

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

---

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

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