

(1) WhatsApp | Introducing C | Census Data | ML Project v | Machine learn | 6.3. Preprocess | Supervised L | DataCamp Ap | EEET2493 Lab | +

localhost:8888/notebooks/ML%20Project%20vinod%20yadav/Machine%20learning%20Project.ipynb

School Manageme... myphAdmin

jupyter Machine learning Project Last Checkpoint 04/11/2023 (unsaved changes)

File Edit View Insert Cell Kernel Widgets Help

Not Trusted Python 3 (ipykernel) O

```
# Initialize a scaler, then apply it to the features
scaler = MinMaxScaler() # default=(0, 1)
numerical = ['age', 'education-num', 'capital-gain', 'capital-loss', 'hours-per-week']
features_raw[numerical] = scaler.fit_transform(data[numerical])
```

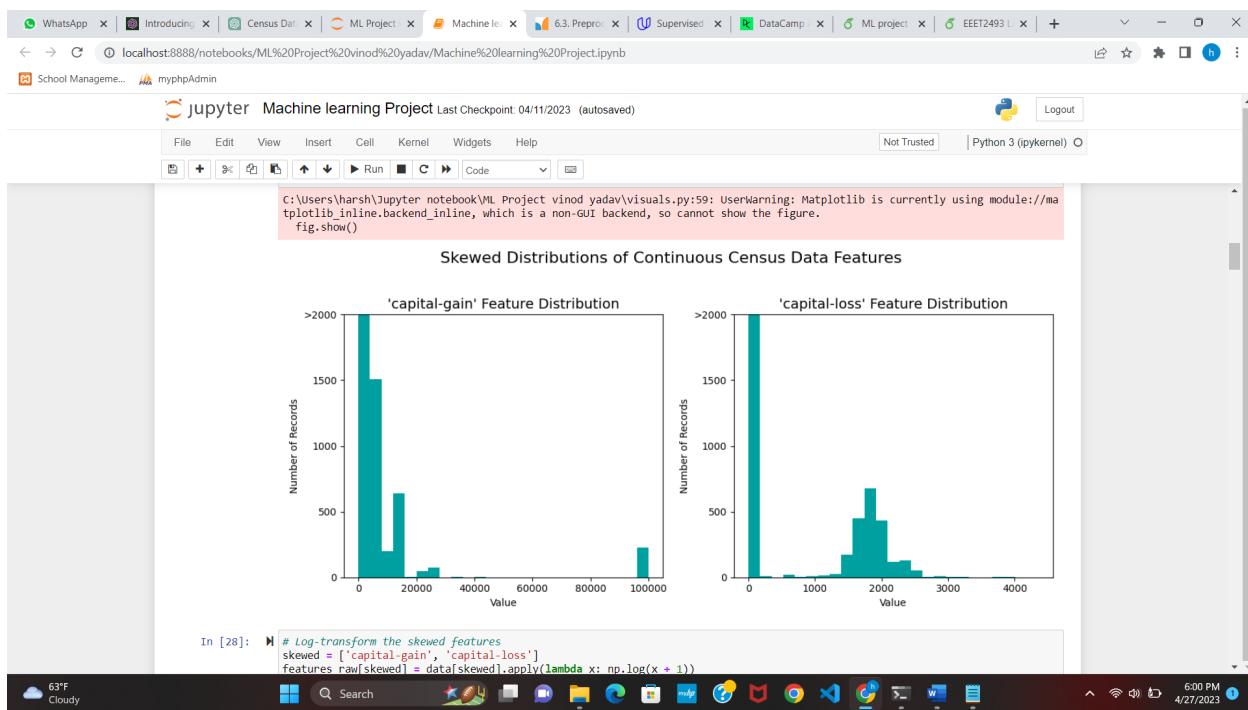
# Show an example of a record with scaling applied
display(features\_raw.head(n = 5))

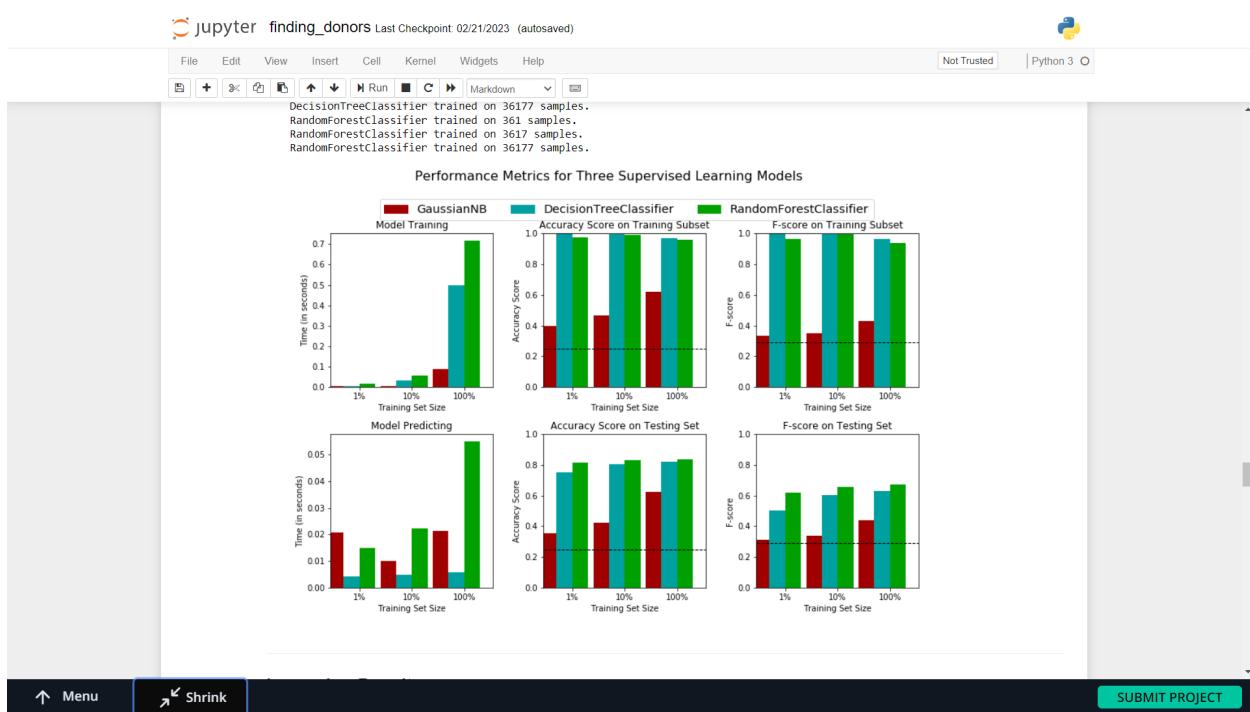
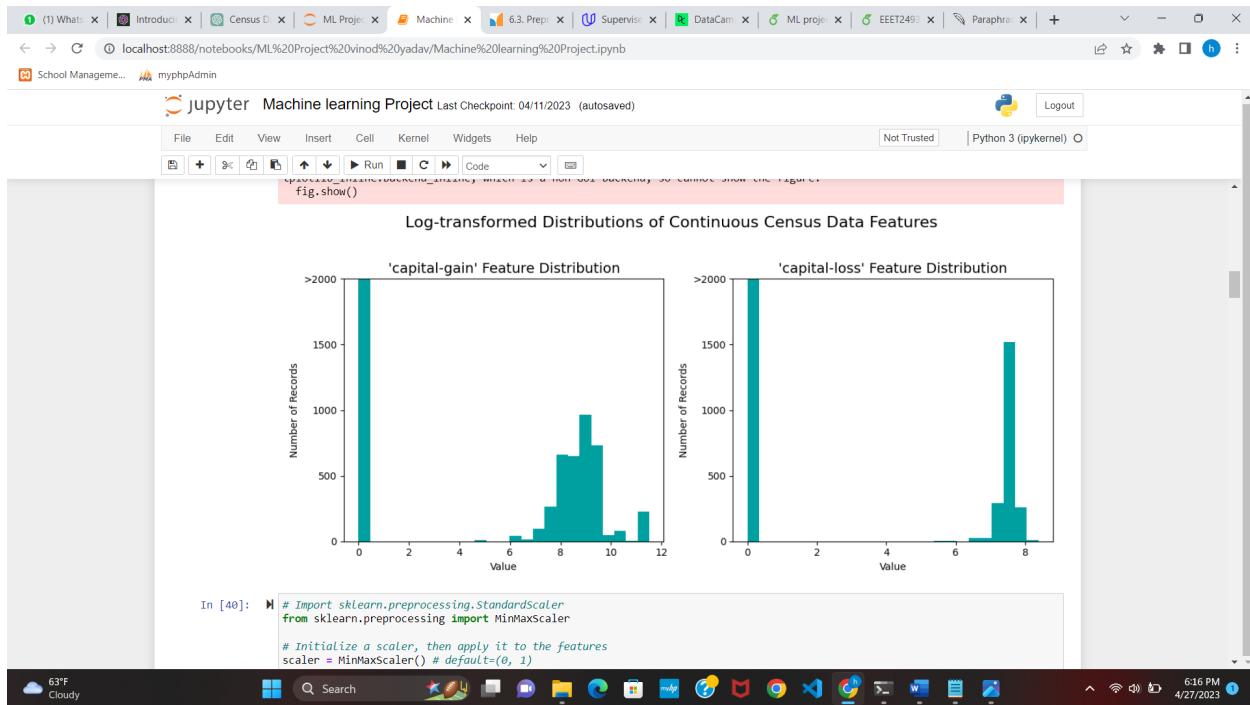
	age	workclass	education_level	education-num	marital-status	occupation	relationship	race	sex	capital-gain	capital-loss	hours-per-week	native-country
0	301370	State-gov	Bachelors	0.800000	Never-married	Adm-clerical	Not-in-family	White	Male	0.02174	0.0	0.397959	United-States
1	0.452055	Self-emp-not-inc	Bachelors	0.800000	Married-civ-spouse	Exec-managerial	Husband	White	Male	0.00000	0.0	0.122449	United-States
2	0.287671	Private	HS-grad	0.533333	Divorced	Handlers-cleaners	Not-in-family	White	Male	0.00000	0.0	0.397959	United-States
3	0.493151	Private	11th	0.400000	Married-civ-spouse	Handlers-cleaners	Husband	Black	Male	0.00000	0.0	0.397959	United-States
4	0.150685	Private	Bachelors	0.800000	Married-civ-spouse	Prof-specialty	Wife	Black	Female	0.00000	0.0	0.397959	Cuba

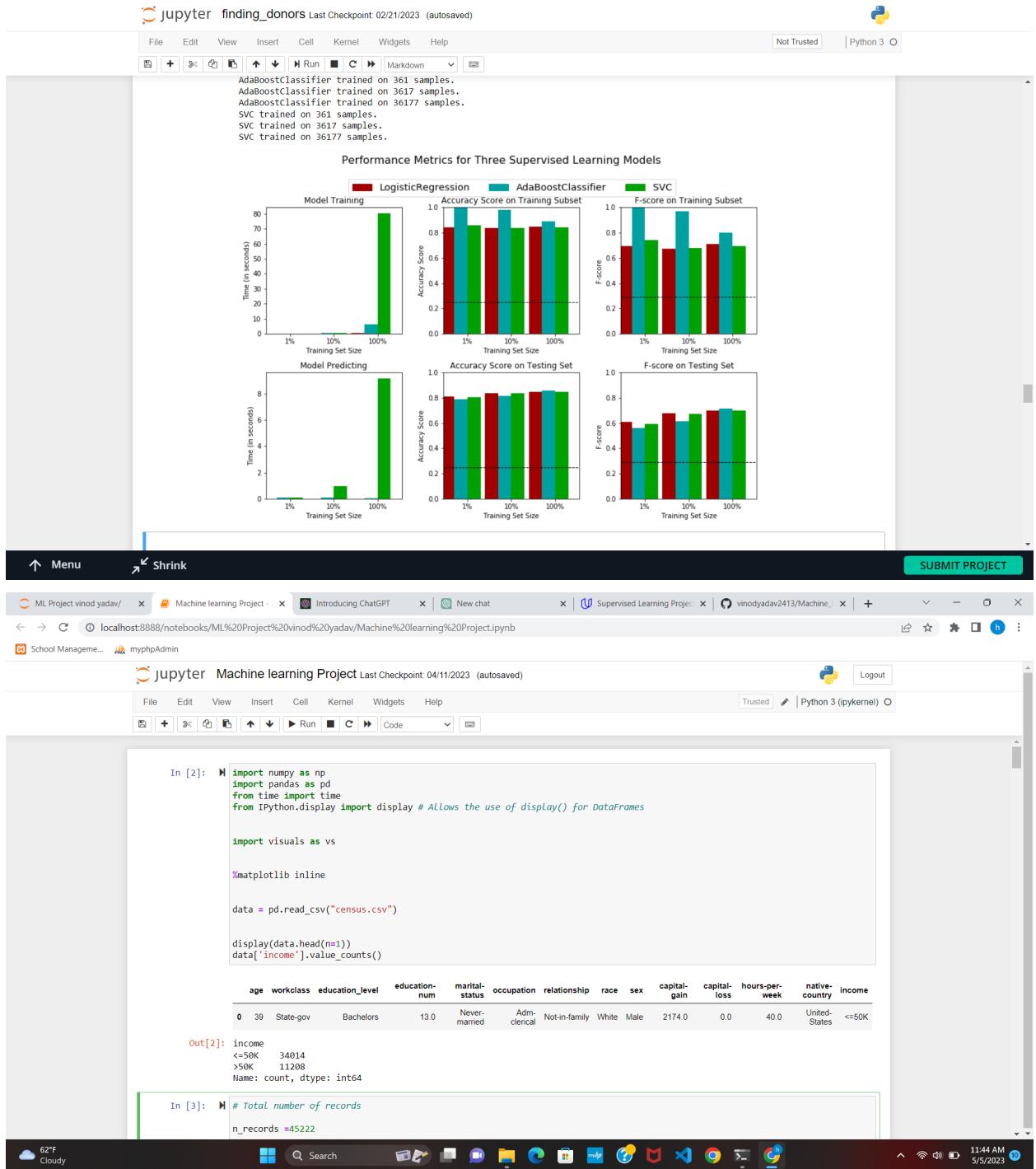
In [30]: # TODO: One-hot encode the 'features\_log\_minmax\_transform' data using pandas.get\_dummies()
features\_final = pd.get\_dummies(features\_raw)

# TODO: Encode the 'income\_raw' data to numerical values
income = income\_raw.apply(lambda x: 0 if x == '<50K' else 1)

# Print the number of features after one-hot encoding
encoded = list(features\_final.columns)
print("{} total features after one-hot encoding.".format(len(encoded)))







ML Project vinod yadav/ Machine learning Project - Introducing ChatGPT New chat Supervised Learning Project vinodyadav2413/Machine... +

localhost:8888/notebooks/ML%20Project%20vinod%20yadav/Machine%20learning%20Project.ipynb

School Manageme... myphAdmin

jupyter Machine learning Project Last Checkpoint 04/11/2023 (autosaved)

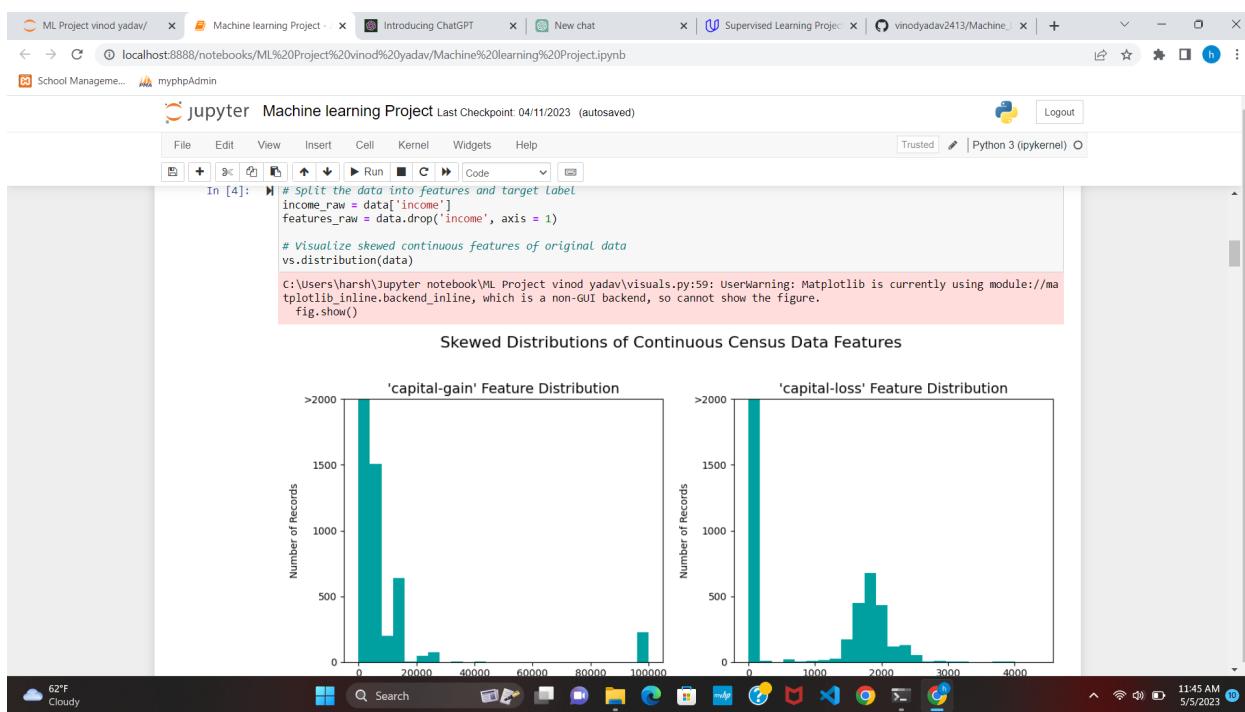
File Edit View Insert Cell Kernel Widgets Help Trusted Python 3 (ipykernel) Logout

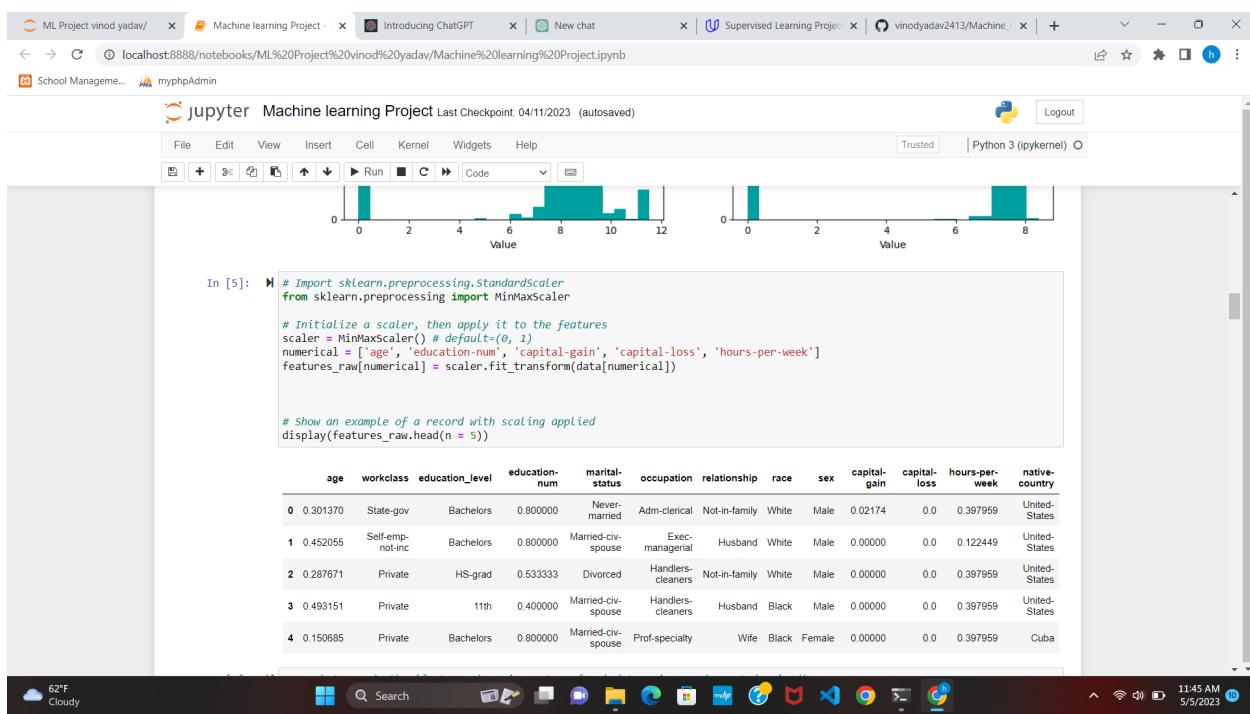
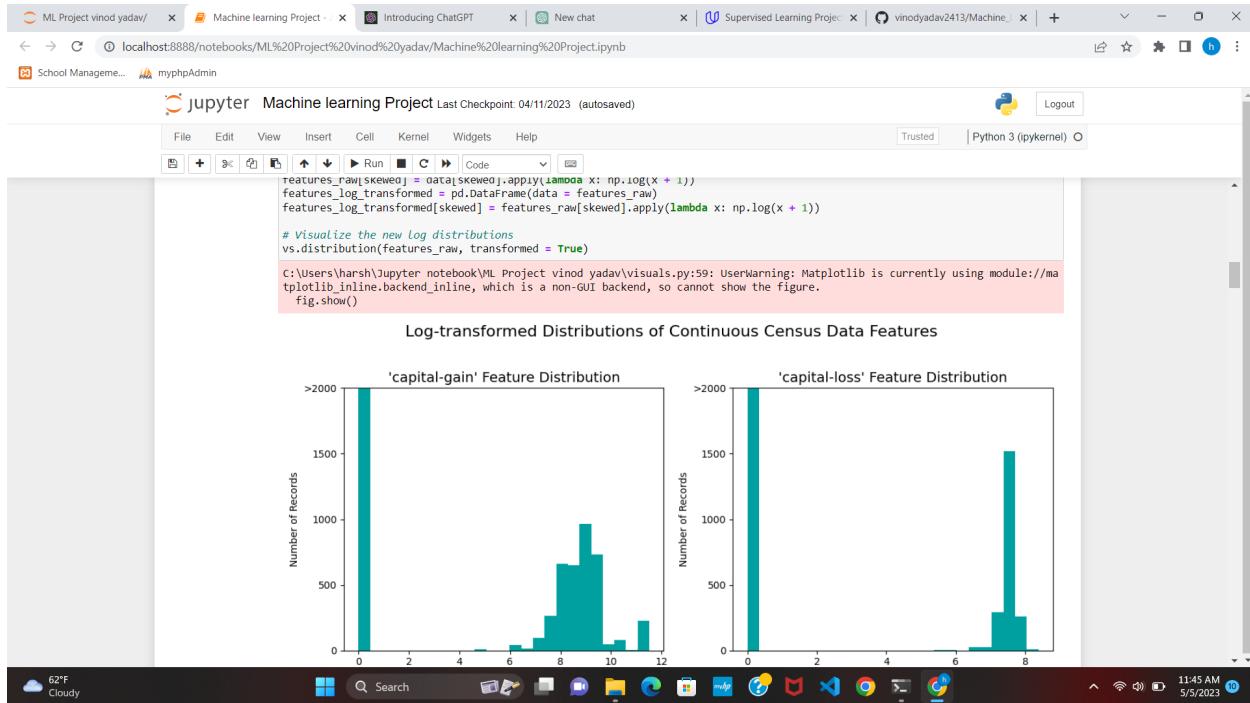
```
<=50K    34014  
>50K    11208  
Name: count, dtype: int64
```

In [3]: # Total number of records  
n\_records = 45222  
  
# Number of records where individual's income is more than \$50,000  
n\_greater\_50k = 11208  
  
# Number of records where individual's income is at most \$50,000  
n\_at\_most\_50k = 34014  
  
# Percentage of individuals whose income is more than \$50,000  
greater\_percent = 24.78  
  
# Print the results  
print("Total number of records: {}".format(n\_records))  
print("Individuals making more than \$50,000: {}".format(n\_greater\_50k))  
print("Individuals making at most \$50,000: {}".format(n\_at\_most\_50k))  
print("Percentage of individuals making more than \$50,000: {}%".format(greater\_percent))

Total number of records: 45222  
Individuals making more than \$50,000: 11208  
Individuals making at most \$50,000: 34014  
Percentage of individuals making more than \$50,000: 24.78%

In [4]: # Split the data into features and target label  
income\_raw = data['income']  
features\_raw = data.drop('income', axis = 1)  
  
# Visualize skewed continuous features of original data  
vs.distribution(data)





ML Project vinod yadav/ Machine learning Project - Introducing ChatGPT New chat Supervised Learning Project vinodyadav2413/Machine... +

localhost:8888/notebooks/ML%20Project%20vinod%20yadav/Machine%20learning%20Project.ipynb

School Manageme... myphAdmin

jupyter Machine learning Project Last Checkpoint 04/11/2023 (autosaved)

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3 (ipykernel) Logout

In [6]: # One-hot encode the 'features\_Log\_minmax\_transform' data using pandas.get\_dummies()  
`features\_final = pd.get\_dummies(features\_raw)  
# Encode the 'income\_raw' data to numerical values  
income = income\_raw.apply(lambda x: 0 if x == '<=50K' else 1)  
for x in income\_raw:  
if x == '<50K':  
income.append(0)  
else:  
income.append(1)  
income = pd.Series(income)  
# Print the number of features after one-hot encoding  
encoded = list(features\_final.columns)  
print("{} total features after one-hot encoding.".format(len(encoded)))  
# Uncomment the following line to see the encoded feature names  
print(encoded)

103 total features after one-hot encoding.  
['age', 'education-num', 'capital-loss', 'hours-per-week', 'workclass\_Federal-gov', 'workclass\_Local-gov', 'workclass\_Private', 'workclass\_Self-emp-inc', 'workclass\_Self-emp-not-inc', 'workclass\_State-gov', 'workclass\_Without-pay', 'education\_level\_10th', 'education\_level\_11th', 'education\_level\_12th', 'education\_level\_1st-4th', 'education\_level\_5th-9th', 'education\_level\_7th-8th', 'education\_level\_9th', 'education\_level\_Ascoc-acdm', 'education\_level\_Ascoc-vo', 'education\_level\_Bachelors', 'education\_level\_Dكتور', 'education\_level\_HS-grad', 'education\_level\_Masters', 'education\_level\_Preschool', 'education\_level\_Prof-school', 'education\_level\_Some-college', 'marital-status\_Divorced', 'marital-status\_Married-AF-spouse', 'marital-status\_Married-civ-spouse', 'marital-status\_Married-spouse-absent', 'marital-status\_Never-married', 'marital-status\_Separated', 'marital-status\_Widowed', 'occupation\_Adm-clerical', 'occupation\_Armed-forces', 'occupation\_Craft-repair', 'occupation\_Exec-managerial', 'occupation\_Farming-fishing', 'occupation\_Handlers-cleaners', 'occupation\_Machin-op-inspct', 'occupation\_Other-service', 'occupation\_Priv-house-serv', 'occupation\_Prof-specialty', 'occupation\_Protective-serv', 'occupation\_Sales', 'occupation\_Tech-support', 'occupation\_Transport-moving', 'relationship\_Husband', 'relationship\_Not-in-family', 'relationship\_other-relative', 'relationship\_Own-child', 'relationship\_Unmarried', 'relationship\_Wife', 'race\_Amer-Indian-Eskimo', 'race\_Asian-Pac-Islander', 'race\_Black', 'race\_Other', 'race\_White', 'sex\_Female', 'sex\_Male', 'native-country\_Cambodia', 'native-country\_China', 'native-country\_Columbia', 'native-country\_Cuba', 'native-country\_Dominican-Republic', 'native-country\_Ecuador', 'native-country\_El-Salvador', 'native-country\_England', 'native-country\_France', 'native-country\_Germany', 'native-country\_Greece', 'native-country\_Guatemala', 'native-country\_Haiti', 'native-country\_Holand-Netherlands', 'native-country\_Honduras', 'native-country\_Hong', 'native-country\_Hungary', 'native-country\_India', 'native-country\_Iran', 'native-country\_Ireland', 'native-country\_Italy', 'native-country\_Jamaica', 'native-country\_Japan', 'native-country\_Laos', 'native-country\_Mexico', 'native-country\_Nicaragua', 'native-country\_Outlying-US-Guam-USVI-etc', 'native-country\_Peru', 'native-country\_Philippines', 'native-country\_Poland', 'native-country\_Portugal', 'native-country\_Puerto-Rico', 'native-country\_Scotland', 'native-country\_South', 'native-country\_Taiwan', 'native-country\_Thailand', 'native-country\_Trinidad&Tobago', 'native-country\_United-States', 'native-country\_Vietnam', 'native-country\_Yugoslavia']

In [7]: # Import train\_test\_split  
from sklearn.model\_selection import train\_test\_split  
# split the 'features' and 'income' data into training and testing sets  
X\_train, X\_test, y\_train, y\_test = train\_test\_split(features\_final,  
income,  
test\_size = 0.2,  
random\_state = 0)  
# Show the results of the split  
print("Training set has {} samples.".format(X\_train.shape[0]))  
print("Testing set has {} samples.".format(X\_test.shape[0]))

Training set has 36177 samples.  
Testing set has 9045 samples.

In [8]: ...  
TP = np.sum(income) # Counting the ones as this is the naive case. Note that 'income' is the 'income\_raw' data  
encoded to numerical values done in the data preprocessing step.  
FP = income.count() - TP # Specific to the naive case  
TN = 0 # No predicted negatives in the naive case  
FN = 0 # No predicted positives in the naive case

Cloudy 11:46 AM 5/5/2023

ML Project vinod yadav/ Machine learning Project - Introducing ChatGPT New chat Supervised Learning Project vinodyadav2413/Machine... +

localhost:8888/notebooks/ML%20Project%20vinod%20yadav/Machine%20learning%20Project.ipynb

School Manageme... myphpAdmin

jupyter Machine learning Project Last Checkpoint 04/11/2023 (autosaved)

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3 (ipykernel) Logout

In [8]:

```
# split the 'features' and 'income' data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(features_final,
                                                    income,
                                                    test_size = 0.2,
                                                    random_state = 0)

# Show the results of the split
print("Training set has {} samples.".format(X_train.shape[0]))
print("Testing set has {} samples.".format(X_test.shape[0]))

Training set has 36177 samples.
Testing set has 9045 samples.
```

In [8]:

```
...
TP = np.sum(income) # Counting the ones as this is the naive case. Note that 'income' is the 'income_raw' data
encoded to numerical values done in the data preprocessing step.
FP = income.count() - TP # Specific to the naive case

TN = 0 # No predicted negatives in the naive case
FN = 0 # No predicted negatives in the naive case
...
# accuracy
accuracy = n_greater_50k/float(n_records)

# F-score using the formula above for beta = 0.5
recall = n_greater_50k/n_greater_50k
precision = n_greater_50k/float(n_records)
fscore = (1+0.5**2) * precision * recall/(0.5**2 * precision + recall)

# Print the results
print("Naive Predictor: [Accuracy score: {:.4f}, F-score: {:.4f}]".format(accuracy, fscore))

Naive Predictor: [Accuracy score: 0.2478, F-score: 0.2917]
```

In [10]:

```
# Import two metrics from sklearn - fbeta score and accuracy score
```

52°F Cloudy

ML Project vinod yadav/ Machine learning Project - Introducing ChatGPT New chat Supervised Learning Project vinodyadav2413/Machine... +

localhost:8888/notebooks/ML%20Project%20vinod%20yadav/Machine%20learning%20Project.ipynb

School Manageme... myphpAdmin

jupyter Machine learning Project Last Checkpoint 04/11/2023 (autosaved)

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3 (ipykernel) Logout

In [11]:

```
# Import the three supervised learning models from sklearn
from sklearn import tree
from sklearn.naive_bayes import GaussianNB
from sklearn.ensemble import RandomForestClassifier
from sklearn.linear_model import LogisticRegression

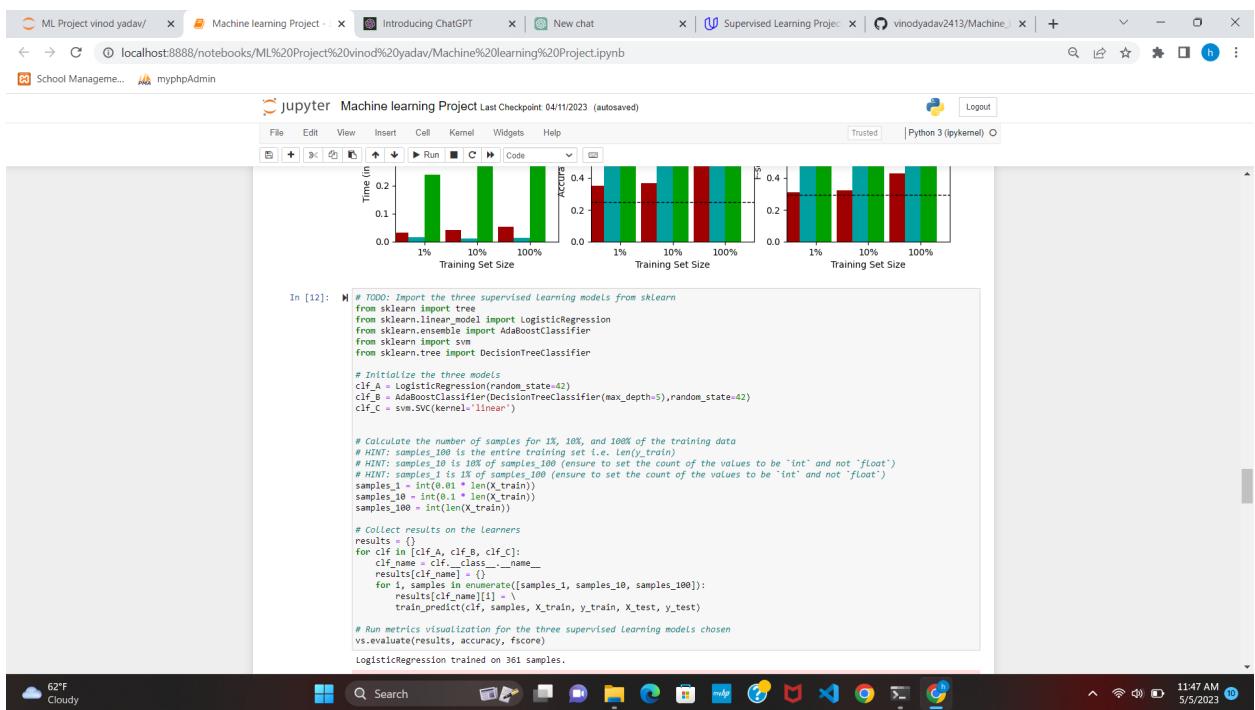
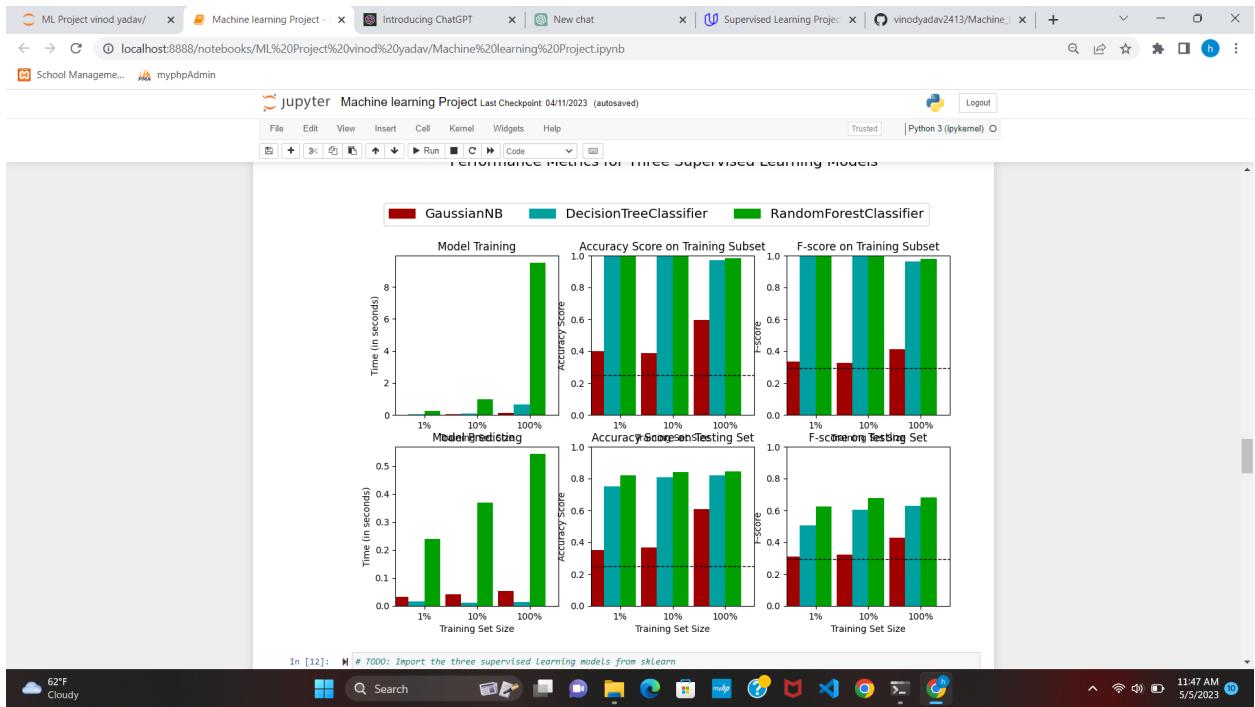
# Initialize the three models
clf_A = GaussianNB()
clf_B = tree.DecisionTreeClassifier(random_state=10)
clf_C = RandomForestClassifier(random_state=10)
#clf_D = LogisticRegression(random_state=42)
#clf_E = AdaBoostClassifier(DecisionTreeClassifier(max_depth=5), random_state=42)
#clf_F = ...

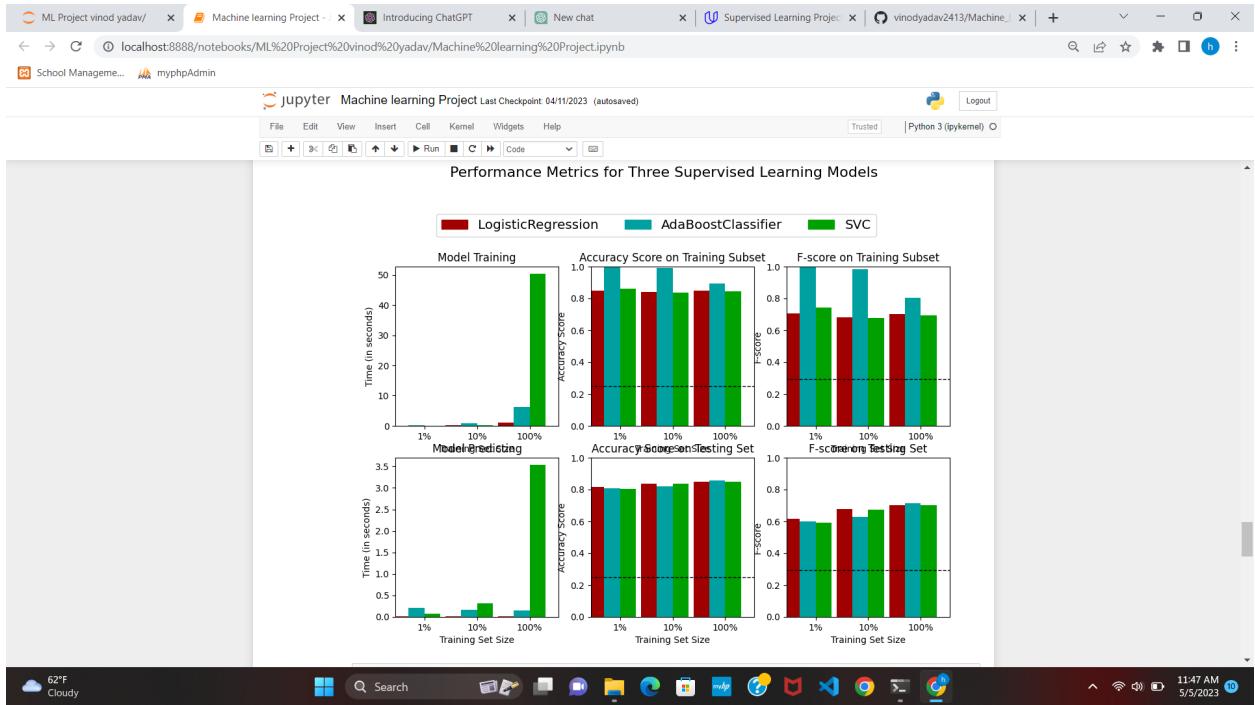
# Calculate the number of samples for 1%, 10%, and 100% of the training data
# HINT: samples_100 is the entire training set i.e. len(y_train)
# HINT: samples_10 is 10% of samples_100 (ensure to set the count of the values to be 'int' and not 'float')
# HINT: samples_1 is 1% of samples_100 (ensure to set the count of the values to be 'int' and not 'float')
samples_1 = int(0.01 * len(X_train))
samples_10 = int(0.1 * len(X_train))
samples_100 = int(len(X_train))

# Collect results on the learners
results = {}
for clf in [clf_A, clf_B, clf_C]:
    clf_name = clf.__class__.__name__
    results[clf_name] = {}
    for i, samples in enumerate([samples_1, samples_10, samples_100]):
        results[clf_name][i] = \
            train_predict(clf, samples, X_train, y_train, X_test, y_test)

# Run metrics visualization for the three supervised learning models chosen
vs.evaluate(results, accuracy, fscore)
```

52°F Cloudy





jupyter Machine learning Project Last Checkpoint: 04/11/2023 (autosaved)

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3 (ipykernel) Logout

```
In [13]: # Import 'GridSearchCV', 'make_scorer', and any other necessary libraries
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import make_scorer

# Initialize the classifier
clf = RandomForestClassifier(random_state = 10)

# Create the parameters list you wish to tune
parameters = {'max_depth': [5,20,2]}

# Make an fbeta_score scoring object
scorer = make_scorer(fbeta_score, beta=0.5)

# Perform grid search on the classifier using 'scorer' as the scoring method
grid_obj = GridSearchCV(clf, parameters)

# Fit the grid search object to the training data and find the optimal parameters
grid_fit = grid_obj.fit(X_train,y_train)

# Get the estimator
best_clf = grid_fit.best_estimator_

# Make predictions using the unoptimized and optimized model
predictions = (clf.fit(X_train, y_train)).predict(X_test)
best_predictions = best_clf.predict(X_test)

# Report the before-and-afterscores
print("Unoptimized model:-----")
print("Accuracy score on testing data: {:.4f}".format(accuracy_score(y_test, predictions)))
print("F-score on testing data: {:.4f}".format(fbeta_score(y_test, predictions, beta = 0.5)))
print("Optimized Model:-----")
print("Final accuracy score on the testing data: {:.4f}".format(accuracy_score(y_test, best_predictions)))
print("Final F-score on the testing data: {:.4f}".format(fbeta_score(y_test, best_predictions, beta = 0.5)))
```

Unoptimized model  
.....  
Accuracy score on testing data: 0.8421  
F-score on testing data: 0.6805  
Optimized Model  
.....  
Final accuracy score on the testing data: 0.8611  
Final F-score on the testing data: 0.7346

ML Project vinod yadav/ Machine learning Project - Introducing ChatGPT New chat Supervised Learning Project vinodyadav2413/Machine... + - ×

localhost:8888/notebooks/ML%20Project%20vinod%20yadav/Machine%20learning%20Project.ipynb

School Manageme... myphAdmin

jupyter Machine learning Project Last Checkpoint: 04/11/2023 (autosaved)

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3 (ipykernel) Logout

In [14]:

```
# Import a supervised learning model that has 'feature_importances_'
from sklearn.ensemble import GradientBoostingClassifier
# Train the supervised model on the training set using .fit(X_train, y_train)
model = GradientBoostingClassifier(n_estimators=100, learning_rate=1.0, max_depth=1, random_state=0).fit(X_train, y_train)

# Extract the feature importances using .feature_importances_
importances = model.feature_importances_

# Plot
vs.feature_plot(importances, X_train, y_train)
```

Normalized Weights for First Five Most Predictive Features

Feature	Feature Weight	Cumulative Feature Weight
marital-status_Married-civ-spouse	~0.42	~0.42
capital-gain	~0.20	~0.62
education-num	~0.15	~0.77
marital-status_Married-AF-spouse	~0.10	~0.87
capital-loss	~0.05	~0.92

52°F Cloudy Search 11:47 AM 5/5/2023

ML Project vinod yadav/ Machine learning Project - Introducing ChatGPT New chat Supervised Learning Project vinodyadav2413/Machine... + - ×

localhost:8888/notebooks/ML%20Project%20vinod%20yadav/Machine%20learning%20Project.ipynb

School Manageme... myphAdmin

jupyter Machine learning Project Last Checkpoint: 04/11/2023 (autosaved)

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3 (ipykernel) Logout

In [15]:

```
# Import functionality for cloning a model
from sklearn.base import clone

# Reduce the feature space
X_train_reduced = X_train[X_train.columns.values[(np.argsort(importances)[::-1])[:5]]]
X_test_reduced = X_test[X_test.columns.values[(np.argsort(importances)[::-1])[:5]]]

# Train on the "best" model found from grid search earlier
clf = (clone(best_clf)).fit(X_train_reduced, y_train)

# Make new predictions
reduced_predictions = clf.predict(X_test_reduced)

# Report scores from the final model using both versions of data
print("Final Model trained on full data")
print("Accuracy on testing data: {:.4f}".format(accuracy_score(y_test, best_predictions)))
print("F-score on testing data: {:.4f}".format(fbeta_score(y_test, best_predictions, beta = 0.5)))
print("\nFinal Model trained on reduced data")
print("Accuracy on testing data: {:.4f}".format(accuracy_score(y_test, reduced_predictions)))
print("F-score on testing data: {:.4f}".format(fbeta_score(y_test, reduced_predictions, beta = 0.5)))
```

Final Model trained on full data  
Accuracy on testing data: 0.8611  
F-score on testing data: 0.7346  
Final Model trained on reduced data  
Accuracy on testing data: 0.8568  
F-score on testing data: 0.7286

In [ ]:

52°F Cloudy Search 11:47 AM 5/5/2023