

Identifying safe loans with decision trees

The LendingClub is a peer-to-peer leading company that directly connects borrowers and potential lenders/investors. In this notebook, you will build a classification model to predict whether or not a loan provided by LendingClub is likely to default.

In this notebook you will use data from the LendingClub to predict whether a loan will be paid off in full or the loan will be charged off and possibly go into default. In this assignment you will:

- Use SFrames to do some feature engineering.
- Train a decision-tree on the LendingClub dataset.
- Visualize the tree.
- Predict whether a loan will default along with prediction probabilities (on a validation set).
- Train a complex tree model and compare it to simple tree model.

If you are doing the assignment with IPython Notebook

An IPython Notebook has been provided below to you for this assignment. This notebook contains the instructions, quiz questions and partially-completed code for you to use as well as some cells to test your code.

Make sure that you are using GraphLab Create 1.8.3. See this post (<https://www.coursera.org/learn/ml-classification/supplement/LgZ3I/installing-correct-version-of-graphlab-create>) for installing the correct version of GraphLab Create.

What you need to download

If you are using GraphLab Create:

- Download the Lending club data in SFrame format: `lending-club-data.gl.zip`
(<https://s3.amazonaws.com/static.dato.com/files/coursera/course-3/lending-club-data.gl.zip>)
- Download the companion IPython Notebook: `module-5-decision-tree-assignment-1-blank.ipynb`
(<https://s3.amazonaws.com/static.dato.com/files/coursera/course-3/module-5-decision-tree-assignment-1-blank.ipynb>)
- Save both of these files in the same directory (where you are calling IPython notebook from) and unzip the data file.
- Follow the instructions contained in the IPython notebook.

If you are not using GraphLab Create:

- If you are using SFrame, download the LendingClub dataset in SFrame format: `lending-club-data.gl.zip`
(<https://s3.amazonaws.com/static.dato.com/files/coursera/course-3/lending-club-data.gl.zip>)
- If you are using a different package, download the LendingClub dataset in CSV format: `lending-club-data.csv.zip`
(<https://s3.amazonaws.com/static.dato.com/files/coursera/course-3/lending-club-data.csv.zip>)

If you are using GraphLab Create and the companion IPython Notebook

Open the companion IPython notebook and follow the instructions in the notebook.

If you are using other tools

This section is designed for people using tools other than GraphLab Create. Even though some instructions are specific to scikit-learn, most part of the assignment should be applicable to other tools as well. However, we highly suggest you use SFrame (<https://github.com/dato-code/SFrame>) since it is open source. In this part of the assignment, we describe general instructions, however we will tailor the instructions for SFrame and scikit-learn.

- If you choose to use SFrame and scikit-learn, you should be able to follow the instructions here and complete the assessment. **All code samples given here will be applicable to SFrame and scikit-learn.**
- You are free to experiment with any tool of your choice, but **some many not produce correct numbers for the quiz questions.**

Load the Lending Club dataset

We will be using a dataset from the LendingClub (<https://www.lendingclub.com/>).

1. Load the dataset into a data frame named **loans**. Using SFrame, this would look like

```
import sframe
loans = sframe.SFrame('lending-club-data.gl/')
```

Note: To install SFrame (without installing GraphLab Create), run

```
pip install sframe
```

Exploring some features

2. Let's quickly explore what the dataset looks like. First, print out the column names to see what features we have in this dataset. On SFrame, you can run this code:

```
loans.column_names()
```

Here, we should see that we have some feature columns that have to do with grade of the loan, annual income, home ownership status, etc.

Exploring the target column

The target column (label column) of the dataset that we are interested in is called `bad_loans`. In this column **1** means a risky (bad) loan **0** means a safe loan.

In order to make this more intuitive and consistent with the lectures, we reassign the target to be:

- **+1** as a safe loan
- **-1** as a risky (bad) loan

3. We put this in a new column called **safe_loans**.

```
# safe_loans = 1 => safe
# safe_loans = -1 => risky
loans['safe_loans'] = loans['bad_loans'].apply(lambda x : +1 if x==0 else -1)
loans = loans.remove_column('bad_loans')
```

4. Now, let us explore the distribution of the column `safe_loans`. This gives us a sense of how many safe and risky loans are present in the dataset. Print out the percentage of safe loans and risky loans in the data frame.

You should have:

- Around 81% safe loans
- Around 19% risky loans

It looks like most of these loans are safe loans (thankfully). But this does make our problem of identifying risky loans challenging.

Features for the classification algorithm

5. In this assignment, we will be using a subset of features (categorical and numeric). The features we will be using are **described in the code comments** below. If you are a finance geek, the LendingClub (<https://www.lendingclub.com/>) website has a lot more details about these features. Extract these feature columns and target column from the dataset. We will only use these features.

```
features = ['grade',           # grade of the loan
            'sub_grade',       # sub-grade of the loan
            'short_emp',       # one year or less of employment
            'emp_length_num',   # number of years of employment
            'home_ownership',   # home_ownership status: own, mortgage or rent
            'dti',              # debt to income ratio
            'purpose',          # the purpose of the loan
            'term',              # the term of the loan
            'last_delinq_none',  # has borrower had a delinquency
            'last_major_derog_none', # has borrower had 90 day or worse rating
            'revol_util',        # percent of available credit being used
            'total_rec_late_fee', # total late fees received to day
            ]

target = 'safe_loans'          # prediction target (y) (+1 means safe, -1 is risky)

# Extract the feature columns and target column
loans = loans[features + [target]]
```

What remains now is a **subset of features** and the **target** that we will use for the rest of this notebook.

Notes to people using other tools

If you are using SFrame, proceed to the section "Sample data to balance classes".

If you are NOT using SFrame, download the list of indices for the training and validation sets: module-5-assignment-1-train-idx.json (<https://s3.amazonaws.com/static.dato.com/files/coursera/course-3/indices-json/module-5-assignment-1-train-idx.json>), module-5-assignment-1-validation-idx.json (<https://s3.amazonaws.com/static.dato.com/files/coursera/course-3/indices-json/module-5-assignment-1-validation-idx.json>). Then follow the following steps:

- Apply one-hot encoding to **loans**. Your tool may have a function for one-hot encoding. Alternatively, see #7 for implementation hints.
- Load the JSON files into the lists **train_idx** and **validation_idx**.

- Perform train/validation split using **train_idx** and **validation_idx**. In Pandas, for instance:

```
train_data = loans.iloc[train_idx]
validation_data = loans.iloc[validation_idx]
```

IMPORTANT: If you are using a programming language with 1-based indexing (e.g. R, Matlab), make sure to increment all indices by 1.

Note. Some elements in loans are included neither in **train_data** nor **validation_data**. This is to perform sampling to achieve class balance.

Now proceed to the section "Build a decision tree classifier", skipping three sections below.

Sample data to balance classes

6. As we explored above, our data is disproportionately full of safe loans. Let's create two datasets: one with just the safe loans (**safe_loans_raw**) and one with just the risky loans (**risky_loans_raw**).

```
safe_loans_raw = loans[loans[target] == +1]
risky_loans_raw = loans[loans[target] == -1]
print "Number of safe loans : %s" % len(safe_loans_raw)
print "Number of risky loans : %s" % len(risky_loans_raw)
```

One way to combat class imbalance is to undersample the larger class until the class distribution is approximately half and half. Here, we will undersample the larger class (safe loans) in order to balance out our dataset. This means we are throwing away many data points. We used seed=1 so everyone gets the same results.

```
# Since there are fewer risky loans than safe loans, find the ratio of the sizes
# and use that percentage to undersample the safe loans.
percentage = len(risky_loans_raw)/float(len(safe_loans_raw))

risky_loans = risky_loans_raw
safe_loans = safe_loans_raw.sample(percentage, seed=1)

# Append the risky_loans with the downsampled version of safe_loans
loans_data = risky_loans.append(safe_loans)
```

You can verify now that **loans_data** is comprised of approximately 50% safe loans and 50% risky loans.

Note: There are many approaches for dealing with imbalanced data, including some where we modify the learning algorithm. These approaches are beyond the scope of this course, but some of them are reviewed in this paper (<http://ieeexplore.ieee.org/xpl/login.jsp?tp=&arnumber=5128907&url=http%3A%2F%2Fieeexplore.ieee.org%2Fiel5%2F69%2F5173046%2F05128907.pdf%3Farnumber%3D5128907>). For this assignment, we use the simplest possible approach, where we subsample the overly represented class to get a more balanced dataset. In general, and especially when the data is highly imbalanced, we recommend using more advanced methods.

One-hot encoding

7. For scikit-learn's decision tree implementation, it requires numerical values for its data matrix. This means you will have to turn categorical variables into binary features via one-hot encoding. **The next assignment has more details about this.**

If you are using SFrame, feel free to use this piece of code as is. Refer to the SFrame API documentation for a deeper understanding. If you are using different machine learning software, make sure you prepare the data to be passed to the learning software.

```

loans_data = risky_loans.append(safe_loans)

categorical_variables = []
for feat_name, feat_type in zip(loans_data.column_names(), loans_data.column_types()):
    if feat_type == 'str':
        categorical_variables.append(feat_name)

for feature in categorical_variables:
    loans_data_one_hot_encoded = loans_data[feature].apply(lambda x: {x: 1})
    loans_data_unpacked = loans_data_one_hot_encoded.unpack(column_name_prefix=feature)

    # Change None's to 0's
    for column in loans_data_unpacked.column_names():
        loans_data_unpacked[column] = loans_data_unpacked[column].fillna(0)

    loans_data.remove_column(feature)
    loans_data.add_columns(loans_data_unpacked)

```

Split data into training and validation

8. We split the data into training and validation sets using an 80/20 split and specifying `seed=1` so everyone gets the same results. Call the training and validation sets **train_data** and **validation_data**, respectively.

Note: In previous assignments, we have called this a **train-test split**. However, the portion of data that we don't train on will be used to help **select model parameters** (this is known as model selection). Thus, this portion of data should be called a **validation set**. Recall that examining performance of various potential models (i.e. models with different parameters) should be on validation set, while evaluation of the final selected model should always be on test data. Typically, we would also save a portion of the data (a real test set) to test our final model on or use cross-validation on the training set to select our final model. But for the learning purposes of this assignment, we won't do that.

```
train_data, validation_data = loans_data.random_split(.8, seed=1)
```

Build a decision tree classifier

9. Now, let's use the built-in scikit learn decision tree learner (`sklearn.tree.DecisionTreeClassifier` (<http://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html>)) to create a loan prediction model on the training data. To do this, you will need to import **sklearn**, **sklearn.tree**, and **numpy**.

Note: You will have to first convert the SFrame into a numpy data matrix, and extract the target labels as a numpy array (Hint: you can use the **.to_numpy()** method call on SFrame to turn SFrames into numpy arrays). See the API (<http://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html>) for more information.

Make sure to set `max_depth=6`.

Call this model **decision_tree_model**.

10. Also train a tree using with **max_depth=2**. Call this model **small_model**.

Visualizing a learned model (Optional)

10a. For this optional section, we would like to see what the small learned tree looks like. If you are using scikit-learn and have the package Graphviz (<http://graphviz.readthedocs.org/en/latest/#>), then you will be able to perform this section. If you are using a different software, try your best to follow along.

Visualize **small_model** in the software of your choice.

Making predictions

Let's consider two positive and two negative examples **from the validation set** and see what the model predicts. We will do the following:

- Predict whether or not a loan is safe.
- Predict the probability that a loan is safe.

11. First, let's grab 2 positive examples and 2 negative examples. In SFrame, that would be:

```
validation_safe_loans = validation_data[validation_data[target] == 1]
validation_risky_loans = validation_data[validation_data[target] == -1]

sample_validation_data_risky = validation_risky_loans[0:2]
sample_validation_data_safe = validation_safe_loans[0:2]

sample_validation_data = sample_validation_data_safe.append(sample_validation_data_risky)
sample_validation_data
```

12. Now, we will use our model to predict whether or not a loan is likely to default. For each row in the **sample_validation_data**, use the **decision_tree_model** to predict whether or not the loan is classified as a **safe loan**. (Hint: if you are using scikit-learn, you can use the **.predict()** method)

Quiz Question: What percentage of the predictions on **sample_validation_data** did **decision_tree_model** get correct?

Explore probability predictions

13. For each row in the **sample_validation_data**, what is the probability (according **decision_tree_model**) of a loan being classified as **safe**? (Hint: if you are using scikit-learn, you can use the **.predict_proba()** method)

Quiz Question: Which loan has the highest probability of being classified as a **safe loan**?

Checkpoint: Can you verify that for all the predictions with probability ≥ 0.5 , the model predicted the label **+1**?

Tricky predictions!

14. Now, we will explore something pretty interesting. For each row in the **sample_validation_data**, what is the probability (according to **small_model**) of a loan being classified as **safe**?

Quiz Question: Notice that the probability predictions are the **exact same** for the 2nd and 3rd loans. Why would this happen?

Visualize the prediction on a tree

14a. Note that you should be able to look at the small tree (of depth 2), traverse it yourself, and visualize the prediction being made. Consider the following point in the **sample_validation_data**:

```
sample_validation_data[1]
```

If you have Graphviz, go ahead and re-visualize **small_model** here to do the traversing for this data point.

Quiz Question: Based on the visualized tree, what prediction would you make for this data point (according to **small_model**)? (If you don't have Graphviz, you can answer this quiz question by executing the next part.)

15. Now, verify your prediction by examining the prediction made using **small_model**.

Evaluating accuracy of the decision tree model

Recall that the accuracy is defined as follows:

$$\text{accuracy} = \frac{\# \text{ correctly classified examples}}{\# \text{ total examples}}$$

16. Evaluate the accuracy of **small_model** and **decision_tree_model** on the training data. (Hint: if you are using scikit-learn, you can use the **.score()** method)

Checkpoint: You should see that the **small_model** performs worse than the **decision_tree_model** on the training data.

17. Now, evaluate the accuracy of the **small_model** and **decision_tree_model** on the entire validation_data, not just the subsample considered above.

Quiz Question: What is the accuracy of decision_tree_model on the validation set, rounded to the nearest .01?

Evaluating accuracy of a complex decision tree model

Here, we will train a large decision tree with **max_depth=10**. This will allow the learned tree to become very deep, and result in a very complex model. Recall that in lecture, we prefer simpler models with similar predictive power. This will be an example of a more complicated model which has similar predictive power, i.e. something we don't want.

18. Using `sklearn.tree.DecisionTreeClassifier` (<http://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html#sklearn.tree.DecisionTreeClassifier>), train a decision tree with maximum depth = 10. Call this model **big_model**.

19. Evaluate the accuracy of **big_model** on the training set and validation set.

Checkpoint: We should see that **big_model** has even better performance on the training set than **decision_tree_model** did on the training set.

Quiz Question: How does the performance of **big_model** on the validation set compare to **decision_tree_model** on the validation set? Is this a sign of overfitting?

Quantifying the cost of mistakes

Every mistake the model makes costs money. In this section, we will try and quantify the cost each mistake made by the model. Assume the following:

- **False negatives:** Loans that were actually safe but were predicted to be risky. This results in an opportunity cost of losing a loan that would have otherwise been accepted.
- **False positives:** Loans that were actually risky but were predicted to be safe. These are much more expensive because it results in a risky loan being given.
- **Correct predictions:** All correct predictions don't typically incur any cost.

Let's write code that can compute the cost of mistakes made by the model. Complete the following 4 steps:

1. First, let us compute the predictions made by the model.

2. Second, compute the number of false positives.
3. Third, compute the number of false negatives.
4. Finally, compute the cost of mistakes made by the model by adding up the costs of true positives and false positives.

Quiz Question: Let's assume that each mistake costs us money: a false negative costs \$10,000, while a false positive costs \$20,000. What is the total cost of mistakes made by `decision_tree_model` on `validation_data`?

