



# Applied Machine Learning in Engineering

**Lecture 05 summer term 2025**

Prof. Merten Stender

Cyber-Physical Systems in Mechanical Engineering, Technische Universität Berlin

[www.tu.berlin/cpsme](http://www.tu.berlin/cpsme)

[merten.stender@tu-berlin.de](mailto:merten.stender@tu-berlin.de)

# Recap: Lecture 04

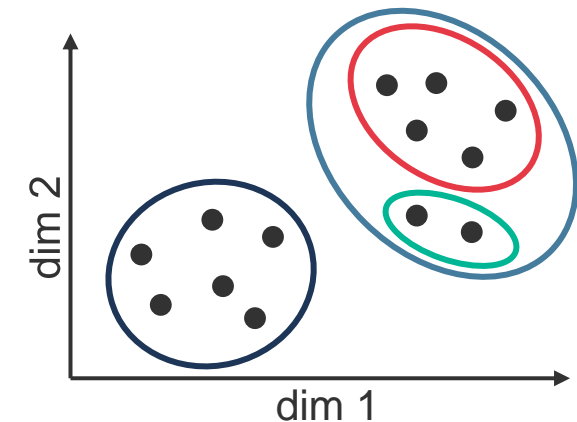
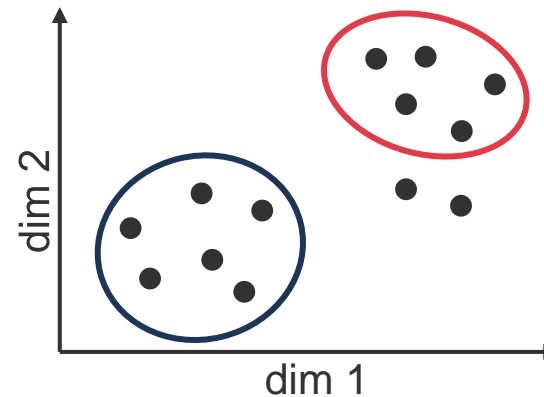
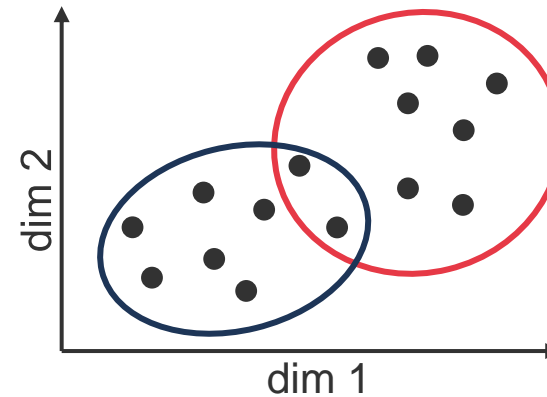


**Clustering** = the entity of clusters = the overall result of a clustering process

1. Nesting

2. Exclusiveness

3. Completeness

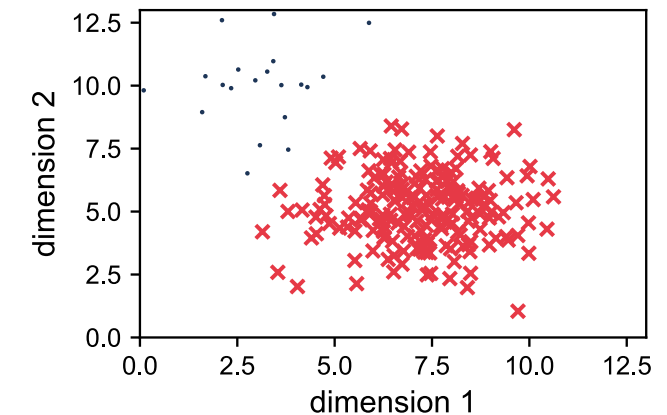
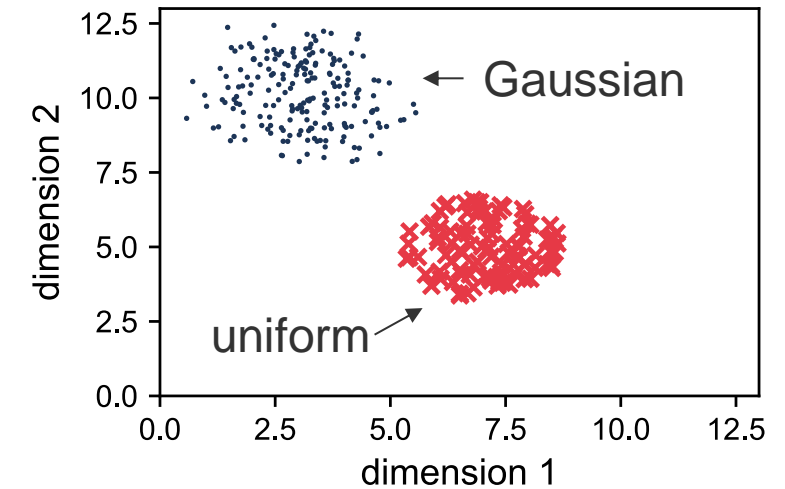
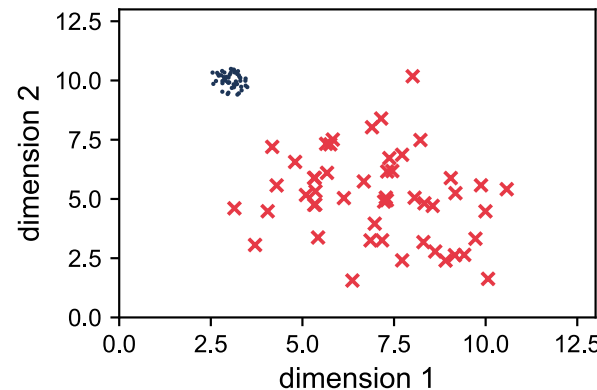


# Recap: Lecture 04



## Characteristics of Clusters

- **Distribution** of points within a cluster
  - Examples:
    - Gaussian distribution
    - Uniform distribution
- **Density**
  - High density clusters
  - Low density clusters
- **Size / variance**
  - Hypervolume consumed by cluster
  - Number of cluster members

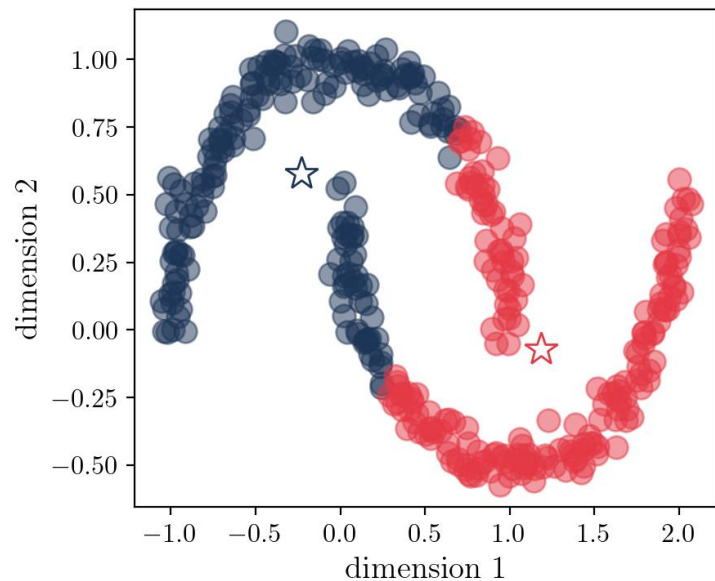


# Recap: Lecture 04



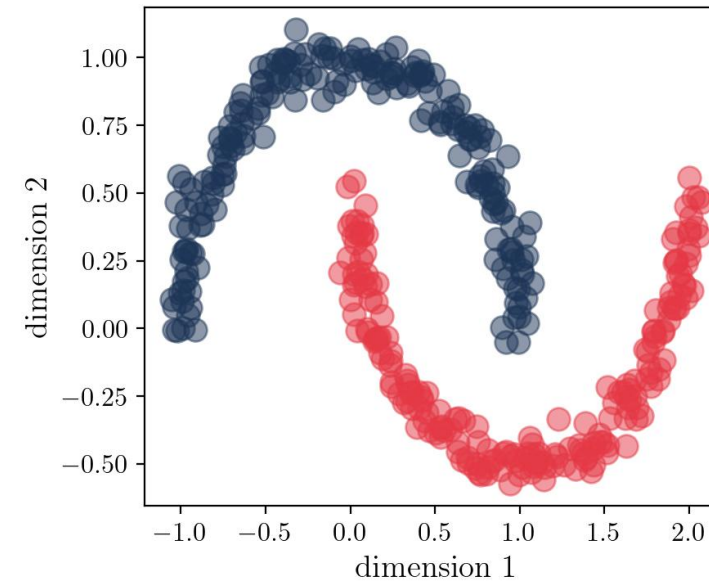
## K-means

- Built for prototype-based clusters (globular shape)
- No outlier handling (*exclusive* and *complete* clustering)



## DBSCAN\*

- Built for density-based clusters (any shape)
- Allows for incomplete clusterings (outliers without cluster assignment)



\* **Density-based spatial clustering of applications with noise**

# DBSCAN: Definition



- DBSCAN builds on 3 different types of data points:
- **Core point**  $\mathbf{x}_{\text{core}}$  has at least  $N_{\text{min}}$  points within  $\epsilon$  neighborhood (incl. itself). Interior of a cluster
- **Edge point**  $\mathbf{x}_{\text{edge}}$  is reachable from a core point within  $\epsilon$ , but is not a core point. Edge of a cluster
- **Outliers**  $\mathbf{x}_{\text{outlier}}$  is no core point and is not  $\epsilon$ -reachable from any core point. Not a cluster member



# Recap: Exercise 04



- From scratch implementation of Z-scoring
- initialization
- .fit method: calculate mean and std. deviation of the given data set
- .transform: z-scoring operation
- .inverse\_transform: reverse z-scoring operation

```
class ZScorer:

    def __init__(self):
        self._means: np.ndarray
        self._stds: np.ndarray

    def fit(self, data: np.ndarray):
        # Fits the class object to the data set,
        # extracts mean and std per feature column

        self._means = np.mean(data, axis=0)
        self._stds = np.std(data, axis=0)

    def transform(self, data: np.ndarray) -> np.ndarray:
        # returns a zscored data set

        # 1. subtract the mean from the data set, columnwise
        data_transformed = data - self._means

        # 2. divide by standart deviations (columnwise)
        data_transformed = data_transformed / self._stds

        return data_transformed

    def inverse_transform(self, data: np.ndarray) -> np.ndarray:
        # reverse the zscoring transformation

        # 1. multiply with std
        data_reversed = data * self._stds

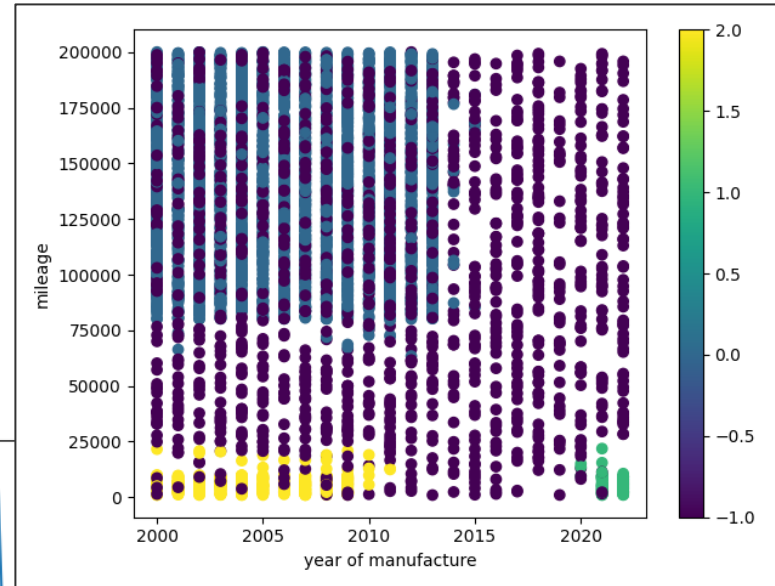
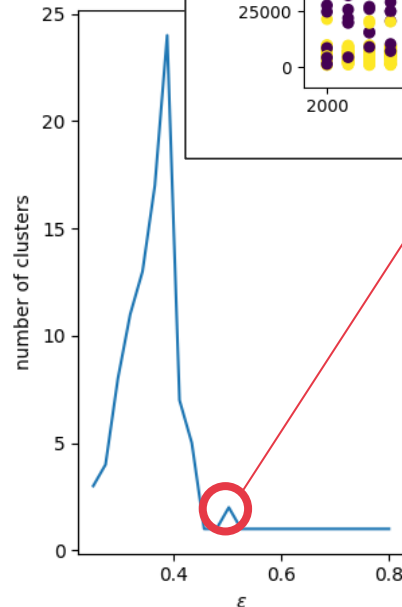
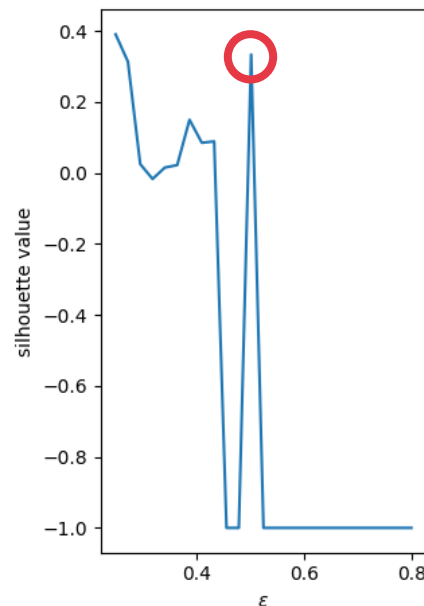
        # 2. add the means
        data_reversed = data_reversed + self._means

        return data_reversed
```

# Recap: Exercise 04



- Clustering data using **DBSCAN**
- **Hyperparameter optimization** through grid search for optimum silhouette value



```
# hyperparameter search for N_min, epsilon
nmin_grid = np.arange(10, 11)
epsilon_grid = np.linspace(0.25, 0.8, num=25)

silhouette_values = []
number_of_clusters = []
for _nmin in nmin_grid:
    for _epsilon in epsilon_grid:

        _dbscan = DBSCAN(eps=_epsilon, min_samples=_nmin)
        _labels = _dbscan.fit_predict(data_scaled)

        _num_clusters = num_clusters(_labels)
        if _num_clusters > 1:
            _silhouette_value = sil_coeff(data_scaled, _labels)
        else:
            _silhouette_value = -1

        silhouette_values.append(_silhouette_value)
        number_of_clusters.append(_num_clusters)

silhouette_values = np.array(silhouette_values)
number_of_clusters = np.array(number_of_clusters)
```



# Questions?



# Agenda



Cyber-Physical Systems  
in Mechanical Engineering TU Berlin

- Introduction to supervised learning
- Decision and Regression Trees
- Entropy and class purity measures

# Learning outcomes



Cyber-Physical Systems  
in Mechanical Engineering TU Berlin

## Learn to ...

- Estimate the value of high-quality labels for supervised learning
- Build a decision tree classifier
- Find optimal binary feature space segmentations

## Know about ...

- Measures of purity in populations of discrete events
- Shannon entropy
- Information gain and Gini Index
- Python: Recursive Functions



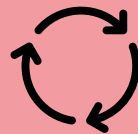
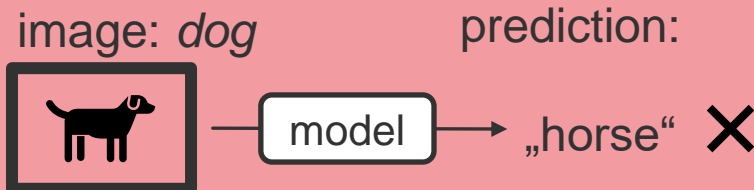
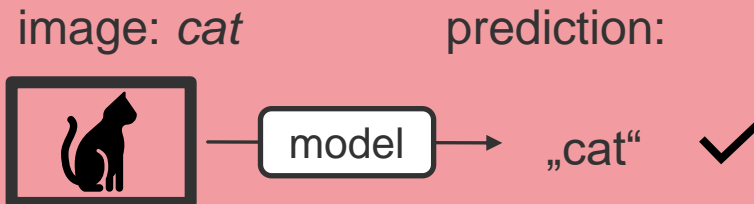
# Supervised Learning

# Supervised and Unsupervised Learning



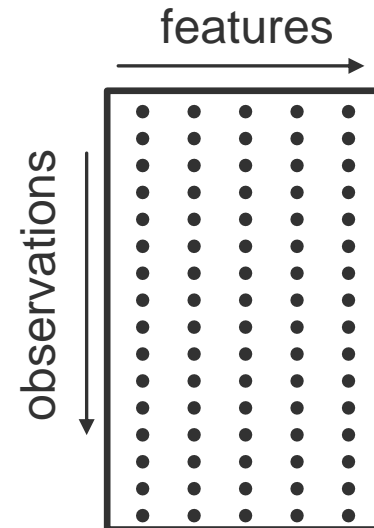
Cyber-Physical Systems  
in Mechanical Engineering TU Berlin

## supervised learning (predictive task)

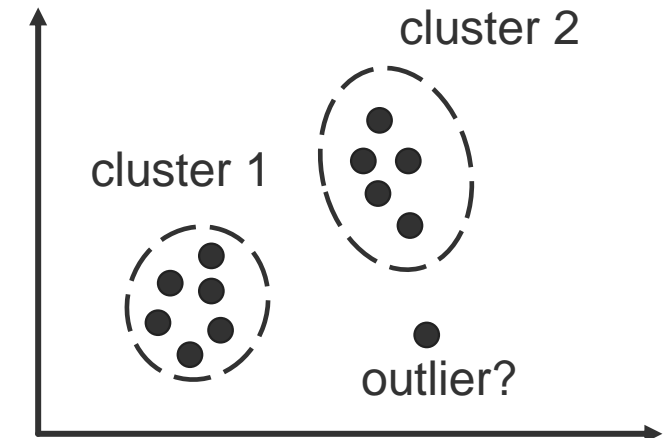


model training = reduce prediction error

## data (tabular)



## unsupervised learning (descriptive task)



finding clusters, groups and anomalies



**Supervised learning** = fitting prediction models to data for which ground truth targets exist

$$\mathcal{M}_{\theta}: \mathbf{X} \mapsto \mathbf{y}, \mathbf{X} \in \mathbb{R}^{N \times n}, \mathbf{y} \in \mathbb{R}^{N \times m}$$

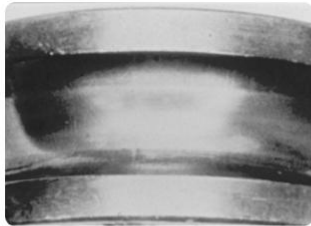
- **Ground truth data** (*'labels'*)
  - Desired target quantities  $y_i$ 
    - Allows comparing  $y_i$  against model predictions  $\hat{y}_i$ , prediction error  $E = \|\mathbf{y} - \hat{\mathbf{y}}\|$
    - Quantitative statements about model prediction quality
- **Model fitting:**
  - Reduction of error on training data set  $\min_{\theta} E(D_{\text{train}}, \theta)$  through optimization of  $\theta$
  - Model validation on hold-out validation data set  $D_{\text{val}}$
  - Under- and overfitting as potential issues
- **Classification and regression tasks**

# Application Cases in Engineering



Cyber-Physical Systems  
in Mechanical Engineering TU Berlin

## Structural Health Monitoring Predictive Maintenance Remaining lifetime prediction



1 Fine roughening or waviness



2 Small cracks



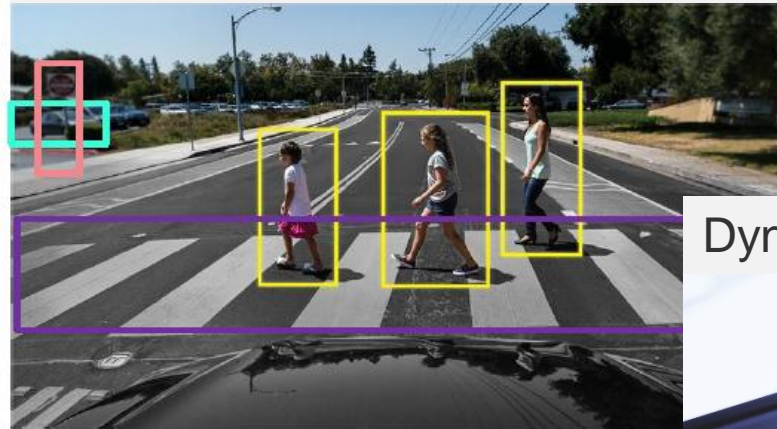
3 Local spalling



4 Spalling over the entire surface

SKF® Bearing damage and failure analysis

## Computer vision



## Dynamical behavior prediction



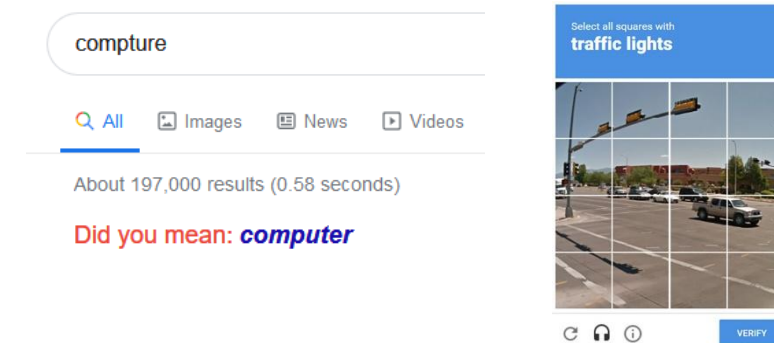
... and many more

# Generation of Targets



- **Extremely important** (*trash in – trash out*), yet tedious and expensive
  - Correct and high-quality ground truth targets are of **crucial importance**
  - Less but high-quality data should always be preferred over large and less-quality data

- Creative ways to generate labels:
  - Did you mean ... ? → grammar / language models
  - reCAPTCHA - *are you a robot?* → computer vision



- Professional data labeling services

- Read newspaper article in ZEIT:  
<https://www.zeit.de/2023/25/bild-annotation-kuenstliche-intelligenz-auto-indien>

aws	
AWS > Documentation > Amazon SageMaker	
Amazon SageMaker Developer Guide	
What Is Amazon SageMaker? Get Started Machine Learning Environments Autopilot: Automated ML Label Data Ground Truth	
Overview Features Pricing FAQs Developer Resources Customers	
Customized quote; fill out the <a href="#">project requirements form</a> .	
Object pricing details	
You are charged for the number of dataset objects that are reviewed. A dataset object is defined as an atomic unit of data across all modalities.	
Reviewed objects (images, video frames, text documents, audio files, etc.)	
Number of reviewed objects per month	
Price per reviewed object	
Less than 50,000 objects	\$0.08
50,000 to 1,000,000 objects	\$0.04
Greater than 1,000,000 objects	\$0.02

Amazon, 06 April 2023



# Decision Trees



# Decision Trees



**input features**  $x = [x_1, x_2, x_3]$

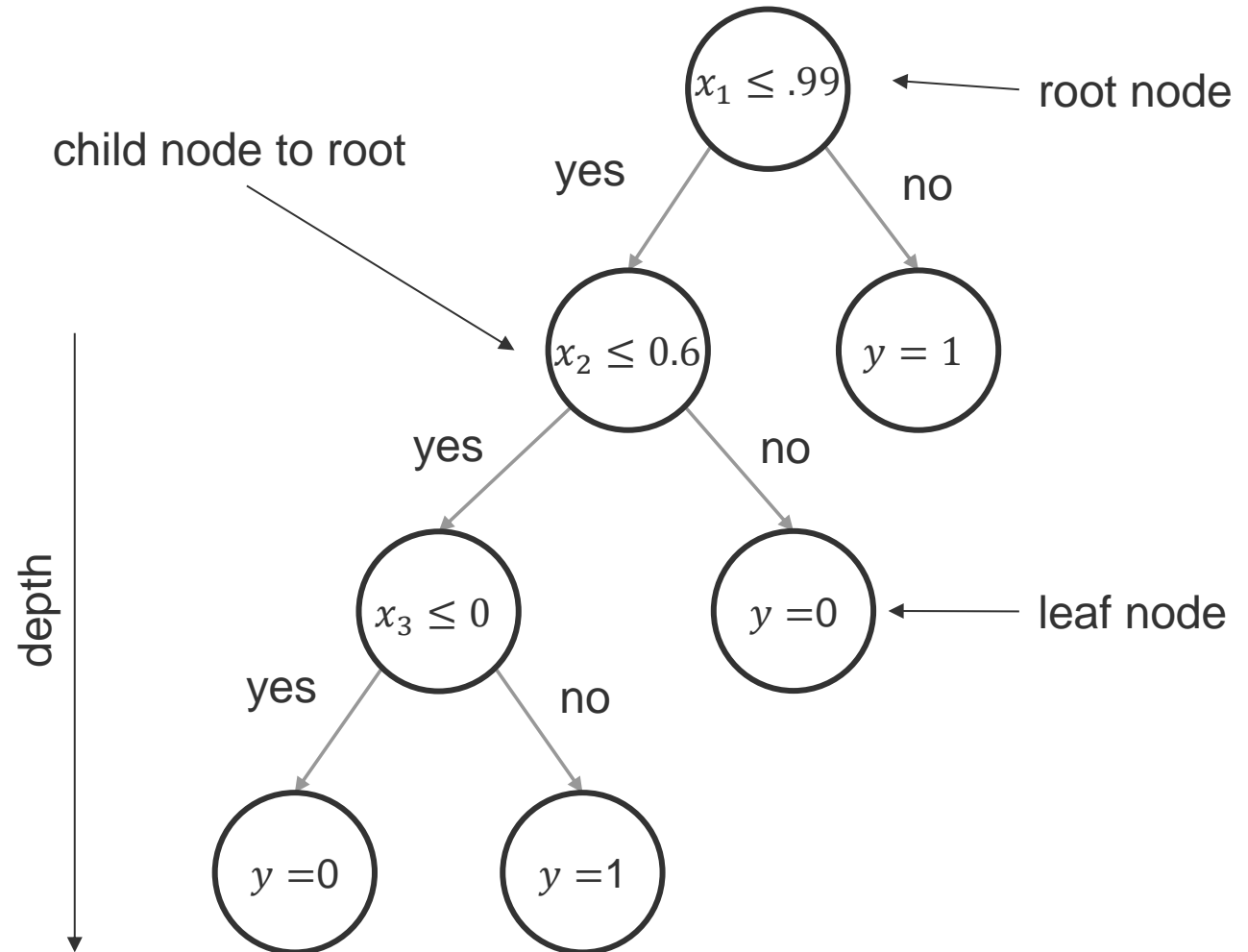
$x_1$ : sun is shining:  $\{0, 1\}$

$x_2$ : probability of rain:  $[0, 1.0]$

$x_3$ : ambient temperature:  $[-20, 40]^\circ\text{C}$

**target**  $y$

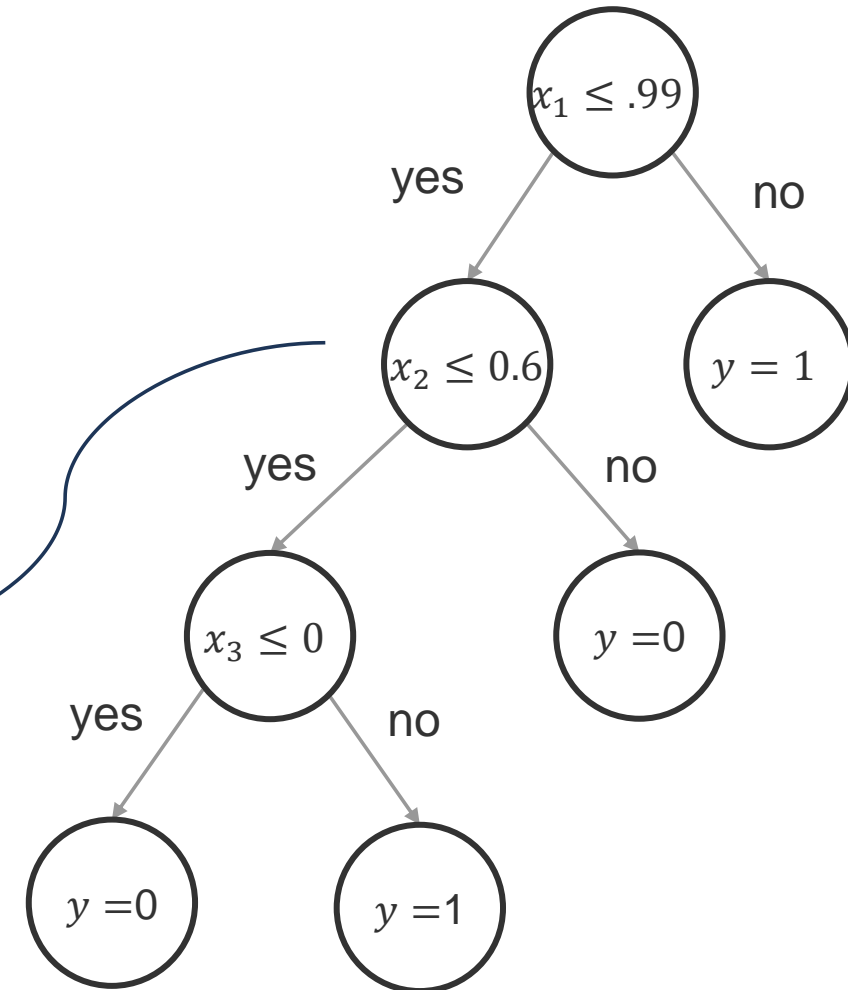
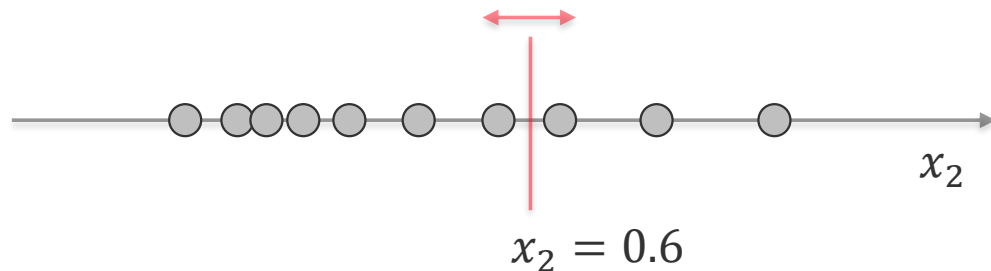
$y$ : ride bike to work?  $\{0, 1\}$



# Decision Trees



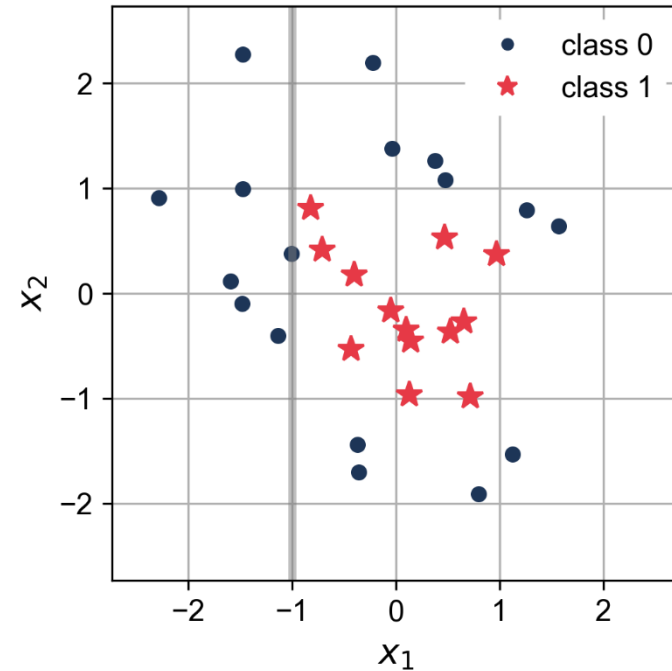
- Sequential decision rules partitioning the feature space
- Greedy algorithm
  - Previous splits not affected by current split
  - Algorithm does not 'look ahead' or back
- Recursive binary **feature space segmentation**



# Decision Trees: Example



- Aim: classify 2-dimensional data set with two classes (binary classification task)



Which split to do?

- Feature dimension ( $x_1, x_2$ )?
- Feature value ( $x^*$ ) ?

- Split condition 1:  $x_1 \leq -1.0$

child 1 (condition true)

7 × ●  
0 × ★

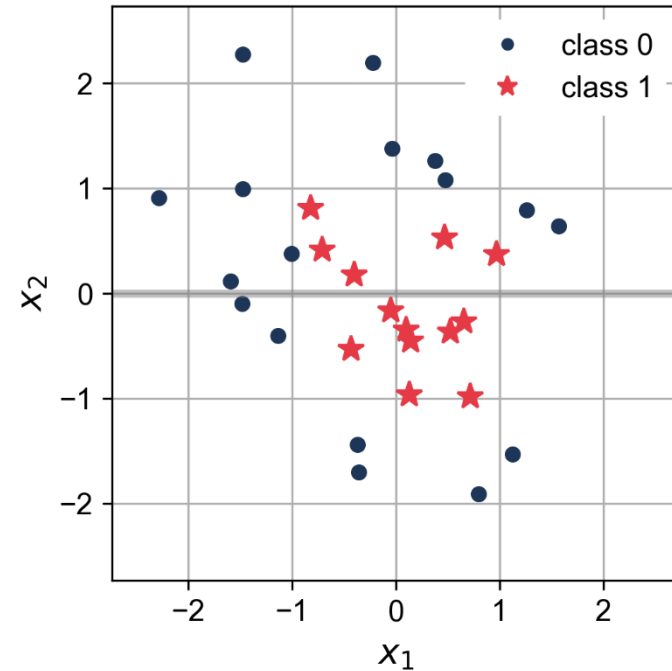
child 2 (condition false)

10 × ●  
13 × ★

# Decision Trees: Example



- Aim: classify 2-dimensional data set with two classes (binary classification task)



Which split to do?

- Feature dimension ( $x_1, x_2$ )?
- Feature value ( $x^*$ ) ?

- Split condition 2:  $x_2 \leq 0$

child 1 (condition true)

6 × ●  
8 × ★

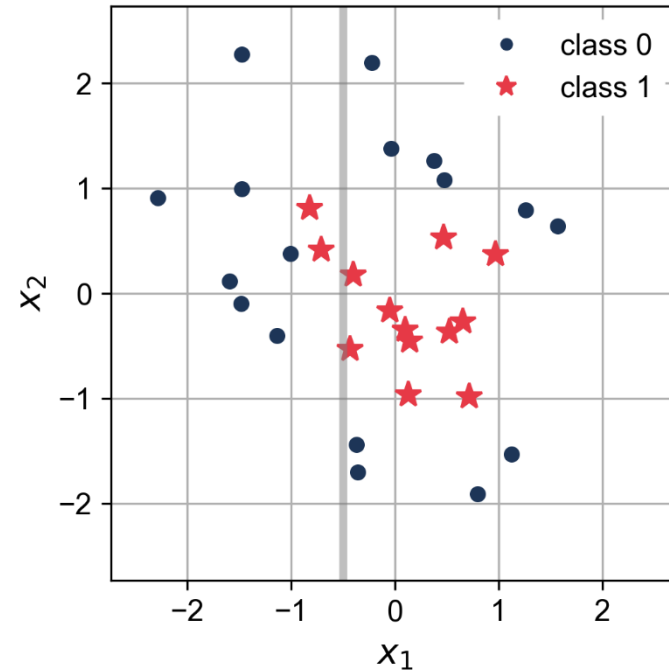
child 2 (condition false)

11 × ●  
5 × ★

# Decision Trees: Example



- Aim: classify 2-dimensional data set with two classes (binary classification task)



Which split to do?

- Feature dimension ( $x_1, x_2$ )?
- Feature value ( $x^*$ ) ?
- Split condition 3:  $x_1 \leq -0.5$

child 1 (condition true)

7 × ●  
2 × ★

child 2 (condition false)

10 × ●  
11 × ★

# Selecting the Best Split



- Full data set



- Data split
  - Which to select?



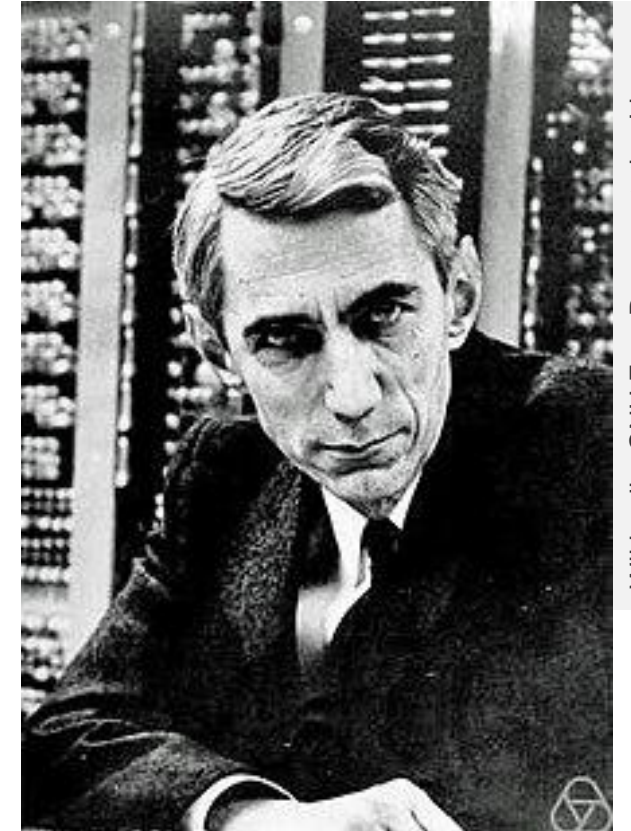


$$H(\mathbf{x}) = - \sum_{i \in \mathcal{C}} P(\mathbf{x}_i) \log P(\mathbf{x}_i)$$

**Claude E. Shannon (1916-2001)**

Founder of information theory

*1948: A Mathematical Theory of Communication*



Wikipedia, GNU Free Documentation License

# Shannon Entropy: Definition



- Also denoted *information entropy* or *entropy index*

$$H(x) = - \sum_{i \in C} P(x_i) \log_a(P(x_i)) = \sum_{i \in C} P(x_i) \log_a \left( \frac{1}{P(x_i)} \right) \in [0, ]$$

- $C = \{c_1, c_2, c_3\}$  set of distinct classes
- $P(x_i)$  probability of a single event  $i$ :
  - fraction of population composed of a single species  $i$
- $H(x)$  amount of information gained by observing an event of probability  $P(x_i)$

$a$  base of the logarithm

The unit of entropy depends on the base of the logarithm

- Computer science:  $a = 2$   
→ unit *bits*
- Euler's number:  $a = e$   
→ unit *nats*



# Shannon Entropy: Intuition



- Intuition and edge cases

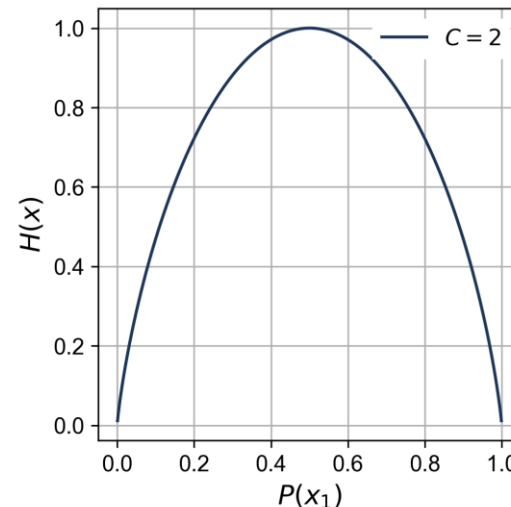
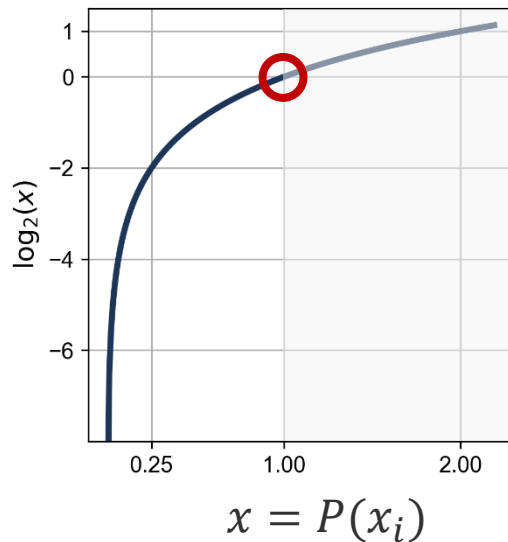
$$H(x) = - \sum_{i \in \mathcal{C}} P(x_i) \cdot \log_2(P(x_i))$$

$C = 2$  (binary classification):  $P(x_1) + P(x_2) = 1$

$$H(x) = -P(x_1) \cdot \log_2(P(x_1)) - P(x_2) \cdot \log_2(P(x_2))$$

$$H(x) = -P(x_1) \cdot \log_2(P(x_1)) - (1 - P(x_1)) \cdot \log_2(1 - P(x_1))$$

Log-2



- $P(x_i) = 1 \rightarrow \text{entropy} = 0$
- $\max(H(x))_{C=2} = 1.0$  for  $P(x_1) = P(x_2)$

Arbitrary  $C$ :

- $\max(H(x))_C = -C \cdot \left(\frac{1}{C} \log \frac{1}{C}\right) = -\log \frac{1}{C}$

# Shannon Entropy: Example



$$H(x) = - \sum_{i \in C} P(x_i) \log_2(P(x_i))$$

- Example: **calculate** the uncertainty coming with a certain character appearing next
  - sequence of numbers  $[2\ 3\ 0\ 2\ 7\ 1]$
  - probabilities:  $P(0) = 1/6, P(1) = 1/6, P(2) = 2/6, P(3) = 1/6, P(7) = 1/6$
  - entropy  $H(x) = 2.2516$
- Edge case 1: Outcome is certain: vanishing entropy  $H(x) = 0$  , e.g.,  $[2\ 2\ 2\ 2\ 2\ 2]$
- Edge case 2: The more proportional the frequencies of occurrence are, the harder it gets to make a prediction, hence the larger the entropy  $H(x) \gg 0$ , e.g.,  $[1\ 2\ 3\ 4\ 5\ 6]$



# Selecting the best split



- Root (parent) data set



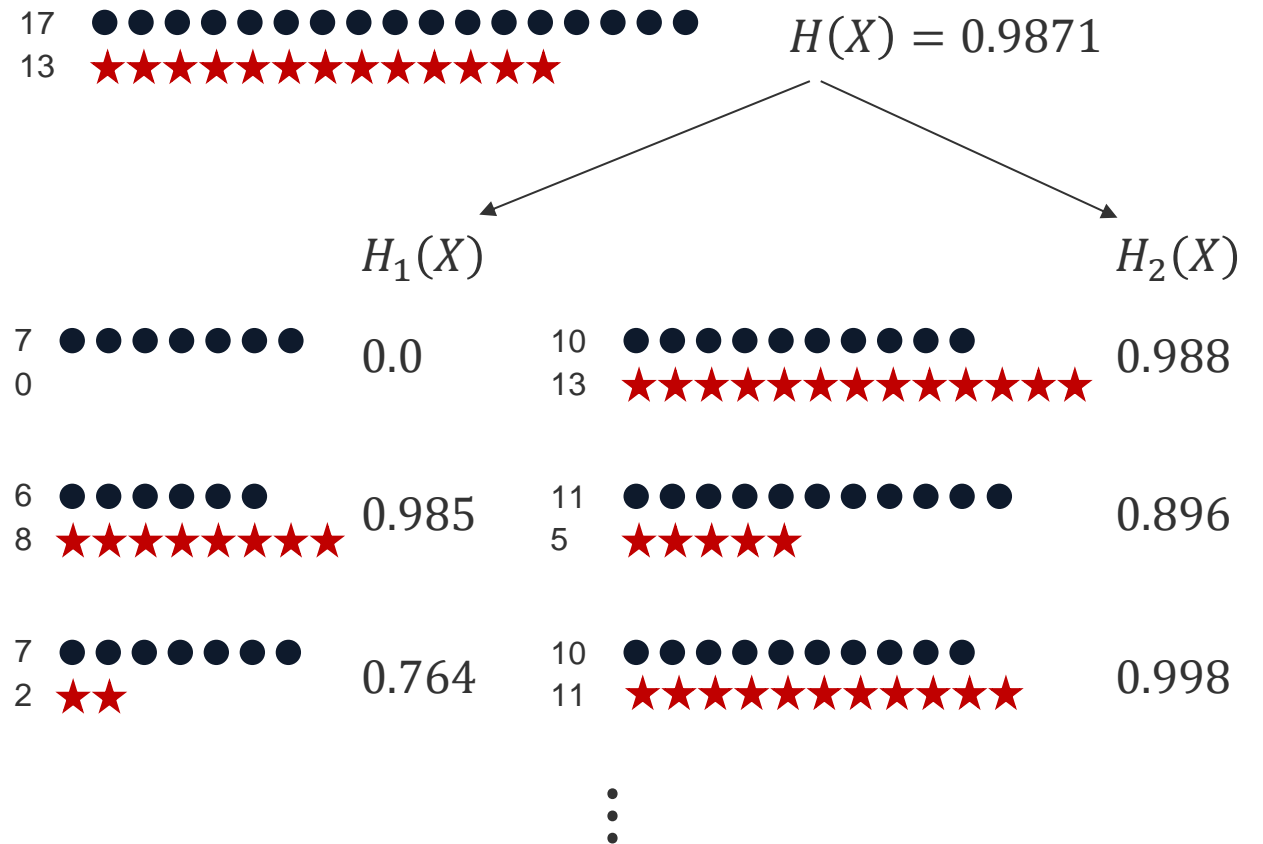
- Compute** the entropy of every possible data split

- What now? Which split to select?  
→ **Information gain**

Split 1

Split 2

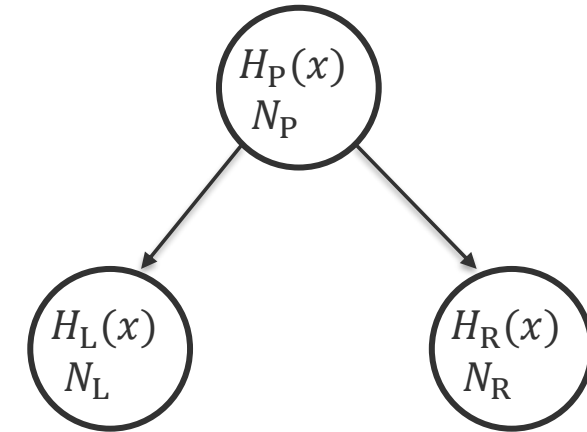
Split 3



# Information Gain



- DT: segmentation of the feature space
- Aim: maximal class purity of the sub-data set
- Purity metric: entropy
  - Entropy at parent node:  $H_P(x)$
  - Entropy at children:  $H_{L,R}(x)$
  - Number of samples:  $N_{P,L,R}$



- Information gain

$$I(x) = H_P(x) - \underbrace{\left( \frac{N_L}{N_P} \cdot H_L(x) + \frac{N_R}{N_P} \cdot H_R(x) \right)}$$

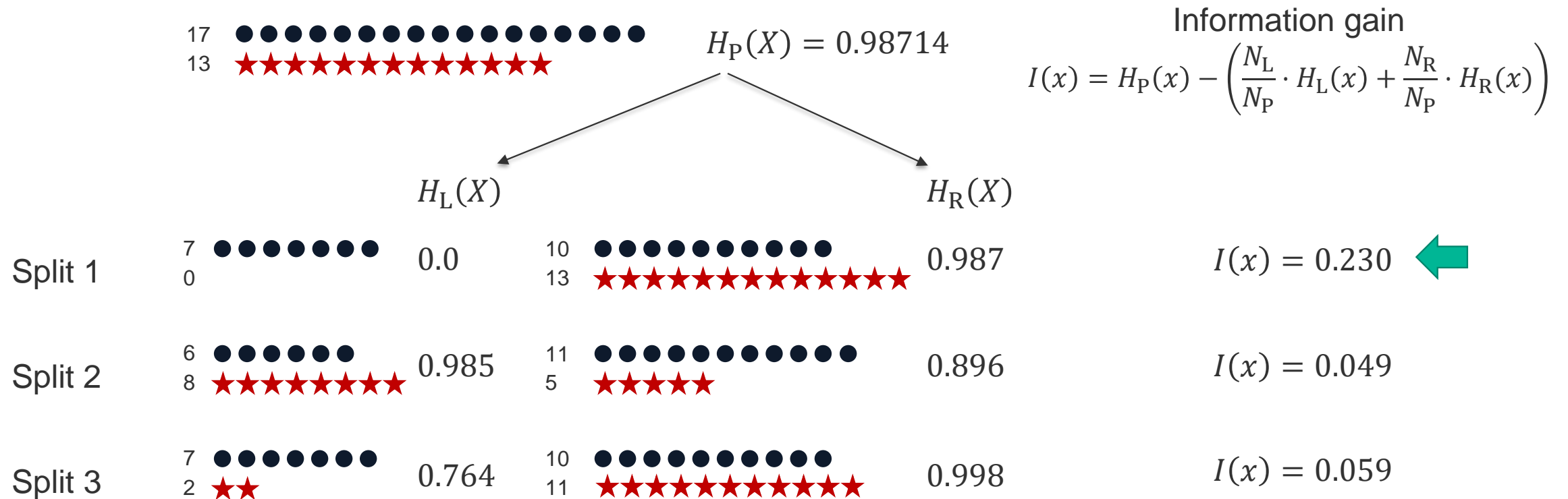
- Split for maximum information gain

$$x^* = \max(I(x))$$

Maximum gain:  $H_L = H_R = 0$

entropy in both child nodes  
vanishes (pure child nodes)

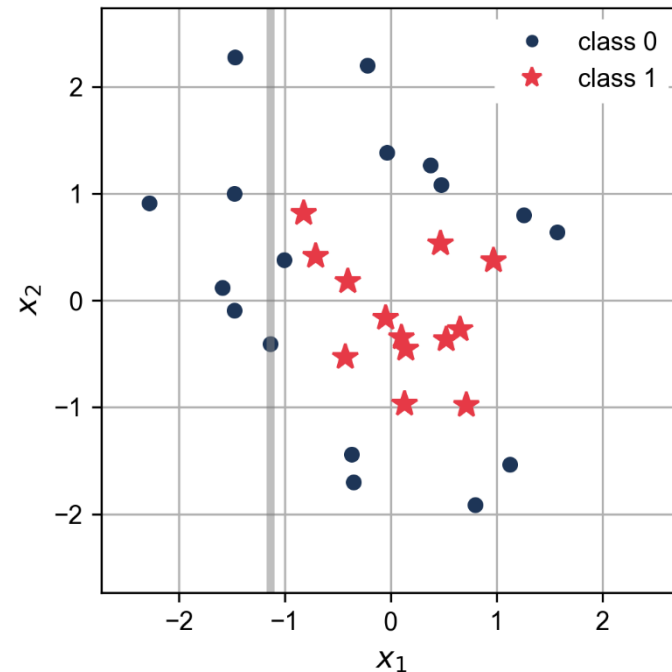
# Best Split: Information Gain



# Decision Trees: Example



- Previous slides: 3 exemplary data splits
- Actual: 58 splits possible (29 for 1<sup>st</sup> feature dimension, 29 for 2<sup>nd</sup> feature dimension)  
[excluding the splits that would put the complete data set into a single child, hence assuming a non-optimal purity in the data set to be split]



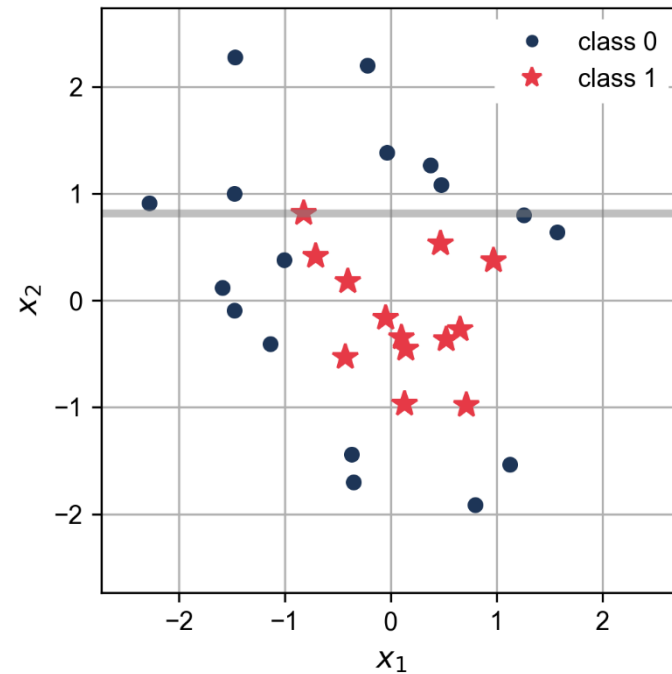
## Greedy approach:

- Computation of all information gains
- Top 3 splits:
  - (3) condition:  $x_1 \leq -1.137 \rightarrow$  information gain  $I(x) = 0.191$

# Decision Trees: Example



- Previous slides: 3 exemplary data splits
- Actual: 58 splits possible (29 for feature dimension 1, 29 for feature dimension 2)



## Greedy approach:

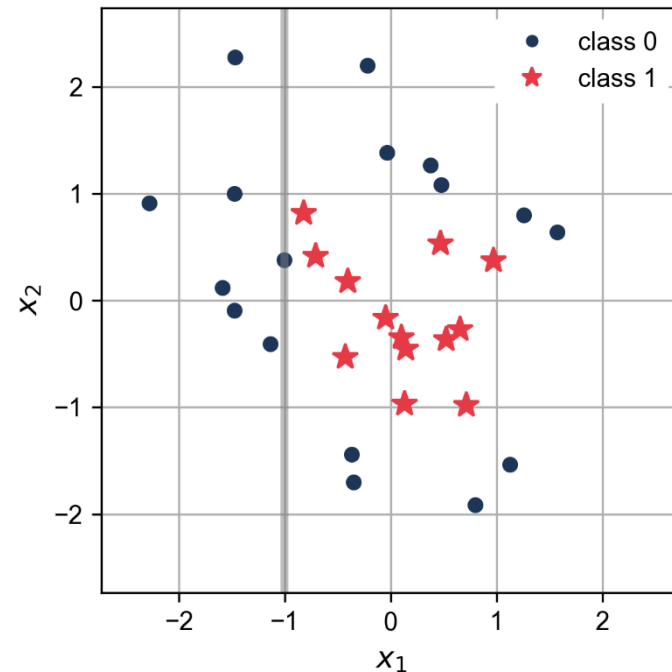
- Computation of all information gains
- Top 3 splits:
  - (3) condition:  $x_1 \leq -1.137 \rightarrow$  information gain  $I(x) = 0.191$
  - (2) condition:  $x_2 \leq 0.813 \rightarrow$  information gain  $I(x) = 0.229$



# Decision Trees: Example



- Previous slides: 3 exemplary data splits
- Actual: 58 splits possible (29 for feature dimension 1, 29 for feature dimension 2)



## Greedy approach:

- Computation of all information gains
- Top 3 splits:
  - (3) condition:  $x_1 \leq -1.137 \rightarrow$  information gain  $I(x) = 0.191$
  - (2) condition:  $x_2 \leq 0.813 \rightarrow$  information gain  $I(x) = 0.229$
  - (1) condition:  $x_1 \leq -1.007 \rightarrow$  information gain  $I(x) = 0.229$

# Decision Trees: Example



- **First split** of the data set:

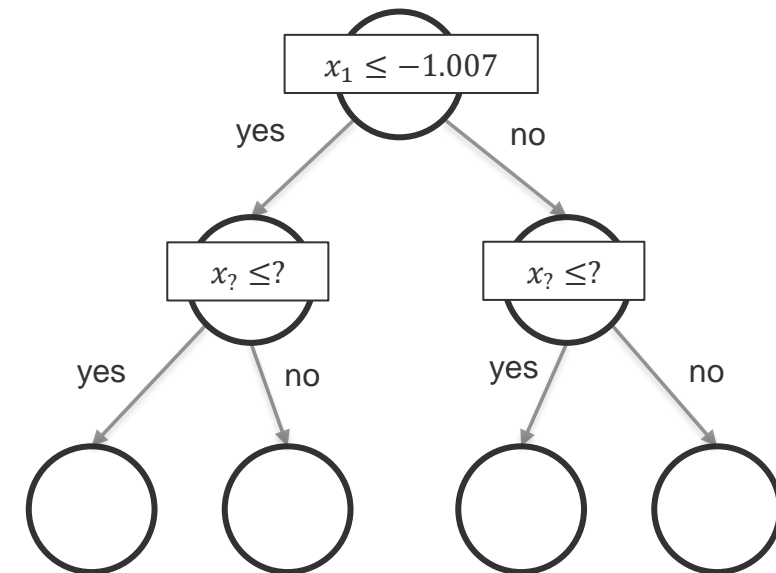
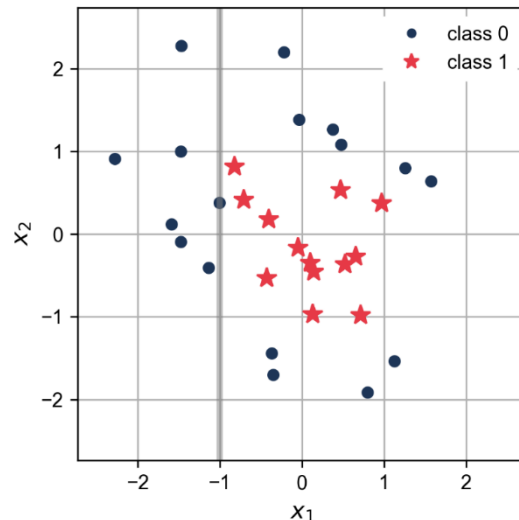
- Left child node:  $x_1 \leq -1.007$
- Right child node:  $x_1 > -1.007$

7 members of class 0

10 members of class 0, 13 members of class 1

- **Second split:**

- Optimal (information gain) split of child nodes
- Here: left child node is pure  $\rightarrow$  no more splitting!



$N_1 = ?$   
 $H_1(x) =$

$N_2 = ?$   
 $H_2(x) =$

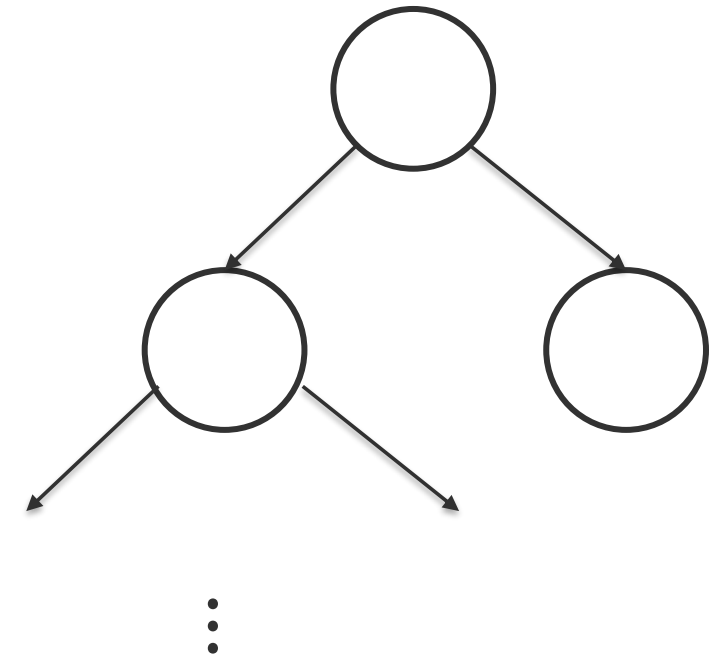
$N_1 = ?$   
 $H_1(x) =$

$N_2 = ?$   
 $H_2(x) =$

# Stopping Criteria




- Aim: maximize purity in leaf nodes
- Without any constraints there is a solution with  $H(x) = 0$  in each leaf node
  - Worst case:  $N = 1$  samples per leaf
  - Overfitting the data
- Excessive splitting leads to **overfitting**:
  - Very strong performance on training data set
  - Weak performance on new (unseen) data
- **Constraints** to tree growing
  - Minimum number of samples per node  $N_{\min}$
  - Maximum depth of tree  $D_{\max}$
- **Impure leaves**: return the most-common class label



# scikit-learn: DecisionTreeClassifier



Cyber-Physical Systems  
in Mechanical Engineering TU Berlin

[Install](#) [User Guide](#) [API](#) [Examples](#) [Community](#) [More](#) 

[Prev](#) [Up](#) [Next](#)

**scikit-learn 1.2.2**  
[Other versions](#)

Please [cite us](#) if you use the software.

[sklearn.tree.DecisionTreeClassifier](#)  
[DecisionTreeClassifier](#)  
Examples using  
[sklearn.tree.DecisionTreeClassifier](#)

## sklearn.tree.DecisionTreeClassifier

```
class sklearn.tree.DecisionTreeClassifier(*, criterion='gini', splitter='best', max_depth=None, min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features=None, random_state=None, max_leaf_nodes=None, min_impurity_decrease=0.0, class_weight=None, ccp_alpha=0.0)
```

[\[source\]](#)

A decision tree classifier.

Read more in the [User Guide](#).

Parameters:	
<b>criterion</b> : {"gini", "entropy", "log_loss"}, default="gini"	The function to measure the quality of a split. Supported criteria are "gini" for the Gini impurity and "log_loss" and "entropy" both for the Shannon information gain, see <a href="#">Mathematical formulation</a> .
<b>splitter</b> : {"best", "random"}, default="best"	The strategy used to choose the split at each node. Supported strategies are "best" to choose the best split and "random" to choose the best random split.
<b>max_depth</b> : int, default=None	The maximum depth of the tree. If None, then nodes are expanded until all leaves are pure or until all leaves contain less than min_samples_split samples.
<b>min_samples_split</b> : int or float, default=2	The minimum number of samples required to split an internal node: <ul style="list-style-type: none"><li>If int, then consider min_samples_split as the minimum number.</li><li>If float, then min_samples_split is a fraction and <code>ceil(min_samples_split * n_samples)</code> are the minimum number of samples for each split.</li></ul>

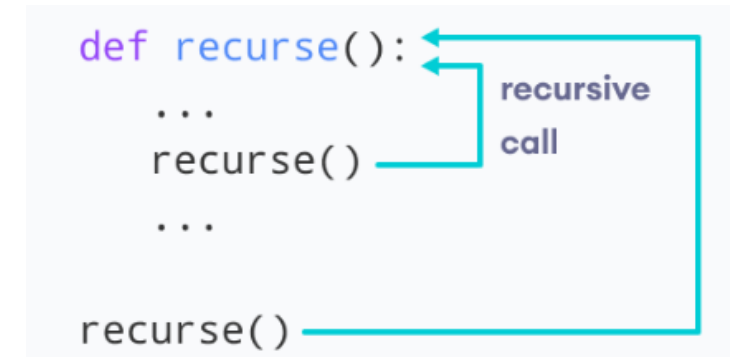


# Python: Recursive Functions

# Recursive Functions

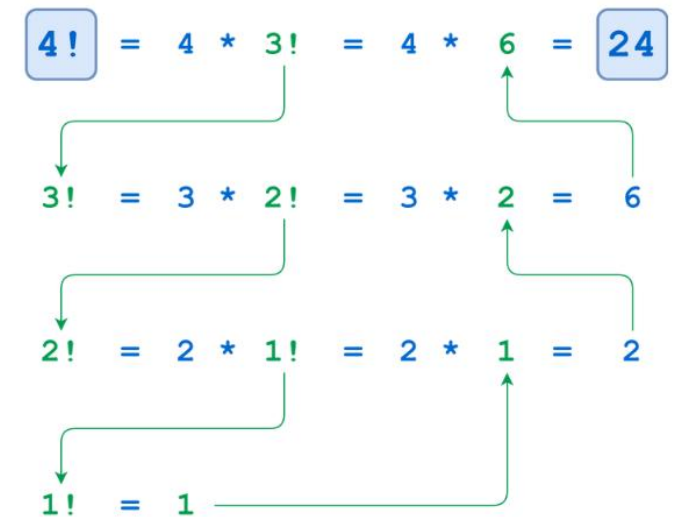


- A recursive function **calls itself** until some break condition (so-called „base condition“) is met
- Let the user created nested function calls or object structures
- Example: computing the factorial of a number



- Get the Python recursion limit (after which Python will stop automatically)

```
from sys import getrecursionlimit  
getrecursionlimit()
```



<https://realpython.com/python-recursion/>

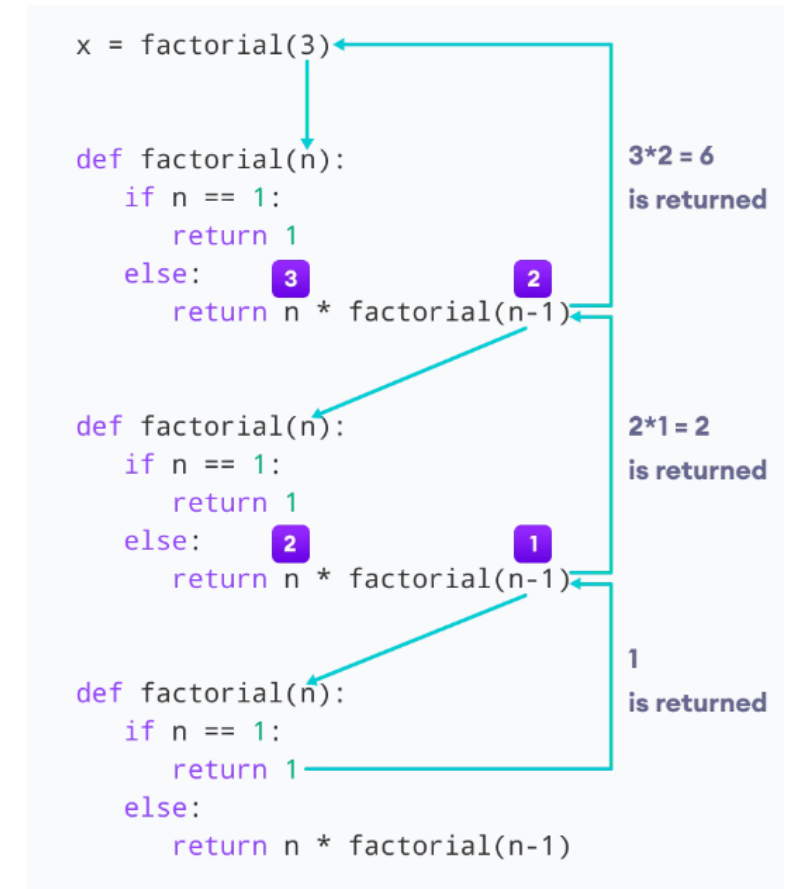
# Recursive Functions



- A recursive function **calls itself** until some break condition is met
- Example: computing the factorial of a number

```
def factorial(x: int) -> int:
    """Compute the factorial of an
    integer using a recursive function"""
    if x == 1: ← stopping condition
                (,base condition')
        return 1
    else:
        return x * (factorial(x-1))

print(f'factorial of x=3 is {factorial(3)}')
```



From <https://www.programiz.com/python-programming/recursion>



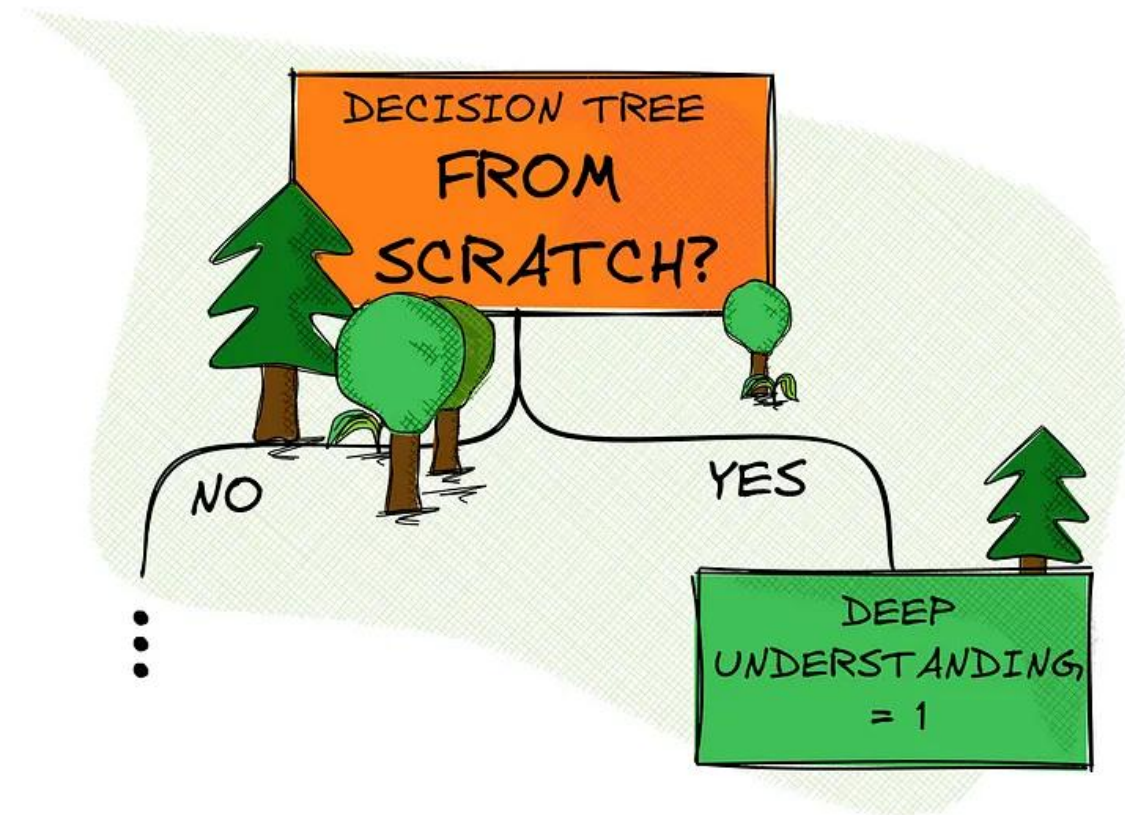
# Exercise 05



# Exercise 05



- Implementation of a decision tree from scratch
- Implement entropy
- Implement information gain
- Implement a splitting routine
- Implement a best split routine



© Marvin Lanhenke, <https://towardsdatascience.com/implementing-a-decision-tree-from-scratch-f5358ff9c4bb>



# Questions?