

Comparative Methods

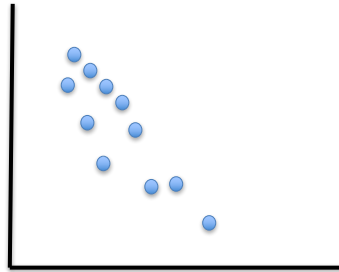
3 February 2017

SESYNC

Comparative Biology

- e.g. taller trees should have bigger seeds
- experimental evolution: often not practical
- interspecific comparison: test whether traits correlate across species
- problem: related species may share the same traits due to shared ancestry = phylogenetic non-independence
- result is that species cannot be taken as independent data points

Traditional Correlation Tests

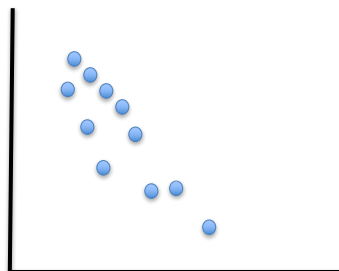


Central Assumptions in Statistics

-Normally distributed data (parametric statistics)

-Independent observations
- Points are independent

Traditional Correlation Tests

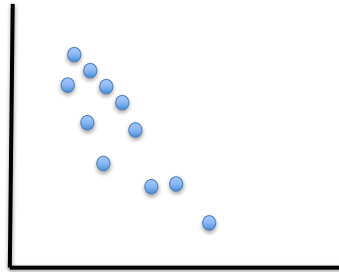


Central Assumptions in Statistics

-Normally distributed data (parametric statistics)

-Independent observations
- Points are independent

Problem with Traditional Correlation Tests



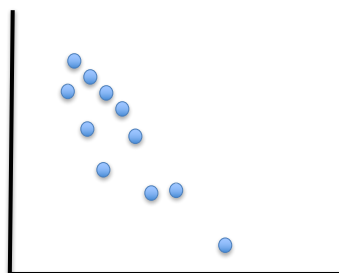
Central Assumptions in Statistics

-Normally distributed data (parametric statistics)

-Independent observations
- Points are independent

Species share evolutionary history and therefore are not independent observations

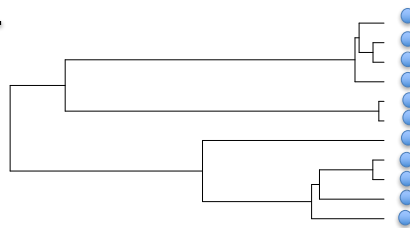
Problem with Traditional Correlation Tests



Central Assumptions in Statistics

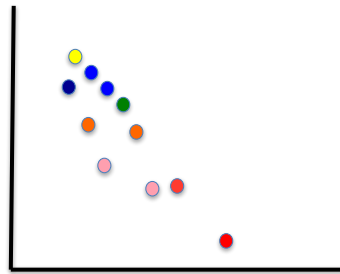
-Normally distributed data (parametric statistics)

-Independent observations
- Points are independent



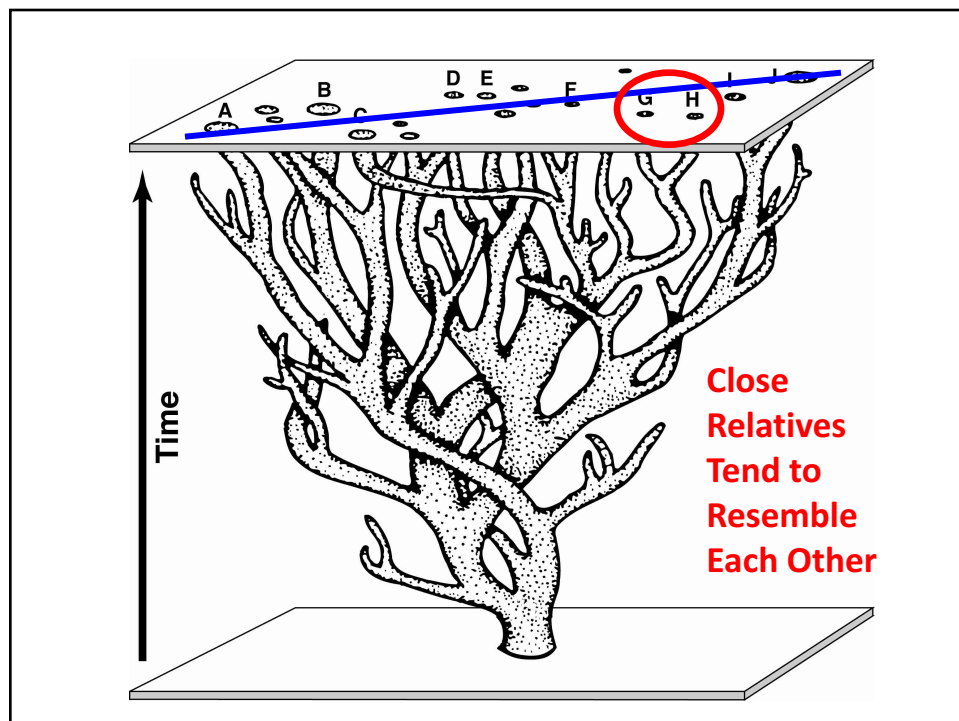
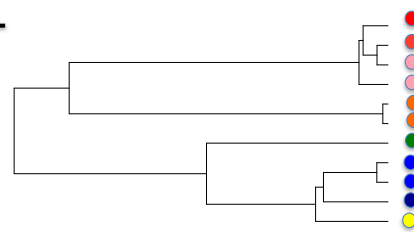
Problem with Traditional Correlation Tests

Central Assumptions in Statistics



-Normally distributed data (parametric statistics)

-Independent observations
- Points are independent



Solution

- Because species share common ancestors they cannot be considered independent observations in a correlation study.
- Therefore a central assumption of correlation statistics is violated
- Is there a solution?

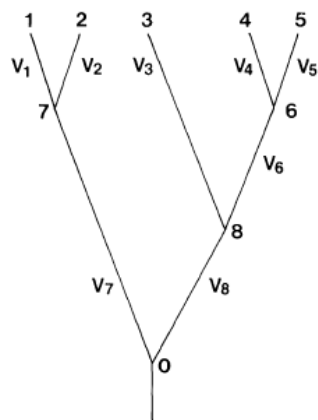
Many Solutions

- Originally using paired species comparisons
- Methods with an underlying model of trait evolution
 - Phylogenetically independent contrasts (PIC)
 - Phylogenetic generalized least squares (pGLS)
- Methods with no underlying model
 - Phylogenetic eigenvectors

Phylogenetically Independent Contrasts

Felsenstein's Solution Step 3

COMPARATIVE METHOD

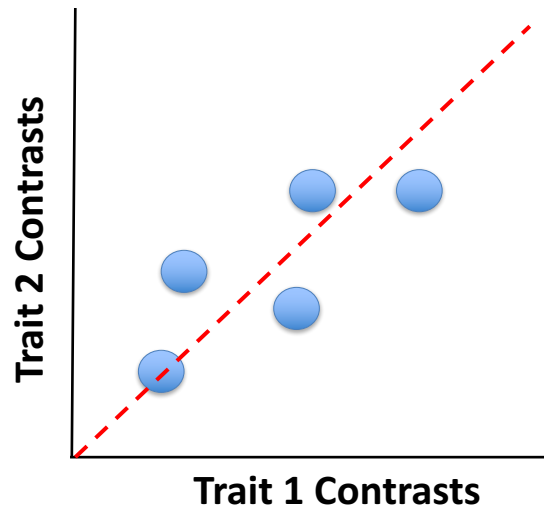


CONTRAST

$X_1 - X_2$
 $X_4 - X_5$
 $X_3 - X_6$
 $X_7 - X_8$

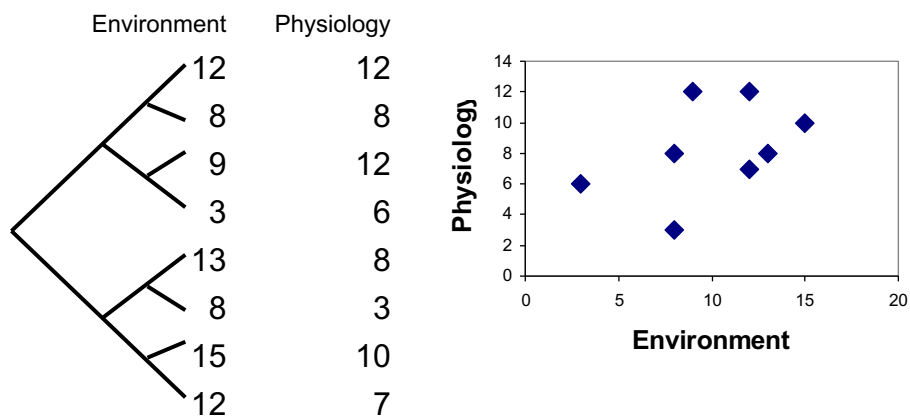
**Calculate
Regression
Through Origin
Between Contrasts**

Felsenstein's Solution Step 3

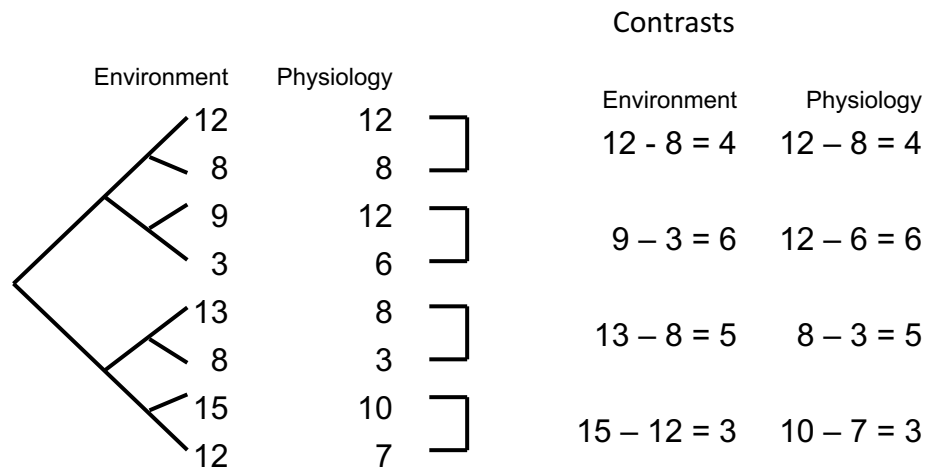


Example of analysis that takes phylogeny into account

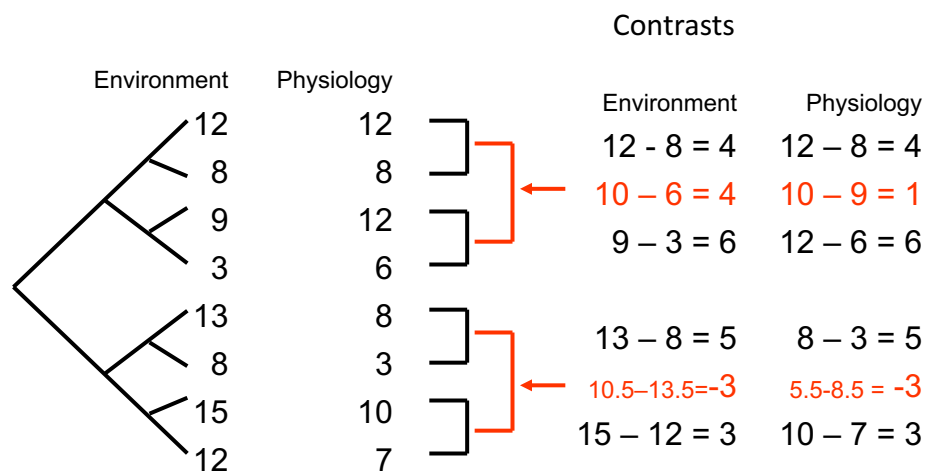
- Phylogenetically independent contrasts



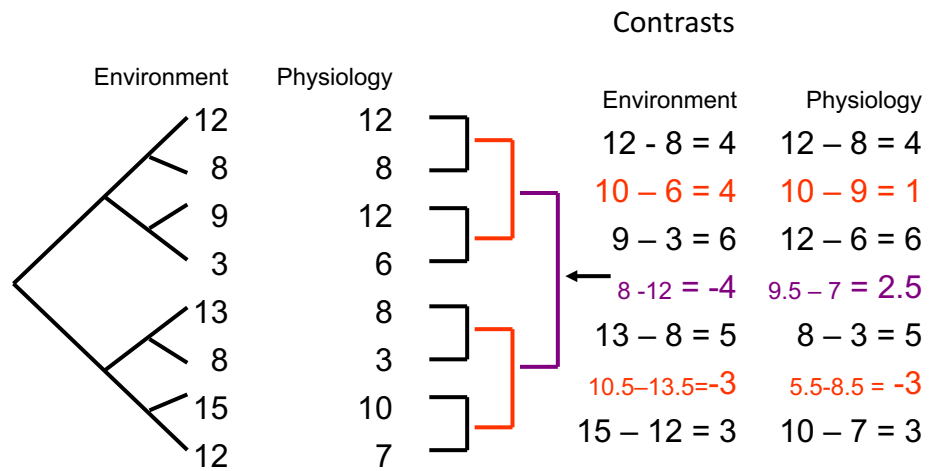
Phylogenetically independent contrasts



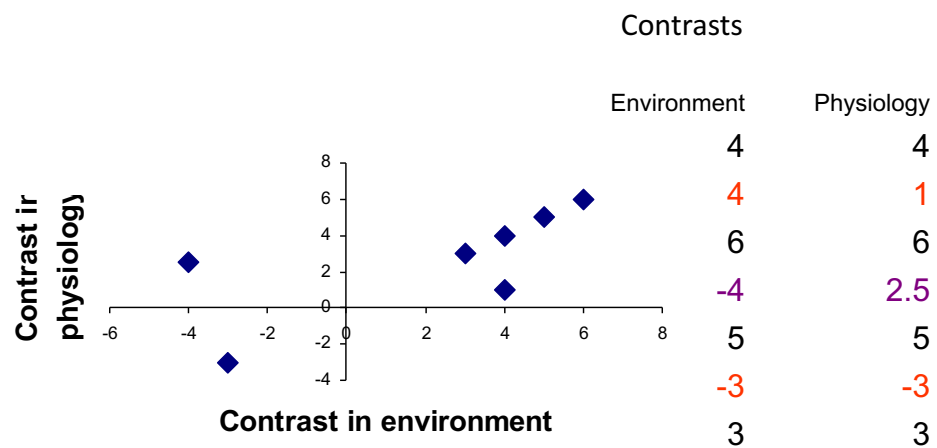
Phylogenetically independent contrasts



Phylogenetically independent contrasts



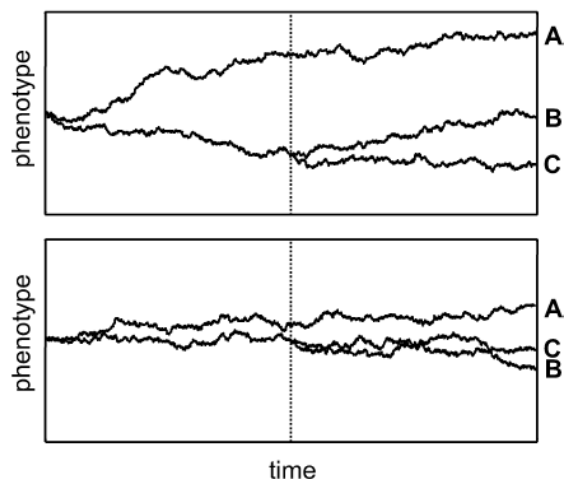
Phylogenetically independent contrasts



PIC ASSUMPTIONS

- The phylogenetic topology is 'correct'. IE we know the correct phylogeny
 - Phylogenies as hypotheses
- Branch lengths are measured in the units expected to correlate with character evolution
- We know the species trait values
- Traits evolve under a Brownian motion model
 - Increasingly, people are developing methods to rely on other potential models or that fit models given the data

Brownian Motion

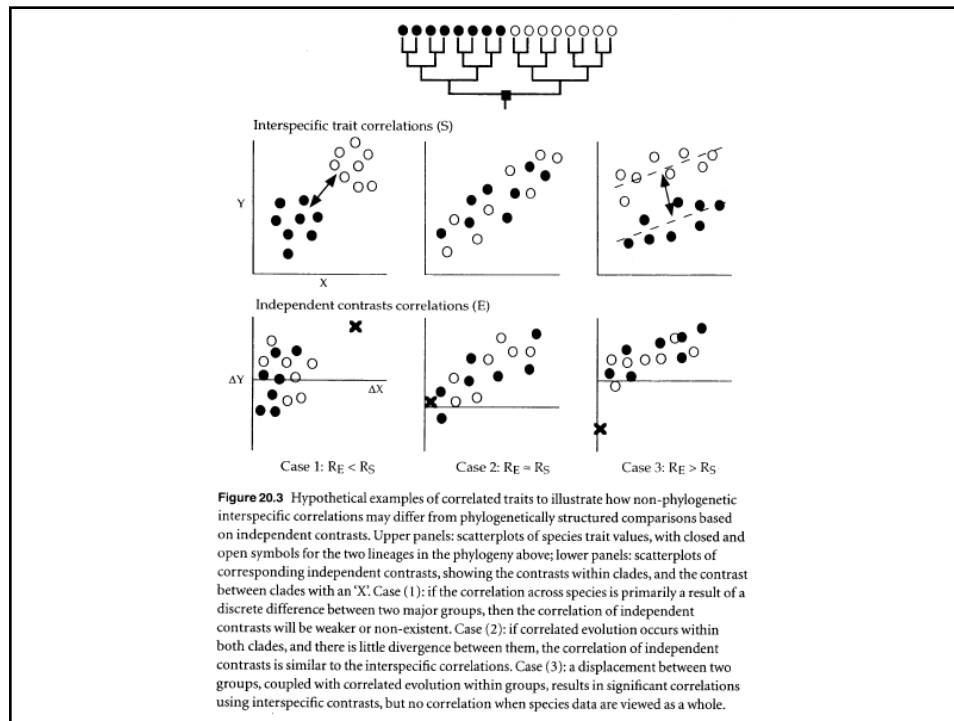


Limits to phylogenetically based analyses

- What if a particular adaptation only occurred once?
 - Don't want to ignore a "neat story" that is supported by other kinds of evidence
- Trait relationships may be non-linear
 - Must transform data to satisfy BM assumption
 - Reduces statistical power

How might this influence my statistical inference?

- This depends on the degree of phylogenetic signal in data
- If there is no phylogenetic signal in trait data, then your result is similar compared to traditional correlation
- If there is high phylogenetic signal in trait data, then your inference can be very different compared to traditional correlation



What kind of data can (or should) be analyze using Phylogenetically Independent Contrasts?

- Where species are units of observation
 - Though, a handful of methods are now available that utilize variation or probability distributions w/i species
- Comparative Questions
 - Do species in drier soils have smaller leaves (data are soil moisture preference and leaf size for many species)
 - Is maximum height of species correlated with their seed size? (data are maximum heights and seed size of many species)

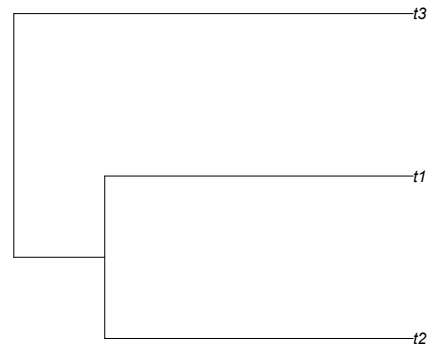
Phylogenetic GLS

Phylogenetic GLS

- In the simplest case equivalent to independent contrast analysis (Grafen 1989; Martins & Hansen 1997) but various extensions.
- One important extension is fitting a model of trait evolution to the residuals

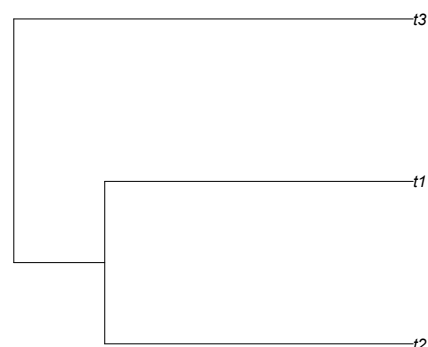
Phylogenetic GLS

- We start with a phylogeny
- Must convert it to a phylogenetic variance-covariance (VCV) matrix



What is a phylogenetic VCV?

- A phylo VCV is constructed assuming a model of trait evolution.
 - Usually Brownian motion
 - Trait variance and co-variance scales with branch length

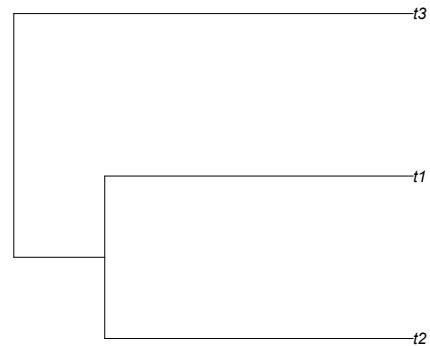


What is a phylogenetic VCV?

```

      t2      t1      t3
t2 0.6174106 0.1405757 0.0000000
t1 0.1405757 0.6174106 0.0000000
t3 0.0000000 0.0000000 0.6174106

```



What is a phylogenetic VCV?

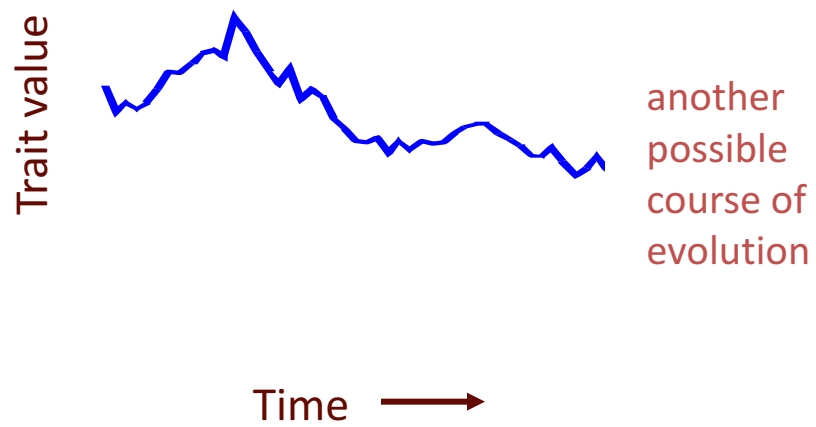
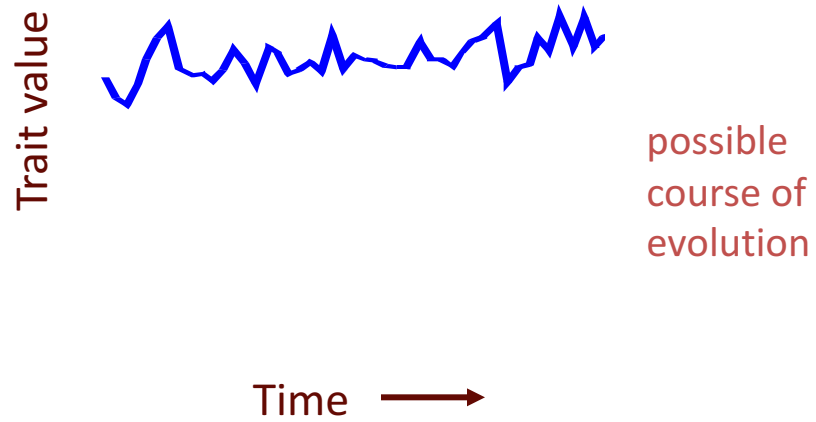
```

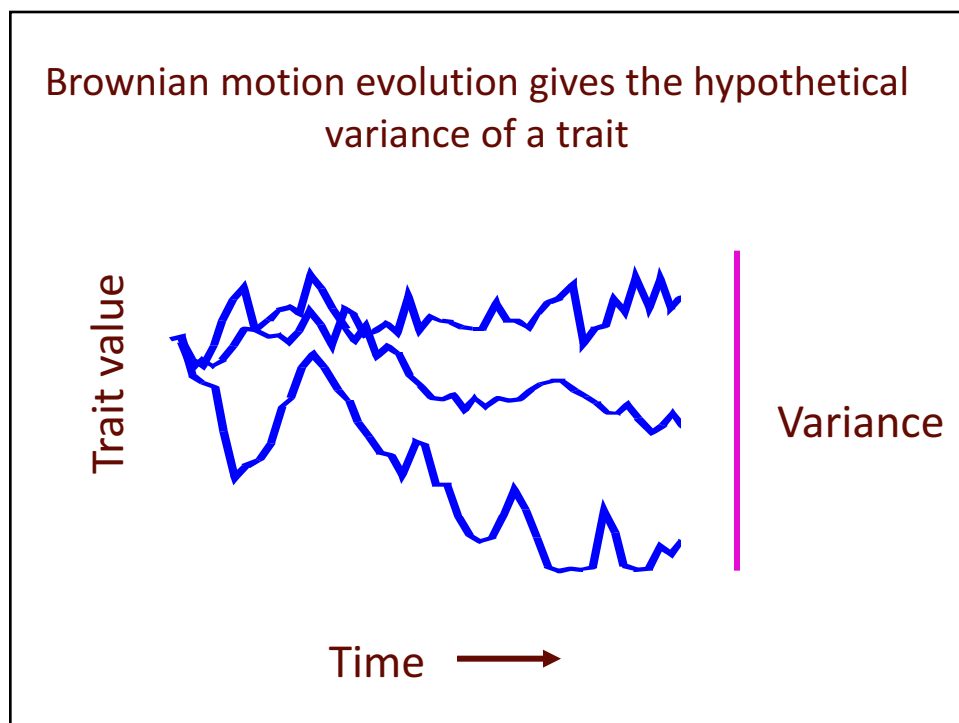
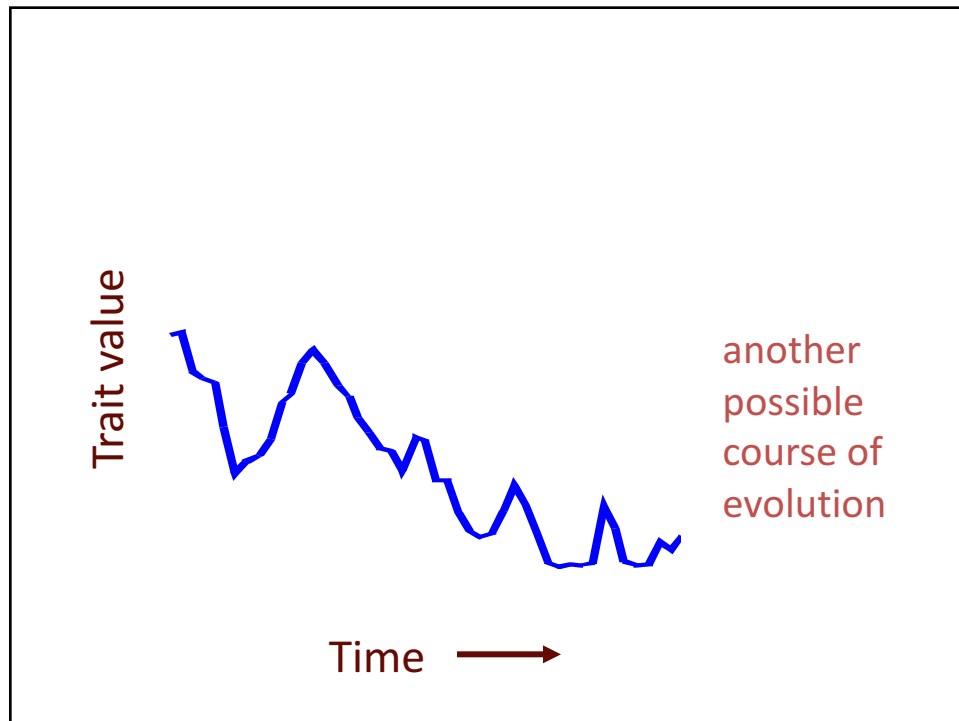
      t3      t1      t2
t3 0.5795273 0.5625709 0.0000000
t1 0.5625709 0.5795273 0.0000000
t2 0.0000000 0.0000000 0.5795273

```

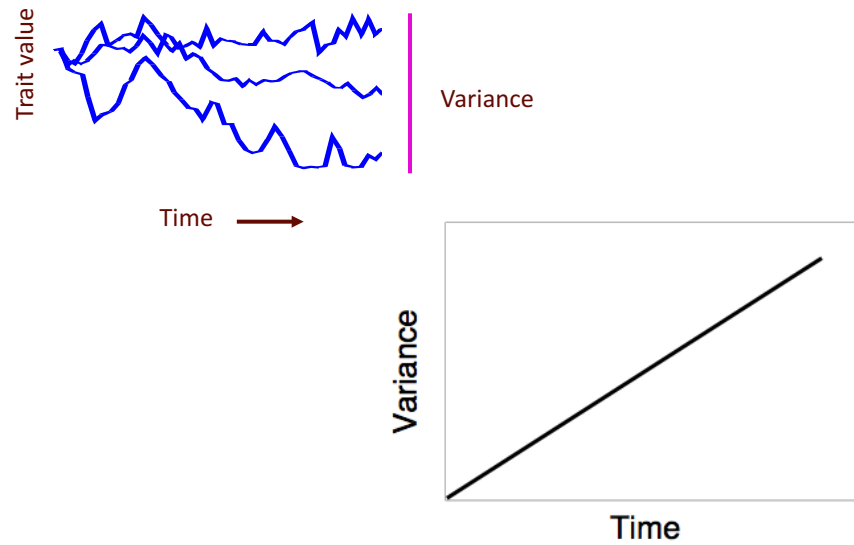


Hypothetical trait for a single species under
Brownian motion evolution

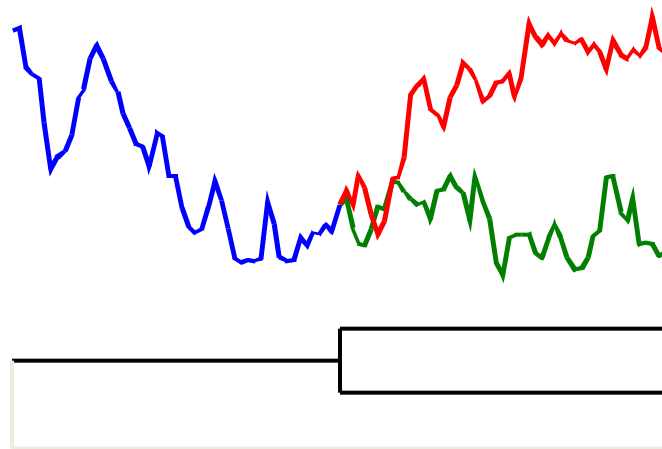


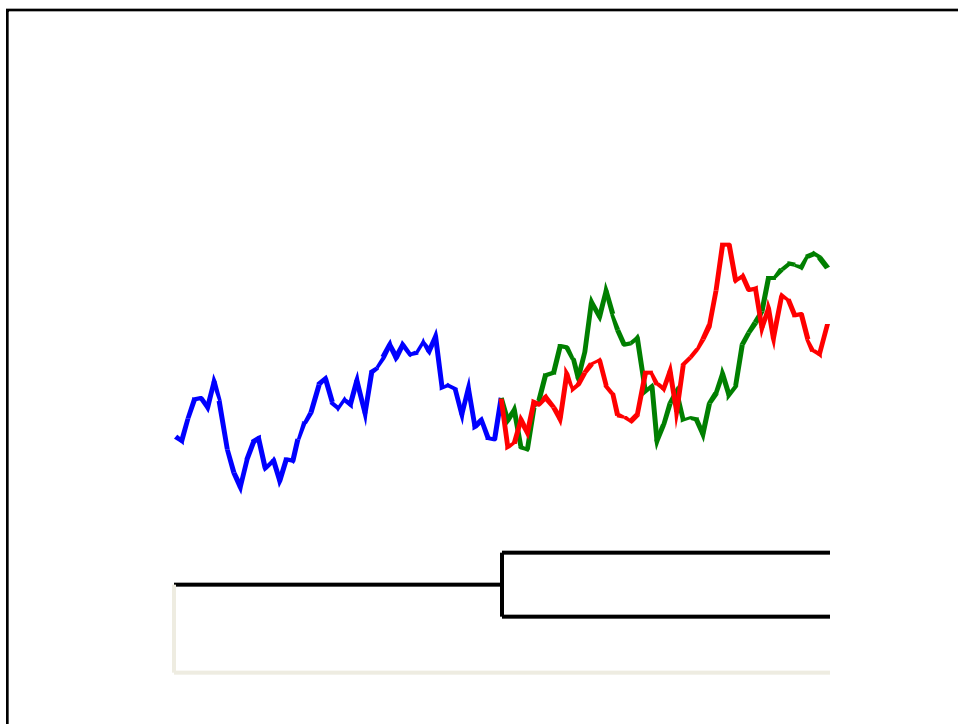
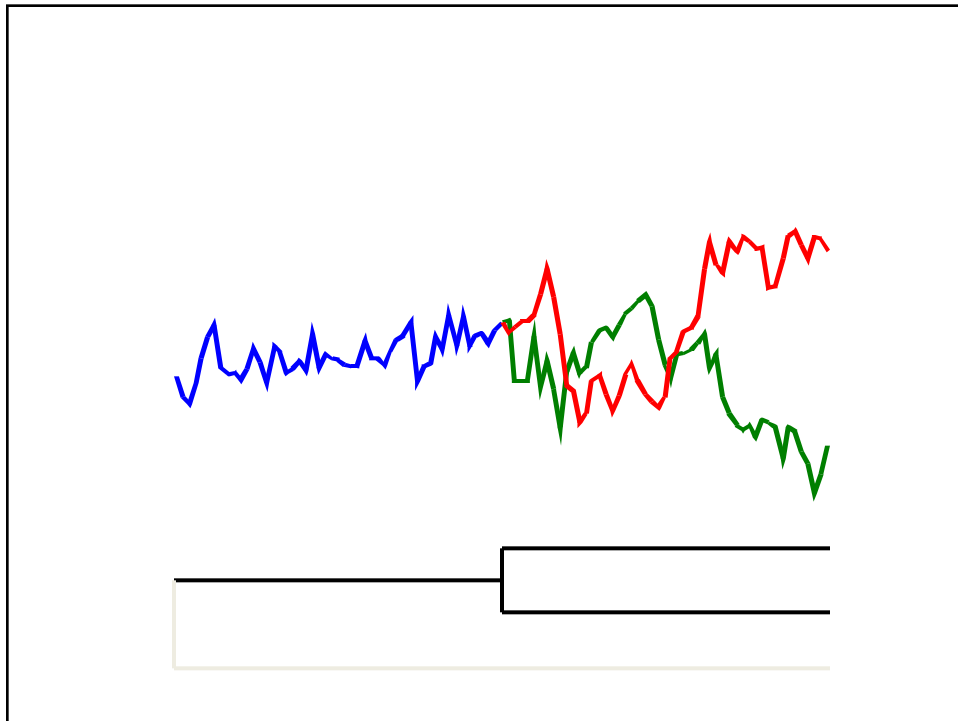


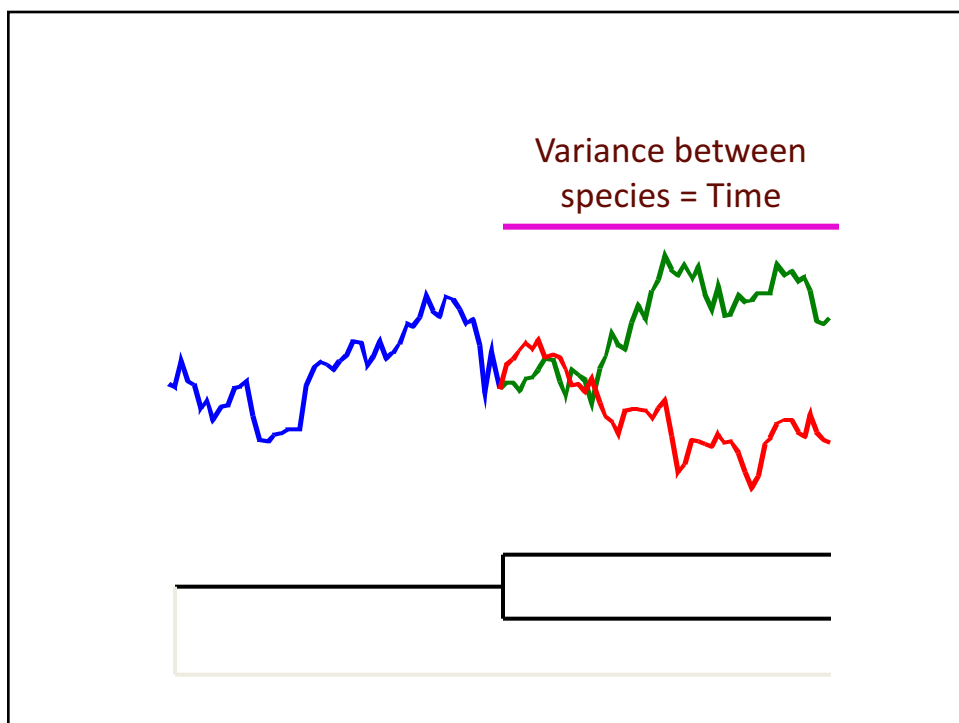
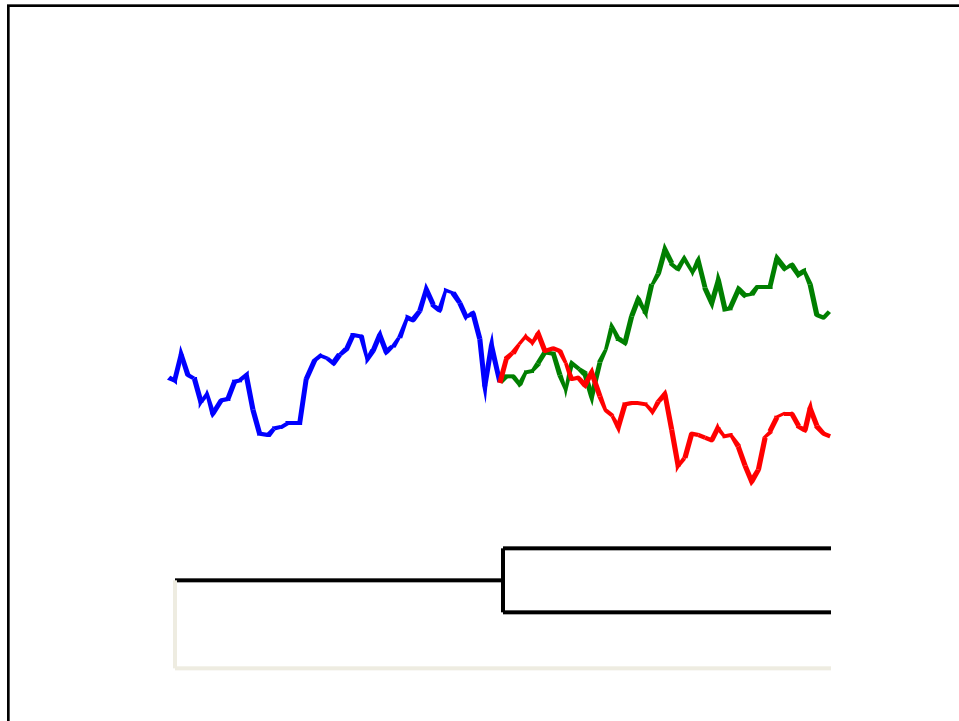
Brownian motion evolution

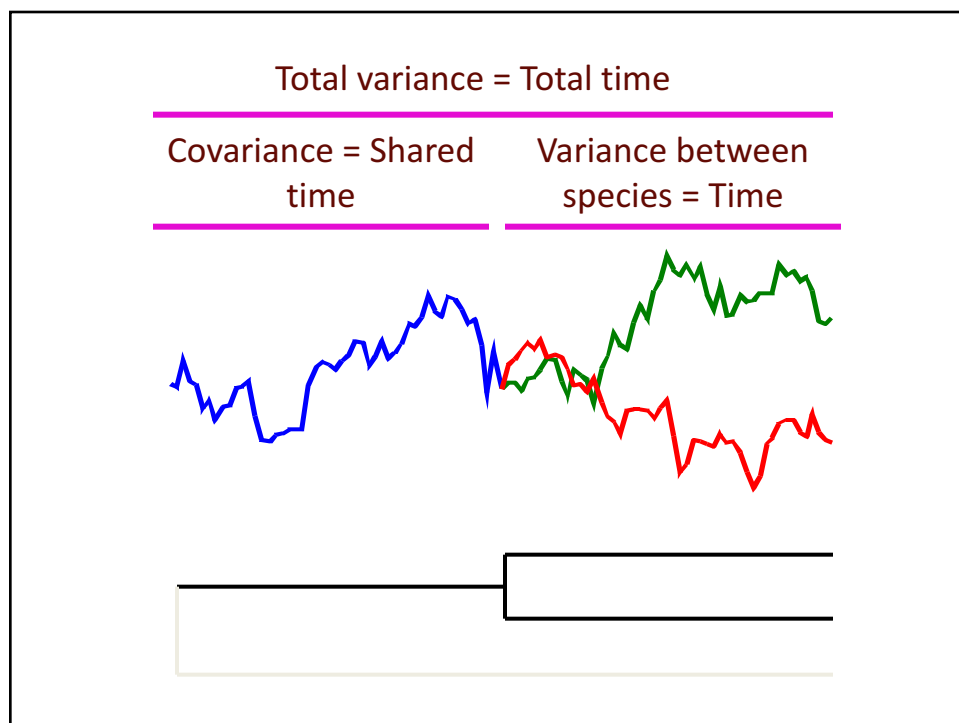
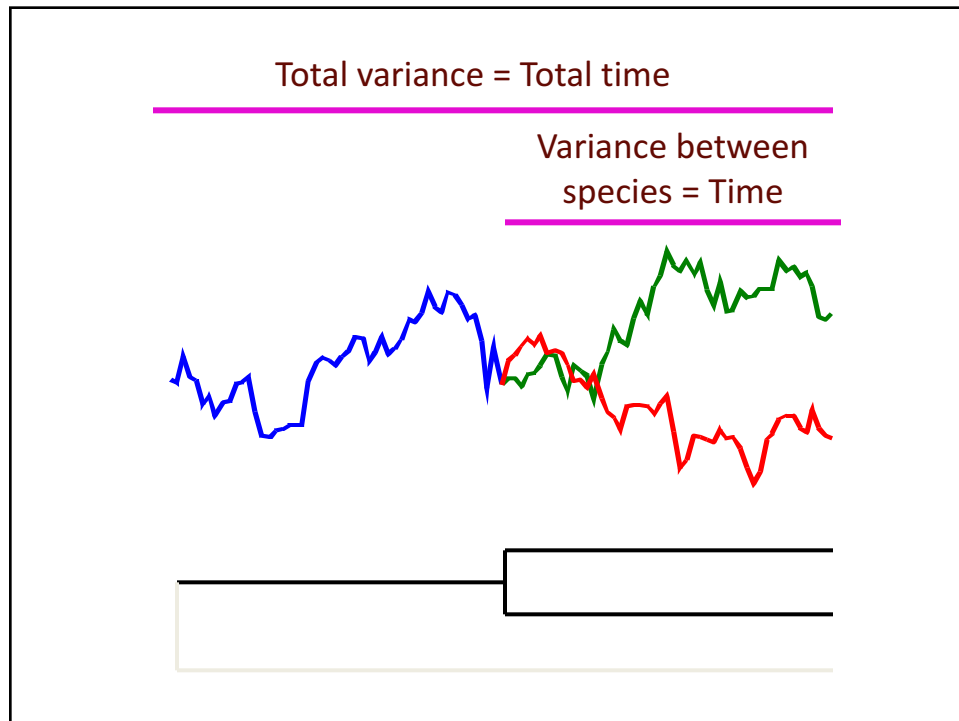


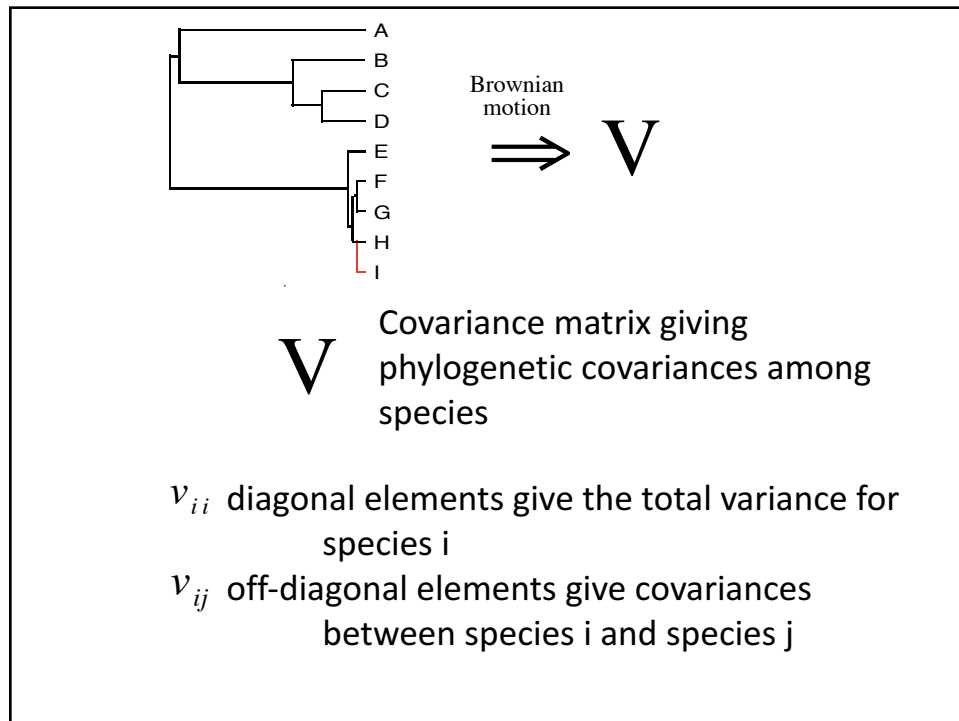
Brownian motion evolution of a hypothetical trait during speciation





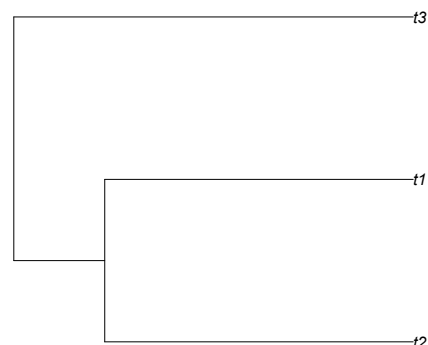






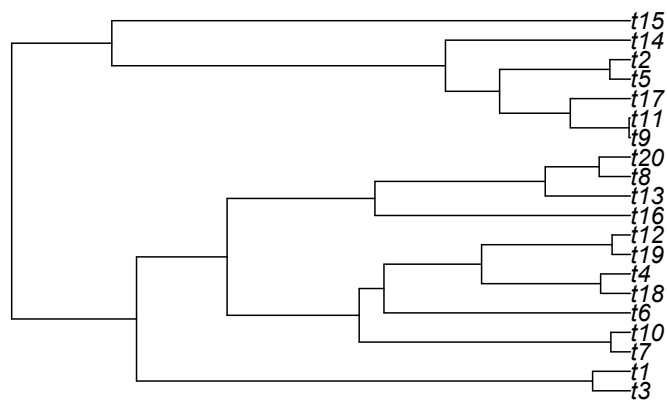
Phylo GLS

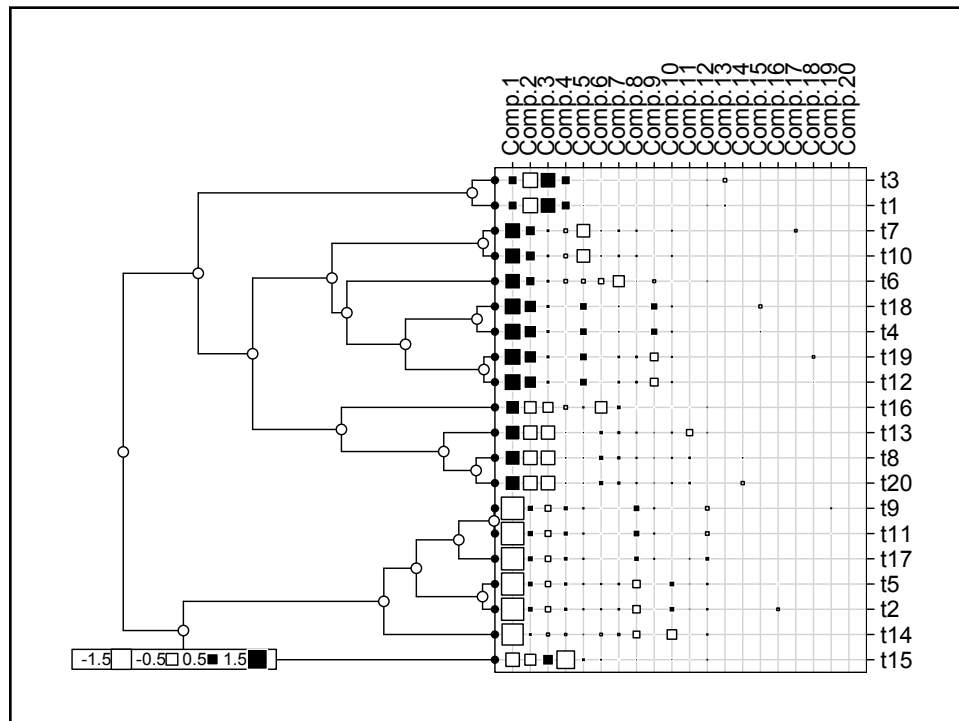
- The phylo VCV is then used in a generalized linear model
 - Can use the `r` command `pGLS()`



Phylogenetic Eigenvectors

- Inputs a phylogenetic distance matrix into a PCA to decompose the phylo into many eigenvectors





Phylogenetic Eigenvectors

- Phylogenetic eigenvectors are then used in a linear model similar
 - E.g. as independent variables in the model
- Some well discussed weaknesses
 - No explicit model of trait evolution
 - How do we pick which eigenvectors to use?

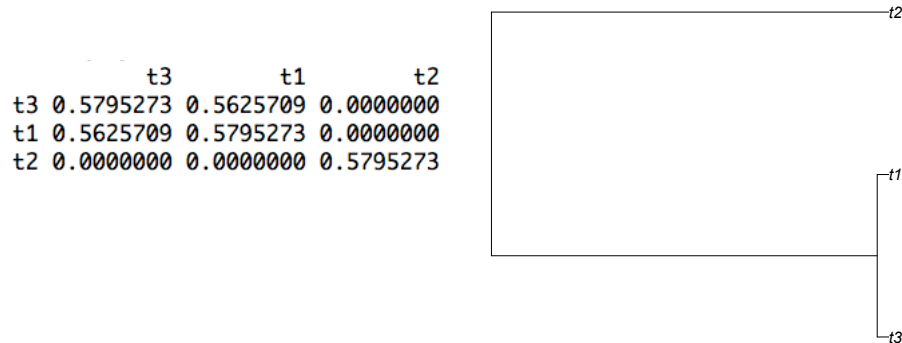
Phylogenetic ANOVAs/t-tests/PCAs/etc

- Phylogenetic ANOVAs and MANOVAS can be conducted (or t-tests for that matter)
- A standard way is to calculate a regular (M)ANOVA
- Next trait values are simulated on the phylogeny under a given model of trait evolution
 - This generates a null distribution

Intra-specific variation

- Error/variance
- Re-sampling from probability distribution

Our old friend the vcv



Ives et al. 2007 sys biol

ANALYSES

Univariate Analyses and Phylogenetic Signal

The problem of finding the best estimator of the expectation of a random variable when there is phylogenetic correlation and measurement error is given by the statistical model

$$\mathbf{X}^* = \mathbf{a} + \boldsymbol{\varepsilon} \quad (1)$$

$$\mathbf{X} = \mathbf{X}^* + \boldsymbol{\eta}$$

where \mathbf{X}^* is a $N \times 1$ dimensional vector containing the true values of a trait in a sample of N species (tips), \mathbf{a} is a scalar giving the expected value of the trait, $\boldsymbol{\varepsilon}$ is a $N \times 1$ vector of zero-mean error terms depicting the evolutionary variance of the trait among species, \mathbf{X} is a $N \times 1$ vector containing the observed values of the trait, and $\boldsymbol{\eta}$ is the $N \times 1$ vector of errors associated with measurement. Note that for notational convenience we have written $\mathbf{X}^* = \mathbf{a} + \boldsymbol{\varepsilon}$ as the sum of a scalar and a vector to represent $\mathbf{X}^* = \mathbf{a}\mathbf{1} + \boldsymbol{\varepsilon}$ where $\mathbf{1}$ is the $N \times 1$ vector of ones.

Because closely related species will likely have similar values of trait \mathbf{x} , values of $\boldsymbol{\varepsilon}$ will be correlated among species. Thus, we assume the covariance matrix for $\boldsymbol{\varepsilon}$ is given by $\text{E}\{\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}'\} = \sigma^2\mathbf{C}$, where σ^2 scales the overall phylogenetically inherited variance (sometimes referred to as the rate of evolution; Garland et al., 1999; Garland and Ives, 2000), and \mathbf{C} gives the correlation structure created by phylogenetic relatedness. The most common assumption in phylogenetic analyses is that evolution proceeds like a "Brownian motion" process; through time, the value of a trait changes in small increments in random directions, like a random walk in continuous time (Felsenstein, 1985). Under this assumption, $\boldsymbol{\varepsilon}$ has a multivariate normal distribution in which the element c_{ij} of \mathbf{C} is proportional to the length of the shared

branch from root to the last common ancestor of

Although we do not consider correlated measurement errors in detail, nonzero off-diagonal elements of \mathbf{M} can be used in all of the methods we derive. Finally, although we will typically assume that $\boldsymbol{\varepsilon}$ and $\boldsymbol{\eta}$ have multivariate normal distributions, for some of the statistical procedures described below, $\boldsymbol{\varepsilon}$ and $\boldsymbol{\eta}$ need not be restricted to being normally distributed.

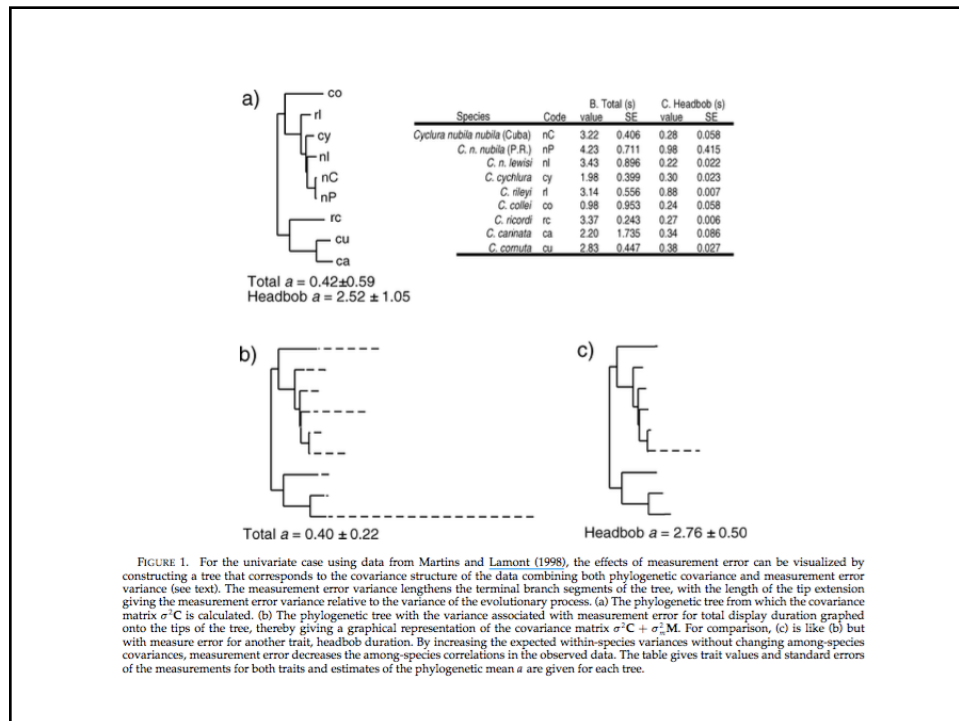
Consider first the case of no measurement error. Equation 1 can be reformulated as a phylogenetic regression problem in which the error terms are correlated, and hence can be analyzed using either independent contrasts or, as we will do here, generalized least squares (Hansen and Martins, 1996; Garland and Ives, 2000; Rohlf, 2001). Because \mathbf{C} is a covariance matrix (and hence real, symmetric, and nonsingular), there exists another matrix \mathbf{D} such that $\mathbf{D}\mathbf{C}\mathbf{D}' = \mathbf{I}$, where the apostrophe denotes transpose and \mathbf{I} is the $N \times N$ identity matrix. Matrix \mathbf{D} can be used to transform values of trait \mathbf{x} by letting $\mathbf{Z} = \mathbf{D}\mathbf{X}$, $\mathbf{U} = \mathbf{D}\mathbf{1}$ (the $N \times 1$ vector of 1's), and $\boldsymbol{\alpha} = \mathbf{D}\mathbf{a}$. From Equation 1 (with $\boldsymbol{\eta} = 0$), this gives

$$\mathbf{Z} = \mathbf{U}\boldsymbol{\alpha} + \boldsymbol{\alpha} \quad (2)$$

The covariance matrix of $\boldsymbol{\alpha}$ is $\text{E}\{\boldsymbol{\alpha}\boldsymbol{\alpha}'\} = \text{E}\{\mathbf{D}\boldsymbol{\varepsilon}(\mathbf{D}\boldsymbol{\varepsilon})'\} = \text{E}\{\mathbf{D}\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}'\mathbf{D}'\} = \text{E}\{\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}'\}\mathbf{D}' = \mathbf{D}(\sigma^2\mathbf{C})\mathbf{D}' = \sigma^2\mathbf{I}$. Thus, no covariance terms appear in the covariance matrix of $\boldsymbol{\alpha}$, so the error terms $\boldsymbol{\alpha}$ are uncorrelated. Equation 2 can, therefore, be analyzed as a standard least-squares regression problem with independent errors. Specifically, the generalized least-squares (GLS) estimator of $\boldsymbol{\alpha}$ is

$$\hat{\boldsymbol{\alpha}} = \frac{\mathbf{U}'\mathbf{Z}}{\mathbf{U}'\mathbf{U}} = \frac{(\mathbf{D}\mathbf{1})'(\mathbf{D}\mathbf{X})}{(\mathbf{D}\mathbf{1})'(\mathbf{D}\mathbf{1})} = \frac{\mathbf{1}'\mathbf{D}'\mathbf{D}\mathbf{X}}{\mathbf{1}'\mathbf{D}'\mathbf{D}\mathbf{1}} = \frac{\mathbf{1}'\mathbf{C}^{-1}\mathbf{X}}{\mathbf{1}'\mathbf{C}^{-1}\mathbf{1}} \quad (3)$$

The corresponding estimate of σ^2 is the mean squared error



Evans et al. 2009 AmNat

VOL. 173, NO. 2 THE AMERICAN NATURALIST FEBRUARY 2009

®

Climate, Niche Evolution, and Diversification of the “Bird-Cage” Evening Primroses (*Oenothera*, Sections *Anogra* and *Kleinia*)

Margaret E. K. Evans,^{1,2,*} Stephen A. Smith,^{1,†} Rachel S. Flynn,^{1,‡} and Michael J. Donoghue^{1,3,§}

Evans et al. 2009 AmNat

- Obtain a probability distribution for a trait
- Re-sample that distribution and estimate ancestral states each time

