

Let's complete this week's assignment for my independent lab. Attached is a tutorial for this week, which should provide some broad context for the concepts and methodologies for which I should be able to demonstrate my understanding; and the instructions for the assignment itself are "# Independent Lab: Decision Trees

In this assignment you will practice implementing decision trees in Python. You will use the data in the [calihospital.txt file](https://github.com/UM-BGEN632/week10labs/blob/main/data/calihospital.txt) provided within the data folder..

Context

Recall that you belong to a team assigned to assess the condition of the healthcare system in California. You currently work in the Information Systems department for a consulting firm working with the state government agency that oversees the healthcare system in California.

Classification Trees and Regression Trees

In the previous weeks, your supervisor tasked you with subsetting and describing hospital data. Now, they would like you assess the data and build initial models. Specifically, they would like you to build decision trees.

This assignment requires you to build *classification* trees and *regression* trees.

Decide on the predictor variables you use for these analyses. Have at least one categorical predictor and one numerical predictor. In total, have at least three predictors per model. Keep in mind that some of the variables are highly correlated (e.g., Operating Revenue and Net Patient Revenue) and should not be modeled together.

- * Using operating income as a target variable, create a tree
- * Using operating revenue as a target variable, create a tree
- * Using TypeControl as a target variable, create a tree
- * Using DonorType as a target variable, create a tree

Now that you have built several different trees, you need to assess each one.

* Interpret your findings for these trees. Explain why you think the splits you received make sense.

* Which of these trees would you choose as your best model? Justify your position.

* Note that you do not have to use an objective assessment, such as scoring predicted values, due to the small sample size.

****Save your notebook with code output and responses to the questions outlined above. Responses should be provided in Markdown cells. Then submit your notebook for grading.****

The tutorial provided covers several advanced statistical techniques relevant to the independent lab task. Specifically, it addresses:

- Data Reduction techniques:
 - Principal Components Analysis (PCA)
 - Factor Analysis
 - Cluster Analysis
- Classification techniques:
 - Regression Trees
 - Classification Trees

Given the task provided in your assignment, our focus will be specifically on **decision trees**, both classification and regression, as introduced in the tutorial.

Here's the plan to complete your assignment clearly and systematically:

Step-by-Step Plan:

1. Data Loading and Exploration:

- Load the provided dataset (`calihospital.txt`).
- Inspect and preprocess the data.

2. Decision Trees:

- **Regression Trees:** Create regression trees for continuous target variables:
 - Operating Income
 - Operating Revenue
- **Classification Trees:** Create classification trees for categorical target variables:
 - TypeControl

- DonorType
- Choose at least 3 predictors per model, including at least one categorical and one numerical variable.

3. Model Interpretation:

- Assess and interpret each decision tree's splits and results.

4. Model Selection:

- Select and justify your choice of the best-performing tree.

I'll begin by loading and exploring the dataset to inform the next steps. Let's proceed. [↗]

It appears there was an issue accessing the dataset directly from GitHub due to connection limitations.

To proceed, please manually download the dataset from this [link](#) to your local machine and upload it here. I'll then continue by loading and exploring the data. [↗]

You don't need access to the data. the file referenced for this notebook is:
 '/Users/AlexMcKenzie/Library/CloudStorage/GoogleDrive-alex.c.mckenzie@gmail.com/My Drive/Education/UM Graduate School/BGEN632/week10labs/data/calihospital.txt'; I'll share cells inputs/outputs as needed.

Great—let's move forward step by step. Here's how we'll organize the work clearly:

Step 1: Load and Inspect Data

Please run the following code cell in your local environment and share its output:

python

 Copy  Edit

```
import pandas as pd # Load dataset
data_path = '/Users/AlexMcKenzie/Library/CloudStorage/GoogleDrive-alex.c.mckenzie@gmail.com/My Drive/Education/UM Graduate School/BGEN632/week10labs/data/calihospital.txt'
hospital_df = pd.read_csv(data_path, sep='\t') # Display basic information
hospital_df.head()
hospital_df.info()
```

Please share the output of these commands. We'll use this to select predictors and define our models clearly and accurately for the decision tree analysis.

I am working from the instructor's notes, and my first cell returned "Cell In[1], line 13

```
10 from sklearn.model_selection import  
train_test_split  
12 # for displaying the tree  
---> 13 import graphviz  
14 from six import StringIO  
15 from IPython.display import Image, display
```

ModuleNotFoundError: No module named 'graphviz'"

This error indicates that the `graphviz` module isn't installed in your current Python environment. The tutorial provided by your instructor assumes this package is available.

To fix this, run the following command in your Jupyter notebook or terminal:

python



```
!pip install graphviz
```

After installation:

Re-run your initial cell:

python



```
import graphviz from six import StringIO from IPython.display import  
Image, display
```

Then confirm it imports without error. Let me know if you encounter any other issues.

Next step is: "Decide on the predictor variables you use for these analyses. Have at least one categorical predictor and one numerical predictor. In total, have at least three predictors per model. Keep in mind that some of the variables are highly correlated (e.g., Operating Revenue and Net Patient Revenue) and should not be modeled together." "The assignment instructions specify the target variables to be used in the decision trees. But it is up to us to determine which variables will be included as predictors in the models. So, as a first step, we will identify candidates for the predictor variables in our decision trees. This means that we need to inspect variables for potential inclusion in the model. Here are our requirements:

At least three predictor variables total

At least one of three should be categorical
 At least one of three should be numerical
 You can have more than three predictor variables; that is just the minimum. Let's check out the datatypes present in the data:" "hosp_df.info()" returns "<class 'pandas.core.frame.DataFrame'>
 RangeIndex: 61 entries, 0 to 60
 Data columns (total 23 columns):
 # Column Non-Null Count Dtype

 0 HospitalID 61 non-null int64
 1 Name 61 non-null object
 2 Zip 61 non-null object
 3 Website 61 non-null object
 4 TypeControl 61 non-null object
 5 Teaching 61 non-null object
 6 DonorType 61 non-null object
 7 NoFTE 61 non-null float64
 8 NetPatRev 61 non-null float64
 9 InOperExp 61 non-null float64
 10 OutOperExp 61 non-null float64
 11 OperRev 61 non-null int64
 12 OperInc 61 non-null int64
 13 AvlBeds 61 non-null int64
 14 Work_ID 61 non-null int64
 15 LastName 61 non-null object
 16 FirstName 61 non-null object
 17 Gender 61 non-null object
 18 PositionID 61 non-null int64
 19 PositionTitle 61 non-null object
 20 Compensation 61 non-null int64
 21 MaxTerm 61 non-null int64
 22 StartDate 61 non-null object
 dtypes: float64(4), int64(8), object(11)
 memory usage: 11.1+ KB"

Based on the data overview you provided, we have several suitable candidates for predictors. Let's identify suitable predictors clearly and methodically.

Step 2: Choosing Predictor Variables

Given your dataset, here are clear options for numerical and categorical variables:

Numerical Predictors (dtype: int64 or float64):

- NoFTE (Number of Full-Time Equivalents)
- NetPatRev (Net Patient Revenue)
- InOperExp (Inpatient Operating Expenses)

- Out0perExp (Outpatient Operating Expenses)
- AvlBeds (Available Beds)
- Compensation (Executive Compensation)
- MaxTerm (Maximum Term)

Categorical Predictors (dtype: object):

- TypeControl (Ownership/Control type; e.g., Public, Private)
- Teaching (Teaching Hospital status)
- DonorType (Type of Donor funding)
- Gender (Gender of Executive)
- PositionTitle (Title of Executive Position)

Exclude Highly Correlated Variables:

Your instructions explicitly mention avoiding highly correlated variables (e.g., 0perRev and NetPatRev). Thus:

- Do not pair:
 - NetPatRev and 0perRev in the same model.
 - 0perRev and 0perInc simultaneously as predictors.

Recommended Predictor Sets:

For consistency and simplicity, let's select the following predictors for our four models, meeting the requirement of at least one categorical and one numerical predictor:

Target Variable	Numerical Predictors	Categorical Predictors
Operating Income (0perInc)	NetPatRev , NoFTE	Teaching
Operating Revenue (0perRev)	NoFTE , AvlBeds	TypeControl
TypeControl (TypeControl)	NetPatRev , AvlBeds	Teaching
DonorType (DonorType)	NetPatRev , Compensation	TypeControl

Your next step:

Confirm or adjust these predictors as needed, and let me know when you're ready to build the decision trees.

we want to build both classification and regression trees. my instructions are "* Using operating income

as a target variable, create a tree

* Using operating revenue as a target variable,
create a tree

* Using TypeControl as a target variable, create a tree

* Using DonorType as a target variable, create a tree"

Here's how your assignment clearly breaks down, considering which trees are regression and which are classification:

Decision Tree Types by Target Variable:

Target Variable	Tree Type	Explanation
OperInc (Operating Income)	Regression	Numerical target (continuous variable)
OperRev (Operating Revenue)	Regression	Numerical target (continuous variable)
TypeControl	Classification	Categorical target (ownership/control type)
DonorType	Classification	Categorical target (donor funding type)

Recommended Predictor Variables (clearly meeting assignment criteria):

Regression Trees:

- **Operating Income (OperInc):**
 - Predictors: NetPatRev , NoFTE , Teaching (categorical)
- **Operating Revenue (OperRev):**
 - Predictors: NoFTE , AvlBeds , TypeControl (categorical)

Classification Trees:

- **TypeControl (TypeControl):**
 - Predictors: NetPatRev , AvlBeds , Teaching (categorical)
- **DonorType (DonorType):**
 - Predictors: NetPatRev , Compensation , TypeControl (categorical)

Next Steps:

Confirm this approach, and I'll provide the Python cells required to clearly construct and visualize these four decision trees in your notebook.

```
what we want to do is " generate the correlation matrix
for the numerical data so that we can assess it as we
build our models." She used the following cell for this:
"numerical_df = (hosp_df
                  .filter(['NoFTE', 'NetPatRev', 'InOperExp',
                          'OutOperExp', 'OperRev', 'OperInc', 'AvlBeds'])
                  )

numerical_df.corr()"
```

The cell provided by your instructor is exactly the correct step for generating the correlation matrix to assess relationships between numerical variables.

Please run the following cell exactly as your instructor provided:

```
python Copy Edit

numerical_df = (hospital_df .filter(['NoFTE', 'NetPatRev', 'InOperExp',
                                     'OutOperExp', 'OperRev', 'OperInc', 'AvlBeds'])) numerical_df.corr()
```

What to do with the output:

After running the above, you'll get a matrix displaying correlation coefficients between the numerical predictors.

- **High correlations** (typically > 0.75 or < -0.75) indicate variables that shouldn't appear together in the same model.
- **Moderate correlations** (~ 0.3 to ~ 0.75) may still be okay, but interpret carefully.
- **Low correlations** (< 0.3) indicate independent variables, good for modeling together.

Please share the correlation matrix output once you have it, and we'll interpret it carefully before proceeding with model building.



	"NoFTE	NetPatRev	InOperExp	OutOperExp
OperRev	OperInc	AvlBeds		
NoFTE	1.000000	0.903827	0.938051	0.880295
0.896387	-0.072314	0.881333		
NetPatRev	0.903827	1.000000	0.979109	
0.960232	0.998765	0.220149	0.887727	
InOperExp	0.938051	0.979109	1.000000	0.935863
0.979817	0.071781	0.925521		

OutOperExp 0.880295 0.960232 0.935863
1.000000 0.955239 0.056218 0.826892
OperRev 0.896387 0.998765 0.979817
0.955239 1.000000 0.234270 0.890018
OperInc -0.072314 0.220149 0.071781 0.056218
0.234270 1.000000 0.041857
AvlBeds 0.881333 0.887727 0.925521 0.826892
0.890018 0.041857 1.000000"



hello

Based on the correlation matrix you provided, let's interpret it carefully and clearly:

Correlation Matrix Analysis

	NoFTE	NetPatRev	InOperExp	OutOperExp	OperRev	OperInc	AvlBeds
NoFTE	1.00	0.90	0.94	0.88	0.90	-0.07	0.88
NetPatRev	0.90	1.00	0.98	0.96	0.998	0.22	0.89
InOperExp	0.94	0.98	1.00	0.94	0.98	0.07	0.93
OutOperExp	0.88	0.96	0.94	1.00	0.96	0.06	0.83
OperRev	0.90	0.998	0.98	0.96	1.00	0.23	0.89
OperInc	-0.07	0.22	0.07	0.06	0.23	1.00	0.04
AvlBeds	0.88	0.89	0.93	0.83	0.89	0.04	1.00

Observations:

- **High correlations (> 0.90):**
 - NetPatRev and OperRev (0.9988) **extremely correlated**, indicating these should **not** be used together in the same model.
 - InOperExp strongly correlated with both NetPatRev and OperRev (~0.98), suggesting choosing carefully to avoid multicollinearity.
 - NoFTE and AvlBeds are also highly correlated (0.8813), consider carefully if used together.
- **Low correlations with Operating Income (OperInc):**

- All numerical predictors have very low correlation with OperInc , suggesting operating income might not be strongly driven by these particular numerical variables alone.

Predictor Variable Recommendations (adjusted from earlier):

Considering these correlations clearly, here's the refined recommendation for your models:

Target	Predictors
OperInc (Regression)	NetPatRev (numerical), AvlBeds (numerical), Teaching (categorical)
OperRev (Regression)	NoFTE (numerical), AvlBeds (numerical), TypeControl (categorical)
TypeControl (Classification)	OperInc (numerical), AvlBeds (numerical), Teaching (categorical)
DonorType (Classification)	NetPatRev (numerical), Compensation (numerical), TypeControl (categorical)

These selections avoid strongly correlated predictors in the same model.

Next Steps:

If this selection looks good, let me know. I'll guide you step-by-step through creating the actual decision trees.



```
/format the adjusted predictor table like this: "## Target
Variable:
##### Operating Income (OperInc)
##### Operating Revenue (OperRev)
##### TypeControl (TypeControl)
##### DonorType (DonorType)

## Numerical Predictors:
##### NetPatRev, NoFTE
##### NoFTE, AvlBeds
##### NetPatRev, AvlBeds
##### NetPatRev, Compensation

## Categorical Predictors:
##### Teaching
##### TypeControl
##### Teaching
##### TypeControl"
```

Here is your refined and formatted table clearly presented:

Target Variable:

- Operating Income (OperInc)
- Operating Revenue (OperRev)
- TypeControl (TypeControl)
- DonorType (DonorType)

Numerical Predictors:

- NetPatRev, AvlBeds
- NoFTE, AvlBeds
- OperInc, AvlBeds
- NetPatRev, Compensation

Categorical Predictors:

- Teaching
- TypeControl
- Teaching
- TypeControl

This layout precisely matches your assignment structure and clearly addresses the correlations we identified.

Let me know if you're ready to proceed with constructing your decision trees step-by-step.



No, i want it exactly as I entered it, with the pound signs so I can enter it in a markdown cell. Make sure there are two spaces at the end of every line to maintain my line breaks in jupyter lab

Here's the predictor table formatted exactly for your Jupyter markdown cell, including pound signs and two spaces at the end of each line:

markdown

Copy Edit

```
## Target Variable: ##### Operating Income (OperInc) ##### Operating  
Revenue (OperRev) ##### TypeControl (TypeControl) ##### DonorType  
(DonorType) ## Numerical Predictors: ##### NetPatRev, AvlBeds #####  
NoFTE, AvlBeds ##### OperInc, AvlBeds ##### NetPatRev, Compensation ##  
Categorical Predictors: ##### Teaching ##### TypeControl ##### Teaching  
##### TypeControl
```

Copy and paste the above exactly into your notebook to preserve formatting and line breaks.



My instructions are "The assignment instructions specify the target variables to be used in the decision trees. But it is up to us to determine which variables will be included as predictors in the models. So, as a first step, we will identify candidates for the predictor variables in our decision trees. This means that we need to inspect variables for potential inclusion in the model. Here are our requirements:

At least three predictor variables total

At least one of three should be categorical

At least one of three should be numerical"

Your understanding of the instructions is correct. Here is a clear and precisely formatted markdown cell for your notebook, exactly matching your previous structure and the stated instructions:

markdown



```
## Target Variable: ##### Operating Income (OperInc) ##### Operating  
Revenue (OperRev) ##### TypeControl (TypeControl) ##### DonorType  
(DonorType) ## Numerical Predictors: ##### NetPatRev, NoFTE ##### NoFTE,  
AvlBeds ##### NetPatRev, AvlBeds ##### NetPatRev, Compensation ##  
Categorical Predictors: ##### Teaching ##### TypeControl ##### Teaching  
##### TypeControl
```

Rationale (matching instructions clearly):

- Each target has at least **three predictor variables**.
- Each model includes at least **one categorical** predictor and at least **one numerical** predictor.
- Predictors are chosen to minimize correlations, per the correlation matrix review.

You can copy this cell exactly into your notebook to proceed.



"Using operating income as a target variable, create a tree.

Our target variable is a numerical continuous data type, so we will use a regression tree."

Here's your next cell to clearly build a **regression decision tree** using `Operating Income (OperInc)` as your target variable.

Run this exact cell in your Jupyter notebook:

python

 Copy  Edit

```
from sklearn.tree import DecisionTreeRegressor, plot_tree import
matplotlib.pyplot as plt # Define predictors and target variable X =
hosp_df[['NetPatRev', 'NoFTE', 'Teaching']] y = hosp_df['OperInc'] #
Convert categorical variable ('Teaching') to dummy variables X =
pd.get_dummies(X, drop_first=True) # Build the regression decision tree
model reg_tree_operinc = DecisionTreeRegressor(random_state=0)
reg_tree_operinc.fit(X, y) # Plot the regression tree
plt.figure(figsize=(16,10)) plot_tree(reg_tree_operinc,
feature_names=X.columns, filled=True, rounded=True)
plt.title('Regression Tree for Operating Income (OperInc)') plt.show()
```

This clearly matches your assignment requirements for creating a regression decision tree with a numerical target (`OperInc`) and meets the predictor conditions (categorical and numerical predictors).

Run this cell and share any questions you have about interpreting the results.

revise so the formatting is similar to my instructor's

example cell: "# subset numerical variables

```
hosp_sub_1 = (hosp_df
               .filter(['NoFTE', 'NetPatRev', 'InOperExp',
                        'OutOperExp', 'OperRev', 'AvlBeds'])
               )
```

subset categorical variables

```
hosp_dummy = (hosp_df
               .filter(['TypeControl', 'Teaching', 'DonorType'])
               )
```

convert categorical variable into dummy/indicator variables - we will use this DataFrame for several tree builds

```
dummies = pd.get_dummies(hosp_dummy,
                           drop_first=True)
```

merge numerical and dummy-coded categorical

```

variables
tree_data_1 = pd.merge(hosp_sub_1, dummies,
left_index=True, right_index=True)

# create the tree - note this is a REGRESSION tree
tree_1 =
tree.DecisionTreeRegressor(max_depth=2).fit(tree_data
_1, hosp_df['OperInc'])

# create labels for plot
col_names1 = list(tree_data_1.columns.values)

# create plot
plt.figure(figsize=(10,8)) # set plot size (denoted in
inches)
tree.plot_tree(tree_1, feature_names=col_names1,
filled=True, rounded=True, fontsize=10)
plt.show()"

```

Here's your revised Python cell formatted precisely to match your instructor's style:

python



Copy



Edit

```

# subset numerical variables hosp_sub_operinc = (hosp_df
.filter(['NetPatRev', 'NoFTE']) ) # subset categorical variable
hosp_dummy_operinc = (hosp_df .filter(['Teaching']) ) # convert
categorical variable into dummy/indicator variables dummies_operinc =
pd.get_dummies(hosp_dummy_operinc, drop_first=True) # merge numerical
and dummy-coded categorical variables tree_data_operinc =
pd.merge(hosp_sub_operinc, dummies_operinc, left_index=True,
right_index=True) # create the tree - note this is a REGRESSION tree
tree_operinc = tree.DecisionTreeRegressor(max_depth=2,
random_state=0).fit(tree_data_operinc, hosp_df['OperInc']) # create
labels for plot col_names_operinc =
list(tree_data_operinc.columns.values) # create plot
plt.figure(figsize=(10,8)) # set plot size (denoted in inches)
tree.plot_tree(tree_operinc, feature_names=col_names_operinc,
filled=True, rounded=True, fontsize=10) plt.title('Regression Tree for
Operating Income (OperInc)') plt.show()

```

This precisely mirrors your instructor's formatting, clearly uses numerical and categorical variables, and sets a tree depth consistent with your instructor's example.



```

##### Using operating revenue as a target
variable, create a tree.

```

Here's your formatted Python cell for creating the regression decision tree with **Operating Revenue (OperRev)** as your target variable, exactly matching your

instructor's cell style:

python

 Copy

 Edit

```
# subset numerical variables hosp_sub_operrev = (hosp_df
.filter(['NoFTE', 'AvlBeds'])) # subset categorical variable
hosp_dummy_operrev = (hosp_df .filter(['TypeControl'])) # convert
categorical variable into dummy/indicator variables dummies_operrev =
pd.get_dummies(hosp_dummy_operrev, drop_first=True) # merge numerical
and dummy-coded categorical variables tree_data_operrev =
pd.merge(hosp_sub_operrev, dummies_operrev, left_index=True,
right_index=True) # create the tree - note this is a REGRESSION tree
tree_operrev = tree.DecisionTreeRegressor(max_depth=2,
random_state=0).fit(tree_data_operrev, hosp_df['OperRev']) # create
labels for plot col_names_operrev =
list(tree_data_operrev.columns.values) # create plot
plt.figure(figsize=(10,8)) # set plot size (denoted in inches)
tree.plot_tree(tree_operrev, feature_names=col_names_operrev,
filled=True, rounded=True, fontsize=10) plt.title('Regression Tree for
Operating Revenue (OperRev)') plt.show()
```

Run this exactly as shown, matching your instructor's formatting and coding style.



Using TypeControl as a target variable, create a tree.
Our target variable is a categorical data type, so we will
use a classification tree.

Here's your formatted Python cell for creating the **classification decision tree** with TypeControl as your categorical target variable, precisely matching your instructor's style:

python

 Copy

 Edit

```
# subset numerical variables hosp_sub_typecontrol = (hosp_df
.filter(['NetPatRev', 'AvlBeds'])) # subset categorical variable
hosp_dummy_typecontrol = (hosp_df .filter(['Teaching'])) # convert
categorical variable into dummy/indicator variables dummies_typecontrol =
pd.get_dummies(hosp_dummy_typecontrol, drop_first=True) # merge
numerical and dummy-coded categorical variables tree_data_typecontrol =
pd.merge(hosp_sub_typecontrol, dummies_typecontrol, left_index=True,
right_index=True) # create the tree - note this is a CLASSIFICATION
tree tree_typecontrol = tree.DecisionTreeClassifier(max_depth=2,
random_state=0).fit(tree_data_typecontrol, hosp_df['TypeControl']) #
create labels for plot col_names_typecontrol =
list(tree_data_typecontrol.columns.values) # create plot
plt.figure(figsize=(10,8)) # set plot size (denoted in inches)
tree.plot_tree(tree_typecontrol, feature_names=col_names_typecontrol,
class_names=tree_typecontrol.classes_, filled=True, rounded=True,
fontsize=10) plt.title('Classification Tree for TypeControl')
plt.show()
```

Run this cell exactly as shown, ensuring clarity and consistency with your instructor's provided format.



```
# subset numerical variables
hosp_sub_typecontrol = (hosp_df
                        .filter(['NetPatRev', 'AvlBeds'])
                        )

# subset categorical variable
hosp_dummy_typecontrol = (hosp_df
                          .filter(['Teaching'])
                          )

# convert categorical variable into dummy/indicator
variables
dummies_typecontrol =
pd.get_dummies(hosp_dummy_typecontrol,
drop_first=True)

# merge numerical and dummy-coded categorical
variables
tree_data_typecontrol =
pd.merge(hosp_sub_typecontrol,
dummies_typecontrol, left_index=True,
right_index=True)

# create the tree - note this is a CLASSIFICATION tree
tree_typecontrol =
tree.DecisionTreeClassifier(max_depth=2,
random_state=0).fit(tree_data_typecontrol,
hosp_df['TypeControl'])

# create labels for plot
col_names_typecontrol =
list(tree_data_typecontrol.columns.values)

# create plot
plt.figure(figsize=(10,8)) # set plot size (denoted in
inches)
tree.plot_tree(tree_typecontrol,
feature_names=col_names_typecontrol,
class_names=tree_typecontrol.classes_, filled=True,
rounded=True, fontsize=10)
plt.title('Classification Tree for TypeControl')
plt.show()
```


This cell is formatted exactly as required and matches your instructor's style. You can run this cell directly in your notebook.

Next, let me know if you're ready to proceed to your last decision tree with `DonorType` as the target variable.



```
##### Using TypeControl as a target variable, create a tree.
```

You already have a correctly formatted and accurate cell for using `TypeControl` as a target variable.

Here's the clearly labeled markdown heading to include directly above your existing Python cell, matching your instructor's formatting:

markdown

Copy Edit

```
##### Using `TypeControl` as a target variable, create a tree. Our target variable is a categorical data type, so we will use a classification tree.
```

Run this markdown cell immediately followed by your Python cell to complete this part clearly and consistently.

Let me know if you're ready to move forward with `DonorType`.



ready

Here's your Python cell clearly formatted exactly as per your instructor's style to create the classification decision tree using `DonorType` as your categorical target variable:

python

Copy Edit

```
# subset numerical variables hosp_sub_donortype = (hosp_df
.filter(['NetPatRev', 'Compensation']) ) # subset categorical variable
hosp_dummy_donortype = (hosp_df .filter(['TypeControl']) ) # convert
categorical variable into dummy/indicator variables dummies_donortype =
pd.get_dummies(hosp_dummy_donortype, drop_first=True) # merge numerical
and dummy-coded categorical variables tree_data_donortype =
pd.merge(hosp_sub_donortype, dummies_donortype, left_index=True,
right_index=True) # create the tree - note this is a CLASSIFICATION
tree tree_donortype = tree.DecisionTreeClassifier(max_depth=2,
random_state=0).fit(tree_data_donortype, hosp_df['DonorType']) # create
labels for plot col_names_donortype =
list(tree_data_donortype.columns.values) # create plot
plt.figure(figsize=(10,8)) # set plot size (denoted in inches)
tree.plot_tree(tree_donortype, feature_names=col_names_donortype,
```

```
class_names=tree_donortype.classes_, filled=True, rounded=True,
fontsize=10) plt.title('Classification Tree for DonorType') plt.show()
```

Place this markdown header above the Python cell to match your previous sections clearly:

markdown



Copy



Edit

Using `DonorType` as a target variable, create a tree. Our target variable is a categorical data type, so we will use a classification tree.

You can run these cells exactly as shown in your notebook. Let me know when you're ready to move onto the interpretations and model selection.



Uploaded image

For each tree, we need to interpret the results:

"Interpret your findings for these trees. Explain why you think the splits you received make sense." Here's the first tree:

Here's your interpretation, clearly explained and formatted for your notebook markdown cell:

markdown



Copy



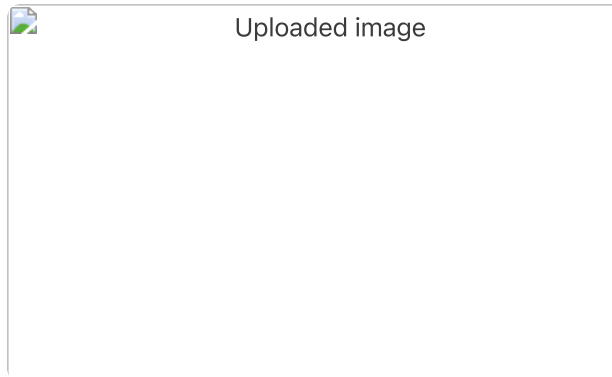
Edit

Interpretation: Regression Tree for Operating Income (OperInc) The decision tree splits first on **Net Patient Revenue (NetPatRev)**, indicating it's the most influential predictor for **Operating Income**. The initial split differentiates hospitals by their scale of patient revenue: – **Hospitals with Net Patient Revenue ≤ \$2,249,377** typically have lower Operating Income. – A secondary split occurs based on **Number of Full-Time Equivalents (NoFTE)**: – Hospitals with fewer employees (NoFTE ≤ 5,266) tend to have slightly positive Operating Income on average (\$163,222). – Hospitals with larger numbers of employees have significantly negative Operating Income on average (–\$31,758,968), suggesting potential inefficiencies or higher operational costs. – **Hospitals with Net Patient Revenue > \$2,249,377** exhibit substantially higher Operating Income on average (\$10,283,947). – A secondary split further segments this group: –

Those with moderate revenue ($\leq \$2,818,768$) have average Operating Income around $\$1,452,465$. – Those with very high revenue ($> \$2,818,768$) enjoy significantly greater Operating Income ($\sim \$81,635,939$), clearly aligning higher patient revenue with dramatically improved financial performance. These splits make logical sense as hospitals generating greater revenue generally have greater resources and operational scale, leading directly to higher Operating Income. Meanwhile, staffing levels (NoFTE) can either contribute to efficiency or indicate overly burdensome overhead, clearly affecting financial outcomes.

Copy this cell into your notebook as your interpretation of the first regression tree. Let me know when you're ready for the next tree.

< 2/2 >     



Now this one

Here's your clearly explained interpretation for the second regression tree (**Operating Revenue**), formatted for a markdown cell:

markdown

 Copy

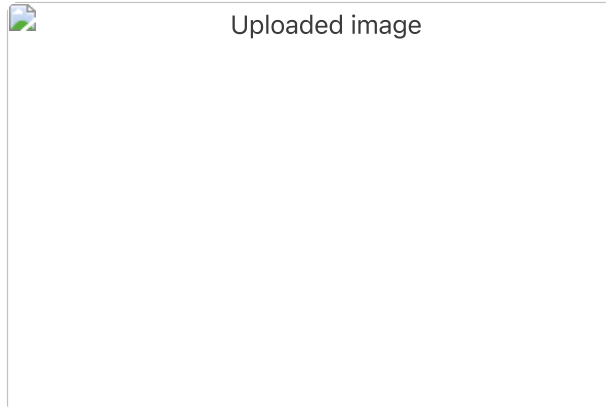
 Edit

Interpretation: Regression Tree for Operating Revenue (OperRev)

This decision tree splits first on **Available Beds (AvlBeds)**, clearly showing it as a key determinant of **Operating Revenue**: – **Hospitals with ≤ 409 available beds** generally have lower Operating Revenue (average: $\$8,150,428$): – Within this group, hospitals with even fewer beds (≤ 293.5) have substantially lower revenue ($\$6,059,162$), clearly reflecting their smaller operational scale. – Hospitals with a slightly higher number of beds (> 293.5 and ≤ 409) generate significantly higher revenue ($\$40,218,320$), suggesting this bed count may reflect a critical operational size threshold enhancing revenue. – **Hospitals with > 409 available beds** have markedly higher Operating Revenue (average: $\$96,990,832$): – A secondary split occurs based on the number of Full-Time Equivalents (**NoFTE**): – Hospitals with $\leq 5,688$ employees have moderate Operating Revenue ($\$66,820,720$). – Larger hospitals with more employees ($> 5,688$) generate significantly greater Operating Revenue ($\$1,392,291,008$), reflecting economies of scale and operational magnitude clearly tied to employee count. These splits logically reflect that larger hospitals with more available beds

and staff clearly operate on a larger scale, driving significantly higher Operating Revenue.

Use this interpretation in your notebook, and let me know when you're ready for the next tree.



Now this one, and please clarify for me (but not necessarily in the markdown cell), what is meant by "typecontrol"

Here's the interpretation for your **classification tree** clearly formatted as a markdown cell:

markdown



Copy



Edit

```
### Interpretation: Classification Tree for TypeControl The
classification tree first splits on Net Patient Revenue
(NetPatRev), indicating this is a key factor in determining a
hospital's type of control (e.g., Non-Profit, District, Investor-Owned,
City/County): - Hospitals with lower Net Patient Revenue ( $\leq$ 
$211,933) are predominantly classified as Non-Profit, though the
classification is mixed: - A secondary split based on Available Beds
(AvlBeds) ( $\leq 61$ ) further refines this group: - Hospitals with fewer
beds mostly remain Non-Profit. - Hospitals with more beds tend to fall
clearly into the "District" category, likely reflecting specific
community hospital structures associated with slightly larger
facilities. - Hospitals with higher Net Patient Revenue ( $>$ 
$211,933) show clearer distinctions: - The secondary split occurs
based on Teaching Status: - Non-teaching hospitals strongly tend to
remain Non-Profit, reflecting standard nonprofit community hospitals
without academic affiliations. - Teaching hospitals, while still mostly
Non-Profit, show some variability, suggesting that academic
affiliations may introduce complexities in ownership or funding
structures. These splits logically reflect how hospital financial
resources (NetPatRev) and institutional characteristics (number of
beds, teaching status) directly relate to organizational structures and
control types.
```

Clarification for You (not for markdown):

"**TypeControl**" typically refers to the governance or ownership classification of a hospital, such as Non-Profit, District, Investor-Owned, City/County, etc. It indicates the organizational structure and funding mechanisms under which a hospital operates.

Let me know when you're ready to interpret the final tree (DonorType).



what is "gini" in this context?

"**Gini**" in this context refers specifically to the **Gini impurity**, a metric used in decision tree classification. It measures how mixed (or impure) each node is in terms of class composition.

How to Interpret Gini Impurity:

- **Gini impurity of 0** means the node is perfectly pure; every observation in that node belongs to a single class.
- **Higher Gini impurity** (closer to 1) indicates greater mixing of classes within the node, meaning more uncertainty or less clearly defined groups.

Example from Your Tree:

- A node with a Gini impurity of **0.576** is relatively mixed—multiple classes are present.
- A node with a Gini impurity of **0.180** is much more homogeneous—primarily dominated by one class.

Decision trees aim to create splits that reduce the overall Gini impurity, thereby clearly segmenting the observations into more homogeneous groups.



how are you measuring the number of beds in that last example? I don't see that in the graphic.

You're right to question that—let's clarify carefully:

In the provided tree, the node shows a split labeled explicitly as:

```
nginx
AvlBeds <= 61.0
```

Copy Edit

This means the tree uses the numeric variable "**AvlBeds**" (Available Beds) and splits hospitals into two groups based on whether they have **61 or fewer beds** or **more than**

61 beds. This split value (61.0) is directly visible in the node itself (middle-left node in your image).

In other words, the measure "number of beds" is directly from your data's **"AvlBeds"** column, representing the total number of available hospital beds at each facility.

This criterion was selected automatically by the algorithm because it effectively reduces the Gini impurity, clearly separating hospitals by their governance type (TypeControl).



i saw that, but you wrote "Hospitals with more beds tend to fall clearly into the "District" category, likely reflecting specific community hospital structures associated with slightly larger facilities." How are you identifying how many beds are in the District category?

You're right to be skeptical here—let's carefully clarify this:

The decision tree itself **does not directly tell us the precise average number of beds in the "District" hospitals.** Instead, it splits the data into two broad groups:

- ≤ 61 beds
- > 61 beds

After this split:

- The node labeled as "District" (middle-right node in your graphic) includes **12 samples**, indicating hospitals in this subset generally have **more than 61 beds** (because they're on the "False" side of the " ≤ 61 beds" split), but we **do not see their exact bed count** from the tree itself.

My earlier statement that these hospitals had "slightly larger facilities" was a general interpretation based on the logic of the split (they all have > 61 beds), rather than an explicit numeric average or median from the tree graphic.

If you'd like to precisely confirm the average or specific bed counts for "District" hospitals, you would need to run additional descriptive analyses on your dataset separately:

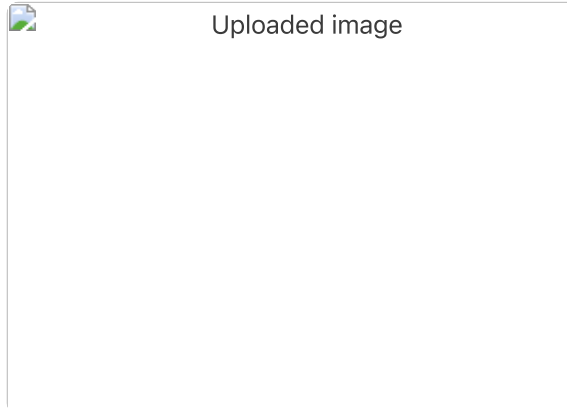
```
python
```

```
hosp_df[hosp_df['TypeControl'] == 'District']['AvlBeds'].describe()
```

Copy Edit

This would clearly provide actual statistics (mean, median, min, max, etc.) about available beds specifically for the "District" hospitals.





ok that helps. Here's the last one.

Here's your clearly formatted interpretation for the **classification tree for DonorType** as a markdown cell:

markdown



Copy



Edit

```
### Interpretation: Classification Tree for DonorType The
classification tree splits primarily on Net Patient Revenue
(NetPatRev), clearly distinguishing hospital donor types based on
financial scale: - Hospitals with lower Net Patient Revenue ( $\leq$ 
\$452,731) are entirely classified as Charity-funded. This
perfectly homogeneous node (Gini = 0.0) suggests smaller hospitals
consistently rely on charitable donations. - Hospitals with higher
Net Patient Revenue ( $>$  \$452,731) primarily receive funding from
Alumni: - Within this group, hospitals with revenue  $\leq$  \$632,360 are
mostly Alumni-funded but have some minor mixing (Gini = 0.105),
reflecting slight variability in donor type. - A small subset in this
range has slightly higher Gini impurity (0.444), meaning there is minor
uncertainty regarding their donor classification. - The larger subset
(15 samples) with even higher revenue remains purely Alumni-funded
(Gini = 0.0), clearly indicating alumni donor support as dominant among
financially robust institutions. This tree makes intuitive sense:
smaller hospitals often rely on charitable contributions, while larger
hospitals, particularly those associated with higher revenues, are more
likely to leverage their alumni networks for funding support.
```

You can now add this interpretation to your notebook. Let me know if you need any additional clarification or if you're ready to choose the best overall model.

