A FULL VERSION OF THE MATHEMATICAL APPENDIX (§10) TO

THE PHYLOGENETIC REGRESSION

BY A. GRAFEN

*Animal Behaviour Research Group, Zoology Department, Oxford University*

MATHEMATICAL APPENDIX

*(a) Preliminary remarks.*

The purpose of this appendix is to state results of importance to the phylogenetic

regression. The mathematics done is all quite simple, but in order to express it

economically it has been necessary to adopt a rather formal approach. A major

notational problem is the formal treatment of an arbitrary phylogeny. In these

preliminary remarks, the meaning and relevance of the four theorems is discussed.

Throughout the appendix, formal remarks are made to explain the direction of the

developing argument. In many cases the actual objects of interest are not mentioned

in the mathematics at all. These objects are statistical tests. The first is the *standard*

*regression*. The variables involved are a y-variable y, a set of x-variables X to be

controlled for, and a set of x-variables Z to be tested for. The regression is defined by

$$\mathbf{E}(y) = \mathbf{1}_t\mu + X\beta + Z\gamma, \qquad (y\text{-}X\beta\text{-}Z\gamma) \propto \mathbf{N}(0,V),$$

where $\mathbf{1}_t$ is the constant term, $V$ is defined by

$$V_{ij}(\rho) = (1 - h_{ij}^{\rho}),$$

and $h_{ij}$ is the height in the initial working phylogeny at which the paths to species i

and j diverge. "$\propto$" is used to mean that the variance-covariance matrix of the error is

assumed only to be proportional to V, not necessarily equal to it. For the purpose of

this appendix the path segment lengths are fixed, so that $\rho$ is considered known.

The first theorem states that the standard regression is equivalent to the *long*

*regression*, defined by

$$\mathbf{E}(Ly) = S\delta + LX\beta + LZ\gamma, \qquad (Ly\text{-}S\delta\text{-}LX\beta\text{-}LZ\gamma) \propto \mathbf{N}(0,C),$$

in which L, S and C are matrices defined formally later. The two regressions are shown to be equivalent in the sense that the residual sum of squares of the long regression, concentrated for $\beta$ and $\delta$, is the same function of y, Z and $\gamma$ as the residual sum of squares of the standard regression concentrated for $\mu$ and $\beta$. This shows that the significance tests for $\gamma$=0, controlling for $\mathbf{1}_t$ and X in the case of the standard regression and for LX and S in the case of the long regression, will yield the same test statistic with the same distribution. Each datapoint in the long regression represents the deviation of a node's value from its parent node's value. The data in this form is suitable for defining the randomization test explained in §3(c). It is important that C is a diagonal matrix, so that this theorem allows the standard regression to be fitted by a package which cannot handle non-diagonal variance-covariance matrices. GLIM is such a package. The reason it is necessary to prove this first theorem is to show that the formulae for L and C are correct - their forms are far from obvious *a priori*. L represents the process of "hanging on the tree" described in §3(a).

The second and third theorems concern the *short regression*, defined by

$$\mathbf{E}(GC^{-1}Ly) = GC^{-1}LX\beta + GC^{-1}LZ\gamma, \qquad (GC^{-1}L(y\text{-}X\beta\text{-}Z\gamma)) \propto \mathbf{N}(0,I),$$

The distribution of $(GC^{-1}L(y\text{-}X\beta\text{-}Z\gamma))$ is understood as a distribution conditional on G, as G is a random matrix because it depends on the value of y. The second theorem states that the process of performing the long regression, defining the random linear contrasts $GC^{-1}$ and forming the elements of the short regression does indeed result in the same, standard, statistical test as the short regression. This is shown by proving that conditional on G, the residual in the short regression after regression of $GC^{-1}Ly$ on $GC^{-1}LX$ has the same probability density whether the randomness arises through $\varepsilon$, the error in the standard regression, as transmuted by construction of G and the formation of the short regression; or whether the randomness is assumed to arise as a $\mathbf{N}(0,I)$ variable in the short regression itself. The first reason it is necessary to prove

this theorem is to show that the formula for G is correct.  The second reason is that G is a random matrix, as it depends on ε.  In general, using contrasts G that depend on ε will violate the standard formulae for the variances and covariances of the contrasts, which rely on fixed G.  As was seen in the simulations in §5, the short regression has high mean square error in its parameter estimate under the null hypothesis, and has biassed estimates under the alternative hypothesis.  It is therefore not at all obvious that the short regression will be valid, but, as the theorem shows, it is.

The third theorem states that the short regression is equivalent to the *long regression with T*, defined by

$$\mathbf{E}(Ly) = S\delta + LX\beta + T\tau + LZ\gamma, \qquad (Ly\text{-}S\delta\text{-}LX\beta\text{-}T\theta\text{-}LZ\gamma) \propto \mathbf{N}(0,C),$$

T is a matrix representing a set of artificial variables added to the long regression to ensure that no matter what value Z may take, the residuals after regression on S, LX, T and LZ will remain proportional, within each radiation separately, to the residuals after regression on S and LX alone.  T therefore depends on y, and like G is a random matrix.  Equivalence means that the residual sum of squares for the long regression with T, concentrated for β, δ and θ, is the same function of y, Z and γ as the residual sum of squares of the short regression concentrated for β.  The theorem is proved to show that the phylogenetic regression can be interpreted as conditioning within the standard regression on the patterns of the residual in each radiation, in the sense of "pattern" explained in §3(c).

The fourth theorem shows that the randomization test explained in §3(c) and defined formally below, is equivalent to the short regression in the sense that the null distribution of the test statistic of the randomization test is also an F-distribution with the required degrees of freedom.

As well as these four results, the mathematical development defines the matrices used to construct the long and short regressions, and so formally defines the phylogenetic regression.

The published version of this appendix contained only the definitions, theorems and remarks. This full version contains in addition all lemmas and proofs needed to prove the theorems.

### *(b) Mathematical development*

A preliminary note on matrix notation. I shall define matrices as $\Delta \times \Theta$, where $\Delta$ and $\Theta$ are finite sets, rather than as m×n, were m and n are integers. An $\Delta \times \Theta$ matrix A will have elements $A_{ij}$, where $i \in \Delta$ and $j \in \Theta$. Where a matrix is defined as m×n or $\Delta \times$n, the integers m and n should be understood as shorthand for the sets $\{1, 2 \ldots m\}$ and $\{1, 2 \ldots n\}$. The advantage of this notation is that if $\Delta'$ and $\Theta'$ are subsets of $\Delta$ and $\Theta$, respectively, then a submatrix $A'$ can be concisely defined as the $\Delta' \times \Theta'$ submatrix of A.

Lemma 1. If A is an $n_A \times n_A$ matrix of full rank, B is an $n_A \times n_B$ matrix of full rank, $\lambda$ is a scalar, $n_B \leq n_A$, I represents the $n_B \times n_B$ identity matrix and $(I + \lambda B^T A^{-1} B)$ is of full rank, then

$$(A + \lambda BB^T)^{-1} = A^{-1} - \lambda A^{-1} B (I + \lambda B^T A^{-1} B)^{-1} B^T A^{-1}$$

Proof. By multiplication of the proposed inverses and collection of terms in $B(...)B^T A^{-1}$.

Lemma 2. If A is an $n_A \times n_A$ positive definite matrix, B and D are $n_A \times n_B$ and $n_A \times n_D$ matrices of full rank, $n_B + n_D = n_A$, $B^T A^{-1} D = 0$, and I is the $n_A \times n_A$ identity matrix, then

$$I - D(D^T A^{-1} D)^{-1} D^T A^{-1} = B(B^T A^{-1} B)^{-1} B^T A^{-1}$$

<u>Proof</u>. Any $n_A \times 1$ vector x can be expressed uniquely as Bb+Dd, for some $n_B \times 1$ vector b and some $n_D \times 1$ vector d, because together the columns of B and D span the whole space. The LHS pre-multiplies x into Bb, and so does the RHS. They are therefore the same matrix.

<u>Lemma</u> <u>3</u>. If A is an $n_A \times n_A$ positive definite matrix, and B is an $n_A \times n_B$ matrix of full rank, $n_B < n_A$, then there exists an $n_A \times n_D$ matrix D such that

  i)   D is of full rank

  ii)   $B^T A^{-1} D = 0$

  iii)   $n_B + n_D = n_A$

<u>Proof</u>. Let the inner product of two vectors a and b of $\mathbf{R}^{n_A}$ be defined by $a^T A^{-1} b$. Then D must be chosen so that its columns minimally span the subspace complementary to that spanned by the columns of B.

<u>Lemma</u> <u>4</u>. If A is an $n_A \times n_A$ positive definite matrix, and B, D, E and F are full rank matrices of size $n_A \times n_B$, $n_A \times n_D$, $n_A \times n_E$, and $n_A \times n_F$, respectively, and rk(B|E)=rk(B)+rk(E), then if

i) $n_D + n_B = n_A$

ii) $D^T A^{-1} B = 0$, and

iii) $F^T A^{-1}(B|E) = 0$,

then there exists an $n_D \times n_F$ matrix H such that H is of full rank and F=DH.

<u>Proof</u>. The columns of F must lie in the subspace spanned by the columns of D. The ith column of H is the linear combination of the columns of D which equals the ith

column of F, $1 \le i \le n_D$. H is of full rank because $rk(F) \le \min\{rk(D), rk(H)\}$, $rk(F) = n_F$, and $n_F \le n_D$.

<u>Definition</u> of $\Pi$, $\Pi_h$, $\Pi_s$, $\Pi_t$, $\Pi_i$, $\Pi_{di}$. These definitions are made with respect to the working phylogeny. Let $\Pi$ be the set of all nodes, $\Pi_t$ the set of species nodes, $\Pi_h$ the set of higher (i.e. non-species) nodes and $\Pi_s$ the set of all nodes except the root. Let $\Pi_i$, $i \in \Pi$, be the set of species nodes which are descendants of (or equal to) node i. Let $\Pi_{di}$, $i \in \Pi_h$, be the set of daughter nodes of node i. Associate each node with a distinct integer, to establish an arbitrary ordering over $\Pi$.

<u>Definition</u> of P, $P_i$. Let P denote the partition $\{\Pi_{di}\}_{\in \Pi_h}$ of $\Pi_s$, and let $P_i$, $i \in \Pi$, denote the partition $\{\Pi_j\}_{j \in \Pi_{di}}$ of $\Pi_i$.

<u>Definition</u> of n, $n_t$, $n_s$, $n_h$. Let n be the number of nodes in the working phylogeny, $n_t$ be the number of species nodes, $n_s$ be n-1, and $n_h$ be the number of higher nodes. Note every species is either a species node or a higher node but not both, so that $n_h + n_t = n$. It follows that $n_s - n_h = n_t - 1$.

<u>Definition</u> of $'$. Let $i' \in \Pi_h$ denote the parent node of i, $i \in \Pi_s$.

<u>Definition</u> of $\kappa_i$, $h_i$. Let $\kappa_i$ be arbitrary non-negative real numbers representing the length of the path segment between i and $i'$, $i \in \Pi_s$, with $\kappa_i > 0$ if $i \in \Pi_t$. Let $h_i$ be the summed length of the path segments between the root and i, $i \in \Pi$.

<u>Definition</u> of $\alpha(i,j)$. Let $\alpha(i,j)$ be the lowest common ancestor of i and j, $i,j \in \Pi$.

<u>Remark</u>. The $h_i$ just defined are related to the $h_{ij}$ used in the body of the paper by the relationship $h_{\alpha(i,j)} = 1 - h_{ij}^{\rho}$. The working phylogeny as used in the appendix is taken as having already undergone transformation by $\rho$.

<u>Definition</u> of $\Omega_t$, $\Omega_s$, $\Omega_{ti}$, $\Omega_{si}$. Let $\Omega_t$ be the set of column vectors with real elements indexed by $\Pi_t$, and let $\Omega_{ti}$, $i \in \Pi$, be the subspace of $\Omega_t$ with only those elements

indexed by $\Pi_i$.  Let $\Omega_s$ be the set of column vectors with real elements indexed by $\Pi_s$, and let $\Omega_{si}$, $i \in \Pi_h$, be the subspace of $\Omega_s$ with only those elements indexed by $\Pi_{di}$.

Remark.  The definitions of the various $\Pi$'s allows means at higher nodes to be dealt with in the same way as species values.  $\Omega_t$ and $\Omega_s$ are the dataspaces of the standard and long regressions, respectively.

Definition of $\mathbf{1}_{ti}$, $\mathbf{I}_{ti}$, $\mathbf{1}_{si}$.  Let $\mathbf{1}_{ti} \in \Omega_{ti}$, $i \in \Pi$, be the vector each of whose elements equals one.  Let $\mathbf{I}_{ti}$ be the identity matrix over $\Omega_{ti}$.  Let $\mathbf{1}_{si} \in \Omega_{si}$ be the vector each of whose elements equals one.

Definition of $U_i$.  Let $U_i$, $i \in \Pi$, be the $\Pi_i \times \Pi_t$ matrix defined by

$$(U_i)_{jk} = \begin{cases} 1 & j=k \\ 0 & j \neq k \end{cases}$$

Remark.  $U_i$ is a matrix which picks out from a vector $x \in \Omega_t$ those elements indexed by elements of $\Pi_i$.  $(U_i x) \in \Omega_i$, and equals x over those elements held in common.  $U_i^T$ transforms a vector $x \in \Omega_i$ into a vector which is an element of $\Omega_t$, equals x in those elements indexed in common, and equals zero elsewhere.

Lemma 5.  If $\alpha(i,k)=i$, $\alpha(j,m)=j$ and $\alpha(i,j) \neq i,j$, then $\alpha(k,m)=\alpha(i,j)$.

Proof.  This is obvious from the nature of a tree.

Definition of V.  Let V be the $\Pi_t \times \Pi_t$ matrix defined by

$$V_{ij} = h_{\alpha(i,j)}, \text{ for } i,j \in \Pi_t$$

Extension of subscript notation for V.  As an extension of the usual subscript notation, let $V_{ij}$ also be defined when i and j are not necessarily species nodes, as the $\Pi_i \times \Pi_j$ submatrix of V.  Further, let $V_i$ denote $V_{ii}$.

<u>Lemma 6</u>.  $V_{ij} = h_{\alpha(i,j)} \mathbf{1}_i \mathbf{1}_j^T$  $\alpha(i,j) \neq i,j$.

<u>Proof</u>.  By definition of V, $V_{km} = h_{\alpha(k,m)}$, $k,m \in \Pi_t$.  Under the conditions of the lemma, $\alpha(k,m) = \alpha(i,j)$ $\forall k,m$, by Lemma 5, proving the result.

<u>Remark</u>.  The restrictions in the following lemma are needed, as $V_i - h_i \mathbf{1}_i \mathbf{1}_i^T$ equals zero for species nodes, and $V_i - h_{i'} \mathbf{1}_i \mathbf{1}_i^T$ is undefined for the root node.

<u>Lemma 7</u>.  $V_i$, $i \in \Pi$, and $V_i - h_i \mathbf{1}_i \mathbf{1}_i^T$ , $i \in \Pi_h$, and $V_i - h_{i'} \mathbf{1}_i \mathbf{1}_i^T$ , $i \in \Pi_s$, are all of full rank and positive definite.

<u>Proof.</u>  For this proof, let PD denote "positive definite".  First it is shown that

$$\text{If } V_i - h_i \mathbf{1}_i \mathbf{1}_i^T \text{ is PD then } V_i - h_{i'} \mathbf{1}_i \mathbf{1}_i^T \text{ is PD, } i \in \Pi_s \tag{A1}$$

If the left hand matrix is A and the right hand is B, then we have

$$B = A + (h_i - h_{i'})\mathbf{1}_i \mathbf{1}_i^T$$

Now $(h_i - h_{i'})$ is non-negative because by definition of h it equals $\kappa_i$.  Hence for any conformable vector x,

$$x^T B x = x^T A x + (h_i - h_{i'})(x^T \mathbf{1}_i)^2$$

The second term on the right is non-negative for arbitrary x.  Hence if A is positive definite, then so is B, as required.  Similar arguments, using the non-negativity of $h_{i'}$ or $h_i$ instead of $(h_i - h_{i'})$, show that

$$\text{If } V_i - h_{i'} \mathbf{1}_i \mathbf{1}_i^T \text{ is PD then } V_i \text{ is PD, } i \in \Pi_s \tag{A2a}$$

$$\text{If } V_i - h_i \mathbf{1}_i \mathbf{1}_i^T \text{ is PD then } V_i \text{ is PD, } i \in \Pi \tag{A2b}$$

Next it is shown that

$$\text{If } V_j - h_i \mathbf{1}_j \mathbf{1}_j^T \text{ is PD for all } j \in \Pi_{di}, \text{ then } V_i - h_i \mathbf{1}_i \mathbf{1}_i^T \text{ is PD, } i \in \Pi_h \tag{A3}$$

This follows because, by Lemma 6, considered as a $\Pi_{di} \times \Pi_{di}$ matrix according to the partition $P_i$ of $\Pi_i$, $V_i - h_i \mathbf{1}_i \mathbf{1}_i^T$ is diagonal and its jth element is $V_j - h_i \mathbf{1}_j \mathbf{1}_j^T$. A block diagonal matrix is positive definite if all of its diagonal blocks are positive definite.

The three results (A1), (A2) and (A3) now allow the lemma to be proved. (A1) and (A3) show that the property that $V_i - h_{i'} \mathbf{1}_i \mathbf{1}_i^T$ is PD is inherited from daughters to parents in the sense that if all the daughters of node i possess it then so does node i, provided it is defined for that node. But for a species node, $i \in \Pi_t$, $V_i - h_{i'} \mathbf{1}_i \mathbf{1}_i^T$ is a $1 \times 1$ matrix whose element equals $\kappa_i$, which is by definition strictly positive for $i \in \Pi_t$. Hence the property that $V_i - h_{i'} \mathbf{1}_i \mathbf{1}_i^T$ is PD is possessed by all species nodes and so is inherited by all $i \in \Pi_s$. Now by (A3), $V_i - h_i \mathbf{1}_i \mathbf{1}_i^T$ is PD for $i \in \Pi_h$. The only case remaining in the statement of the lemma is $V_i$. We have now shown for every node either $V_i - h_{i'} \mathbf{1}_i \mathbf{1}_i^T$ is PD, or $V_i - h_i \mathbf{1}_i \mathbf{1}_i^T$ is PD; and by (A2a) and (A2b) this is sufficient to show that $V_i$ is PD, $i \in \Pi$. Positive definiteness has been established for all the cases in the statement of the lemma.

Finally, it is sufficient to note that a positive definite matrix must be of full rank. This completes the proof.

<u>Definition of $\sigma_i^2$</u>. Let $\sigma_i^2 = (\mathbf{1}_i^T V_i^{-1} \mathbf{1}_i)^{-1}$, $i \in \Pi$. Lemma 7 shows that both the inverses exist.

<u>Remark</u>. $\sigma_i^2$ is the sampling variance of the mean of all the species below node i, of a variable whose variance-covariance matrix is V.

<u>Lemma 8</u>.

$$\mathbf{1}_i^T (V_i - h_i \mathbf{1}_i \mathbf{1}_i^T)^{-1} \mathbf{1}_i = \sum_{j \in \Pi_{di}} (\sigma_j^2 - h_i)^{-1}, \; i \in \Pi_h$$

$$\sigma_i^2 = \left( \sum_{j \in \Pi_{di}} (\sigma_j^2 - h_i)^{-1} \right)^{-1} + h_i \, , \; i \in \Pi_h$$

Proof. The inverses employed throughout this proof are shown to exist by Lemma 7. According to Lemma 6, $V_i - h_i \mathbf{1}_i \mathbf{1}_i^T$, $i \in \Pi_h$, is diagonal when considered as a $\Pi_{di} \times \Pi_{di}$ matrix according to the partition $P_i$. Hence

$$\mathbf{1}_i^T (V_i - h_i \mathbf{1}_i \mathbf{1}_i^T)^{-1} \mathbf{1}_i = \sum_{j \in \Pi_{di}} \mathbf{1}_j^T (V_j - h_i \mathbf{1}_j \mathbf{1}_j^T)^{-1} \mathbf{1}_j$$

But $V_j - h_i \mathbf{1}_j \mathbf{1}_j^T$ can be inverted by Lemma 1, pre- and post-multiplied by $\mathbf{1}_j$, and re-arranged to give

$$\mathbf{1}_j^T (V_j - h_i \mathbf{1}_j \mathbf{1}_j^T)^{-1} \mathbf{1}_j = \frac{\mathbf{1}_j^T V_j^{-1} \mathbf{1}_j}{1 - h_i \mathbf{1}_j^T V_j^{-1} \mathbf{1}_j} = \frac{1}{\sigma_j^2 - h_i}$$

and so

$$\mathbf{1}_i^T (V_i - h_i \mathbf{1}_i \mathbf{1}_i^T)^{-1} \mathbf{1}_i = \sum_{j \in \Pi_{di}} (\sigma_j^2 - h_i)^{-1}$$

establishing the first part of the lemma. $V_i$ in the form $(V_i - h_i \mathbf{1}_i \mathbf{1}_i^T) + h_i \mathbf{1}_i \mathbf{1}_i^T$ can be inverted by Lemma 1, pre- and post-multiplied by $\mathbf{1}_i$, and then rearranged to give

$$(\mathbf{1}_i^T V_i^{-1} \mathbf{1}_i)^{-1} = (\mathbf{1}_i^T (V_i - h_i \mathbf{1}_i \mathbf{1}_i^T)^{-1} \mathbf{1}_i)^{-1} + h_i$$

Together, these last two equations establish the second part of the lemma, completing the proof.

Definition of $f_i$. Let $f_i \in \Omega_t$, $i \in \Pi$, be defined by $f_i = U_i^T V_i^{-1} \mathbf{1}_i (\mathbf{1}_i^T V_i^{-1} \mathbf{1}_i)^{-1}$.

<u>Lemma 9</u>.    $f_i^T V f_j =$    $h_{\alpha(i,j)}$    $\alpha(i,j) \neq i,j.$

$\sigma_i^2$    $\alpha(i,j)=i.$

<u>Proof</u>.  After expanding the $f_i$ and $f_j$, it is necessary to notice that $U_i\,VU_j^T = V_{ij}$, and that when $\alpha(i,j)=i$, $U_i^T\,U_i\,U_j^T = U_j^T$ and $U_j\,U_i^T\,\mathbf{1}_i = \mathbf{1}_j$ .  With the definition of $\sigma_i^2$ and Lemma 6 on the form of $V_{ij}$, the results then follow immediately by direct computation.

<u>Definition of</u> L, $L_i$, W <u>and</u> K.  Let L be a $\Pi_s \times \Pi_t$ matrix whose ith row is denoted by $L_i$, and defined by

$$L_i = [\,f_i^T \; - \; f_{i'}^T\,]$$

and let W be a $\Pi_s \times \Pi_s$ matrix defined by $W = LVL^T$.  Let K be a $\Pi_t \times \Pi_s$ matrix defined by

$$K_{ij} \quad = \quad 1 \qquad \alpha(i,j)=j$$

$$0 \qquad \text{o.w.}$$

<u>Remark</u>.  L is the matrix of linear contrasts which transforms the variables of the standard regression into the corresponding variables of the long regression.  $f_i$ is a vector which maps (by taking the inner product) a vector of species values into the mean value for species below node i.  So $f_i - f_{i'}$ produces the deviation of the mean of the species below node i from the mean of the species below the parent of node i.  If the variance covariance matrix of a random vector x is V, then that of Lx is W.  K is a matrix with a row for every species, and a column for every node except the root.  An element equal to 1 indicates that the column-node is an ancestor of (or is equal to) the row-node.

<u>Notational convention of bracketed subscripts</u>.  Any array dimension indexed by $\Pi_s$ can also be considered to be indexed by $\Pi_h$, according to the partition P of $\Pi_s$.  It is

convenient to be able to use both forms of indexing explicitly. Accordingly unbracketed subscripts will refer in the usual way to indexing by $\Pi_s$, while bracketed subscripts will refer to the partitional indexing. Thus $W_{ij}$ is a single element of the matrix W, defined for $i,j \in \Pi_s$. $W_{(i)j}$ is a $\Pi_{di} \times 1$ vector defined for $i \in \Pi_h$, $j \in \Pi_s$. $W_{(ij)}$ is the $\Pi_{di} \times \Pi_{dj}$ submatrix of W, defined for $i,j \in \Pi_h$.

<u>Lemma 10</u>.

$$W_{ij} = \quad 0 \qquad\qquad\qquad \text{if } i' \neq j'$$

$$-(\sigma_{i'}^2 - h_{i'}) \qquad\quad \text{if } i' = j', \, i \neq j$$

$$\sigma_i^2 - \sigma_{i'}^2 \qquad\qquad \text{if } i' = j', \, i = j$$

Equivalently,

$$W_{(ij)} = \quad 0 \qquad\qquad\qquad\qquad\qquad\qquad i \neq j$$

$$\text{diag}_{k \in \Pi_{di}}(\sigma_k^2 - h_i) \; - (\sigma_i^2 - h_i)\mathbf{1}_{si}\,\mathbf{1}_{si}^T \qquad\qquad i = j$$

<u>Proof</u>. It is convenient to assume without loss of generality that if $i'$ and $j'$ are ancestor and descendant, then it is $i'$ that is the ancestor. Formally, if $\alpha(i',j')=j'$, then $i'=j'$. The proof considers in turn five distinct and mutually exhaustive cases. By definition, $W_{ij}=L_i\,VL_j^T$. Expanding $L_i\,VL_j^T$ using Lemma 9 yields:

Case 1, $\alpha(i',j') \neq i'$: $L_i\,VL_j^T = h_{\alpha(i,j)} - h_{\alpha(i,j)} - h_{\alpha(i,j)} + h_{\alpha(i,j)} = 0$.

Case 2, $\alpha(i',j')=i'$, $i' \neq j'$, $\alpha(i,j')=i$: $L_i\,VL_j^T = \sigma_i^2 - \sigma_{i'}^2 - \sigma_i^2 + \sigma_{i'}^2 = 0$.

Case 3, $\alpha(i',j')=i'$, $i' \neq j'$, $\alpha(i,j') \neq i$: $L_i\,VL_j^T = h_{\alpha(i,j')} - \sigma_{i'}^2 - h_{\alpha(i,j')} + \sigma_{i'}^2 = 0$.

Case 4, $\alpha(i',j')=i'$, $i'=j'$, $i \neq j$: $L_i\,VL_j^T = h_{i'} - \sigma_{i'}^2 - \sigma_{i'}^2 + \sigma_{i'}^2 = h_{i'} - \sigma_{i'}^2$.

Case 5, $\alpha(i',j')=i'$, $i'=j'$, $i=j$: $L_i\,VL_j^T = \sigma_i^2 - \sigma_{i'}^2 - \sigma_{i'}^2 + \sigma_{i'}^2 = \sigma_i^2 - \sigma_{i'}^2$.

Cases 1, 2 and 3 all have i′≠j′, and so show the first part of the proposition. Case 4 and Case 5 demonstrate the second and third parts respectively.

Definition of C. Let C be a $\Pi_s \times \Pi_s$ matrix defined by $C = \text{diag}_{i \in \Pi_s}(\sigma_i^2 - h_{i'})$.

Definition of |. If A and B are two matrices with the same number of rows, then let A|B denote the matrix formed by juxtaposing the columns of A and B.

Definition of $M^t$, $N^t$, $M^s$, $N^s$. If A is a $\Pi_t \times n_A$ matrix of full rank, $n_A \leq n_t$, then let $M_A^t = A(A^T V^{-1} A)^{-1} A^T V^{-1}$, and let $N_A^t = I_t - M_A^t$. If A is a $\Pi_s \times n_A$ matrix of full rank, $n_A \leq n_s$, then let $M_A^s = A(A^T C^{-1} A)^{-1} A^T C^{-1}$, and let $N_A^s = I_s - M_A^s$. In each case, $rk(M_A) = rk(A)$. In each case, if A is a null matrix, then let $M_A = 0$, and $N_A = I$.

Remark. The M's and N's are orthogonal projection matrices in the $\Omega$ space indicated by their superscript. M projects onto the columns of the subscripted matrix, while N projects onto the space orthogonal to them. Orthogonality in $\Omega_t$ is taken with respect to $V^{-1}$, and in $\Omega_s$ is taken with respect to $C^{-1}$. The principal properties of projection matrices, which will be used without comment, are that $M_{A|B}A = A$, and $N_{A|B}A = 0$; that $M_A M_A = M_A$, $N_A N_A = N_A$ and $M_A N_A = 0$; that if the columns of A and B span the same subspace, then $M_A = M_B$ and $N_A = N_B$; and that if the columns of A are orthogonal to the columns of B then $M_A B = 0$ and $N_A B = B$.

Lemma 11.    $KL = N_{\mathbf{1}_t}^t$

Proof. The ith row of KL is

$$([f_i^T - f_{i'}^T] + [f_{i'}^T - f_{i''}^T] + [f_{i''}^T - f_{i'''}^T] + \ldots + - f_r^T]),$$

where r represents the root node. This equals $[f_i^T - f_r^T]$. $f_i^T$, $i \in \Pi_t$, by definition contains a 1 in position i, and zeroes elsewhere. $f_r^T$ is constant for all rows, and as $U_r = I_t$, $f_r^T$ equals $(\mathbf{1}_i^T V_i^{-1} \mathbf{1}_i)^{-1} \mathbf{1}_i^T V_i^{-1}$. Hence

$$KL = \mathbf{I}_t - \mathbf{1}_t(\mathbf{1}_t^T V^{-1} \mathbf{1}_t)^{-1} \mathbf{1}_t^T V^{-1},$$

which equals $N\overset{t}{\mathbf{1}_t}$ as required.

Lemma 12. $L\mathbf{1}_t = 0$

Proof. The ith row of L is $[f_i^T - f_{i'}^T]$, but calculation from the definition of $f_i$ using $U_i \mathbf{1}_t = \mathbf{1}_i$ shows that $f_i^T \mathbf{1}_t = 1$, for all $i \in \Pi$. Hence each element of $L\mathbf{1}_t$ equals $1-1=0$, as required.

Definition of S. Let S be a $\Pi_s \times \Pi_h$ matrix defined by

$$S_{ij} = \quad 1 \qquad\qquad i \in \Pi_{dj}$$

$$0 \qquad\qquad i \notin \Pi_{dj}$$

Equivalently, when considered as a $\Pi_h \times \Pi_h$ matrix according to the partition P of $\Pi_s$, S is diagonal with $S_{(i)i} = \mathbf{1}_{si}$, $i \in \Pi_h$.

Lemma 13. $W = C - S(S^T C^{-1} S)^{-1} S^T$.

Proof. C and S can be considered as $\Pi_h \times \Pi_h$ matrices according to the partition P of $\Pi_s$. The off-diagonal elements of both C and S all equal zero by definition. Hence it suffices to prove that

$$W_{(ii)} = C_{(ii)} - \mathbf{1}_{si}(\mathbf{1}_{si}^T C_{(ii)}^{-1} \mathbf{1}_{si})^{-1} \mathbf{1}_{si}^T, \quad i \in \Pi_h.$$

Consider first the off-diagonal elements. The off-diagonal elements of the left hand side are all equal to $h_i - \sigma_i^2$. The off-diagonal elements of the right hand side all equal

$$-(\mathbf{1}_{si}^T C_{(ii)}^{-1} \mathbf{1}_{si})^{-1} = -\left(\sum_{j \in \Pi_{ti}} (\sigma_j^2 - h_{j'})^{-1}\right)^{-1}$$

so the off-diagonal elements of the two sides are equal according to Lemma 8.

Consider next the diagonal elements. The diagonal element of the left hand side indexed by j, $j \in \Pi_{di}$, is $\sigma_j^2 - \sigma_i^2$ by Lemma 10. The corresponding diagonal element of the right hand side is $C_{jj}$ plus the off-diagonal element. But as these are equal in the right and left hand sides, we can write

$$c_j = \sigma_i^2 - h_{i'} - (h_i - \sigma_i^2) = \sigma_j^2 - \sigma_i^2$$

which proves the result.

Lemma 14. $S^T C^{-1} L = 0$

Proof. Let $S^T$ be partitioned into a row of submatrices $s_i$, $i \in \Pi_h$, in which $s_i$ contains (as columns) the rows of S indexed by $\Pi_{di}$. Let L be partitioned into a column of submatrices $m_i$, $i \in \Pi_h$, in which $m_i$ contains the rows of L indexed by $\Pi_{di}$. We can now write

$$S^T C^{-1} L = \sum_{i \in \Pi_h} s_i C_{(ii)}^{-1} m_i$$

We proceed by showing that each element of the sum equals zero. The element indexed by i looks like this:

$$
\begin{pmatrix}
00000... \\
. . . . . . \\
. . . . . . \\
00000... \\
11111... \\
00000... \\
. . . . . .
\end{pmatrix}
\begin{pmatrix}
C_{j_1 j_1}^{-1} (f_{j_1}^T - f_i^T) \\
C_{j_2 j_2}^{-1} (f_{j_2}^T - f_i^T) \\
C_{j_3 j_3}^{-1} (f_{j_3}^T - f_i^T) \\
C_{j_4 j_4}^{-1} (f_{j_4}^T - f_i^T) \\
............
\end{pmatrix}
$$

where $j_1$, $j_2$ etc represent the elements of $\Pi_{di}$. The rows of the matrix product corresponding to the zero rows in the left hand factor will be zero. This leaves only

the row indexed by i, which is the row of ones. The product of this row with the right hand factor is

$$\sum_{j \in \Pi_{di}} C_{jj}^{-1} (f_j - f_i)^T ,$$ (A4)

and the lemma will be proved if this can be shown to equal zero for all $i \in \Pi_h$. Expanding the $f_i$ from its definition, and using the definition of C, we obtain that (A4) equals

$$\sum_{j \in \Pi_{di}} \frac{1}{\sigma_j^2 - h_i} (\sigma_j^2 \mathbf{1}_j^T V_j^{-1} U_j - \sigma_i^2 \mathbf{1}_i^T V_i^{-1} U_i)$$

By Lemma 8,

$$\sum_{j \in \Pi_{di}} \frac{1}{\sigma_j^2 - h_i} = \frac{1}{\sigma_i^2 - h_i} ,$$

so the lemma will be proved if we can prove the equality

$$\sum_{j \in \Pi_{di}} \frac{\sigma_j^2}{\sigma_j^2 - h_i} \mathbf{1}_j^T V_j^{-1} U_j = \frac{\sigma_i^2}{\sigma_i^2 - h_i} \mathbf{1}_i^T V_i^{-1} U_i$$ (A5)

which we now proceed to do.

$V_i$ in the form $(V_i - h_i \mathbf{1}_i \mathbf{1}_i^T) + h_i \mathbf{1}_i \mathbf{1}_i^T$ can be inverted by Lemma 1, pre-multiplied by $\mathbf{1}_i$, and re-arranged to yield

$$\mathbf{1}_i^T V_i^{-1} = \frac{1}{1 + h_i \mathbf{1}_i^T (V_i - h_i \mathbf{1}_i \mathbf{1}_i^T)^{-1} \mathbf{1}_i} \mathbf{1}_i^T (V_i - h_i \mathbf{1}_i \mathbf{1}_i^T)^{-1}$$ (A6)

The scalar factor on the RHS of (A6) equals, by the first part of Lemma 8,

$$\frac{1}{1 + h_i \displaystyle\sum_{j \in \Pi_{di}} (\sigma_j^2 - h_i)^{-1}}$$

which using the second part of Lemma 8 becomes

$$\frac{\sigma_i^2 - h_i}{\sigma_i^2}$$

Using

$$(V_i - h_i \mathbf{1}_i \mathbf{1}_i^T) = \text{diag}_{j \in \Pi_{di}}(V_j - h_i \mathbf{1}_j \mathbf{1}_j^T)$$

from Lemma 6, the matrix inverse on the RHS of (A6) can be expressed using Lemma 1 as

$$\text{diag}_{j \in \Pi_{di}}(V_j^{-1} + \frac{h_i}{1 - h_i \mathbf{1}_j^T V_j^{-1} \mathbf{1}_j} V_j^{-1} \mathbf{1}_j \mathbf{1}_j^T V_j^{-1})$$

Hence the RHS as a whole equals

$$\frac{\sigma_i^2 - h_i}{\sigma_i^2} \ \text{row}_{j \in \Pi_{di}}(\frac{\sigma_i^2}{\sigma_j^2 - h_i} \mathbf{1}_j^T V_j^{-1})$$

We can therefore conclude from (A6) that

$$\frac{\sigma_i^2}{\sigma_i^2 - h_i} \ \mathbf{1}_i^T V_i^{-1} = \text{row}_{j \in \Pi_{di}}(\frac{\sigma_i^2}{\sigma_j^2 - h_i} \mathbf{1}_j^T V_j^{-1})$$

This equality differs only notationally from the equality (A5) we had to prove, and so completes the proof.

Lemma 15. If A is a $\Pi_t \times n_A$ matrix of full rank, B is a $\Pi_t \times n_B$ matrix of full rank, and A and B are linearly independent, then

$$N_{A|B}^t = N_{(N_A^t B)}^t N_A^t$$

Proof. The lemma states that residualizing on (A|B) is equivalent to residualizing on A, and then on B residualized on A. Let D be a matrix of full rank such that (D|A|B) spans $\Omega_t$ and such that $D^T V^{-1}(A|B)=0$. The existence of D is guaranteed by Lemma 3. Then there is a unique decomposition of a vector $x \in \Pi_t$ into $x = Dd + (A|B)\binom{a}{b}$ The LHS of the statement of the lemma pre-multiplies x into Dd. $N_A^t x = Dd + (0|N_A^t B)\binom{a}{b}$, and the second term is annihilated on pre-multiplication by $N_{(N_A^t B)}^t$ because $N_{(N_A^t B)}^t N_A^t B = 0$. To prove the lemma, it remains to show that $N_{(N_A^t B)}^t D = D$. By properties of projection matrices, this will be true if $N_A^t B$ and D are orthogonal. But

$$B^T N_A^t{}^T V^{-1} D = B^T (I_t - V^{-1}A(A^T V^{-1}A)^{-1}A^T)V^{-1}D,$$

and as $B^T V^{-1}D=0$ and $A^T V^{-1}D=0$, the whole expression equals zero as required. Hence D and $N_A^t B$ are orthogonal and the lemma is proved.

Lemma 16. If A is a $\Pi_t \times n_A$ matrix of full rank, and if $1_t$ and A are linearly independent, then

$$N_{1_t|A}^t{}^T V^{-1} N_{1_t|A}^t = L^T N_{LA|S}^s{}^T C^{-1} N_{LA|S}^s L$$

Proof. Expanding the projection matrices, the RHS becomes

$$L^T C^{-1} L - L^T C^{-1}(LA|S)\begin{pmatrix} A^T L^T C^{-1}LA & A^T L^T C^{-1}S \\ S^T C^{-1}LA & S^T C^{-1}S \end{pmatrix}^{-1}(LA|S)^T C^{-1}L$$

By Lemma 14, $S^T C^{-1}L=0$, and so this simplifies to

$$L^T C^{-1} L - (L^T C^{-1}LA|0)\begin{pmatrix} A^T L^T C^{-1}LA & 0 \\ 0 & S^T C^{-1}S \end{pmatrix}^{-1}(L^T C^{-1}LA|0)^T$$

and so to

$$L^T C^{-1} L - L^T C^{-1}LA(A^T L^T C^{-1}LA)^{-1}A^T L^T C^{-1}L$$

By Lemma 15, we have $N_{1_t|A}^t = N_{(N_{1_t}^t A)}^t N_{1_t}^t = (I_t - M_{(N_{1_t}^t A)}^t)N_{1_t}^t$, which is defined because the condition of the theorem that $1_t$ and A are linearly independent ensures

that $N_{1_t}^t A$ is of full rank. We can use this to expand the LHS as follows

$$N_{1_t}^t TV^{-1}N_{1_t}^t - N_{1_t}^t TV^{-1}N_{1_t}^t A(A^T N_{1_t}^t TV^{-1}N_{1_t}^t A)^{-1}A^T N_{1_t}^t TV^{-1}N_{1_t}^t$$

To show that the RHS and LHS are equal, it will therefore suffice to show that

$L^T C^{-1}L = N_{1_t}^t TV^{-1}N_{1_t}^t$. Lemma 13 implies that $WC^{-1}W = W$, and W by definition

equals $LVL^T$. Hence

$$LVL^T C^{-1}LVL^T = LVL^T$$

Pre-multiplying by K and postmultiplying by $K^T$, and applying Lemma 11 we obtain

$$N_{1_t}^t VL^T C^{-1}LVN_{1_t}^t {}^T = N_{1_t}^t VN_{1_t}^t {}^T$$

Substituting $(I_t - M_{1_t}^t)$ for $N_{1_t}^t$ in the LHS, and using $L1_t = 0$ from Lemma 12 yields

$$VL^T C^{-1}LV = N_{1_t}^t VN_{1_t}^t {}^T$$

Calculation shows that $V^{-1}N_{1_t}^t V = N_{1_t}^t {}^T$, so pre- and post-multiplication by $V^{-1}$ gives

$$L^T C^{-1}L = N_{1_t}^t TV^{-1}N_{1_t}^t$$

as required. This completes the proof.

Remark. The following theorem says that the long regression, which has C as its variance-covariance matrix, is equivalent to the standard regression. C is a diagonal matrix. This allows GLIM, for example, to handle the standard regression, even though it does not allow covariances among the errors. The LHS of the statement of the theorem is the residual sum of squares in the standard regression, concentrated for the parameter vectors for $1_t$ and X, as a function of y, Z and $\gamma$. The RHS is the residual sum of squares in the long regression, concentrated for the parameter vectors of S and LX, as a function of y, Z and $\gamma$. Before the theorem we formally define the data of the analysis. Note that the null hypothesis is implicit in the definition of y.

<u>Definition</u> of $\varepsilon$, y, $\mu$, X, $\beta$, Z, $n_X$, $n_Z$. Let $\varepsilon$ be a random variable over $\Omega_t$, distributed as $\mathbf{N}(0,V)$. Let $\mu$ be a scalar, X a $\Pi_t \times n_X$ matrix of full rank which is linearly independent of $\mathbf{1}_t$, $\beta$ an $n_X \times 1$ vector, and Z a $\Pi_t \times n_Z$ matrix. Let y be a random variable defined by $y = (\mathbf{1}_t|X)\begin{pmatrix}\mu\\\beta\end{pmatrix} + \varepsilon$.

<u>Theorem</u> 1. If $\gamma$ is an $n_Z \times 1$ vector, then

$$(y\text{-}Z\gamma)^T \, N^t_{\mathbf{1}_t|X}{}^T \, V^{-1} N^t_{\mathbf{1}_t|X}(y\text{-}Z\gamma) \; = (Ly\text{-}LZ\gamma)^T \, N^s_{LX|S}{}^T \, C^{-1} N^s_{LX|S}(Ly\text{-}LZ\gamma)$$

<u>Proof</u>. $N^s_{LX|S}$ is well defined only if $(LX|S)$ is of full rank, and this is established first. S is of full rank by construction, and by Lemma 14 $S^T C^{-1} L = 0$ so to show $(LX|S)$ of full rank it remains to show that LX is of full rank. Suppose not, then there exists a vector $a \neq 0$ such that $LXa = 0$ and so by Lemma 11 we have $KLXa = N^t_{\mathbf{1}_t} Xa = 0$. As X is of full rank, if $a \neq 0$ then $Xa \neq 0$. But only multiples of $\mathbf{1}_t$ are annihilated by $N^t_{\mathbf{1}_t}$, hence Xa is a multiple of $\mathbf{1}_t$. However, this contradicts the definition of X, which states that $\mathbf{1}_t$ and X are linearly independent.

The statement of the theorem will be true if

$$N^t_{\mathbf{1}_t|X}{}^T \, V^{-1} N^t_{\mathbf{1}_t|X} \;\; = \;\; L^T \, N^s_{LX|S}{}^T \, C^{-1} N^s_{LX|S} \, L$$

but this is the statement of Lemma 16, with X in the place of A. This completes the proof.

<u>Definition</u> of e, $\Pi_g$, $n_g$, $\Omega_g$, $\mathbf{I}_g$, $\tau$, $\lambda$. Let e be a random variable over $\Omega_s$ defined by

$$e = N^s_{S|LX} \, L\varepsilon$$

Let $\Pi_g = \{i|i \in \Pi_h, e_{(i)} \neq 0\}$. Let $n_g$ be the number of elements of $\Pi_g$. Let $\Omega_g$ be the set of column vectors with real elements indexed by $\Pi_g$. Let $\mathbf{I}_g$ be the identity matrix over $\Omega_g$. Let $j_i = \min\{j|j \in \Pi_{di}, \tau_j \neq 0\}$, $i \in \Pi_g$. Let $\tau$ be a random variable over $\Omega_s$ defined by

$$\tau_{(i)} \propto e_{(i)}, \quad \tau_{(i)}^{T} C_{(ii)}^{-1} \tau_{(i)} = 1, \quad \tau_{j_i} > 0 \qquad \qquad i \in \Pi_g$$

$$\tau_{(i)} = 0 \qquad \qquad i \notin \Pi_g$$

For $i \in \Pi_g$, the conditions define in turn the relative values of the elements of $\tau_{(i)}$, the magnitude of $\tau_{(i)}$ and the sign of $\tau_{(i)}$. Let $\lambda$ be a random variable over $\Omega_g$ defined by $e_{(i)} = \lambda_i \tau_{(i)}, i \in \Pi_g$.

Remark. It is formally possible that $\Pi_g = \{\}$, if all of the variability in y has been explained by X. In what follows I shall tacitly assume that this is not the case. In practical terms, this situation will be obvious because of a zero sum of squares in the standard regression, and in theoretical terms it has no particular interest. There is no possibility of discovering if Z explains variability in y from such a dataset.

Definition of $M^g$, $N^g$. If A is a $\Pi_g \times n_A$ matrix of full rank, $n_A \leq n_g$, then let $M_A^g = A(A^T A)^{-1} A^T$, and let $N_A^g = I_g - M_A^g$. Note that $rk(M_A) = rk(A)$. In the case that A is a null matrix, let $M_A^g = 0$, and $N_A^g = I_g$. $M^g$ and $N^g$ are orthogonal projection matrices over $\Omega_g$, and orthogonality is taken with respect to $I_g$.

Definition of G. Let G be a $\Pi_g \times \Pi_s$ matrix defined by

$$G_{i(j)} = \begin{cases} \tau_{(i)}^{T} & i = j \\ 0 & i \neq j \end{cases}$$

Remark. e is the residual in the long regression after regression of y on X. $\Pi_g$ is the set of higher nodes at which these residuals are not identically zero. The circumstances in which some of the residuals are identically zero is discussed in §3(e). Usually, $\Pi_g = \Pi_h$. $\Omega_g$ is the dataspace of the short regression. $\tau$ is a vector containing the "pattern" of the residuals, and $\lambda$ contains the "magnitudes" in the sense of §3(c). G is a matrix which in combination with C will form the linear contrasts $GC^{-1}$ which transform the long regression into the short regression. $G^T G C^{-1}$ is a projection matrix, as the following lemma shows. The short regression can therefore

be seen as a projection of the long regression onto the columns of $G^T$. The ith column of $G^T$ has zero everywhere except in the radiation of node i, and there it is proportional to $e_{(i)}$, the residuals in the long regression of y on X. This projection ensures that all the residuals after regression on $GC^{-1}LZ$ must lie in the same space, and so must be proportional, within each radiation, to $e_{(i)}$.

Lemma 17. $GC^{-1}G^T = \mathbf{I}_g$ and $GC^{-1}S = 0$

Proof. The first part is obvious in view of the diagonality of C and the definition of G. C and S are diagonal matrices according to the partition P of $\Pi_s$. From the definition of G,

$$G_{i(j)} = 0 \qquad\qquad i \neq j$$
$$(GC^{-1}S)_{ii} \propto e_{(i)}^T C^{-1}S_{(i)i} \qquad\qquad i \in \Pi_g$$

Hence if $e^T C^{-1}S = 0$, then so does $GC^{-1}S$. But from the definition of e, $e^T C^{-1}S$ equals $(Ly)^T N_{S|LX}^s {}^T C^{-1}S = (Ly)^T C^{-1} N_{S|LX}^s S$. However, $N_{S|LX}^s S = 0$ and so the second part of the lemma is proved.

Definition of $X^g$. Let $X^g$ be a $\Pi_g \times rk(GC^{-1}LX)$ matrix defined such that the columns of $X^g$ span the same subspace as those of $GC^{-1}LX$. $X^g$ may be a null matrix.

Remark. This definition is needed in case $GC^{-1}LX$ is not of full rank even though LX is. See §3(e). $X^g = GC^{-1}LX$ will satisfy the definition when $GC^{-1}LX$ is of full rank.

Lemma 18. If A is a $\Pi_t \times n_A$ matrix of full rank, and $\mathbf{1}_t$ and A are linearly independent, then $LN_{\mathbf{1}_t|A}^t = N_{S|LA}^s L$.

Proof. The lemma will be true if

$$LM_{\mathbf{1}_t|A}^t = M_{S|LA}^s L$$

which we now prove. $(\mathbf{1}_t|A)$ and $(\mathbf{1}_t|N^t_{\mathbf{1}_t} A)$ span the same subspace of $\Omega_t$, so it follows that $M^t_{\mathbf{1}_t|A} = M^t_{\mathbf{1}_t|N^t_{\mathbf{1}_t}A}$ . The LHS therefore equals

$$L(\mathbf{1}_t|N^t_{\mathbf{1}_t} A)\begin{pmatrix} \mathbf{1}^T_t V^{-1}\mathbf{1}_t & 0 \\ 0 & A^T N^t_{\mathbf{1}_t} TV^{-1}N^t_{\mathbf{1}_t}A \end{pmatrix}^{-1}(\mathbf{1}_t|N^t_{\mathbf{1}_t}A)\ TV^{-1}$$

which using $L\mathbf{1}_t = 0$ from Lemma 12, its consequence that $LN^t_{\mathbf{1}_t} = L$, and the block diagonality of the inverse of the partitioned matrix, reduces to

$$LA(A^T N^t_{\mathbf{1}_t} TV^{-1}N^t_{\mathbf{1}_t} A)^{-1}(N^t_{\mathbf{1}_t} A)^T V^{-1} \qquad (A7)$$

Because $(S|LA)$ and $(S|N^s_S LA)$ span the same subspace of $\Omega_s$, $M^s_{S|LA} = M^s_{S|N^s_S LA}$ , and so the RHS equals

$$(S|N^s_S LA)\begin{pmatrix} S^T C^{-1}S & 0 \\ 0 & A^T L^T N^s_S TC^{-1}N^s_S LA \end{pmatrix}^{-1}(S|N^s_S LA)\ TC^{-1}L$$

Using $S^T C^{-1}L = 0$ from Lemma 14 , its consequence that $N^s_S L = L$, and the block diagonality of the inverse of the partitioned matrix, this equals

$$LA(A^T L^T C^{-1}LA)^{-1}(LA)^T C^{-1}L \qquad (A8)$$

It will now be shown piecewise that (A7) equals (A8) thus proving the lemma. Both formulae begin with LA. The matrix inverses that follow are equal because $L^T C^{-1}L = N^t_{\mathbf{1}_t} TV^{-1}N^t_{\mathbf{1}_t}$ by Lemma 16, and the remaining portions are equal for the same reason and because $N^t_{\mathbf{1}_t} TV^{-1}N^t_{\mathbf{1}_t} = N^t_{\mathbf{1}_t} TV^{-1}$. This completes the proof.

<u>Lemma 19</u>. Conditional on $\tau$, $\lambda \sim \mathbf{N}(0, N^g_{Xg})$.

<u>Proof</u>. First it is established that the support of the probability distribution of e is the subspace of $\Omega_s$ orthogonal to $(S|LX)$, and that there the density is proportional to

$$\exp(-\tfrac{1}{2}\ e^T C^{-1}e)$$

Let $\eta = N_{\mathbf{1}_t|X}^t \, \varepsilon$, and $\varphi = M_{\mathbf{1}_t|X}^t \, \varepsilon$. As the cross-product $M_{\mathbf{1}_t|X}^t {}^T V^{-1} N_{\mathbf{1}_t|X}^t$ equals zero, the density of $\varepsilon$ is proportional to

$$\exp(-\tfrac{1}{2}\, \varepsilon^T V^{-1} \varepsilon) \;=\; \exp\big(-\tfrac{1}{2}\, \varepsilon^T N_{\mathbf{1}_t|X}^t {}^T V^{-1} N_{\mathbf{1}_t|X}^t \, \varepsilon - \tfrac{1}{2}\, \varepsilon^T M_{\mathbf{1}_t|X}^t {}^T V^{-1} M_{\mathbf{1}_t|X}^t \, \varepsilon\big)$$

which in turn yields

$$= \exp\big(-\tfrac{1}{2}\, \eta^T V^{-1}\eta - \tfrac{1}{2}\, \varphi^T V^{-1}\varphi\big)$$

Hence the density of $\eta$ is zero on the subspace orthogonal to $(\mathbf{1}_t|X)$, where it is proportional to

$$\exp\big(-\tfrac{1}{2}\, \eta^T V^{-1}\eta\big) \int \exp\big(-\tfrac{1}{2}\, \varphi^T V^{-1}\varphi\big)\, d\varphi$$

and so it is also proportional simply to

$$\exp\big(-\tfrac{1}{2}\, \eta^T V^{-1}\eta\big)$$

$e = L\eta$, so the density of $e$ will be as claimed if

$$\eta^T V^{-1}\eta = e^T C^{-1} e \quad \text{where } N_{\mathbf{1}_t|X}^t \, \eta = \eta \text{ and } e = L\eta$$

which we now prove. Over the relevant subspace, the LHS is equal to

$$\eta^T N_{\mathbf{1}_t|X}^t {}^T V^{-1} N_{\mathbf{1}_t|X}^t \, \eta$$

and the RHS equals

$$\eta^T N_{\mathbf{1}_t|X}^t {}^T L^T C^{-1} L N_{\mathbf{1}_t|X}^t \, \eta$$

By Lemma 18, $L N_{\mathbf{1}_t|X}^t = N_{S|LX}^s \, L$, so the RHS equals

$$\eta^T L^T N_{S|LX}^s {}^T C^{-1} N_{S|LX}^s \, L\eta$$

But Lemma 16 shows that

$$N^{t}_{\mathbf{1}_{t}|X}{}^{T} V^{-1} N^{t}_{\mathbf{1}_{t}|X} = L^{T} N^{s}_{S|LX} C^{-1} N^{s}_{S|LX} L$$

and so the two sides are equal, completing the proof of the distribution of e.

We now proceed to find the distribution of $\lambda$. From the definition of $\tau$ and $\lambda$, we have

$$e^{T}_{(i)} C^{-1}_{(ii)} e_{(i)} = \lambda_{i}^{2} \tau^{T}_{(i)} C^{-1}_{(ii)} \tau_{(i)} = \lambda_{i}^{2} \quad i \in \Pi_{g}$$

hence

$$e^{T}C^{-1}e = \lambda^{T}\lambda, \text{ and so } \exp(-\tfrac{1}{2} e^{T}C^{-1}e) = \exp(-\tfrac{1}{2} \lambda^{T}\lambda)$$

To derive the distribution of $\lambda$, it remains to establish over what subset of $\Omega_{g}$ the density is non-zero, and what the Jacobian of the transformation is. The subset of $\Omega_{g}$ is those values of $\lambda$ corresponding to an e belonging to the subspace of $\Omega_{s}$ defined by $M^{s}_{S|LX} e=0$. As $e_{(i)}=\lambda_{i}\tau_{(i)}$ this condition becomes

$$(S|LX)^{T}C^{-1} \text{col}_{i \in \Pi_{g}} \{\lambda_{i}\tau_{(i)}\} = 0$$

But $S^{T}C^{-1}G^{T}=0$ by Lemma 17, so $(S^{T}C^{-1})^{T}_{(i)} \tau_{(i)}=0$ and therefore S may be dropped. Writing $(LX)_{(i)}$ for the submatrix of LX containing only those rows indexed by $\Pi_{di}$, this condition is equivalent to

$$\sum_{i \in \Pi_{h}} \lambda_{i}(LX)^{T}_{(i)}C^{-1}_{(ii)}\tau_{(i)} = 0, \text{ or } (GC^{-1}LX)^{T}\lambda = 0, \text{ hence to } (X^{g})^{T}\lambda = 0$$

The subset of $\Omega_{g}$ over which the density of $\lambda$ is defined conditional on $\tau$ is therefore a subspace defined by $(X^{g})^{T}\lambda=0$. Conditional on $\tau$, $e_{(i)}=\lambda_{(i)}\tau_{(i)}$ implies that $de \propto d\lambda$, so there is no non-constant term in the Jacobian. We have now established that the density of $\lambda$ is proportional to

$$\exp(-\tfrac{1}{2} \lambda^{T}\lambda)$$

where $M_{Xg}^g \lambda = 0$.  A random variable with distribution $\mathbf{N}(0, N_{Xg}^g)$ has density

proportional to $\exp(-\frac{1}{2}\xi^T\xi)$ where $M_{Xg}^g \xi = 0$, and zero elsewhere.  But this is the

same density as $\lambda$, and so $\lambda$ too is distributed as $\mathbf{N}(0, N_{Xg}^g)$.  This completes the proof.

Remark.  The following theorem shows that the short regression is analytically valid

in the case where the working phylogeny is the true phylogeny, and $\rho$ is known and

taken as fixed at its true value.  It does this by giving the probability density of $N_{Xg}^g$

$GC^{-1}Ly$, the residual vector after regression of $GC^{-1}Ly$ on $GC^{-1}LX$, based on the

whole process of computing the long regression, conditioning on $\tau$, and using the

random linear contrasts $GC^{-1}$ to form the short regression.  The theorem shows that

conditional on $\tau$, the probability density is the same as it would have been if $GC^{-1}LX$

were taken as fixed, and the residual's density calculated on the basis of an error $\psi$

distributed as $\mathbf{N}(0, \mathbf{I}_g)$ in the regression $y^g = GC^{-1}LX\beta + \psi$.  This equivalence of the

residual density in the two cases establishes the exactness of the short regression for

testing for the addition of $GC^{-1}LZ$.  Note that conditional on $\tau$, $GC^{-1}LZ$ is fixed and

not random.

Theorem 2.  Conditional on $\tau$

$$N_{Xg}^g\ GC^{-1}Ly\ \sim\ \mathbf{N}(0, N_{Xg}^g)$$

Proof.  We show that $N_{Xg}^g \lambda = \lambda$, and then the theorem follows immediately from

Lemma 19.  We begin with the identity

$$Ly = L\left((\mathbf{1}_t|X)\binom{\mu}{\beta} + \varepsilon\right) = M_{S|LX}^s\ L\left((\mathbf{1}_t|X)\binom{\mu}{\beta} + \varepsilon\right) + N_{S|LX}^s\ L\left((\mathbf{1}_t|X)\binom{\mu}{\beta} + \varepsilon\right)$$

$L\mathbf{1}_t = 0$ from Lemma 12.  We also know that $N_{S|LX}^s\ L\varepsilon = e$ by definition, and $N_{S|LX}^s$

annihilates $LX$ while $M_{S|LX}^s$ preserves it.  Hence these expressions also equal

$$M_{S|LX}^s\ L(X\beta + \varepsilon) + e$$

Pre-multiplying by $GC^{-1}$ yields

$$GC^{-1}Ly = GC^{-1}M^s_{S|LX} L(X\beta+\varepsilon) + GC^{-1}e$$

$N^g_{Xg} GC^{-1}(S|LX)=0$ because $GC^{-1}S=0$ by Lemma 17, and because the columns of $X^g$ and the columns of $GC^{-1}LX$ span the same subspace of $\Omega_g$. Expansion of the projection matrix therefore shows that $N^g_{Xg}$ annihilates the first term of the RHS. Further, calculation shows that $GC^{-1}e=\lambda$. Hence pre-multiplying by $N^g_{Xg}$ gives

$$N^g_{Xg} GC^{-1}Ly = N^g_{Xg} \lambda$$

But $N^g_{Xg} \lambda=\lambda$ because by Lemma 19 the density of $\lambda$ is zero except over the subspace where this is true. This completes the proof.

Definition of T, $n_T$. Let $n_T=n_s-n_h-n_g$. Let T be a $\Pi_s \times n_T$ matrix of full rank which satisfies $T^TC^{-1}(S|G^T)=0$, and $rk(T|S|G^T)=n_s$. If $n_s-n_h-n_g=0$, then T will be a null matrix.

Remark. T will be null only when the working phylogeny is binary and as a consequence the phylogenetic regression and the standard regression are the same.

Remark. The following theorem shows that the long regression with T is equivalent to the short regression. The LHS is the sum of squares of the long regression, concentrated for the parameter vectors associated with LX, S and T. The RHS is the residual sum of squares of the short regression, concentrated for the parameter vectors associated with $GC^{-1}LX$. The introduction of A allows for possible collinearity between T and LX, or in other words that $rk(X^g)<rk(LX)$.

Theorem 3. Let A be a $\Pi_s \times rk(S|LX|T)$ matrix of full rank whose columns span the same subspace of $\Omega_s$ as the columns of $(S|LX|T)$. Then conditional on $\tau$,

$$(Ly-LZ\gamma)^T N^s_A{}^T C^{-1} N^s_A (Ly-LZ\gamma)$$

$$= (GC^{-1}Ly - GC^{-1}LZ\gamma)^T \, N_{X^g}^{g\,T} \, \mathbf{I}_g^{-1} \, N_{X^g}^g (GC^{-1}Ly - GC^{-1}LZ\gamma)$$

<u>Proof</u>. The statement of the theorem will be true if

$$N_A^{s\,T} C^{-1} N_A^s = C^{-1} G^T N_{X^g}^{g\,T} \mathbf{I}_g^{-1} N_{X^g}^g GC^{-1}$$

which by definition of the projection matrices simplifies to

$$C^{-1} N_A^s = C^{-1} G^T N_{X^g}^g GC^{-1}$$

which we now prove. Let $n_R = n_s - rk(S|T|LX) = n_s - rk(A)$. By Lemma 3, there exists a $\Pi_s \times n_R$ matrix R of full rank, which satisfies $rk(A|R) = rk(S|T|LX|R) = n_s$, and $A^T C^{-1} R = 0$ so that $(S|T|LX)^T C^{-1} R = 0$. As $rk(S|T|G^T) = n_s$, and $G^T C^{-1}(S|T) = 0$ from Lemma 17 and construction of T, by Lemma 4 there exists a $\Pi_g \times n_R$ matrix Q of full rank such that $R = G^T Q$.

By Lemma 2 the LHS equals $C^{-1} M_R^s$, which can in turn be written as

$$C^{-1} R (R^T C^{-1} R)^{-1} R^T C^{-1}$$

Substituting $G^T Q$ for R, and as $GC^{-1}G^T = \mathbf{I}_g$ by Lemma 17, this equals

$$C^{-1} G^T Q (Q^T GC^{-1}G^T Q)^{-1} Q^T GC^{-1}$$

$$= C^{-1} G^T Q (Q^T Q)^{-1} Q^T GC^{-1}$$

$$= C^{-1} G^T M_Q^g GC^{-1}$$

To prove the theorem, it therefore remains to show that

$$C^{-1} G^T M_Q^g GC^{-1} = C^{-1} G^T N_{X^g}^g GC^{-1}$$

By Lemma 2, $M_Q^g = N_{X^g}^g$ and so this equality will hold, only provided $rk(X^g|Q) = n_g$, and $Q^T X^g = 0$. $Q^T(GC^{-1}LX) = R^T C^{-1} LX$, which equals zero by the definition of R, and because the columns of $X^g$ and those of $GC^{-1}LX$ span the same subspace of $\Omega_g$,

$Q^T X^g = 0$ too. Now $GC^{-1}$ is a $\Pi_g \times \Pi_s$ matrix of full rank, and by construction of R, $rk(S|T|LX|R) = n_s$. Hence

$$rk(GC^{-1}(S|T|LX|R)) = rk(GC^{-1}S|GC^{-1}T|GC^{-1}LX|GC^{-1}R) = n_g$$

As $GC^{-1}S = 0$ by Lemma 17 and $GC^{-1}T = 0$ by construction of T, and $GC^{-1}R = GC^{-1}G^TQ = Q$ by Lemma 17, it follows that

$$rk(GC^{-1}LX|Q) = n_g,$$

and so $rk(X^g|Q) = n_g$ too, as required. This completes the proof of the theorem.

<u>Definition</u> of F. If a is a $\Pi_g \times 1$ vector, B is an $\Pi_g \times n_B$ matrix of full rank, $n_B < n_g$, D is a $\Pi_g \times n_D$ matrix of full rank, $n_B + n_D < n_g$, $rk(B|D) = rk(B) + rk(D)$, and $a^T N^g_{B|D} a \neq 0$, then let F(a,B,D) equal

$$\frac{a^T M^g_{N_B D} a}{a^T N^g_{B|D} a} = \frac{a^T M^g_{N_B D} a}{a^T N^g_B a - a^T M^g_{N_B D} a}$$

<u>Definition</u> of $Z^c$. Let $Z^c$ be a matrix of full rank formed by deleting columns from Z in such a way that $(X^g|GC^{-1}LZ^c)$ is of full rank and that $rk(X^g|GC^{-1}LZ^c) = rk(X^g|GC^{-1}LZ)$. $Z^c$ may be a null matrix.

<u>Definition</u> of $m_X$, $m_Z$. Let $m_X = rk(X^g)$, and $m_Z = rk(Z^c)$.

<u>Definition</u> of $\Gamma$. Let $\Gamma$ be defined almost surely as a $\Pi_g \times m_Z$ random matrix whose elements are independently distributed in normal distributions, with zero mean, and

$$Var(\Gamma_{ij}) = \frac{e_{(i)}^T C_{(i)}^{-1} e_{(i)}}{(e_{(i)}^T C_{(i)}^{-1} (LZ^c)_{(i)j})^2}$$

This definition fails when the denominator is zero for any i,j.

<u>Definition</u> of $\Psi$. Let $\Psi$ be a $\Pi_s \times m_Z$ random matrix defined almost surely by

$$\Psi_{(i)j} = \quad \Gamma_{ij}(LZ^c)_{(i)j} \qquad i\in\Pi_g, j=1..m_Z$$

$$0 \qquad\qquad i\in\Pi_h\text{-}\Pi_g, j=1..m_Z$$

Remark. $\Psi$ is the random alternative to Z of the randomization test described in §3(c).

Definition of $\Phi_1$, $\Phi_2$, $\Phi$. Let $\Phi_1$ be a random $\Pi_g\times n_Z$ matrix defined by $\Phi_1=GC^{-1}\Psi$, and let $\Phi_2$ be a random $\Pi_g\times(n_g\text{-}m_X\text{-}m_Z)$ matrix whose elements have normal distributions with zero mean and unit variance, independent of each other and of $\Gamma$. Let $\Phi=(\Phi_1|\Phi_2)$.

Lemma 20. $\Phi_{ij}$ have independent normal distributions with zero mean and unit variance.

Proof. Independence and zero mean follow for $\Phi_1$ because $\Phi_{ij}$ equals a constant times $\Gamma_{ij}$. Unit variance follows because that constant equals the inverse of the square root of the variance of $\Gamma_{ij}$. $\Phi_2$ has the required distribution by definition.

Remark. The following theorem shows that the randomization test described in §3(c) is the same test as the short regression. Formally, the randomization test is to find the p-value for the null hypothesis that $\gamma=0$ by finding

$$\Pr\{F(GC^{-1}Ly,X^g,GC^{-1}\Psi)>F(GC^{-1}Ly,X^g,GC^{-1}LZ^c)\},$$

in which the substitution of $\Psi$ for $LZ^c$ is the only difference between the two sides of the inequality. If this probability is low, it implies that randomly selected explanatory variables would rarely explain as much of the remaining variation in Ly as $LZ^c$ does. The short regression's test statistic would also have an F-distribution with $m_Z$ and $n_g-m_X-m_Z$ degrees of freedom.

Theorem 4. Conditional on y, $F(GC^{-1}Ly,X^g,GC^{-1}\Psi)$ has an F-distribution with $m_Z$ and $n_g-m_X-m_Z$ degrees of freedom

Proof. For brevity, let $a = GC^{-1}Ly$, $B = X^g$. By the definition of F and the relationship between the F-distribution and the Beta distribution, and noting that $GC^{-1}\Psi = \Phi_1$, it suffices to show that

$$q = \frac{a^T M^g_{N_B \Phi_1} a}{a^T N^g_B a}$$

has a Beta distribution with $m_Z$ and $n_g - m_X - m_Z$ degrees of freedom. For a definition of the Beta distribution, and its relationship to the F distribution, see Abramowitz & Stegun (1965, 26.5.1, 26.5.2 and 26.5.28). We apply Gram-Schmidt orthogonalization (see for example Kingman and Taylor 1966, page 201) to the columns of $N^g_B(\Phi_1|\Phi_2)$ in order, yielding a $\Pi_g \times (n_g - m_X - m_Z)$ random matrix $\Lambda$ conformally partitioned into $(\Lambda_1|\Lambda_2)$. There is still ambiguity in the definition of $\Lambda$, up to multiplication of each column by $\pm 1$. We complete the definition of $\Lambda$ by requiring that the inner product of each pair of corresponding columns from $N^g_B(\Phi_1|\Phi_2)$ and $\Lambda$ should be positive. Gram-Schmidt orthogonalization ensures that the columns of $\Lambda$ are mutually orthogonal, and are of unit length. Further, that the first k columns of $\Lambda$, k=1,2, .. $rk(\Lambda)$ span the same subspace as the first k columns of $N^g_B(\Phi_1|\Phi_2)$. By virtue of its construction $\Lambda$ therefore satisfies

$$\Lambda^T\Lambda = I, \ \Lambda^T B = 0,$$

Because of symmetry between the columns of $\Phi$, and the spherical symmetry of the normal distribution which $\Phi$ has by virtue of Lemma 20, the distribution of $\Lambda$ satisfies

$$\text{freq}(\Lambda) \quad \propto \quad 1 \qquad \Lambda^T\Lambda = I, \ \Lambda^T B = 0$$
$$0 \qquad \text{o.w.}$$

By construction, $\Lambda_1$ spans the same subspace as $N^g_B \Phi_1$, so there exists a random $m_Z \times m_Z$ matrix $\varsigma$ defined by $\Lambda_1 = N^g_B \Phi_1 \varsigma$. Now

$$M^g_{N^g_B \Phi_1} = M^g_{\Lambda_1 \varsigma^{-1}} = \Lambda_1 \varsigma^{-1}((\varsigma^{-1})^T \Lambda_1^T \Lambda_1 \varsigma^{-1})^{-1}(\varsigma^{-1})^T \Lambda_1^T = \Lambda_1 \Lambda_1^T$$

and so

$$q = \frac{a^T \Lambda_1 \Lambda_1^T a}{a^T N^g_B a}$$

Let $d$ be the random $(n_g-m_X) \times 1$ vector defined by $d = \Lambda^T a$, and let $d_1 = \Lambda_1^T a$. The distribution of $d$ in $\mathbf{R}^{n_g-m_X}$ is spherically symmetrical with fixed length $a^T N^g_B a$. The volume element of $d$ for which $d_1$ is a constant is the surface of an $(n_g - m_X - m_Z)$ sphere of radius $\sqrt{a^T N^g_B a - d^T d}$, and the volume element of vectors $d_1$ such that $d_1^T d_1$ equals a constant is the surface of an $m_Z$ sphere of radius $\sqrt{d_1^T d_1}$. Hence the density of $d_1^T d_1$, which equals $a^T \Lambda_1 \Lambda_1^T a$, is proportional to

$$\left(\sqrt{d_1^T d_1}\right)^{m_Z} \left(\sqrt{a^T N^g_B a - d_1^T d_1}\right)^{n_g-m_X-m_Z}$$

and so the density of $q$ is proportional to

$$q^{m_Z/2}(1-q)^{(n_g-m_X-m_Z)/2}$$

which is the Beta distribution with $m_Z$ and $n_g - m_X - m_Z$ degrees of freedom (Abramowitz & Stegun 1965, 26.5.1 and 26.5.2), as required.

### REFERENCES

Abramowitz, M. & Stegun, I.A. (ed.s) 1965 *Handbook of mathematical functions*. New York: Dover.

Kingman, J.F.C. & Taylor, S.J. 1966 *Introduction to measure and probability*. Cambridge: Cambridge University Press.