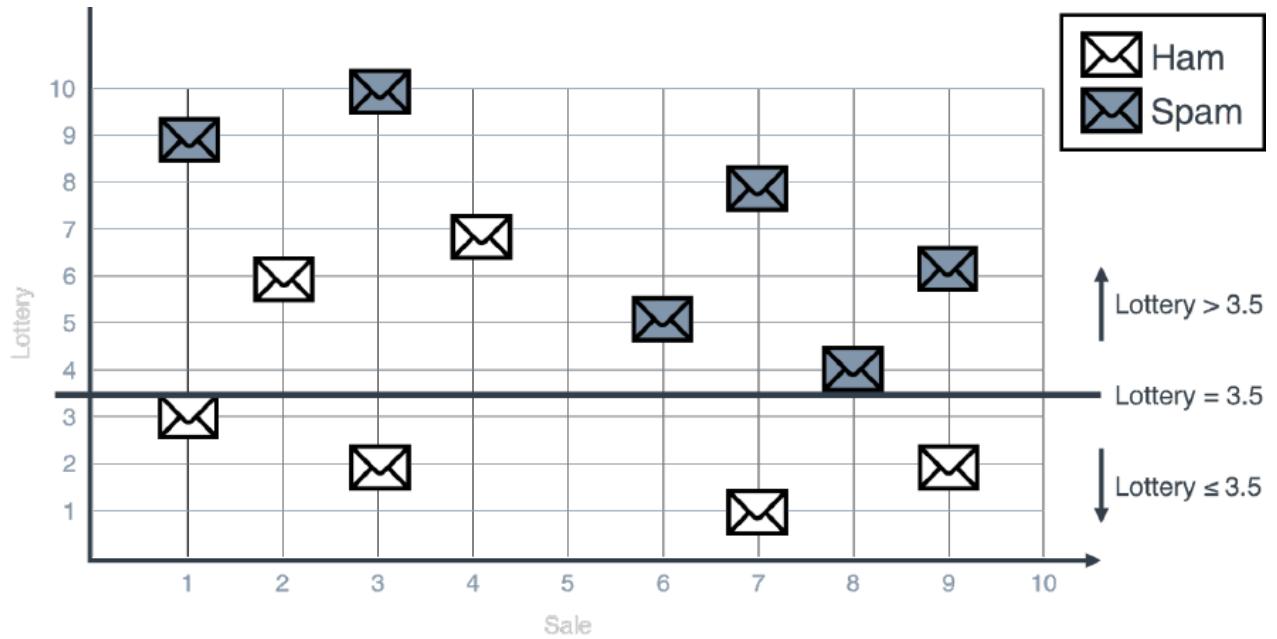
Graphical Example: Spam Detector

Lottery	Sale	Spam
7	1	No
3	2	No
3	9	No
1	3	No
2	6	No
4	7	No
1	9	Yes
3	10	Yes
6	5	Yes
7	8	Yes
8	4	Yes
9	6	Yes

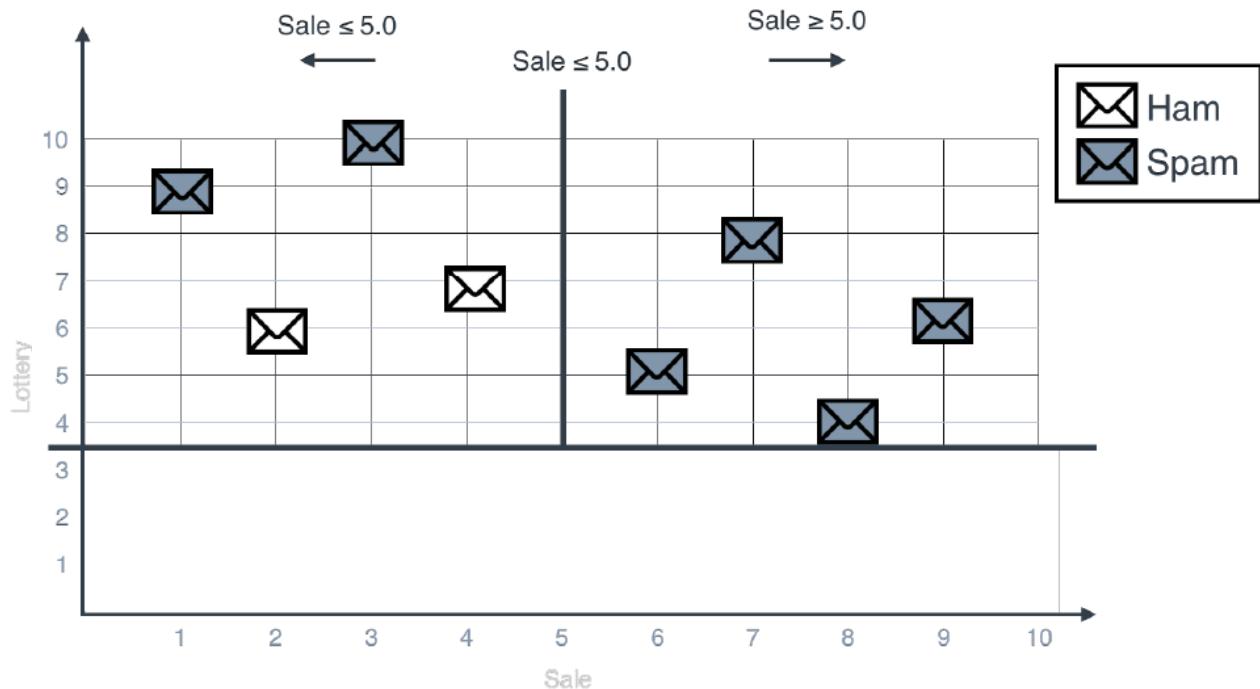
Graphical Example: Spam Detector



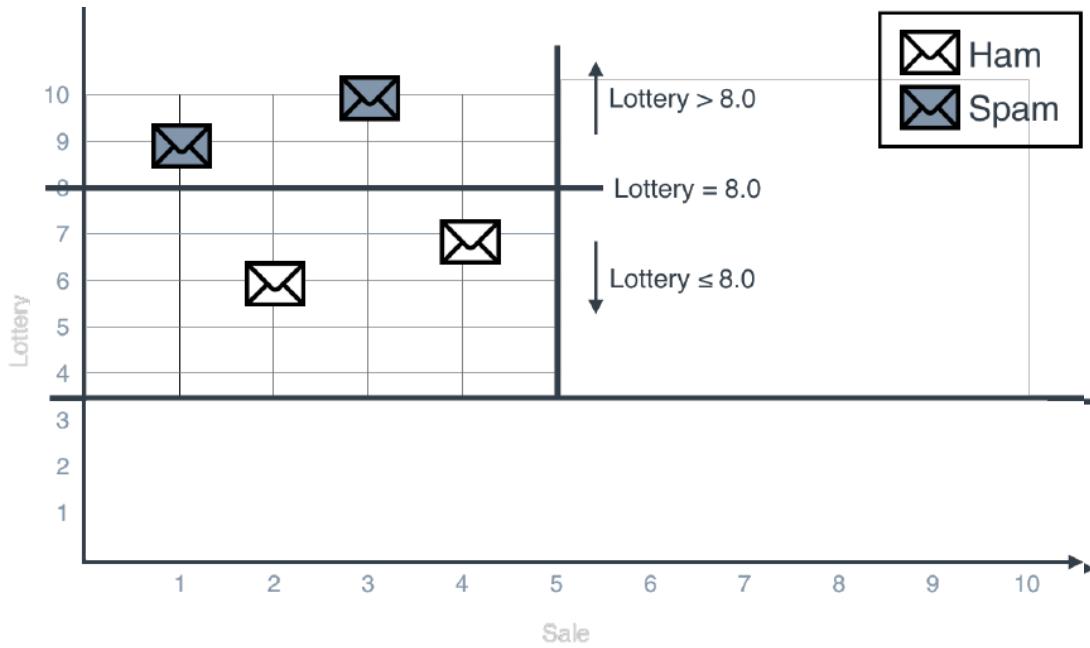
Graphical Example: Spam Detector



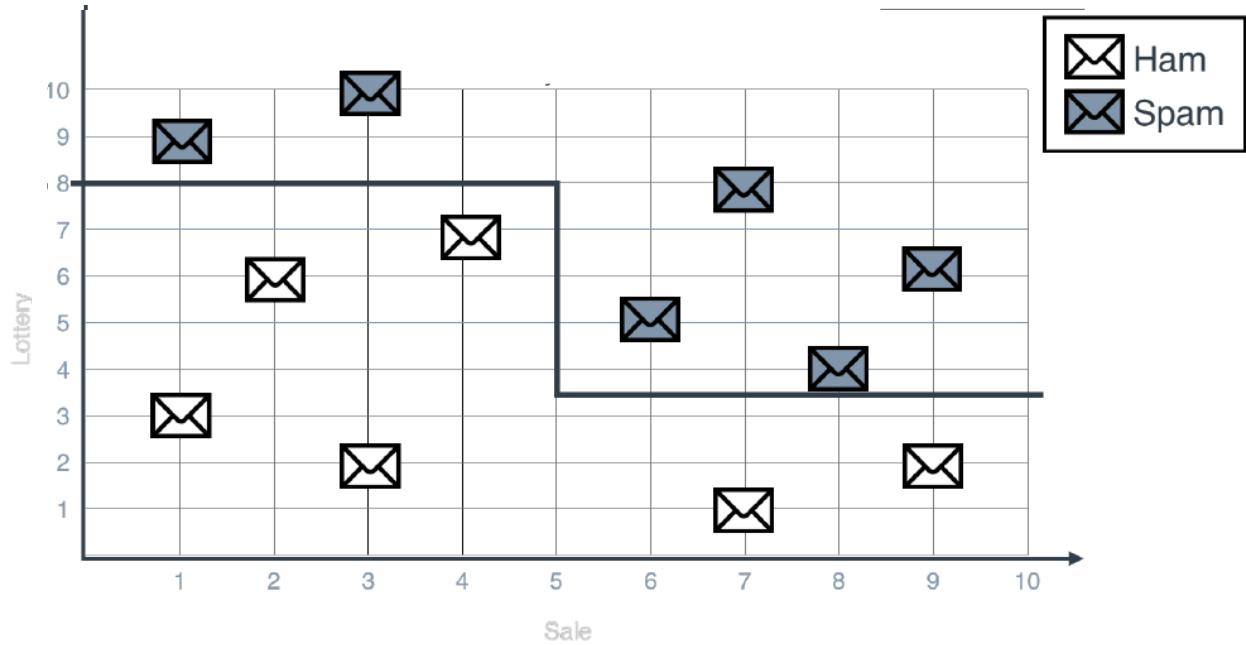
Graphical Example: Spam Detector



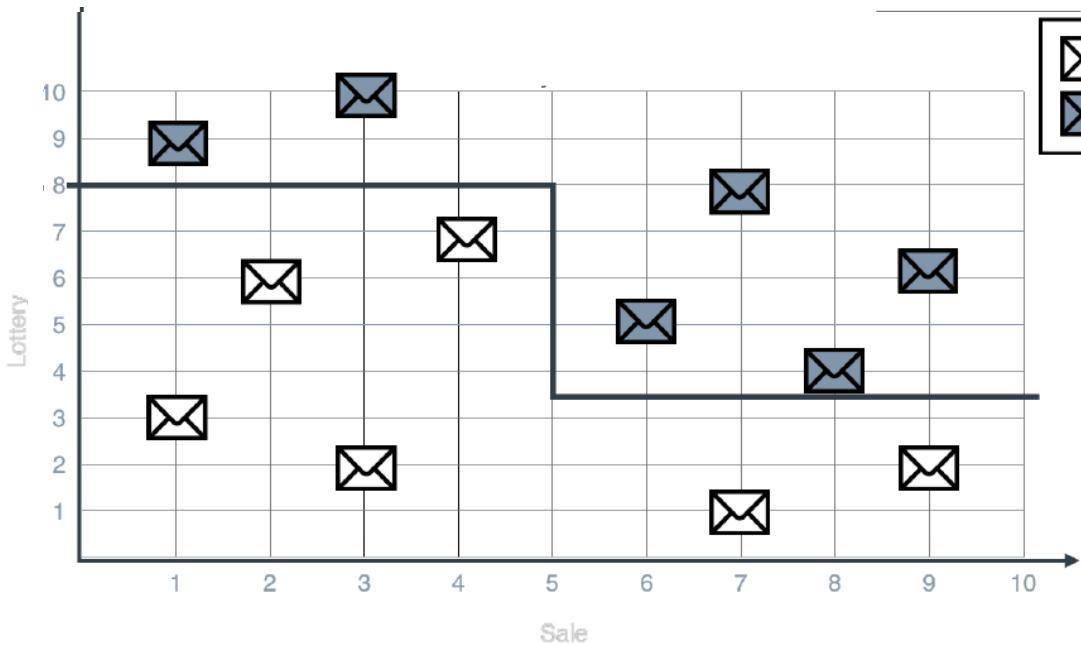
Graphical Example: Spam Detector



Graphical Example: Spam Detector



Graphical Example: Spam Detector



*Each Stub
represents
vertical or
horizontal line*

Decision Tree For Classification



Use colab to open this github notebook:

[“s7s/machine_learning_1/decision_trees/decision_trees.ipynb”](https://colab.research.google.com/github/s7s/machine_learning_1/blob/main/decision_trees/decision_trees.ipynb)

Decision Trees For Regression

Decision Trees For Regression

For regression,

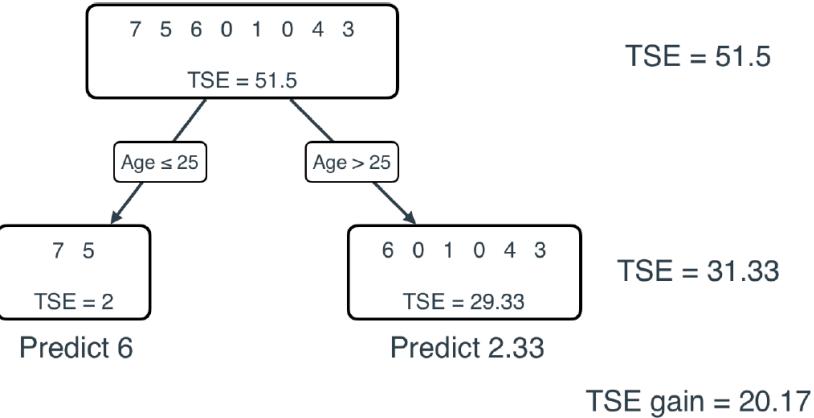
- The decision is **the average** of the set.
- The error is **Total Square Error(TSE)** from the average.
- The best question that gives **the least error**.

Decision Trees For Regression

Decision Trees For Regression

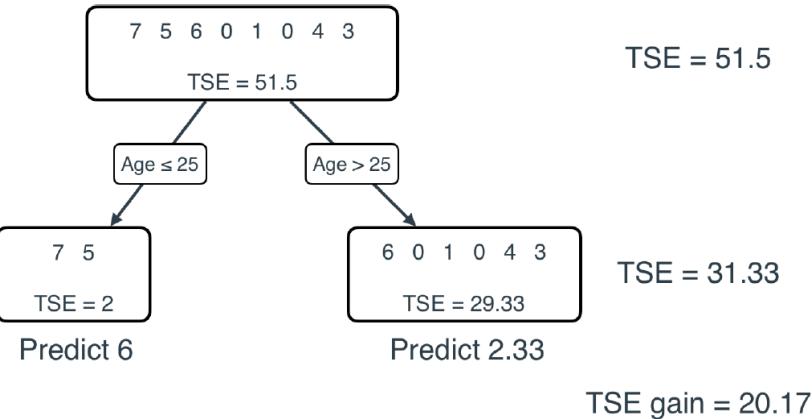
Age	Engagement
10	7
20	5
30	6
40	0
50	1
60	0
70	4
80	3

Decision Trees For Regression



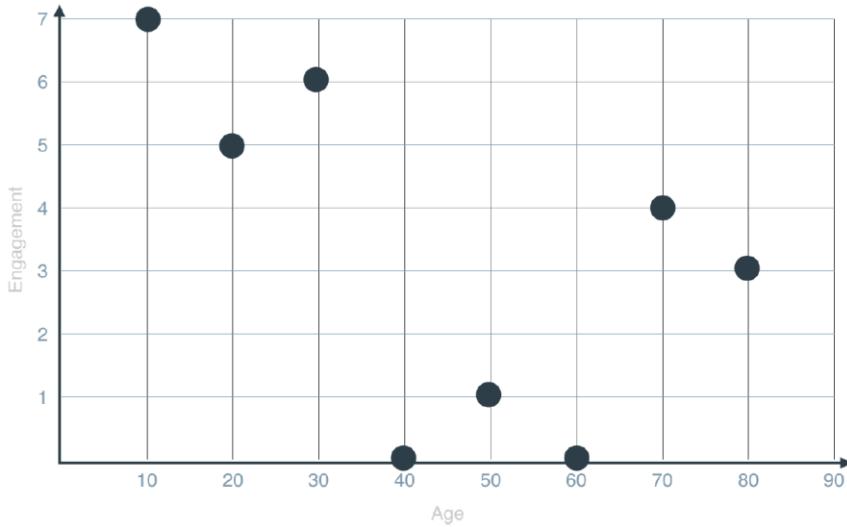
Age	Engagement
10	7
20	5
30	6
40	0
50	1
60	0
70	4
80	3

Decision Trees For Regression

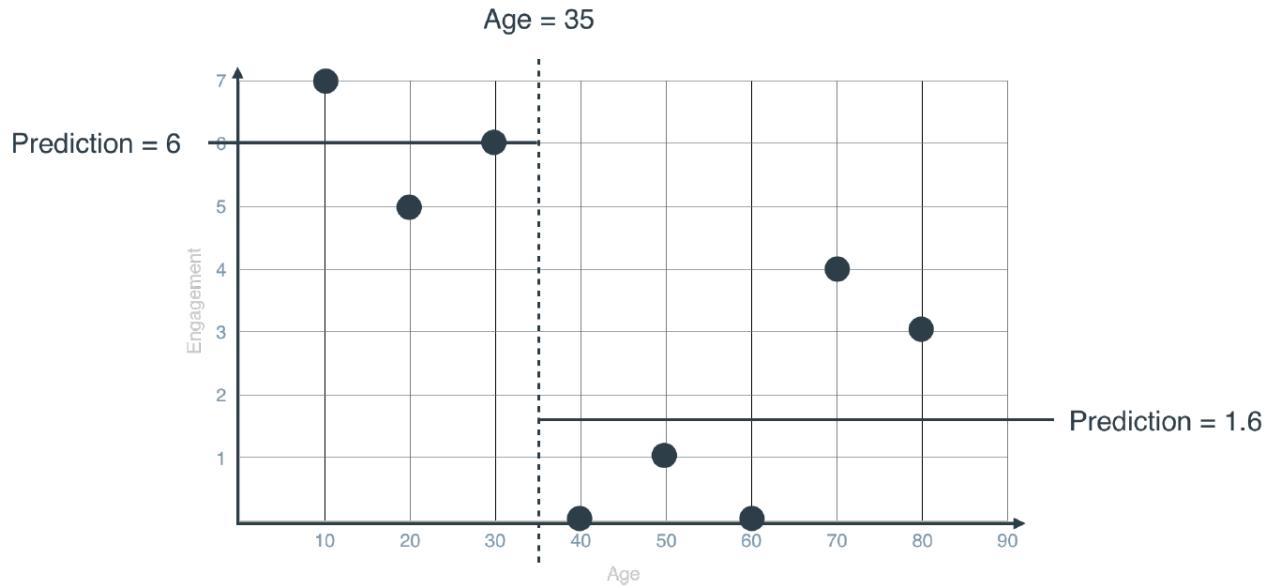


Age	Engagement
10	7
20	5
30	6
40	0
50	1
60	0
70	4
80	3

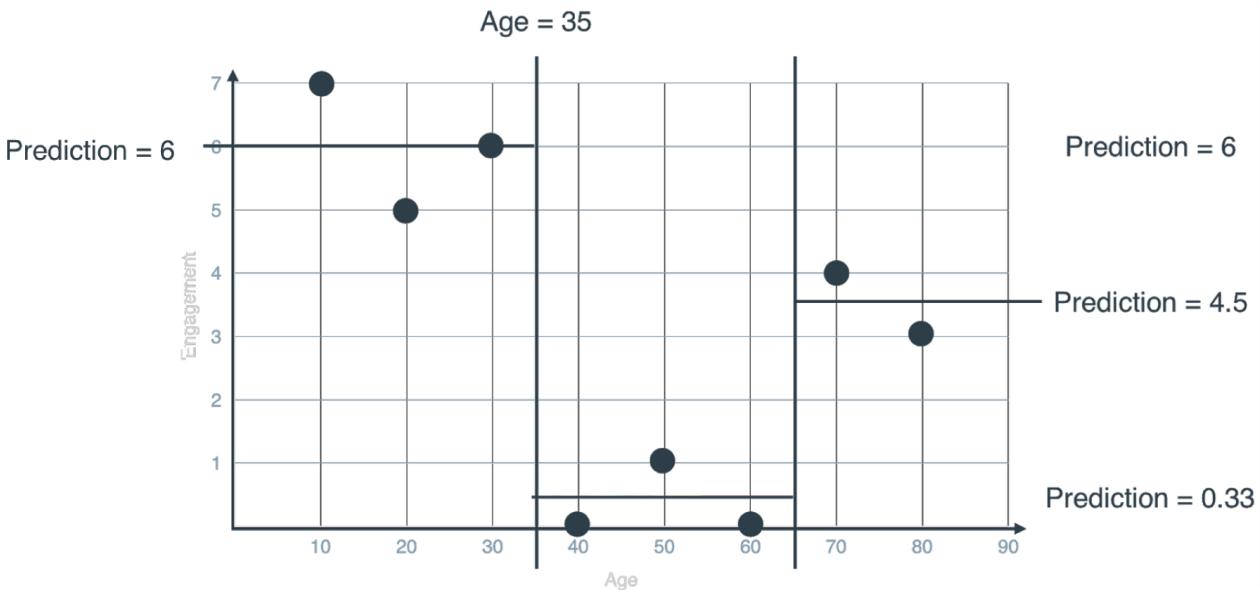
Decision Trees For Regression



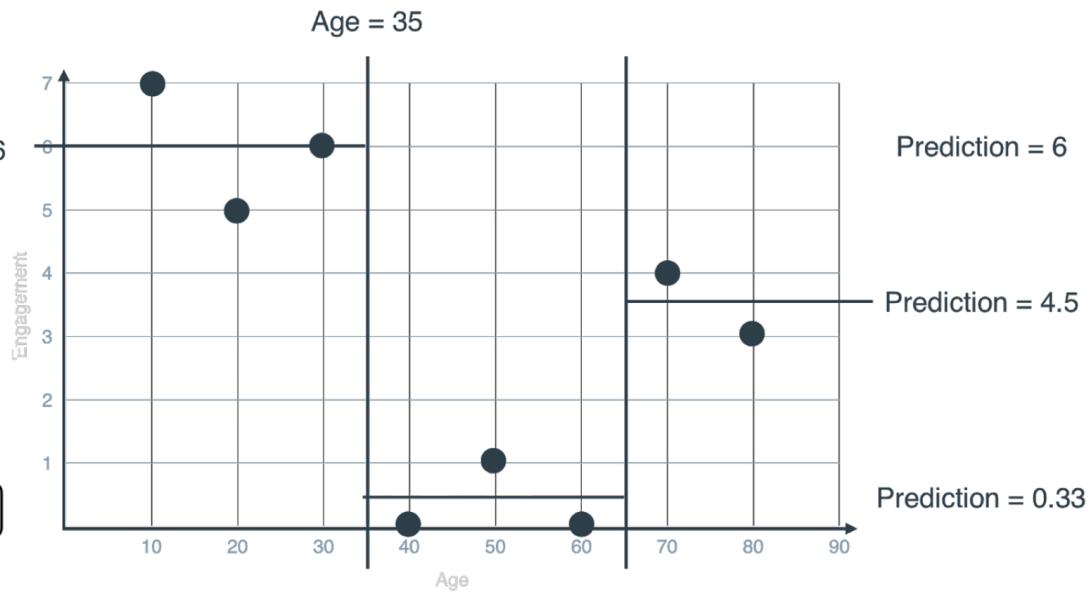
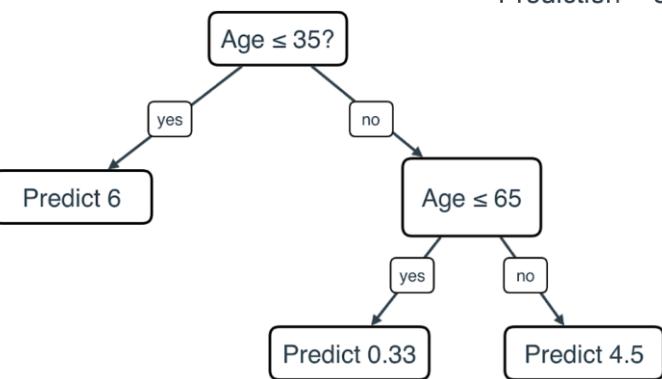
Decision Trees For Regression



Decision Trees For Regression



Decision Trees For Regression



Decision Tree For Regression



Use colab to open this github notebook:

[“s7s/machine_learning_1/decision_trees/decision_trees.ipynb”](https://colab.research.google.com/github/s7s/machine_learning_1/blob/main/decision_trees/decision_trees.ipynb)

Lecture Overview

Simple Decision Tree

Mobile Apps Recommendation Example

Split Using Accuracy

Split Using Gini Impurity

Split Using Entropy

Tree Hyperparameters

One-hot Encoded Features

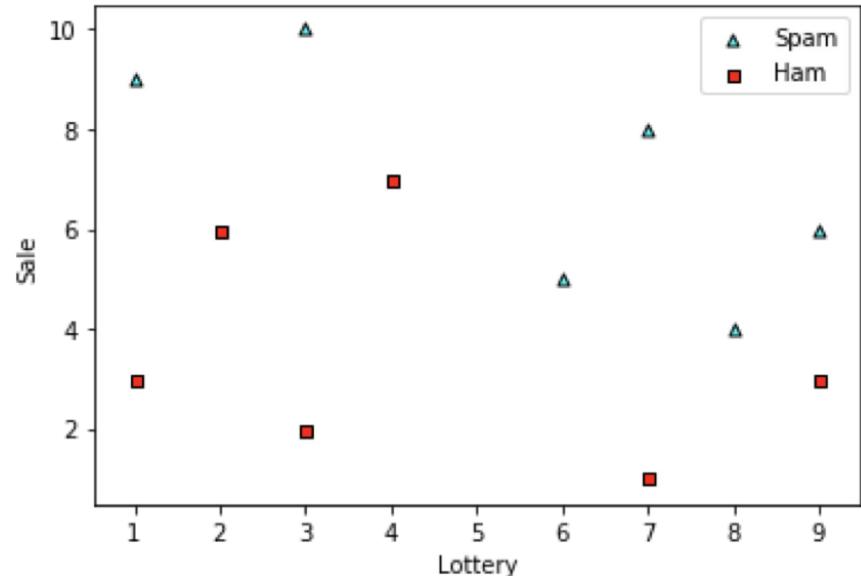
Continues Features

Decision Tree For Regression

Random Forest

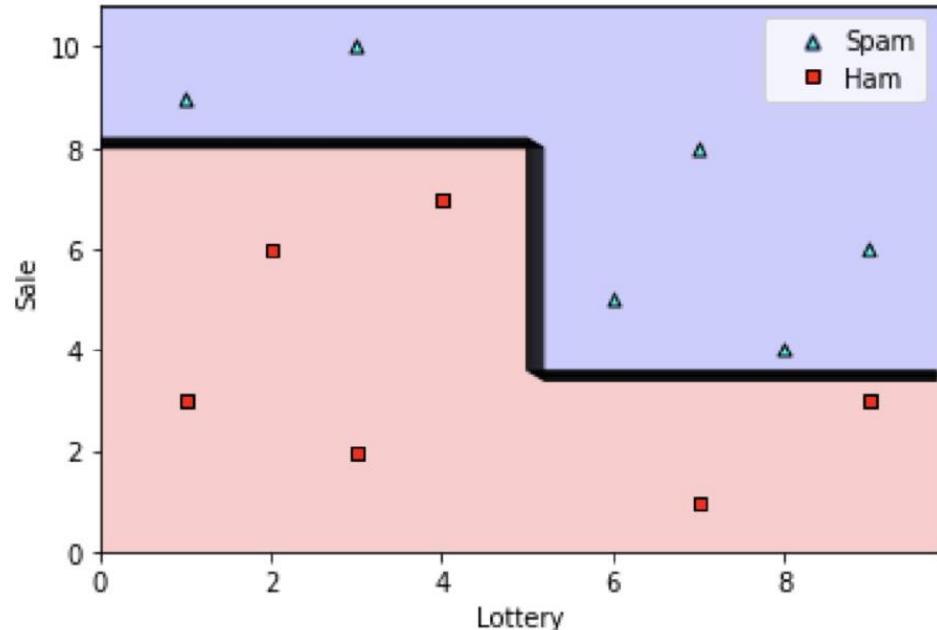
Random Forest

Using normal Decision Tree



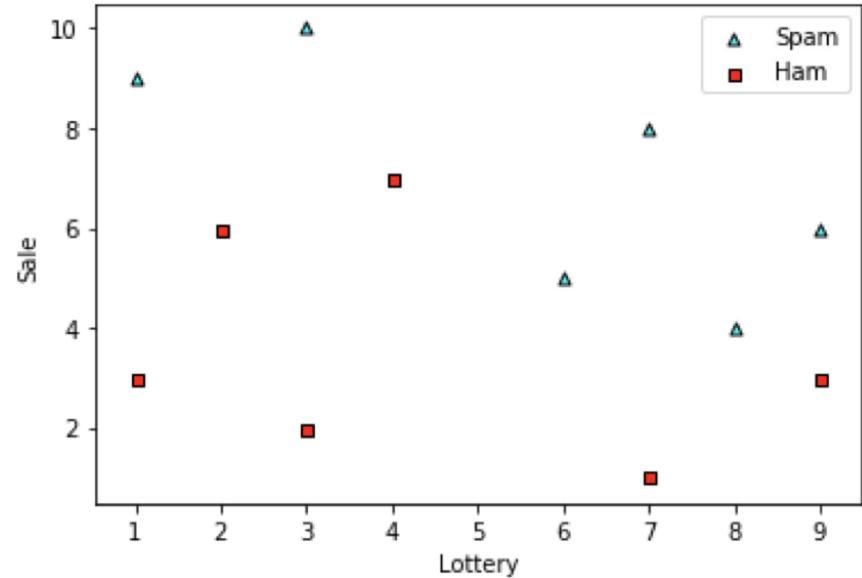
Random Forest

Using normal Decision Tree



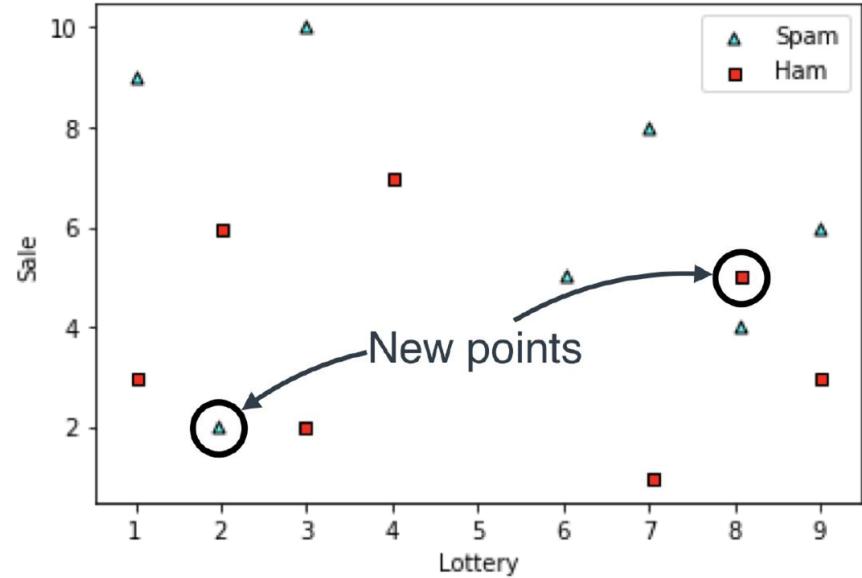
Random Forest

Using normal Decision Tree



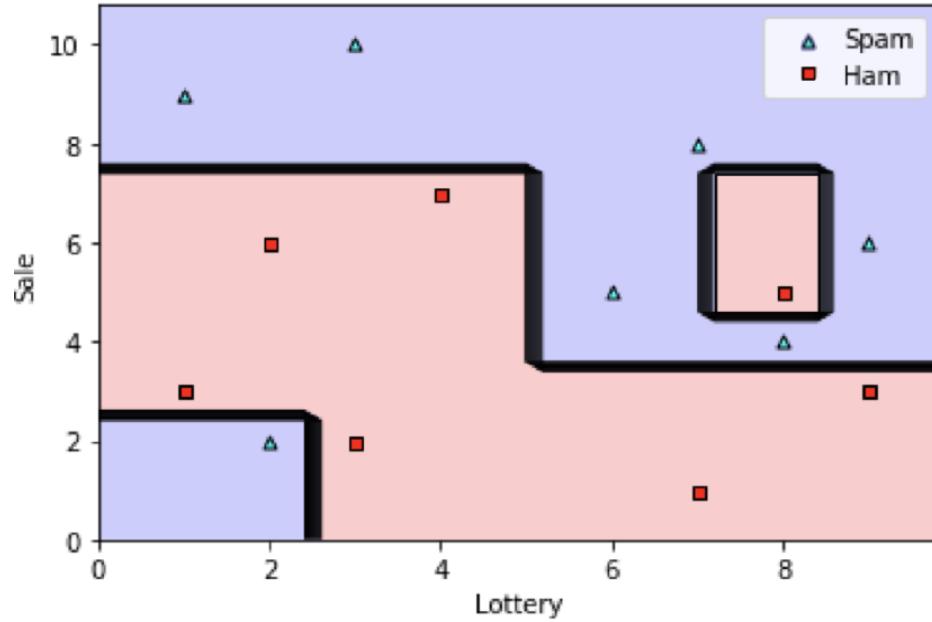
Random Forest

Using normal Decision Tree



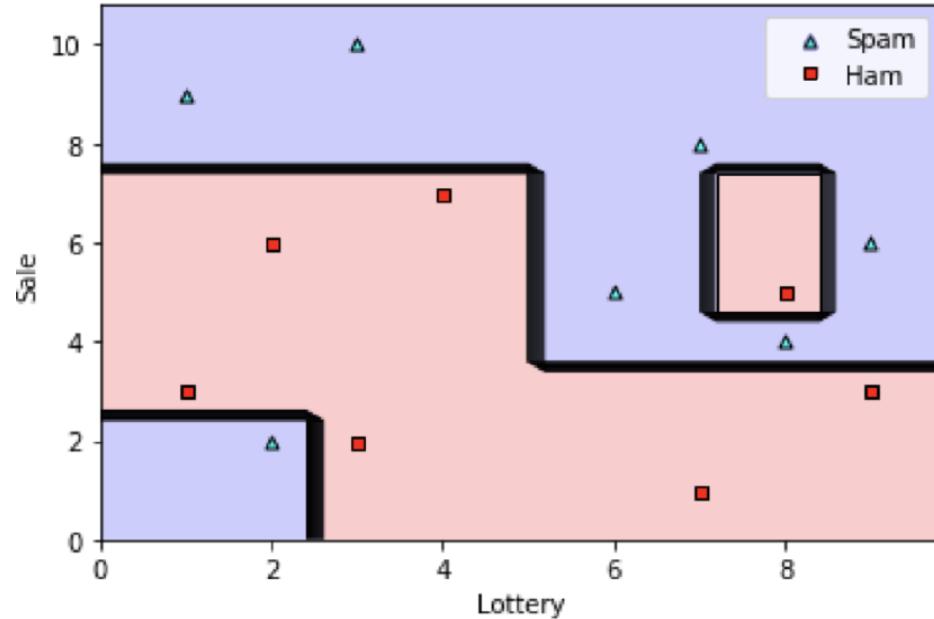
Random Forest

Using normal Decision Tree



Random Forest

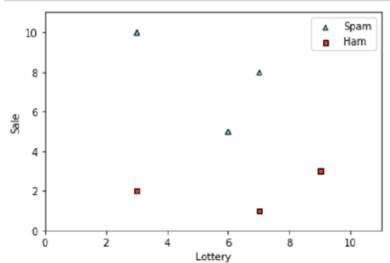
Using normal Decision Tree



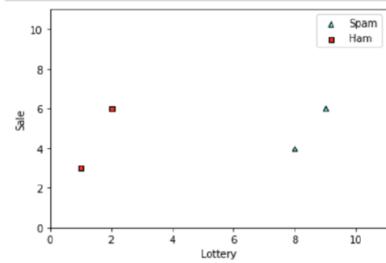
Do you notice the overfitting?

Random Forest

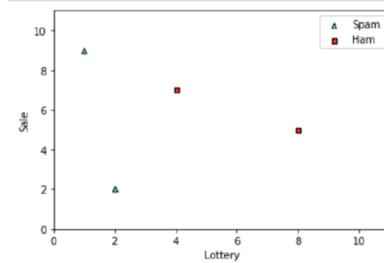
Subset 1



Subset 2



Subset 3

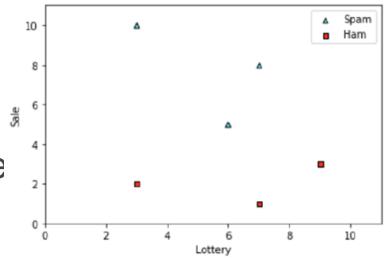


Random Forest

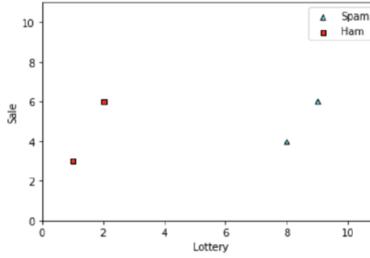
Random Forest Steps:

- Split data randomly into sets:
- Train each decision tree at one of these sets.
- Take decision based on vote between these models.

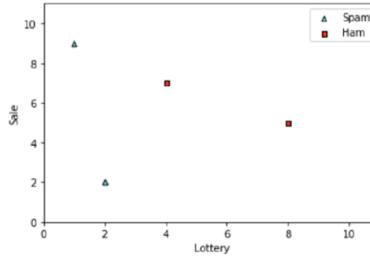
Subset 1



Subset 2



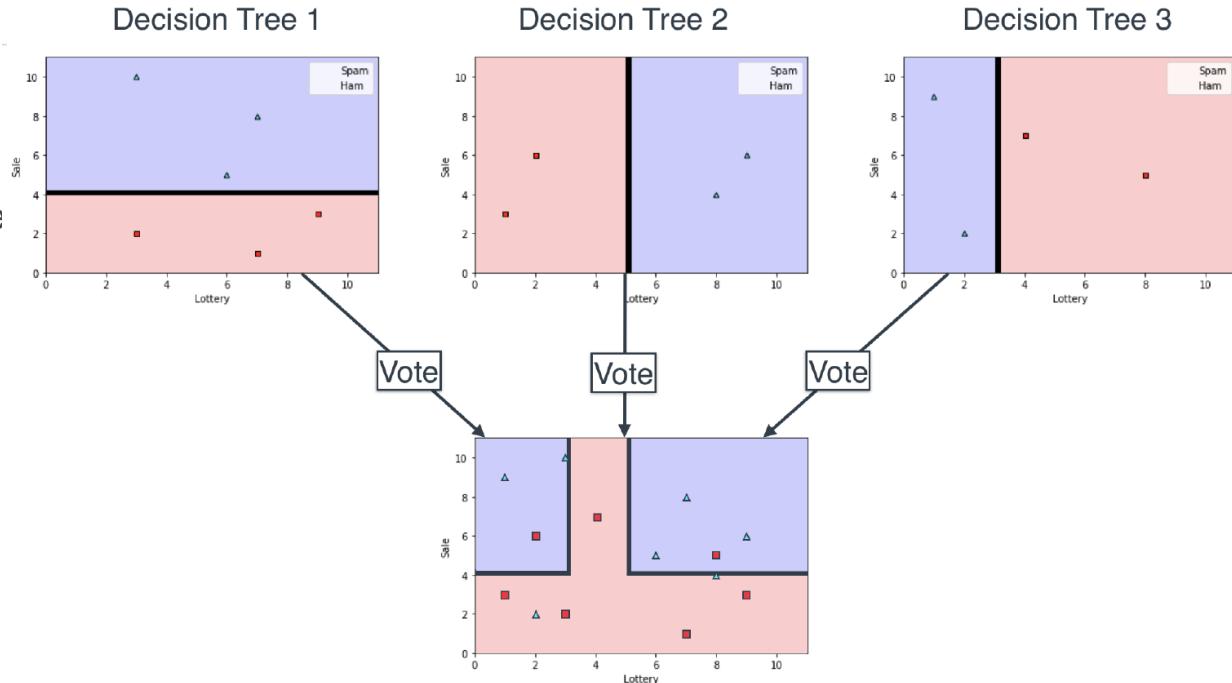
Subset 3



Random Forest

Random Forest Steps:

- Split data randomly into sets:
- Train each decision tree at one of these sets.
- Take decision based on vote between these models.



Random Forest



Use colab to open this github notebook:

[“s7s/machine_learning_1/decision_trees/random_forest.ipynb”](https://colab.research.google.com/github/s7s/machine_learning_1/blob/main/decision_trees/random_forest.ipynb)

Lecture Overview

Simple Decision Tree

Mobile Apps Recommendation Example

Split Using Accuracy

Split Using Gini Impurity

Split Using Entropy

Tree Hyperparameters

One-hot Encoded Features

Continues Features

Decision Tree For Regression

Random Forest



Feedback