

## Decision-trees for classification

## Reflections on Al Al Ethics

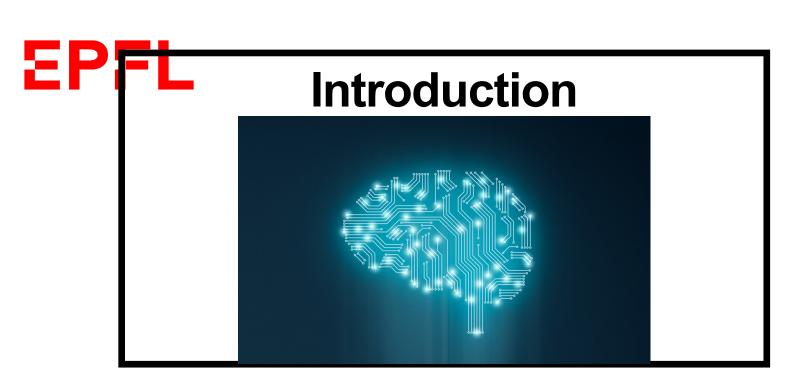
02.12.2024

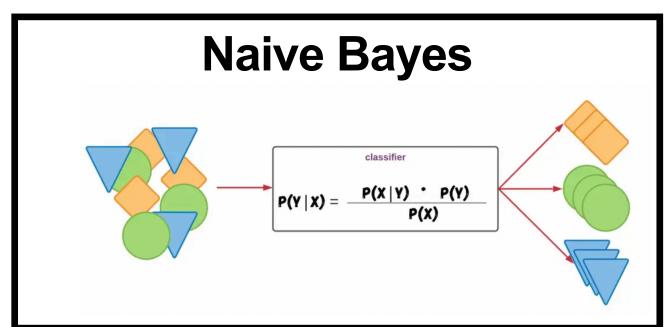


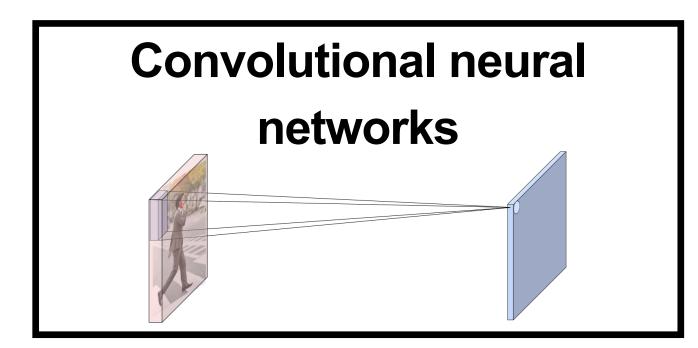
#### Outline

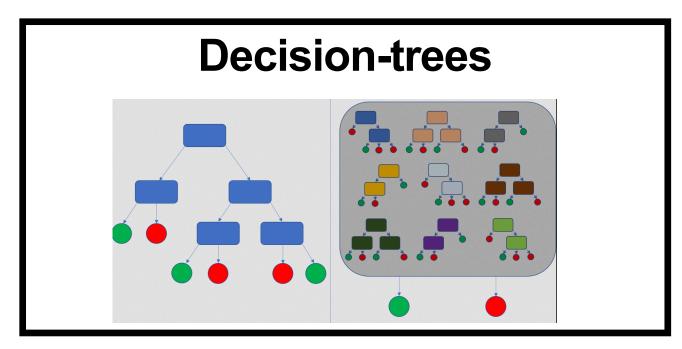
- Hour 1
  - Introduction on AI reflections/ethics
  - Decision-tree continued

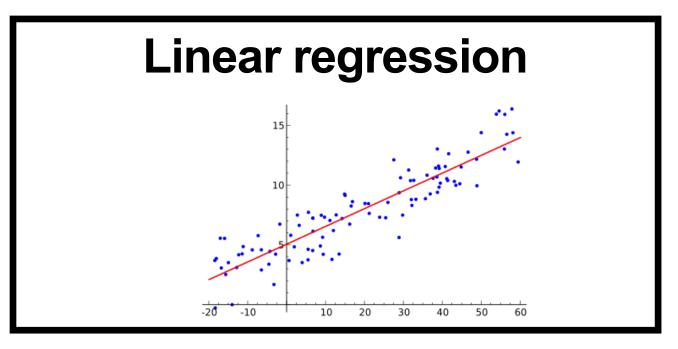
- Hour 2: Reflections on Al
  - Societal challenges
  - How could AI help/hurt?

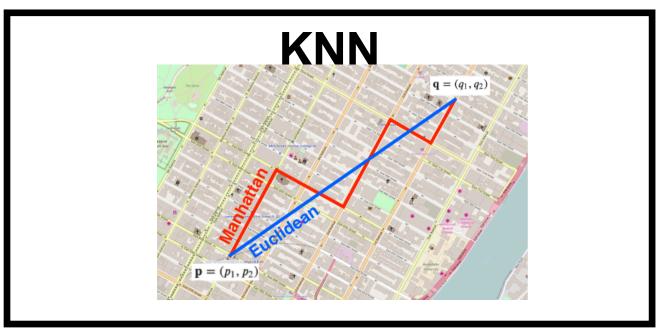


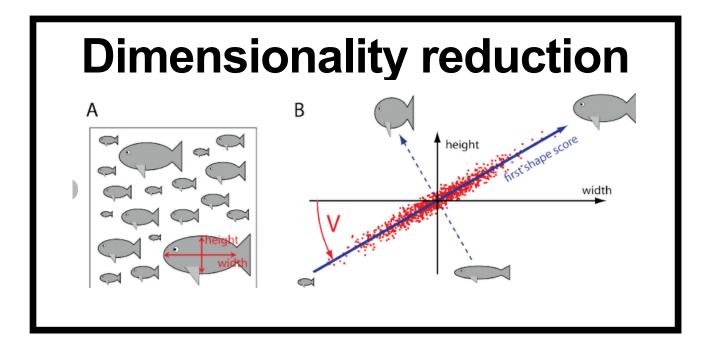


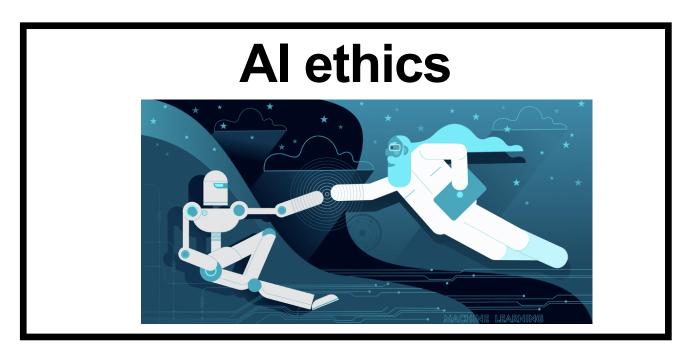


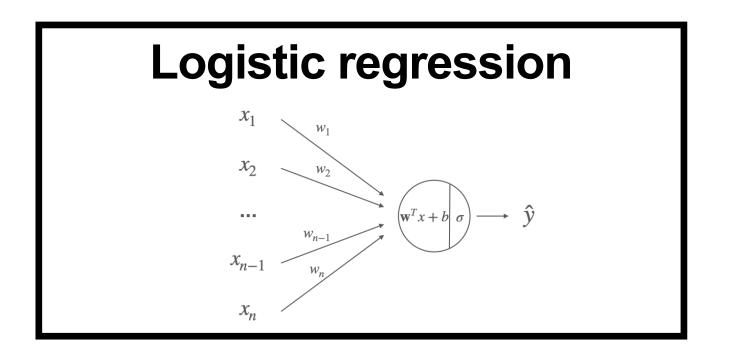


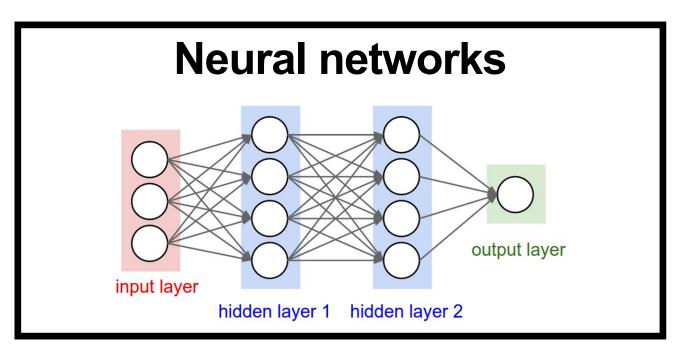


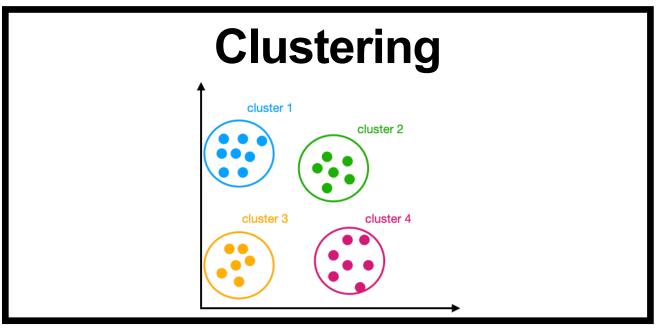


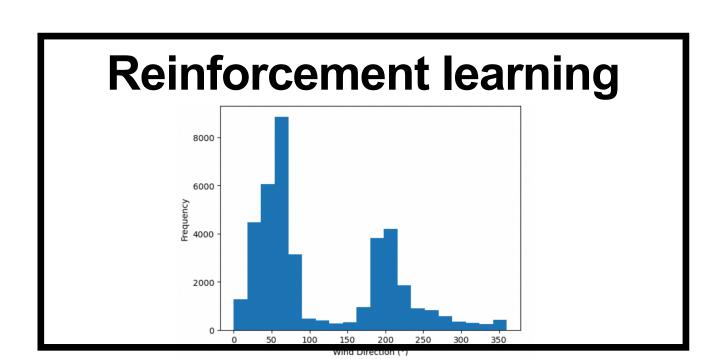












### EPFL Al and its societal impacts

As engineers in this era we might use or develop Al-based tools

We need to understand the impact of our choices in designing a tool, using a dataset, and deploying a tool, examples:

- Why should I use a ML-based approach?
- Who will benefit from my product?
- Are there biases in the data, or the methods I have used?
- Can I explain the decisions of the approach I have used?
- What are some harmful consequences?
- What are the environmental impacts of my choice?



https://www.nature.com/articles/d41586-024-01184-4



https://keefeandkeefe.com/

### **EPFL** Critical thinking

Even broader than our engineering applications...

Who is deciding that the Al narrative should dominate our discussions, courses, approaches?

 Who is deciding how I should be spending my moments, what news I should read, how I should interact with my friends, what movies/clips to watch, where to shop, what food to eat, etc?

What are some benefits and potential dangers of using Al in my courses and work?

How can I use/develop AI to empower myself and others, while being aware of its challenges?

#### **EPFL**

#### Plan for reflection and ethics in our course

#### 1. Discussion on conditions for AI to benefit our societies

Today in class with Prof. Sascha Nick

https://forms.gle/qrWcMXppq9kZ4EK6A

#### 2. Videos to watch at home (link also on our Moodle page) https://mediaspace.epfl.ch/channel/MOOCS+Ethique+num%C3%A9rique/56510

Prepared by Dr. Johan Rochel

#### 3. Bonus assignment (2 points)

To be posted on 04.12.2024 and to return on Moodle by 20.12.2024

### **EPFL** Decision Trees for classification

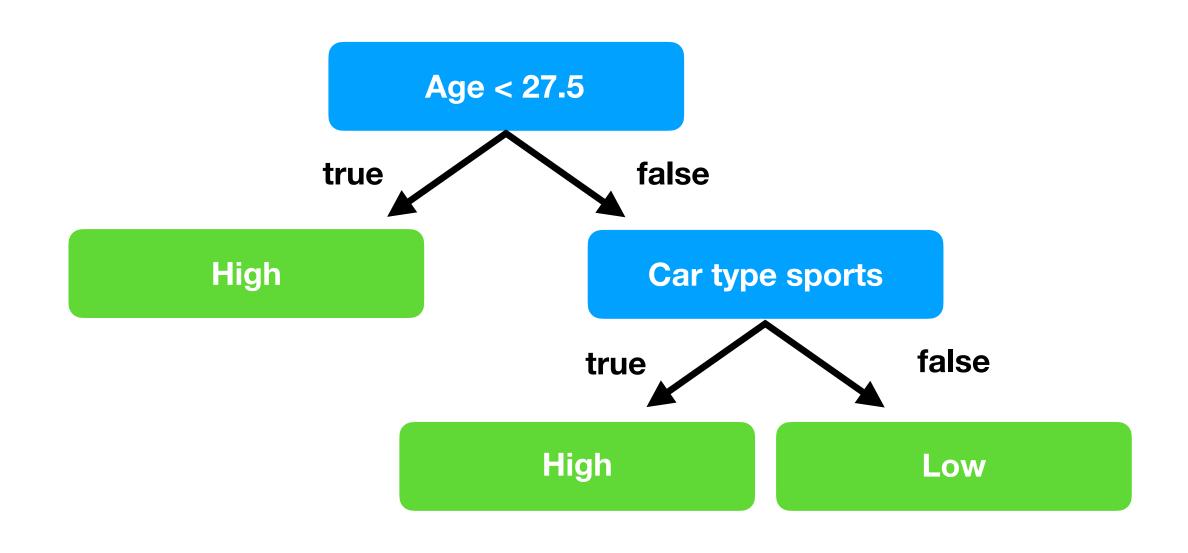
ation  $\{x', y'\}$   $x' \in \mathbb{R}^d, y' \in \{0, 1\}$ 

Which feature to use at each depth to do a split?

For the continuous feature, at what value to do a split?

For the categorical feature, which category to use for the split?

Continuous Feature	Categorical Feature	Class label
Age	Car type	Risk
23	family	high
17	sports	high
43	sports	high
68	family	low
32	family	low
20	family	high



#### **EPFL** Classification trees

Greedy approach: choose a feature and the split sequentially based on minimising a performance metric, for example, the **Gini impurity** of a node

Gini impurity of a leaf node: based on empirical probability of class

$$\sum_{l=1}^{K} P_{l} \sum_{j\neq l} P_{l}, = \sum_{l=1}^{K} P_{l} (1-P_{l})$$

$$e \times : K = 2 \quad (binary classification) : P_{l}P_{2} + P_{l}P_{2} = 2P_{l}P_{2}$$

Gini impurity of a node

Weighted sum of the gini impurity of the two leaf nodes associated with the node

Note: Criteria other than gini index (such as entropy) are also used for node split



#### Classification Trees

Example - constructing the tree

Gini impurity to construct the classification tree

- 1. compute Gini impurity for different features and at different split values
- 2. choose feature and split with lowest Gini impurity

Continuous Feature	Categorical Feature	Class label
Age	Car type	Risk
23	family	high
17	sports	high
43	sports	high
68	family	low
32	family	low
20	family	high

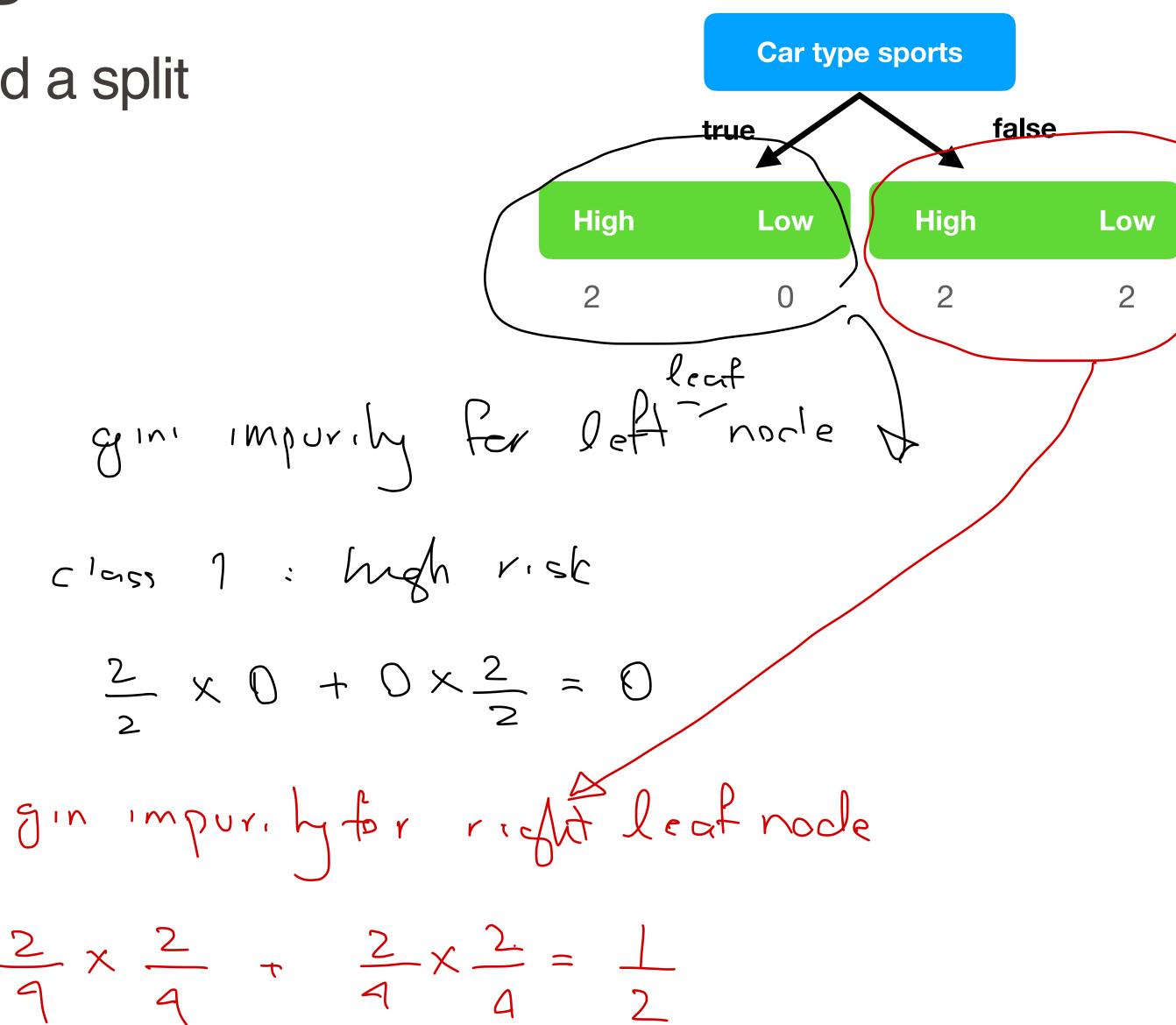
### **EPFL** Classification trees

Criteria for choosing a feature and a split

$$\frac{2}{6} \times 0 + \frac{4}{6} \times \frac{1}{2} = \frac{1}{3}$$

Continuous	Categorical	Class
<b>Feature</b>	<b>Feature</b>	label

Age	Car type	Risk
23	family	high
17	sports	high
43	sports	high
68	family	low
32	family	low
20	family	high





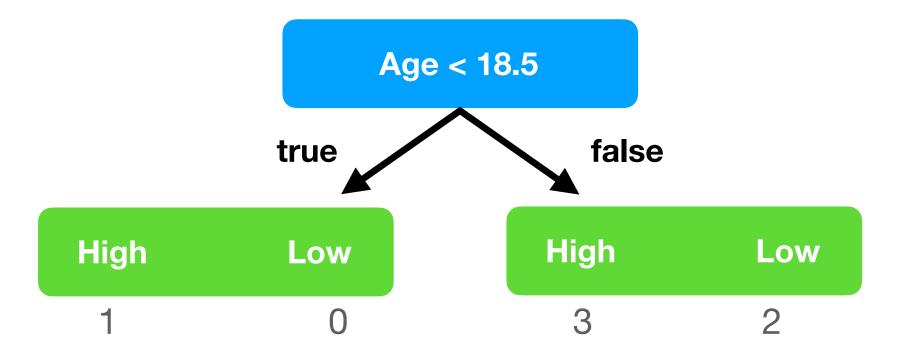
#### **Decision Trees**

#### Splitting continuous features

tid	Age	Risk
0	23	high
1	17	high
2	43	high
3	68	low
4	32	low
5	20	high

tid	Age	Risk	
1	17	high	3 181
5	20	high	7 /2
0	23	high	21/2
4	32	low	27/
2	43	high	37/2
3	68	low	$\begin{bmatrix} 1 & 1/2 \\ 1 & 1/2 \end{bmatrix}$
			7 22 /

$$\frac{1}{6}(0) + \frac{S}{6} \times \frac{12}{2S} = \frac{2}{S}$$



$$5\left(\frac{1}{1}\times 0\right) = 0$$

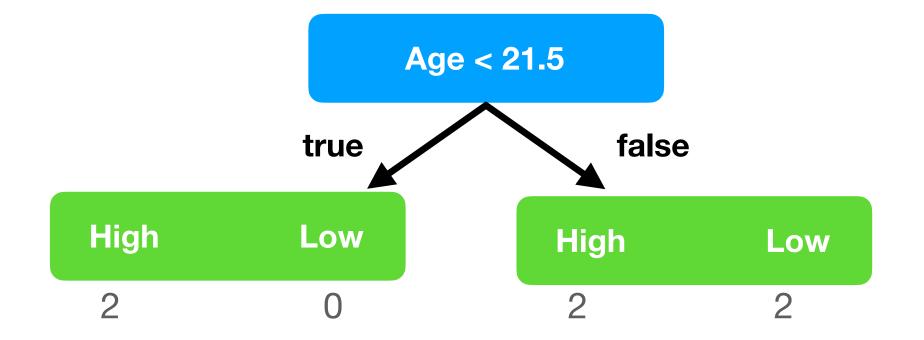
$$\frac{3}{5} \times \frac{2}{5} + \frac{2}{5} \times \frac{3}{5} = \frac{12}{25}$$



### Decision Trees Splitting continuous features

tid	Age	Risk
0	23	high
1	17	high
2	43	high
3	68	low
4	32	low
5	20	high

tid	Age	Risk
1	17	high
5	20	high
0	23	high
4	32	low
2	43	high
3	68	low



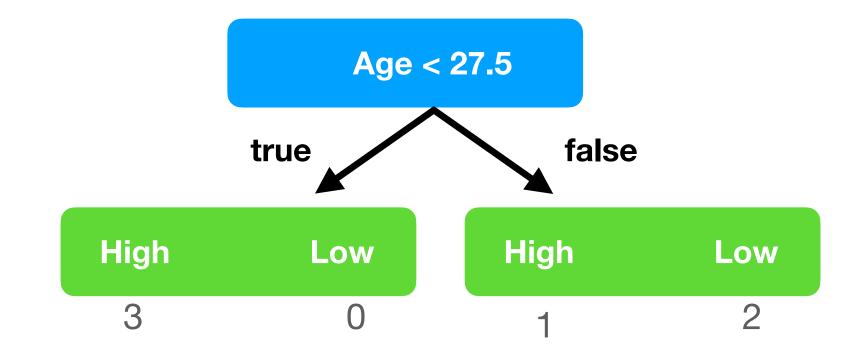


#### **Decision Trees**

#### Splitting continuous features

tid	Age	Risk
0	23	high
1	17	high
2	43	high
3	68	low
4	32	low
5	20	high

tid	Age	Risk
1	17	high
5	20	high
0	23	high
4	32	low
2	43	high
3	68	low





#### **Decision Trees**

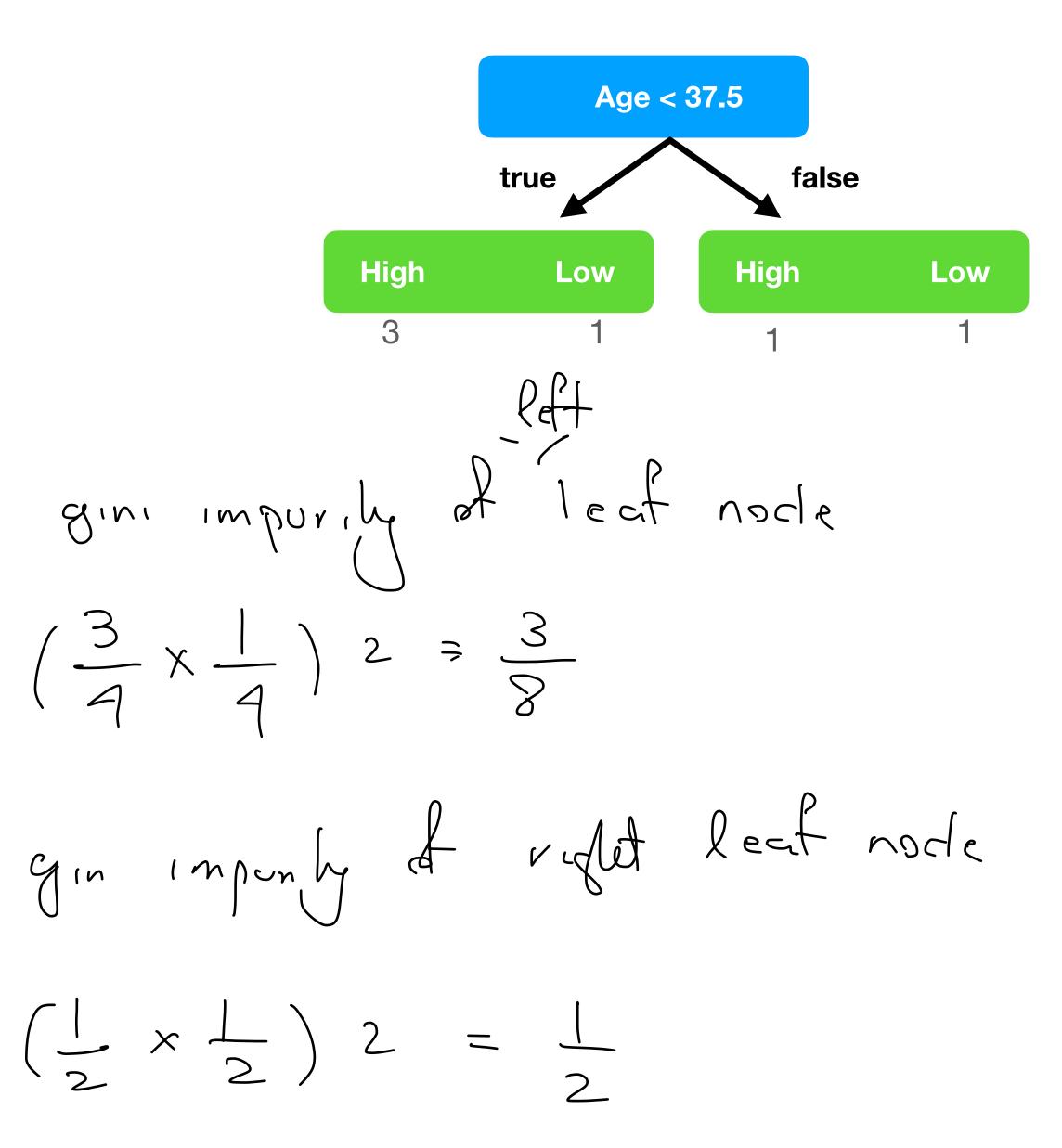
#### Splitting continuous features

tid	Age	Risk
0	23	high
1	17	high
2	43	high
3	68	low
4	32	low
5	20	high

Reorder the data depending on Age

tid	Age	Risk
1	17	high
5	20	high
0	23	high
4	32	low
2	43	high
3	68	low

5/2



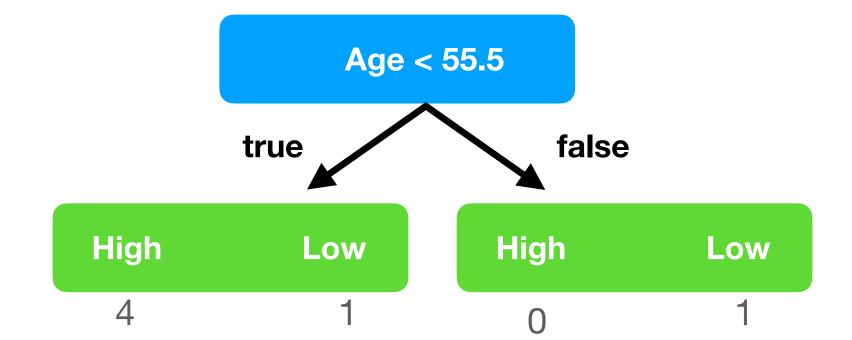


### Decision Trees Splitting continuous features

tid	Age	Risk
0	23	high
1	17	high
2	43	high
3	68	low
4	32	low
5	20	high

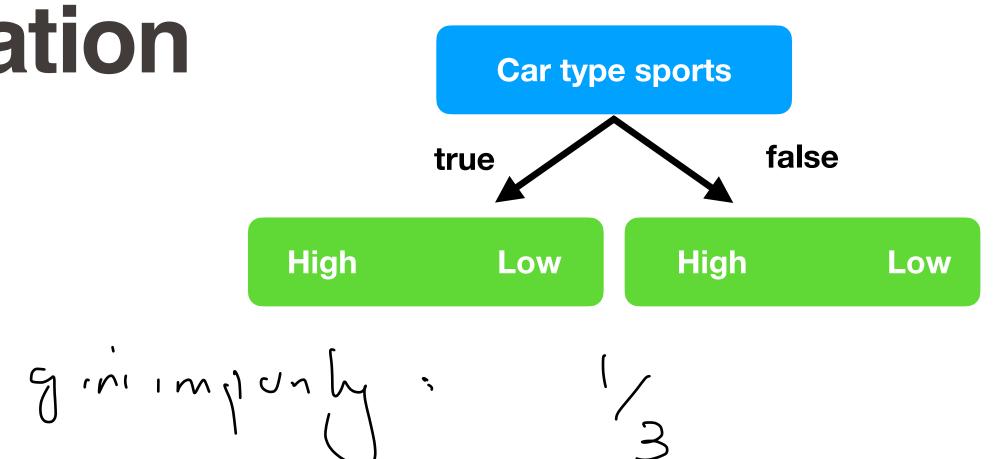
tid	Age	Risk
1	17	high
5	20	high
0	23	high
4	32	low
2	43	high
3	68	low





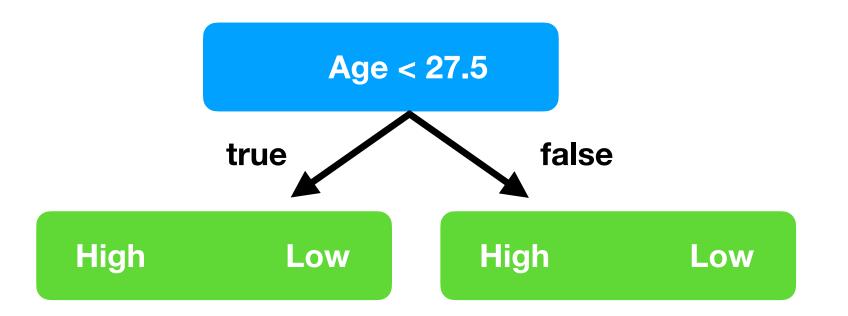


#### Decision Trees for classification



Continuous	Categorical	Class
<b>Feature</b>	Feature	label

Age	Car type	Risk
23	family	high
17	sports	high
43	sports	high
68	family	low
32	family	low
20	family	high



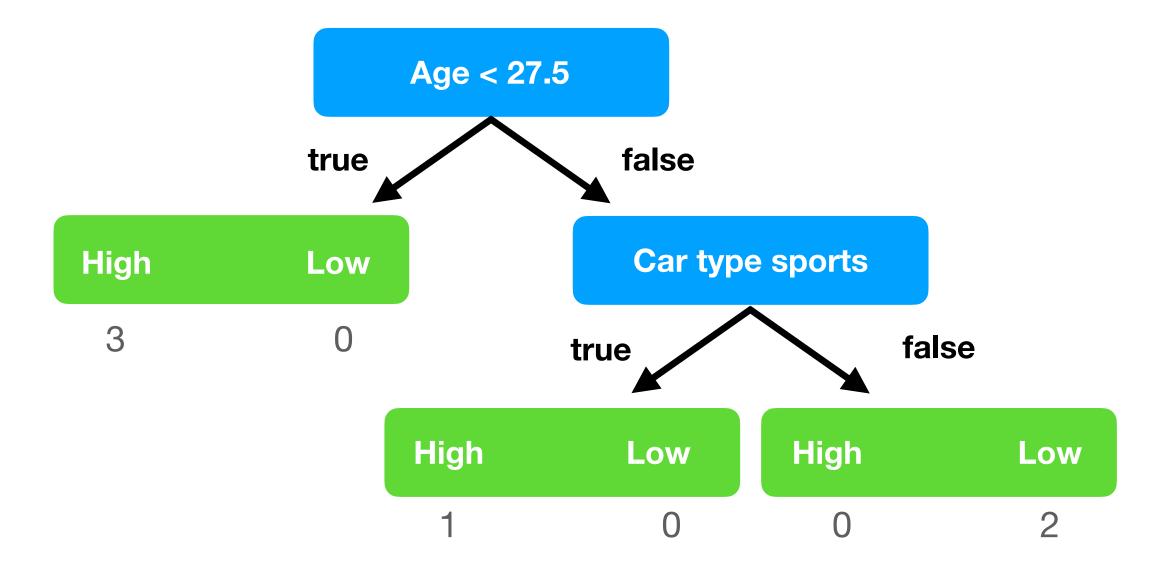
gini impunh:

#### **EPFL** Classification tree

### Example - Final tree

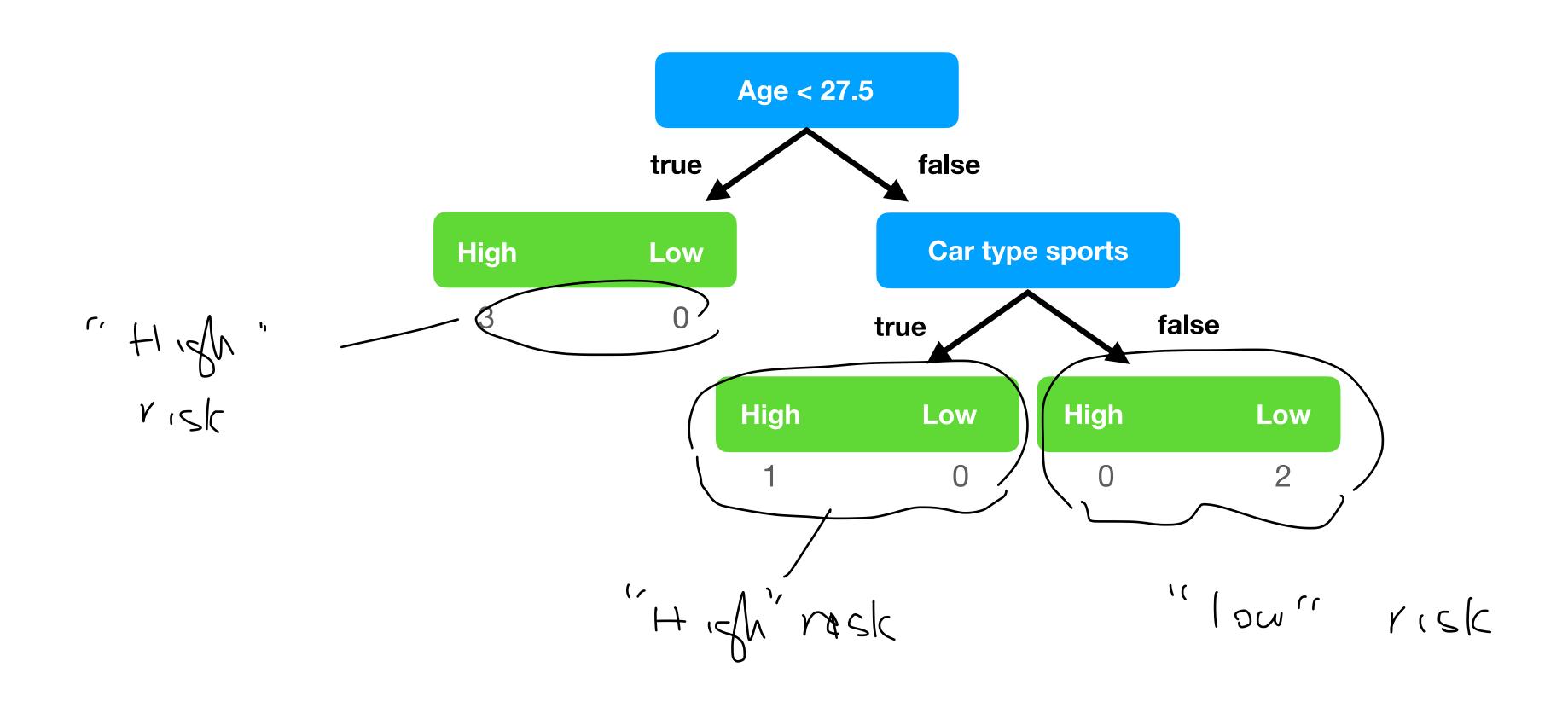
Continuous	Categorical	Class
Feature	Feature	label

Age	Car type	Risk
23	family	high
17	sports	high
43	sports	high
68	family	low
32	family	low
20	family	high



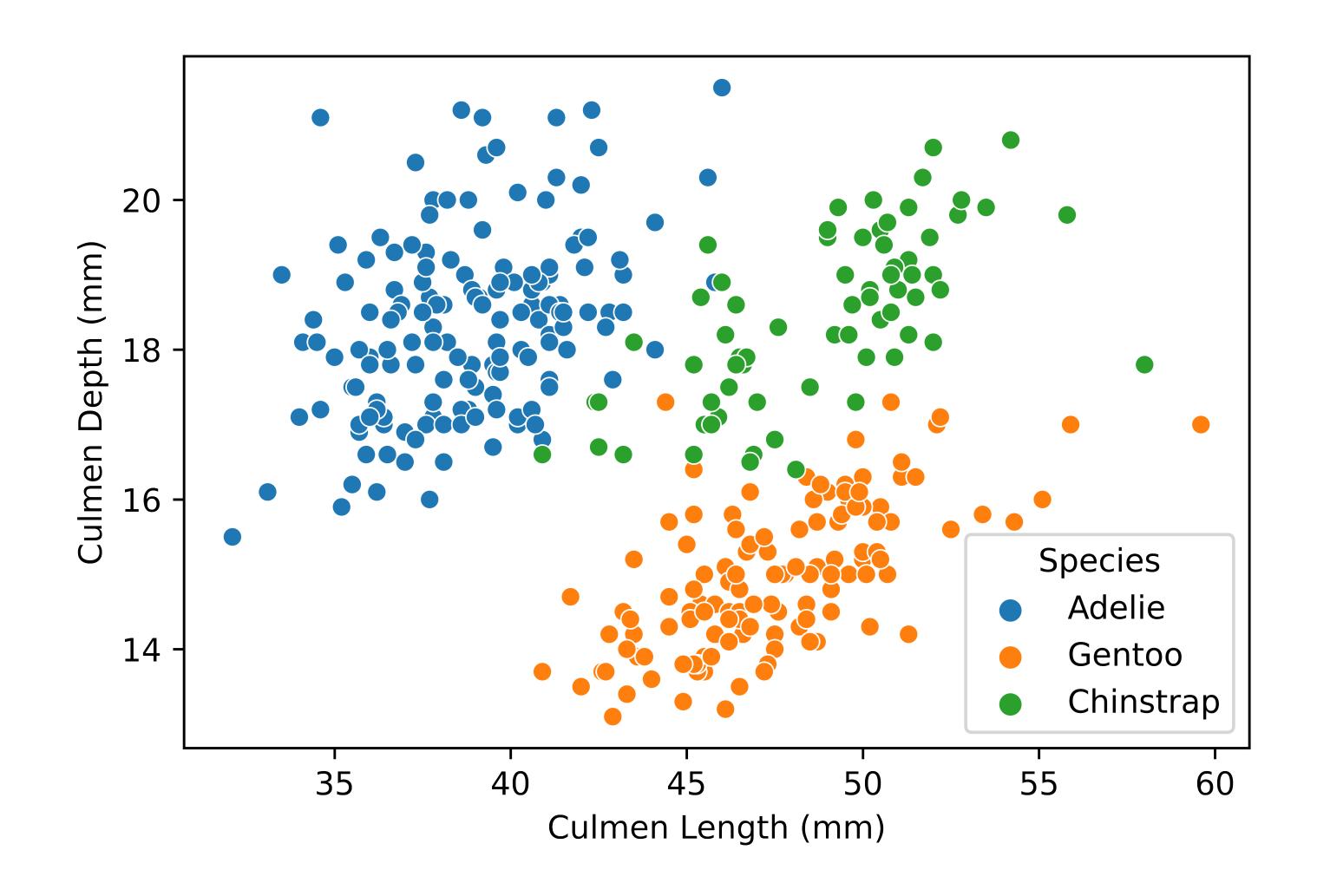
#### **EPFL** Classification tree

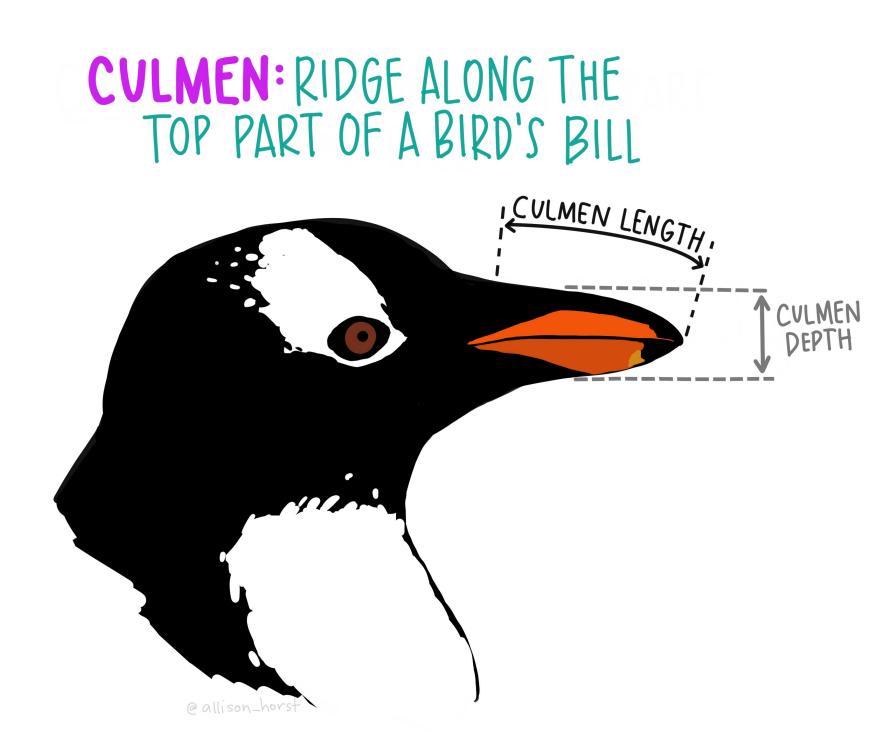
How do we use this for prediction?





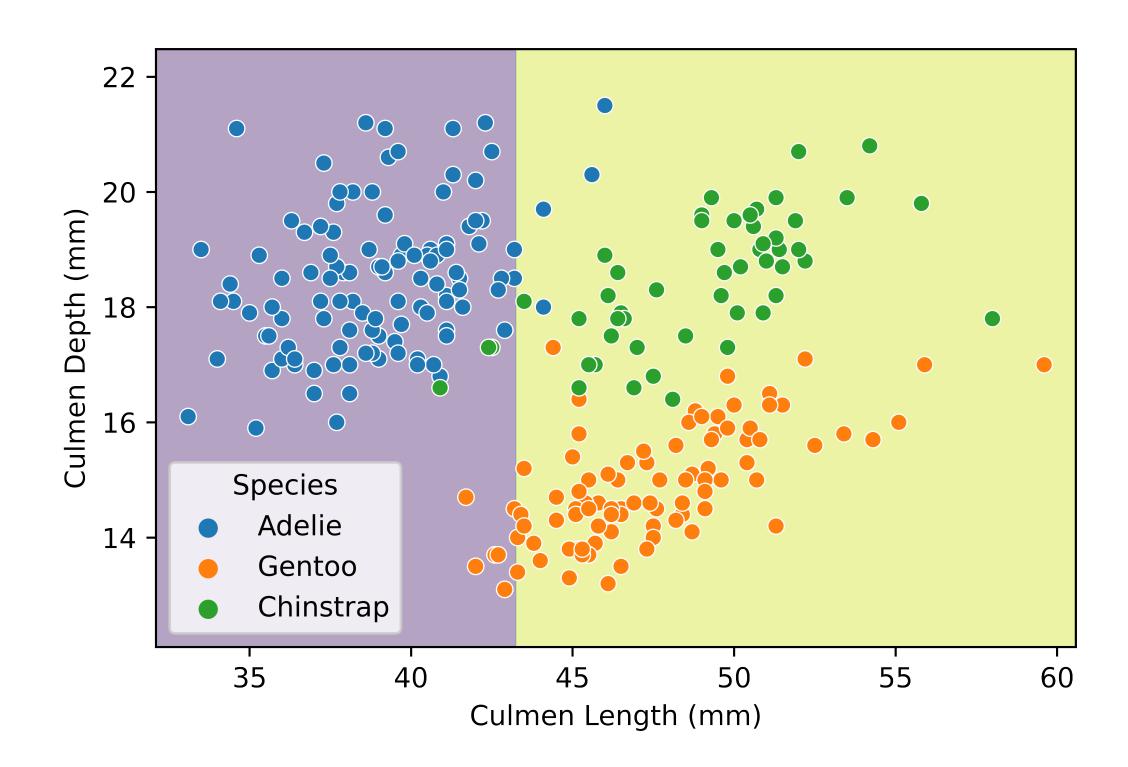
# Classification tree example penguin data set

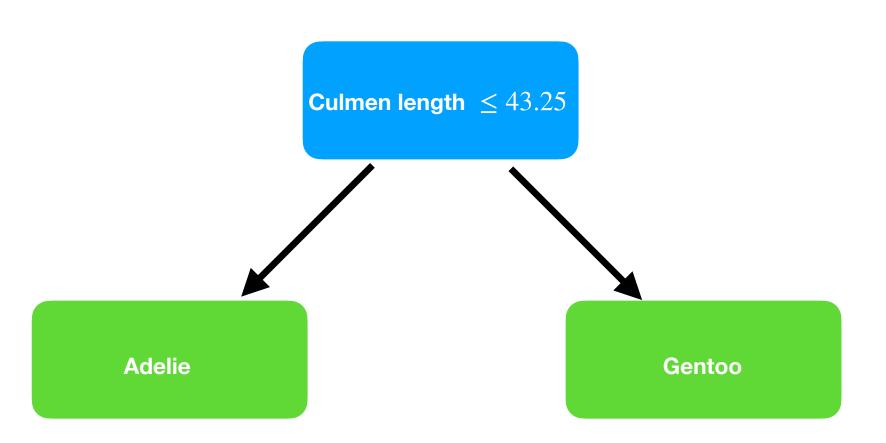






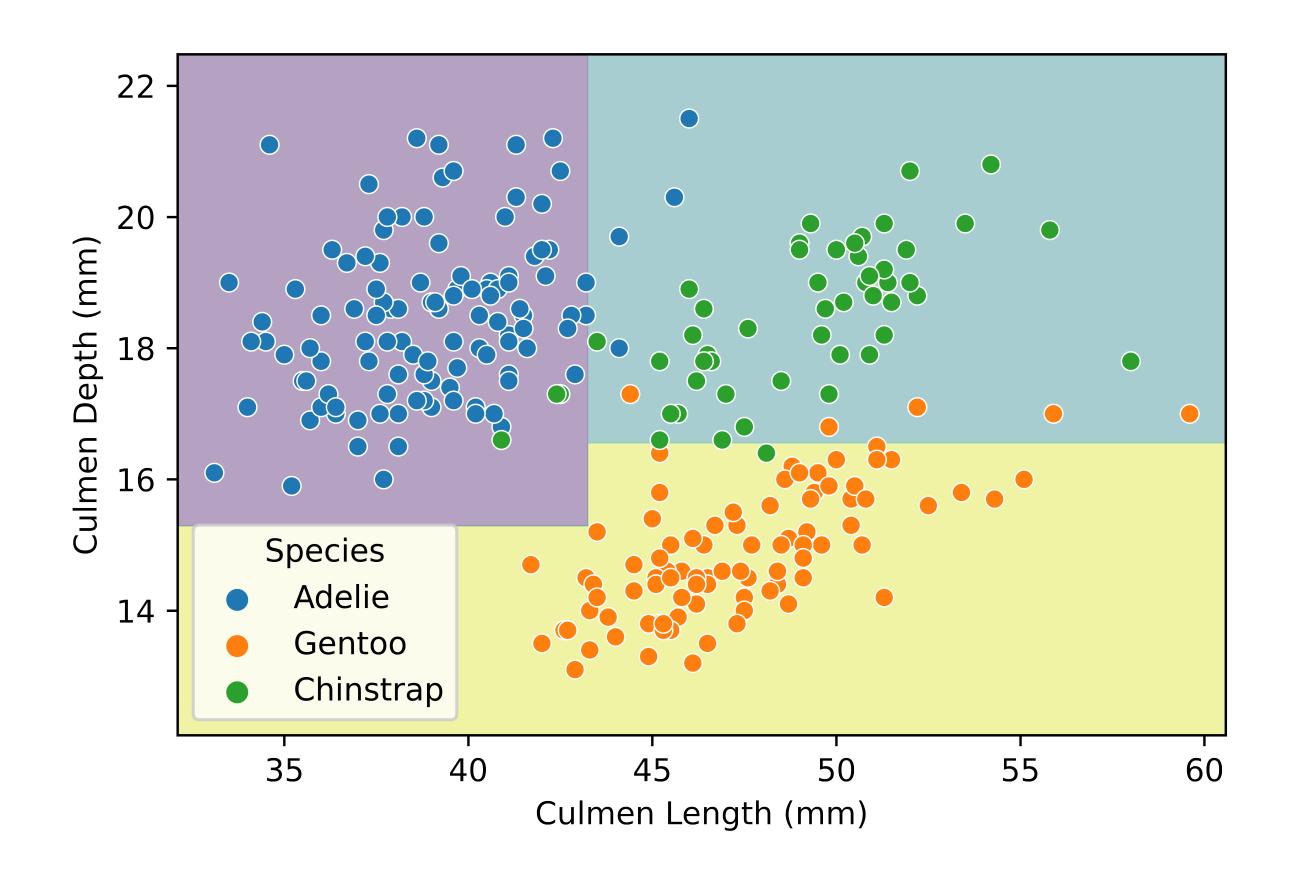
# Decision Trees Penguins example

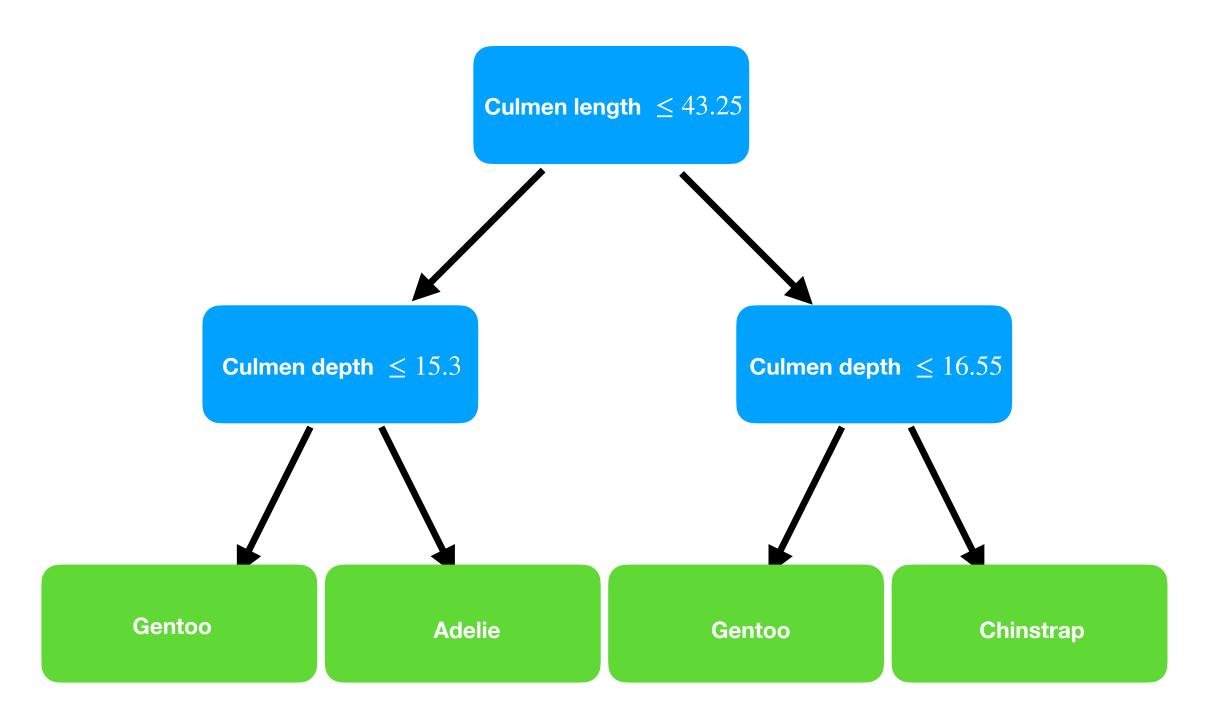






# Decision Trees Penguins example

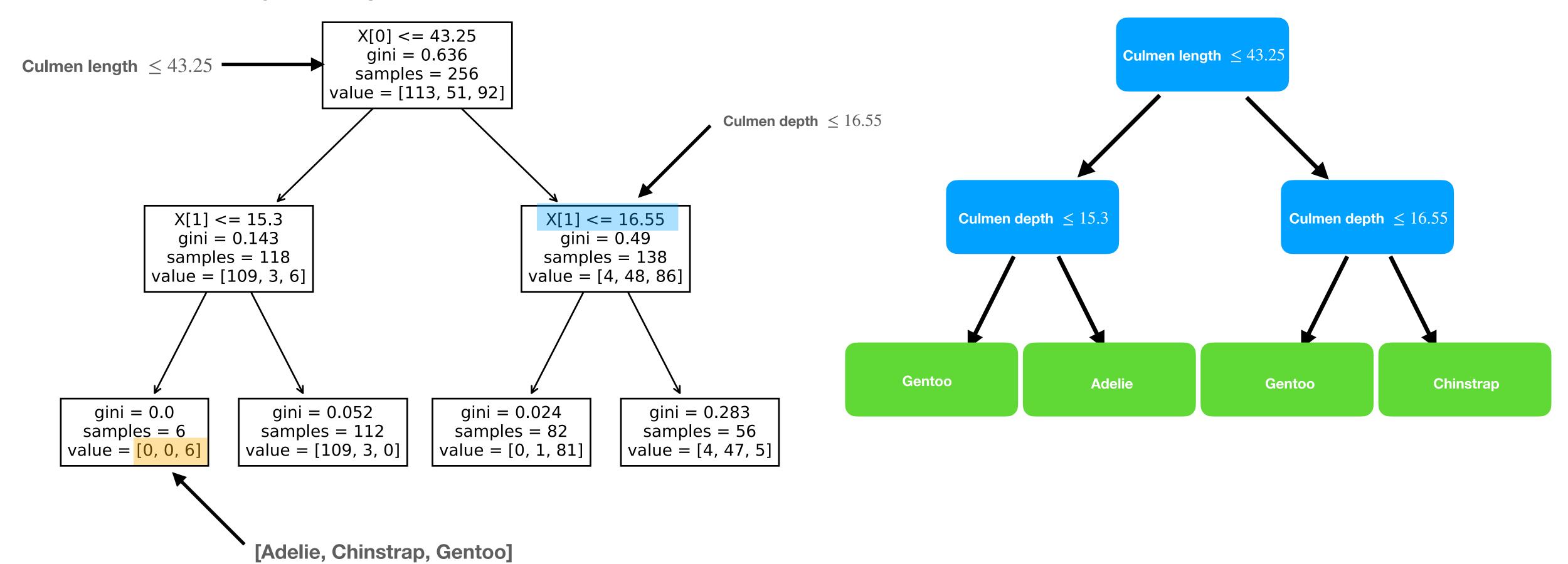






# Decision Trees Penguins example

Python output for *sklearn.tree* 





### Summary - decision tree construction

- We select the most discriminative Feature and Threshold
  - Discriminative power based on a criteria:
    - Regression tree: Average Squared Loss
    - Classification tree: Gini impurity
    - We create a node based on this feature
- We repeat for each new branch until a stopping criteria
- To ensure the tree-depth is not too high (avoid overfitting) "pruning" is done...

More examples on StatQuestion!!!: https://www.youtube.com/watch?v=\_L39rN6gz7Y