## Some formulæ for ENGI 4421 Probability and Statistics

On these pages are many more formulae than will fit comfortably on your allocation of two sheets in the final examination. Examine these suggestions together with your own notes, in order to construct a more concise set of formula sheets that will work best for *you*.

**Descriptive Statistics**

median $=$ middle value in ordered set

mode $=$ most frequently occurring value

$$\text{mean} = \bar{x} = \frac{1}{n}\sum x$$

$$s^2 = \frac{1}{n-1}\sum (x-\bar{x})^2 \quad \text{or} \quad s^2 = \frac{n\sum x^2 - \left(\sum x\right)^2}{n(n-1)}$$

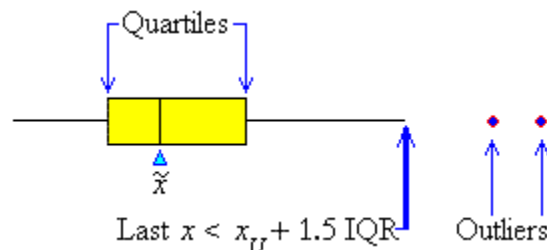From a frequency table ( $f_i$ occurrences of distinct value $x_i$ observed):

$$\text{mean} = \bar{x} = \frac{\sum f_i x_i}{\sum f_i}, \qquad \left(n = \sum f_i\right)$$

$$s^2 = \frac{1}{n-1}\sum f_i (x_i - \bar{x})^2 \quad \text{or} \quad s^2 = \frac{n\sum f_i x_i^2 - \left(\sum f_i x_i\right)^2}{n(n-1)}$$

Quartiles: in ordered set $\{\, x_1, x_2, \ldots, x_n \,\}$ are $x_{(n+1)/4}$ and $x_{3(n+1)/4}$

(interpolate as necessary)

Histogram: (area of a bar) = (rel. freq. of that interval)

Boxplot:



**Decision Tree**

[See examples in problem sets, past exams and additional exercises.]

**Counting Techniques:**
Number of distinct ways of selecting $r$ objects from $n$ objects:

with replacement, order of selection matters $= n^r$

without replacement, order of selection matters $= {}^nP_r = n!/(n-r)!$

without replacement, order of selection doesn't matter $= {}^nC_r = n!/(r!\,(n-r)!)$

$${}^nP_0 = {}^nC_0 = {}^nC_n = 1, \qquad {}^nP_1 = {}^nC_1 = {}^nC_{n-1} = n, \qquad {}^nP_{n-1} = {}^nP_n = n!, \qquad {}^nC_{n-r} = {}^nC_r$$

Number of distinct ways to partition $n$ objects into piles of $r_1, r_2, \ldots, r_k$ objects:

$$\begin{pmatrix} n \\ r_1\ r_2\ \cdots\ r_k \end{pmatrix} = \frac{n!}{r_1!\,r_2!\cdots r_k!}$$

**Laws of Probability**

Odds [on] $= r$ and probability $= p$:

$$r = \frac{p}{1-p}, \quad p = \frac{r}{1+r}$$

deMorgan's laws:     $\sim(A \cup B) = \sim A \cap \sim B$     and     $\sim(A \cap B) = \sim A \cup \sim B$

General addition law of probability:          $P[A \vee B] = P[A] + P[B] - P[A \wedge B]$

General multiplication law of probability:     $P[A \wedge B] = P[A] \times P[B \mid A] = P[B] \times P[A|B]$

Events $A, B$ are independent iff $P[A \wedge B] = P[A] \times P[B]$

Events $A, B$ are incompatible (mutually exclusive) iff $P[A \wedge B] = 0$

If $\{ E_1, E_2, \ldots, E_n \}$ is a partition, then **Total Probability Law:**

$$P[A] = \sum_i P[A \cap E_i] = \sum_i P[A \mid E_i]P[E_i]$$

and **Bayes' Theorem:**

$$P[E_k \mid A] = \frac{P[A \mid E_k]P[E_k]}{\sum_i P[A \mid E_i]P[E_i]}$$

---

**Conditions for $p(x)$ to be a probability mass function (p.m.f.):**

$$p(x) \geq 0 \quad \forall x \qquad \text{and} \quad \sum_{\text{all } x} p(x) = 1 \ \left(\text{coherence condition}\right)$$

**Conditions for $f(x)$ to be a probability density function (p.d.f.):**

$$f(x) \geq 0 \quad \forall x \qquad \text{and} \quad \int_{-\infty}^{\infty} f(x)\,dx = 1 \ \left(\text{coherence condition}\right)$$

**Cumulative Distribution Function (c.d.f.):**

Discrete:          Continuous:

$$F(x) = \sum_{y \leq x} p(x) \qquad F(x) = \int_{-\infty}^{x} f(x)\,dx \qquad P[a < X \leq b] = F(b) - F(a)$$

(staircase)          (ogive)

Median: $F(\tilde{\mu}) = \frac{1}{2}$          Lower quartile: $F(x_L) = \frac{1}{4}$   Upper quartile: $F(x_U) = \frac{3}{4}$

**Expected value:**

$$\mu = \mathrm{E}[X] = \sum_x x \cdot p(x) \quad \text{or} \quad \int_{-\infty}^{\infty} x f(x) \, dx$$

$$\mathrm{E}[h(X)] = \sum_x h(x) \cdot p(x) \quad \text{or} \quad \int_{-\infty}^{\infty} h(x) f(x) \, dx$$

**Variance:** $\qquad \sigma^2 = \mathrm{V}[X] = \mathrm{E}\left[(X - \mu)^2\right] = \mathrm{E}\left[X^2\right] - \left(\mathrm{E}[X]\right)^2$

---

**Joint Probability Mass Functions and Joint Probability Density Functions:**

$$p(x, y) = \mathrm{P}\left[(X = x) \wedge (Y = y)\right]$$

Marginal p.m.f.'s:                                     Marginal p.d.f.'s:

$$p_X(x) = \sum_y p(x, y) \quad \text{and} \quad p_Y(y) = \sum_x p(x, y) \qquad f_X(x) = \int_{-\infty}^{\infty} f(x, y) \, dy$$

Conditional p.m.f.'s:                            and $f_Y(y) = \int_{-\infty}^{\infty} f(x, y) \, dx$

$$p_{Y|X}(y \mid x) = \frac{p(x, y)}{p_X(x)}, \qquad p_{X|Y}(x \mid y) = \frac{p(x, y)}{p_Y(y)}$$

Conditional p.d.f.'s: $\qquad f_{Y|X}(y \mid x) = \dfrac{f(x, y)}{f_X(x)} \quad \text{and} \quad f_{X|Y}(x \mid y) = \dfrac{f(x, y)}{f_Y(y)}$

**Covariance:** $\; \mathrm{Cov}[X, Y] = \mathrm{E}\left[(X - \mu_X)(Y - \mu_Y)\right] = \mathrm{E}[XY] - \mathrm{E}[X]\mathrm{E}[Y]$

$$\mathrm{V}[X] = \mathrm{Cov}[X, X] \geq 0$$

**Correlation:** $\; \rho_{X,Y} = \dfrac{\mathrm{Cov}[X, Y]}{\sqrt{\mathrm{V}[X] \cdot \mathrm{V}[Y]}} ; \qquad -1 \leq \rho \leq +1$

Independence $\Rightarrow \rho = 0$ ;    (but $\rho = 0 \not\Rightarrow$ independence)

---

**Combinations of Random Quantities:**

$$\mathrm{E}\left[\sum_{i=1}^{n} a_i X_i\right] = \sum_{i=1}^{n} a_i \mathrm{E}[X_i], \qquad \mathrm{V}\left[\sum_{i=1}^{n} a_i X_i\right] = \sum_{i=1}^{n}\sum_{j=1}^{n} a_i a_j \mathrm{Cov}[X_i, X_j]$$

$\mathrm{E}[aX + bY] = a\,\mathrm{E}[X] + b\,\mathrm{E}[Y]$   and

$\mathrm{V}[aX + bY] = a^2\,\mathrm{V}[X] + 2ab\,\mathrm{Cov}[X, Y] + b^2\,\mathrm{V}[Y]$

for all constants $a, b$ and all random quantities $X, Y$.

Special case, when $X$ and $Y$ are independent (or at least uncorrelated):

$\mathrm{E}[X \pm Y] = \mathrm{E}[X] \pm \mathrm{E}[Y]$   and   $\mathrm{V}[X \pm Y] = \mathrm{V}[X] + \mathrm{V}[Y]$

$$\mathrm{E}\left[\overline{X}\right] = \mu, \qquad \mathrm{V}\left[\overline{X}\right] = \frac{\sigma^2}{n}$$

**Propagation of Error**

Sample size $n$ leads to estimate $\left( \bar{x} \pm \dfrac{s}{\sqrt{n}} \right)$

Pooling $\left( x_1 \pm \sigma_1 \right)$ with $\left( x_2 \pm \sigma_2 \right)$:   The most precise estimate is $\left( x \pm \sigma \right)$   where

$$x = c_1 x_1 + c_2 x_2 \, , \quad \sigma^2 = c_1^{\,2} \sigma_1^{\,2} + c_2^{\,2} \sigma_2^{\,2} \, , \quad c_1 = \dfrac{\sigma_2^{\,2}}{\sigma_1^{\,2} + \sigma_2^{\,2}} \quad \text{and} \quad c_2 = \dfrac{\sigma_1^{\,2}}{\sigma_1^{\,2} + \sigma_2^{\,2}}$$

If not independent, $Y = \displaystyle\sum_{i=1}^{n} a_i X_i \quad \Rightarrow \quad \sigma_Y \leq \left| a_1 \right| \sigma_1 + \left| a_2 \right| \sigma_2 + \ldots + \left| a_n \right| \sigma_n$

Non-linear: $U = U(X) \quad \Rightarrow \quad \sigma_U \approx \left| \left( \left. \dfrac{dU}{dX} \right|_{X = x_{\text{obs}}} \right) \right| \sigma_X$

Multivariate and independent:

$$U = U\left( X_1, X_2, \ldots, X_n \right) \quad \Rightarrow \quad \sigma_U \approx \sqrt{ \left( \dfrac{\partial U}{\partial X_1} \right)^2 \sigma_1^{\,2} + \left( \dfrac{\partial U}{\partial X_2} \right)^2 \sigma_2^{\,2} + \ldots + \left( \dfrac{\partial U}{\partial X_n} \right)^2 \sigma_n^{\,2} }$$

Multivariate and not independent:

$$\sigma_U \leq \left| \dfrac{\partial U}{\partial X_1} \right| \sigma_1 + \left| \dfrac{\partial U}{\partial X_2} \right| \sigma_2 + \ldots + \left| \dfrac{\partial U}{\partial X_n} \right| \sigma_n$$

**Bernoulli distribution:**

$$p(x) = \begin{cases} p & (x = 1) \\ 1 - p & (x = 0) \end{cases} , \qquad \mu = p \, , \qquad \sigma^2 = p(1 - p)$$

**Discrete Uniform distribution:**

$$p(x) = \begin{cases} \dfrac{1}{n} & (x = x_1, x_2, \cdots, x_n) \\ 0 & (\text{otherwise}) \end{cases} , \qquad \mu = \dfrac{1}{n} \sum x$$

**Binomial distribution   bin($n, p$):**

$$\mathrm{P}\big[ X = x \big] = \mathrm{b}(x; n, p) = {}^{n}C_x \, p^x \left( 1 - p \right)^{n-x} , \quad \mu = np, \quad \sigma^2 = np(1 - p)$$

$$\mathrm{P}\big[ X \leq x \big] = \mathrm{B}(x; n, p) = \mathrm{b}(0; n, p) + \mathrm{b}(1; n, p) + \ldots + \mathrm{b}(x; n, p)$$

Four conditions for $X$ to be binomial:
1) Each trial has a complementary pair of outcomes;
2) $p = \mathrm{P}[\text{success}]$ is the same for all trials;
3) Trials are all independent of each other;
4) The number of trials, $n$, is fixed.

**Hypergeometric distribution $H(n, R, N)$:**

$$P[X = x] = p(x) = \frac{{}^{R}C_x \times {}^{N-R}C_{n-x}}{{}^{N}C_n} \quad \text{(may use binomial instead if } n < 5\% \text{ of } N)$$

**Poisson distribution  Poisson($\lambda$):**

$$p(x; \lambda) = \frac{e^{-\lambda}\lambda^x}{x!} \quad (x = 0, 1, 2, \cdots); \quad \sigma^2 = \mu = \lambda$$

---

**Continuous Uniform distribution  U($a, b$):**

$$f(x) = \frac{1}{b-a} \quad (a \le x \le b); \quad P[X \le x] = F(x) = \begin{cases} 0 & (x < a) \\ \dfrac{x-a}{b-a} & (a \le x \le b) \\ 1 & (x > b) \end{cases}$$

$$\mu = \frac{a+b}{2}, \qquad \sigma^2 = \frac{(b-a)^2}{12}$$

**Normal distribution  N($\mu, \sigma^2$):**

To convert $X \sim N(\mu, \sigma^2)$ to $Z \sim N(0, 1)$:

$$Z = \frac{X - \mu}{\sigma}; \qquad \text{then} \quad P[a < X < b] = P[z_a < Z < z_b] = \Phi(z_b) - \Phi(z_a)$$

**Exponential distribution:**

$$f(x) = \begin{cases} 0 & (x < 0) \\ \lambda e^{-\lambda x} & (x \ge 0) \end{cases}; \quad P[X \le x] = F(x) = \begin{cases} 0 & (x < 0) \\ 1 - e^{-\lambda x} & (x \ge 0) \end{cases}$$

$$P[X \ge x] = e^{-\lambda x} \quad (x \ge 0); \quad \sigma = \mu = \frac{1}{\lambda}$$

**Gamma distribution (Erlang if $r \in \mathbb{N}$ ):**

$$f(t) = \frac{\lambda^r t^{r-1}}{\Gamma(r)} e^{-\lambda t} \quad (t > 0), \quad \mu = \frac{r}{\lambda}, \quad \sigma^2 = \frac{r}{\lambda^2}; \text{ for Erlang, } \Gamma(r) = (r-1)!$$

---

**Central Limit Theorem:**

For any random quantity $X$ with $E[X] = \mu$, $V[X] = \sigma^2$ and sufficiently large sample size $n$

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right) \quad \text{and} \quad Z = \frac{\bar{X} - \mu}{\left(\dfrac{\sigma}{\sqrt{n}}\right)} \sim N(0,1)$$

**Bayesian Confidence Intervals:**

If prior information is $X \sim N\left(\mu_o, \sigma_o^2\right)$ and

data from random sample:    size $n$, mean $= \bar{x}$, standard deviation $= s$, then
Calculate weights

$$w_o = \frac{1}{\sigma_o^2}, \quad w_d = \frac{n}{s^2} \quad \rightarrow \quad w^* = w_o + w_d$$

and updated mean and variance

$$\mu^* = \frac{w_d \bar{x} + w_o \mu_o}{w^*}, \quad \left(\sigma^2\right)^* = \frac{1}{w^*} \quad \rightarrow \quad \text{Posterior dist'n } N\left(\mu^*, \left(\sigma^2\right)^*\right)$$

Then the $(1-\alpha) \times 100\%$ confidence interval estimate for $\mu$ is

$$\mu^* - t_{\alpha/2, n-1} \cdot \sigma^* < \mu < \mu^* + t_{\alpha/2, n-1} \cdot \sigma^*$$

[**Note:** if $\sigma$ is known, then it replaces $s$, and $t$ becomes $z$.
Otherwise, the true number of degrees of freedom on $t$ is actually a number between $n-1$ and $\infty$. The interval shown above is a conservative approximation.]

**Classical Confidence Intervals on $\mu$:**

$$\sigma_o^2 \rightarrow \infty \quad \Rightarrow \quad w_o \rightarrow 0 \quad \Rightarrow \quad \mu^* = \bar{x} \quad \text{and} \quad \left(\sigma^2\right)^* = \frac{s^2}{n}$$

Then the $(1-\alpha) \times 100\%$ confidence interval estimate for $\mu$ is

$$\bar{x} - t_{\alpha/2, n-1} \cdot \frac{s}{\sqrt{n}} < \mu < \bar{x} + t_{\alpha/2, n-1} \cdot \frac{s}{\sqrt{n}}$$

**Classical Confidence Intervals on $p$:**

$$p^* \pm z_{\alpha/2} \sqrt{\frac{p^* q^*}{n^*}} \quad \text{where} \quad p^* = \frac{x^*}{n^*} = \frac{x+2}{n+4}$$

**Classical Confidence Intervals on $\mu_1 - \mu_2$ (large sample sizes & independent):**

$$(\bar{x}_1 - \bar{x}_2) - z_{\alpha/2} \cdot \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} < \mu < (\bar{x}_1 - \bar{x}_2) + z_{\alpha/2} \cdot \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

**Classical Confidence Intervals on $p_1 - p_2$:**

$$p_X^* - p_Y^* \pm z_{\alpha/2} \sqrt{\frac{p_X^* q_X^*}{n_X^*} + \frac{p_Y^* q_Y^*}{n_Y^*}}$$

where $x^* = x+1$, $y^* = y+1$, $n_X^* = n_X + 2$ and $n_Y^* = n_Y + 2$

**Classical Confidence Intervals on $\mu_1 - \mu_2$ (small sample sizes & independent):**

$$\left(\bar{x}_1 - \bar{x}_2\right) - t_{\alpha/2,\,\nu} \cdot \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} \;<\; \mu \;<\; \left(\bar{x}_1 - \bar{x}_2\right) + t_{\alpha/2,\,\nu} \cdot \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

$$\text{where} \quad \nu = \text{INT}\left( \frac{\left(\dfrac{s_1^2}{n_1} + \dfrac{s_2^2}{n_2}\right)^2}{\dfrac{1}{n_1-1}\left(\dfrac{s_1^2}{n_1}\right)^2 + \dfrac{1}{n_2-1}\left(\dfrac{s_2^2}{n_2}\right)^2} \right)$$

**Classical Confidence Intervals on $\mu_1 - \mu_2$ (not independent):**

Paired if equal-size samples are pairs of observations on the *same set of individuals*.
Unpaired otherwise – but must be uncorrelated.

---

**Classical Hypothesis Tests:**

$\mathcal{H}_o : \mu = \mu_o$. Select the appropriate $\mathcal{H}_A$ and level of significance $\alpha$.

The burden of proof is on $\mathcal{H}_A$.

Two tailed test, $\mathcal{H}_A : \mu \neq \mu_o$:

Method 1: Evaluate $\quad c_L, c_U = \mu_o \pm t_{\alpha/2,\,n-1} \cdot \dfrac{s}{\sqrt{n}}$

Reject $\mathcal{H}_o$ iff $\bar{x} < c_L$ or $\bar{x} > c_U$

Method 2: Reject $\mathcal{H}_o$ iff

$$\left| t_{\text{obs}} \right| > t_{\alpha/2,\,n-1}, \quad \text{where} \quad t_{\text{obs}} = \frac{\bar{x} - \mu_o}{\left(\dfrac{s}{\sqrt{n}}\right)}$$

Method 3: Reject $\mathcal{H}_o$ iff $\text{P}\left[\,\left|T\right| > \left|t_{\text{obs}}\right|\,\right] < \alpha$, where $t_{\text{obs}} = \dfrac{\bar{x} - \mu_o}{\left(\dfrac{s}{\sqrt{n}}\right)}$

Upper tailed test, $\mathcal{H}_A : \mu > \mu_o$:

Method 1: Evaluate $\quad c = \mu_o + t_{\alpha,\,n-1} \cdot \dfrac{s}{\sqrt{n}}$

Reject $\mathcal{H}_o$ iff $\bar{x} > c$

Method 2: Reject $\mathcal{H}_o$ iff $t_{\text{obs}} > t_{\alpha,n-1}$, where $t_{\text{obs}} = \dfrac{\bar{x} - \mu_o}{\left(\dfrac{s}{\sqrt{n}}\right)}$

Method 3: Reject $\mathcal{H}_o$ iff $\text{P}\left[T > t_{\text{obs}}\right] < \alpha$, where $t_{\text{obs}} = \dfrac{\bar{x} - \mu_o}{\left(\dfrac{s}{\sqrt{n}}\right)}$

Lower tailed test,   $\mathcal{H}_A : \mu < \mu_o$ :

Method 1:     Evaluate       $c = \mu_o - t_{\alpha, n-1} \cdot \dfrac{s}{\sqrt{n}}$

               Reject $\mathcal{H}_o$ iff   $\bar{x} < c$

Method 2:     Reject $\mathcal{H}_o$ iff   $t_{obs} < -t_{\alpha, n-1}$ ,   where   $t_{obs} = \dfrac{\bar{x} - \mu_o}{\left(\dfrac{s}{\sqrt{n}}\right)}$

Method 3:     Reject $\mathcal{H}_o$ iff   $P\left[T < t_{obs}\right] < \alpha$,   where   $t_{obs} = \dfrac{\bar{x} - \mu_o}{\left(\dfrac{s}{\sqrt{n}}\right)}$

Modify the above if the sample size is large and/or $\sigma^2$ is known.
Modify the above for two sample unpaired tests.

---

**Chi-Square Goodness-of-Fit Test:**

$\mathcal{H}_o : p_1 = p_{1o}, \; p_2 = p_{2o}, \ldots, \; p_k = p_{ko}$   vs.   $\mathcal{H}_A :$ not $\mathcal{H}_o$

Observed values $\{O_1, O_2, \ldots, O_k\}$.   Calculate expected values $\{E_1, E_2, \ldots, E_k\}$, then

compare $\chi^2 = \displaystyle\sum_{i=1}^{k} \frac{(O_i - E_i)^2}{E_i}$   to $c = \chi^2_{\alpha, k-1}$.   Reject $\mathcal{H}_o$ iff   $\chi^2 > c$

---

**Chi-Square Test for Independence:**

Observations $\{o_{ij}\}$ in $I$ rows and $J$ columns.   $\mathcal{H}_o :$ factors are independent.

Calculate $e_{ij} = \dfrac{o_{i\bullet} \times o_{\bullet j}}{o_{\bullet\bullet}}$.   Iff $\chi^2 = \displaystyle\sum_{i=1}^{I}\sum_{j=1}^{J} \frac{\left(o_{ij} - e_{ij}\right)^2}{e_{ij}} > \chi^2_{\alpha, (I-1)(J-1)}$ then reject $\mathcal{H}_o$.

---

**Simple Linear Regression:**

Model:       $Y = \beta_1 x + \beta_0 + \varepsilon$,   $\varepsilon \sim N\left(0, \sigma^2\right)$

Check for linear trend with constant error variance and normally distributed residuals.
Find summary statistics:

      $n$      $\sum x$      $\sum y$      $\sum x^2$      $\sum xy$      $\sum y^2$

Evaluate

$$n S_{xx} = n \sum x^2 - \left(\sum x\right)^2 \;,\quad n S_{xy} = n \sum xy - \sum x \sum y \;,\quad n S_{yy} = n \sum y^2 - \left(\sum y\right)^2$$

then

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} \;,\quad \hat{\beta}_0 = \frac{1}{n}\left(\sum y - \hat{\beta}_1 \sum x\right)$$

The line of best fit in the least squares sense is   $y = \hat{\beta}_1 x + \hat{\beta}_0$

Entries in the ANOVA table:

$$SSR = \sum(\hat{y}-\bar{y})^2 = \frac{\left(nS_{xy}\right)^2}{n\left(nS_{xx}\right)} \quad , \quad SST = \sum(y-\bar{y})^2 = S_{yy} = \frac{\left(nS_{yy}\right)}{n}$$

$$SSE = \sum(y-\hat{y})^2 = SST - SSR = \frac{\left(nS_{xx}\right)\left(nS_{yy}\right) - \left(nS_{xy}\right)^2}{n\left(nS_{xx}\right)}$$

$$MSR = \frac{SSR}{\nu_R} = SSR \quad , \qquad s^2 = MSE = \frac{SSE}{\nu_E} = \frac{SSE}{n-2} \quad , \qquad f = t^2 = \frac{MSR}{MSE}$$

| Source | d.f. | Sum Sq. | Mean Sq. | $f$ |
|--------|------|---------|----------|-----|
| R | 1 | SSR | MSR | $f$ |
| E | $n-2$ | SSE | MSE | |
| T | $n-1$ | SST | | |

Variance of slope $\hat{\beta}_1$:  $s_b^{\,2} = \dfrac{MSE}{S_{xx}} = \dfrac{\left(nS_{xx}\right)\left(nS_{yy}\right)-\left(nS_{xy}\right)^2}{(n-2)\left(nS_{xx}\right)^2}$

Coefficient of Determination:

$$r^2 = \frac{SSR}{SST} = \frac{\left(nS_{xy}\right)^2}{\left(nS_{xx}\right)\left(nS_{yy}\right)}$$

Correlation coefficient $= r = \operatorname{sign}\left(\hat{\beta}_1\right)\cdot\sqrt{r^2}$

Hypothesis tests on the linear association of $X$ and $Y$  (all have $\nu = \nu_E = n-2$):

To test $\mathcal{H}_o$: $\rho = 0$ vs. $\mathcal{H}_A$: $\rho \neq 0$  (or, equivalently, $\mathcal{H}_o$: $\beta_1 = 0$ vs. $\mathcal{H}_A$: $\beta_1 \neq 0$):
Use *any* of

$$t = \frac{r\sqrt{n-2}}{\sqrt{1-r^2}}, \quad t = \frac{\hat{\beta}_1 - 0}{s_b}, \quad t = \sqrt{\frac{MSR}{MSE}}, \quad t = \sqrt{\frac{\left(nS_{xy}\right)^2(n-2)}{\left(nS_{xx}\right)\left(nS_{yy}\right)-\left(nS_{xy}\right)^2}}$$

To test $\mathcal{H}_o$: $\beta_1 = \beta_{1o}$ vs. $\mathcal{H}_A$: $\beta_1 > \beta_{1o}$  (when $\beta_{1o} \neq 0$) use

$$t = \frac{\hat{\beta}_1 - \beta_{1o}}{s_b} \qquad \text{or} \qquad t = \left(\left(nS_{xy}\right) - \beta_{1o}\left(nS_{xx}\right)\right)\sqrt{\frac{(n-2)}{\left(nS_{xx}\right)\left(nS_{yy}\right)-\left(nS_{xy}\right)^2}}$$

The $(1-\alpha)\times 100\%$ **confidence interval** estimate for $\mu = \mathrm{E}\left[Y \mid x = x_o\right]$ is

$$\left(\hat{\beta}_0 + \hat{\beta}_1 x_o\right) \pm t_{\alpha/2,\,(n-2)}\, s \sqrt{\frac{1}{n} + \frac{n\left(x_o - \bar{x}\right)^2}{\left(nS_{xx}\right)}}$$

The $(1-\alpha)\times 100\%$ **prediction interval** estimate for $Y \mid x = x_o$ is

$$\left(\hat{\beta}_0 + \hat{\beta}_1 x_o\right) \pm t_{\alpha/2,\,(n-2)}\, s \sqrt{1 + \frac{1}{n} + \frac{n\left(x_o - \bar{x}\right)^2}{\left(nS_{xx}\right)}}$$

[Space for Additional Notes]