

```
# Unbalanced dataset
# Applied Logistic Regression, Naive Bayes Classifier, and Decision Tree
# Compared changes in accuracy rates when 10-fold cross-validation applied
```

```
import sys
!{sys.executable} -m pip install -U pandas-profiling[notebook]
!jupyter nbextension enable --py widgetsnbextension
!pip install matplotlib
!pip install graphviz
```

```
Requirement already satisfied: pandas-profiling[notebook] in /usr/local/lib/python3.7/dist-packages (1.4.1)
WARNING: pandas-profiling 1.4.1 does not provide the extra 'notebook'
Requirement already satisfied: six>=1.9 in /usr/local/lib/python3.7/dist-packages (from pandas-profiling[notebook]) (1.16.0)
Requirement already satisfied: Jinja2>=2.8 in /usr/local/lib/python3.7/dist-packages (from pandas-profiling[notebook]) (2.11.3)
Requirement already satisfied: MarkupSafe>=0.23 in /usr/local/lib/python3.7/dist-packages (from Jinja2>=2.8->pandas-profiling[notebook]) (2.0.1)
Requirement already satisfied: matplotlib>=1.4 in /usr/local/lib/python3.7/dist-packages (from pandas-profiling[notebook]) (3.3.4)
Requirement already satisfied: pandas>=0.19 in /usr/local/lib/python3.7/dist-packages (from pandas-profiling[notebook]) (1.1.5)
Requirement already satisfied: numpy>=1.11 in /usr/local/lib/python3.7/dist-packages (from pandas>=0.19->pandas-profiling[notebook]) (1.19.5)
Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.7/dist-packages (from matplotlib>=1.4->pandas-profiling[notebook]) (0.11.0)
Requirement already satisfied: python-dateutil>=2.1 in /usr/local/lib/python3.7/dist-packages (from pandas>=0.19->pandas-profiling[notebook]) (2.8.2)
Requirement already satisfied: pyparsing!=2.0.4,!=2.1.2,!=2.1.6,>=2.0.1 in /usr/local/lib/python3.7/dist-packages (from python-dateutil>=2.1->pandas>=0.19->pandas-profiling[notebook]) (3.0.7)
Requirement already satisfied: kiwisolver>=1.0.1 in /usr/local/lib/python3.7/dist-packages (from matplotlib>=1.4->pandas-profiling[notebook]) (1.3.2)
Requirement already satisfied: typing-extensions in /usr/local/lib/python3.7/dist-packages (from pandas>=0.19->pandas-profiling[notebook]) (4.1.1)
Requirement already satisfied: pytz>=2017.3 in /usr/local/lib/python3.7/dist-packages (from pandas>=0.19->pandas-profiling[notebook]) (2021.3)
Enabling notebook extension jupyter-js-widgets/extension...
Paths used for configuration of notebook:
  /root/.jupyter/nbconfig/notebook.json
- Validating: OK
Paths used for configuration of notebook:
  /root/.jupyter/nbconfig/notebook.json
Requirement already satisfied: matplotlib in /usr/local/lib/python3.7/dist-packages (3.3.4)
Requirement already satisfied: pyparsing!=2.0.4,!=2.1.2,!=2.1.6,>=2.0.1 in /usr/local/lib/python3.7/dist-packages (3.0.7)
Requirement already satisfied: numpy>=1.11 in /usr/local/lib/python3.7/dist-packages (1.19.5)
Requirement already satisfied: python-dateutil>=2.1 in /usr/local/lib/python3.7/dist-packages (2.8.2)
Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.7/dist-packages (0.11.0)
Requirement already satisfied: kiwisolver>=1.0.1 in /usr/local/lib/python3.7/dist-packages (1.3.2)
Requirement already satisfied: typing-extensions in /usr/local/lib/python3.7/dist-packages (4.1.1)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.7/dist-packages (1.16.0)
Requirement already satisfied: graphviz in /usr/local/lib/python3.7/dist-packages (0.10)
```

```
from google.colab import files
uploaded = files.upload()
```

Choose files stroke_preprocessed.arff

- **stroke_preprocessed.arff**(n/a) - 368681 bytes, last modified: 06/03/2022 - 100% done
- Saving stroke_preprocessed.arff to stroke_preprocessed (1).arff

```
import pandas as pd
from scipy.io import arff
```

```

import numpy as np

# Timing how long predictors take to run for efficiency calculations
# Import libraries
import time

data_file = "stroke_preprocessed.arff"
data = arff.loadarff(data_file)

df = pd.DataFrame(data[0])
for col in df.columns:
    if df[col].dtype == 'object':
        # Ensure data isn't read as bytes but rather as strings from file
        df[col] = df[col].str.decode('utf-8')
# Examine data types
print(df.dtypes)

    "id"                float64
    "gender"            object
    "age"               float64
    "hypertension"      object
    "heart_disease"     object
    "ever_married"      object
    "work_type"         object
    "residence_type"    object
    "avg_glucose_level" float64
    "bmi"               float64
    "smoking_status"    object
    "stroke"            object
    dtype: object

# Display first 10 rows
df.head(10)

```

	"id"	"gender"	"age"	"hypertension"	"heart_disease"	"ever_married"	"work_type"
0	9046.0	Male	67.0	0	1	Yes	Private
1	51676.0	Female	61.0	0	0	Yes	Self-employed
2	31112.0	Male	80.0	0	1	Yes	Private

```
# Examine meta info about data
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5110 entries, 0 to 5109
Data columns (total 12 columns):
#   Column                Non-Null Count  Dtype
---  -
0   "id"                  5110 non-null   float64
1   "gender"              5110 non-null   object
2   "age"                 5110 non-null   float64
3   "hypertension"        5110 non-null   object
4   "heart_disease"       5110 non-null   object
5   "ever_married"        5110 non-null   object
6   "work_type"           5110 non-null   object
7   "residence_type"      5110 non-null   object
8   "avg_glucose_level"   5110 non-null   float64
9   "bmi"                 5110 non-null   float64
10  "smoking_status"      5110 non-null   object
11  "stroke"              5110 non-null   object
dtypes: float64(4), object(8)
memory usage: 479.2+ KB
```

```
# The original 201 null values were all from bmi column, and they have been replaced by place
# Convert the 5000 values back into null values
df = df.replace(5000.0, np.nan)
```

```
# Check head of dataset again
df.head(10)
```

	"id"	"gender"	"age"	"hypertension"	"heart_disease"	"ever_married"	"work_type"
0	9046.0	Male	67.0	0	1	Yes	Private
1	51676.0	Female	61.0	0	0	Yes	Self-employed
2	31112.0	Male	80.0	0	1	Yes	Private
3	60182.0	Female	49.0	0	0	Yes	Private
4	1665.0	Female	70.0	1	0	Yes	Self-employed

```
# Check structure of data types to ensure bmi remains float
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5110 entries, 0 to 5109
Data columns (total 12 columns):
#   Column                      Non-Null Count  Dtype
---  -
0   "id"                        5110 non-null   float64
1   "gender"                    5110 non-null   object
2   "age"                        5110 non-null   float64
3   "hypertension"              5110 non-null   object
4   "heart_disease"             5110 non-null   object
5   "ever_married"              5110 non-null   object
6   "work_type"                  5110 non-null   object
7   "residence_type"            5110 non-null   object
8   "avg_glucose_level"         5110 non-null   float64
9   "bmi"                        4909 non-null   float64
10  "smoking_status"            5110 non-null   object
11  "stroke"                     5110 non-null   object
dtypes: float64(4), object(8)
memory usage: 479.2+ KB
```

```
# Remove records with NAs from dataset
df_noNA = df
df_noNA = df_noNA.dropna()
df_noNA.head(10)
```

	"id"	"gender"	"age"	"hypertension"	"heart_disease"	"ever_married"	"work_type"
0	9046.0	Male	67.0	0	1	Yes	Priv:
2	31112.0	Male	80.0	0	1	Yes	Priv:
3	60182.0	Female	49.0	0	0	Yes	Priv:
4	1665.0	Female	79.0	1	0	Yes	St employ
5	56669.0	Male	81.0	0	0	Yes	Priv:

```
# Change 'stroke' attribute into data type float
df_noNA["stroke"] = df_noNA["stroke"].astype(float)
df_noNA.head(10)
```

/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:2: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_

	"id"	"gender"	"age"	"hypertension"	"heart_disease"	"ever_married"	"work_type"
0	9046.0	Male	67.0	0	1	Yes	Priv:
2	31112.0	Male	80.0	0	1	Yes	Priv:
3	60182.0	Female	49.0	0	0	Yes	Priv:
4	1665.0	Female	79.0	1	0	Yes	St employ
5	56669.0	Male	81.0	0	0	Yes	Priv:
6	53882.0	Male	74.0	1	1	Yes	Priv:
7	10434.0	Female	69.0	0	0	No	Priv:
9	60491.0	Female	78.0	0	0	Yes	Priv:
10	12109.0	Female	81.0	1	0	Yes	Priv:
11	12095.0	Female	61.0	0	1	Yes	Govt_



```
print(df_noNA.dtypes)
```

```
"id"                float64
"gender"            object
"age"              float64
"hypertension"      object
"heart_disease"     object
```

```

"ever_married"      object
"work_type"         object
"residence_type"    object
"avg_glucose_level" float64
"bmi"               float64
"smoking_status"    object
"stroke"            float64
dtype: object

```

```

# See if there are any extreme values in numeric data
df_noNA.describe()

```

	"id"	"age"	"avg_glucose_level"	"bmi"	"stroke"
count	4909.000000	4909.000000	4909.000000	4909.000000	4909.000000
mean	37064.313506	42.865374	105.305150	28.893237	0.042575
std	20995.098457	22.555115	44.424341	7.854067	0.201917
min	77.000000	0.080000	55.120000	10.300000	0.000000
25%	18605.000000	25.000000	77.070000	23.500000	0.000000
50%	37608.000000	44.000000	91.680000	28.100000	0.000000
75%	55220.000000	60.000000	113.570000	33.100000	0.000000
max	72940.000000	82.000000	271.740000	97.600000	1.000000



```

# Normalize continuous numeric variables
# Such as age, avg_glucose_level, and bmi
# Using z-score methods

```

```

# Import libraries for normalization
from sklearn import preprocessing
scaler = preprocessing.MinMaxScaler()

```

```

# Only need to normalize continuous numeric variables
var_to_norm = ["age", "avg_glucose_level", "bmi"]
df_noNA[var_to_norm] = scaler.fit_transform(df_noNA[var_to_norm])

```

```

# Examine first 10 rows of normalized dataset
df_noNA.head()

```

```

# The 3 columns are now standardized to values between 0-1

```

/usr/local/lib/python3.7/dist-packages/pandas/core/frame.py:3678: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stable/user_](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html)
self[col] = igetitem(value, i)

	"id"	"gender"	"age"	"hypertension"	"heart_disease"	"ever_married"	"work_t"
0	9046.0	Male	0.816895	0	1	Yes	Pr
2	31112.0	Male	0.975586	0	1	Yes	Pr
3	60182.0	Female	0.597168	0	0	Yes	Pr
4	1665.0	Female	0.963379	1	0	Yes	empl

Create list of categorical columns - removed ID as it is not relevant to prediction
cat_cols = ["gender", "hypertension", "heart_disease", "ever_married", "work_type",



```
# Create copy of a data frame in memory w/ a different name
df_dummy = df_noNA.copy()
# Convert only categorical feature into dummy/one-hot features
df_dummy = pd.get_dummies(df_noNA, columns = cat_cols, prefix = cat_cols)
# Print dataset
df_dummy
```

```

# Create train test set split
from sklearn.model_selection import train_test_split

# Set class name as "stroke", all else will be used as features
class_col_name = "stroke"
# Obtain necessary dummy feature names
dummy_feature_name = df_dummy.columns.values.tolist()
dummy_feature_names = dummy_feature_name[5:]
# 70% training, 30% test set split
x_train, x_test, y_train, y_test = train_test_split(df_dummy.loc[:, dummy_feature_names], df_

5404 141800 0.157715 0.221402 0.005074 0.0 1
start.=.time.time()

# Import needed libraries for Logistic Regression Model
from sklearn.linear_model import LogisticRegression

# Begin to Implement Logistic Regression Model
log_regr = LogisticRegression()

# Apply data into Logistic Regression Model
log_regr.fit(x_train, y_train)
y_pred = log_regr.predict(x_test)

# Obtain Confusion Matrix and Evaluation Metrics for the Logistic Regression Model
from sklearn import metrics
cnf_matrix = metrics.confusion_matrix(y_test, y_pred)
cnf_matrix

array([[1409,  0],
       [ 64,  0]])

# Display Evaluation Metrics for Logistic Regression Model
print("Logistic Regression Accuracy:\t", metrics.accuracy_score(y_test, y_pred))
print("Logistic Regression Precision:\t", metrics.precision_score(y_test, y_pred))
print("Logistic Regression Recall:\t", metrics.recall_score(y_test, y_pred))

Logistic Regression Accuracy: 0.956551255940258
Logistic Regression Precision: 0.0
Logistic Regression Recall: 0.0
/usr/local/lib/python3.7/dist-packages/sklearn/metrics/_classification.py:1318: Undefined
_warn_prf(average, modifier, msg_start, len(result))

# Import libraries for cross-validation
from sklearn.model_selection import cross_val_score, cross_val_predict

# 10-Fold Cross Validation for Logistic Regression

```



```
cv_lr = cross_val_score(log_regr, df_dummy, df_dummy[class_col_name], cv=10)
print("Cross-validated scores:\t", cv_lr)
```

Increased Accuracy score from 0.957 to 1

```
/usr/local/lib/python3.7/dist-packages/sklearn/linear_model/_logistic.py:818: Converge
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.
```

Increase the number of iterations (max_iter) or scale the data as shown in:

<https://scikit-learn.org/stable/modules/preprocessing.html>

Please also refer to the documentation for alternative solver options:

https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression

```
extra_warning_msg=_LOGISTIC_SOLVER_CONVERGENCE_MSG,
/usr/local/lib/python3.7/dist-packages/sklearn/linear_model/_logistic.py:818: Converge
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.
```

Increase the number of iterations (max_iter) or scale the data as shown in:

<https://scikit-learn.org/stable/modules/preprocessing.html>

Please also refer to the documentation for alternative solver options:

https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression

```
extra_warning_msg=_LOGISTIC_SOLVER_CONVERGENCE_MSG,
Cross-validated scores: [1.          1.          1.          1.          1.          1.
 1.          1.          0.95723014 0.95918367]
```



Cross validation accuracy for Logistic Regression (R2 score)

```
predictions = cross_val_predict(log_regr, df_dummy, df_dummy[class_col_name], cv=10)
accuracy = metrics.r2_score(df_dummy[class_col_name], predictions)
print("Cross-Predicted Accuracy for Logistic Regression: ", accuracy)
```

```
/usr/local/lib/python3.7/dist-packages/sklearn/linear_model/_logistic.py:818: Converge
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.
```

Increase the number of iterations (max_iter) or scale the data as shown in:

<https://scikit-learn.org/stable/modules/preprocessing.html>

Please also refer to the documentation for alternative solver options:

https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression

```
extra_warning_msg=_LOGISTIC_SOLVER_CONVERGENCE_MSG,
/usr/local/lib/python3.7/dist-packages/sklearn/linear_model/_logistic.py:818: Converge
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.
```

Increase the number of iterations (max_iter) or scale the data as shown in:

<https://scikit-learn.org/stable/modules/preprocessing.html>

Please also refer to the documentation for alternative solver options:

https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression

```
extra_warning_msg=_LOGISTIC_SOLVER_CONVERGENCE_MSG,
Cross-Predicted Accuracy for Logistic Regression: 0.7951043469408531
```



```
end=time.time()
```

```
print("Time to run Logistic Regression: ", end-start)
```

Time to run Logistic Regression: 2.4975674152374268

```
start = time.time()

# Naive Bayes modeling
from sklearn.naive_bayes import MultinomialNB

# Create Multinomial NB Classifier
nb = MultinomialNB()

# Train model using training sets
nb.fit(x_train, y_train)

MultinomialNB()

# Predict response for test dataset
y_pred = nb.predict(x_test)

# Print Naive Bayes output
print("Number of features used: ", nb.n_features_)
print("Classes: ", nb.classes_)
print("Number of records for classes: ", nb.class_count_)
print("Log prior probability for classes: ", nb.class_log_prior_)
print("Log conditional probability for each feature given a class: ", nb.feature_log_prob_)

Number of features used: 20
Classes: [0. 1.]
Number of records for classes: [3291. 145.]
Log prior probability for classes: [-0.04311653 -3.16532954]
Log conditional probability for each feature given a class: [[-2.47270743 -2.84008951 -
-5.08288006 -2.97245497 -2.38986067 -3.98660149 -7.15535293 -2.49664198
-3.8615758 -3.90153906 -2.67927954 -2.60089142 -3.11914766 -3.7690812
-2.91883388 -3.83311859]
[-2.49950545 -2.81502232 -6.94215671 -2.28819636 -3.20448709 -2.16303321
-3.6099522 -4.2341065 -2.05935478 -3.85111425 -6.94215671 -2.4312972
-3.38680864 -6.94215671 -2.75250196 -2.54770755 -3.94642443 -3.18095659
-2.93482352 -3.5081695 ]]
/usr/local/lib/python3.7/dist-packages/sklearn/utils/deprecation.py:103: FutureWarning:
warnings.warn(msg, category=FutureWarning)
```

```
# Confusion matrix and Evaluation metrics
from sklearn.metrics import confusion_matrix
cf = confusion_matrix(y_test, y_pred)
print("Confusion Matrix")
print(cf)
tn, fp, fn, tp = cf.ravel()
print("TP: ", tp, ", FP: ", fp, ", TN: ", tn, ", FN: ", fn)
```

```
Confusion Matrix
[[1404    5]
 [  64    0]]
TP:  0 , FP:  5 , TN: 1404 , FN:  64
```

```
# Classifier report
from sklearn.metrics import classification_report
from sklearn import metrics
```

```
print(classification_report(y_test, y_pred))
```

	precision	recall	f1-score	support
0.0	0.96	1.00	0.98	1409
1.0	0.00	0.00	0.00	64
accuracy			0.95	1473
macro avg	0.48	0.50	0.49	1473
weighted avg	0.91	0.95	0.93	1473

```
# Display Evaluation Metrics for Naive Bayes Classifier
print("Naive Bayes Classifier Accuracy:\t", metrics.accuracy_score(y_test, y_pred))
print("Naive Bayes Classifier Precision:\t", metrics.precision_score(y_test, y_pred))
print("Naive Bayes Classifier Recall:\t\t", metrics.recall_score(y_test, y_pred))
```

```
Naive Bayes Classifier Accuracy:      0.9531568228105907
Naive Bayes Classifier Precision:      0.0
Naive Bayes Classifier Recall:         0.0
```

```
# 10-Fold Cross Validation for Naive Bayes Classifier
cv_nb = cross_val_score(nb, df_dummy, df_dummy[class_col_name], cv=10)
print("Cross-validated scores:\t", cv_nb)
```

```
# Increased Accuracy score from 0.953 to 0.998
```

```
Cross-validated scores: [0.99796334 0.99592668 0.99796334 0.99796334 0.99592668 0.99389
0.99592668 0.99796334 0.99796334 0.99795918]
```



```
# Cross validation accuracy for Naive Bayes Classifier (R2 score)
predictions = cross_val_predict(nb, df_dummy, df_dummy[class_col_name], cv=10)
accuracy = metrics.r2_score(df_dummy[class_col_name], predictions)
print("Cross-Predicted Accuracy for Naive Bayes Classifier: ", accuracy)
```

```
Cross-Predicted Accuracy for Naive Bayes Classifier: 0.9250381757100682
```

```
end.=.time.time()
```

```
print("Time to run Naive Bayes Classifier: ", end--start)
```

Time to run Naive Bayes Classifier: 0.33214259147644043

```
start = time.time()

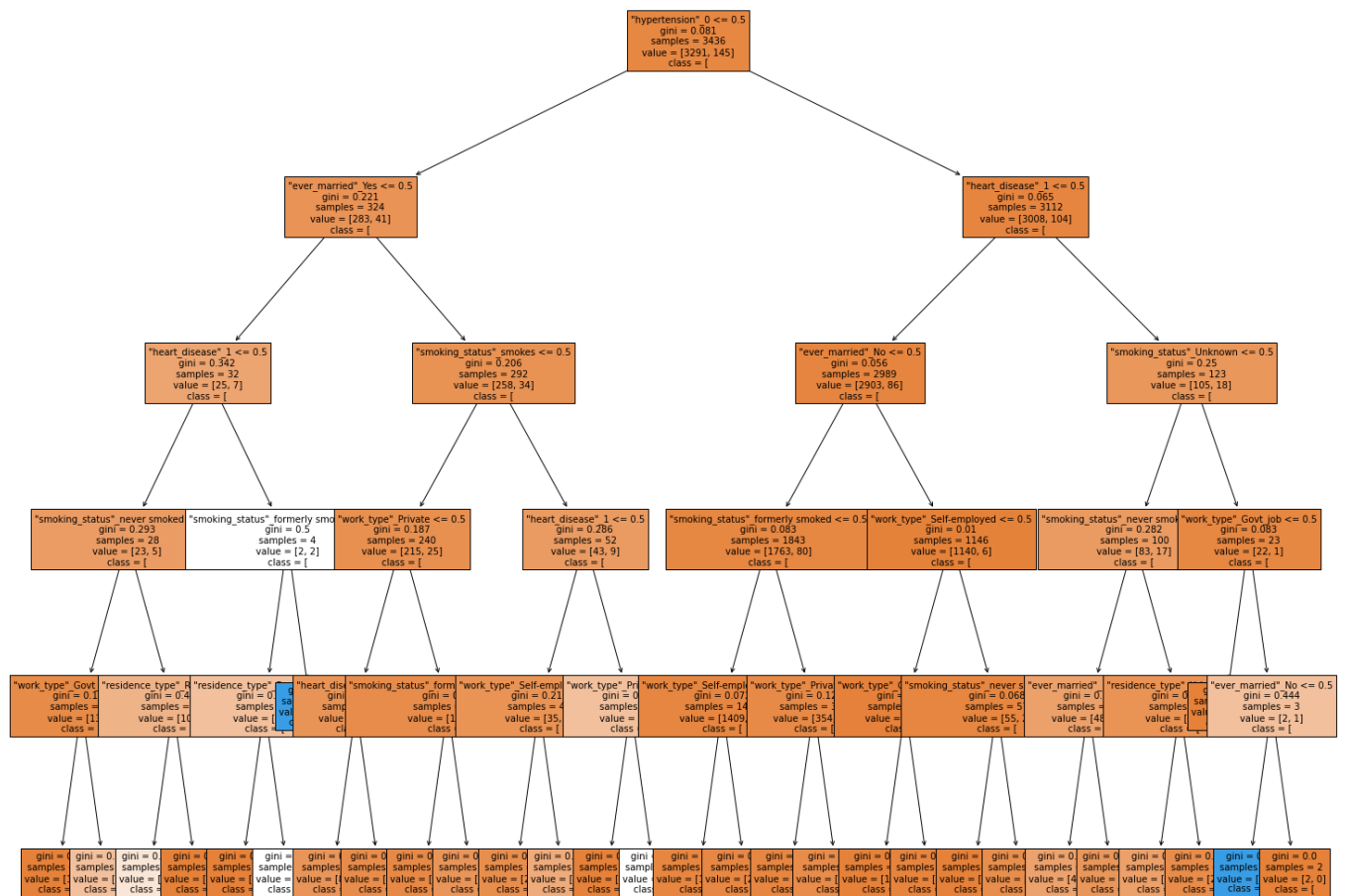
# Decision tree on dummy encoded data
from sklearn import tree
clf = tree.DecisionTreeClassifier(max_depth = 5) # 5 levels set
clf = clf.fit(x_train, y_train)

import graphviz
# Obtain unique class values to show on tree
class_values = df_dummy[class_col_name].unique()
print("class names: ", class_values)

class names:  [1. 0.]

# Import libraries for plotting the decision tree
import matplotlib
from matplotlib import pyplot as plt

# Plot decision tree
fig = plt.figure(figsize=(25,20))
_ = tree.plot_tree(clf, feature_names = dummy_feature_names, class_names = str(class_values),
```



```
# Save decision tree figure
fig.savefig("decision_tree6.png")
```

```
# Perform prediction on test set
y_pred = clf.predict(x_test)
```

```
# Get decision tree confusion matrix
cf = confusion_matrix(y_test, y_pred)
print("Confusion Matrix")
print(cf)
tn, fp, fn, tp = cf.ravel()
print("TP: ", tp, ", FP: ", fp, ", TN: ", tn, ", FN: ", fn)
```

```
Confusion Matrix
[[1407   2]
 [ 64   0]]
TP:  0 , FP:  2 , TN: 1407 , FN: 64
```

```
# Get decision tree report
from sklearn.metrics import classification_report
from sklearn import metrics
print(classification_report(y_test, y_pred))
```

	precision	recall	f1-score	support
0.0	0.96	1.00	0.98	1409
1.0	0.00	0.00	0.00	64
accuracy			0.96	1473
macro avg	0.48	0.50	0.49	1473
weighted avg	0.91	0.96	0.93	1473

Display Evaluation Metrics for Decision Tree

```
print("Decision Tree Accuracy:\t\t", metrics.accuracy_score(y_test, y_pred))
print("Decision Tree Precision:\t", metrics.precision_score(y_test, y_pred))
print("Decision Tree Recall:\t\t", metrics.recall_score(y_test, y_pred))
```

```
Decision Tree Accuracy:      0.955193482688391
Decision Tree Precision:     0.0
Decision Tree Recall:        0.0
```

10-Fold Cross Validation for Decision Tree

```
cv_dt = cross_val_score(clf, df_dummy, df_dummy[class_col_name], cv=10)
print("Cross-validated scores:\t", cv_dt)
```

Increased Accuracy score from 0.954 to 1

```
Cross-validated scores: [1. 1. 1. 1. 1. 1. 1. 1. 1. 1.]
```

Cross validation accuracy for Decision Tree (R2 score)

```
predictions = cross_val_predict(clf, df_dummy, df_dummy[class_col_name], cv=10)
accuracy = metrics.r2_score(df_dummy[class_col_name], predictions)
print("Cross-Predicted Accuracy for Decision Tree: ", accuracy)
```

```
Cross-Predicted Accuracy for Decision Tree: 1.0
```

```
end.=.time.time()
```

```
print("Time to run Decision Tree: ", end.-.start)
```

```
Time to run Decision Tree: 4.365260124206543
```

✓ 0s completed at 6:02 AM

