

understanding_decision_trees.ipynb

Iris Dataset composed of:

4 features = petal length, petal width, sepal length & sepal width

Target Variable = Iris Species

↳ 3: ⁰setosa, ¹versicolour, ²virginica

DATA: 150 observations (Rows), 50 for each iris = **Balanced DATA**

For easier understanding work w/ just petal width & sepal width

• **Confusion Matrix**: y-axis is target
x-axis is species decision tree predicted

• 2 feature decision tree 93% accurate

• **How was the Decision Tree Built?**

- Split data into 2

1. Training Set

2. Testing Set

↳ used to measure performance

★ Goal of Decision Tree is to split training set into homogenous areas where only 1 iris species is present according to given features

Node 0: Root Node

- First decision boundary of tree:

petal width = 0.8cm
↳ easily visible on graph

↳ decided by testing all possible decision splitting dataset
→ choose 1 that minimizes Gini Impurity of 2 splits

↳ measures probability from randomly chosen element to be incorrectly classified

all possible classes 3

$$\begin{aligned} \Rightarrow \text{probability choosing an element times probability misclassified} \\ \sum_{i=1}^3 p_i \sum_{k \neq i} p_k &= \sum_{i=1}^3 p_i (1 - p_i) = \sum_{i=1}^3 (p_i - p_i^2) = \sum_{i=1}^3 p_i - \sum_{i=1}^3 p_i^2 \\ &= 1 - \sum_{i=1}^3 p_i^2 \end{aligned}$$

used to perform Gini Test

③ How does the algo split?

- Tries all possible boundaries along all features
- compute Gini impurity of 2 groups
- Choose boundary w/ lowest Gini impurity for each Group

(only Se's)

- Left Node $Gini = 0$ so becomes a leaf & doesn't get split again
- Right node $GI = 2.5$ cuz there's about equal # of V_0 & V_1
- creating group w/ 2 not always best cuz you might not have all of em
- Overfitting can happen if we don't limit size of Tree