

Due data: 3/13/2019, end of the day. **Please submit the following 2 files via Canvas:**

- 1) For question 1 – 5, please submit in a word file or a PDF file;
- 2) For question 6, please submit a .ipynb file (Python jupyter notebook file).

I. Please answer the following questions related to Machine Learning concepts:

1. [3 points] Explain what is the bias-variance trade-off?
2. [3 points] Describe few techniques to reduce bias and variance respectively.
3. [4 points] Describe following regression measures:
 - 1) RMSE,
 - 2) R^2 measure.
4. [10 points] Explain following concepts related to classification measures:
 - 1) confusion matrix,
 - 2) precision,
 - 3) recall,
 - 4) F_1 score,
 - 5) ROC curve.

II. Please answer the following questions related to Concept Learning:

5. [10 points] Consider the following training examples (which is similar to EnjoySport but with slightly different attributes) and the hypothesis space H that we described in lecture 3-2 (i.e., hypothesis is conjunction of attributes).

	Temp	Humidity	Water	Sky	EnjoySport
1	Warm	Normal	Warm	Sunny	Yes
2	Warm	Normal	Cold	Sunny	Yes
3	Cold	Normal	Warm	Rainy	No
4	Warm	High	Warm	Sunny	Yes

Trace the Candidate-Elimination algorithm to show the sequence of S and G boundary sets.

III. Programming problem (40 points):

6. In this programming problem, you will get familiar with building a decision tree, using cross validation to prune a tree, evaluating the tree performance, and interpreting the result.

Potential packages to use and short tutorials:

(1)<http://scikit-learn.org/stable/modules/tree.html>

(2)http://chrisstrelhoff.ws/sandbox/2015/06/25/decision_trees_in_python_again_cross_validation.html

```
from sklearn import tree # tree library
tree.DecisionTreeClassifier() # for classification tree
tree.DecisionTreeRegressor() # for regression tree
# X: design matrix; Y: labels
fit(X, Y) # fit a tree
predict(X) # make prediction on test data
tree.export_graphviz(model) # visualize tree
from sklearn.model_selection import KFold # K-fold cross validation
```

```
from sklearn.grid_search import GridSearchCV
```

In python, you may have to do gridsearch and cross validation using

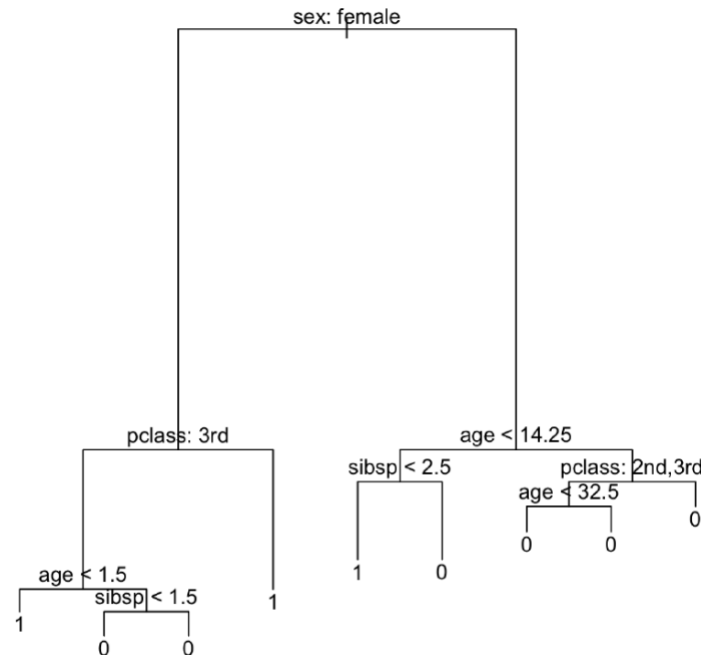
GridSearchCV() to choose the best parameters. Try use different values for "max_leaf_nodes": [None, 1,2,3,4,5,6,7,8,9], (see reference 2).

classification tree

Use the titanic.csv dataset included in the assignment.

Step 1: read in Titanic.csv and observe a few samples, some features are categorical and others are numerical. Take a random 70% samples for training and the rest 30% for test.

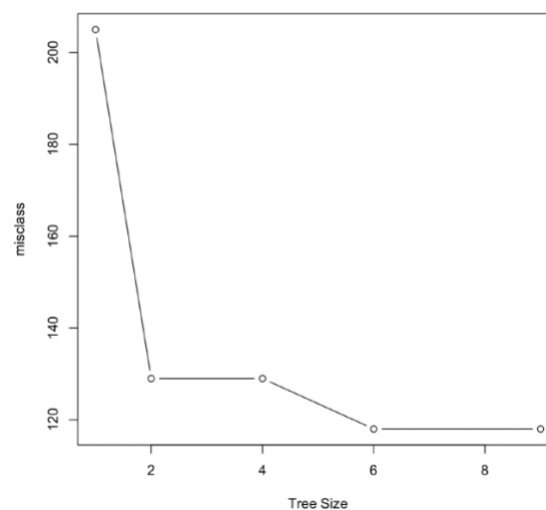
Step 2: fit a decision tree model using independent variables 'pclass + sex + age + sibsp' and dependent variable 'survived'. Plot the full tree. Make sure 'survived' is a qualitative variable taking 1 (yes) or 0 (no) in your code. You may see a tree similar to this one:



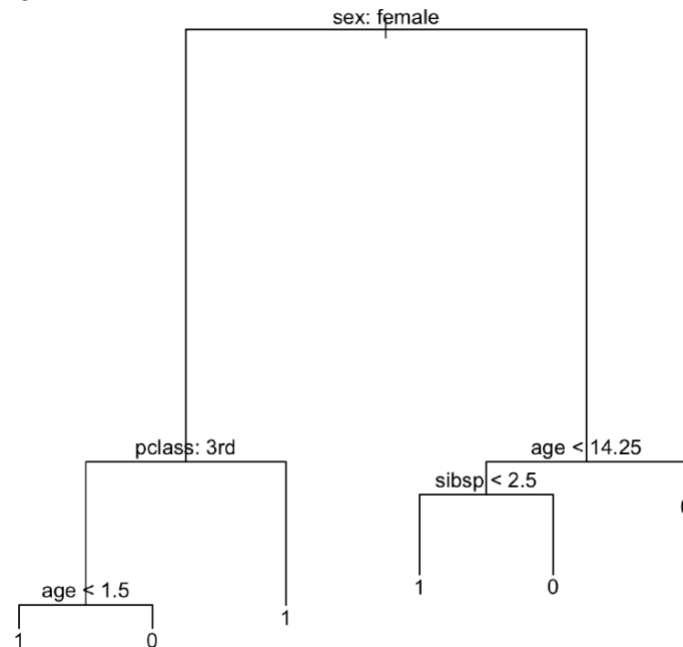
Step 3: check the performance of the full model: in-sample and out-of-sample accuracy, defined as:

- in-sample percent survivors correctly predicted (on training set)
- in-sample percent fatalities correctly predicted (on training set)
- out-of-sample percent survivors correctly predicted (on test set)
- out-of-sample percent fatalities correctly predicted (on test set)

Step 4: use cross-validation to find the best parameter to prune the tree. You should be able to plot a graph with the 'tree size' as the x-axis and 'number of misclassification' as the Y-axis. Find the minimum number of misclassification and choose the corresponding tree size to prune the tree. You may have a plot similar to:



Step 5: prune the tree with the optimal tree size. Plot the pruned tree. You may see a similar tree like this:



Step 6: Report as many details as you can on the final pruned tree.

Required reports on: in-sample and out-of-sample accuracy, defined as

- in-sample percent survivors correctly predicted (on training set)
- in-sample percent fatalities correctly predicted (on training set)
- out-of-sample percent survivors correctly predicted (on test set)
- out-of-sample percent fatalities correctly predicted (on test set)

Check whether there is improvement in out-of-sample for the full tree (bigger model) and the pruned tree (smaller model).