

<b>Code</b>	<b>LOGISTIC_REGRESSION.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 logistic regression to predict probability of passing an exam using number of hours studied (using both sklearn and statsmodels)
<b>Methods/ Process</b>	<p><a href="#">Logistic regression</a></p> <ul style="list-style-type: none"> <li>- Extends concepts of linear regression (OLS) to model a binary (0/1) rather than continuous quantity</li> <li>- Basis is the <a href="#">logistic</a> function, which takes the form of an s-curve, with dual asymptotes at zero and one</li> <li>- Therefore, functions as binary classifier, outputting the probability that <math>Y = 1</math></li> <li>- Argument (X) of the logistic function is a linear equation containing one or multiple variables (continuous or categorical)</li> <li>- Accordingly, the logistic function maps an unbounded quantity (X) onto the bounded probability space (0-1)</li> <li>- Coefficients for linear component (X) often estimated using <a href="#">maximum likelihood</a><sup>1</sup></li> <li>- Model fit typically assessed by examining rates of true/false positives</li> <li>- Extensions of logistic regression exist for non-binary outcome variables (Y), or those that can take on more than two possible levels<sup>2</sup></li> </ul>
<b>Training Data</b>	<a href="#">Exam data</a> – synthetic data on whether 20 students passed an exam and number of hours studied (from <i>Wikipedia</i> )
<b>Results</b>	Predictions from the two logistic regression modules (sklearn and statsmodels) are similar, although the statsmodels fit is slightly better (based on log-likelihood)

---

<sup>1</sup> Informally, this means choosing the model parameters so as to maximize the likelihood that the particular training data would be observed. More rigorously, it means selecting the s-curve parameters (location, steepness) so as to minimize the aggregate distance between the s-curve and the true outcomes (0 or 1).

<sup>2</sup> If there is a “natural” order (monotonicity) to the groups (e.g., small/medium/large), [ordinal logistic regression](#) is useful. Otherwise, if no such natural order exists (e.g., favorite color), [multinomial logistic regression](#) relaxes the monotonicity assumption, and allows the regression to select the order of the groups.