What the Code Actually Does (The Step-by-Step Process)

Think of this process as teaching a student for an exam. You give them study materials, they learn, and then you test them on questions they haven't seen before.

1. Loading and Preparing the "Study Material" 🤚

First, the code loads the bank.csv file. This file is our historical data—our textbook. It contains thousands of records of past customers.

However, computers don't understand words like 'married', 'technician', or 'yes'. They only understand numbers.

- **Feature Engineering**: We convert all the categorical text columns (like job, marital, education) into numbers using a technique called **one-hot encoding**. It turns each category into a new column with a 1 (if the customer belongs to that category) or a 0 (if they don't).
- Label Encoding: We convert our target variable, the 'deposit' column, from 'yes'/'no' into 1/0. This is the "answer key" the model will learn from.

2. Splitting the Data (Training vs. Testing)

You wouldn't test a student on the exact same questions they studied, right? You want to know if they actually *learned the concepts*. We do the same thing here.

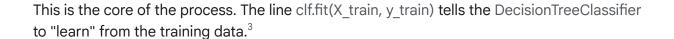
We split our dataset into two parts:

• Training Set (80% of the data): This is the "study material." The model will analyze this

data to find patterns and create its decision rules.²

• Testing Set (20% of the data): This is the "final exam." The model has never seen this data before. We use it at the very end to check how well the model learned to make predictions on new, unseen customers.

3. Training the Model (The "Learning" Phase)



The algorithm sifts through all the customer features (age, job, balance, etc.) and asks thousands of questions, like:

- "Is the duration of the last call less than 500 seconds?"
- "Is the customer's age over 40?"
- "Does the customer have a housing loan?"

For each question, it measures how well that question splits the customers into two groups: those who subscribed (1) and those who didn't (0). It selects the question that creates the "purest" splits (i.e., the split that does the best job of separating the 1s from the Os). It repeats this process over and over, building a tree of questions.⁴

4. Evaluating the Model (Checking the "Exam" Score)

After the model is trained, we use the testing set (the 20% of data it never saw) to evaluate it.

- **Prediction**: The model looks at each customer in the test set and follows its decision tree rules to predict whether they will subscribe (1) or not (0).⁵
- **Comparison**: We then compare the model's predictions to the actual answers in the test set.
- Accuracy: The final accuracy score (e.g., 88.48%) tells us what percentage of the predictions on the unseen data were correct.

What is the Decision Tree? (The "Flowchart")



The decision tree is the final output of the learning process. 6 It's a model of the decisions the computer learned. It's just like a flowchart. You start at the top and work your way down.

Let's look at the first few lines of the tree rules you generated:

```
I--- duration <= 699.50
| |--- duration <= 321.50
```

Here's how to read it for a new customer:

- 1. Start at the Root Node (the first line): The most important question is: "Was the duration of the last call less than or equal to 699.5 seconds?"
- 2. Follow the Branch:
 - o If YES, you go down the first path.⁸ The next question is: "Was the duration also less than or equal to 321.5 seconds?"
 - If NO (meaning the duration was greater than 699.5), you would skip to the other main branch that starts with I--- duration > 699.50.
- 3. Continue Down the Tree: You keep answering guestions about the customer (housing yes <= 0.50 means "Does the customer NOT have a housing loan?"), following the branches until you reach a final conclusion, which is a class: 0 (predicts 'no') or class: 1 (predicts 'yes').9

So, the tree is simply the set of rules the model will use to classify any new customer you give it. 10 The features closer to the top of the tree (like duration) are the most powerful predictors.