


## 6.1 Multiple Linear Regression

Wednesday, 26 October 2022 13:07

Summary	<ul style="list-style-type: none"> <li>Multiple Linear Regression</li> <li><math>\bar{R}^2</math> and <math>s_e</math></li> <li>Calibration Plot</li> <li>Marginal and Partial Slopes</li> <li>Path Diagram and Collinearity</li> <li>F-statistic</li> <li>t-statistic</li> <li>Prediction interval</li> </ul>
<ul style="list-style-type: none"> <li>What's the basic difference between multiple and simple linear regression?</li> </ul>	<h3>Multiple Regression</h3> <ul style="list-style-type: none"> <li>In simple linear regression, there's <b>one explanatory variable</b> and one response (dependent) variable.</li> <li>In multiple linear regression, we have <b>multiple explanatory variables</b> and one response variable.</li> </ul>
	<h3>The Multiple Regression Model</h3> <ul style="list-style-type: none"> <li>Use multiple regression to describe the relationship between several explanatory variables and the response.</li> <li>Multiple regression separates the effects of each explanatory variable on the response and reveals which really matter.</li> <li>We will study the effect of each of these explanatory variable on the response variable.</li> </ul>
<ul style="list-style-type: none"> <li>Full form of MLR and MRM.</li> <li>What is <math>k</math> in MRM?</li> </ul>	<h3>The Multiple Regression Model</h3> <ul style="list-style-type: none"> <li>Multiple regression model (MRM): model for the association in the population between multiple explanatory variables and a response.</li> <li><b><math>k</math>: the number of explanatory variables in the multiple regression (<math>k = 1</math> in simple regression).</b></li> <li>We can call it either: MLR / MRM.             <ul style="list-style-type: none"> <li>MLR : Multiple Linear Regression</li> <li>MRM: Multiple Regression Model</li> </ul> </li> </ul>
<ul style="list-style-type: none"> <li>Equation of MRM</li> <li>Assumptions about the error term</li> </ul>	<h3>The Multiple Regression Model</h3> <p>The <b>response <math>Y</math></b> is linearly related to <b><math>k</math> explanatory variables <math>X_1, X_2, \dots, X_k</math></b> by the equation</p> $Y = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k + \epsilon$ $\epsilon \sim N(0, \sigma_\epsilon^2)$ <ul style="list-style-type: none"> <li><i>The unobserved errors in the model</i> <ol style="list-style-type: none"> <li>are independent of one another,</li> <li>have equal variance, and</li> <li>are normally distributed around the regression equation.</li> </ol> </li> <li>We will estimate <math>\beta_0, \beta_1, \dots, \beta_k</math>.</li> </ul>
<ul style="list-style-type: none"> <li>Difference in the error terms between MLR and SLR</li> </ul>	<h3>The Multiple Regression Model</h3> <ul style="list-style-type: none"> <li><b>While the SRM bundles all but one explanatory variable into the error term, multiple regression allows for the inclusion of several variables in the</b></li> </ul>

	<p>model.</p> <ul style="list-style-type: none"> <li>In the MRM, residuals departing from normality may suggest that an important explanatory variable has been omitted.</li> </ul> $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k + \varepsilon$ <ul style="list-style-type: none"> <li>The impact of all other explanatory variables, which was not considered in SLR, will get lumped into the error term of SLR.</li> </ul>
<ul style="list-style-type: none"> <li><math>\bar{R}^2</math></li> <li><math>s_e</math></li> <li>How do they behave when an explanatory variable is added?</li> </ul> <p>□? If <math>\bar{R}^2 &lt; R^2</math>, and <math>\bar{R}^2</math> keeps increasing as you add more explanatory variables, then does this <math>\bar{R}^2 &lt; R^2</math> relationship always hold? Does it mean that no matter how many explanatory variables you add, you'll never cross <math>R^2</math>?</p>	<h2>Interpreting Multiple Regression</h2> <ul style="list-style-type: none"> <li>R-squared and <math>s_e</math></li> <li><math>\bar{R}^2</math> is known as the adjusted R-squared. It adjusts for both sample size <math>n</math> and model size <math>k</math>. It is always smaller than <math>R^2</math>.</li> <li>The residual degrees of freedom (<math>n-k-1</math>) is the divisor of <math>s_e</math>. <math>\bar{R}^2</math> and <math>s_e</math> move in opposite directions when an explanatory variable is added to the model (<math>\bar{R}^2</math> goes up while <math>s_e</math> goes down).</li> </ul> <p>Adjusted <math>R^2 &lt; R^2</math></p> <ul style="list-style-type: none"> <li>Adjusted <math>R^2</math> gives us a slightly more realistic picture of what is the combined explanatory power of all these explanatory variables.</li> </ul> <p><math>s_e</math> is the estimate of [Equation] Standard Deviation of the error term, standard error in estimating the slopes</p> <ul style="list-style-type: none"> <li>In general, we want a large value of [Equation], and we want a smaller value of <math>s_e</math>, i.e., if you add [Equation] and [Equation]</li> </ul>
<ul style="list-style-type: none"> <li>State Adjusted R squared formula</li> </ul>	<p>Adjusted R squared formula (not in the lectures but is asked in the questions)</p> <p>[Equation]</p>
<ul style="list-style-type: none"> <li>Define Calibration Plot</li> <li>What does [Equation] represent in: <ul style="list-style-type: none"> <li>SLR</li> <li>MLR</li> </ul> </li> </ul>	<h2>Interpreting Multiple Regression</h2> <ul style="list-style-type: none"> <li>Calibration Plot</li> <li>Calibration plot: scatterplot of the response <math>\mathcal{Y}</math> on the fitted values <math>\hat{\mathcal{Y}}</math></li> <li><math>R^2</math> is the correlation between <math>\hat{\mathcal{Y}}</math> and <math>\mathcal{Y}</math>; the tighter data cluster along the diagonal line in the calibration plot, the larger the <math>R^2</math> value.</li> <li>In SLR, <math>R</math> represented the correlation between [Equation] and <math>\mathcal{Y}</math>.</li> <li>In MLR, <math>R</math> represents the coefficient of correlation between observed value(<math>y</math>) of the response variable and the fitted/predicted value( [Equation] ) of the response variable.</li> </ul>
<ul style="list-style-type: none"> <li>Define <ul style="list-style-type: none"> <li>Marginal slope</li> <li>Partial slope</li> </ul> </li> <li>When is Marginal slope [Equation] Partial slope?</li> </ul>	<h2>Interpreting Multiple Regression</h2> <ul style="list-style-type: none"> <li>Marginal and Partial Slopes</li> <li>Partial slope: slope of an explanatory variable in a multiple regression that statistically excludes the effects of other explanatory variables.</li> <li>Marginal slope: slope of an explanatory variable in a simple regression.</li> </ul>

	<ul style="list-style-type: none"> <li>Partial and marginal slopes only agree when the explanatory variables are uncorrelated.</li> <li>In SLR: [Equation] Marginal Slope. <ul style="list-style-type: none"> <li>Change in Y variable with one unit change in X variable.</li> </ul> </li> <li>In MLR: [Equation] Partial Slope. <ul style="list-style-type: none"> <li>Change in Y variable with one unit change in X variable keeping all the other X variables constant.</li> </ul> </li> <li>Ideally, we want explanatory variables which are orthogonal to each other, <i>i.e.</i>, which are independent of each other. <ul style="list-style-type: none"> <li>If the variables are truly independent of each other, then the marginal slope and the partial slope will have the same value. But it happens rarely.</li> </ul> </li> </ul>
<ul style="list-style-type: none"> <li>What is a path diagram? <ul style="list-style-type: none"> <li>What is the total effect of <math>X</math> on <math>Y</math>?</li> <li>It is represented by the _____ slope?</li> </ul> </li> <li>Define collinearity</li> </ul>	<h2>Interpreting Multiple Regression</h2> <ul style="list-style-type: none"> <li>Path Diagram <ul style="list-style-type: none"> <li>Path diagram: schematic drawing of the relationships among the explanatory variables and the response.</li> <li>Collinearity: very high correlations among the explanatory variables that make the estimates in multiple regression uninterpretable.</li> <li>Collinearity: When the explanatory variables are not independent, <i>i.e.</i>, when they are correlated. It is a situation that represents very high correlation amongst the explanatory variables. Sometimes it is so severe that it actually makes the estimates of MLR very difficult to interpret.</li> <li>Path Diagram:  </li> <li>We say that [Equation] impacts <math>Y</math> in two different ways: <ul style="list-style-type: none"> <li>There's a <b>direct effect</b> of [Equation] on <math>Y</math>, and</li> <li>There's also an <b>indirect effect</b> of [Equation] on [Equation] as it impacts <math>X_2</math>, because they are correlated, that in turn impacts <math>Y</math>.</li> <li>[Equation]</li> </ul> </li> <li>Total effect of [Equation] on [Equation] <u>is represented in the Marginal Slope.</u></li> </ul> </li> </ul>
Slides that were not covered in lectures	
<ul style="list-style-type: none"> <li>Errors in MRM satisfy what conditions?</li> </ul>	<h2>Checking Conditions</h2> <p>Conditions for Inference</p> <ul style="list-style-type: none"> <li>Use the residuals from the fitted MRM to check that the errors in the model <ul style="list-style-type: none"> <li>are independent;</li> <li>have equal variance; and</li> <li>follow a normal distribution.</li> </ul> </li> </ul>
<ul style="list-style-type: none"> <li>How do you calculate F-Statistic in MRM?</li> <li>What null hypothesis do we assume in MRM?</li> </ul>	<h2>Inference in Multiple Regression</h2> <ul style="list-style-type: none"> <li>Inference for the Model: <i>F</i>-test</li> <li>The <i>F</i>-Statistic <math display="block">F = \frac{R^2}{1 - R^2} \frac{n - k - 1}{k}</math> </li> </ul>

	<p>is used to test the null hypothesis that all slopes are equal to zero,</p> $H_0: \beta_1 = \beta_2 = 0$ <div>[Equation]</div> <div>[Equation]</div>
<ul style="list-style-type: none"><li>• What is the purpose of t-statistic in MRM?</li><li>• How is t-statistic calculated?</li></ul>	<h2>Inference in Multiple Regression</h2> <h3>Inference for One Coefficient</h3> <ul style="list-style-type: none"><li>▪ The <math>t</math>-statistic is used to test each slope using the null hypothesis <math>H_0: \beta_j = 0</math>.</li><li>▪ The <math>t</math>-statistic is calculated as</li></ul> $t_j = \frac{b_j - 0}{se(b_j)}$ <div>[Equation]</div>
<ul style="list-style-type: none"><li>• An approximate <div>[Equation]</div> prediction interval is given by?</li></ul>	<h2>Inference in Multiple Regression</h2> <ul style="list-style-type: none"><li>• <b>Prediction Intervals</b></li><li>▪ An approximate 95% prediction interval is given by <math>\hat{y} \pm 2s_e</math>.</li><li>▪ For example, the 95% prediction interval for price.</li></ul>