

# Machine Learning Engineer Nanodegree

## Supervised Learning

### Project: Finding Donors for *CharityML*

Welcome to the second project of the Machine Learning Engineer Nanodegree! In this notebook, some template code has already been provided for you, and it will be your job to implement the additional functionality necessary to successfully complete this project. Sections that begin with **'Implementation'** in the header indicate that the following block of code will require additional functionality which you must provide. Instructions will be provided for each section and the specifics of the implementation are marked in the code block with a `'TODO'` statement. Please be sure to read the instructions carefully!

In addition to implementing code, there will be questions that you must answer which relate to the project and your implementation. Each section where you will answer a question is preceded by a **'Question X'** header. Carefully read each question and provide thorough answers in the following text boxes that begin with **'Answer:'**. Your project submission will be evaluated based on your answers to each of the questions and the implementation you provide.

**Note:** Please specify WHICH VERSION OF PYTHON you are using when submitting this notebook. Code and Markdown cells can be executed using the **Shift + Enter** keyboard shortcut. In addition, Markdown cells can be edited by typically double-clicking the cell to enter edit mode.

## 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. This sort of task can arise in a non-profit setting, where organizations survive on donations. Understanding an individual's income can help a non-profit better understand how large of a donation to request, or whether or not they should reach out to begin with. While it can be difficult to determine an individual's general income bracket directly from public sources, we can (as we will see) infer this value from other publically available features.

The dataset for this project originates from the [UCI Machine Learning Repository](https://archive.ics.uci.edu/ml/datasets/Census+Income) (<https://archive.ics.uci.edu/ml/datasets/Census+Income>). The dataset was donated by Ron Kohavi and Barry Becker, after being published in the article "*Scaling Up the Accuracy of Naive-Bayes Classifiers: A Decision-Tree Hybrid*". You can find the article by Ron Kohavi online (<https://www.aaai.org/Papers/KDD/1996/KDD96-033.pdf>). The data we investigate here consists of small changes to the original dataset, such as removing the 'fnlwgt' feature and records with missing or ill-formatted entries.

---

## Exploring the Data

Run the code cell below to load necessary Python libraries and load the census data. Note that the last column from this dataset, 'income', will be our target label (whether an individual makes more than, or at most, \$50,000 annually). All other columns are features about each individual in the census database.

```
In [105]: # 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 1

# Import supplementary visualization code visuals.py
import visuals as vs

# Pretty display for notebooks
%matplotlib inline

# Load the Census dataset
data = pd.read_csv("census.csv")

# Success - Display the first record
```

	age	workclass	education_level	education-num	marital-status	occupation	relationship	race
0	39	State-gov	Bachelors	13.0	Never-married	Adm-clerical	Not-in-family	White

## 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'.

**HINT:** You may need to look at the table above to understand how the 'income' entries are formatted.

```
In [106]: # 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.groupby('income').size()[1]

# TODO: Number of records where individual's income is at most $50,000
n_at_most_50k = data.groupby('income').size()[0]

# 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))
```

```
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: 24.7843969749237
1%
```

## 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-spouse.
- **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, Trinidad&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.

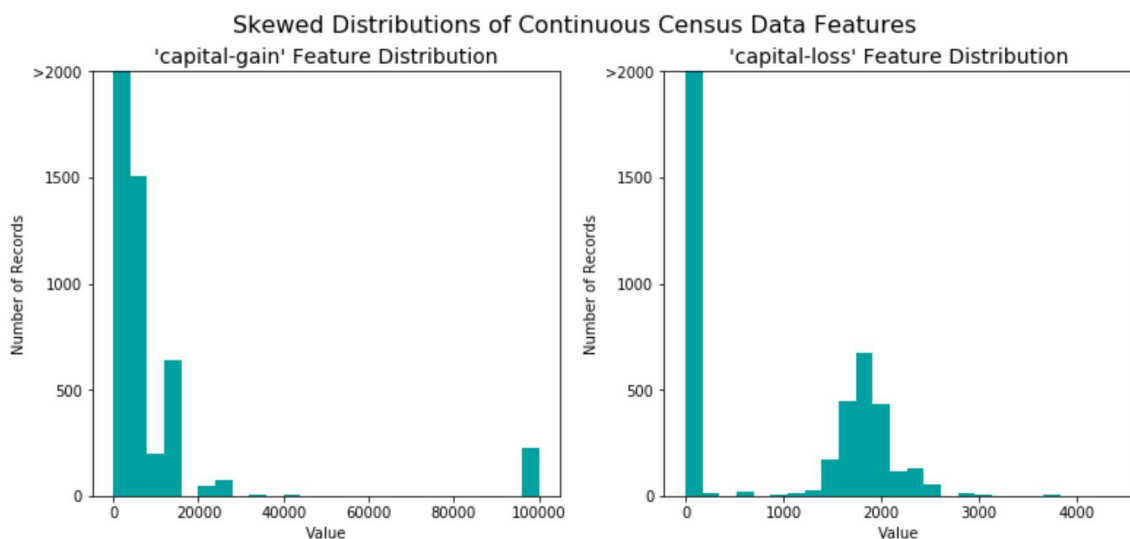
### Transforming Skewed Continuous Features

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'.

Run the code cell below to plot a histogram of these two features. Note the range of the values present and how they are distributed.

```
In [107]: # Split the data into features and target label
income_raw = data['income']
features_raw = data.drop('income', axis = 1)

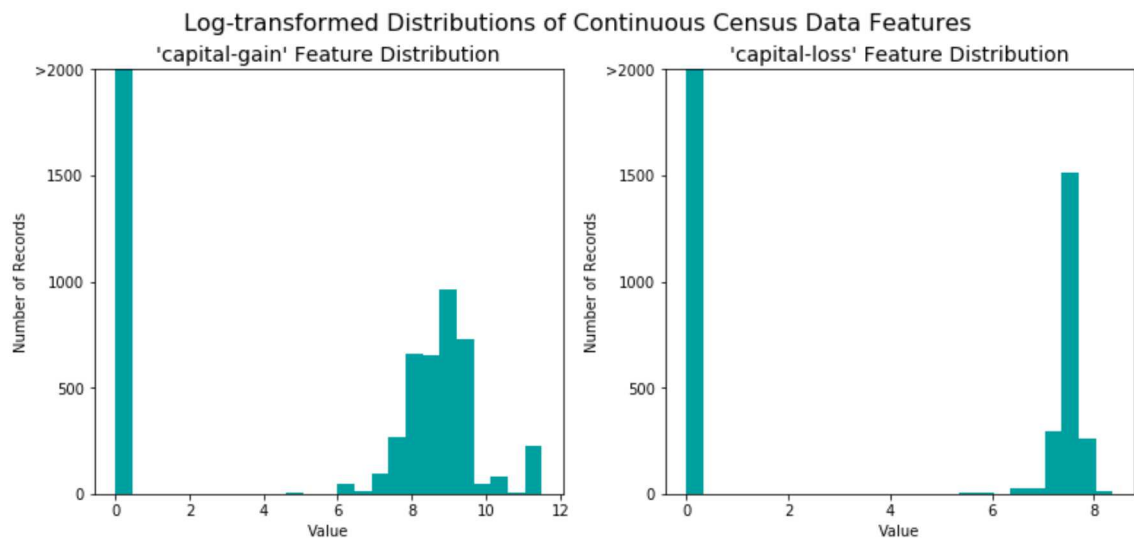
# Visualize skewed continuous features of original data
```



For highly-skewed feature distributions such as 'capital-gain' and 'capital-loss', it is common practice to apply a logarithmic transformation ([https://en.wikipedia.org/wiki/Data\\_transformation\\_\(statistics\)](https://en.wikipedia.org/wiki/Data_transformation_(statistics))) 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.

Run the code cell below to perform a transformation on the data and visualize the results. Again, note the range of values and how they are distributed.

```
In [108]: # 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:
# Visualize the new log distributions
```



## 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); however, normalization ensures that each feature is treated equally when applying supervised learners. Note that once scaling is applied, observing the data in its raw form will no longer have the same original meaning, as exemplified below.

Run the code cell below to normalize each numerical feature. We will use `sklearn.preprocessing.MinMaxScaler` (<http://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.MinMaxScaler.html>) for this.

```
In [109]: # 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', 'l

          features_log_minmax_transform = pd.DataFrame(data = features_log_transf
          features_log_minmax_transform[numerical] = scaler.fit_transform(feature

          # Show an example of a record with scaling applied
```

	age	workclass	education_level	education-num	marital-status	occupation	relationship
0	0.301370	State-gov	Bachelors	0.800000	Never-married	Adm-clerical	Not-in-family
1	0.452055	Self-emp-not-inc	Bachelors	0.800000	Married-civ-spouse	Exec-managerial	Husband
2	0.287671	Private	HS-grad	0.533333	Divorced	Handlers-cleaners	Not-in-family
3	0.493151	Private	11th	0.400000	Married-civ-spouse	Handlers-cleaners	Husband
4	0.150685	Private	Bachelors	0.800000	Married-civ-spouse	Prof-specialty	Wife

## 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`.

	<b>someFeature</b>		<b>someFeature_A</b>	<b>someFeature_B</b>	<b>someFeature_C</b>
0	B		0	1	0
1	C	----> one-hot encode ---->	0	0	1
2	A		1	0	0

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 (" $\leq 50K$ " 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()` ([http://pandas.pydata.org/pandas-docs/stable/generated/pandas.get\\_dummies.html?highlight=get\\_dummies#pandas.get\\_dummies](http://pandas.pydata.org/pandas-docs/stable/generated/pandas.get_dummies.html?highlight=get_dummies#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 " $\leq 50K$ " to 0 and records with " $> 50K$ " to 1.



```
In [110]: # TODO: One-hot encode the 'features_log_minmax_transform' data using
features_final = pd.get_dummies(features_log_minmax_transform)

# TODO: Encode the 'income_raw' data to numerical values
income = (income_raw.replace('<=50K', 0)).replace('>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

103 total features after one-hot encoding.
['age', 'education-num', 'capital-gain', '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-6th', 'education_level_ 7th-8
th', 'education_level_ 9th', 'education_level_ Assoc-acdm', 'educati
on_level_ Assoc-voc', 'education_level_ Bachelors', 'education_level_
Doctorate', 'education_level_ HS-grad', 'education_level_ Masters'
, 'education_level_ Preschool', 'education_level_ Prof-school', 'edu
cation_level_ Some-college', 'marital-status_ Divorced', 'marital-st
atus_ Married-AF-spouse', 'marital-status_ Married-civ-spouse', 'mar
ital-status_ Married-spouse-absent', 'marital-status_ Never-married'
, 'marital-status_ Separated', 'marital-status_ Widowed', 'occupatio
n_ Adm-clerical', 'occupation_ Armed-Forces', 'occupation_ Craft-rep
air', 'occupation_ Exec-managerial', 'occupation_ Farming-fishing',
'occupation_ Handlers-cleaners', 'occupation_ Machine-op-inspct', 'o
ccupation_ Other-service', 'occupation_ Priv-house-serv', 'occupatio
n_ Prof-specialty', 'occupation_ Protective-serv', 'occupation_ Sale
```

## 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 [111]: # 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]))
```

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 performance 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:

$$F_{\beta} = (1 + \beta^2) \cdot \frac{\text{precision} \cdot \text{recall}}{(\beta^2 \cdot \text{precision}) + \text{recall}}$$

In particular, when  $\beta = 0.5$ , more emphasis is placed on precision. This is called the **F<sub>0.5</sub> 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.

### Note: Recap of accuracy, precision, recall

**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

$$[\text{True Positives} / (\text{True Positives} + \text{False Positives})]$$

**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

$$[\text{True Positives} / (\text{True Positives} + \text{False Negatives})]$$

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

## Question 1 - Naive Predictor Performance

- 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.

**Please note** that the the purpose of generating a naive predictor is simply to show what a base model without any intelligence would look like. In the real world, ideally your base model would be either the results of a previous model or could be based on a research paper upon which you are looking to improve. When there is no benchmark model set, getting a result better than random choice is a place you could start from.

### HINT:

- When we have a model that always predicts '1' (i.e. the individual makes more than 50k) then our model will have no True Negatives(TN) or False Negatives(FN) as we are not making any negative('0' value) predictions. Therefore our Accuracy in this case becomes the same as our Precision(True Positives/(True Positives + False Positives)) as every prediction that we have made with value '1' that should have '0' becomes a False Positive; therefore our denominator in this case is the total number of records we have in total.
- Our Recall score(True Positives/(True Positives + False Negatives)) in this setting becomes 1 as we have no False Negatives.

```
In [112]: TP = np.sum(income) # Counting the ones as this is the naive case. Note
#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 = (TP + TN)/n_records
recall = TP/(TP + FN)
precision = TP/(TP + FP)

# TODO: Calculate F-score using the formula above for beta = 0.5 and c
fscore = (1 + 0.25)*(precision*recall)/(0.25*precision + recall)

# Print the results
```

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 scikit-learn ([http://scikit-learn.org/stable/supervised\\_learning.html](http://scikit-learn.org/stable/supervised_learning.html)) 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

## Question 2 - Model Application

List three of the supervised learning models above that are appropriate for this problem that you will test on the census data. For each model chosen

- Describe one real-world application in industry where the model can be applied.
- What are the strengths of the model; when does it perform well?
- What are the weaknesses of the model; when does it perform poorly?
- What makes this model a good candidate for the problem, given what you know about the data?

### HINT:

Structure your answer in the same format as above^, with 4 parts for each of the three models you pick. Please include references with your answer.

**Answer:**

- Decision Trees:
  - can be applied in many problems from business and industry, such as energy consumption or healthcare management (e.g. <http://what-when-how.com/artificial-intelligence/decision-tree-applications-for-data-modelling-artificial-intelligence/> (<http://what-when-how.com/artificial-intelligence/decision-tree-applications-for-data-modelling-artificial-intelligence/>))
  - are very fast, result in simple and understandable rules (<https://de.wikipedia.org/wiki/Entscheidungsbaum> (<https://de.wikipedia.org/wiki/Entscheidungsbaum>))
  - are quite simple, therefore they cannot be used for complex problems, and also tend to overfit
  - I chose decision trees because they will result in simple rules that CharityML will easily understand, e.g. "Someone earns probably more than \$50,000 if he is older than 50 years and has a doctoral degree".
- AdaBoost:
  - can be applied for image classification in autonomous driving (e.g. <https://ieeexplore.ieee.org/abstract/document/6728254> (<https://ieeexplore.ieee.org/abstract/document/6728254>))
  - performs well when computational resources are limited
  - but can be sensitive to noisy data and outliers (see <http://www.nickgillian.com/wiki/pmwiki.php/GRT/AdaBoost#Disadvantages> (<http://www.nickgillian.com/wiki/pmwiki.php/GRT/AdaBoost#Disadvantages>))
  - I chose AdaBoost because it combines the advantages of decision trees with fast computation time and only little danger of overfitting. Also, the data is cleaned and doesn't contain many outliers and noise, so AdaBoost should work well.
- Support Vector Machines (SVMs):
  - can be applied for pedestrian detection in autonomous driving (e.g. [https://link.springer.com/chapter/10.1007/978-3-540-76386-4\\_86](https://link.springer.com/chapter/10.1007/978-3-540-76386-4_86) ([https://link.springer.com/chapter/10.1007/978-3-540-76386-4\\_86](https://link.springer.com/chapter/10.1007/978-3-540-76386-4_86)))
  - Its strength lies in the kernel trick, with a good kernel extremely challenging problems can be solved, and SVMs will not find local minima like neural networks (see <https://statinfer.com/204-6-8-svm-advantages-disadvantages-applications/> (<https://statinfer.com/204-6-8-svm-advantages-disadvantages-applications/>)). Also, the margin will be maximized ([https://de.wikipedia.org/wiki/Support\\_Vector\\_Machine](https://de.wikipedia.org/wiki/Support_Vector_Machine) ([https://de.wikipedia.org/wiki/Support\\_Vector\\_Machine](https://de.wikipedia.org/wiki/Support_Vector_Machine))).
  - can only be directly applied with fully-labeled data and two-class problems (see [https://en.wikipedia.org/wiki/Support\\_vector\\_machine#Issues](https://en.wikipedia.org/wiki/Support_vector_machine#Issues) ([https://en.wikipedia.org/wiki/Support\\_vector\\_machine#Issues](https://en.wikipedia.org/wiki/Support_vector_machine#Issues)))
  - I chose SVM because the data is linearly separable and SVM will not only find a hyperplane that separates the data, but will maximize the margin as well.

## 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. Your implementation here will be used in the following section. In the code block below, you will need to implement the following:

- Import `fbeta_score` and `accuracy_score` from `sklearn.metrics` (<http://scikit-learn.org/stable/modules/classes.html#sklearn-metrics-metrics>).
- Fit the learner to the sampled training data and record the training time.
- Perform predictions on the test data `x_test`, and also on the first 300 training points `x_train[:300]`.
  - Record the total prediction time.
- Calculate the accuracy score for both the training subset and testing set.
- Calculate the F-score for both the training subset and testing set.
  - Make sure that you set the `beta` parameter!

```

In [128]: # TODO: Import two metrics from sklearn - fbeta_score and accuracy_score
from sklearn.metrics import fbeta_score, accuracy_score

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'
    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[:300])
    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)

    # 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, beta=0.5)

    # TODO: Compute F-score on the test set which is y_test
    results['f_test'] = fbeta_score(y_test, predictions_test, beta=0.5)

    # Success
    print("{} trained on {} samples.".format(learner.__class__.__name__, sample_size))

    # Return the results
    return results

```



## Implementation: Initial Model Evaluation

In the code cell, you will need to implement the following:

- Import the three supervised learning models you've discussed in the previous section.
- Initialize the three models and store them in `'clf_A'`, `'clf_B'`, and `'clf_C'`.
  - Use a `'random_state'` for each model you use, if provided.
  - **Note:** Use the default settings for each model — you will tune one specific model in a later section.
- Calculate the number of records equal to 1%, 10%, and 100% of the training data.
  - Store those values in `'samples_1'`, `'samples_10'`, and `'samples_100'` respectively.

**Note:** Depending on which algorithms you chose, the following implementation may take some time to run!

```
In [114]: # TODO: Import the three supervised learning models from sklearn
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import AdaBoostClassifier
from sklearn.svm import SVC

# TODO: Initialize the three models
clf_A = DecisionTreeClassifier(random_state=411)
clf_B = AdaBoostClassifier(random_state=411)
clf_C = SVC(random_state=411)

# TODO: Calculate the number of samples for 1%, 10%, and 100% of the training set
# 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 samples_100)
# HINT: samples_1 is 1% of samples_100 (ensure to set the count of the samples_100)
samples_100 = len(y_train)
samples_10 = int(samples_100/10)
samples_1 = int(samples_100/100)

# 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 called
```

DecisionTreeClassifier trained on 361 samples.

DecisionTreeClassifier trained on 3617 samples.

DecisionTreeClassifier trained on 36177 samples.

AdaBoostClassifier trained on 361 samples.

AdaBoostClassifier trained on 3617 samples.

AdaBoostClassifier trained on 36177 samples.

C:\Users\Antje\Anaconda3\lib\site-packages\sklearn\metrics\classification.py:1135: UndefinedMetricWarning: F-score is ill-defined and being set to 0.0 due to no predicted samples.

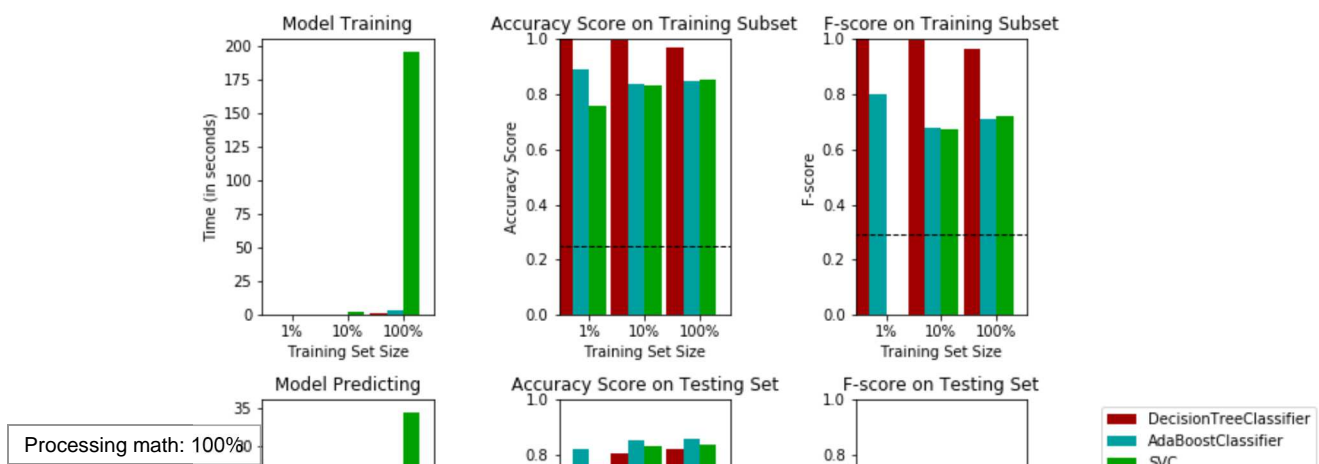
'precision', 'predicted', average, warn\_for)

SVC trained on 361 samples.

SVC trained on 3617 samples.

SVC trained on 36177 samples.

### Performance Metrics for Three Supervised Learning Models



## 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.

### Question 3 - Choosing the Best Model

- Based on the evaluation you performed earlier, in one to two paragraphs, explain to *CharityML* which of the three models you believe to be most appropriate for the task of identifying individuals that make more than \$50,000.

**HINT:** Look at the graph at the bottom left from the cell above (the visualization created by `vs.evaluate(results, accuracy, fscore)`) and check the F score for the testing set when 100% of the training set is used. Which model has the highest score? Your answer should include discussion of the:

- metrics - F score on the testing when 100% of the training data is used,
- prediction/training time
- the algorithm's suitability for the data.

**Answer:** I would choose AdaBoost classifier. F-score and accuracy score on testing set are highest for AdaBoost, no matter which training set size. Also, the decision tree classifier tends to overfit: F-score and accuracy are nearly 1 on the training subset, but much lower on the testing set. This is not the case for AdaBoost. Calculation time is also reasonable and much smaller than for SVM. AdaBoost is suitable for the data because it is a supervised learning algorithm that combines several weak learners to one strong learner, therefore saving calculation time on large data sets.

### Question 4 - Describing the Model in Layman's Terms

- In one to two paragraphs, explain to *CharityML*, in layman's terms, how the final model chosen is supposed to work. Be sure that you are describing the major qualities of the model, such as how the model is trained and how the model makes a prediction. Avoid using advanced mathematical jargon, such as describing equations.

**HINT:**

When explaining your model, if using external resources please include all citations.

**Answer:** The method chosen consists of different weak learners. A weak learner is a simple classification algorithm that can classify some data well, but will perform poorly on other data. Therefore, it is only slightly better than chance or simply guessing. For example, one could use very simple decision trees with depth 1 as weak learners. A decision tree of depth 1 would for example only separate between people of  $\text{age} > 50$  and  $\text{age} \leq 50$  and ignore all other features. Another tree would separate between people who have a university degree and people who don't.

AdaBoost is an ensemble method combining several weak learners to one strong learner. After the first weak learner is applied on the training data, the misclassified points are given a higher weight, therefore ensuring that the next weak classifier will correct the mistakes of the first. A weight is simply a multiplying factor. This is repeated several times. Therefore, each weak learner solves a subproblem, leaving the remaining problem for the next weak learner. Finally, the average over all weak learners is calculated at each data point as result.

## Implementation: Model Tuning

Fine tune the chosen model. Use grid search (`GridSearchCV`) with at least one important parameter tuned with at least 3 different values. You will need to use the entire training set for this. In the code cell below, you will need to implement the following:

- Import `sklearn.grid_search.GridSearchCV` ([http://scikit-learn.org/0.17/modules/generated/sklearn.grid\\_search.GridSearchCV.html](http://scikit-learn.org/0.17/modules/generated/sklearn.grid_search.GridSearchCV.html)) and `sklearn.metrics.make_scorer` ([http://scikit-learn.org/stable/modules/generated/sklearn.metrics.make\\_scorer.html](http://scikit-learn.org/stable/modules/generated/sklearn.metrics.make_scorer.html)).
- Initialize the classifier you've chosen and store it in `clf`.
  - Set a `random_state` if one is available to the same state you set before.
- Create a dictionary of parameters you wish to tune for the chosen model.
  - Example: `parameters = {'parameter' : [list of values]}`.
  - **Note:** Avoid tuning the `max_features` parameter of your learner if that parameter is available!
- Use `make_scorer` to create an `fbeta_score` scoring object (with  $\beta = 0.5$ ).
- Perform grid search on the classifier `clf` using the '`scorer`', and store it in `grid_obj`.
- Fit the grid search object to the training data (`X_train`, `y_train`), and store it in `grid_fit`.

**Note:** Depending on the algorithm chosen and the parameter list, the following implementation may take some time to run!

```
In [124]: # TODO: Import 'GridSearchCV', 'make_scorer', and any other necessary
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import make_scorer

# TODO: Initialize the classifier
clf = AdaBoostClassifier(random_state=411)

# TODO: Create the parameters list you wish to tune, using a dictionary
# HINT: parameters = {'parameter_1': [value1, value2], 'parameter_2': [value1, value2]}
parameters = {'base_estimator': [DecisionTreeClassifier(max_depth=1), DecisionTreeClassifier(max_depth=2)]}

# 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
grid_obj = GridSearchCV(clf, parameters, scoring=scorer)

# TODO: 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 model
predictions = (clf.fit(X_train, y_train)).predict(X_test)
best_predictions = best_clf.predict(X_test)

# Report the before-and-after scores
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_predictions, beta=0.5)))
```

Unoptimized model

-----

Accuracy score on testing data: 0.8576

F-score on testing data: 0.7246

Optimized Model

-----

Final accuracy score on the testing data: 0.8684

Final F-score on the testing data: 0.7423

## Question 5 - Final Model Evaluation

- What is your optimized model's accuracy and F-score on the testing data?
- Are these scores better or worse than the unoptimized model?
- How do the results from your optimized model compare to the naive predictor benchmarks you found earlier in **Question 1**?

**Note:** Fill in the table below with your results, and then provide discussion in the **Answer** box.

**Results:**

Metric	Unoptimized Model	Optimized Model
Accuracy Score	0.8576	0.8684
F-score	0.7246	0.7423

**Answer:** My predicted accuracy and F-score are slightly better than the unoptimized model. These results are much better than the naive predictor, which had an accuracy score of 0.2478 and F-score of 0.2917. The naive predictor was about 25% correct (as about 25% of people earn more than \$50,000), the optimized model is about 87% correct.

---

## Feature Importance

An important task when performing supervised learning on a dataset like the census data we study here is determining which features provide the most predictive power. By focusing on the relationship between only a few crucial features and the target label we simplify our understanding of the phenomenon, which is most always a useful thing to do. In the case of this project, that means we wish to identify a small number of features that most strongly predict whether an individual makes at most or more than \$50,000.

Choose a scikit-learn classifier (e.g., adaboost, random forests) that has a `feature_importance_` attribute, which is a function that ranks the importance of features according to the chosen classifier. In the next python cell fit this classifier to training set and use this attribute to determine the top 5 most important features for the census dataset.

## Question 6 - Feature Relevance Observation

When **Exploring the Data**, it was shown there are thirteen available features for each individual on record in the census data. Of these thirteen records, which five features do you believe to be most important for prediction, and in what order would you rank them and why?

**Answer:** I suppose capital-gain is the most important as it shows the current income. I would choose occupation second and workclass third, because the income highly depends on the current occupation. I would suggest education-num and education\_level forth and fifth, because education has a high impact on income as well.

## Implementation - Extracting Feature Importance

Choose a `scikit-learn` supervised learning algorithm that has a `feature_importance_` attribute available for it. This attribute is a function that ranks the importance of each feature when making predictions based on the chosen algorithm.

In the code cell below, you will need to implement the following:

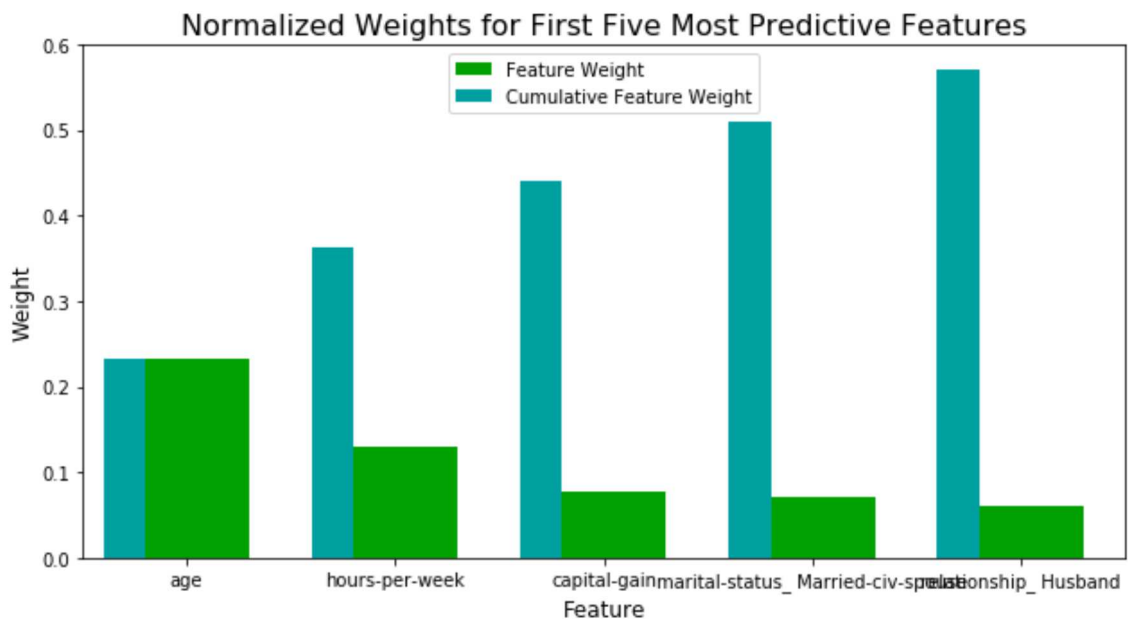
- Import a supervised learning model from `sklearn` if it is different from the three used earlier.
- Train the supervised model on the entire training set.
- Extract the feature importances using `'.feature_importances_'`.

```
In [126]: # TODO: Import a supervised learning model that has 'feature_importance_
from sklearn.ensemble import ExtraTreesClassifier

# TODO: Train the supervised model on the training set using .fit(X_train, y_train)
model = ExtraTreesClassifier(n_estimators=250, random_state=411)
model.fit(X_train, y_train)

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

# Plot
```



## Question 7 - Extracting Feature Importance

Observe the visualization created above which displays the five most relevant features for predicting if an individual makes at most or above \$50,000.

- How do these five features compare to the five features you discussed in **Question 6**?
- If you were close to the same answer, how does this visualization confirm your thoughts?
- If you were not close, why do you think these features are more relevant?

Processing math: 100%

**Answer:** The result shows that the most important features are age, hours-per-week, capital gain, being married and being a husband. My suggestions in Question 6 were wrong, obviously profession and education are not as important as I thought. Age is very important, which makes sense, because usually people reach top positions only after gaining lots of experience. Hours-per-week are important because the higher the position, the more time-consuming the job is. Capital-gain is important, just as I suggested above. Marital status and being a husband is probably important because successful people often have a spouse taking care of household, children etc.

## Feature Selection

How does a model perform if we only use a subset of all the available features in the data? With less features required to train, the expectation is that training and prediction time is much lower — at the cost of performance metrics. From the visualization above, we see that the top five most important features contribute more than half of the importance of **all** features present in the data. This hints that we can attempt to *reduce the feature space* and simplify the information required for the model to learn. The code cell below will use the same optimized model you found earlier, and train it on the same training set *with only the top five important features*.

```
In [127]: # 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(importance
          X_test_reduced = X_test[X_test.columns.values[(np.argsort(importances)

          # 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,
          print("F-score on testing data: {:.4f}".format(fbeta_score(y_test, best
          print("\nFinal Model trained on reduced data\n-----")
          print("Accuracy on testing data: {:.4f}".format(accuracy_score(y_test,

Final Model trained on full data
-----
Accuracy on testing data: 0.8684
F-score on testing data: 0.7423

Final Model trained on reduced data
-----
Accuracy on testing data: 0.8197
F-score on testing data: 0.6321
```



## Question 8 - Effects of Feature Selection

- How does the final model's F-score and accuracy score on the reduced data using only five features compare to those same scores when all features are used?
- If training time was a factor, would you consider using the reduced data as your training set?

**Answer:** The accuracy of the reduced model is slightly lower, 82% instead of 87%. Also, the F-score is 63% instead of 74%. But if training time was a factor, I would definitely consider using this reduced training set because it finished immediately, whereas the full model took much longer to finish.

**Note:** Once you have completed all of the code implementations and successfully answered each question above, you may finalize your work by exporting the iPython Notebook as an HTML document. You can do this by using the menu above and navigating to

**File -> Download as -> HTML (.html).** Include the finished document along with this notebook as your submission.