

Code	RANDOM_FOREST.PY
Author	Nathaniel Heatwole, PhD (heatwolen@gmail.com) (GitHub) (LinkedIn)
Summary	Uses a random forest to predict survival for passengers in the Titanic disaster.
Methods/ process	<p><u>Random forest</u>:</p> <ul style="list-style-type: none"> - Supervised learning method that fits many <u>decision trees</u> (“forest”) and aggregates the results. - Combines benefits of decision tree learning, while mitigating their tendency to overfit to their training data. - Each decision tree fitted on: 1) random subset of features; and 2) random selection of training data observations (with replacement). - Randomly withholding some information (that would otherwise be available to fit the model) reduces correlations between the trees. - Trees can be split using various measures, including entropy¹ and Gini impurity² (minimum sought in either case). - <i>Root node</i> (top of tree): quantity/threshold yielding the best split. - Predictions of the many individual trees are homogenized using either plurality vote (classification) or average (regression). - Out-of-bag testing can also be used (if entire dataset not used to generate tree). <p>Steps:</p> <ol style="list-style-type: none"> 1. Import and clean training data. 2. Conduct feature engineering (prepare data for use in a random forest model and maximize useful information to be extracted from it). 3. Assess feature importance (Pearson correlation matrix, chi-squared, and coefficient of variation). 4. Fit random forest to training data to predict survival.
Training data	<u>Titanic dataset</u> – containing data for 891 Titanic passengers (from Kaggle).
Output	<p>Plot</p> <ul style="list-style-type: none"> - Actual versus predicted survival (with group-level summary stats)
Result	<p>The predictions align with expectations: high predicted survival probabilities for survivors (mean = 0.87), and the opposite for non-survivors (mean = 0.09). Additionally, the predicted probability IQRs for the two groups are non-overlapping.</p>

¹ *Entropy* is a measure of disorder. It is defined as the average (expected value) of information, and is equal to: $\sum [-p * \ln(p)]$.

² *Gini impurity* measures how often a randomly chosen element would be incorrectly labeled (if group labels were assigned randomly and independently using the distribution of labels in the training set). It is equal to: $\sum [p * (1 - p)] = 1 - \sum (p^2)$.