



# PROJECT 3: ASSESS LEARNERS

## DUE DATE

09/20/2020 11:59PM [Anywhere on Earth time](#)

## REVISIONS

This assignment is subject to change up until 3 weeks prior to the due date. We do not anticipate changes; any changes will be logged in this section.

## OVERVIEW

**This assignment counts towards 15% of your overall grade.**

You are to implement and evaluate three learning algorithms as Python classes: a “classic” Decision Tree learner, a Random Tree learner, and a Bootstrap Aggregating learner, and an Insane Learner. Note that a Linear Regression learner is provided for you in the assess learners zip file, name LinRegLearner. The new classes should be named DTLearner, RTLearner, and BagLearner respectively.

Note that we are considering this a **regression** problem (not classification). So the goal for your learner is to return a continuous numerical result (not a discrete result). In this project we are ignoring the time order aspect of the data and treating it as if it is static data and time does not matter. In a later project we will make the transition to considering time series data.

You must write your own code for the Decision Tree learner, Random Tree learner, and Bagging. You are NOT allowed to use other people’s code to implement these learners.

This project has two main components: the code for your learners, which will be auto graded, and your report, `report.pdf`, that should include the components listed in the Report section below.

NOTE: you are given a grading script, named `grade_learners.py`, which gives you a rough indication on how you are doing with the project. The score you get on the provided grading script does not necessarily translate to your final grade with the private grader. You are encouraged to test edge cases as needed.

Your learners should be able to handle any number of dimensions in  $X$  from 2 to  $N$ .

NOTE: All of the the decision tree learners you create should be implemented using a matrix data representation (numpy ndarray **\*\*\*no casting\*\*\***).

## TEMPLATE

Instructions:

- Download the appropriate zip file [File:Assess\\_Learners\\_2020Fall.zip](#)
- You should see these files:
  - `assess_learners` the assignment directory
  - `assess_learners/Data/`: Contains data for you to test your learning code on.
  - `assess_learners/LinRegLearner.py`: An implementation of the `LinRegLearner` class. You can use it as a template for implementing your learner classes.
  - `assess_learners/testlearner.py`: Simple testing scaffold that you can use to test your learners. Useful for debugging.
  - `assess_learners/grade_learners.py`: The grading script; more details here: [ML4T\\_Software\\_Setup#Running\\_the\\_grading\\_scripts](#)

In the `assess_learners/Data/` directory you will find these files:

- `3_groups.csv`
- `ripple.csv`
- `simple.csv`
- `winequality-red.csv`
- `winequality-white.csv`
- `winequality.names.txt`
- `Istanbul.csv`

In these files, we've provided test data for you to use in determining the correctness of the output of your learners. Each data file contains  $N+1$  columns:  $X_1, X_2, \dots, X_N$ , and  $Y$ . For the task below, you will mainly be working with the `Istanbul` datafile. This file includes the returns of multiple worldwide indexes for a number of days in history. In this task, the overall objective is to predict what the return for the MSCI Emerging Markets (EM) index

will be on the basis of the other index returns. Y in this case is the last column to the right of the Istanbul.csv file while the X values are the remaining columns to the left (except the first column). The first column of data in this file is the date, **which you should ignore**. Note that the local test script does this automatically for you, but you will have to handle it yourself when working on your report.

When the grading script tests your code it randomly selects 60% of the data to train on and uses the other 40% for testing.

The other files, besides Istanbul.csv are there as alternative sets for you to test your code on. Each data file contains N+1 columns: X1, X2, ... XN, and Y.

The Istanbul data is also available here: [File:Istanbul.csv](#)

## HINTS & RESOURCES

“Official” course-based materials:

- [How to use a decision tree if you have one \(Balch Youtube video\)](#)
- [How to build a decision tree & Random Trees \(Balch Youtube video\)](#)
- [Media:How-to-learn-a-decision-tree.pdf](#) Balch slides on decision trees
- [Media:Decision-tree-example.xlsx](#) Example tabular version of decision tree

Additional supporting materials:

- You may be interested in looking at Andrew Moore’s slides on [instance based learning](#).
- A definition of [correlation](#) which is used to assess the quality of the learning.
- [Bootstrap Aggregating](#)
- [AdaBoost](#)
- [numpy corrcoef](#)
- [numpy argsort](#)
- [RMS error](#)

## TASKS

### Implement DTLearner (15 points)

Implement a Decision Tree learner class named DTLearner in the file DTLearner.py. You should follow the algorithm outlined in the presentation here [decision tree slides](#).

- We define “best feature to split on” as the feature ( $X_i$ ) that has the highest absolute value correlation with Y.

The algorithm outlined in those slides is based on the paper by [JR Quinlan](#) which you may also want to review as a reference. Note that Quinlan's paper is focused on creating classification trees, while we're creating regression trees here, so you'll need to consider the differences.

For this part of the project, your code should build a single tree only (not a forest). We'll get to forests later in the project. Your code should support exactly the API defined below. DO NOT import any modules besides those listed in the allowed section below. You should implement the following functions/methods:

```
import DTLearner as dt
learner = dt.DTLearner(leaf_size = 1, verbose = False) # constructor
learner.add_evidence(Xtrain, Ytrain) # training step
Y = learner.query(Xtest) # query
```

"Leaf\_size" is a hyperparameter that defines the maximum number of samples to be aggregated at a leaf. While the tree is being constructed recursively, if there are leaf\_size or fewer elements at the time of the recursive call, the data should be aggregated into a leaf. Xtrain and Xtest should be ndarrays (numpy objects) where each row represents an X1, X2, X3... XN set of feature values. The columns are the features and the rows are the individual example instances. Y and Ytrain are single dimension ndarrays that indicate the value we are attempting to predict with X.

If "verbose" is True, your code can print out information for debugging. If verbose = False your code should not generate ANY output. When we test your code, verbose will be False.

This code should not generate statistics or charts.

NOTE: casting the end result, or anywhere in-between, is not implementing with an ndarray (numpy object).

## Implement RTLearner (15 points)

Implement a Random Tree learner class named RTLearner in the file RTLearner.py. This learner should behave exactly like your DTLearner, except that the choice of feature to split on should be made randomly. You should be able to accomplish this by removing a

few lines from DTLearner (the ones that compute the correlation) and replacing the line that selects the feature with a call to a random number generator.

You should implement the following functions/methods:

```
import RTLearner as rt
learner = rt.RTLearner(leaf_size = 1, verbose = False) # constructor
learner.add_evidence(Xtrain, Ytrain) # training step
Y = learner.query(Xtest) # query
```

## Implement BagLearner (20 points)

Implement Bootstrap Aggregating as a Python class named BagLearner. Your BagLearner class should be implemented in the file BagLearner.py. It should support EXACTLY the API defined below. This API is designed so that BagLearner can accept any learner (e.g., RTLearner, LinRegLearner, even another BagLearner) as input and use it to generate a learner ensemble. Your BagLearner should support the following function/method prototypes:

```
import BagLearner as bl
learner = bl.BagLearner(learner = al.ArbitraryLearner, kwargs = {"argument1": 1})
learner.add_evidence(Xtrain, Ytrain)
Y = learner.query(Xtest)
```

Where learner is the learning class to use with bagging. You should be able to support any learning class that obeys the API defined above for DTLearner and RTLearner. kwargs are keyword arguments to be passed on to the learner's constructor and they vary according to the learner (see example below). The "bags" argument is the number of learners you should train using Bootstrap Aggregation. If boost is true, then you should implement boosting (optional). If verbose is True, your code can generate output. Otherwise the code should be silent.

As an example, if we wanted to make a random forest of 20 Decision Trees with leaf\_size 1 we might call BagLearner as follows

```
import BagLearner as bl
learner = bl.BagLearner(learner = dt.DTLearner, kwargs = {"leaf_size": 1}, b
```

```
learner.add_evidence(Xtrain, Ytrain)
Y = learner.query(Xtest)
```

As another example, if we wanted to build a bagged learner composed of 10 LinRegLearners we might call BagLearner as follows

```
import BagLearner as bl
learner = bl.BagLearner(learner = lrl.LinRegLearner, kwargs = {}, bags = 10)
learner.add_evidence(Xtrain, Ytrain)
Y = learner.query(Xtest)
```

Note that each bag should be trained on a different subset of the data. You will be penalized if this is not the case.

Boosting is an optional topic and not required. There's a citation in the Resources section that outlines a method of implementing boosting.

If the training set contains  $n$  data items, each bag should contain  $n$  items as well. Note that because you should sample with replacement, some of the data items will be repeated.

This code should not generate statistics or charts. If you want create charts and statistics, you can modify `testlearner.py`.

You can use code like the below to instantiate several learners with the parameters listed in `kwargs`:

```
learners = []
kwargs = {"k": 10}
for i in range(0, bags):
    learners.append(learner(**kwargs))
```

## Implement InsaneLearner (Up to 10 point penalty)

Your BagLearner should be able to accept any learner object so long as the learner obeys the API defined above. We will test this in two ways: 1) By calling your BagLearner with an arbitrarily named class and 2) By having you implement InsaneLearner as described below. If your code dies in either case, you will lose 10 points. Note, grading script only does a

rudimentary check thus we will also manually inspect your code for correct implementation and grade accordingly.

Using your BagLearner class and the provided LinRegLearner class, implement InsaneLearner as follows: InsaneLearner should contain 20 BagLearner instances where each instance is composed of 20 LinRegLearner instances. We should be able to call your InsaneLearner using the following API:

```
import InsaneLearner as it
learner = it.InsaneLearner(verbose = False) # constructor
learner.add_evidence(Xtrain, Ytrain) # training step
Y = learner.query(Xtest) # query
```

The code for InsaneLearner should be 20 lines or less. Each ";" in the code counts as one line. All lines (except comments and empty lines) will be counted. There is no credit for this, but a penalty if it is not implemented correctly.

The entire InsaneLearner.py should be 20 lines or less. Please do not include unnecessary blank lines or comments. Have your file under 20 lines total to ease grading.

## Implement author() (Up to 10 point penalty)

For all learners you submit (DT, RT, Bag, Insane) should implement a method called `author()` that returns your Georgia Tech user ID as a string. This must be implemented within **each individual file** even if using inheritance. It is not your 9 digit student number. Here is an example of how you might implement `author()` within a learner object:

```
class LinRegLearner(object):

    def __init__(self):
        pass # move along, these aren't the drones you're looking for

    def author(self):
        return 'tb34' # replace tb34 with your Georgia Tech user id.
```

And here's an example of how it could be called from a testing program:

```
# create a learner and train it
learner = lr1.LinRegLearner() # create a LinRegLearner
```

```
learner.add_evidence(trainX, trainY) # train it
print(learner.author())
```

## Extra Credit (0 points)

Implement boosting as part of BagLearner. How does boosting affect performance compared to not boosting? Does overfitting occur as the number of bags with boosting increases? Create your own dataset for which overfitting occurs as the number of bags with boosting increases.

- Submit your report regarding boosting as report-boosting.pdf

## Report (50 points)

Answer the following prompt in a maximum of 7 pages (excluding references) in [JDF format](#). Any content beyond 7 pages will not be considered for a grade. Seven pages is a maximum, not a target; our recommended per-section lengths intentionally add to less than seven pages to leave you room to decide where to delve into more detail. This length is intentionally set expecting that your submission will include diagrams, drawings, pictures, etc. These should be incorporated into the body of the paper unless specifically required to be included in an appendix.

The [JDF format](#) specifies font sizes and margins, which should not be altered. Include charts (not tables) to support each of your answers. Charts should be generated by the code and saved to files. Charts should be properly annotated with legible and appropriately named labels, titles, and legends.

Please address each of these points / questions, the questions asked in the Project 3 wiki, and the items stated in the Project 3 rubric in your report. The report is to be submitted as **report.pdf**.

### **Abstract: ~0.25 pages**

First, include an abstract that briefly introduces your work and gives context behind your investigation. Ideally, the abstract will fit into 50 words, but should not be more than 100 words.

### **Introduction: ~0.5 pages**

The report should briefly describe the paper's justification. While the introduction may assume that the reader has some domain knowledge, it should assume that the reader is unfamiliar with the specifics of the assignment. The introduction should also present an initial hypothesis (or hypotheses).

### **Methods: ~0.5 pages**



Discuss the setup of the experiment(s) in sufficient detail that an informed reader (someone with familiarity of the field, but not necessarily the assignment) could setup and repeat the experiment(s) you performed.

### **Discussion: ~ 3 pages**

#### **Experiment 1**

Research and discuss overfitting as observed in the experiment. (Use the dataset Istanbul.csv with DTLearner). Support your assertion with graphs/charts. (Do not use bagging in Experiment 1). At a minimum, the following question(s) that must be answered in the discussion:

- Does overfitting occur with respect to leaf\_size?
- For which values of leaf\_size does overfitting occur? Indicate the starting point and the direction of overfitting. Support your answer in the discussion or analysis. Use RMSE as your metric for assessing overfitting.

#### **Experiment 2**

Research and discuss the use of bagging and its effect on overfitting. (Again, use the dataset Istanbul.csv with DTLearner.) Provide charts to validate your conclusions. Use RMSE as your metric. At a minimum, the following questions(s) must be answered in the discussion.

- Can bagging reduce overfitting with respect to leaf\_size? Can bagging eliminate overfitting with respect to leaf\_size? To investigate this, choose a fixed number of bags to use and vary leaf\_size to evaluate. If there is overfitting, indicate the starting point and the direction of overfitting. Support your answer in the discussion or analysis.

#### **Experiment 3**

Quantitatively compare “classic” decision trees (DTLearner) versus random trees (RTLearner). Provide at least two new quantitative measures in the comparison. Using two similar measures that illustrate the same broader metric does not count as two separate measures. (For example, do not use two measures for accuracy.) Note for this part of the report you must conduct new experiments, don’t use the results of the experiments above for this. Importantly, RMSE and correlation are not allowed as a new experiment. Provide charts to support your conclusions. At a minimum, the following question(s) must be answered in the discussion.

- In which ways is one method better than the other?

### **Summary: ~0.5 pages**

The summary is an opportunity to synthesize and summarize the experiments. Ideally it presents key findings and insights discovered during the research. It may also identify interesting areas for future investigation.

**References: ~0.25 pages**

### (Optional)

References should be placed at the end of the paper in a dedicated section. Reference lists should be numbered and organized alphabetically by first author's last name. If multiple papers have the same author(s) and year, you may append a letter to the end of the year to allow differentiated in-line text (e.g. Joyner, 2018a and Joyner, 2018b in the section above). If multiple papers have the same author(s), list them in chronological order starting with the older paper. Only works that are cited in-line should be included in the reference list. The reference list does not count against the length requirements.

### Submission Instructions

Complete your assignment using JDF, then save your submission as a PDF. Assignments should be submitted to the corresponding assignment submission page in Canvas. You should submit a single PDF for this assignment. Be sure to following the instructions in Canvas for the report and code submissions.

This is an individual assignment. All work you submit should be your own. Make sure to cite any sources you reference and use quotes and in-line citations to mark any direct quotes.

Late work is not accepted without advanced agreement except in cases of medical or family emergencies. In the case of such an emergency, please immediately contact the Dean of Students followed by contacting the "Instructors" through a Piazza private post.

## WHAT TO TURN IN

Be sure to follow these instructions diligently!

### Canvas:

Submit the following files (only) via Canvas before the deadline:

- Project 3: Assess Learners (Report)
  - Your report as report.pdf.

Unlimited resubmissions are allowed up to the deadline for the project.

### Gradescope:

- (SUBMISSION) Project 3: Assess Learners
  - Your code as `RTLearner.py`, `DTLearner.py`, `InsaneLearner.py`, `BagLearner.py`, and `testlearner.py`.

Do not submit any other files.

You are only allowed 3 submissions to **(SUBMISSION) Project 3: Assess Learners** but unlimited resubmissions are allowed on **(TESTING) Project 3: Assess Learners**.

Note that Gradescope does **not** grade your assignment live; instead, it pre-validates that it will run against our batch autograder that we will run after the deadline. There will be **no** credit given for coding assignments that do not pass this pre-validation.

Refer to the [Gradescope Instructions](#) for more information.

## RUBRIC

### Report [50 Points]

- Is the report neat and well organized? (-5 points if not)
- Is the experimental methodology and setup well described (up to -5 points per question if not)
- Does the chart properly reflect the intent of the assignment? (up to -10 points if not)
- Does the chart include properly labeled axis and legend? (up to -5 points if not)
- Overfitting / leaf\_size question:
  - Is one or more charts provided to support the argument? (up to -5 points if not)
  - Does the student state where the region of overfitting occurs (or state that there is no overfitting)? (up to -5 points if not)
  - Are the starting point and direction of overfitting identified supported by the data (or if the student states that there is no overfitting, is that supported by the data)? (up to -5 points if not)
- Does bagging reduce or eliminate overfitting?:
  - Is a chart provided to support the argument? (-5 points if not)
  - Does the student state where the region of overfitting occurs (or state that there is no overfitting)? (up to -5 points if not)
  - Are the starting point and direction of overfitting identified supported by the data (or if the student states that there is no overfitting, is that supported by the data)? (up to

-5 points if not)

- Comparison of DT and RT learning
  - Is each quantitative experiment explained well enough that someone else could reproduce it (up to -5 points if not)
  - Are there at least two new quantitative properties that are compared (do not use RMSE or correlation)? (-5 points if only one, -10 if none)
  - Is each conclusion regarding each comparison supported well with charts? (up to -10 points if not)
- Was the report exceptionally well done? (up to +2 points)
- Does the student's response indicate a lack of understanding of overfitting? (up to -10 points)
- Were all charts provided generated in Python? (up to -20 points if not)

## Code

- Is the `author()` method correctly implemented for all files: `DTLearner`, `InsaneLearner`, `BagLearner` and `RTLearner`? (-10 points for each if not)
- Is `InsaneLearner` correctly implemented in 20 lines or less **\*\*Please remove all blank lines and comments\*\*** (-10 points if not)
- Does `BagLearner` work correctly with an arbitrarily named class (-10 points if not)
- Does `BagLearner` generate a different learner in each bag? (-10 points if not)
- Are the decision tree learners implemented using a numpy array (`ndarray`) **\*\*NO CASTING\*\***? (-20 points if not)
- Does the implemented code properly reflect the intent of the assignment? (-20 points if not)
- Does the `testlearner.py` file run using `PYTHONPATH=../:. python testlearner.py Data/Istanbul.csv`? (-20 points if not)
- Does `testlearner.py` run in under 10 minutes? (-20 points if not)
- Does the code generate appropriate charts written to png files? **DO NOT use `plt.show()` and manually save your charts. The charts should be created and saved using Python code** (-20 points if not)

## Auto-Grader [50 Points]

- `DTLearner` in sample/out of sample test, auto grade 5 test cases (4 using `istanbul.csv`, 1 using another data set), 3 points each: 15 points.
  - For each test 60% of the data will be selected at random for training and 40% will be selected for testing.

- Success criteria for each of the 5 tests:
  - 1) Does the correlation between predicted and actual results for **in sample data** exceed 0.95 with leaf\_size = 1?
  - 2) Does the correlation between predicted and actual results for **out of sample** data exceed 0.15 with leaf\_size=1?
  - 3) Is the correlation between predicted and actual results for **in sample** data below 0.95 with leaf\_size = 50?
  - 4) Does the test complete in less than 10 seconds (i.e. 50 seconds for all 5 tests)?
- RTLearner in sample/out of sample test, auto grade 5 test cases (4 using istanbul.csv, 1 using another data set), 3 points each: 15 points.
  - For each test 60% of the data will be selected at random for training and 40% will be selected for testing.
  - Success criteria for each of the 5 tests:
    - 1) Does the correlation between predicted and actual results for **in sample data** exceed 0.95 with leaf\_size = 1?
    - 2) Does the correlation between predicted and actual results for **out of sample** data exceed 0.15 with leaf\_size=1?
    - 3) Is the correlation between predicted and actual results for **in sample** data below 0.95 with leaf\_size = 50?
    - 4) Does the test complete in less than 3 seconds (i.e. 15 seconds for all 5 tests)?
- BagLearner, auto grade 10 test cases (8 using istanbul.csv, 2 using another data set), 2 points each 20 points
  - For each test 60% of the data will be selected at random for training and 40% will be selected for testing.
  - leaf\_size = 20
- Success criteria for each run of the 10 tests:
  - 1) For out of sample data is correlation with 1 bag lower than correlation for 20 bags?
  - 2) Does the test complete in less than 10 seconds (i.e. 100 seconds for all 10 tests)?

## REQUIRED, ALLOWED & PROHIBITED

Required:

- Your code must implement a Random Tree learner.
- Your project must be coded in Python 3.6.x.
- Your code must be submitted to Gradescope in the appropriate Gradescope assignment.
- Your code must run in less than 10 seconds per test case.
- The code you submit should NOT include any data reading routines. The provided `testlearner.py` code reads data for you.
- The learner code files (DT, RT, etc.) you submit should NOT generate any output: No prints, no charts, etc. `testlearner.py` is the only file that should generate charts. \*Note that all charts you provide must be generated in Python
- Reference any code used in the “Allowed” section in your code. At minimum it should have the link/filename/video name of where it came from.

Allowed:

- You can develop your code on your personal machine, but it must also run successfully on Gradescope.
- Your code may use standard Python libraries.
- You may use the NumPy, SciPy, matplotlib and Pandas libraries. Be sure you are using the correct versions.
- Code provided by the instructor, or allowed by the instructor to be shared.
- Cheese.

Prohibited:

- Any other method of reading data besides `testlearner.py`
- Any libraries not listed in the “allowed” section above.
- Any code you did not write yourself.
- Any Classes (other than Random) that create their own instance variables for later use (e.g., learners like `kdtree`).
- Code that includes any data reading routines. The provided `testlearner.py` code reads data for you.
- Code that generates any output when `verbose = False`: No prints, no charts, etc.
- Any use of `plot.show()`
- Absolute import statements of the **current** project folder such as `from assess_learners import XXXX` or `import assess_learners.XXXX`
- Extra directories (manually or code created)
- Extra files not listed in “WHAT TO TURN IN”

## FAQ

- **Q:** Can I use an ML library or do I have to write the code myself?  
**A:** You must write the decision tree and bagging code yourself. The LinRegLearner is provided to you. Do not use other libraries or your code will fail the auto grading test cases.
- **Q:** Which libraries am I allowed to use? Which library calls are prohibited?  
**A:** The use of classes that create and maintain their own data structures are prohibited. So for instance, use of `scipy.spatial.KDTree` is not allowed because it builds a tree and keeps that data structure around for reference later. The intent for this project is that YOU should be building and maintaining the data structures necessary. You can, however, use methods that return immediate results and do not retain data structures
  - Examples of things that are allowed: `sqrt()`, `sort()`, `argsort()` — note that these methods return an immediate value and do not retain data structures for later use.
  - Examples of things that are prohibited: any scikit add on library, `scipy.spatial.KDTree`, importing things from libraries other than pandas, numpy or scipy.
- **Q:** How should I read in the data?  
**A:** Your code does not need to read in data, that is handled for you in the `testlearner.py` and `grade_learners.py` code. For testing your code you can modify `testlearner.py` to read in different datasets, but your solution should NOT depend on any special code in `testlearner.py`
- **Q:** How many data items should be in each bag?  
**A:** If the training set is of size  $N$ , each bag should contain  $N$  items. Note that since sampling is with replacement some of the data items will be repeated.
- **Q:** How do I ignore the first (date) column?  
**A:** Please review the grading script to see a method of performing this action.

## Acknowledgements and Citations

The data used in this assignment was provided by [UCI's ML Datasets](#).