

## Lecture 12

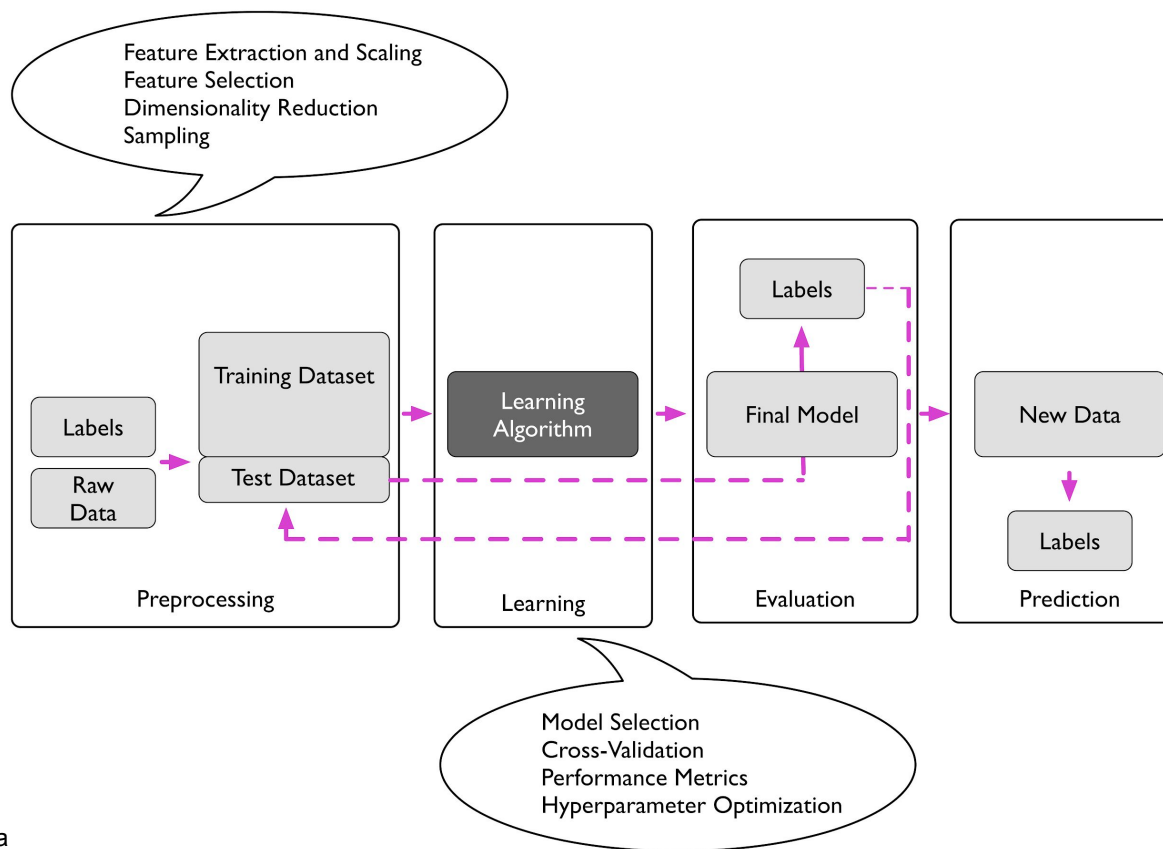
# Decision trees

<https://github.com/dalcimar/MA28CP-Intro-to-Machine-Learning>

UTFPR - Federal University of Technology - Paraná

<https://www.dalcimar.com/>

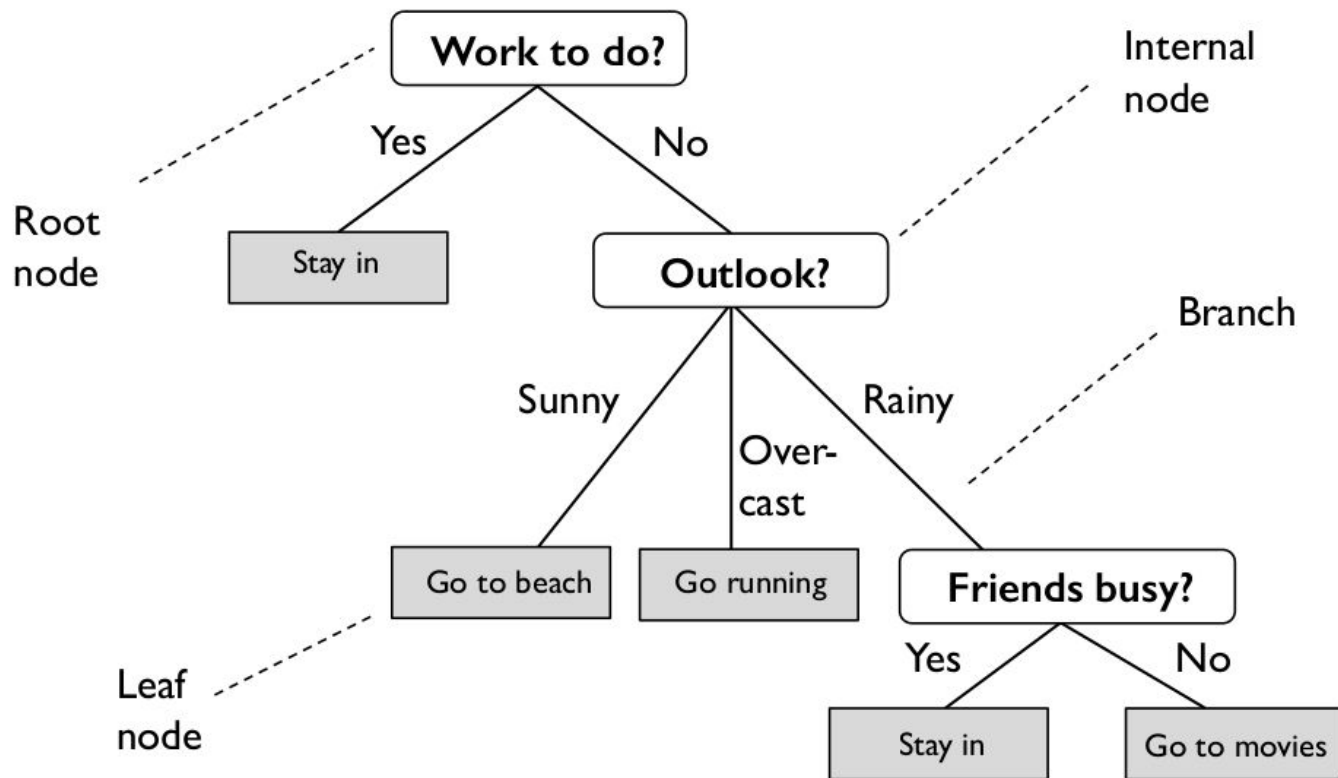
# Machine learning pipeline



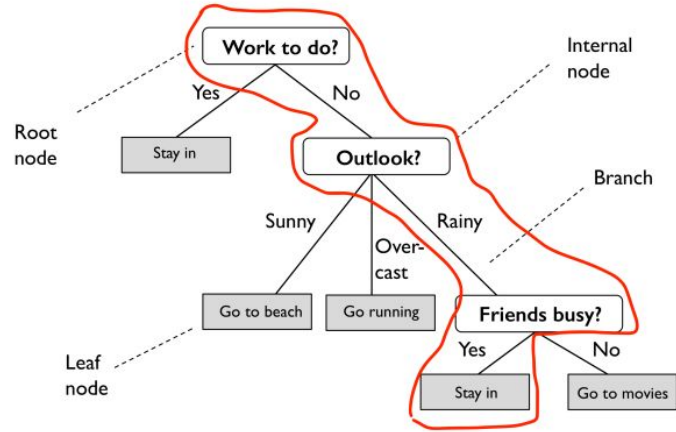
# Topics

1. **Intro to decision trees**
2. Recursive and divide & conquer strategy
3. Types of decision trees
4. Splitting criteria
5. Gini & Entropy vs misclassification error
6. Improvements & dealing with overfitting
7. Code example

# Decision tree terminology



# Decision tree as rulesets



**IF**

\_\_\_\_\_

\_\_\_\_\_

**THEN**

\_\_\_\_\_

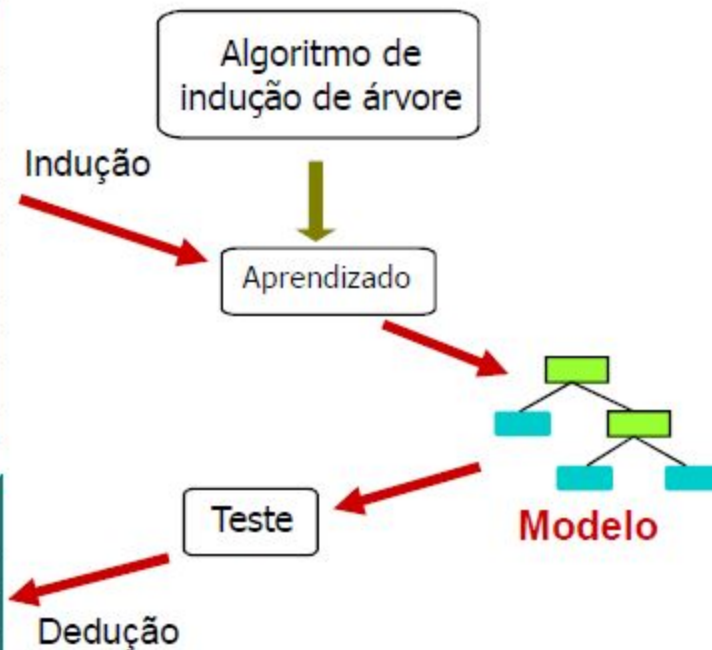
# Decision tree pipeline

Conjunto de  
treinamento

<i>Id</i>	<i>E Credor</i>	<i>Estado Civil</i>	<i>Salário</i>	<i>Calote</i>
1	Sim	Solteiro	125K	Não
2	Não	Casado	100K	Não
3	Não	Solteiro	70K	Não
4	Sim	Casado	120K	Não
5	Não	Divorciado	95K	Sim
6	Não	Casado	60K	Não
7	Sim	Divorciado	220K	Não
8	Não	Solteiro	85K	Sim
9	Não	Casado	75K	Não
10	Não	Solteiro	90K	Sim

Conjunto de  
teste

<i>Id</i>	<i>E Credor</i>	<i>Estado Civil</i>	<i>Salário</i>	<i>Calote</i>
11	Não	Casado	80K	?
12	Não	Solteiro	100K	?
13	Sim	Solteiro	100K	?
14	Não	Casado	120K	?
15	Sim	Solteiro	80K	?

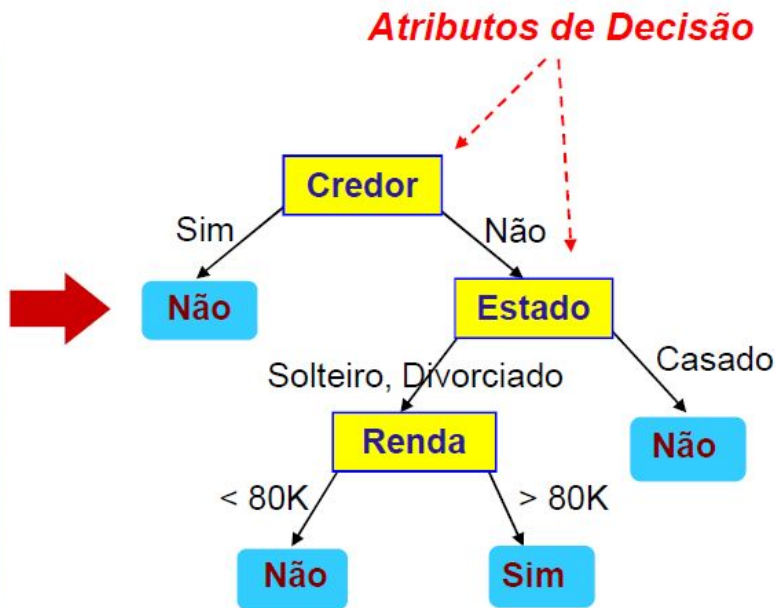


# Training

- Works on categorical (binary, nominal or ordinal) and
- Works also with numeric (continuos) attributes

<i><b>Id</b></i>	<i><b>E Credor</b></i>	<i><b>Estado Civil</b></i>	<i><b>Salário</b></i>	<i><b>Calote</b></i>
1	Sim	Solteiro	125K	Não
2	Não	Casado	100K	Não
3	Não	Solteiro	70K	Não
4	Sim	Casado	120K	Não
5	Não	Divorciado	95K	Sim
6	Não	Casado	60K	Não
7	Sim	Divorciado	220K	Não
8	Não	Solteiro	85K	Sim
9	Não	Casado	75K	Não
10	Não	Solteiro	90K	Sim

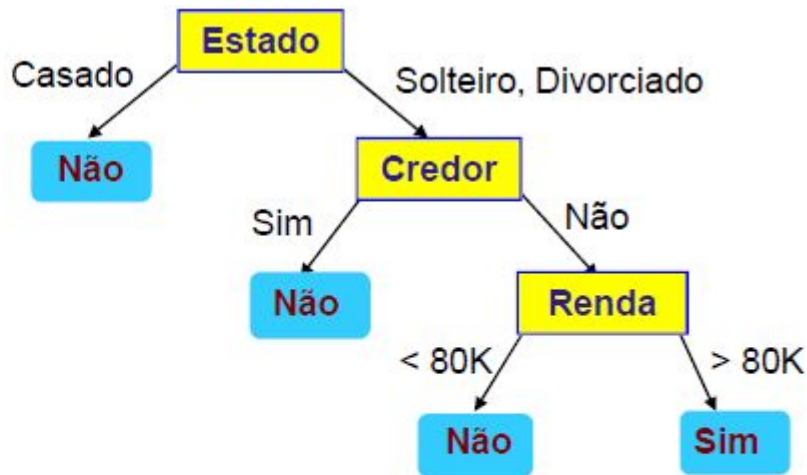
**Dados de Treinamento**



**Modelo: Árvore de Decisão**

# Training

<i>Id</i>	<i>E Credor</i>	<i>Estado Civil</i>	<i>Salário</i>	<i>Calote</i>
1	Sim	Solteiro	125K	Não
2	Não	Casado	100K	Não
3	Não	Solteiro	70K	Não
4	Sim	Casado	120K	Não
5	Não	Divorciado	95K	Sim
6	Não	Casado	60K	Não
7	Sim	Divorciado	220K	Não
8	Não	Solteiro	85K	Sim
9	Não	Casado	75K	Não
10	Não	Solteiro	90K	Sim



Diferentes árvores podem ser ajustadas  
para os mesmos dados !



# Decision tree on Iris dataset (2 features)

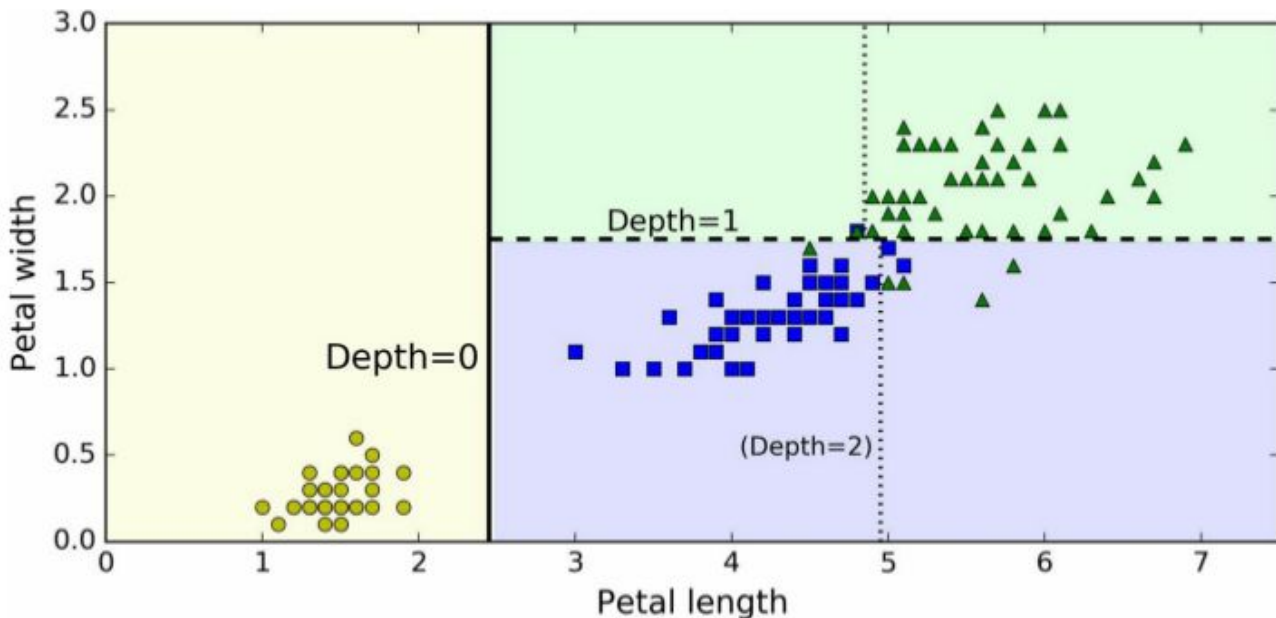
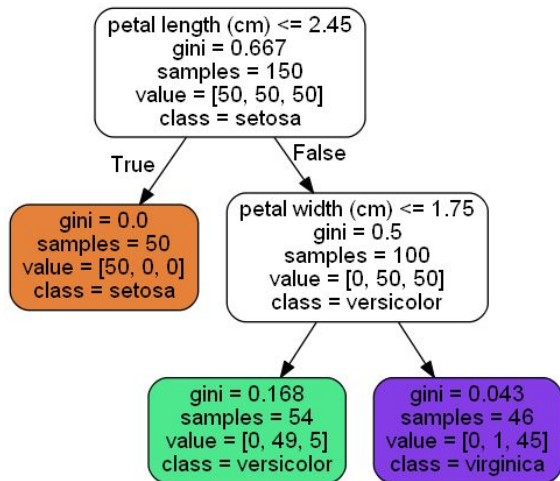
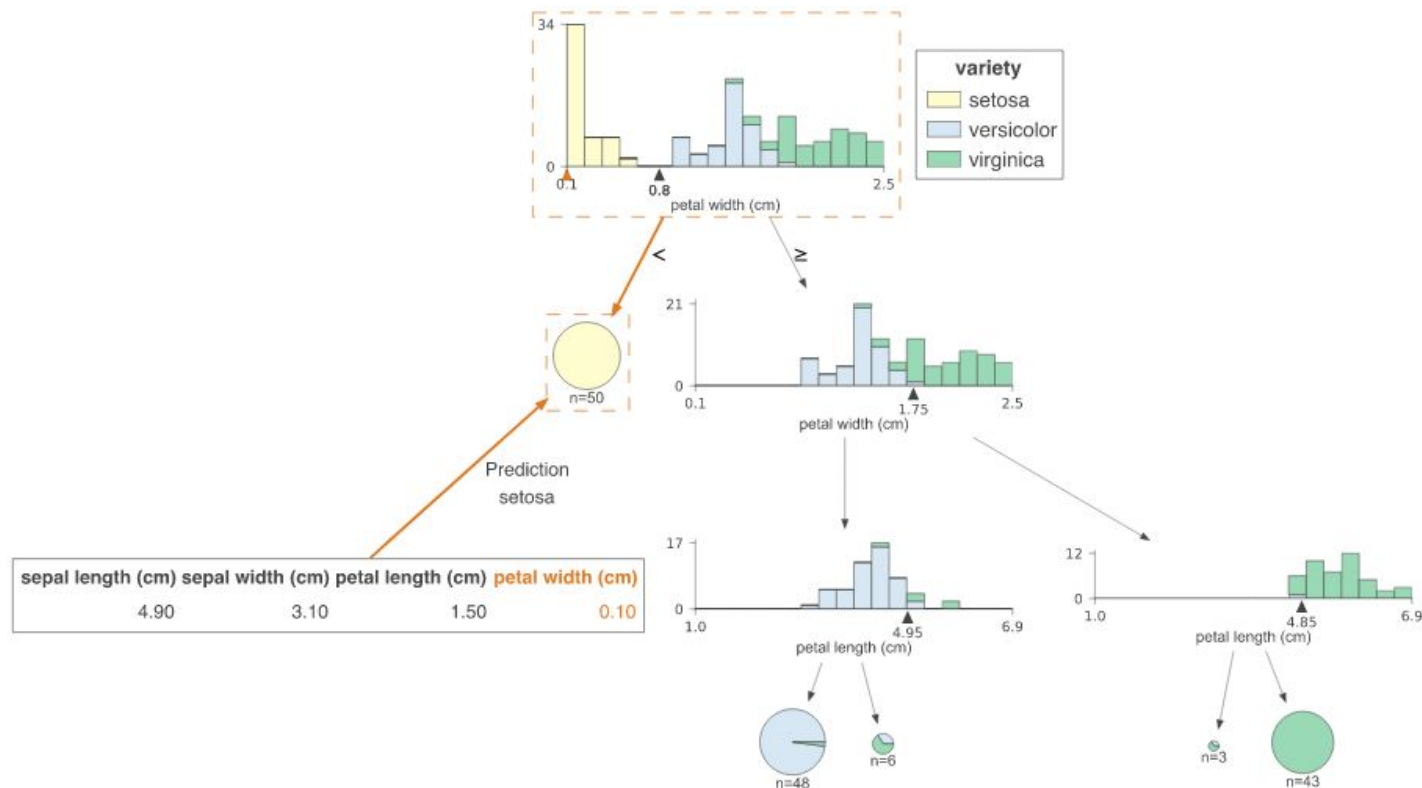
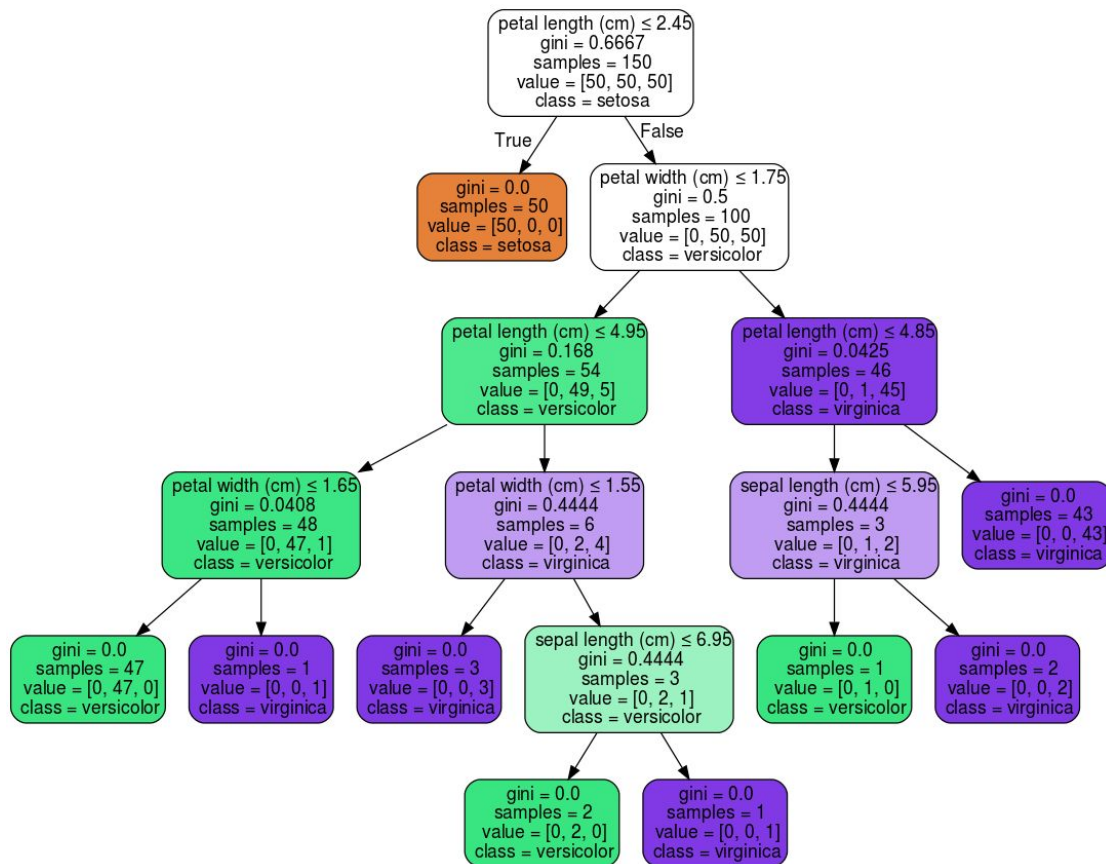


Figure 6-2. Decision Tree decision boundaries

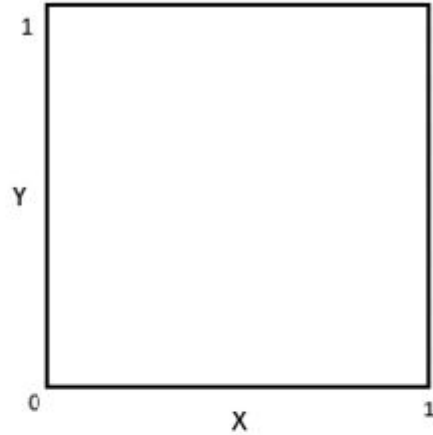
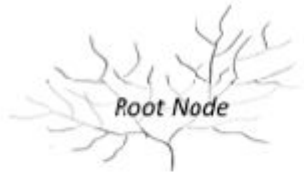
# Decision tree on Iris dataset (2 features)



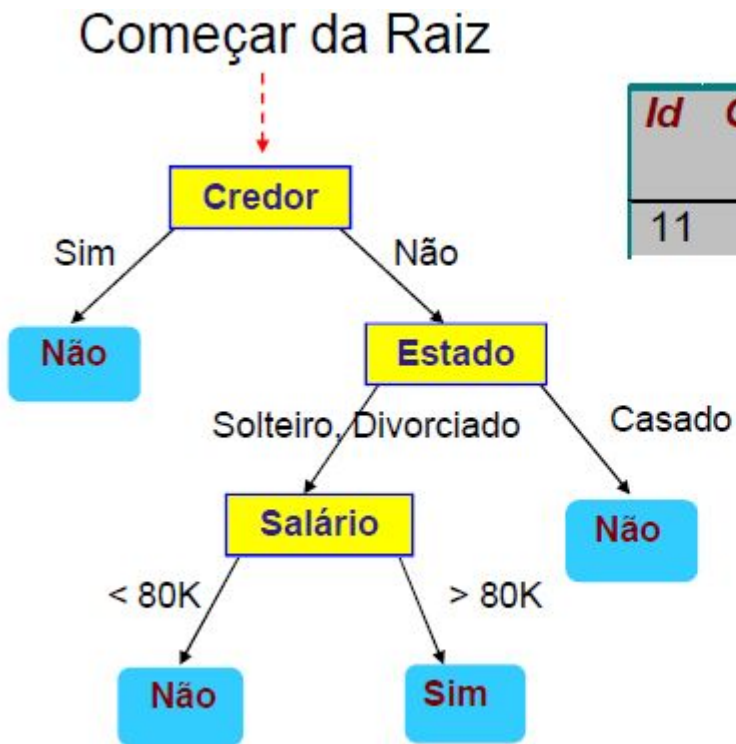
# Decision tree on Iris dataset (4 features)



# Growing depth process



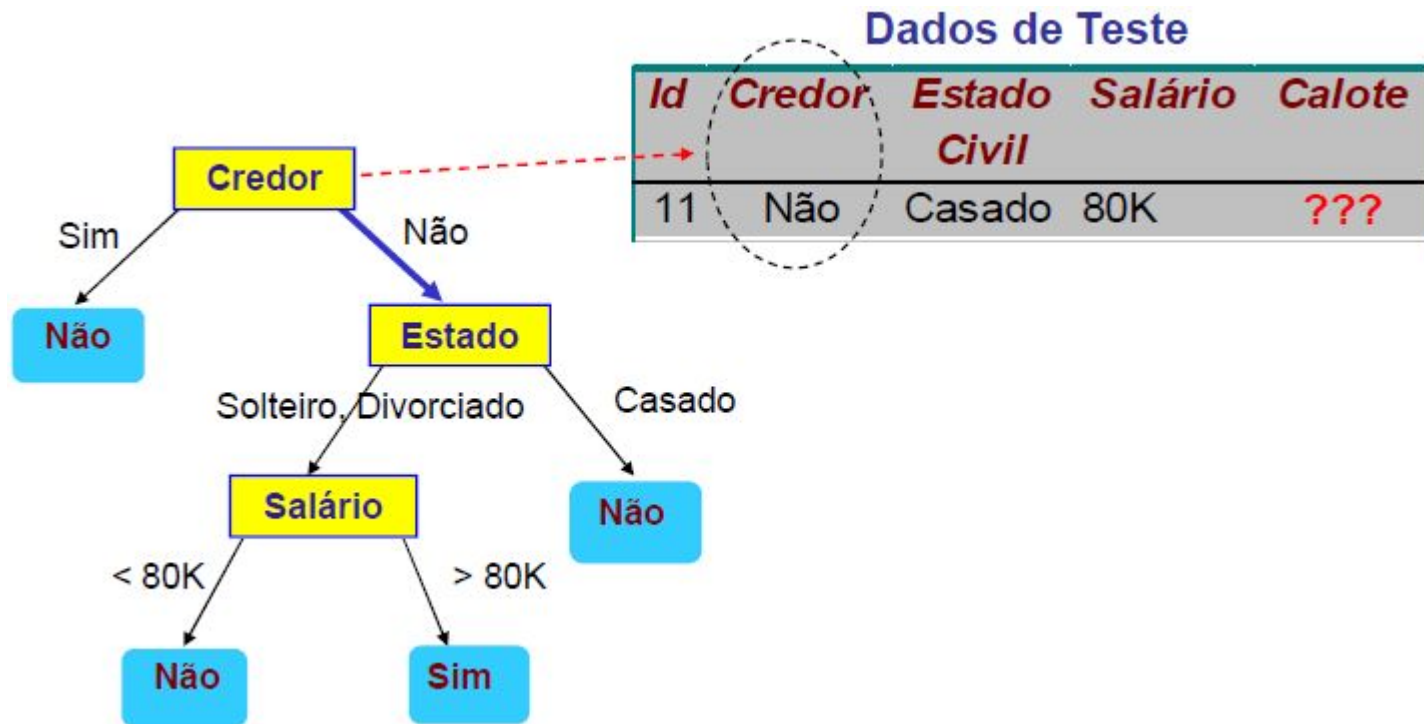
# Inference



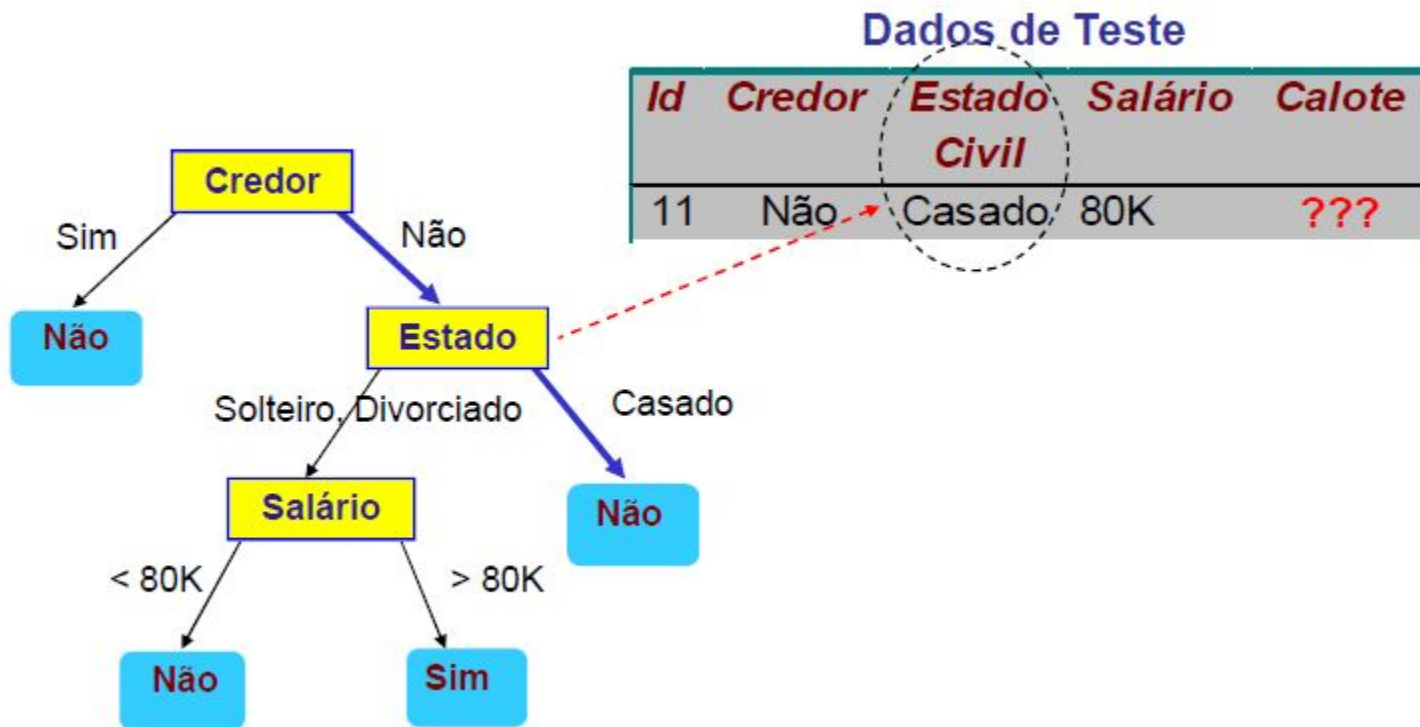
Dados de Teste

<i>Id</i>	<i>Credor</i>	<i>Estado Civil</i>	<i>Salário</i>	<i>Calote</i>
11	Não	Casado	80K	???

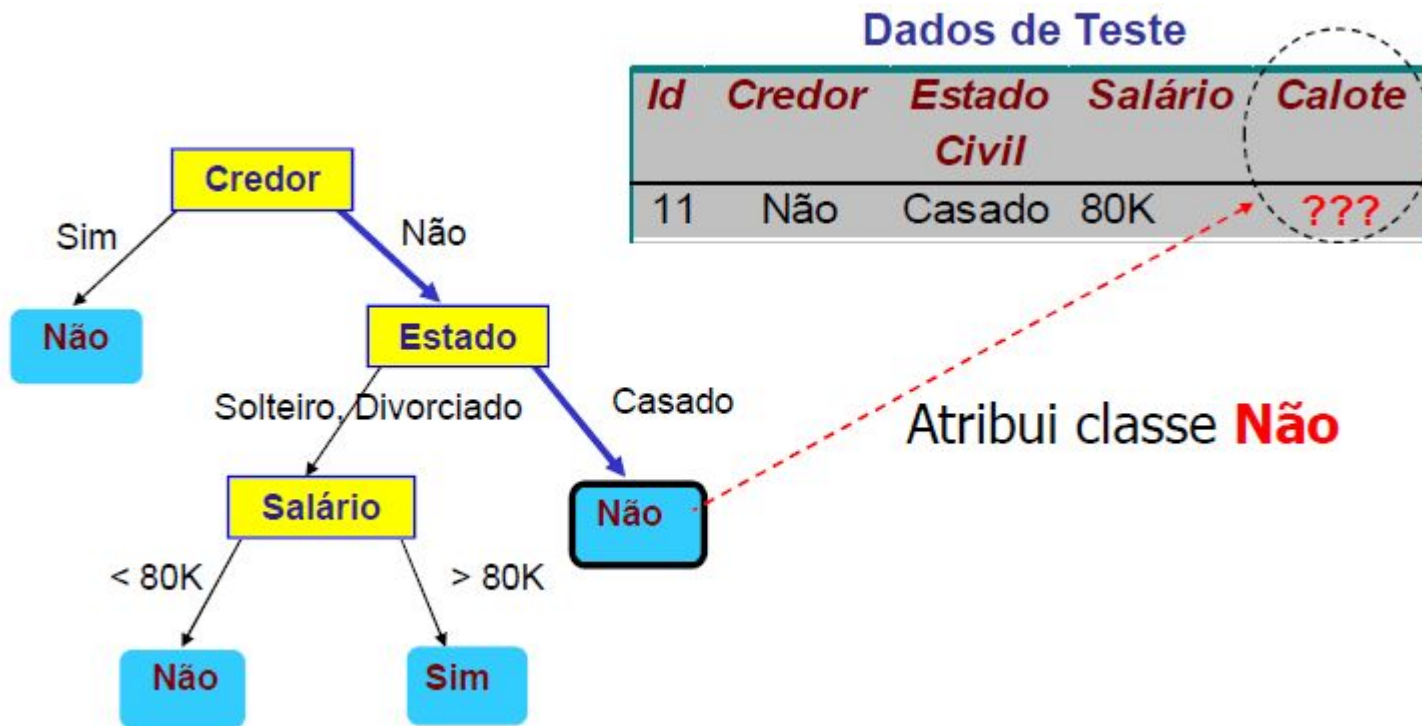
# Inference



# Inference



# Inference





# Tree depth overfitting

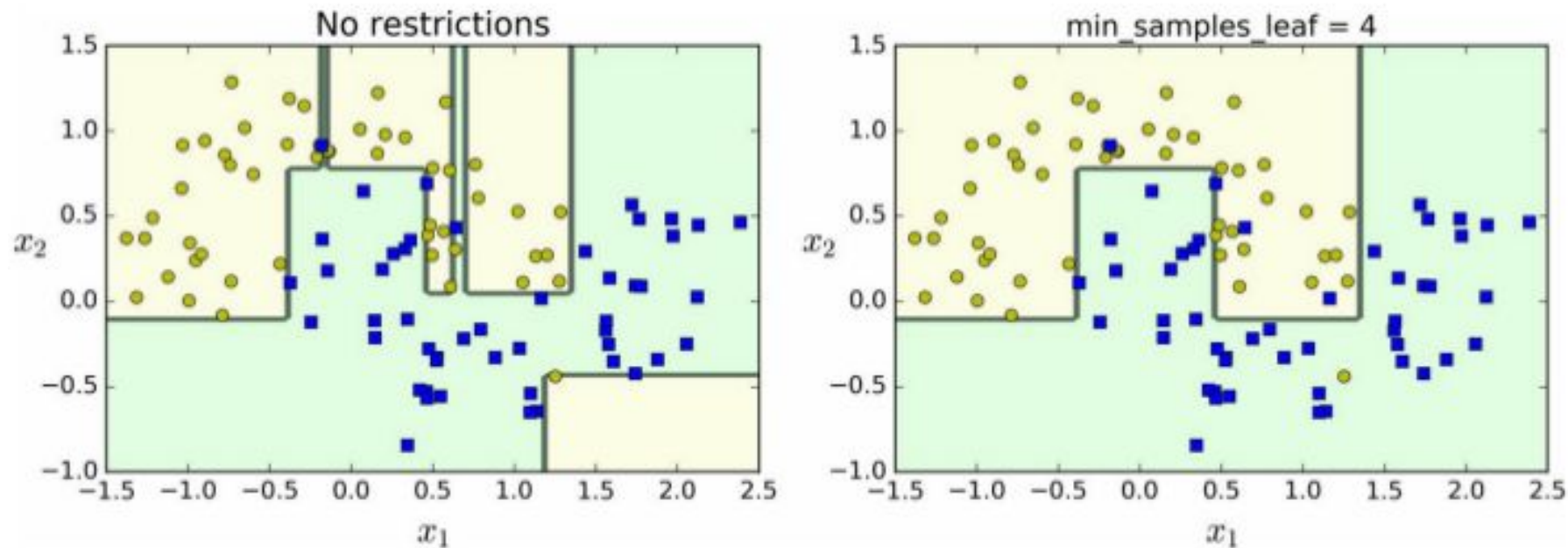


Figure 6-3. Regularization using `min_samples_leaf`

# Animations

Amazing animation of decision tree

<http://www.r2d3.us/visual-intro-to-machine-learning-part-1/>

# Topics

1. Intro to decision trees
2. **Recursive and divide & conquer strategy**
3. Types of decision trees
4. Splitting criteria
5. Gini & Entropy vs misclassification error
6. Improvements & dealing with overfitting
7. Code example

# Recursion / Recursive Algorithms

What does this function do?

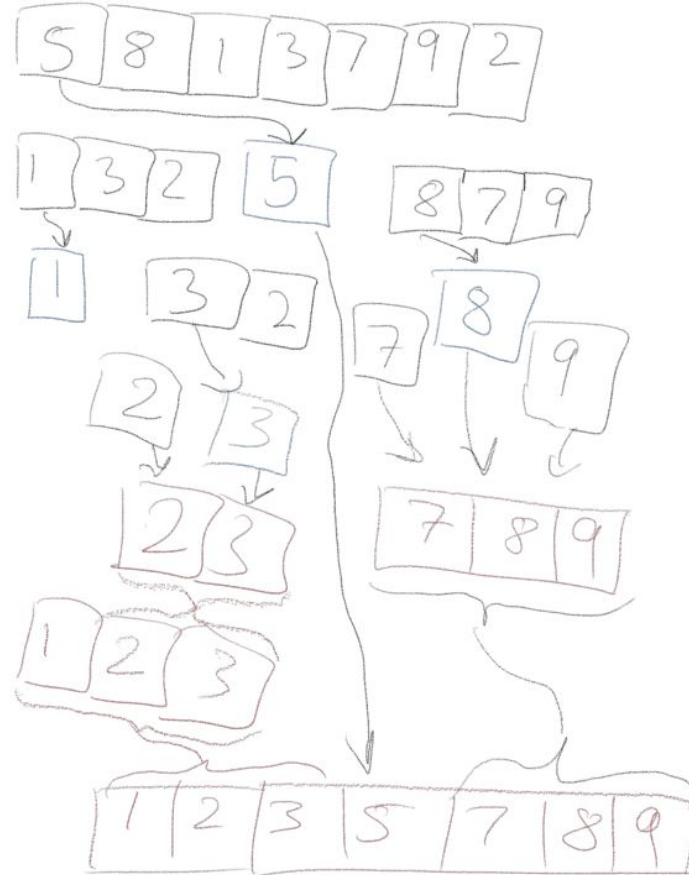
```
1 def some_fun(x):  
2     if x == []:  
3         return 0  
4     else:  
5         return 1 + some_fun(x[1:])
```

# Divide & Conquer Algorithms: Quicksort

```
1 def quicksort(array):
2     if len(array) < 2:
3         return array
4     else:
5         pivot = array[0]
6         smaller, bigger = [], []
7         for ele in array[1:]:
8             if ele <= pivot:
9                 smaller.append(ele)
10            else:
11                bigger.append(ele)
12        return quicksort(smaller) + [pivot] + quicksort(bigger)
```

# Divide & Conquer Algorithms: Quicksort

```
1 def quicksort(array):
2     if len(array) < 2:
3         return array
4     else:
5         pivot = array[0]
6         smaller, bigger = [], []
7         for ele in array[1:]:
8             if ele <= pivot:
9                 smaller.append(ele)
10            else:
11                bigger.append(ele)
12     return quicksort(smaller) + [pivot] + quicksort(bigger)
```



# Time complexity of quicksort

$O(n \log n)$

```
1 def quicksort(array):
2     if len(array) < 2:
3         return array
4     else:
5         pivot = array[0]
6         smaller, bigger = [], []
7         for ele in array[1:]:
8             if ele <= pivot:
9                 smaller.append(ele)
10            else:
11                bigger.append(ele)
12        return quicksort(smaller) + [pivot] + quicksort(bigger)
```

# Time and space complexity of sorting algorithms

Algorithm	Time Complexity			Space Complexity
	Best	Average	Worst	Worst
<u>Quicksort</u>	$\Omega(n \log(n))$	$\Theta(n \log(n))$	$O(n^2)$	$O(\log(n))$
<u>Mergesort</u>	$\Omega(n \log(n))$	$\Theta(n \log(n))$	$O(n \log(n))$	$O(n)$
<u>Timsort</u>	$\Omega(n)$	$\Theta(n \log(n))$	$O(n \log(n))$	$O(n)$
<u>Heapsort</u>	$\Omega(n \log(n))$	$\Theta(n \log(n))$	$O(n \log(n))$	$O(1)$
<u>Bubble Sort</u>	$\Omega(n)$	$\Theta(n^2)$	$O(n^2)$	$O(1)$
<u>Insertion Sort</u>	$\Omega(n)$	$\Theta(n^2)$	$O(n^2)$	$O(1)$
<u>Selection Sort</u>	$\Omega(n^2)$	$\Theta(n^2)$	$O(n^2)$	$O(1)$
<u>Tree Sort</u>	$\Omega(n \log(n))$	$\Theta(n \log(n))$	$O(n^2)$	$O(n)$
<u>Shell Sort</u>	$\Omega(n \log(n))$	$\Theta(n(\log(n))^2)$	$O(n(\log(n))^2)$	$O(1)$
<u>Bucket Sort</u>	$\Omega(n+k)$	$\Theta(n+k)$	$O(n^2)$	$O(n)$
<u>Radix Sort</u>	$\Omega(nk)$	$\Theta(nk)$	$O(nk)$	$O(n+k)$
<u>Counting Sort</u>	$\Omega(n+k)$	$\Theta(n+k)$	$O(n+k)$	$O(k)$
<u>Cubesort</u>	$\Omega(n)$	$\Theta(n \log(n))$	$O(n \log(n))$	$O(n)$



# Decision Tree in Pseudocode

GenerateTree( $\mathcal{D}$ ):

- if  $y = 1 \forall \langle \mathbf{x}, \mathbf{y} \rangle \in \mathcal{D}$  or  $y = 0 \forall \langle \mathbf{x}, \mathbf{y} \rangle \in \mathcal{D}$ :
  - return Tree
- else:
  - Pick best feature  $x_j$ :
    - $\mathcal{D}_0$  at Child<sub>0</sub> :  $x_j = 0 \forall \langle \mathbf{x}, \mathbf{y} \rangle \in \mathcal{D}$
    - $\mathcal{D}_1$  at Child<sub>1</sub> :  $x_j = 1 \forall \langle \mathbf{x}, \mathbf{y} \rangle \in \mathcal{D}$

return Node( $x_j$ , GenerateTree( $\mathcal{D}_0$ ), GenerateTree( $\mathcal{D}_1$ ))

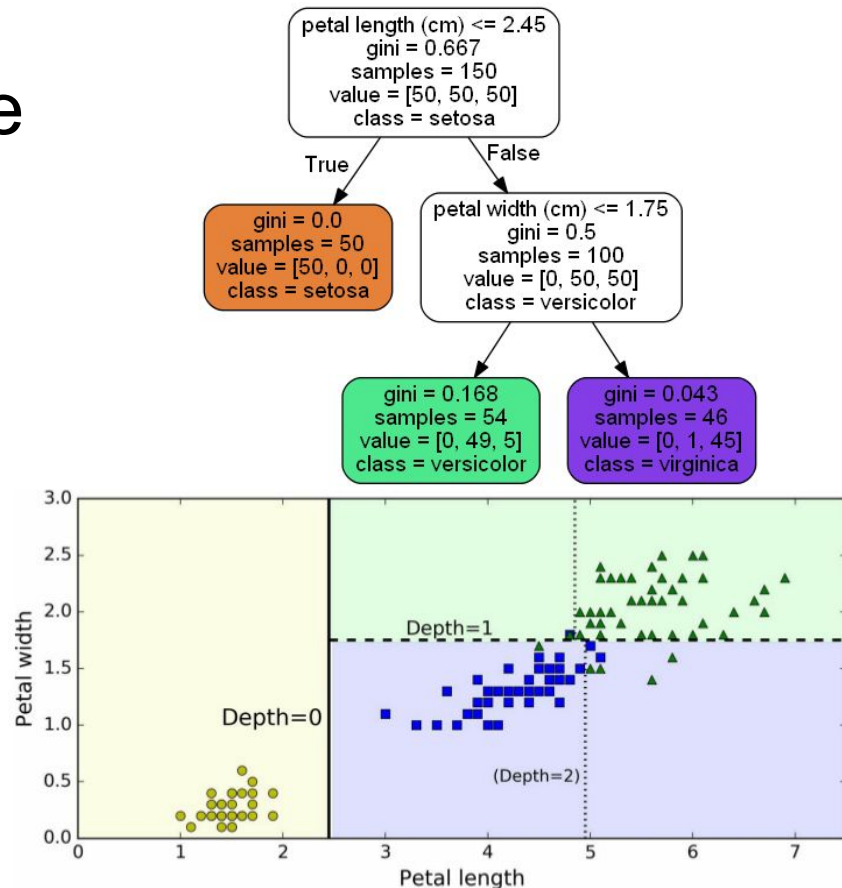


Figure 6-2. Decision Tree decision boundaries

# Time complexity decision tree

Growing the tree (**fit()**):  $O(m \cdot n^2 \log n)$

- Assuming we have continuous features and perform binary splits, the runtime of the decision tree construction is
- Sorting the values of continuous features helps with determining a decision threshold. If we have  $n$  examples, the sorting has time complexity  $O(n \log n)$
- If we have to compare sort  $m$  features, this becomes  $O(m \cdot n \log n)$
- Sorting step up to  $n/2$  times  $O(m \cdot n^2 \log n)$

Querying the tree (**predict()**):  $O(\log n)$

- depth of  $\log_2 n$

# Topics

1. Intro to decision trees
2. Recursive and divide & conquer strategy
3. **Types of decision trees**
4. Splitting criteria
5. Gini & Entropy vs misclassification error
6. Improvements & dealing with overfitting
7. Code example