

# Advanced Analyses Calculations

## Introduction

This document is a collection of the calculations used in the advanced analyses in STATGRAPHICS *Plus* Quality and Design, and STATGRAPHICS *Plus* Professional.

The calculations are arranged by analysis in the order that they are listed on the Special menu:

Analysis	Page
Quality Control	2
Design of Experiments	27
Time-Series	32
Multivariate Methods	40
Advanced Regression	45

# Quality Control Calculations

## Acceptance Chart

The Acceptance Chart shows the measurements or subgroup averages with modified control limits. If Sigma Multiple is selected:

$$UCL = (USL - Z_{\delta}\sigma) + \frac{3\sigma}{\sqrt{n}}$$

$$LCL = (LSL + Z_{\delta}\sigma) - \frac{3\sigma}{\sqrt{n}}$$

where

$\delta$  is the maximum allowable proportion of non-conforming items.

If Beta Risk is specified, then the control limits are set at:

$$UCL = (USL - Z_{\delta}\sigma) - \frac{Z_{\beta}\sigma}{\sqrt{n}}$$

$$LCL = (LSL + Z_{\delta}\sigma) + \frac{Z_{\beta}\sigma}{\sqrt{n}}$$

where

$\beta$  is the probability of being within the control limits when the proportion of non-conforming items equals  $\delta$ .

## Acceptance Sampling

Acceptance Sampling is the process for determining sample sizes and decision rules for incoming or outgoing products. STATGRAPHICS *Plus* has procedures for generating two types of acceptance sampling plans: attributes and variables. For each of the two basic plan types there are three variations: OC Plans, AOQL Plans and LTPD Plans.

Notation:

$N$  = the number of items in the lot

$n$  = the number of items sampled from the lot and inspected

AQL = Acceptable Quality Level, the poorest level of quality which the consumer finds acceptable on average

LTPD = Lot Tolerance Percent Defective, the poorest level of quality that the consumer is willing to tolerate in any given lot

AOQL = Average Outgoing Quality Limit, the maximum percent of defective items accepted by a given sampling plan

### AOQL Plans

The focus is on the maximum percent of non-conforming items shipped after verification.

$$AOQL = \max_{\theta} \left[ \theta P(\text{accept} | \theta) \left( \frac{N-n}{N} \right) \right]$$

### LTPD Plans

The LTPD Plans minimize the number of inspections while controlling the probability of accepting a "bad" lot. The consumer's risk is calculated according to:

$$\beta = P(\text{accept} | \text{LTPD})$$

## OC Plans

The OC Plan controls the probability of accepting a lot at both the the AQL and the LTPD. All compute acceptance probabilities according to:

$$P(\text{accept} | \theta) = \sum_{x=0}^c P_H(x | \theta, N) = \sum_{x=0}^c \frac{\binom{N\theta}{x} \binom{N(1-\theta)}{n-x}}{\binom{N}{n}}$$

## ARIMA Charts

Underlying the ARIMA control charts is the general AutoRegressive Integrated Moving Average Model:

$$z_j = \theta_0 + \Phi_1 z_{t-1} + \Phi_2 z_{t-2} + \dots + \Phi_p z_{t-p} \\ + a - \theta_1 a_{t-1} - \theta_2 a_{t-2} - \dots - \theta_q a_{t-q}$$

where

$$z_j = \nabla^d \bar{x}_j$$

and

$$a_j \sim \text{NID}(0, \sigma_a)$$

Parameters which define this model:

$\theta_0 = \text{constant}$

$\Phi_1, \Phi_2, \dots, \Phi_p = \text{autoregressive parameters}$

$\theta_1, \theta_2, \dots, \theta_q = \text{moving average parameters}$

## Chart Type

Data with Long-Term Limits - plots the subgroup averages with limits defined by:

$$\hat{\mu} \pm 3\hat{\sigma}_z$$

Data with One-Step Limits - plots the subgroup averages with limits defined by:

$$\hat{z}_j (j-1) \pm 3\sigma_a$$

where

$\hat{z}_j (j - 1)$  is the conditional expectation of the subgroup average at period  $j$  given all information through period  $j-1$ .

Residuals - plots the residuals defined by:

$$\hat{a}_j = z_j - \hat{z}_j (j - 1)$$

Normalized Residuals - plots standardized residuals defined by:

$$\frac{\hat{a}_j}{\hat{\sigma}_a}$$

## Process Capability

### Normal Capability Indices

$$C_p = \frac{USL - LSL}{6\hat{\sigma}}$$

$$C_{pk} = \text{Min} \left\{ \frac{(USL - \bar{X})}{3\hat{\sigma}}, \frac{(\bar{X} - LSL)}{3\hat{\sigma}} \right\}$$

$$C_{pk(\text{upper})} = \frac{USL - \bar{X}}{3\hat{\sigma}}$$

$$C_{pk(\text{lower})} = \frac{\bar{X} - LSL}{3\hat{\sigma}}$$

$$C_r = \frac{1}{C_p}$$

$$C_{pm} = \frac{USL - LSL}{6\hat{\sigma}}$$

where

$$\hat{\sigma} = \sqrt{\sum_{i=1}^n \frac{(\bar{X}_i - \text{Nominal})^2}{n-1}}$$

$$K = \frac{\bar{\bar{X}} - \text{Nominal}}{(\text{USL} - \text{LSL})/2}$$

Note: For variables control charts, STATGRAPHICS uses the appropriate  $\hat{\sigma}$  to calculate the capability indices. For example, with X-bar and R Charts,

$$\hat{\sigma} = \frac{\bar{R}}{d_2}$$

### Confidence Limits for Capability Indices

$$C_{pk} \pm Z_{\alpha/2} \sqrt{\frac{1}{9n} + \frac{C_{pk}^2}{2n-2}}$$

$$\left[ \sqrt{\frac{C_1}{v}} C_{pm}, \sqrt{\frac{C_2}{v}} C_{pm} \right]$$

where

$$v = \frac{(n + \lambda)^2}{n + 2\lambda}$$

and

$$\lambda = n \left( \frac{\bar{X} - \text{Target}}{s} \right)^2$$

and

$$C_1 = \frac{\alpha}{2} \text{ lower tail of the central } \chi_v^2$$

and

$$C_2 = 1 - \frac{\alpha}{2} \text{ of the central } \chi_v^2$$

$$\left[ C_p \sqrt{\frac{\chi_{1-\alpha/2}^2}{n-1}}, C_p \sqrt{\frac{\chi_{\alpha/2}^2}{n-1}} \right]$$

### Non-normal Capability Indices

$$C_p = \frac{USL - LSL}{U_p - L_p}$$

$$C_r = \frac{U_p - L_p}{USL - LSL} = \frac{1}{C_p}$$

$$C_{pk(\text{upper})} = \frac{USL - \text{Median}_{\text{est}}}{U_p - \text{Median}_{\text{est}}}$$

$$C_{pk(\text{lower})} = \frac{\text{Median}_{\text{est}} - LSL}{\text{Median}_{\text{est}} - L_p}$$

$$C_{pk} = \min [C_{pk(\text{upper})}, C_{pk(\text{lower})}]$$

$$K = \frac{\text{Median}_{\text{est}} - \text{Nominal}}{\left( \frac{USL - LSL}{2} \right)}$$

$$U_p = \bar{x} + s U_p^*$$

$$L_p = \bar{x} - s L_p^*$$

$$\text{Median}_{\text{est}} = \bar{x} + s \text{Median}^*$$

The system uses the kurtosis and skewness statistics to locate the  $L_p$ , Median, and  $U_p$  values in the Pearson Curve tables that Gruska, Mirkhani, and Lamberson display in their work.

## Pareto

For information on the calculations STATGRAPHICS uses in Pareto Analysis, see Duncan (1974).

## All Variables Control Charts

In the following formulas, coefficients like  $A_2$ ,  $c_4$ ,  $d_2$ ,  $d_3$ , and so on have their usual meaning in the quality control literature. You can look up their values in tables found in almost any quality control textbook.

$N$  = number of subgroups

$n_j$  = number of observations in subgroup  $j$

$j$  = 1, 2, ...,  $N$

$x_{ij}$  =  $i$ th observation in subgroup  $j$

### Subgroup

Means:

$$\bar{x}_j = \frac{\sum_{i=1}^{n_j} x_{ij}}{n_j}$$

Standard Deviations:

$$s_j = \sqrt{\frac{\sum_{i=1}^{n_j} (x_{ij} - \bar{x}_j)^2}{(n_j - 1)}}$$

Ranges:

$$R_j = \max \{ x_{ij} \mid 1 \leq i \leq n_j \} - \min \{ x_{ij} \mid 1 \leq i \leq n_j \}$$



## Grand Mean

Mean Range:  $\bar{R} = \frac{1}{N} \sum_{j=1}^N R_j$

Let

CL = centerline

UCL = upper control limit

LCL = lower control limit

N = number of subgroups

$n_j$  = size for the jth subgroup

$\bar{\bar{X}}$  = grand average of all subgroups

$\bar{R}$  = average (mean) range

$\mu$  = standard process mean

$\sigma$  = standard process sigma

$\hat{\sigma}$  = estimated process sigma

$k_u$  = multiple of sigma for upper control limit

$k_l$  = multiple of sigma for lower control limit

$d_2$  = expected value of  $R/\sigma$  (depends on n)

$d_3$  = standard deviation of  $R/\sigma$  (depends on n)

$\bar{s}$  = average of subgroup standard deviations

$c_4$  = expected value of  $s/\sigma$  (depends on n)

pooled s =  $\sqrt{\text{pooled variance}}$

$$\text{pooled variance} = \frac{\sum_{j=1}^N (n_j - 1) s_j^2}{\sum_{j=1}^N (n_j - 1)}$$

$\alpha_u$  = alpha for upper control limit  
(probability of a Type I error beyond limit)

$\alpha_l$  = alpha for lower control limit  
(probability of a Type I error beyond limit)

$\chi^2_{1-\alpha, v}$  =  $1 - \alpha$  probability point of chi-square distribution  
with v degrees of freedom

## X-bar and R Charts

**Estimated:**

process mean =  $\bar{\bar{X}}$

process sigma =  $\frac{\bar{R}}{d_2} = \hat{\sigma}$

mean range =  $\bar{R}$

**Standard:**

process mean =  $\mu$

process sigma =  $\sigma$

mean range =  $d_2 \sigma$

## Initial Study

X-bar Chart:

CL =  $\bar{\bar{X}}$

UCL =  $\bar{\bar{X}} + k_u \frac{\hat{\sigma}}{\sqrt{n}}$

$$LCL = \bar{\bar{X}} - k_1 \frac{\hat{\sigma}}{\sqrt{n}}$$

Range Chart:

$$CL = \bar{R}$$

$$UCL = \bar{R} + k_u d_3 \frac{\bar{R}}{d_2}$$

$$LCL = \text{Max} \left\{ 0, \bar{R} - k_l d_3 \frac{\bar{R}}{d_2} \right\}$$

### Control-to-Standard

X-bar Chart:

$$CL = \mu$$

$$UCL = \mu + k_u \frac{\sigma}{\sqrt{n}}$$

$$LCL = \mu - k_l \frac{\sigma}{\sqrt{n}}$$

Range Chart:

$$CL = d_2 \sigma$$

$$UCL = (d_2 \sigma) + k_u d_3 \sigma$$

$$LCL = \text{Max} \{0, (d_2 \sigma) - k_l d_3 \sigma\}$$

### Normalization: X-bar (i = 1,...,N)

Initial Study:

$$\frac{(\bar{x}_i - \bar{\bar{X}})}{\hat{\sigma}}$$

Control-to-Standard:

$$\frac{(\bar{x}_i - \mu)}{\sigma}$$

## X-bar and S Charts

<b>Estimated:</b>	process mean = $\bar{\bar{X}}$
	process sigma = $\frac{\bar{s}}{c_4} = \hat{\sigma}$
	mean sigma = $\bar{s}$
<b>Standard:</b>	process mean = $\mu$
	process sigma = $\sigma$
	mean sigma = $c_4 \sigma$

### Initial Study

X-bar Chart:	CL = $\bar{\bar{X}}$
	UCL = $\bar{\bar{X}} + k_u \frac{\hat{\sigma}}{\sqrt{n}}$
	LCL = $\bar{\bar{X}} - k_l \frac{\hat{\sigma}}{\sqrt{n}}$
S Chart:	CL = $\bar{s}$
	UCL = $\bar{s} + k_u \bar{s} \frac{\sqrt{1 - (c_4)^2}}{c_4}$
	LCL = $\text{Max} \left\{ 0, \bar{s} - k_l \bar{s} \frac{\sqrt{1 - (c_4)^2}}{c_4} \right\}$

### Control-to-Standard

X-bar Chart:	CL = $\mu$
--------------	------------

$$UCL = \mu + k_u \frac{\sigma}{\sqrt{n}}$$

$$LCL = \mu - k_l \frac{\sigma}{\sqrt{n}}$$

## X-bar and S-squared Charts

### Estimated:

$$\text{process mean} = \bar{\bar{X}}$$

$$\text{process sigma} = \sqrt{\text{pooled } s^2} = \hat{\sigma}$$

### Standard:

$$\text{process mean} = \mu$$

$$\text{process sigma} = \sigma$$

$$\text{variance} = \sigma^2$$

### Initial Study:

X-bar Chart

$$CL = \bar{\bar{X}}$$

$$UCL = \bar{\bar{X}} + k_u \frac{\hat{\sigma}}{\sqrt{n}}$$

$$LCL = \bar{\bar{X}} - k_l \frac{\hat{\sigma}}{\sqrt{n}}$$

S-squared Chart:

$$CL = \text{pooled } s^2$$

$$UCL = \frac{\text{pooled } s^2}{n-1} \chi^2_{\alpha_u, n-1}$$

$$LCL = \frac{\text{pooled } s^2}{n-1} \chi^2_{1-\alpha, n-1}$$

### Control-to-Standard:

X-bar Chart:

$$CL = \mu$$

$$UCL = \mu + k_u \frac{\sigma}{\sqrt{n}}$$

$$LCL = \mu - k_l \frac{\sigma}{\sqrt{n}}$$

S-squared Chart:

$$CL = \sigma^2$$

$$UCL = \frac{\sigma^2}{n-1} \chi^2_{\alpha, n-1}$$

$$LCL = \frac{\sigma^2}{n-1} \chi^2_{\alpha, n-1}$$

## Individuals Charts

**Ranges are estimated by:**

$$MR(2)_j = |X_j - X_{j-1}|, \quad j = 2, \dots, N$$

$$MR(2)_1 = \text{missing}$$

STATGRAPHICS calculates  $d_2$  and  $d_3$  on a subgroup size of 2.

Estimated:

$$\text{process mean} = \bar{\bar{X}}$$

$$\text{process sigma} = \frac{\overline{MR(2)}}{d_2} = \hat{\sigma}$$

$$\text{mean } MR(2) = \overline{MR(2)}$$

Standard:                      process mean =  $\mu$   
                                      process sigma =  $\sigma$   
                                      mean MR(2) =  $d_2 \sigma$

### Initial Study

X Chart:                      CL =  $\bar{\bar{X}}$   
                                      UCL =  $\bar{\bar{X}} + k_u \hat{\sigma}$   
                                      LCL =  $\bar{\bar{X}} - k_l \hat{\sigma}$

MR(2) Chart:                      CL =  $\overline{MR(2)}$   
                                      UCL =  $\overline{MR(2)} + k_u d_3 \frac{\overline{MR(2)}}{d_2}$   
                                      LCL =  $\text{Max} \left\{ 0, \overline{MR(2)} - k_l d_3 \frac{\overline{MR(2)}}{d_2} \right\}$

### Control-to-Standard

X Chart:                      CL =  $\mu$   
                                      UCL =  $\mu + k_u \sigma$   
                                      LCL =  $\mu - k_l \sigma$

Range Chart:                      CL =  $d_2 \sigma$   
                                      UCL =  $(d_2 \sigma) + k_u d_3 \sigma$   
                                      LCL =  $\text{Max} \{0, (d_2 \sigma) - k_l d_3 \sigma\}$

**Normalization: X**

Initial Study: 
$$\frac{(x_j - \bar{x})}{\hat{\sigma}}$$

Control to Standard: 
$$\frac{(x_j - \mu)}{\sigma}$$

**Normalization: Range**

Initial Study: 
$$(MR(2)_j - \overline{MR(2)}) \frac{d_2}{(MR(2) d_3)}$$

Control to Standard: 
$$\frac{(MR(2)_j - d_2 \sigma)}{(\sigma d_3)}$$

**Moving Average Chart**

Moving averages of order  $q$ ,  $q = 1, 2, \dots, n$

Display Limits: 
$$CL = \mu$$

$$UCL = \mu + \frac{k_u \sigma}{\sqrt{nq}}$$

$$LCL = \mu - \frac{k_l \sigma}{\sqrt{nq}}$$

Chart Limits: 
$$CL = \mu$$

$$UCL[t] = \mu + \frac{k_u \sigma}{\sqrt{nq[t]}}$$

$$LCL[t] = \mu - \frac{k_l \sigma}{\sqrt{nq[t]}}$$



where

$$q[t] = \begin{cases} n & \text{for } t < 8 \\ t & \text{for } t \geq q \text{ for } t = 1, \dots, n \end{cases}$$

Chart Data:

$$Z[t] = \begin{cases} \frac{(x[t] + x[t-1] + \dots + x[t-q+1])}{q} & \text{for } t \geq q \\ \frac{(x[t] + x[t-1] + \dots + x[1])}{t} & \text{for } t < q \end{cases}$$

## Exponentially Weighted Moving Average (EWMA) Chart

Display Limits:

$$CL = \mu$$

$$UCL = \mu + k_u \sigma \sqrt{\frac{\lambda}{(2-\lambda)n}}$$

$$LCL = \mu - k_l \sigma \sqrt{\frac{\lambda}{(2-\lambda)n}}$$

Chart Limits:

$$CL = \mu$$

$$UCL = \mu + k_u \sigma \sqrt{\left(\frac{\lambda}{2-\lambda}\right) (1 - (1-\lambda)^{2t}) \frac{1}{n}}$$

$$LCL = \mu - k_l \sigma \sqrt{\left(\frac{\lambda}{2-\lambda}\right) (1 - (1-\lambda)^{2t}) \frac{1}{n}}$$

Chart Data:

$$z[0] = \mu$$

$$z[t] = (\lambda x[t]) + (1-\lambda) z[t-1] \quad \text{for } t=1,2,\dots,N$$

Moving Range: Moving Range of order q for  $q=1,2,\dots,N$

Control Limits:

$$CL = d_2 \sigma$$

$$UCL = (d_2 \sigma) + k_u d_3 \sigma$$

$$LCL = (d_2 \sigma) - k_l d_3 \sigma$$

where  $d_2$  and  $d_3$  are based on a subgroup size of  $q$

Chart Data:

$$R[t] = \begin{aligned} &\text{Max} (x[t] + x[t - 1] + \dots + x[t - q + 1]) \\ &\quad - \text{Min} (x[t] + x[t - 1] \\ &\quad + \dots + x[t - q + 1]) \quad \text{for } t \geq q \\ &= \text{missing value for } t < q \end{aligned}$$

## Cumulative Sum (CuSum) Chart

$\mu$  = control mean  
 $\sigma$  = standard deviation  
 $\Delta$  = difference to detect  
 $\alpha$  = Type I error  
 $\beta$  = Type II error  
 $n$  = number of data points (subgroups)  
 $\bar{x}$  = mean of subgroup  $i$

Plot cumulative sums  $C_t$  versus  $t$  where:

$$C_t = \sum_{i=1}^t (\bar{x}_i - \mu) \quad \text{for } t = 1, 2, \dots, n$$

The V-mask is located at distance

$$d = \frac{2}{\Delta} \left[ \frac{\frac{\sigma^2}{n}}{\Delta} \ln \frac{1-\beta}{\frac{\alpha}{2}} \right]$$

in front of the last data point.

$$\text{Angle of mask} = 2 \tan^{-1} \left( \frac{\Delta}{2} \right)$$

$$\text{Slope of the line} = \pm \frac{\Delta}{2}$$

Let

$\mu$  = process mean

$\sigma$  = process sigma

$x[i]$  =  $i$ th observation (if individual data are input)  
 $i$ th  $\bar{x}$  (if subgroup statistics)

$R[i]$  =  $i$ th range (if subgroup statistics)

$n$  = 1 (if individual data are input)  
 average subgroup size (if subgroup statistics)

$N$  = number of subgroups

$k_u$  = sigma multiple for upper control

$k_l$  = sigma multiple for lower control

## Multivariate Control Chart

$S$  = sample covariance matrix

$\tilde{X}_t$  = observation vector at time  $t$

$\tilde{\bar{X}}$  = vector of column averages

Then Hotelling's T-squared is

$$T_t^2 = (\underline{X}_t - \bar{\underline{X}})' S^{-1} (\underline{X}_t - \bar{\underline{X}})$$

and the upper control limit is

$$UCL = \left( \frac{k(n-1)}{(n-k)} \right) F_{k, n-k, \alpha}$$

Points are plotted at  $T_t^2$ ; glyphs are centered at  $T_t^2$ .

## All Attributes Control Charts

$N$  = number of subgroups

$n_j$  = number of observations in subgroup  $j$

$j = 1, 2, \dots, N$

$$\bar{n} = \frac{1}{N} \sum_{j=1}^N n_j$$

$u_j$  = number of defects in sample  $j$

$$\bar{u} = \frac{\text{number of defects in all samples}}{\text{number of units in all samples}} = \frac{\sum u_j}{\sum n_j}$$

$u'$  = standard  $u$

$$p_j = \frac{\text{number of defective units in subgroup } j}{\text{number of units inspected in subgroup } j}$$

$$\bar{p} = \frac{\text{number of defectives in all samples}}{\text{number of units in all samples}} = \frac{\sum p_j n_j}{\sum n_j}$$

$p'$  = standard  $p$

$$\bar{f} = \bar{p} N$$

$$f' = p' N$$

$c_j$  = number of defects in item j

$$\bar{c} = \frac{\text{number of defects in all items}}{\text{number of items}} = \frac{\sum c_j}{n}$$

$c'$  = standard c = standard number of defects

CL = centerline

UCL = upper control limit

LCL = lower control limit

$k_u$  = multiplier for upper control limit

$k_l$  = multiplier for lower control limit

$n[j]$  =  $\bar{n}$  (when "Avg. subgroup size" is selected)  
 $n_j$  (when "Avg. subgroup size" is not selected)

## p Chart

Estimated:

$$\text{mean } p = \bar{p}$$

$$\text{sigma} = \sqrt{\frac{\bar{p}(1-\bar{p})}{\bar{n}}}$$

Standard:

$$\text{mean } p = p'$$

$$\text{sigma} = \sqrt{\frac{p'(1-p')}{\bar{n}}}$$

Initial:

$$\text{CL} = \bar{p}$$

$$\text{UCL} = \bar{p} + k_u \sqrt{\frac{\bar{p}(1-\bar{p})}{n[j]}}$$

$$LCL = \bar{p} - k_l \sqrt{\frac{\bar{p}(1-\bar{p})}{n[j]}}$$

Control-to-Standard:

$$CL = p'$$

$$UCL = p' + k_u \sqrt{\frac{p'(1-p')}{n[j]}}$$

$$LCL = p' - k_l \sqrt{\frac{p'(1-p')}{n[j]}}$$

### Normalization

Initial:

$$\frac{(p_j - \bar{p})}{\sqrt{\frac{\bar{p}(1-\bar{p})}{n[j]}}}$$

Control-to-Standard

$$\frac{(p_j - p')}{\sqrt{\frac{p'(1-p')}{n[j]}}}$$

### np Chart

Estimated:

$$\text{mean np} = \bar{f}$$

$$\text{sigma} = \sqrt{\bar{f}(1-\bar{f})}$$

Standard:

$$\text{mean np} = f'$$

$$\text{sigma} = \sqrt{f'(1-f')}$$

Initial:

$$CL = \bar{f}$$

$$UCL = \bar{f} + k_u \sqrt{n[j] \bar{p} (1 - \bar{p})}$$

$$LCL = \bar{f} - k_l \sqrt{n[j] \bar{p} (1 - \bar{p})}$$

Control-to-Standard:

$$CL = f'$$

$$UCL = f' + k_u \sqrt{n[j] p' (1 - p')}$$

$$LCL = f' - k_l \sqrt{n[j] p' (1 - p')}$$

### Normalization

Initial:

$$\frac{(f_j - \bar{f})}{\sqrt{n[j] \bar{p} (1 - \bar{p})}}$$

Control-to-Standard:

$$\frac{(f_j - f')}{\sqrt{n[j] p' (1 - p')}}}$$

### c Chart

Estimated:

$$\text{mean } c = \bar{c}$$

$$\text{sigma} = \sqrt{\bar{c}}$$

Standard:                      mean  $c = c'$

                                     sigma    =    $\sqrt{c'}$

Initial:                              CL =  $\bar{c}$

   UCL =  $\bar{c} + k_u \sqrt{\bar{c}}$

   LCL =  $\bar{c} - k_l \sqrt{\bar{c}}$

Control-to-Standard:            CL =  $c'$

   UCL =  $c' + k_u \sqrt{c'}$

   LCL =  $c' - k_l \sqrt{c'}$

### **Normalization**

Initial:                               $\frac{(c_j - \bar{c})}{\sqrt{\bar{c}}}$

Control-to-Standard:             $\frac{(c_j - c')}{\sqrt{c'}}$

### **u Chart**

Estimated:                        mean  $u = \bar{u}$

   sigma    =    $\sqrt{\frac{\bar{u}}{n}}$



Standard:  $\text{mean } u = u'$

$$\text{sigma} = \sqrt{\frac{u'}{n}}$$

Initial:  $CL = \bar{u}$

$$UCL = \bar{u} + k_u \sqrt{\frac{\bar{u}}{n[j]}}$$

$$LCL = \bar{u} - k_l \sqrt{\frac{\bar{u}}{n[j]}}$$

Control-to-Standard:  $CL = u'$

$$UCL = u' + k_u \sqrt{\frac{u'}{n[j]}}$$

$$LCL = u' - k_l \sqrt{\frac{u'}{n[j]}}$$

### Normalization

Initial:  $\frac{(u_j - \bar{u})}{\sqrt{\frac{\bar{u}}{n[j]}}}$

Control-to-Standard:  $\frac{(u_j - u')}{\sqrt{\frac{u'}{n[j]}}}$

## Gage R&R

For information on the calculations STATGRAPHICS uses in the Gage R&R Analysis, see AIAG (1990), *Measurement Systems Analysis Reference Manual*.

## Toolwear Chart

The toolwear chart was originally developed to monitor the wear of tools where the measurements were expected to follow a natural trend. That linear trend is defined by:

$$\mu_j = A + B_j$$

where

$\mu_j$  = the mean wear at period  $j$

$A$  = the y-axis intercept

$B$  = the slope

Control limits are placed by default at:

$$\hat{A} + \hat{B}_j \pm 3 \frac{\hat{\sigma}}{\sqrt{n}}$$

For initial studies the model estimation error can be included:

$$\hat{A} + \hat{B}_j \pm 3 \hat{\sigma} \sqrt{\frac{1}{\bar{n}} \left( 1 + \frac{(j - \bar{j})^2}{\sum_{i=1}^m n_i (i - \bar{j})^2} \right)}$$

where

$\bar{n}$  = average subgroup size ( $\bar{n} = 1$  for individuals data)

$\bar{j}$  = average of the period numbers used to fit the line

# Design of Experiments Calculations

## D-Optimal Design

The optimize design procedure seeks to find the design which maximizes the determinant of the M matrix defined by:

$$|M| = \frac{|X'X|}{N^p}$$

where

N = total number of runs

p = the number of estimated coefficients

X = the normalized design matrix.

Optimize design also shows various measures of efficiency:

D-efficiency - compare the determinant of M to the value of the best possible design according to:

$$\text{D-efficiency} = 100 \left( \frac{1}{N} |X'X|^{\frac{1}{p}} \right) \%$$

A-efficiency - compares the sum of the variances of the estimated regression coefficients without considering their covariances:

$$\text{A-efficiency} = \left( \frac{p}{\text{trace}(N(X'X)^{-1})} \right) \%$$

G-efficiency - compares the maximum prediction standard error over the design points  $\sigma_m$  through:

$$\text{G-efficiency} = 100 \left( \frac{\sqrt{p/N}}{\sigma_m} \right) \%$$

## Inner/Outer Arrays

For information on the STATGRAPHICS *Plus* implementation of arrays, see Taguchi (1987), *The System of Experimental Designs: Engineering Methods to Optimize Quality and Minimize Costs*.

## Mixture Designs

The Scheffe polynomial formulas for the Linear, Quadratic, Special Cubic, or Cubic models have the following forms.

### **Linear Model:**

$$Y = b_1X_1 + b_2X_2 + \dots + b_NX_N$$

### **Quadratic Model:**

$$Y = \text{Linear Model} + b_{12}X_1X_2 + b_{13}X_1X_3 + \dots + b_{N-1, N}X_{N-1}X_N$$

### **Special Cubic Model:**

$$Y = \text{Quadratic Model} + b_{123}X_1X_2X_3 + \dots + b_{N-2, N-1, N}X_{N-2}X_{N-1}X_N$$

### **Cubic Model:**

$$\begin{aligned} Y = & \text{Special Cubic Model} + d_{12}X_1X_2(X_1 - X_2) + d_{13}X_1X_3(X_1 - X_3) \\ & + \dots \\ & + d_{N-1, N}X_{N-1}X_N(X_{N-1} - X_N) \end{aligned}$$

## Multi-Factor Categorical Designs

Multi-Factor Categorical Designs are analyzed using the Multifactor ANOVA procedure. Additional information on the calculations can be found in Neter, et al (1996), *Applied Linear Statistical Models*, fourth edition.

## Multilevel Factorial Designs

This design fits the general second order model such as the following for three factors from Polhemus (1999):

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_{12} X_1 X_2 + \beta_{13} X_1 X_3 + \beta_{23} X_2 X_3 + \beta_{11} X_1^2 + \beta_{22} X_2^2 + \beta_{33} X_3^2 + \epsilon$$

## Response Surface Designs

### ***Axial Distance***

A design is orthogonally blocked if the axial distance is

$$\alpha = \left[ k \times \frac{(1 + n_{so}/n_s)}{(1 + n_{co}/n_c)} \right]^{1/2}$$

where

$k$  = the number of factors

$n_{so}$  = the number of centerpoints in the star portion of the design

$n_{co}$  = the number of centerpoints in the cube portion of the design

$n_s$  = the number of base points in the star portion

$n_c$  = the number of base points in the cube portion

In a Rotatable design, the variance of the predicted response is constant on a hyperspherical surface centered at the design origin. The axial distance is

$$\alpha = (n_c)^{1/4}$$

where  $n_c$  is the number of points in the cube portion of the design, excluding centerpoints.

An orthogonal design is one in which the estimates of all terms in the second-order model are uncorrelated. To produce an orthogonal design, the axial distance is set to

$$\alpha = \{[(n_c + n_s + n_o)^{1/2} - n_c^{1/2}]^2 \times \frac{n_c}{4}\}^{1/4}$$

where

$n_c$  = the number of cube points

$n_s$  = the number of star points

$n_o$  = the number of centerpoints

You can make a design both rotatable and approximately orthogonal if you set the axial distance to the value used for rotatable designs, and then set the number of centerpoints to the integer closest to

$$n_o = 4 \times n_c^{1/2} + 4 - 2k$$

where

$n_c$  = the number of cube points

$k$  = the number of factors

## Screening Designs

See Chapters 10 and 12 in Box, Hunter, and Hunter (1978) for the calculations used to estimate main effects, interactions, variances, and normal plots.

The design generators and block generators for the factorial designs are listed in Box, Hunter, and Hunter (1978) in Tables 10.B.1 and 12.15.

The error degrees of freedom shown on the Design Selection List screen is computed using the formula

$$\text{d.f.} = n - (k + f + (b - 1)) - 1$$

where

$n$  = the number of runs

$k$  = the number of factors

$f$  = the number of two-factor interactions

$b$  = the number of blocks.

If this number is negative, the system displays a 0. The equation is exactly correct if main effects, two-factor interactions, and block effects are clear of each other. If confounding exists, the degrees of

freedom are increased by adding back the number of nonestimable effects. For Plackett-Burman designs, the system considers only main factor effects to be estimable when computing the degrees of freedom for error.

## Single Factor Categorical Designs

The Single Factor Categorical Designs are analyzed using the One-Way ANOVA analysis. The calculations are in Appendix E of the *Standard Edition Manual*.

## Variance Components Designs

The general class of hierarchical experiments in which factors are nested one within the other. Frequently these experiments are performed to determine what level of a process needs further investigation. The model is of the form

$$Y = \mu + \varepsilon_1 + \varepsilon_2 + \dots + \varepsilon_k$$

where

$Y$  = response variable

$\mu$  = process mean

$\varepsilon_j$  = experimental error due to component  $j$

$k$  = number of variance components

The process standard deviation is generally:

$$\sigma_y = \sqrt{\sigma_1^2 + \sigma_2^2 + \dots + \sigma_k^2}$$

where

$\sigma_y$  = process standard deviation

$\sigma_j$  = standard deviation of error component  $j$

# Time-Series Calculations

## ARIMA Model

The basic form of the model to be fitted is:

$$W_t = \mu + \frac{\theta(B) \Theta_s(B)}{\varphi(B) \Phi_s(B)} a_t$$

This model expresses the data as a combination of the series' past values and the past values of the random input, where:

$t$  = time

$B$  = the backshift operation; that is,  
 $B^l W(t) = W(t-l)$

$W_t$  = the original data or a difference of that data

$\mu$  = the mean

$\theta(B)$  = the nonseasonal moving-average operator,  
 $1 - \theta_1 B - \theta_2 B^2 - \dots - \theta_q B^q$

$\varphi(B)$  = the nonseasonal autoregressive operator,  
 $1 - \varphi_1 B - \varphi_2 B^2 - \dots - \varphi_p B^p$

$\Theta_s(B)$  = the seasonal moving-average operator,  
 $1 - \Theta_1 B^s - \Theta_2 B^{2s} - \dots - \Theta_Q B^{Qs}$

$\Phi_s(B)$  = the seasonal autoregressive operator,  
 $1 - \Phi_1 B^s - \Phi_2 B^{2s} - \dots - \Phi_P B^{Ps}$

$a_t$  = the random error.



This is typically denoted as a  $(p,d,q) = (P,D,Q)^s$  model, where:

$p$  = the order of the nonseasonal autoregressive term  
 $d$  = the order of nonseasonal differencing  
 $q$  = the order of the nonseasonal moving-average term  
 $P$  = the order of the seasonal autoregressive term  
 $D$  = the order of seasonal differencing  
 $Q$  = the order of the seasonal moving-average term  
 $s$  = the length of seasonality

## Autocorrelations

The autocorrelation  $r_k$  at lag  $k$  is calculated as follows:

$$r_k = \frac{c_k}{c_0}$$

where

$$c_k = \frac{1}{n} \sum_{t=1}^{n-k} (y_t - \bar{y})(y_{t+k} - \bar{y})$$

and

$$\bar{y} = \frac{\left( \sum_{t=1}^n y_t \right)}{n}$$

and

$$c_0 = \sum_{t=1}^n (y_t - \bar{y})^2$$

$$\text{standard error at lag } k = \sqrt{\frac{1}{n} \left\{ 1 + 2 \sum_{v=1}^{k-1} r_v^2 \right\}}$$

where

$y_t$  = observation at time  $t$   
 $n$  = number of observations

## Box-Cox Transformations

For the original value  $Z$ , the transformation  $Z_T$  is given by:

$$Z_T = \frac{(Z + \lambda_2)^{\lambda_1 - 1}}{\lambda_1 g^{(\lambda_1 - 1)}} \quad \text{if } \lambda_1 \neq 0$$

$$Z_T = g \ln (Z + \lambda_2) \quad \text{if } \lambda_1 = 0$$

where  $g$  is the sample geometric mean of  $Z + \lambda_2$ . The first parameter  $\lambda_1$  governs the strength of the transformation.

$\lambda_1 = 1$  corresponds to the original data;  $\lambda_1 = 0$  to a logarithm. The system adds  $\lambda_2$  to the data before it applies  $\lambda_1$ .

## Brown's Linear and Quadratic Exponential Smoothing

Let

$\alpha$  = the smoothing constant

$Y_t$  = the observed value at time  $t$

$S_t$  = the smoothed value at time  $t$

Then

$$S_t = \alpha Y_t + (1 - \alpha) S_{t-1}$$

## Crosscorrelations

Crosscorrelation at lag k:  $x$  = input time series

$y$  = output time series

$$r_{xy}(k) = \frac{c_{xy}(k)}{s_x s_y} \quad k = 0, \pm 1, \pm 2, \dots$$

where

$$c_{xy}(k) = \begin{cases} \frac{1}{n} \sum_{t=1}^{n-k} (x_t - \bar{x})(y_{t+k} - \bar{y}) & k = 0, 1, 2, \dots \\ \frac{1}{n} \sum_{t=1}^{n+k} (y_t - \bar{y})(x_{t-k} - \bar{x}) & k = 0, -1, -2, \dots \end{cases}$$

and

$$s_x = \sqrt{c_{xx}(0)}$$

$$s_y = \sqrt{c_{yy}(0)}$$

## Holt's Linear Exponential Smoothing

$N$  = the number of observations

$X_t$  = the  $t$ th observation in the time series

$S_t$  = the smoothed value of the time series at time  $t$

$b_t$  = the estimated trend at time  $t$

$\alpha$  = the smoothing constant for the level of the time series

$\beta$  = the smoothing constant for the trend of the time series

$F_{t+m}$  = the forecast for the time period  $t+m$

The initial values for the estimates are:

$$S_1 = X_1$$

$$b_1 = X_2 - X_1$$

The forecasted values of the time series are given by:

$$F_{t+m} \text{ for } t+m \geq 3$$

$$S_t = \alpha X_t + (1 - \alpha)(S_{t-1} + b_{t-1})$$

$$b_t = \beta(S_t - S_{t-1}) + (1 - \beta)b_{t-1}$$

$$F_{t+m} = S_t + b_t m$$

## Partial Autocorrelation at lag k

The partial autocorrelation at lag k,  $\hat{\theta}_{kk}$ , is obtained by solving the Yule-Walker equations:

$$3r_j =$$

$$\hat{\theta}_{k1} r_{j-1} + \hat{\theta}_{k2} r_{j-2} + \dots + \hat{\theta}_{k(k-1)} r_{j-k+1} + \hat{\theta}_{kk} r_{j-k}$$

$$j = 1, 2, \dots, k$$

$$\text{standard error} = \sqrt{\frac{1}{n}}$$

## Periodogram and Integrated Periodogram

Periodograms are computed using Fourier transforms. The value of the ordinate at each frequency  $f_i$  is given by the following formulas.

If n is odd, where n is the number of observations

$$I(f_i) = \frac{n}{2} (a_i^2 + b_i^2) \quad i=1, 2, \dots, \left[ \frac{n-1}{2} \right]$$

where

$$a_i = \frac{2}{n} \sum_{t=1}^n y_t \cos(2\pi f_i t)$$

$$b_i = \frac{2}{n} \sum_{t=1}^n y_t \sin(2\pi f_i t)$$

$$f_i = \frac{i}{n}$$

If  $n$  is even, an additional term is added:

$$I(0.5) = n \left( \frac{1}{n} \sum_{t=1}^n (1)^t Y_t \right)^2$$

## Resistant Nonlinear Smoothing

For the calculations used in resistant nonlinear smoothing, see Tukey (1977).

## Seasonal Decomposition

For information on the calculations used in seasonal decomposition, see Makridakis, Wheelwright, and McGee (1983).

## Smoothing

For information on the calculations used in resistant nonlinear smoothing, see Tukey (1977).

For information on the calculations used in weighted moving averages, see Kendall and Stuart (1961).

## Tests for Randomness

Tests are variants of the basic runs test. The calculation formulas are the same, but the value of the parameters the formulas use are different.

For the test of the number of runs above and below the median:

$n_1$  = number of observations above the median

$n_2$  = number of observations below the median

For the test of the number of runs up and down:

$n_1$  = number of times the value rises

$n_2$  = number of times the value decreases

The system ignores adjacent pairs (two consecutive numbers that are equal).

For both tests, the expected number of runs is

$$E(R) = \frac{2n_1n_2}{n_1 + n_2} + 1$$

The estimated variance for the runs is

$$V(R) = \frac{2n_1n_2 (2n_1n_2 - n_1 - n_2)}{(n_1 + n_2)^2 (n_1 + n_2 - 1)}$$

The test statistic for the runs test is

$$Z = \frac{|R - E(R)| - 0.5}{\sqrt{V(R)}}$$

## Vertical Time Sequence Plot

For information on the calculations the system uses to produce a vertical time sequence plot, see Box and Jenkins (1976).

## Winter's Exponential Smoothing

$x_t$  = actual demand at time period  $t$

$L$  = length of seasonality

$I_t$  = seasonal factor for period  $t$

$\alpha$  = smoothing constant for level

$\beta$  = smoothing constant for trend

$\gamma$  = smoothing constant for seasonality

$k$  = number of periods into the future ( $k=1,2, \dots$ )

$$\text{level}_t = \alpha \frac{X_t}{I_{t-L}} + (1 - \alpha) (\text{level}_{t-1} + \text{trend}_{t-1})$$

$$\text{trend}_t = \beta (\text{level}_t - \text{level}_{t-1}) + (1 - \beta) \text{trend}_{t-1}$$

$$I_t = \gamma \left( \frac{X_t}{\text{level}_t} \right) + (1 - \gamma) I_{t-L}$$

$$\text{Forecast}_{t+k} = (\text{level}_t + k \times \text{trend}_k) I_{t-L+k}$$

# Multivariate Methods Calculations

## Cluster Analysis

### ***Lance and William Flexible Method***

The formula for the distance measures between groups is:

$$d_{k(ij)} = \alpha_i d_{ki} + \alpha_j d_{kj} + \beta d_{ij} + \gamma |d_{ki} - d_{kj}|$$

where

$d_{ij}$  = the distance between groups  $i$  and  $j$

$\alpha$  = parameter1

$\beta$  = parameter2

$\gamma$  = parameter3

### ***Hierarchical Clustering Scheme***

#### **Quadruple Constraint**

$$\alpha_i + \alpha_j + \beta = 1$$

$$\alpha_i = \alpha_j$$

$$\beta < 1$$

$$\gamma = 0$$



### **Nearest Neighbor**

$$\alpha_i = \alpha_j = \frac{1}{2}$$

$$\beta = 0$$

$$\gamma = -\frac{1}{2}$$

### **Furthest Neighbor**

$$\alpha_i = \alpha_j = \frac{1}{2}$$

$$\beta = 0$$

$$\gamma = \frac{1}{2}$$

### **Centroid**

$$\alpha_i = \frac{n_i}{n_i + n_j}$$

$$\alpha_j = \frac{n_j}{n_i + n_j}$$

$$\beta = -\alpha_i \alpha_j$$

$$\gamma = 0$$

### **Median**

$$\alpha_i = \alpha_j = \frac{1}{2}$$

$$\beta = -\frac{1}{4}$$

$$\gamma = 0$$

### **Group Average**

$$\alpha_i = \frac{n_i}{n_i + n_J}$$

$$\alpha_j = \frac{n_j}{n_i + n_J}$$

$$\beta = \gamma = 0$$

### **Ward's Method**

$$\alpha_i = \frac{n_k + n_i}{n_k + n_i + n_J}$$

$$\alpha_j = \frac{n_k + n_j}{n_k + n_i + n_J}$$

$$\beta = \frac{-n_k}{n_k + n_i + n_J}$$

$$\gamma = 0$$

## **Factor Analysis**

### ***Factor Scores for Unstandardized Variables***

$$F_{jk} = \sum_{i=1}^p W_{ji} X_i = W_{j1} X_1 + W_{j2} X_2 + \dots + W_{jp} X_p$$

where

$p$  = number of variables

$X_i$  = unstandardized variables

$W_j$  = factor weights

### ***Factor Scores for Standardized Variables***

$$F_{jk} = \sum_{i=1}^p W_{ji} Z_{ik} = W_{j1} Z_{1k} + W_{j2} Z_{2k} + \dots + W_{jp} Z_{pk}$$

where

$p$  = number of variables

$W_j$  = factor weights

$Z_i$  = standardized variables

$k$  = observations

$$\text{Percent of Variance} = \frac{\text{Eigenvalue}}{p} \times 100$$

## **Discriminant Analysis**

### ***Discriminant Function***

$$D_i = d_{i1} Z_1 + d_{i2} Z_2 + \dots + d_{ip} Z_p$$

where

$p$  = number of variables

$Z_1, Z_2, \dots, Z_p$  = standardized variables

$d_i$  = standardized classification function coefficients

$$\text{Eigenvalues} = \frac{\text{Between-groups sums of squares}}{\text{Within-groups sums of squares}}$$

$$\text{Wilks' Lambda} = \frac{\text{Within-groups sums of squares}}{\text{Total sums of squares}}$$

$$\text{Percent of Variance} = \frac{\text{Eigenvalue}}{p} \times 100$$

## Principal Components

### *Principal Components for Unstandardized Variables*

$$P_{jk} = \sum_{i=1}^p W_{ji} X_i = W_{j1} X_1 + W_{j2} X_2 + \dots + W_{jp} X_p$$

where

$p$  = number of variables

$X_i$  = unstandardized variables

$W_j$  = component weights

### *Principal Components for Standardized Variables*

$$P_{jk} = \sum_{i=1}^p W_{ji} Z_i = W_{j1} Z_1 + W_{j2} Z_2 + \dots + W_{jp} Z_p$$

where

$p$  = number of variables

$Z_i$  = standardized variables

$W_j$  = component weights

$$\text{Percent of Variance} = \frac{\text{Eigenvalue}}{p} \times 100$$

# Advanced Regression Calculations

## Notation Common to All Analyses

$\underline{Y}$  = vector of  $n$  observations for the dependent variable

$\underline{X}$  =  $n$ -by- $p$  matrix of observations for  $p-1$  independent variables and the constant term, if any

$\underline{W}$  = vector of weights (default =  $\underline{1}$  )

$\underline{\beta}$  =  $p \times 1$  vector of unknown model coefficients

$\underline{\varepsilon}$  =  $n \times 1$  vector of random errors

$\underline{b}$  =  $p \times 1$  vector of estimated model coefficients

$s^2 \{b\}$  = estimated variance-covariance matrix

MSE = mean squared error

$\hat{\underline{Y}}$  = vector of predicted values for  $\underline{Y}$

$\underline{e}$  = vector of residuals

$\underline{h}$  = vector of leverages

$\underline{d}$  = vector of studentized residuals

**Mean:**

$$\bar{Y} = \frac{\sum_{i=1}^n Y_i w_i}{\sum_{i=1}^n w_i}$$

$\bar{\underline{X}}$  = vector of column means for  $\underline{X}$

**General Linear Model:**

$$\underline{Y} = \underline{X} \underline{\beta} + \underline{\varepsilon}$$

**Least Squares Estimates:**

$$\underline{b} = (\underline{X}' \underline{W} \underline{X})^{-1} \underline{X}' \underline{W} \underline{Y}$$

**Estimated Variance-Covariance Matrix:**

$$s^2[\underline{b}] = \text{MSE} (\underline{X}' \underline{W} \underline{X})^{-1}$$

where

$$\text{MSE} = \frac{\sum_{i=1}^n w_i (Y_i - \hat{Y}_i)^2}{n - p} = \frac{\text{SSE}}{n - p}$$

**Predicted Values:**

$$\hat{\underline{Y}} = \underline{X} \underline{b}$$

**Residuals:**

$$\underline{e} = \underline{Y} - \hat{\underline{Y}}$$

**R-Squared:**

$$R^2 = \frac{SSTO - SSE}{SSTO}$$

where

$$SSTO = \begin{cases} \sum_{i=1}^n w_i (Y_i - \bar{Y})^2 & \text{if constant in model} \\ \sum_{i=1}^n w_i Y_i^2 & \text{if no constant} \end{cases}$$

**Adjusted R-Squared:**

$$1 - \left( \frac{n-1}{n-p} \right) (1 - R^2)$$

**Standard Error of Estimate:**

$$SE = \sqrt{MSE}$$

**Durbin-Watson Statistic:**

$$D = \frac{\sum_{i=1}^{n-1} (e_{i+1} - e_i)^2}{\sum_{i=1}^n e_i^2}$$

**Mean Absolute Error:**

$$MAE = \frac{\sum_{i=1}^n |e_i| w_i}{\sum_{i=1}^n w_i}$$

**Mean Error:**

$$ME = \frac{\sum_{i=1}^n e_i w_i}{\sum_{i=1}^n w_i}$$

**Mean Percentage Error:**

$$MPE = \frac{100 \sum_{i=1}^n \frac{e_i}{Y_i} w_i}{\sum_{i=1}^n w_i}$$

**Mean Absolute Percentage Error:**

$$MAPE = \frac{100 \sum_{i=1}^n \left| \frac{e_i}{Y_i} \right| w_i}{\sum_{i=1}^n w_i}$$

**Forecasts:**

$\underline{\tilde{X}}_h$  = m-by-p matrix of independent variables  
for m predictions

**Predicted Value:**

$$\underline{\hat{Y}}_h = \underline{\tilde{X}}_h \underline{\hat{b}}$$

**Standard Error of Prediction:**

$$S(\underline{\hat{Y}}_{h(new)}) = \sqrt{\text{diagonal elements of MSE } (1 + \underline{\tilde{X}}_h (\underline{\tilde{X}}' \underline{\tilde{W}} \underline{\tilde{X}})^{-1} \underline{\tilde{X}}_h')}$$



**Standard Error of Mean Response:**

$$S(\hat{Y}_{\tilde{h}}) = \sqrt{\text{diagonal elements of MSE } (\tilde{X}_h (\tilde{X}' \tilde{W} \tilde{X})^{-1} \tilde{X}_h')}$$

**Variance Inflation Factors:**

$$\text{VIF}_i = \frac{1}{1 - R_i^2}$$

where

$R_i^2$  = coefficient of multiple determination when regressing  $x_i$  against other regressor variables

**Leverage:**

$$h_{\tilde{}} = \text{diagonal elements of } \tilde{W} \tilde{X} (\tilde{X}' \tilde{W} \tilde{X})^{-1} \tilde{X}'$$

**Studentized Residuals:**

$$d_i = e_i \sqrt{w_i} \left[ \frac{n - p - 1}{\text{SSE}(1 - h_i) - e_i^2 w_i} \right]^{\frac{1}{2}}$$

where

$$\text{SSE} = \sum_{i=1}^n w_i (Y_i - \hat{Y})^2$$

**DFITS:**

$$\text{DFITS}_i = d_i \sqrt{\frac{h_i}{1 - h_i}}$$

**Mahalanobis Distance:**

$$\text{MD}_i = \frac{n \begin{pmatrix} n - 2 \\ \sum_{i=1}^n w_i \end{pmatrix} \left( h_i - \frac{w_i}{n} \right)}{(n - 1) (1 - h_i)}$$

## Calibration Models

A 100 (1 -  $\alpha$ )% **Prediction Interval for X given Y:**

$$\frac{Y - a}{b} \pm \frac{t(1 - \frac{\alpha}{2}, n - 2) \sqrt{MSE}}{b} \sqrt{\frac{1}{m} + \frac{1}{n} + \left(\frac{Y - \bar{Y}_w}{b}\right)^2 \frac{1}{S_{wXX}}}$$

where

$S_{wXX}$  is the weighted sum of squares =  $\sum w_i (X - \bar{X}_w)^2$ ,

$\bar{X}_w$ ,  $\bar{Y}_w$  denote weighted averages,

and m is the Mean Size or Weight for the new observation Y.

Refer to Caulcutt & Boddy (1983) for details.

## General Linear Models

**Indicator Variables for Categorical Factors:**

k = number of levels

Construct k - 1 indicator variable where

$$I_1 = \begin{cases} 1 & \text{if level} = 1 \\ -1 & \text{if level} = k \\ 0 & \text{otherwise} \end{cases}$$

$$I_2 = \begin{cases} 1 & \text{if level} = 2 \\ -1 & \text{if level} = k \\ 0 & \text{otherwise} \end{cases}$$

...

**Multiple Comparison Intervals:**

ith contrast

$$L_i = \sum_{j=1}^k c_{ij} \mu_j$$

k = number of factor levels

$$\hat{L}_i = \sum_{j=1}^k c_{ij} \hat{\mu}_j$$

where

$\hat{\mu}_j$  = predicted value of Y for level j when all quantitative factors are set at their mean levels

$s(\hat{L}_i)$  = standard error of estimated contrast

$v$  = number of degrees of freedom for error

$w$  = number of contrasts

#### **Least Significant Differences (LSD) Intervals:**

$$\hat{L}_i \pm t s(\hat{L}_i)$$

where

$t$  is the upper  $\alpha/2$  critical value of Student's t distribution with  $v$  degrees of freedom.

#### **Tukey Intervals:**

$$\hat{L}_i \pm \frac{Q}{\sqrt{2}} s(\hat{L}_i)$$

where

$Q$  is the upper  $\alpha$  critical value of studentized range distribution with parameters  $k$  and  $v$ .

#### **Dunnnett Intervals:**

$$\hat{L}_i \pm |d| s(\hat{L}_i)$$

where

$|d|$  is the upper  $\alpha$  critical value of  $(k - 1)$  - variate Student's t distribution with  $v$  degrees of freedom

**Scheffé Intervals:**

$$\hat{L} \pm S s(\hat{L})$$

where

$S^2 = (k - 1) F$  and  $F$  is the upper  $\alpha$  critical value of Snedecor's  $F$  distribution with  $k - 1$  and  $v$  degrees of freedom

**Bonferroni Intervals:**

$$\hat{L} \pm B s(\hat{L})$$

where

$B =$  upper  $\alpha/2w$  critical value of student's  $t$  distribution with  $v$  degrees of freedom

**Multivariate t Intervals:**

$$\hat{L} \pm P s(\hat{L})$$

where

$P =$  upper  $\alpha/2$  critical value of  $w$ -variate  $t$  distribution with  $v$  degrees of freedom.

**MANOVA Tests:**

$H =$  the hypothesis matrix

$E =$  the error matrix

$p =$  the number of dependent variables

$df_h =$  the degrees of freedom for  $H$

$df_e =$  the degrees of freedom for  $E$

$S = \min(df_h, p)$

$$m = \frac{|df_h - p| - 1}{2}$$

$$n = \frac{df_e - p - 1}{2}$$

$$r = df_e - \frac{p - df_h + 1}{2}$$

$$u = \frac{p \times df_h - 2}{4}$$

$$t = \sqrt{\frac{p^2 df_h^2 - 4}{p^2 + df_h^2 - 5}}, \quad \text{if } p^2 + df_h^2 - 5 > 0$$

$t = 1$ , otherwise,

eigenvalues of  $E^{-1} H$  are  $\lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \dots \geq \lambda_p$

**Wilk's Lambda:**

$$A = \frac{1}{1 + \lambda_1} \times \frac{1}{1 + \lambda_2} \times \dots \times \frac{1}{1 + \lambda_p}$$

$$F = \frac{1 - \Lambda^{1/2}}{\Lambda^{1/2}} \times \frac{rt - 2u}{p \times df_h}, \text{ with } p \times df_h \text{ and } (rt - 2u) \text{ df}$$

**Hotelling-Lawley Trace:**

$$U = \lambda_1 + \lambda_2 + \dots + \lambda_p$$

$$F = 2(Sn + 1) \frac{U}{S^2(2m + S + 1)}, \text{ with } S(2m + S + 1) \text{ and } 2(Sn + 1) \text{ df}$$

**Pillai's Trace:**

$$V = \frac{\lambda_1}{1 + \lambda_1} \times \frac{\lambda_2}{1 + \lambda_2} \times \dots \times \frac{\lambda_p}{1 + \lambda_p}$$

$$F = \frac{2n + S + 1}{2m + S + 1} \frac{V}{S - V}, \text{ with } S(2m + S + 1) \text{ and } S(2n + S + 1) \text{ df}$$

**Roy's Greatest Root Test:**

$$\max(\lambda_1, \lambda_2, \dots, \lambda_p)$$

# Logistic Regression

**Model:**

$$E(\tilde{Y}) = \frac{1}{1 + \exp(-\tilde{X} \tilde{\beta})}$$

**Logistic Transformation:**

$$\log\left(\frac{E(\tilde{Y})}{1 - E(\tilde{Y})}\right) = \tilde{X} \tilde{\beta}$$

**Weighted Least Squares:**

$Y_i = p_i$  = (sample proportions)

$n_i$  = samples sizes

$$w_i = \frac{n_i p_i}{1 - p_i}$$

by default, all  $p_i$  restricted to the range (user controlled)

$$\left(\frac{1}{2n_i}, 1 - \frac{1}{2n_i}\right)$$

$n$  = number of samples

**Maximum Likelihood (Proportions):**

Maximize

$$L(\hat{\beta}) = \prod_{i=1}^n \left( \frac{1}{1 + e^{-x_i \beta}} \right)^{r_i} \left( \frac{e^{-x_i \beta}}{1 + e^{-x_i \beta}} \right)^{n_i - r_i}$$

where

$$r_i = n_i p_i$$

**Maximum Likelihood (0's and 1's):**

Maximize

$$L(\hat{\beta}) = \left( \frac{\prod_{i=1}^{n_i} e^{x_i \beta}}{\prod_{i=1}^n (1 + e^{x_i \beta})} \right)$$

where

$n_i$  = number of successes

**Likelihood Ratio Test:**

$$\lambda(\beta) = -2 \ln \left[ \frac{L(\hat{\beta})}{L(\hat{p})} \right]$$

where

$$L(\hat{p}) = \prod_{i=1}^n (Y_i)^{Y_i} (1 - Y_i)^{(1 - Y_i)}$$

**Percentage of Deviance:**

$$R^2 = \frac{\lambda(\beta_1, \beta_2, \dots, \beta_{p-1} | \beta_0)}{\lambda(\beta_0)}$$

**Adjusted Percentage of Deviance:**

$$\text{Adjusted } R^2 = \frac{\lambda(\beta_1, \beta_2, \dots, \beta_{p-1} | \beta_0) - 2p}{\lambda(\beta_0)}$$

## Nonlinear Regression

**Notation:**

$F(X, \hat{\beta})$  = values of nonlinear function using parameter estimates

**Marquardt Algorithm:**

Estimated coefficients are obtained by minimizing the residual sum of squares using a search procedure suggested by Marquardt. This is a

compromise between the Gauss-Newton linearization and steepest descent methods.

The user specifies:

1. Initial value of Marquardt parameter  $\lambda$ .
2. Multiple  $p$  used to modify Marquardt parameter after each iteration, where

$$\lambda_i = \frac{\lambda_{i-1}}{p}$$

provided SSE decreases, or

$$\lambda_i = p \lambda_{i-1} \text{ if SSE increases.}$$

3. Convergence Criterion #1 — estimation stops if relative change in SSE is less than this criterion.
4. Convergence Criterion # 2 — estimation stops if relative change in all parameter estimates is less than this criterion.
5. Initial estimates  $\beta_0$

Find revised estimates where

$$\beta_i = \beta_{i-1} + \Delta$$

such that

$$\text{SSE}(\beta_i) < \text{SSE}(\beta_{i-1})$$

where

$$\Delta = (X'X + \lambda \text{diag}(X'X))^{-1} X'e$$

An extended discussion of nonlinear estimation can be found in Draper and Smith, 1981.

## Regression Model Selection

$$C_p = \frac{\text{SSE}_p}{\text{MSE}_F} - (n - 2p)$$

where

$F$  = the full model

$n$  = the number of observations

$p$  = 1 + the number of independent variables



## Ridge Regression

### Standardized Coefficient Estimates:

If  $\tilde{Z} = n \times p$  matrix of independent variables standardized  
so that  $\tilde{Z}'\tilde{Z}$  equals the correlation matrix

and  $\Theta$  = value of the ridge parameter

then  $\tilde{b}'(\Theta) = (\tilde{Z}'\tilde{Z} + \Theta I_p)^{-1} \tilde{Z}'\tilde{Y}$

where

$I_p$  is a  $p \times p$  identity matrix.

### Natural Coefficients:

For the  $i$ th independent variable,

$$b_i = \frac{b'_i}{S_i} \text{ for } i = 1 \dots \text{number of independent variables}$$

for the constant,

$$b_o = b'_o - \frac{b'_1 \bar{x}_1}{S_1} - \frac{b'_2 \bar{x}_2}{S_2} - \dots - \frac{b'_k \bar{x}_k}{S_k}$$

### Variance Inflation Factors:

$$VIF_1 = \frac{1}{1 - R_1^2}$$

where

$R_1^2$  = the coefficient of multiple determination  $R^2$  when regressing  $X_i$   
against the other independent variables