# LINEAR REGRESSION FORMULA SHEET
*Note: This is not a cheat sheet. You cannot use this during the exam.*

## EQUATIONS TO MEASURE RELATIONSHIPS BETWEEN VARIABLES

| | | | |
|---|---|---|---|
| Sample mean of x | $\overline{x}$ | $\dfrac{1}{n}\sum x_i$ | Average $x_i$ <br> $E(x_i) = \overline{x}, \forall_i$ |
| Sample variance of x | $s_x^2$ <br> var(x) | $\dfrac{\sum(x_i - \overline{x})^2}{n-1} = \dfrac{S_{xx}}{n-1}$ | Variation in $x_i$ about $\overline{x}$ |
| Sample standard deviation of x | $s_x$ <br> sd(x) | $\sqrt{s_x^2}$ | Average distance of $x_i$ from $\overline{x}$ |
| Sample covariance of x and y | $s_{xy}$ <br> cov(x,y) | $\dfrac{\sum(x_i - \overline{x})(y_i - \overline{y})}{n-1} = \dfrac{S_{xy}}{n-1}$ | Measures direction of relationship: if y goes up or down when x goes up or down. Is <u>not</u> predictive! |
| Sample estimate of the correlation between x and y | $r_{xy}$ <br> corr(x,y) | $\dfrac{s_{xy}}{s_x s_y} = \dfrac{\text{cov}(x,y)}{s_x s_y}$ | Correlation is the standardized (unit-less) version of covariance. The correlation ranges between -1 and +1, inclusively. |
| Not quite the variance (uppercase 'S') | $S_{xx}$ | $\sum(x_i - \overline{x})^2 = \sum(x_i - \overline{x})\cdot x_i$ | Used in equations below. This is a mathematical "shorthand"; it doesn't have meaning in and of itself. |
| Not quite the covariance (uppercase 'S') | $S_{xy}$ | $\sum(x_i - \overline{x})(y_i - \overline{y}) = \sum(x_i - \overline{x})\cdot y_i$ | Used in equations below. This is a mathematical "shorthand"; it doesn't have meaning in and of itself. |

## BIVARIATE REGRESSION EQUATIONS

$y_i$ = sample value; $\overline{y}$ = average of the sample values; $\hat{y}$ = value predicted by model; $e_i = y_i - \hat{y}$ = residual

| | | | |
|---|---|---|---|
| Two-variable linear model | | Population model: $y = \alpha + \beta x + \varepsilon$ <br> Sample model: $y = a + bx + e$ <br> Estimated model: $\hat{y} = a + bx$ | ε = unobservable errors <br> e = residual errors <br> No error term in the estimate! |
| OLS estimate for α | a | $\overline{y} - b\overline{x}$ | a is the y-intercept |
| OLS estimate for β | b | $\dfrac{S_{xy}}{S_{xx}} = \dfrac{s_{xy}}{s_x^2}$ | b is the slope of the regression line. A one-unit change in *x* causes a change of *b* in *y*. *b* is predictive! |
| Total sum of squared deviations <br> **SST = SSE + SSR** | SST <br> (TSS) | $\sum(y_i - \overline{y})^2$ | How much the actual sample points vary from the sample average. |
| Explained sum of squared deviations <br> **SSE = SST - SSR** | SSE <br> (ESS) | $\sum(\hat{y}_i - \overline{y})^2 = b^2 \cdot \sum(x_i - \overline{x})^2 = b^2 \cdot S_{xx}$ | How much the predicted points vary from the sample average. |
| Residual sum of squared deviations. <br> **SSR = SST - SSE** | SSR <br> (RSS) | $\sum e_i^2 = \sum(y_i - \hat{y}_i)^2$ | How much the actual sample points vary from the predicted points. |
| Coefficient of determination | $R^2$ | $\dfrac{SSE}{SST} = 1 - \dfrac{SSR}{SST}$ | Proportion of the total sum of squares that is explained by *x*; how well *x* predicts *y*. Convert to percentage. |
| Standard error of the entire regression <br><br> *An <u>estimate</u> of the population std error.* | se <br><br> root MSE | $\sqrt{\dfrac{SSR}{N-(k+1)}}$ | N = total observations <br> k = # of independent variables <br><br> <u>Note</u>: This is **not** the standard error of the population, but an estimate of the population standard error. |
| Estimated variance of the error | $se^2$ | $\dfrac{SSR}{N-(k+1)}$ | Measure of variance in *y* not explained by the variance in *x* for the entire model. |

| Variance of the estimator b | Var(b) | $$\frac{se^2}{S_{xx}} = \frac{se^2}{(n-1)s_x^2}$$ | |
|---|---|---|---|
| Standard error of the estimator b | $se_b$ | $\sqrt{\operatorname{var}(b)}$ | Standard error of the estimator for the x coefficient – i.e. the width of the estimator's sampling distribution. |
| Variance of the estimator a | Var(a) | $$se^2 \cdot \left( \frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)$$ | |
| Standard error of the estimator a | $se_a$ | $\sqrt{\operatorname{var}(a)}$ | Standard error of the estimator for the y intercept – i.e. the width of the estimator's sampling distribution. |
| Using statistical inference, how to tell if b is significantly different from β | t | $$\frac{b - \beta}{se_b}$$ | Test statistic for testing the significance of the predicted x effect on y. Usually, set β=0 |

## MULTIVARIATE REGRESSION EQUATIONS

| Multivariate linear model | | Population model: $y = \alpha + \beta x + \gamma z + \varepsilon$ <br> Sample model: $y = a + bx + gz + e$ <br> Estimated model: $\hat{y} = a + bx + gz$ | ε = unobservable errors <br> e = residual errors <br> No error term in the estimate! |
|---|---|---|---|
| OLS estimate for α | a | $\bar{y} - b\bar{x} - g\bar{z}$ | |
| OLS estimate for β | b | $$\frac{S_{xy}S_{zz} - S_{zy}S_{xz}}{S_{xx}S_{zz} - S_{xz}^2}$$ | Holding constant all other explanatory variables in the model, a one-unit change in x causes a change of b in y. Depends on the relationship between x and z, and z and y. |
| OLS estimate for γ | g | $$\frac{S_{zy}S_{xx} - S_{xy}S_{xz}}{S_{xx}S_{zz} - S_{xz}^2}$$ | Ceteris paribus, a one-unit change in z causes a change of g in y. Depends on the relationship between z and x, and x and y |
| Variance of the estimator b <br><br> **Precision** <br> This is also known as precision or efficiency: smaller var(b) means higher precision. | Var(b) | $$\frac{se^2}{n \cdot s_x^2} \cdot \frac{1}{1 - r_{x,z}^2}$$ | $r_{x,z}^2$ is the sample correlation between x & z. <br><br> If x and z are not correlated, we recover the bivariate formula (almost – need 'n-1' instead of 'n'). If x and z are perfectly correlated, the standard error blows up. <br><br> In general, as x and z approach perfect correlation, se blows up. |
| Standard error of the estimator b | $se_b$ | $\sqrt{\operatorname{var}(b)}$ | See comments above. |
| Omitted variable bias | E(b) | $$\beta + \gamma \cdot \left( \frac{\operatorname{cov}(x,z)}{\operatorname{var}(x)} \right)$$ | The second term is the effect of x on z. It's magnitude of the bias. <br><br> E(b) = β if γ = 0 or if cov(x,z) = 0 <br><br> Not solving for numbers, just positive or negative signs which tell you "understated" or "overstated" effect of the omitted variable <br><br> Note: var(x) is always positive. |