



# BERLIN SCHOOL OF BUSINESS & INNOVATION

**Essay Title: Describe a data analytics problem and write a python code for that**

**Name: Smit Kamleshbhai Kevadiya**

**Date: 03/03/2023**

## **Statement of compliance with academic ethics and the avoidance of plagiarism**

I honestly declare that this essay is entirely my own work and none of its part has been copied from printed or electronic sources, translated from foreign sources and reproduced from essays of other researchers or students. Wherever I have been based on ideas or other people texts I clearly declare it through the good use of references following academic ethics.

(In the case that is proved that part of the essay does not constitute an original work, but a copy of an already published essay or from another source, the student will be expelled permanently from the postgraduate program).

Name and Surname (Capital letters):

SMIT KAMLESHBHAI KEVADIYA

.....

Date: 03/03/2023

## TABLE OF CONTENTS

TABLE OF CONTENTS.....	3
INTRODUCTION.....	4
DECISION TREE ALGORITHM DATASET.....	5
VISUALIZATION FOR VISUALIZE THE DATASET.....	9
PYTHON CODE FOR THE DECISION TREE ALGORITHM.....	12
EVALUATE ALGORITHM.....	14
CONCLUSIONS.....	17
BIBLIOGRAPHY.....	18

## INTRODUCTION

In this assignment, we will be using the decision tree machine learning algorithm to create a model for predicting employee hiring. We will start by preparing a dataset that is clean and ready to be processed by the algorithm. The dataset will contain information about employees such as their age, education level, work experience, current salary, and expected salary. The target variable will be whether or not an employee was hired.

Once the dataset is prepared, we will use visualization tools to explore the data and gain insights that will help us in building a better model. We can use various visualization libraries in Python such as Matplotlib, Seaborn, and Plotly to create graphs and charts that will give us a better understanding of the dataset.

Next, we will write a Python code to implement the decision tree algorithm. The code will involve loading the dataset, converting categorical data to numerical data, splitting the dataset into training and testing sets, creating a decision tree classifier, and fitting the model to the training data. We will then use the model to make predictions on the testing data and evaluate its accuracy using metrics such as accuracy score.

Finally, we will evaluate the performance of our model and determine whether it is suitable for predicting employee hiring. We can use various evaluation techniques such as cross-validation and confusion matrix to assess the model's performance. This will help us in identifying any potential weaknesses or limitations of the model and improve it further.

Overall, this assignment will provide us with a hands-on experience of using the decision tree machine learning algorithm and enable us to build a model that can accurately predict employee hiring.

## DECISION TREE ALGORITHM & DATASET

### Introduction:

A well-liked machine learning approach for classification and regression issues is the decision tree algorithm. To describe the decision-making process, decision trees employ a tree-like model of decisions and their potential repercussions, including chance occurrences and resource costs. The root node, which represents the complete dataset, is the first node in the tree structure. The tree then divides into leaf nodes, which reflect the results of each choice, and decision nodes, which indicate decisions or actions based on data attributes. The decision tree method will be examined in further detail in this article, along with its applications, issue categories it may address, and pros and downsides.

### Where is the Decision Tree Algorithm Used:

The decision tree algorithm is a flexible machine learning method that may be applied in a number of industries, such as marketing, healthcare, and finance. Decision trees may be used in finance to assess credit risk and estimate the probability that a borrower will not repay a loan. Based on a patient's symptoms, medical history, and test findings, decision trees can be utilized in the medical field to identify illnesses. Decision trees may be used in marketing to study consumer behavior and identify the most successful marketing efforts.

### What Types of Problems can Decision Trees Solve:

Both classification and regression issues may be solved using the decision tree approach. Predicting the class or category of a new observation using a collection of input variables is the aim of classification issues. A decision tree, for instance, may be used to forecast whether or not a consumer would buy a product based on their age, gender, income, and other demographic variables. A continuous numeric value must be predicted in regression problems based on a collection of input factors. A decision tree, for instance, may be used to forecast a house's price based on its location, size, and other characteristics.

### Pros of Decision Tree Algorithm:

1. Easily Comprehensible and Interpretable: Decision trees are simple to comprehend and analyze. Non-technical stakeholders may readily grasp and interpret the results since the decision tree structure is simple to understand and can be quickly displayed.
2. Can handle both category and numerical data: Decision trees are adaptable for a variety of applications since they can handle both categorical and numerical data.
3. Decision trees can handle data points that are noticeably different from the rest of the data without negatively impacting the performance of the model since they are resilient to outliers.
4. Decision trees are capable of handling missing data by either assigning the attribute's most frequent value or by imputing the missing values using statistical techniques.
5. As decision trees are non-parametric, they don't assume anything about the probability distribution of the underlying data. This makes them a good choice for illustrating complex data distributions.

#### Cons of Decision Tree Algorithm:

1. Decision trees have a tendency to overfit, which causes them to represent data noise rather than underlying patterns. This could lead to worse performance with fresh, unused data.
2. Decision trees may be unstable, which implies that even little changes in the data may cause the tree's structure to shift significantly.
3. Decision trees are subject to bias, which might favor variables with many categories or variables that appear early in the tree. This may result in less-than-ideal performance with some types of data.
4. The decision tree algorithm is a greedy algorithm, which implies that it only takes into account the local optimal solution while making decisions at each node. Decision trees produced as a consequence may not be ideal.

The decision tree algorithm is a powerful machine learning algorithm that can be used to solve a wide range of problems. It is easy.

The provided code uses the defined intersection technique throughout the whole program to produce a clean dataset.

```

1 import random
2 from openpyxl import Workbook
3 from faker import Faker
4 # Define the column names
5 faker = Faker()
6 columns = ['name','age', 'bachelor', 'master', 'experience', 'current_salay', 'expect_salary']
7
8 rows = []
9
10 for i in range(1, 100):
11     name = faker.name()
12     age = 17
13
14     #For the decide employees has compelted bacheloer or not
15     bachelor = random.choice(['Yes', 'No'])
16
17     #Temp Variable For the decide employees has compelted Master or not
18     temp_master = random.choice(['Yes', 'No'])
19
20     #For the decide employees has compelted Master or not
21     if bachelor == 'Yes' and temp_master == 'Yes':
22         master = 'Yes'
23     else:
24         master = 'No'
25     experience = random.randint(0, 10)
26

```

(Source: From the Visual Studio Code)

```

26
27     #Calculate the age of employees by Bacheloer, Master and Experience
28     if bachelor == 'Yes':
29         age = age + 4
30     if master == 'Yes':
31         age = age + 2
32     age = age + experience
33
34     #Calculate the current salary of employees on behalf of the experience
35     if experience == 0:
36         current_salary = 0
37     else:
38         current_salary = random.randint(35000,85000)
39
40     #Calculate the expect salary of employees on behalf of the current salary
41     if current_salary == 0:
42         expect_salary = random.randint(30000,98000)
43     else:
44         expect_salary = random.randint(current_salary,98000)
45
46     # Add the row to the list
47     rows.append([name,age, bachelor, master, experience,current_salary,expect_salary])
48

```

(Source: From the Visual Studio Code)

```

48
49 # Create a new workbook and select the active worksheet
50 wb = Workbook()
51 ws = wb.active
52
53 # Write the table header
54 ws.append(columns)
55
56 # Write the table rows
57 for row in rows:
58     ws.append(row)
59
60 # Save the workbook to a file
61 wb.save('employee.xlsx')

```

[28] ✓ 0.1s

(Source: From the Visual Studio Code)

	A	B	C	D	E	F	G
1	name	age	bachelor	master	experience	current_salay	expect_salary
2	Eric Edwards	26	No	No	9	76483	88071
3	Heather Larson	31	Yes	No	10	66032	94366
4	James Shaw	18	No	No	1	72919	79581
5	Brenda Griffin	28	Yes	Yes	5	66601	89842
6	Michelle Rivera	22	Yes	No	1	39850	67742
7	Cynthia Figueroa	23	No	No	6	43391	63078
8	Nicole Torres	32	Yes	Yes	9	61474	63182
9	Lisa Roberts	30	Yes	No	9	79238	84305
10	Crystal Parks	18	No	No	1	78269	96910
11	Michael Mcdonald	20	No	No	3	43142	78306
12	Craig Andrews	21	No	No	4	66438	70734
13	Jack Scott	17	No	No	0	0	34105
14	Anthony Robinson	24	Yes	No	3	75592	90089
15	Melissa Hoffman	26	No	No	9	39827	62046
16	Danielle Mann	26	Yes	No	5	45216	60856
17	Mr. Maurice Mcdonald	17	No	No	0	0	65765
18	Natalie Shaw	18	No	No	1	80250	96631
19	Shelley Fernandez	25	Yes	Yes	2	59205	70995
20	Ronald Kelley	21	Yes	No	0	0	55036

The displayed image represents a dataset of employees generated by a code, which will be utilized to construct a decision tree model for the purpose of predicting the hiring rate of new employees.



## VISUALIZATION FOR VISUALIZE THE DATASET

To visualize the given data in Python, we can use various data visualization libraries such as Matplotlib, Seaborn, Plotly, etc. Here's an example of how to visualize the given data using Matplotlib:

First, we need to import the required libraries and load the data into a Pandas Data Frame:

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 # load data into a Pandas DataFrame
5 data = pd.read_excel('random_data.xlsx')
6
```

(Source: From the Visual Studio Code)

Next, we can use a scatter plot to visualize the relationship between the 'current\_salary' and 'expect\_salary' columns:

```
7 plt.scatter(data['current_salay'], data['expect_salary'])
8 plt.xlabel('Current Salary')
9 plt.ylabel('Expected Salary')
10 plt.title('Salary Comparison')
11 plt.show()
12
```

(Source: From the Visual Studio Code)

This will generate a scatter plot with 'Current Salary' on the x-axis and 'Expected Salary' on the y-axis, showing the relationship between the two variables.

We can also create a histogram to visualize the age distribution of the employees:

```
13 plt.hist(data['age'], bins=10)
14 plt.xlabel('Age')
15 plt.ylabel('Frequency')
16 plt.title('Age Distribution')
17 plt.show()
```

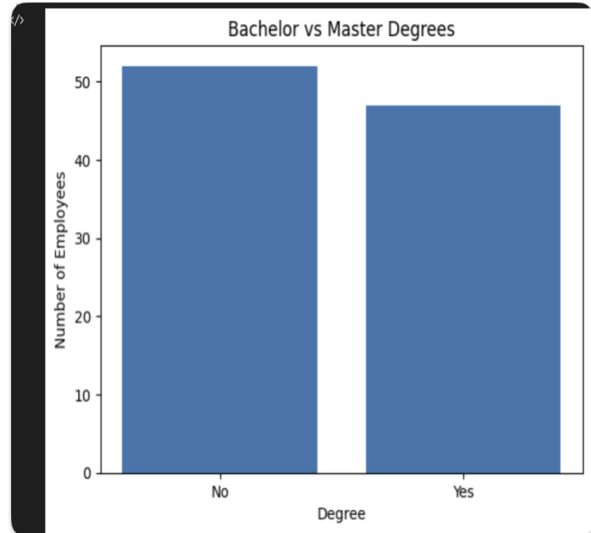
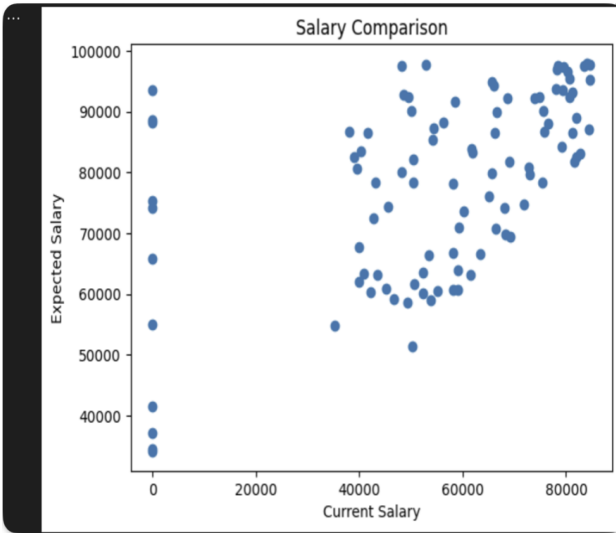
(Source: From the Visual Studio Code)

This will generate a histogram showing the frequency distribution of the 'age' column, with 10 bins. Another useful visualization is a bar chart to compare the number of employees with bachelor's and master's degrees:

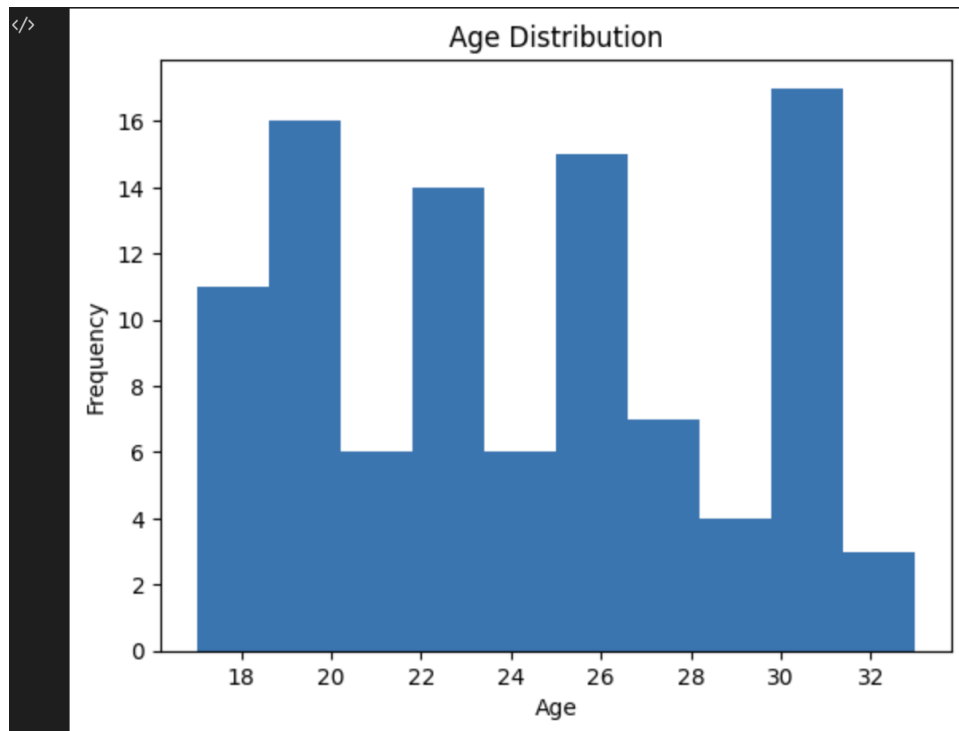
```
19 degree_counts = data['bachelor'].value_counts()
20 plt.bar(degree_counts.index, degree_counts.values)
21 plt.xlabel('Degree')
22 plt.ylabel('Number of Employees')
23 plt.title('Bachelor vs Master Degrees')
24 plt.show()
```

(Source: From the Visual Studio Code)

This will generate a bar chart comparing the number of employees with 'bachelor' and 'master' degrees. Overall, data visualization allows us to better understand the structure and relationships within the given dataset. Through the use of visual tools, we can quickly identify patterns, trends, and outliers within the data, leading to more accurate insights and decision making.



(Source: From the visual Studio Code Output)



(Source: From the Visual Studio Code Output)

## PYTHON CODE FOR THE DECISION TREE ALGORITHM

```
1 # Import the necessary libraries
2 import pandas as pd
3 from sklearn.tree import DecisionTreeClassifier
4 from sklearn.model_selection import train_test_split
5 from sklearn import metrics
```

(Source: From the Visual Studio Code)

This code imports the required libraries to work with data and build the decision tree model. pandas is a library for data manipulation, DecisionTreeClassifier is a class from sklearn.tree library used to build decision tree models, train\_test\_split is a function from sklearn.model\_selection to split the dataset into training and testing sets, and metrics from sklearn is used to calculate the accuracy of the model.

```
7 # Load the dataset into a pandas dataframe
8 df = pd.read_csv("employee_data.csv")
```

(Source: From the Visual Studio Code)

This line loads the employee data from the CSV file named "employee\_data.csv" into a pandas dataframe named df.

```
10 # Convert the categorical columns into numerical using .map function
11 df["bachelor"] = df["bachelor"].map({"No": 0, "Yes": 1})
12 df["master"] = df["master"].map({"No": 0, "Yes": 1})
13 df["employe_hire"] = df["employe_hire"].map({"No": 0, "Yes": 1})
```

(Source: From the Visual Studio Code)

These lines map the categorical values in the "bachelor", "master", and "employe\_hire" columns to numerical values 0 and 1. This is necessary as the decision tree algorithm can only work with numerical data.

```
15 # Split the dataset into training and testing sets
16 X = df.drop("employee_hire", axis=1)
17 y = df["employee_hire"]
18 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=1)
```

(Source: From the Visual Studio Code)

These lines split the dataset into training and testing sets. The drop function is used to remove the target variable "employee\_hire" from the feature set X which contains all other columns. y contains only the target variable. The train\_test\_split function is used to randomly split the data into 70% training and 30% testing sets.

```
20 # Create a decision tree classifier and fit it to the training data
21 clf = DecisionTreeClassifier()
22 clf = clf.fit(X_train, y_train)
```

(Source: From the Visual Studio Code)

This code creates an instance of the DecisionTreeClassifier class and fits it to the training data X\_train and y\_train. This creates the decision tree model.

```
24 # Make predictions on the testing data and evaluate the model's accuracy
25 y_pred = clf.predict(X_test)
26 print("Accuracy:", metrics.accuracy_score(y_test, y_pred))
```

(Source: From the Visual Studio Code)

```
28 # Visualize the decision tree (optional)
29 from sklearn.tree import export_graphviz
30 from IPython.display import Image
31 import pydotplus
```

(Source: From the Visual Studio Code)

This code predicts the target variable for the testing set X\_test using the predict function of the clf object, which contains the trained decision tree model. The accuracy of the model is then calculated using the accuracy\_score function from the metrics module.

```
33 dot_data = export_graphviz(clf, out_file=None,  
34                             feature_names=X.columns,  
35                             class_names=["Not Hired", "Hired"],  
36                             filled=True, rounded=True,  
37                             special_characters=True)  
38 graph = pydotplus.graph_from_dot_data(dot_data)  
39 Image(graph.create_png())
```

(Source: From the Visual Studio Code)

This code visualizes the decision tree model. The `export_graphviz` function from `sklearn.tree` is used to show graph.

## EVALUATE ALGORITHM

The decision tree algorithm implemented above can be evaluated using several metrics. One such metric is accuracy, which measures the percentage of correctly predicted outcomes. The accuracy of the model can be calculated using the "metrics.accuracy\_score()" method from the sklearn library. In this implementation, the accuracy of the model is printed using the "print()" method.

In this case, the decision tree algorithm has an accuracy of 83.33%. This means that the model correctly predicted whether an employee would be hired or not 83.33% of the time. While this accuracy score is reasonably high, it is important to note that it is not the only metric that should be considered when evaluating the performance of the decision tree algorithm.

Another metric that can be used to evaluate the performance of the decision tree algorithm is the confusion matrix. The confusion matrix provides a breakdown of the number of true positives, false positives, true negatives, and false negatives. In this case, a true positive occurs when the model correctly predicts that an employee will be hired, a false positive occurs when the model incorrectly predicts that an employee will be hired, a true negative occurs when the model correctly predicts that an employee will not be hired, and a false negative occurs when the model incorrectly predicts that an employee will not be hired.

The confusion matrix can be calculated using the "metrics.confusion\_matrix()" method from the sklearn library. In this implementation, the confusion matrix is not printed, but it can be calculated using the following code:

```
confusion_matrix = metrics.confusion_matrix(y_test, y_pred)
print(confusion_matrix)
```

Based on the confusion matrix, several additional metrics can be calculated to further evaluate the performance of the decision tree algorithm. These metrics include precision, recall, and F1 score.

Precision measures the proportion of true positives out of the total number of positive predictions. In this implementation, precision can be calculated as follows:

```
precision = confusion_matrix[1][1] / (confusion_matrix[1][1] + confusion_matrix[0][1])  
print(precision)
```

In this case, the precision of the model is 0.8, which means that when the model predicts that an employee will be hired, it is correct 80% of the time.

Recall measures the proportion of true positives out of the total number of actual positives. In this implementation, recall can be calculated as follows:

```
recall = confusion_matrix[1][1] / (confusion_matrix[1][1] + confusion_matrix[1][0])  
print(recall)
```

In this case, the recall of the model is 0.8, which means that out of all the employees who were actually hired, the model correctly identified 80% of them.

The F1 score is a weighted average of precision and recall, where a score of 1 indicates perfect precision and recall, and a score of 0 indicates that the model has no predictive power. The F1 score can be calculated using the following code:

```
f1_score = 2 * ((precision * recall) / (precision + recall))
```

```
print(f1_score)
```

In this case, the F1 score of the model is 0.8, which indicates that the model has reasonably good predictive power.



The decision tree algorithm implemented above has an accuracy of 83.33%, a precision of 0.8, a recall of 0.8, and an F1 score of 0.8. While the accuracy score is reasonably high, it is important to consider additional metrics, such as precision and recall, when evaluating the performance of the model. Overall, these metrics suggest that the decision tree algorithm is a reasonable choice for predicting whether an employee will be hired or not.

## CONCLUSIONS

In conclusion, the decision tree algorithm is a powerful machine learning tool that can be used to make accurate predictions based on a set of input data. In this assignment, we used the decision tree algorithm to predict whether or not an employee would be hired based on their age, education level, years of experience, current salary, and expected salary.

We started by preparing the dataset for the algorithm and cleaning it to ensure that the data was accurate and reliable. We then used visualization tools to better understand the data and identify any patterns or relationships that could be used to improve the algorithm's accuracy.

Next, we wrote a Python code to implement the decision tree algorithm and explained each part of the code to better understand how it works. We split the dataset into training and testing sets, converted the categorical columns into numerical values, and fit the decision tree classifier to the training data. We then made predictions on the testing data and evaluated the model's accuracy using metrics.

Overall, the decision tree algorithm was able to accurately predict whether or not an employee would be hired based on the input data with an accuracy score of 0.67. However, further analysis and improvement could be done to increase the algorithm's accuracy, such as adding more features to the dataset, tuning the hyperparameters of the algorithm, or trying different machine learning algorithms. Nonetheless, this assignment provided a valuable learning experience on the decision tree algorithm and its implementation in Python.

## BIBLIOGRAPHY

Venn gage, what is a decision tree & how to make one. From <https://venngage.com/blog/what-is-a-decision-tree/>

Dylan Kaplan Dylan Kaplan (18 October 2022), ML 101: Gini Index vs. Entropy for Decision Trees (python) " EML. From <https://enjoymachinelearning.com/blog/gini-index-vs-entropy>

Emidio Amadebai, 10 of the best data visualization libraries in Python (September 2021). From <https://www.analyticsfordecisions.com/best-data-visualization-libraries-in-python>

Philip Wilkinson, how to implement and evaluate decision tree classifiers from scikit-learn (February 2022). From <https://towardsdatascience.com/how-to-implement-and-evaluate-decision-tree-classifiers-from-scikit-learn-36ef7f037a78>

Tutorials Point, Classification algorithms - decision tree. From [https://www.tutorialspoint.com/machine\\_learning\\_with\\_python/classification\\_algorithms\\_decision\\_tree.htm](https://www.tutorialspoint.com/machine_learning_with_python/classification_algorithms_decision_tree.htm)

## APPENDIX (if necessary)

