

# Decision Tree-2

## Assignment Questions



You are a data scientist working for a healthcare company, and you have been tasked with creating a decision tree to help identify patients with diabetes based on a set of clinical variables. You have been given a dataset (diabetes.csv) with the following variables:

1. **Pregnancies:** Number of times pregnant (integer)
2. **Glucose:** Plasma glucose concentration a 2 hours in an oral glucose tolerance test (integer)
3. **BloodPressure:** Diastolic blood pressure (mm Hg) (integer)
4. **SkinThickness:** Triceps skin fold thickness (mm) (integer)
5. **Insulin:** 2-Hour serum insulin (mu U/ml) (integer)
6. **BMI:** Body mass index (weight in kg/(height in m)<sup>2</sup>) (float)
7. **DiabetesPedigreeFunction:** Diabetes pedigree function (a function which scores likelihood of diabetes based on family history) (float)
8. **Age:** Age in years (integer)
9. **Outcome:** Class variable (0 if non-diabetic, 1 if diabetic) (integer)

Here's the dataset link: [https://drive.google.com/file/d/1Q4J8KS1wm4-\\_YTuc389enPh6O-eTNcx2/view?usp=sharing](https://drive.google.com/file/d/1Q4J8KS1wm4-_YTuc389enPh6O-eTNcx2/view?usp=sharing)

Your goal is to create a decision tree to predict whether a patient has diabetes based on the other variables. Here are the steps you can follow:

**Q1. Import the dataset and examine the variables. Use descriptive statistics and visualizations to understand the distribution and relationships between the variables.**

**Q2. Preprocess the data by cleaning missing values, removing outliers, and transforming categorical variables into dummy variables if necessary.**

**Q3. Split the dataset into a training set and a test set. Use a random seed to ensure reproducibility.**

**Q4. Use a decision tree algorithm, such as ID3 or C4.5, to train a decision tree model on the training set. Use cross-validation to optimize the hyperparameters and avoid overfitting.**

**Q5. Evaluate the performance of the decision tree model on the test set using metrics such as accuracy, precision, recall, and F1 score. Use confusion matrices and ROC curves to visualize the results.**

**Q6. Interpret the decision tree by examining the splits, branches, and leaves. Identify the most important variables and their thresholds. Use domain knowledge and common sense to explain the patterns and trends.**

**Q7. Validate the decision tree model by applying it to new data or testing its robustness to changes in the dataset or the environment. Use sensitivity analysis and scenario testing to explore the uncertainty and risks.**

By following these steps, you can develop a comprehensive understanding of decision tree modeling and its applications to real-world healthcare problems. Good luck!

**Note:** Create your assignment in Jupyter notebook and upload it to GitHub & share that github repository link through your dashboard. Make sure the repository is public.