

# NAG Library Function Document

## nag\_blgm\_lm\_design\_matrix (g22ycc)

**Note:** please be advised that this function is classed as ‘experimental’ and its interface may be developed further in the future. Please see Section 3.1.1 in How to Use the NAG Library and its Documentation for further information.

### 1 Purpose

**nag\_blgm\_lm\_design\_matrix (g22ycc)** generates a design matrix from a data matrix and model description.

### 2 Specification

```
#include <nag.h>
#include <nagg22.h>

void nag_blgm_lm_design_matrix (void *hform, void *hddesc,
    const double dat[], Integer pddat, Integer sddat, void **hxdesc,
    double x[], Integer pdx, Integer sdx, Integer *mx, NagError *fail)
```

### 3 Description

**nag\_blgm\_lm\_design\_matrix (g22ycc)** generates a design matrix from a data matrix and a model description. Design matrices encapsulate the observed values of the independent variables and the required model in a form that can be used by many of the model fitting functions available in the NAG Library, for example those in Chapter g02.

#### 3.1 Notation

Let  $D$  denote a data matrix with  $n$  observations on  $m_d$  independent variables, denoted by  $V_j$ , for  $j = 1, 2, \dots, m_d$ . If  $V_j$  is a categorical variable, let  $L_j$  denote the number of levels associated with it. If  $V_j$  is a binary, ordinal or continuous variable, let  $L_j = 1$ .

Let  $V_{ji}$  denote the  $i$ th value of  $V_j$ .

Let  $\mathcal{M}$  denote a model made up of one or more terms, denoted by  $T_i$ . Each term consists of either a main effect or an interaction and hence can be described using one or more variable names  $V_j$  and the interaction operator ‘.’. The operator ‘+’ is used to denote the addition of a term to the model. Therefore,  $\mathcal{M} = T_1 + T_2 + T_3 = V_1 + V_2 + V_1.V_2$  denotes a model with three terms, the first two terms being the main effects for variables  $V_1$  and  $V_2$  and the last term the interaction between them. For simplicity we reorder the terms of the model by the number of variables in them, so main effects come first, then two-way interactions, then three-way interactions etc. By default it is assumed that the model  $\mathcal{M}$  contains a mean effect (or intercept term), if the mean effect is excluded, this will be denoted by ‘-1’, so  $\mathcal{M} = T_1$  is a model with one term and a mean effect and  $\mathcal{M} = T_1 - 1$  is the same model with the mean effect dropped.

**nag\_blgm\_lm\_design\_matrix (g22ycc)** generates an  $n$  by  $m_x$  design matrix,  $X$ , from  $D$  and  $\mathcal{M}$ .

#### 3.2 Dummy Variables

When constructing a design matrix, we cannot work directly with categorical variables. Categorical variables must first be recoded into dummy variables. A categorical variable  $V_j$  requires  $L_j$  dummy variables. Let  $\mathcal{D}^j$  denote an  $n \times L_j$  matrix of dummy variables for  $V_j$  defined as

$$\mathcal{D}_{li}^j = \begin{cases} 1; & \text{if } V_{ji} = l, \\ 0; & \text{otherwise} \end{cases}$$

where  $\mathcal{D}_l^j$  is the  $l$ th column of  $\mathcal{D}^j$  and  $\mathcal{D}_{li}^j$  is the  $i$ th element of  $\mathcal{D}_l^j$ .

For a binary, ordinal or continuous variable,  $\mathcal{D}_{1i}^j = V_{ji}$ .

### 3.3 Full Design Matrix

Given a model,  $\mathcal{M}$ , and the matrices of dummy variables constructing the full design matrix  $X_F$  is trivial. Each term is processed in order and

1. If term  $i$  is a main effect, that is  $T_i = V_j$  for some  $j$ ,  $\mathcal{D}^j$  is copied into  $X_F$ .
2. If term  $i$  is a two-way interaction, that is  $T_i = V_j.V_k$ , for some  $j \neq k$ , then
  - (i) Loop over  $l_j = 1, 2, \dots, L_j$ .
  - (ii) Loop over  $l_k = 1, 2, \dots, L_k$ .
  - (iii) Add a column to  $X_F$  corresponding to the element-wise product of  $\mathcal{D}_{l_j}^j$  and  $\mathcal{D}_{l_k}^k$ .
3. Higher interaction terms are handled in a similar manner as the two-way interactions by adding columns constructed from multiplying all combinations of the columns of the corresponding  $\mathcal{D}$ s that correspond to the variables involved. In all cases, the variables towards the right hand side of a term are iterated over the quickest.

### 3.4 Contrasts

Using the full design matrix  $X_F$  in an analysis can result in an overparameterized model. This is due to  $X_F$  often not being of full rank as the sum of all the dummy variables for a particular variable is a vector of ones. This source of overparameterization can be alleviated by using a design matrix  $X$  where (some) dummy variables are replaced by contrasts. For a categorical variable  $V_j$  the contrasts are a set of  $L_j - 1$  functionally independent linear combinations of the dummy variables.

Whilst the choice of contrasts used in term  $T_i$  will affect the individual model coefficients (parameters), it has no effect on the overall contribution of  $T_i$ .

For a given variable  $V_j$ , the contrasts can be represented by an  $L_j$  by  $L_j - 1$  matrix,  $C_j$ . The rows of  $C_j$  correspond to a particular value of  $V_j$  and the columns correspond to the values to use in the design matrix.

Six types of contrast are available in **nag\_blgm\_lm\_design\_matrix (g22ycc)**; two types of treatment contrasts, two types of sum contrasts, Helmert contrasts and polynomial contrasts. Unless specified otherwise, the contrasts used by **nag\_blgm\_lm\_design\_matrix (g22ycc)** are treatment contrasts relative to the first level. See the description of the optional parameter **Contrast** in **nag\_blgm\_lm\_formula (g22yac)** for ways of changing the contrasts used.

#### 3.4.1 Treatment Contrasts

Treatment contrasts are taken relative to either the first or last level of the variable. For example, if  $L_j = 4$ ,

$$C_j = \begin{pmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

would be the contrast matrix for  $V_j$  using treatment contrasts relative to the first level. The contrast matrix obtained when using treatment contrasts relative to the last level is similar, but the row of zeros appears at the bottom and all other rows are shifted up one.

Strictly speaking, the term *contrast* implies that each row in the contrast matrix sums to zero. That is not the case for treatment contrasts, however they are included as this coding is commonly used in practice.

### 3.4.2 Sum Contrasts

Sum contrasts are similar to treatment contrasts and again can be taken relative to the first or last level of the variable. Unlike treatment contrasts, sum contrasts effectively constrain the coefficients related to the variable to sum to zero. For example, if  $L_j = 4$ ,

$$C_j = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ -1 & -1 & -1 \end{pmatrix}$$

would be the contrast matrix for  $V_j$  using treatment contrasts relative to the last level. The contrast matrix obtained when using treatment contrasts relative to the first level is similar, but the row of  $-1$ s appears at the top and all other rows are shifted down one.

### 3.4.3 Helmert Contrasts

With Helmert contrasts level  $l$  of the variable is compared with the average effect of all previous levels. For example, if  $L_j = 4$ ,

$$C_j = \begin{pmatrix} -1 & -1 & -1 \\ 1 & -1 & -1 \\ 0 & 2 & -1 \\ 0 & 0 & 3 \end{pmatrix}$$

would be the contrast matrix for  $V_j$  using Helmert contrasts.

### 3.4.4 Polynomial Contrasts

With polynomial contrasts the entries in the columns of  $C_j$  correspond in linear, quadratic, cubic, quartic, etc. terms to a hypothetical underlying numeric variable that takes equally spaced values at each level. For example, if  $L_j = 4$ ,

$$C_j = \begin{pmatrix} -0.67 & 0.50 & -0.22 \\ -0.22 & -0.50 & 0.67 \\ 0.22 & -0.50 & -0.67 \\ 0.67 & 0.50 & 0.22 \end{pmatrix}$$

would be the contrast matrix for  $V_j$  using polynomial contrasts.

### 3.4.5 When Contrasts Can Be Used

Depending on the specifics of the model,  $\mathcal{M}$ , it may not be possible to always replace the  $L_j$  dummy variables with  $L_j - 1$  contrasts for all variables in all terms and retain the same model. A simple example of this is a data matrix,  $D$ , with four observations and two variables which have two and three levels respectively. This data matrix might look something like:

$$D = \begin{pmatrix} 1 & 1 \\ 2 & 3 \\ 1 & 2 \\ 2 & 2 \end{pmatrix}$$

For the sake of argument, assume that our model contains the main effect for each variable, but does not contain a mean effect (or intercept term). So using the notation established earlier,  $\mathcal{M} = V_1 + V_2 - 1$ . The full design matrix,  $X_F$ , for this data matrix and model would be

$$X_F = \begin{pmatrix} 1 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 \end{pmatrix}$$

However,  $X_F$  is not of full rank (and hence  $\mathcal{M}$  is overparameterized) because the sum of the first two columns is a vector of ones as is the sum of the last three columns.

In order to alleviate this we might try constructing  $X_C$  where the dummy variables have been replaced by contrasts. Assuming treatment contrasts, relative to the first level, we would have

$$X_C = \begin{pmatrix} 0 & 0 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 1 & 0 \end{pmatrix}$$

However, using  $X_C$  makes an implicit assumption that the expected value of the dependent variable (the quantity being modelled) is zero when  $V_1 = 1$  and  $V_2 = 1$ . This assumption was not made when we used  $X_F$  and hence the two design matrices are not equivalent. One solution would be to use dummy variables for  $V_1$  and contrasts for  $V_2$ , which would result in a design matrix,  $X$  of

$$X = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 \end{pmatrix}$$

Using  $X$  would give an equivalent model to using  $X_F$ .

The algorithm used by **nag\_blgm\_lm\_design\_matrix (g22ycc)** to decide which variables, in which terms, can be coded as contrasts and which need to be coded as dummy variables is described below.

Suppose  $V_j$  is any variable that appears in term  $T_i$ , let  $T_{i(j)}$  denote the term obtained by dropping  $V_j$  from  $T_i$ . For example, if  $T_3 = V_1.V_2.V_3$ ,  $T_{3(2)} = V_1.V_3$ . In this context, the empty term is taken to be the mean effect (or intercept term). We say that  $T_{i(j)}$  appears in  $\mathcal{M}$  if there exists a term  $T_k$ ,  $k < i$ , that contains all of the variables appearing in  $T_{i(j)}$ . In most cases  $T_k = T_{i(j)}$ , but this is not required. Note, as stated earlier, the terms in  $\mathcal{M}$  are ordered by the number of variables in them.

A variable,  $V_j$  in term  $T_i$  is coded by contrasts if  $T_{i(j)}$  appears in  $\mathcal{M}$  and by dummy variables otherwise. It is therefore possible for variable  $V_j$  to be coded by contrasts in some terms and dummy variables in others within the same  $X$ .

The above rule assumes the presence of a mean effect. If no such effect is present in the model, the main effect of the first categorical variable is coded by dummy variables to compensate. If no main effects appear in the model, the warning **fail.code** = NW\_POTENTIAL\_PROBLEM is returned.

A longer description and informal proof that the resulting  $X$  is a suitable design matrix for the model of interest can be found in chapter two of Chambers and Hastie (1992).

### 3.5 Mean Effect

The mean effect (or intercept term) is included in a design matrix by adding a column of ones as the first column of  $X$ . However, many model fitting functions in the NAG Library handle the mean effect as a special case and do not require it to be explicitly added to the design matrix. Therefore, by default, **nag\_blgm\_lm\_design\_matrix (g22ycc)** does not explicitly add the mean effect to the design matrix. This behaviour can be changed via the optional parameter **Explicit Mean** in **nag\_blgm\_lm\_formula (g22yac)**.

## 4 References

Chambers J M and Hastie T J (1992) *Statistical Models in S* Wadsworth and Brooks/Cole Computer Science Series

## 5 Arguments

1: **hform** – void \*

*Input*

*On entry:* a G22 handle to the internal data structure containing a description of the model  $\mathcal{M}$  as returned in **hform** by **nag\_blgm\_lm\_formula (g22yac)**.

- 2: **hddesc** – void \* *Input*  
*On entry:* a G22 handle to the internal data structure containing a description of the data matrix,  $D$  as returned in **hddesc** by **nag\_blgm\_lm\_describe\_data (g22ybc)**.
- 3: **dat**[**pddat** × **sddat**] – const double *Input*  
**Note:** the  $(i, j)$ th element of the matrix is stored in **dat**[( $j - 1$ ) × **pddat** +  $i - 1$ ].  
*On entry:* the data matrix,  $D$ . By default  $D_{ij}$ , the  $i$ th value for the  $j$ th variable, for  $i = 1, 2, \dots, n$  and  $j = 1, 2, \dots, m_d$ , should be supplied in **dat**[( $j - 1$ ) × **pddat** +  $i - 1$ ].  
If the optional parameter **Storage Order**, described in **nag\_blgm\_lm\_describe\_data (g22ybc)**, is set to VAROBS,  $D_{ij}$  should be supplied in **dat**[( $i - 1$ ) × **pddat** +  $j - 1$ ].
- 4: **pddat** – Integer *Input*  
*On entry:* the stride separating matrix row elements in the array **dat**.  
*Constraints:*  
if the optional parameter **Storage Order**, described in **nag\_blgm\_lm\_describe\_data (g22ybc)**, is set to VAROBS, **pddat** ≥  $m_d$ ;  
otherwise **pddat** ≥  $n$ .
- 5: **sddat** – Integer *Input*  
*On entry:* the secondary dimension of **dat**.  
*Constraints:*  
if the optional parameter **Storage Order**, described in **nag\_blgm\_lm\_describe\_data (g22ybc)**, is set to VAROBS, **sddat** ≥  $n$ ;  
otherwise **sddat** ≥  $m_d$ .
- 6: **hxdesc** – void \*\* *Input/Output*  
*On entry:* must be set to **NULL**.  
As an alternative an existing G22 handle may be supplied in which case this function will destroy the supplied G22 handle as if **nag\_blgm\_handle\_free (g22zac)** had been called.  
*On exit:* holds a G22 handle to the internal data structure containing a description of the design matrix,  $X$ . You **must not** change the G22 handle other than through the functions in Chapter g22.
- 7: **x**[**pdx** × **sdx**] – double *Output*  
**Note:** the  $(i, j)$ th element of the matrix  $X$  is stored in **x**[( $j - 1$ ) × **pdx** +  $i - 1$ ].  
*On exit:* the design matrix,  $X$ . By default  $X_{ij}$ , the  $i$ th value for the  $j$ th column, for  $i = 1, 2, \dots, n$  and  $j = 1, 2, \dots, m_x$ , is returned in **x**[( $j - 1$ ) × **pdx** +  $i - 1$ ].  
If the optional parameter **Storage Order**, described in **nag\_blgm\_lm\_formula (g22yac)**, is set to VAROBS,  $X_{ij}$  is returned in **x**[( $i - 1$ ) × **pdx** +  $j - 1$ ].  
If **pdx** or **sdx** are too small to hold **x**, the number of columns required to hold the design matrix is returned in **mx**.  
Under some conditions it is possible to use the data matrix in place of the design matrix. Specifically, if  $D$  has no categorical variables,  $\mathcal{M}$  has only main effects and either has no mean effect or the mean effect does not need to be explicitly added to the design matrix. If **pdx** or **sdx** are too small under such circumstances, **fail.code** = NW\_ALTERNATIVE is returned and **hxdesc** is set up in such a way as to allow **dat** to be used as the design matrix. If **pdx** and **sdx** are both zero, **x** is not referenced and may be **NULL**.

- 8: **pdx** – Integer *Input*  
*On entry:* the stride separating matrix row elements in the array **x**.  
*Constraints:*  
 if the optional parameter **Storage Order**, described in **nag\_blgm\_lm\_formula (g22yac)**, is set to VAROBS, **pdx**  $\geq m_x$ ;  
 otherwise **pdx**  $\geq n$ .
- 9: **sdx** – Integer *Input*  
*On entry:* the secondary dimension of **x**.  
*Constraints:*  
 if the optional parameter **Storage Order**, described in **nag\_blgm\_lm\_formula (g22yac)**, is set to VAROBS, **sdx**  $\geq n$ ;  
 otherwise **sdx**  $\geq m_x$ .
- 10: **mx** – Integer \* *Output*  
*On exit:* the minimum number of columns required to hold the design matrix.  
 In most cases **mx** =  $m_x$ . The one exception is when **fail.code** = NW\_ALTERNATIVE, that is the size of **x** was too small but the data matrix given in **dat** can be used as the design matrix. In this case **mx** holds the number of columns that would be required if only the relevant parts of **dat** were copied into a new array.
- 11: **fail** – NagError \* *Input/Output*  
 The NAG error argument (see Section 3.7 in How to Use the NAG Library and its Documentation).

## 6 Error Indicators and Warnings

### NE\_ALLOC\_FAIL

Dynamic memory allocation failed.

See Section 2.3.1.2 in How to Use the NAG Library and its Documentation for further information.

### NE\_ARRAY\_SIZE

On entry,  $m_d = \langle value \rangle$  and **pddat** =  $\langle value \rangle$ .

Constraint: **pddat**  $\geq m_d$ .

On entry,  $m_d = \langle value \rangle$  and **sddat** =  $\langle value \rangle$ .

Constraint: **sddat**  $\geq m_d$ .

On entry,  $n = \langle value \rangle$  and **pddat** =  $\langle value \rangle$ .

Constraint: **pddat**  $\geq n$ .

On entry,  $n = \langle value \rangle$  and **sddat** =  $\langle value \rangle$ .

Constraint: **sddat**  $\geq n$ .

### NE\_BAD\_PARAM

On entry, argument  $\langle value \rangle$  had an illegal value.

### NE\_FIELD\_UNKNOWN

A variable name used when creating **hform** is not present in **hddesc**.

Variable name:  $\langle value \rangle$ .

**NE\_HANDLE**

**hddesc** has not been initialized or is corrupt.

**hddesc** is not a G22 handle as generated by **nag\_blgm\_lm\_describe\_data (g22ybc)**.

**hform** has not been initialized or is corrupt.

**hform** is not a G22 handle as generated by **nag\_blgm\_lm\_formula (g22yac)**.

On entry, **hxdesc** is not **NULL** or a recognised G22 handle.

**NE\_INTERNAL\_ERROR**

An internal error has occurred in this function. Check the function call and any array sizes. If the call is correct then please contact NAG for assistance.

See Section 2.7.6 in How to Use the NAG Library and its Documentation for further information.

**NE\_NO\_LICENCE**

Your licence key may have expired or may not have been installed correctly.

See Section 2.7.5 in How to Use the NAG Library and its Documentation for further information.

**NE\_REAL\_ARRAY**

On entry, column  $j$  of the data matrix,  $D$ , is not consistent with information supplied in **hddesc**,  $j = \langle value \rangle$ .

**NW\_ALTERNATIVE**

On entry, the size of  $x$  is too small to hold the design matrix. **dat** can be used instead.

**NW\_ARRAY\_SIZE**

On entry,  $m_x = \langle value \rangle$  and **pdx** =  $\langle value \rangle$ .

Constraint: **pdx**  $\geq m_x$ .

On entry,  $m_x = \langle value \rangle$  and **sdx** =  $\langle value \rangle$ .

Constraint: **sdx**  $\geq m_x$ .

On entry,  $n = \langle value \rangle$  and **pdx** =  $\langle value \rangle$ .

Constraint: **pdx**  $\geq n$ .

On entry,  $n = \langle value \rangle$  and **sdx** =  $\langle value \rangle$ .

Constraint: **sdx**  $\geq n$ .

**NW\_POTENTIAL\_PROBLEM**

The model contains categorical variables, but no intercept or main effects terms have been requested.

Please check the design matrix returned matches the model you require.

**7 Accuracy**

Not applicable.

**8 Parallelism and Performance**

**nag\_blgm\_lm\_design\_matrix (g22ycc)** is threaded by NAG for parallel execution in multithreaded implementations of the NAG Library.

**nag\_blgm\_lm\_design\_matrix (g22ycc)** makes calls to BLAS and/or LAPACK routines, which may be threaded within the vendor library used by this implementation. Consult the documentation for the vendor library for further information.

Please consult the x06 Chapter Introduction for information on how to control and interrogate the OpenMP environment used within this function. Please also consult the Users' Note for your implementation for any additional implementation-specific information.

## 9 Further Comments

**nag\_blgm\_lm\_submodel (g22ydc)** can be used to obtain labels for the columns of the design matrix  $X$ .

Many of the analysis functions that require a design matrix to be supplied allow submodels to be defined through the use of a vector of ones or zeros indicating whether a column of  $X$  should be included or excluded from the analyses (see for example **sx** in **nag\_regsn\_mult\_linear (g02dac)** or **nag\_glm\_normal (g02gac)**). This allows nested models to be fit without having to reconstruct the design matrix for each analysis. **nag\_blgm\_lm\_submodel (g22ydc)** offers a mechanism for constructing these vectors using submodels specified using **nag\_blgm\_lm\_formula (g22yac)**.

## 10 Example

This example creates and outputs two design matrices for a simple linear regression model. The first design matrix uses sum contrasts for all variables and the second uses a combination of polynomial and Helmert contrasts. Column labels are generated using **nag\_blgm\_lm\_submodel (g22ydc)**.

See also the examples for **nag\_blgm\_lm\_formula (g22yac)**, **nag\_blgm\_lm\_describe\_data (g22ybc)** and **nag\_blgm\_lm\_submodel (g22ydc)**.

### 10.1 Program Text

```

/* nag_blgm_lm_design_matrix (g22ycc) Example Program.
 *
 * Copyright 2017 Numerical Algorithms Group.
 *
 * Mark 26.1, 2016.
 */
/* Pre-processor includes */
#include <stdio.h>
#include <string.h>
#include <nag.h>
#include <nag_stdlib.h>
#include <nagg02.h>
#include <nagg22.h>
#include <nagx04.h>

#define MAX_FORMULA_LEN 200
#define MAX_VNAME_LEN 200
#define MAX_PLAB_LEN 200
#define MAX_OPTION_LEN 200

#define DAT(I,J) dat[j*lddat+i]
#define X(I,J) x[j*ldx+i]

char *read_line(char formula[],Integer nchar);
void print_x(Nag_IncludeIntercept intcpt,char *const plab[],Integer nobs,
            Integer mx,const double x[],Integer ldx,const char *text);

int main(void)
{
    /* Integer scalar and array declarations */
    Integer i, j, ip = 0, lddat, ldx, lisx, lplab = 0, lvinfo,
        lvnames = 0, mx, nobs, nvar, sddat, sdx, lenlab, mncat;
    Integer exit_status = 0;
    Integer *isx = 0, *levels = 0, *vinfo = 0;
    Integer tvinfo[3];

    /* Nag Types */
    NagError fail;
    Nag_IncludeIntercept intcpt;

```



```

/* Double scalar and array declarations */
double *dat = 0, *x = 0;

/* Character scalar and array declarations */
char formula[MAX_FORMULA_LEN], tcontrast[MAX_OPTION_LEN],
    optstr[MAX_OPTION_LEN];
char *pch;
char **vnames = 0, **plab = 0;

/* Void pointers */
void *hform = 0, *hddesc = 0, *hxdesc = 0;

/* Initialize the error structure */
INIT_FAIL(fail);

printf("nag_blgm_lm_design_matrix (g22ycc) Example Program Results\n\n");

/* Skip heading in data file */
#ifdef _WIN32
scanf_s("%*[\n] ");
#else
scanf("%*[\n] ");
#endif

/* Read in the formula for the full model, remove comments and */
/* call nag_blgm_lm_formula (g22yac) to parse it */
read_line(formula,MAX_FORMULA_LEN);
nag_blgm_lm_formula(&hform,formula,&fail);
if (fail.code != NE_NOERROR) {
    printf("Error from nag_blgm_lm_formula (g22yac).\n%s\n", fail.message);
    exit_status = 1;
    goto END;
}

/* Read in the contrast to use for all parameters */
read_line(tcontrast,MAX_OPTION_LEN);

/* Set up the option string */
mncat = MAX_OPTION_LEN;
strcpy(optstr, "Contrast = ");
mncat -= strlen(optstr) + 1;
strncat(optstr, tcontrast, mncat);

/* Call nag_blgm_optset (g22zmc) to set the contrast optional argument */
nag_blgm_optset(hform,optstr,&fail);
if (fail.code != NE_NOERROR) {
    printf("Error from nag_blgm_optset (g22zmc).\n%s\n", fail.message);
    exit_status = 1;
    goto END;
}

/* Read in size of the data matrix and number of variable labels supplied */
#ifdef _WIN32
scanf_s("%" NAG_IFMT "%" NAG_IFMT "%" NAG_IFMT "%*[\n] ", &nobs, &nvar,
    &lvnames);
#else
scanf("%" NAG_IFMT "%" NAG_IFMT "%" NAG_IFMT "%*[\n] ", &nobs, &nvar,
    &lvnames);
#endif

/* Allocate memory */
if (!(levels = NAG_ALLOC(nvar, Integer)) ||
    !(vnames = NAG_ALLOC(lvnames, char *))) {
    printf("Allocation failure\n");
    exit_status = -1;
    goto END;
}
for (i = 0; i < lvnames; i++)
    if (!(vnames[i] = NAG_ALLOC(MAX_VNAME_LEN, char))) {
        printf("Allocation failure\n");
        exit_status = -1;
    }

```

```

        goto END;
    }

    /* Read in number of levels and names for the variables */
    for (i = 0; i < nvar; i++) {
#ifdef _WIN32
        scanf_s("%" NAG_IFMT " ", &levels[i]);
#else
        scanf("%" NAG_IFMT " ", &levels[i]);
#endif
    }
#ifdef _WIN32
    scanf_s("%*[\n] ");
#else
    scanf("%*[\n] ");
#endif

    if (lvnames > 0) {
        for (i = 0; i < lvnames; i++)
#ifdef _WIN32
            scanf_s("%50s", vnames[i], 51);
#else
            scanf("%50s", vnames[i]);
#endif
#ifdef _WIN32
        scanf_s("%*[\n] ");
#else
        scanf("%*[\n] ");
#endif
    }

    /* Call nag_blgm_lm_describe_data (g22ybc) to get a description of */
    /* the data matrix */
    nag_blgm_lm_describe_data(&hddesc,nobs,nvar,levels,lvnames,vnames,&fail);
    if (fail.code != NE_NOERROR) {
        printf("Error from nag_blgm_lm_describe_data (g22ybc).\n%s\n",
            fail.message);
        exit_status = 1;
        goto END;
    }

    /* Read in the data matrix */
    lddat = nobs;
    sddat = nvar;
    if (!(dat = NAG_ALLOC(lddat*sddat, double))) {
        printf("Allocation failure\n");
        exit_status = -1;
        goto END;
    }
    for (i = 0; i < nobs; i++) {
        for (j = 0; j < nvar; j++)
#ifdef _WIN32
            scanf_s("%lf", &DAT(i, j));
#else
            scanf("%lf", &DAT(i, j));
#endif
    }
#ifdef _WIN32
    scanf_s("%*[\n] ");
#else
    scanf("%*[\n] ");
#endif

    /* Start of constructing the design matrix ... */

    /* Call nag_blgm_lm_design_matrix (g22ycc) to get the size of */
    /* the design matrix */
    ldx = 0;
    sdx = 0;
    nag_blgm_lm_design_matrix(hform,hddesc,dat,lddat,sddat,&hxdesc,
        x,ldx,sdx,&mx,&fail);

```

```

if (fail.code != NW_ARRAY_SIZE && fail.code != NW_ALTERNATIVE) {
    printf("Error from nag_blgm_lm_design_matrix (g22ycc).\n%s\n",
        fail.message);
    exit_status = 1;
    goto END;
}

/* Allocate design matrix */
ldx = nob;
sdx = mx;
if (!(x = NAG_ALLOC(ldx*sdx, double))) {
    printf("Allocation failure\n");
    exit_status = -1;
    goto END;
}

/* Call nag_blgm_lm_design_matrix (g22ycc) to generate the design matrix */
nag_blgm_lm_design_matrix(hform,hddesc,dat,lldat,sddat,&hxdesc,
    x,ldx,sdx,&mx,&fail);
if (fail.code != NE_NOERROR) {
    printf("Error from nag_blgm_lm_design_matrix (g22ycc).\n%s\n",
        fail.message);
    exit_status = 1;
    goto END;
}
/* ... End of constructing the design matrix */

/* Start of getting the column labels for X ... */
/* Get size of output arrays used by nag_blgm_lm_submodel (g22ydc) */
lvinfo = 3;
lisx = lplab = lenlab = 0;
nag_blgm_lm_submodel(hform,hxdesc,&intcpt,&ip,lisx,ix,lplab,plab,lenlab,
    lvinfo,tvinfo, &fail);
if (fail.code != NW_ARRAY_SIZE) {
    printf("Error from nag_blgm_lm_submodel (g22ydc).\n%s\n", fail.message);
    exit_status = 1;
    goto END;
}

/* Allocate output arrays (we only want plab) */
lisx = lvinfo = 0;
lplab = tvinfo[1];
if (!(plab = NAG_ALLOC(ip, char *))) {
    printf("Allocation failure\n");
    exit_status = -1;
    goto END;
}
lenlab = MAX_PLAB_LEN;
for (i = 0; i < ip; i++)
    if (!(plab[i] = NAG_ALLOC(lenlab, char))) {
        printf("Allocation failure\n");
        exit_status = -1;
        goto END;
    }

/* Call nag_blgm_lm_submodel (g22ydc) to generate the labels */
nag_blgm_lm_submodel(hform,hxdesc,&intcpt,&ip,lisx,ix,lplab,plab,lenlab,
    lvinfo,vinfo,&fail);
if (fail.code != NE_NOERROR) {
    printf("Error from nag_blgm_lm_submodel (g22ydc).\n%s\n", fail.message);
    exit_status = 1;
    goto END;
}
/* ... End of getting the column labels for X */

/* Display the design matrix */
print_x(intcpt,plab,nob,mx,x,ldx,"First");
printf("\n");

/* Read in the name of the variable whose contrasts need to be changed, */
/* the value to change them to, remove comments and set the contrast */

```

```

/* optional argument for the specified variable */
for (;read_line(tcontrast,MAX_OPTION_LEN);) {
    /* Add an equals sign between variable name and contrast name */
    pch = strstr(tcontrast," ");
    *pch = '=';

    /* Set up the option string */
    mncat = MAX_OPTION_LEN;
    strncpy(optstr, "Contrast:",mncat);
    mncat -= strlen(optstr) + 1;
    strcat(optstr, tcontrast, mncat);

    /* Call nag_blgm_optset (g22zmc) to set the contrast optional argument */
    nag_blgm_optset(hform,optstr,&fail);
    if (fail.code != NE_NOERROR) {
        printf("Error from nag_blgm_optset (g22zmc).\n%s\n", fail.message);
        exit_status = 1;
        goto END;
    }
}

/* Call nag_blgm_lm_design_matrix (g22ycc) to regenerate */
/* the design matrix */
nag_blgm_lm_design_matrix(hform,hddesc,dat,lldat,sddat,&hxdesc,
                        x,ldx,sdx,&mx,&fail);
if (fail.code != NE_NOERROR) {
    printf("Error from nag_blgm_lm_design_matrix (g22ycc).\n%s\n",
        fail.message);
    exit_status = 1;
    goto END;
}

/* Call nag_blgm_lm_submodel (g22ydc) to regenerate the column labels */
nag_blgm_lm_submodel(hform,hxdesc,&intcpt,&ip,lisx,isx,lplab,plab,lenlab,
                    lvinfo,vinfo,&fail);
if (fail.code != NE_NOERROR) {
    printf("Error from nag_blgm_lm_submodel (g22ydc).\n%s\n", fail.message);
    exit_status = 1;
    goto END;
}

/* Display the design matrix */
print_x(intcpt,plab,nobs,mx,x,ldx,"Second");

END:
/* Call nag_blgm_handle_free (g22zac) to clean-up the g22 handles */
nag_blgm_handle_free(&hform,&fail);
nag_blgm_handle_free(&hddesc,&fail);
nag_blgm_handle_free(&hxdesc,&fail);

NAG_FREE(dat);
NAG_FREE(x);
NAG_FREE(levels);
for (i = 0; i < lvnames; i++)
    NAG_FREE(vnames[i]);
NAG_FREE(vnames);
for (i = 0; i < ip; i++)
    NAG_FREE(plab[i]);
NAG_FREE(plab);
NAG_FREE(isx);
NAG_FREE(vinfo);
return (exit_status);
}

char *read_line(char formula[],Integer nchar) {
    /* Read in a line from stdin and remove any comments */
    char *pch;

    /* Read in the model formula */
    if (fgets(formula,nchar,stdin)) {
        /* Strip comments from formula */

```

```

    pch = strstr(formula, "::");
    if (pch) *pch = '\\setminus 0';
    return formula;
} else {
    return 0;
}
}

void print_x(Nag_IncludeIntercept intcpt, char *const plab[], Integer nobs,
            Integer mx, const double x[], Integer ldx, const char *text) {
    /* Print the transpose of the first 10 rows of the design matrix */

    /* Integer scalar and array declarations */
    Integer max_rows = 10, i, j, pnoobs, si;

    /* plab holds the labels for the model parameters, so includes a label */
    /* for the mean effect, if one is present. As the mean effect is not */
    /* being explicitly included in the design matrix, we may need to skip */
    /* the first element of plab (which will always be the label for the */
    /* mean effect if one is present) */
    si = (intcpt == Nag_Intercept) ? 1 : 0;

    /* Printing the first max_rows rows of the design matrix */
    pnoobs = MIN(max_rows, nobs);

    /* Display the design matrix */
    printf(" Transpose of First %"NAG_IFMT" Rows of the %s Design Matrix (X)\n",
           pnoobs, text);
    printf(" Column Name ");
    for (i = 0; i < pnoobs; i++)
        printf(" %2"NAG_IFMT, i+1);
    printf("\n ");
    for (i = 0; i < 15+pnoobs*5; i++)
        printf("-");
    printf("\n");
    for (j = 0; j < mx; j++) {
        printf(" %-15s", plab[j+si]);
        for (i = 0; i < pnoobs; i++)
            printf(" %4.1f", X(i,j));
        printf("\n");
    }

    /* Display the intercept flag using nag_enum_value_to_name (x04nbc) */
    printf(" Intercept flag = %s\n", nag_enum_value_to_name(intcpt));
}

```

## 10.2 Program Data

```

nag_blgm_lm_design_matrix (g22ycc) Example Program Data
F1*F2*Con - F1.F2.Con      :: formula
Sum First                 :: contrast to use
25 3 3                    :: nobs,nvar,lvnames
3 3 1                     :: levels
F1 F2 Con                 :: vnames
3 1 -2.4
3 3 0.2
1 3 -1.4
2 1 -5.4
3 3 0.2
3 2 1.4
1 2 6.8
1 2 6.7
1 1 5.3
2 3 -1.3
3 2 -3.6
3 2 -0.7
1 1 5.7
3 3 2.3
1 2 3.3
2 3 -0.5

```

```

1 1 -2.6
1 2  3.7
1 2  0.9
3 1 -1.1
2 2  2.1
1 3  4.6
2 3  4.6
1 2  5.1
1 3  0.9                                :: dat
F1 Helmert
F2 Polynomial                            :: new contrasts to use

```

### 10.3 Program Results

nag\_blgm\_lm\_design\_matrix (g22ycc) Example Program Results

```

Transpose of First 10 Rows of the First Design Matrix (X)
Column Name      1      2      3      4      5      6      7      8      9     10
-----
F1_SF1           0.0    0.0   -1.0    1.0    0.0    0.0   -1.0   -1.0   -1.0    1.0
F1_SF2           1.0    1.0   -1.0    0.0    1.0    1.0   -1.0   -1.0   -1.0    0.0
F2_SF1          -1.0    0.0    0.0   -1.0    0.0    1.0    1.0    1.0   -1.0    0.0
F2_SF2          -1.0    1.0    1.0   -1.0    1.0    0.0    0.0    0.0   -1.0    1.0
CON             -2.4    0.2   -1.4   -5.4    0.2    1.4    6.8    6.7    5.3   -1.3
F1_SF1.F2_SF1   -0.0    0.0   -0.0   -1.0    0.0    0.0   -1.0   -1.0    1.0    0.0
F1_SF1.F2_SF2   -0.0    0.0   -1.0   -1.0    0.0    0.0   -0.0   -0.0    1.0    1.0
F1_SF2.F2_SF1   -1.0    0.0   -0.0   -0.0    0.0    1.0   -1.0   -1.0    1.0    0.0
F1_SF2.F2_SF2   -1.0    1.0   -1.0   -0.0    1.0    0.0   -0.0   -0.0    1.0    0.0
F1_SF1.CON      -0.0    0.0    1.4   -5.4    0.0    0.0   -6.8   -6.7   -5.3   -1.3
F1_SF2.CON      -2.4    0.2    1.4   -0.0    0.2    1.4   -6.8   -6.7   -5.3   -0.0
F2_SF1.CON       2.4    0.0   -0.0    5.4    0.0    1.4    6.8    6.7   -5.3   -0.0
F2_SF2.CON       2.4    0.2   -1.4    5.4    0.2    0.0    0.0    0.0   -5.3   -1.3
Intercept flag = Nag_Intercept

```

```

Transpose of First 10 Rows of the Second Design Matrix (X)
Column Name      1      2      3      4      5      6      7      8      9     10
-----
F1_H1           0.0    0.0   -1.0    1.0    0.0    0.0   -1.0   -1.0   -1.0    1.0
F1_H2           2.0    2.0   -1.0   -1.0    2.0    2.0   -1.0   -1.0   -1.0   -1.0
F2_P1          -0.7    0.7    0.7   -0.7    0.7    0.0    0.0    0.0   -0.7    0.7
F2_P2           0.4    0.4    0.4    0.4    0.4   -0.8   -0.8   -0.8    0.4    0.4
CON             -2.4    0.2   -1.4   -5.4    0.2    1.4    6.8    6.7    5.3   -1.3
F1_H1.F2_P1     -0.0    0.0   -0.7   -0.7    0.0    0.0   -0.0   -0.0    0.7    0.7
F1_H1.F2_P2     0.0    0.0   -0.4    0.4    0.0   -0.0    0.8    0.8   -0.4    0.4
F1_H2.F2_P1     -1.4    1.4   -0.7    0.7    1.4    0.0   -0.0   -0.0    0.7   -0.7
F1_H2.F2_P2     0.8    0.8   -0.4   -0.4    0.8   -1.6    0.8    0.8   -0.4   -0.4
F1_H1.CON      -0.0    0.0    1.4   -5.4    0.0    0.0   -6.8   -6.7   -5.3   -1.3
F1_H2.CON      -4.8    0.4    1.4    5.4    0.4    2.8   -6.8   -6.7   -5.3    1.3
F2_P1.CON       1.7    0.1   -1.0    3.8    0.1    0.0    0.0    0.0   -3.7   -0.9
F2_P2.CON      -1.0    0.1   -0.6   -2.2    0.1   -1.1   -5.6   -5.5    2.2   -0.5
Intercept flag = Nag_Intercept

```

## 11 Optional Parameters

As well as the optional parameters common to all G22 handles described in **nag\_blgm\_optset (g22zmc)** and **nag\_blgm\_optget (g22znc)**, a number of additional optional parameters can be specified for a G22 handle holding the description of a design matrix as returned by **nag\_blgm\_lm\_design\_matrix (g22ycc)** in **hxdesc**.

The value of an optional parameter can be queried using **nag\_blgm\_optget (g22znc)**.

The remainder of this section can be skipped if you wish to use the default values for all optional parameters.

The following is a list of the optional parameters available. A full description of each optional parameter is provided in Section 11.1.

**Formula****Min Number of Columns****Number of Columns****Number of Observations****Storage Order****11.1 Description of the Optional Parameters**

For each option, we give a summary line, a description of the optional parameter and details of constraints.

The summary line contains:

a parameter value, where the letters  $a$ ,  $i$  and  $r$  denote options that take character, integer and real values respectively;

Keywords and character values are case and white space insensitive.

**Formula** $a$ 

This optional parameter returns a verbose formula string describing the model,  $\mathcal{M}$ , used to create the design matrix. This formula will only contain variable names, the operators ‘+’ and ‘.’ and any contrast identifiers present.

**Min Number of Columns** $i$ 

This optional parameter returns the minimum number of columns required to hold the design matrix,  $X$ . In most cases **Min Number of Columns** = **Number of Columns**. The one exception is when **fail.code** = NW\_ALTERNATIVE, that is the size of  $\mathbf{x}$  was too small but the data matrix given in **dat** can be used as the design matrix. In this case, **Number of Columns** =  $m_x = m_d$  and **Min Number of Columns** holds the number of columns that would be required if only the relevant parts of **dat** were copied into a new array.

**Number of Columns** $i$ 

This optional parameter returns  $m_x$ , the number of columns in the design matrix.

**Number of Observations** $i$ 

This optional parameter returns  $n$ , the number of observations in the design matrix.

**Storage Order** $a$ 

This optional parameter returns how the design matrix,  $X$ , is stored in  $\mathbf{x}$ .

If **Storage Order** = OBSVAR,  $X_{ij}$ , the value for the  $j$ th variable of the  $i$ th observation of the design matrix is stored in  $\mathbf{x}[(j-1) \times \mathbf{pdx} + i - 1]$ .

If **Storage Order** = VAROBS,  $X_{ij}$ , the value for the  $j$ th variable of the  $i$ th observation of the design matrix is stored in  $\mathbf{x}[(i-1) \times \mathbf{pdx} + j - 1]$ .

It should be noted that **Storage Order** is not writeable. If you wish to change the storage order of the design matrix you need to change **Storage Order** in **hform** as described in Section 11 in **nag\_blgm\_lm\_formula** (g22yac) prior to calling **nag\_blgm\_lm\_design\_matrix** (g22ycc).