

Calculation of marginal probability distributions and probabilities:

$$f_X(x_1) = \sum_y f_{X,Y}(x_1, y) = f_{X,Y}(x_1, y_1) + f_{X,Y}(x_1, y_2) + \dots \text{ [or] } f_X(x_1) = \int_y f_{X,Y}(x_1, y) dy$$

independence: For A, B: when $P(A|B) = P(A)$ and $P(B|A) = P(B)$ $\left[P(A|B) = \frac{P(A \cap B)}{P(B)} \right]$

$$Cov(X, Y) = E((X - \mu_X)(Y - \mu_Y)) \quad Cov > 0: \text{ on average } X, Y \text{ above (below) mean together.}$$

$$Cov < 0: \text{ on average, when } X \text{ is above (below) mean, } Y \text{ is below (above).}$$

$$(X, Y): \text{ independent} \Rightarrow Cov(X, Y) = 0 \quad Cov(X, Y) = 0 \not\Rightarrow (X, Y): \text{ independent}$$

$$Cov((a_1X + b_1), (a_2Y + b_2)) = a_1a_2Cov(X, Y) \quad Var(aX + bY) = a^2Var(X) + b^2Var(Y) + 2abCov(X, Y)$$

$$Corr(X, Y) = \rho_{XY} = \frac{Cov(X, Y)}{\sigma_X\sigma_Y} \quad -1 \leq \rho \leq 1$$

$$(X, Y): \text{ independent} \Rightarrow Corr(X, Y) = 0 \quad Corr(X, Y) = 0 \not\Rightarrow (X, Y): \text{ independent}$$

$$E(X) = \sum_{i=1}^k x_i f_i(x_i) \text{ [or] } \int_{-\infty}^{\infty} xf(x)dx \quad E(g(x)) = \sum_{i=1}^k g(x_i) f_i(x_i) \text{ [or] } \int_{-\infty}^{\infty} g(x)xf(x)dx$$

$$E(c) = c, \quad E(aX + b) = aE(X) + b, \quad E(X + Y) = E(X) + E(Y)$$

$$Var(X) = E(X - \mu)^2 = \sigma^2 = E(X^2) - (E(X))^2 \quad Var(c) = 0, \quad Var(aX + b) = a^2Var(X)$$

$$Z = \frac{X - \mu}{\sigma} \sim N(0, 1) \quad \Phi(c) = z - \text{table reading for number } c$$

$$P(Z > z) = 1 - \Phi(z), \quad P(Z < -z) = P(Z > z), \quad P(a \leq Z \leq b) = \Phi(b) - \Phi(a)$$

$$P(|Z| > c) = P(Z > c) + P(Z < -c) = 2(1 - \Phi(c))$$

$$E(Y|X = x) = \sum_j y_j f_{Y|X}(y_j|x) \quad \text{or} \quad \int y_i f_{Y|X}(y_i|x) dy$$

Simple linear regression model. $y = \beta_0 + \beta_1 x + u$ assuming $E\langle u|x \rangle = 0 \Rightarrow E\langle y|x \rangle = \beta_0 + \beta_1 x$
 $E(u) = 0 \Rightarrow E(y - \beta_0 - \beta_1 x) = 0$ $Cov(x, u) = E(xu) = 0 \Rightarrow E(x(y - \beta_0 - \beta_1 x)) = 0$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad \hat{\beta}_1 = \frac{\sum_{i=1}^n x_i(y_i - \bar{y})}{\sum_{i=1}^n x_i(x_i - \bar{x})} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad \hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x \quad \hat{u}_i = y_i - \hat{y}_i$$

$$\sum_{i=1}^n \hat{u}_i = 0 \quad \sum_{i=1}^n x_i \hat{u}_i = 0 \quad \text{The point } (\bar{x}, \bar{y}) \text{ is always on the regression line.}$$

$$\text{bias}(\hat{\beta}_1) = E(\hat{\beta}_1) - \beta_1 \quad E(\hat{\beta}_1) = \beta_1 + E\left(\frac{\sum u_i x_i}{\sum x_i^2}\right) = \beta_1 + \frac{cov(u, x)}{var(x)} = \beta_1 + 0$$

$$\therefore \text{bias}(\hat{\beta}_1) = \beta_1 - \beta_1 = 0 \quad \blacksquare$$

$$SST = \text{total sum of squares} = \text{total variation of } y \text{ about its mean } \sum (y_i - \bar{y})^2$$

$$SST_x = \text{total variation of } x \text{ about its mean } \sum (x_i - \bar{x})^2 \quad SSR = \sum (y_i - \hat{y}_i)^2$$

$$SSE = \text{explained sum of squares} = \text{predicted variation of } y \text{ about } \bar{y} = \sum (\hat{y}_i - \bar{y})^2$$

$$R^2 = \frac{SSE}{SST} = 1 - \frac{SSR}{SST} \quad \text{shows fraction of variation in } y \text{ explained by } x$$

$$E(\hat{\beta}_0) = \beta_0 \quad E(\hat{\beta}_1) = \beta_1 \quad Var(\hat{\beta}_0) = \frac{\sigma^2 \sum x_i^2}{nSST_x} \quad Var(\hat{\beta}_1) = \frac{\sigma^2}{SST_x}$$

$$[SST_x = \sum (x_i - \bar{x})^2] \text{ Unbiased estimator of population variance: } \hat{\sigma}^2 = \frac{1}{n-2} \sum \hat{u}_i^2 = \frac{SSR}{n-2}$$

$$SER = \hat{\sigma} = \sqrt{\hat{\sigma}^2} \quad SE(\hat{\beta}_0) = \sqrt{Var(\hat{\beta}_0)} = \sqrt{\frac{\sigma^2 \sum x_i^2}{nSST_x}} \quad SE(\hat{\beta}_1) = \sqrt{Var(\hat{\beta}_1)} = \sqrt{\frac{\sigma^2}{SST_x}}$$

modify dependent: $cy = c\beta_0 + c\beta_1x$ modify independent: $y = \beta_0 + \left(\frac{\beta_1}{k}\right)(kx)$

lin – lin: $y = \beta_0 + \beta_1x$ $\beta_1 =$ unit change in y for a 1 unit change in x

lin – log: $y = \beta_0 + \beta_1 \ln x$ $\frac{\beta_1}{100} =$ unit change in y for a 1% change in x

log – lin: $\ln y = \beta_0 + \beta_1x$ $100\beta_1 =$ % change in y for a 1 unit change in x

log – log: $\ln y = \beta_0 + \beta_1 \ln x$ $\beta_1 =$ % change in y for a 1% change in x

Partialling: For regression: $y = \beta_0 + \beta_1x_1 + \beta_2x_2$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n \hat{r}_{i1}y_i}{\sum_{i=1}^n \hat{r}_{i1}^2}, \text{ where } \hat{r}_{i1} = \text{residual from OLS regression of } x_1 \text{ on } x_2$$

These residuals are uncorrelated with x_2 , and are estimates of x_1 after effects of x_2 have been 'partialled out.' Therefore, $\hat{\beta}_1$ is simple relationship between y and x_1 after x_2 has been partialled out.

Stat prop mult: Still linear in parameters, still random sampling, still zero conditional mean, still homoscedastic. Specific to multiple: no perfect multicollinearity in independent variables

OVB: $y = \beta_0 + \beta_1x_1 + \beta_2x_2$ but we estimate $y = \widetilde{\beta}_0 + \widetilde{\beta}_1x_1$ and $\widetilde{x}_2 = \widetilde{\delta}_0 + \widetilde{\delta}_1x_1$ then $\widetilde{\beta}_1 = \widehat{\beta}_1 + \widehat{\beta}_2\widetilde{\delta}_1$ then $E(\widetilde{\beta}_1) = \beta_1 + \beta_2\widetilde{\delta}_1$

Including irrelevant vars maintains unbiasedness, but precision of estimates is bad.

Multicollinearity: correlation among expl vars leads to bias in coefficients

$$Var(\widehat{\beta}_j) = \frac{\sigma^2}{SST_j(1-R_j^2)} \quad \text{where } SST_j = \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2 \quad \text{and}$$

$R_j^2 = R^2$ of regression of x_j on all other expl vars

$$\hat{y} = \widehat{\beta}_0 + \widehat{\beta}_1x_1 + \widehat{\beta}_2x_2 \quad \text{and} \quad \tilde{y} = \widetilde{\beta}_0 + \widetilde{\beta}_1x_1 \quad Var(\widetilde{\beta}_1) = \frac{\sigma^2}{SST_1(1-R_1^2)} \quad Var(\widehat{\beta}_1) = \frac{\sigma^2}{SST_1}$$

when $\beta_2 \neq 0$ $\widetilde{\beta}_1$ is biased $\widehat{\beta}_1$ is unbiased and $Var(\widetilde{\beta}_1) \leq Var(\widehat{\beta}_1)$

when $\beta_2 = 0$ $\widetilde{\beta}_1, \widehat{\beta}_1$ both unbiased and $Var(\widetilde{\beta}_1) \leq Var(\widehat{\beta}_1)$

Gauss-Markov Theorem: Given these five assumptions: 1 In population regression, y is linearly related to x and u . 2. Random sample of size n according to population model. 3. Sample outcomes of x are not all the same. 4. Expectation of error term is 0 given any value of x . 5. Error term has the same variance for any value of x . In a linear regression model, with expected errors of zero, uncorrelated errors, and equal variances, the Best Linear Unbiased Estimators are given by the OLS estimators.

Non-linear relationships in linear regressions: add squared term: $y = \beta_0 + \beta_1x + \beta_2x^2 + u$