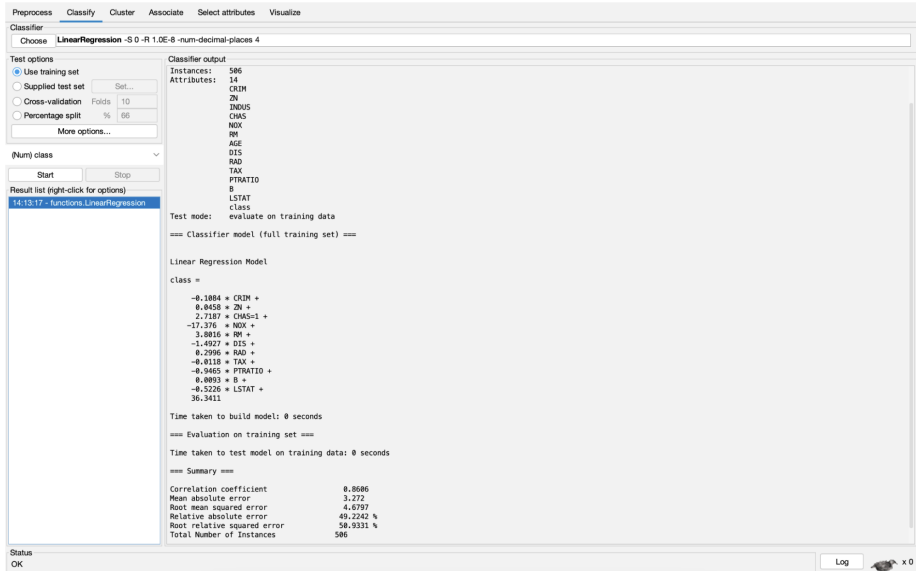


Data Mining Rubrics

WEKA

Q1 (2 points). Build a linear regression equation to predict MEDV. Include a screenshot that shows the linear equation. How many terms are in the equation, and 'why'? In other words, discuss the resulting equation.

POINTS	DESCRIPTION
1 point	<p>For the screenshot: The regression model should be in the screenshot.</p> <p><i>If not, 0 points.</i></p> 
1 point	<p>For the description of results, <i>each following point is worth 0.5 points</i>:</p> <ol style="list-style-type: none"> Number of terms. By default, if we use all 14 attributes (including the class attributes), the number of terms should be 12 (including the constant) or 11 (except the constant). If the student has deliberately deleted some attribute, it is acceptable if they have explained their reasons in the README. <i>If not, 0 points.</i> A reasonable discussion about any one of the following points: <ol style="list-style-type: none"> The “missing” attributes INDUS and AGE; one acceptable reason is that WEKA discards statistically insignificant attributes as measured by their R-squared

	<p>value</p> <p>ii) The weight of each attribute (negative weight means a higher value of the attribute usually lowers the median home price whereas the one with a positive weight raises it)</p> <p>iii) Other rational interpretations of the equation</p>
--	---

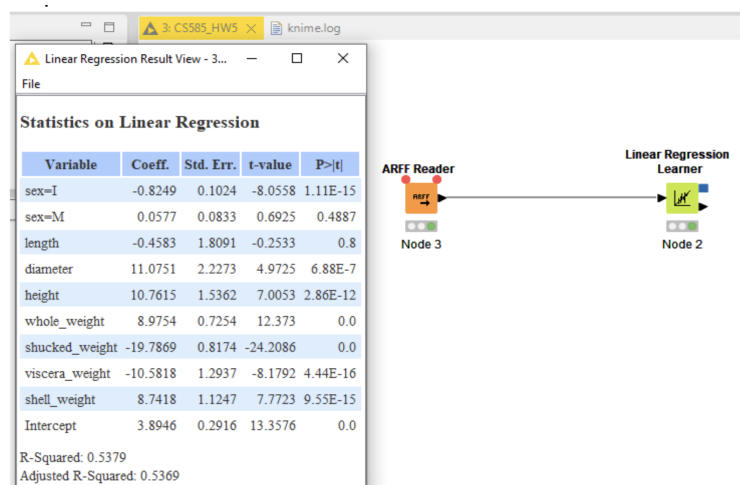
*** if the information on the screenshot does not match the discussion, [0 points for Q1](#)

KNIME

Q2 (2 points). Use KNIME to perform linear regression [on all parameters, not a subset]. You need these nodes: AARF Reader, Linear Regression Learner. Create and connect the nodes, and execute each. What is the linear equation? Include a screenshot.


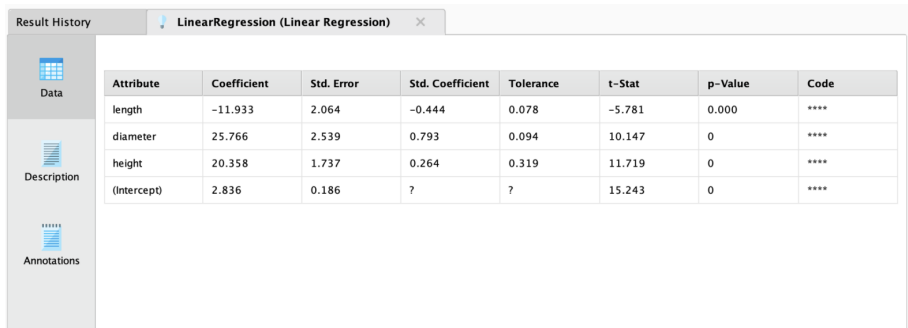
POINTS	DESCRIPTION
1 point	Screenshot shows workflow (students may use a different file reader but not a different learner)
1 point	<p>Screenshot shows statistics</p> <ul style="list-style-type: none"> - At least the first decimal of each coefficient must match the values below eg. the coefficient of length may be 0.46 but not 0.5. If the coefficients vary drastically and the student has not explained the modifications they made, 0 points - If any one variable is missing, 0 points

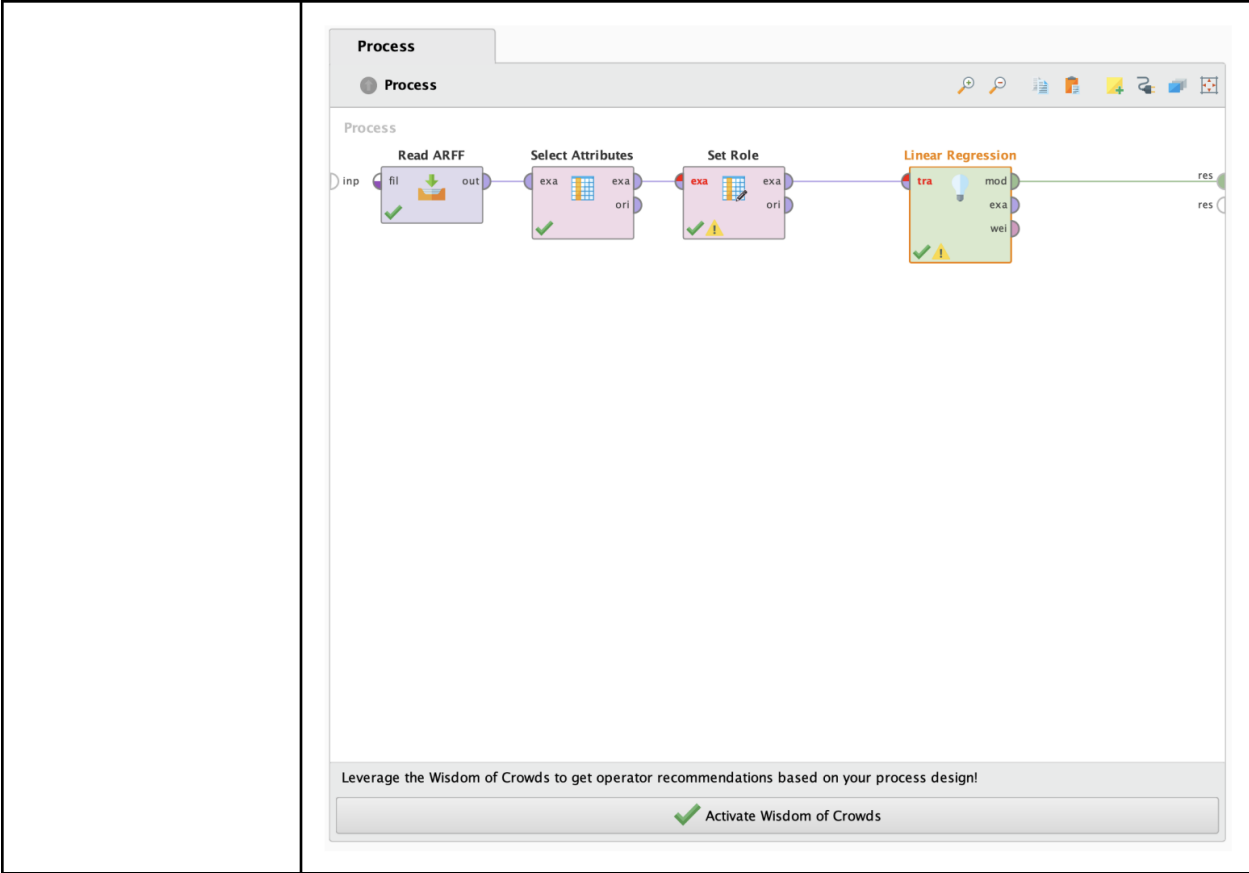
Sample screenshot:



RapidMiner

Q3 (1.5 points). Bring in the shells.arff data, and only work with these 4 params: length, diameter, height,num_rings. Do a linear regression to predict num_rings, from length, diameter, and height. Question: what is the equation? Include a screenshot. Note that you need a 'Set Role' node where you would set num_rings to be a "label", before doing the regression. The regression itself would be done using a 'Linear Regression' operator.

POINTS	DESCRIPTION																																								
1 point	<p>The equation in README.txt:</p> $\text{num_rings} = -11.933 * \text{length} + 25.766 * \text{diameter} + 20.358 * \text{height} + 2.836$ <p>OR a screenshot of it:</p> 																																								
0.5 point	<p>Screenshot showing the student's work (any of the following counts):</p>  <table><thead><tr><th>Attribute</th><th>Coefficient</th><th>Std. Error</th><th>Std. Coefficient</th><th>Tolerance</th><th>t-Stat</th><th>p-Value</th><th>Code</th></tr></thead><tbody><tr><td>length</td><td>-11.933</td><td>2.064</td><td>-0.444</td><td>0.078</td><td>-5.781</td><td>0.000</td><td>****</td></tr><tr><td>diameter</td><td>25.766</td><td>2.539</td><td>0.793</td><td>0.094</td><td>10.147</td><td>0</td><td>****</td></tr><tr><td>height</td><td>20.358</td><td>1.737</td><td>0.264</td><td>0.319</td><td>11.719</td><td>0</td><td>****</td></tr><tr><td>(Intercept)</td><td>2.836</td><td>0.186</td><td>?</td><td>?</td><td>15.243</td><td>0</td><td>****</td></tr></tbody></table>	Attribute	Coefficient	Std. Error	Std. Coefficient	Tolerance	t-Stat	p-Value	Code	length	-11.933	2.064	-0.444	0.078	-5.781	0.000	****	diameter	25.766	2.539	0.793	0.094	10.147	0	****	height	20.358	1.737	0.264	0.319	11.719	0	****	(Intercept)	2.836	0.186	?	?	15.243	0	****
Attribute	Coefficient	Std. Error	Std. Coefficient	Tolerance	t-Stat	p-Value	Code																																		
length	-11.933	2.064	-0.444	0.078	-5.781	0.000	****																																		
diameter	25.766	2.539	0.793	0.094	10.147	0	****																																		
height	20.358	1.737	0.264	0.319	11.719	0	****																																		
(Intercept)	2.836	0.186	?	?	15.243	0	****																																		



*** Incorrect number / missing attributes, [-1 point for Q3](#)

R

Q4 (0.5 point). Take a screenshot of the entire RStudio IDE that shows the code, console output and the map (sufficiently zoomed in) with the hull and centroid markers visible, for submission.

POINTS	DESCRIPTION
0.5 point	<p>Screenshot shows code, console, sufficiently zoomed in map, with visible hull (circles) and centroid markers (dark blue flag icon)</p> <p>Plotting just the convex hull points is what's expected - not all the collected cords. [no point deduction]</p> <p>Also, no need to indicate the area - since this question is about centroid (of the hull) [no point deduction]</p> <p>.</p>

