```
Supervised Learning
           Project: Finding Donors for CharityML
           Getting Started
           In this project, you will employ several supervised algorithms of your choice to accurately model individuals' income using data
           collected from the 1994 U.S. Census. You will then choose the best candidate algorithm from preliminary results and further
           optimize this algorithm to best model the data. Your goal with this implementation is to construct a model that accurately
           predicts whether an individual makes more than $50,000.
           Exploring the Data
 In [1]: # Import libraries necessary for this project
            import numpy as np
           import pandas as pd
            from time import time
            from IPython.display import display # Allows the use of display() for DataFrames
            # Import supplementary visualization code visuals.py
            import visuals as vs
            # Pretty display for notebooks
            %matplotlib inline
            # Load the Census dataset
            data = pd.read_csv(r"C:\Users\Carnival\Downloads\census (1).csv")
            # Success - Display the first record
            display(data.head(n=1))
                                                                                                                        hours-
                                               education- marital-
                                                                                                        capital-
                                                                                                               capital-
                                                                                                                                 nati
               age workclass education_level
                                                                   occupation relationship race sex
                                                                                                                           per-
                                                                                                          gain
                                                           status
                                                     num
                                                                                                                   loss
                                                                                                                                coun
                                                                                                                          week
                                                           Never-
                                                                                                                                 Unit
                                                                               Not-in-family White Male
                                                                                                                           40.0
               39 State-gov
                                     Bachelors
                                                     13.0
                                                                                                       2174.0
                                                           married
                                                                       clerical
           Implementation: Data Exploration
           A cursory investigation of the dataset will determine how many individuals fit into either group, and will tell us about the
           percentage of these individuals making more than \$50,000. In the code cell below, you will need to compute the following:

    The total number of records, 'n_records'

             • The number of individuals making more than \$50,000 annually, 'n_greater_50k'.
             • The number of individuals making at most \$50,000 annually, 'n_at_most_50k'.
             • The percentage of individuals making more than \$50,000 annually, 'greater_percent'.
 In [2]: data.describe()
 Out[2]:
                            age education-num
                                                 capital-gain
                                                               capital-loss hours-per-week
                                  45222.000000 45222.000000
             count 45222.000000
                                                             45222.000000
                                                                              45222.000000
                      38.547941
                                     10.118460 1101.430344
                                                                 88.595418
                                                                                40.938017
             mean
                      17.000000
                                      1.000000
                                                    0.000000
                                                                  0.000000
                                                                                 1.000000
              min
              25%
                      28.000000
                                       9.000000
                                                    0.000000
                                                                  0.000000
                                                                                 40.000000
                      37.000000
                                     10.000000
                                                    0.000000
                                                                  0.000000
                                                                                40.000000
              50%
                      47.000000
                                     13.000000
                                                    0.000000
                                                                 0.000000
                                                                                45.000000
              75%
                      90.000000
                                     16.000000 99999.000000
                                                              4356.000000
                                                                                99.000000
              max
 In [3]: data.income.value_counts()
 Out[3]: <=50K
                       34014
           >50K
                       11208
           Name: income, dtype: int64
 In [4]: # TODO: Total number of records
            n_records = data.shape[0]
            # TODO: Number of records where individual's income is more than $50,000
            n_greater_50k = data.query("income=='>50K'")['income'].count()
            # TODO: Number of records where individual's income is at most $50,000
            n_at_most_50k = data.query("income=='<=50K'")['income'].count()</pre>
            # TODO: Percentage of individuals whose income is more than $50,000
            greater_percent = (n_greater_50k / n_records)*100
            # 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: {:.3f}%".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: 24.784%
           Featureset Exploration
             • age: continuous.

    workclass: Private, Self-emp-not-inc, Self-emp-inc, Federal-gov, Local-gov, State-gov, Without-pay, Never-worked.

    education: Bachelors, Some-college, 11th, HS-grad, Prof-school, Assoc-acdm, Assoc-voc, 9th, 7th-8th, 12th, Masters,

                1st-4th, 10th, Doctorate, 5th-6th, Preschool.
              · education-num: continuous.

    marital-status: Married-civ-spouse, Divorced, Never-married, Separated, Widowed, Married-spouse-absent, Married-AF-

             • occupation: Tech-support, Craft-repair, Other-service, Sales, Exec-managerial, Prof-specialty, Handlers-cleaners,
                Machine-op-inspct, Adm-clerical, Farming-fishing, Transport-moving, Priv-house-serv, Protective-serv, Armed-Forces.

    relationship: Wife, Own-child, Husband, Not-in-family, Other-relative, Unmarried.

    race: Black, White, Asian-Pac-Islander, Amer-Indian-Eskimo, Other,

              • sex: Female, Male.

    capital-gain: continuous.

              • capital-loss: continuous.

    hours-per-week: continuous.

    native-country: United-States, Cambodia, England, Puerto-Rico, Canada, Germany, Outlying-US(Guam-USVI-etc),

                India, Japan, Greece, South, China, Cuba, Iran, Honduras, Philippines, Italy, Poland, Jamaica, Vietnam, Mexico,
                Portugal, Ireland, France, Dominican-Republic, Laos, Ecuador, Taiwan, Haiti, Columbia, Hungary, Guatemala, Nicaragua,
                Scotland, Thailand, Yugoslavia, El-Salvador, Trinadad&Tobago, Peru, Hong, Holand-Netherlands.
           Preparing the Data
           Before data can be used as input for machine learning algorithms, it often must be cleaned, formatted, and restructured — this
           is typically known as preprocessing. Fortunately, for this dataset, there are no invalid or missing entries we must deal with,
           however, there are some qualities about certain features that must be adjusted. This preprocessing can help tremendously
           with the outcome and predictive power of nearly all learning algorithms.
           Transforming Skewed Continuous Features
           A dataset may sometimes contain at least one feature whose values tend to lie near a single number, but will also have a non-
           trivial number of vastly larger or smaller values than that single number. Algorithms can be sensitive to such distributions of
           values and can underperform if the range is not properly normalized. With the census dataset two features fit this description:
            'capital-gain' and 'capital-loss'.
           # 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)
           C:\Users\Carnival\visuals.py:48: UserWarning: Matplotlib is currently using module://ipykerne
           l.pylab.backend_inline, which is a non-GUI backend, so cannot show the figure.
              fig.show()
                                    Skewed Distributions of Continuous Census Data Features
                            'capital-gain' Feature Distribution
                                                                                        'capital-loss' Feature Distribution
               >2000
                                                                           >2000
                1500
                                                                            1500
            of Records
                                                                         of Records
                1000
                                                                            1000
                 500
                                                                             500
                              20000
                                       40000
                                                60000
                                                         80000
                                                                 100000
                                                                                            1000
                                                                                                      2000
                                                                                                                           4000
           For highly-skewed feature distributions such as 'capital-gain' and 'capital-loss', it is common practice to apply
           a <u>logarithmic transformation</u> on the data so that the very large and very small values do not negatively affect the performance
           of a learning algorithm. Using a logarithmic transformation significantly reduces the range of values caused by outliers. Care
           must be taken when applying this transformation however: The logarithm of 0 is undefined, so we must translate the values
           by a small amount above 0 to apply the the logarithm successfully.
 In [6]: # Log-transform the skewed features
            skewed = ['capital-gain', 'capital-loss']
            features log transformed = pd.DataFrame(data = features_raw)
            features_log_transformed[skewed] = features_raw[skewed].apply(lambda x: np.log(x + 1))
            # Visualize the new log distributions
            vs.distribution(features_log_transformed, transformed = True)
                               Log-transformed Distributions of Continuous Census Data Features
                            'capital-gain' Feature Distribution
                                                                                        'capital-loss' Feature Distribution
               >2000
                                                                           >2000
                1500
                                                                            1500
            of Records
                                                                         oţ
                                                                            1000
                1000
                 500
                                                                             500
           Normalizing Numerical Features
           In addition to performing transformations on features that are highly skewed, it is often good practice to perform some type of
           scaling on numerical features. Applying a scaling to the data does not change the shape of each feature's distribution (such as
            'capital-gain' or 'capital-loss' above).
 In [7]: # Import sklearn.preprocessing.StandardScaler
            from sklearn.preprocessing import MinMaxScaler
            # 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_log_minmax_transform = pd.DataFrame(data = features_log_transformed)
            features_log_minmax_transform[numerical] = scaler.fit_transform(features_log_transformed[num
            erical])
            # Show an example of a record with scaling applied
            display(features_log_minmax_transform.head(n = 5))
                                                    education-
                                                                                                                capital- capital-
                        workclass education level
                                                                        occupation relationship
                                                                                                race
                                                                 status
                                                                                                                            loss
                                                                                                                   gain
                                                         num
                                                                 Never-
            0 0.301370
                          State-gov
                                                     0.800000
                                                                                                               0.667492
                                                                                                                             0.0 0.39
                                          Bachelors
                                                                                    Not-in-family White
                                                                                                          Male
                                                                married
                                                                            clerical
                                                                Married-
                          Self-emp-
                                                                             Exec
            1 0.452055
                                          Bachelors
                                                     0.800000
                                                                                        Husband White
                                                                                                         Male
                                                                                                               0.000000
                                                                                                                             0.0 0.12
                                                                   civ-
                                                                         managerial
                                                                 spouse
                                                                          Handlers-
            2 0.287671
                                           HS-grad
                                                     0.533333
                                                               Divorced
                                                                                    Not-in-family White
                                                                                                               0.000000
                                                                                                                             0.0 0.39
                                                                           cleaners
                                                                Married-
                                                                          Handlers-
            3 0.493151
                                                     0.400000
                                                                                                         Male 0.000000
                                                                                                                             0.0 0.39
                            Private
                                              11th
                                                                   civ-
                                                                                        Husband Black
                                                                           cleaners
                                                                 spouse
                                                                Married-
                                                                              Prof-
            4 0.150685
                            Private
                                          Bachelors
                                                     0.800000
                                                                                           Wife Black Female 0.000000
                                                                                                                             0.0 0.39
                                                                   civ-
                                                                           specialty
                                                                 spouse
           Implementation: Data Preprocessing
           From the table in Exploring the Data above, we can see there are several features for each record that are non-numeric.
           Typically, learning algorithms expect input to be numeric, which requires that non-numeric features (called categorical
            variables) be converted. One popular way to convert categorical variables is by using the one-hot encoding scheme. One-
            hot encoding creates a "dummy" variable for each possible category of each non-numeric feature. For example, assume
            someFeature has three possible entries: A , B , or C . We then encode this feature into someFeature_A ,
            someFeature_B and someFeature_C.
                                                                    someFeature_A someFeature_B someFeature C
                                someFeature
                                              ----> one-hot encode ---->
           Additionally, as with the non-numeric features, we need to convert the non-numeric target label, 'income' to numerical
           values for the learning algorithm to work. Since there are only two possible categories for this label ("<=50K" and ">50K"), we
           can avoid using one-hot encoding and simply encode these two categories as 0 and 1, respectively. In code cell below,
           you will need to implement the following:
             • Use pandas.get_dummies() to perform one-hot encoding on the 'features_log_minmax_transform' data.

    Convert the target label 'income_raw' to numerical entries.

                  Set records with "<=50K" to 0 and records with ">50K" to 1.
 In [8]: # TODO: One-hot encode the 'features_log_minmax_transform' data using pandas.get_dummies()
            features_final = pd.get_dummies(data=features_log_minmax_transform, columns=['workclass','ed
            ucation_level', 'marital-status', 'occupation', 'relationship', 'race', 'sex', 'native-country'])
            # TODO: Encode the 'income raw' data to numerical values
            income = income_raw.replace({'<=50K':0, '>50K':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)))
            # Uncomment the following line to see the encoded feature names
            print (encoded)
           103 total features after one-hot encoding.
            ['age', 'education-num', 'capital-gain', 'capital-loss', 'hours-per-week', 'workclass_ Federa
           l-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_leve
           l_ 5th-6th', 'education_level_ 7th-8th', 'education_level_ 9th', 'education_level_ Assoc-acd
m', 'education_level_ Assoc-voc', 'education_level_ Bachelors', 'education_level_ Doctorate',
            'education_level_ HS-grad', 'education_level_ Masters', 'education_level_ Preschool', 'educat
           ion_level_ Prof-school', 'education_level_ Some-college', 'marital-status_ Divorced', 'marita
           l-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-repai
           r', 'occupation_ Exec-managerial', 'occupation_ Farming-fishing', 'occupation_ Handlers-clean
           ers', 'occupation_ Machine-op-inspct', 'occupation_ Other-service', 'occupation_ Priv-house-s
           erv', 'occupation_ Prof-specialty', 'occupation_ Protective-serv', 'occupation_ Sales', 'occu
           pation_ Tech-support', 'occupation_ Transport-moving', 'relationship_ Husband', 'relationship
            _ Not-in-family', 'relationship_ Other-relative', 'relationship_ Own-child', 'relationship_ U
           nmarried', 'relationship_ Wife', 'race_ Amer-Indian-Eskimo', 'race_ Asian-Pac-Islander', 'rac
           e_ Black', 'race_ Other', 'race_ White', 'sex_ Female', 'sex_ Male', 'native-country_ Cambodi
           a', 'native-country_ Canada', 'native-country_ China', 'native-country_ Columbia', 'native-co
           untry_ Cuba', 'native-country_ Dominican-Republic', 'native-country_ Ecuador', 'native-countr
           y_ El-Salvador', 'native-country_ England', 'native-country_ France', 'native-country_ German
           y', 'native-country_ Greece', 'native-country_ Guatemala', 'native-country_ Haiti', 'native-c
           ountry_ Holand-Netherlands', 'native-country_ Honduras', 'native-country_ Hong', 'native-coun
           try_ Hungary', 'native-country_ India', 'native-country_ Iran', 'native-country_ Ireland', 'n
           ative-country_ Italy', 'native-country_ Jamaica', 'native-country_ Japan', 'native-country_ L
           aos', 'native-country_ Mexico', 'native-country_ Nicaragua', 'native-country_ Outlying-US(Gua
           m-USVI-etc)', 'native-country_ Peru', 'native-country_ Philippines', 'native-country_ Polan
           d', 'native-country_ Portugal', 'native-country_ Puerto-Rico', 'native-country_ Scotland', 'n
           ative-country_ South', 'native-country_ Taiwan', 'native-country_ Thailand', 'native-country_
           Trinadad&Tobago', 'native-country_ United-States', 'native-country_ Vietnam', 'native-country
            _ Yugoslavia']
           Shuffle and Split Data
           Now all categorical variables have been converted into numerical features, and all numerical features have been normalized.
           As always, we will now split the data (both features and their labels) into training and test sets. 80% of the data will be used for
           training and 20% for testing.
           Run the code cell below to perform this split.
 In [9]: # 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 = 2)
            # 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.
           Evaluating Model Performance
           In this section, we will investigate four different algorithms, and determine which is best at modeling the data. Three of these
           algorithms will be supervised learners of your choice, and the fourth algorithm is known as a naive predictor.
           Metrics and the Naive Predictor
           CharityML, equipped with their research, knows individuals that make more than \$50,000 are most likely to donate to their
           charity. Because of this, *CharityML* is particularly interested in predicting who makes more than \$50,000 accurately. It would
           seem that using accuracy as a metric for evaluating a particular model's performace would be appropriate. Additionally,
           identifying someone that does not make more than \$50,000 as someone who does would be detrimental to *CharityML*,
           since they are looking to find individuals willing to donate. Therefore, a model's ability to precisely predict those that make
           more than \$50,000 is more important than the model's ability to recall those individuals. We can use F-beta score as a
           metric that considers both precision and recall:
           F \left( \frac{1 + \beta^2 \cdot \frac{1}{\hbar}}{\left( \frac{1 + \beta^2 \cdot \frac{1}{\hbar}}{\left(
           In particular, when \theta = 0.5, more emphasis is placed on precision. This is called the F_{0.5} score (or F-score for
           simplicity).
           Looking at the distribution of classes (those who make at most \$50,000, and those who make more), it's clear most
           individuals do not make more than \$50,000. This can greatly affect accuracy, since we could simply say "this person does
           not make more than \$50,000" and generally be right, without ever looking at the data! Making such a statement would be
           called naive, since we have not considered any information to substantiate the claim. It is always important to consider the
           naive prediction for your data, to help establish a benchmark for whether a model is performing well. That been said, using
           that prediction would be pointless: If we predicted all people made less than \$50,000, CharityML would identify no one as
           donors.
           Accuracy measures how often the classifier makes the correct prediction. It's the ratio of the number of correct predictions to
           the total number of predictions (the number of test data points).
           Precision tells us what proportion of messages we classified as spam, actually were spam. It is a ratio of true positives (words
           classified as spam, and which are actually spam) to all positives(all words classified as spam, irrespective of whether that was
           the correct classificatio), in other words it is the ratio of
            [True Positives/(True Positives + False Positives)]
           Recall(sensitivity) tells us what proportion of messages that actually were spam were classified by us as spam. It is a ratio of
           true positives(words classified as spam, and which are actually spam) to all the words that were actually spam, in other words
           it is the ratio of
            [True Positives/(True Positives + False Negatives)]
           For classification problems that are skewed in their classification distributions like in our case, for example if we had a 100 text
           messages and only 2 were spam and the rest 98 weren't, accuracy by itself is not a very good metric. We could classify 90
            messages as not spam(including the 2 that were spam but we classify them as not spam, hence they would be false
           negatives) and 10 as spam(all 10 false positives) and still get a reasonably good accuracy score. For such cases, precision
           and recall come in very handy. These two metrics can be combined to get the F1 score, which is weighted average(harmonic
           mean) of the precision and recall scores. This score can range from 0 to 1, with 1 being the best possible F1 score(we take
           the harmonic mean as we are dealing with ratios).
           Naive Predictor Performace
             • If we chose a model that always predicted an individual made more than $50,000, what would that model's accuracy and
                F-score be on this dataset? You must use the code cell below and assign your results to 'accuracy' and 'fscore'
                to be used later.
In [10]: '''
            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
            # TODO: Calculate accuracy, precision and recall
            accuracy = (np.sum(income)+0)/(np.sum(income)+income.count() - np.sum(income)+0+0)
            recall = (np.sum(income))/ (np.sum(income)+0)
            precision = (np.sum(income))/ (np.sum(income)+income.count() - np.sum(income))
            # TODO: Calculate F-score using the formula above for beta = 0.5 and correct values for prec
            ision and recall.
            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]
           Supervised Learning Models
           The following are some of the supervised learning models that are currently available in <u>scikit-learn</u> that you
            may choose from:

    Gaussian Naive Bayes (GaussianNB)

    Decision Trees

             • Ensemble Methods (Bagging, AdaBoost, Random Forest, Gradient Boosting)
             • K-Nearest Neighbors (KNeighbors)

    Stochastic Gradient Descent Classifier (SGDC)

    Support Vector Machines (SVM)

              · Logistic Regression
           Implementation - Creating a Training and Predicting Pipeline
           To properly evaluate the performance of each model you've chosen, it's important that you create a training and predicting
           pipeline that allows you to quickly and effectively train models using various sizes of training data and perform predictions on
           the testing data.
In [11]: # TODO: Import two metrics from sklearn - fbeta_score and accuracy_score
            from sklearn.metrics import fbeta_score
            from sklearn.metrics import accuracy_score
            from sklearn import svm
            from sklearn.tree import DecisionTreeClassifier
            from sklearn.linear_model import LogisticRegression
            def train_predict(learner, sample_size, X_train, y_train, X_test, y_test):
                 inputs:
                     - learner: the learning algorithm to be trained and predicted on
                     - sample size: the size of samples (number) to be drawn from training set
                     - X_train: features training set
                    - y_train: income training set
                     - X_test: features testing set
                    - y_test: income testing set
                 results = \{\}
                 # TODO: Fit the learner to the training data using slicing with 'sample_size' using .fit
            (training_features[:], training_labels[:])
                 start = time() # Get start time
                 learner = learner.fit( X_train[:sample_size], y_train[:sample_size] )
                 end = time() # Get end time
                 # TODO: Calculate the training time
                 results['train_time'] = end-start
                 # TODO: Get the predictions on the test set(X_test),
                           then get predictions on the first 300 training samples(X_train) using .predict()
                 start = time() # Get start time
                 predictions_test = learner.predict(X_test)
                 predictions_train = learner.predict(X_train[:300])
                 end = time() # Get end time
                 # TODO: Calculate the total prediction time
                 results['pred_time'] = end-start
                 # TODO: Compute accuracy on the first 300 training samples which is y_train[:300]
                 results['acc_train'] = accuracy_score(y_train[:300], predictions_train[:300])
                 # TODO: Compute accuracy on test set using accuracy_score()
                 results['acc_test'] = accuracy_score(y_test, predictions_test)
                 # TODO: Compute F-score on the the first 300 training samples using fbeta_score()
                 results['f_train'] = fbeta_score(y_train[:300], predictions_train[:300], average=None, b
            eta=0.5)
                 # TODO: Compute F-score on the test set which is y_test
                 results['f_test'] = fbeta_score(y_test, predictions_test, average=None, beta=0.5)
                 # Success
                 print("{} trained on {} samples.".format(learner.__class__.__name__, sample_size))
                 # Return the results
                 return results
           Implementation: Initial Model Evaluation
In [12]: # TODO: Import the three supervised learning models from sklearn
            from sklearn.neighbors import KNeighborsClassifier
            # TODO: Initialize the three models
            clf_A = LogisticRegression(random_state=1)
           clf_B = svm.SVC(random_state=1)
           clf_C = KNeighborsClassifier()
            # TODO: 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` an
            d not `float`)
            samples_100 = int(len(y_train))
            samples_10 = int(len(y_train)*0.1)
            samples_1 = int(len(y_train)*0.01)
            # 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)
           LogisticRegression trained on 361 samples.
           LogisticRegression trained on 3617 samples.
           C:\Users\Carnival\anaconda3\lib\site-packages\sklearn\linear_model\_logistic.py:762: Converge
           nceWarning: lbfgs failed to converge (status=1):
           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
              n iter i = check optimize result(
           C:\Users\Carnival\anaconda3\lib\site-packages\sklearn\linear_model\_logistic.py:762: Converge
           nceWarning: lbfgs failed to converge (status=1):
           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
              n_iter_i = _check_optimize_result(
           LogisticRegression trained on 36177 samples.
           SVC trained on 361 samples.
           SVC trained on 3617 samples.
           SVC trained on 36177 samples.
           KNeighborsClassifier trained on 361 samples.
           KNeighborsClassifier trained on 3617 samples.
           KNeighborsClassifier trained on 36177 samples.
           C:\Users\Carnival\visuals.py:118: UserWarning: Tight layout not applied. tight_layout cannot
           make axes width small enough to accommodate all axes decorations
              pl.tight_layout()
                            Performance Metrics for Three Supervised Learning Models
                           LogisticRegression
                                                           SVC
                                                                          KNeighborsClassifier
                          Model Training
                                                                                      F-score on Training Subset
                                                  Accuracy Score on Training Subset
               100
                80
                60
                                                 0.6
                                                                                  0.6
                40
                                                 0.4
                20
                                                 0.2
                                                                                  0.2
                 0
                               10%
                                                                10%
                                                                         100%
                         MosieinBrestistieg
                                                                                         F-scomenong Testsing Set
                                                    AccuracyTraining Set Stesting Set
                                                1.0
                20
                15
                                                 0.6
              € 10
                                                 0.4
                                                 0.2
                 0
                               10%
                                       100%
                                                                10%
                                                                        100%
                      1%
                                                        1%
                                                                                         1%
                                                                                                 10%
                          Training Set Size
                                                            Training Set Size
                                                                                             Training Set Size
           Improving Results
           In this final section, you will choose from the three supervised learning models the best model to use on the student data. You
           will then perform a grid search optimization for the model over the entire training set ( X_train and y_train ) by tuning at
           least one parameter to improve upon the untuned model's F-score.
           Implementation: Model Tuning
In [13]: # TODO: Import 'GridSearchCV', 'make_scorer', and any other necessary libraries
            from sklearn.model_selection import GridSearchCV
            from sklearn.metrics import make_scorer
            # TODO: Initialize the classifier
            clfc = LogisticRegression(random_state=1)
            # TODO: Create the parameters list you wish to tune, using a dictionary if needed.
            # HINT: parameters = {'parameter_1': [value1, value2], 'parameter_2': [value1, value2]}
           parameters = {'penalty':['11','12'], 'C':[1,2,3,4,5,6]}
            # TODO: Make an fbeta_score scoring object using make_scorer()
            scorer = make_scorer(fbeta_score, beta = 0.5)
            # TODO: Perform grid search on the classifier using 'scorer' as the scoring method using Gri
            grid_obj = GridSearchCV(estimator=clfc, param_grid=parameters, scoring=scorer, n_jobs = -1)
            # TODO: Fit the grid search object to the training data and find the optimal parameters usin
            g fit()
            grid_fit = grid_obj.fit(X_train, y_train)
            # Get the estimator
            best_clf = grid_fit.best_estimator_
            # Make predictions using the unoptimized and 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\n----")
           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("\nOptimized Model\n----")
            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_prediction
            s, beta = 0.5)))
           C:\Users\Carnival\anaconda3\lib\site-packages\sklearn\linear_model\_logistic.py:762: Converge
           nceWarning: lbfgs failed to converge (status=1):
           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
              n_iter_i = _check_optimize_result(
           Unoptimized model
           Accuracy score on testing data: 0.8188
           F-score on testing data: 0.6344
           Optimized Model
           Final accuracy score on the testing data: 0.8422
           Final F-score on the testing data: 0.6902
           Final Model Evaluation
           Results:
                                                              Unoptimized Model Optimized Model
                                                   Metric
                                                                                      0.8423
                                               Accuracy Score
                                                                    0.8188
                                                                    0.6344
                                                                                     0.6907
                                                   F-score
           Feature Importance
           Feature Relevance Observation
           Implementation - Extracting Feature Importance
In [14]: # TODO: Import a supervised learning model that has 'feature_importances_'
            from sklearn.ensemble import RandomForestClassifier
            # TODO: Train the supervised model on the training set using .fit(X_train, y_train)
            model = RandomForestClassifier().fit(X_train, y_train)
            # TODO: Extract the feature importances using .feature_importances_
            importances = model.feature_importances_
            vs.feature_plot(importances, X_train, y_train)
                            Normalized Weights for First Five Most Predictive Features
               0.6

    Feature Weight
```

0.1 0.0 hours-per-week capital-gain marital-status Married-civ-spouseducation-num Feature **Extracting Feature Importance Feature Selection** 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\n----") 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\n-----") 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: 0.8422 F-score on testing data: 0.6902 Final Model trained on reduced data

Cumulative Feature Weight

0.5

0.4

Weight 6.0

0.2

Accuracy on testing data: 0.8265 F-score on testing data: 0.6542