# Phylogenetic Comparative Methods for Multivariate Data

# The (Incomplete) Road to Comparative Methods
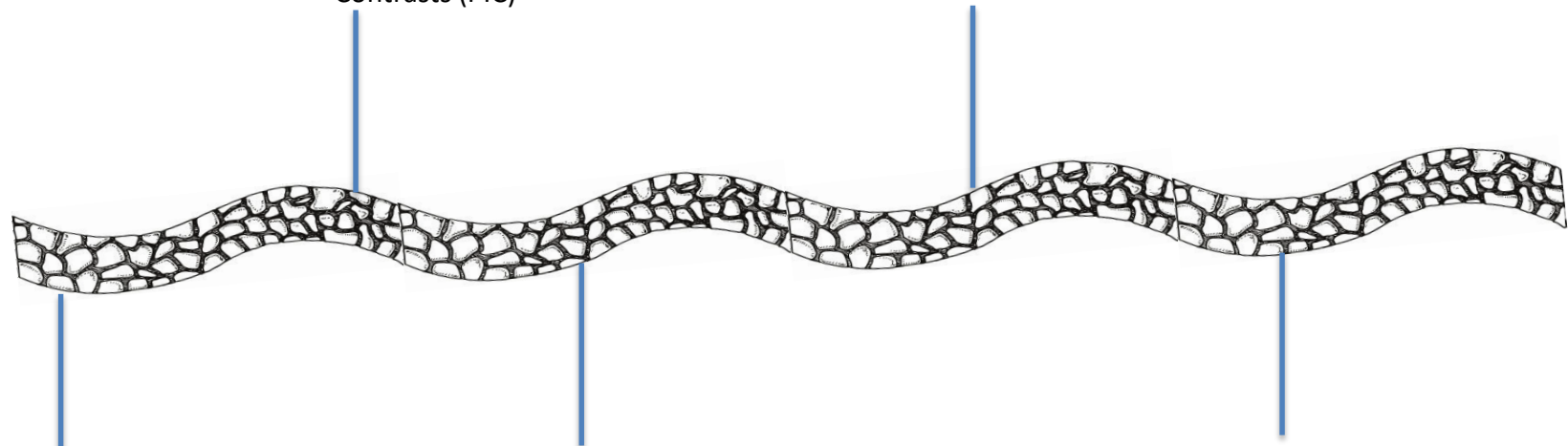
**2000s: Maturation phase**
Synthesis: PIC, PGLS, Phylo-transform
Complex model comparison (BM1, BMM, OU1, OMM)
Bayesian methods
Parameter-shift methods (e.g., MEDUSA, BAMM)
Discrete diversification associations (BiSSe family)

**1985: The Breakthrough**
Phylogenetic Independent
Contrasts (PIC)

**70s – early 80s: early attempts**
Nested ANOVA
Phylogenetic autocorrelation

**80s – 90s: 'niche expansion'**
PGLS
Phylogenetic signal ($\lambda$, $K$)
Phylogenetic ANOVA
Evolutionary models (BM1, OU1, ACDC, $\lambda$)
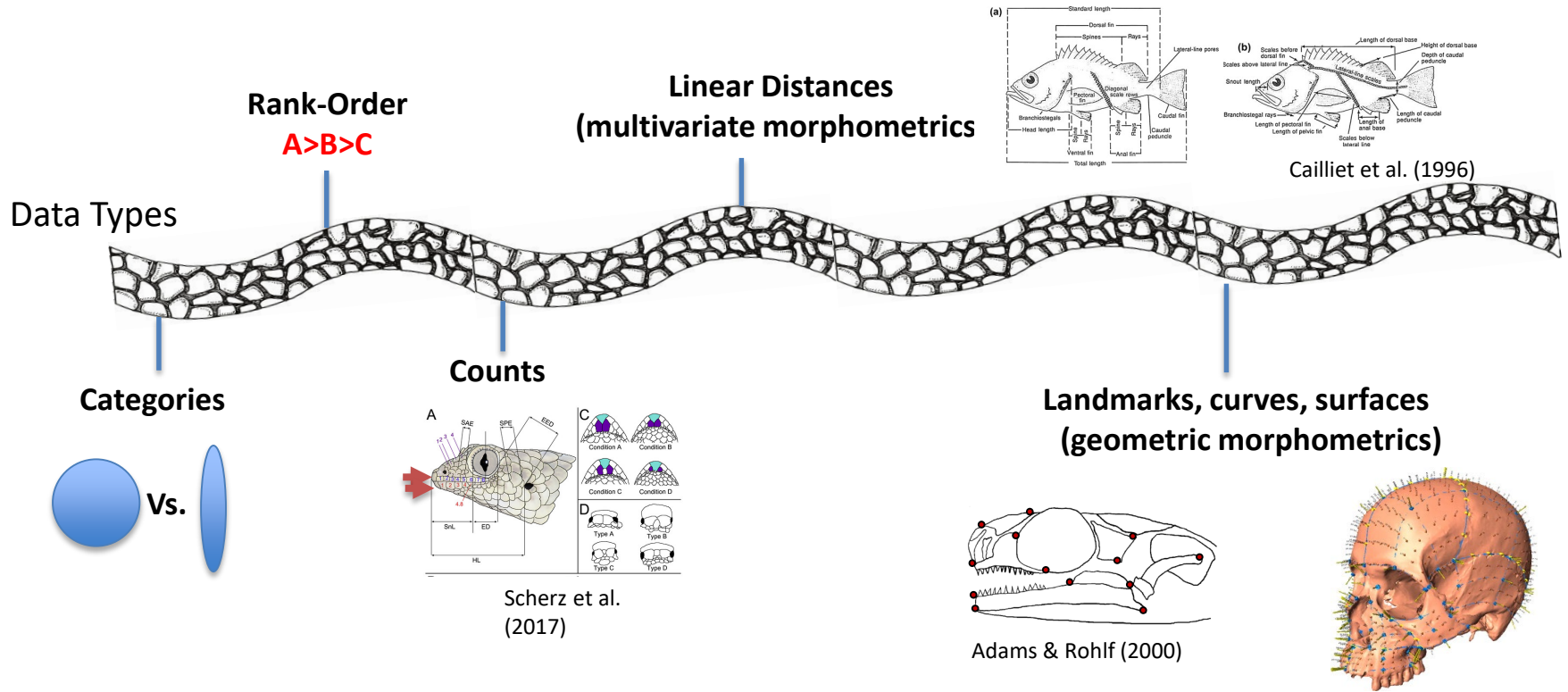Diversity plots (LTT & DTT)
Diversification rates
Discrete trait change models

**~2010s: Multivariate +GMM**

Present day: PCMs: A diverse toolkit for evaluating evolutionary hypotheses

## Morphological quantification has advanced dramatically*

**Data Types**

**Rank-Order**
**A>B>C**

**Linear Distances**
**(multivariate morphometrics**

Cailliet et al. (1996)

**Categories**

**Vs.**

**Counts**

Scherz et al. (2017)

**Landmarks, curves, surfaces (geometric morphometrics)**
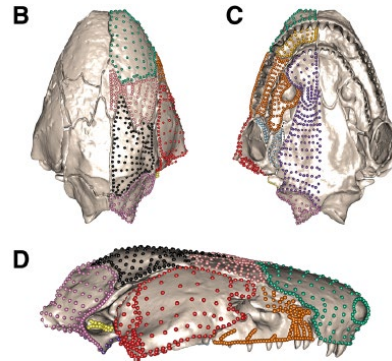
Adams & Rohlf (2000)

Friess 2010

GMM provides greater biological realism, but…
   -greater data complexity
   -requires new mathematical theory
   -analytical and statistical challenges

*See historical treatments in: Reyment, 1996; Bookstein 1998, Adams et al. 2013; Bookstein 2014, 2018, 2019,  among others

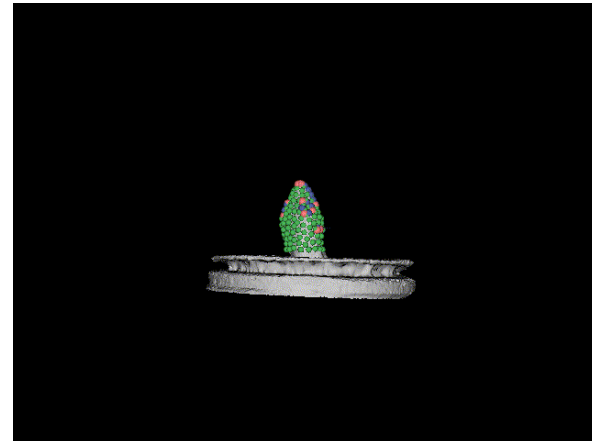## GMM (+ new technology) leads to ever-complex & HD datasets



1469 landmarks =
4407 variables

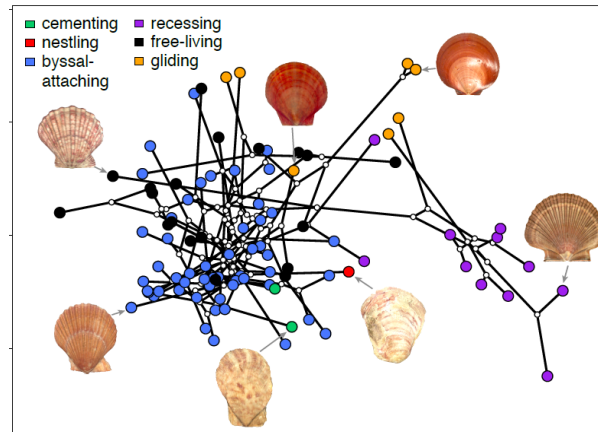Bardua et al. 2019

From this ...                    Obtain this





**How do we handle such phenotypes with statistical rigor?**

PCMs *condition* the data on the phylogeny during the analysis

Empirical Goal: Evaluate evolutionary hypotheses while accounting for (phylogenetic) non-independence



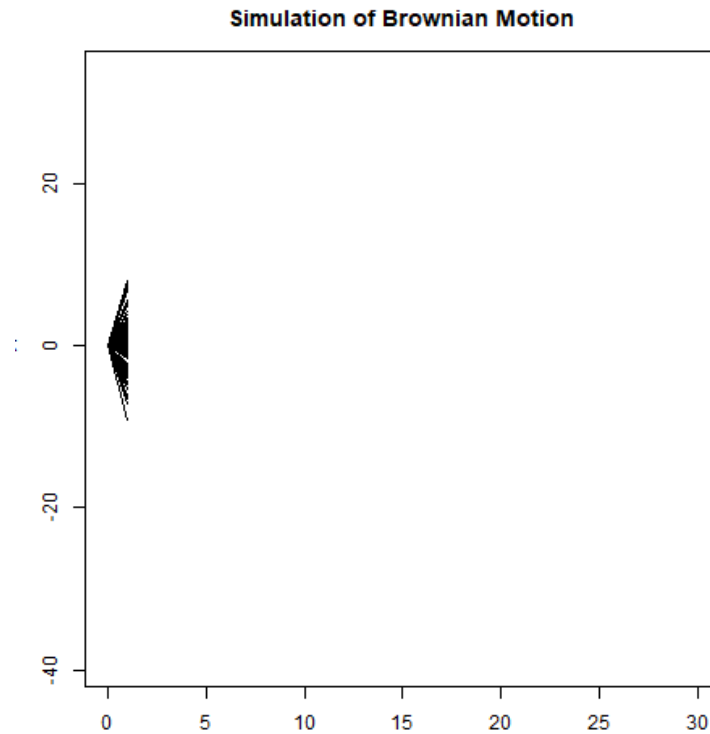Sherratt, Alejandrino, Kraemer, Serb, & Adams (2016)

Requires an evolutionary model of how trait variation is expected to accumulate

## Brownian motion (BM): a *null model* of trait change

Trait changes are independent from time step to time step

Results in: $\Delta\mu=0$, but $\sigma^2 \uparrow \propto$ time
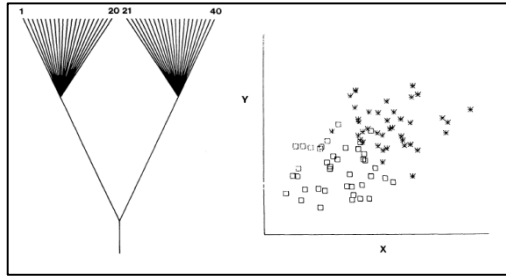
**Simulation of Brownian Motion**



Side-note: this is the continuous-trait model equivalent of the Markov process, and is intimately related to Gaussian theory and the normal distribution
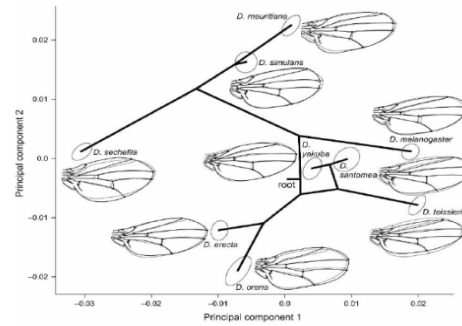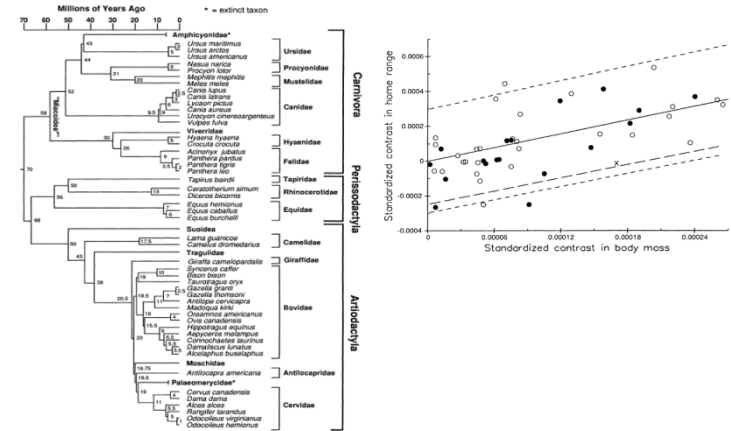
Felsenstein (1973; 1985)
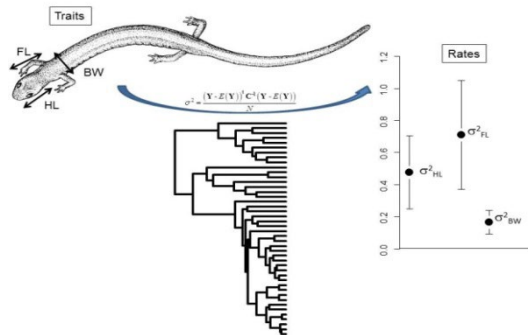
## Phylogenetic Signal
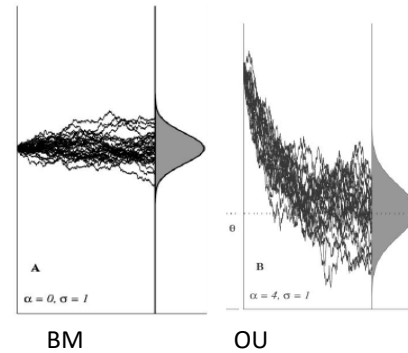


## Phylomorphospace



## Phylogenetic Regression (PIC & PGLS)
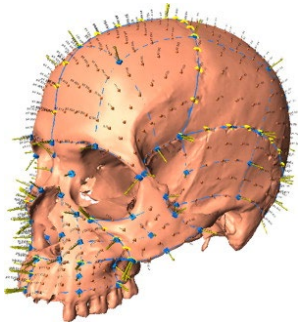


## Evolutionary Rates



## Evolutionary Models



BM        OU

# All are derived from the general PCM model (PGLS)

# The primary statistical model of PCM: GLS (generalized least squares)

$$\mathbf{Y} = \mathbf{X}\widehat{\boldsymbol{\beta}} + \mathbf{E}$$
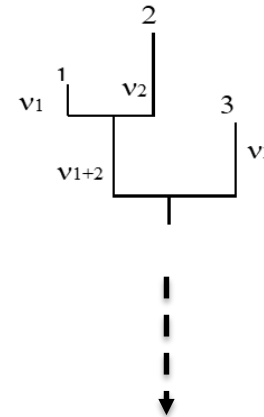
Continuous
Response Data

The Design

Error: $\mathcal{N}(0, \mathbf{V})$
(as described by phylogeny)



Friess 2010

$$\begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 1 & 1 \\ 1 & 1 \\ 1 & 1 \\ 1 & 1 \end{bmatrix}$$



Island vs. Mainland

shape

region

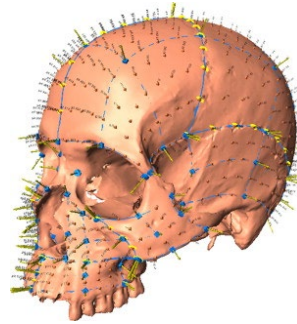$$\mathbf{V} = \mathbf{R} \otimes \begin{pmatrix} v_1 + v_{1+2} & v_{1+2} & 0 \\ v_{1+2} & v_2 + v_{1+2} & 0 \\ 0 & 0 & v_3 \end{pmatrix}$$

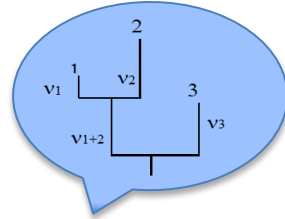*$\mathbf{V}$ can have other formulations for alternative evolutionary models

Shape ~ Region | phylogeny



Friess 2010

$$= \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 1 & 1 \\ 1 & 1 \\ 1 & 1 \\ 1 & 1 \end{bmatrix} * \beta + \varepsilon$$

This GMM/PCM approach requires that one:

    1: Condition multivariate data on phylogeny & fit model parameters

    2: Obtain robust summary statistics

    3: Evaluate significance and effect sizes in reliable manner

**These were rather significant analytical challenges to overcome!**
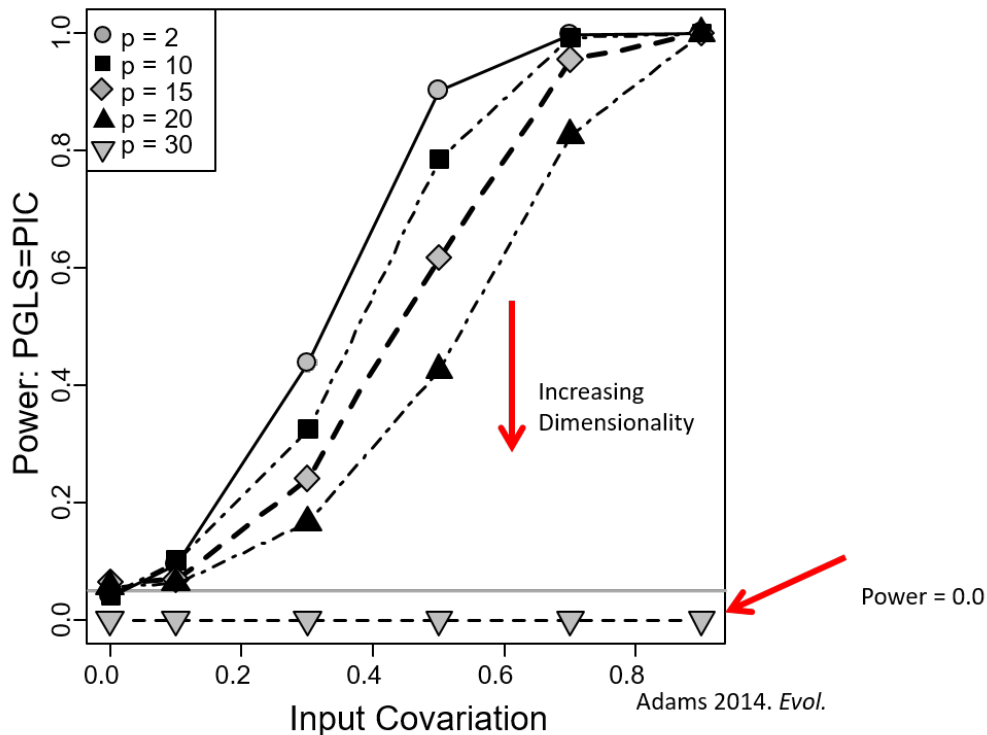
Why not just 'scale up' standard PCMs for GM-data?

Example: phylogenetic regression

$$\mathbf{Y} = \mathbf{X}\widehat{\boldsymbol{\beta}} + \mathbf{E} \qquad \mathbf{E} \sim \mathcal{N}(0, \mathbf{V})$$

Power decreases as $p$-dimensions increases



Adams 2014. *Evol.*

**Why does this happen?**

Adams (2014)
Adams and Collyer (2018)

# The Curse of Parametric Hypothesis Testing

Standard PCMs are rooted in likelihood-based statistical theory

$$\log L = \log \left[ \frac{\exp\left(-0.5(\mathbf{Y} - E(\mathbf{Y}))^t \mathbf{V}^{-1}(\mathbf{Y} - E(\mathbf{Y}))\right)}{\sqrt{(2\pi)^n |\mathbf{V}|}} \right]$$

*Lots of math here!*

$$\log L = \frac{\cdots}{\sqrt{\ldots + |\mathbf{V}|}}$$   The problem?   As  $|\mathbf{V}| \to \mathbf{0}$ as $\boldsymbol{p} \to \boldsymbol{n}$

Translation: divide by zero!

**We need another solution for highly multivariate data!**

See: Adams (2014), (2015)
Adams and Collyer (2018a), (2018b); Adams and Collyer (2019)

Forgo standard ML and parametric approaches for statistical evaluation, and use robust methods
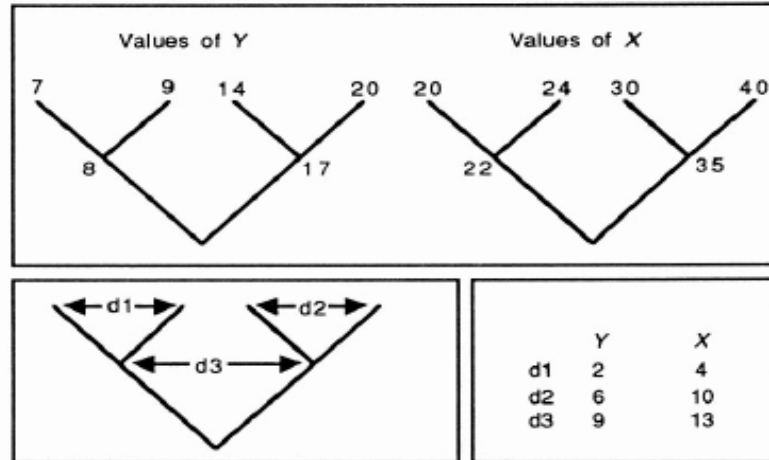
New GMM/PCM Approach:

    1: Condition data on phylogeny & fit model parameters

    2: Obtain robust summary measures (avoid $|\mathbf{V}| = 0$)

    3: Evaluate significance and effect sizes *NOT* using log$L$

See: Adams (2014a), (2014b) (2014c)
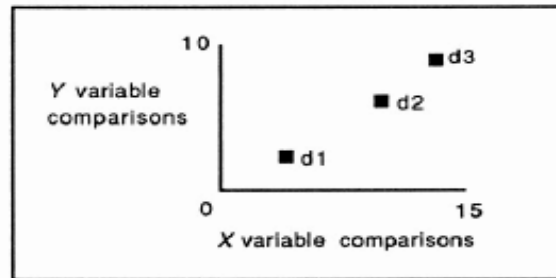Adams and Collyer (2015), (2018a), (2018b); Adams and Collyer (2019)

Three equivalent algebraic implementations

## 1: Phylogenetically Independent Contrasts

Calculate PICs

Analysis of PICs

From Harvey & Pagel (1991)

$$\widehat{\boldsymbol{\beta}} = \left(\mathbf{X}_{pic}^{t}\mathbf{X}_{pic}\right)^{-1}\mathbf{X}_{pic}^{t}\mathbf{Y}_{pic}$$

Felsenstein (1985)

See: Garland and Ives (2000)
Rohlf (2006); Blomberg et al (2012); Adams (2014a)
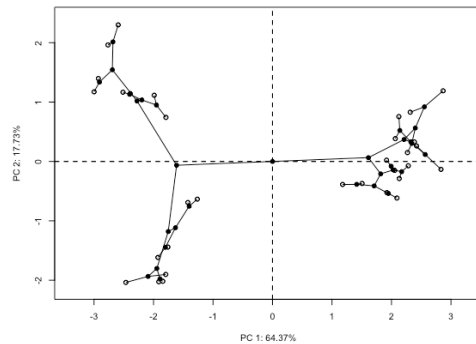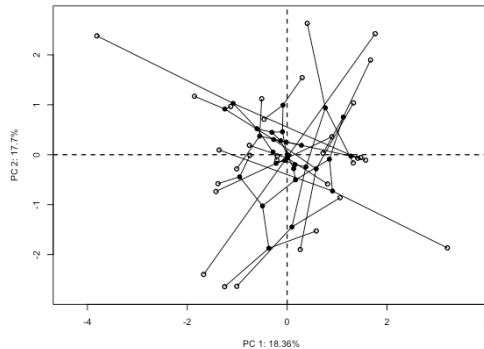
Three equivalent algebraic implementations

2: Phylogenetic (GLS) Regression ($\mathbf{Y} = \mathbf{X}\widehat{\boldsymbol{\beta}} \mid \boldsymbol{phy}$)

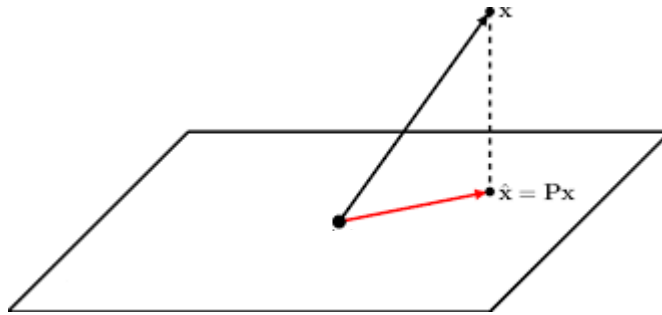Accounts for phylogeny during analysis

Before

After:
$E_{resid}$ independent
of phylogeny

$$\widehat{\boldsymbol{\beta}} = (\mathbf{X}^t\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}^t\mathbf{V}^{-1}\mathbf{Y}$$

Images from Collyer & Adams (2021)

Grafen (1989)

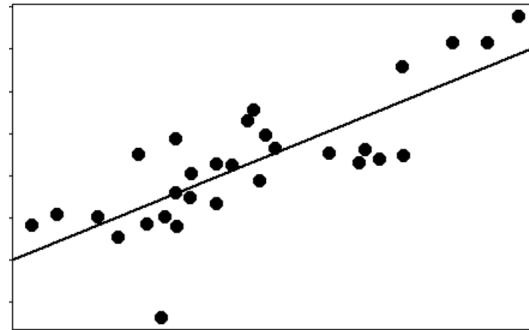See: Garland and Ives (2000)
Rohlf (2006); Blomberg et al (2012); Adams (2014a)

Three equivalent algebraic implementations

3: Phylogenetic Transformation (GLS→OLS)



Project data to phylogenetically-transformed space

Analysis of **Py** vs. **Px**

$$\widehat{\boldsymbol{\beta}} = \left( \widetilde{X}^t \ \widetilde{X} \right)^{-1} \widetilde{X}^t \ \widetilde{Y}$$

We utilize this procedure!

Garland and Ives (2000)
Adams (2014)

See: Garland and Ives (2000)
Rohlf (2006); Blomberg et al (2012); Adams (2014a)

Leverage geometry to obtain robust summary statistics



Images from Anderson (2001)
See: Gower (1966); Goodall (1991); Anderson (2001)

One way: Sums-of-squares from object distances*

Avoids $|\mathbf{V}| = 0$, but still obtains: SS, MS, F, $R^2$, etc.

*Note: approach also used for Goodall's F-test

See: Adams (2014a), (2014b) (2014c)
Adams and Collyer (2015), (2018a), (2018b); Adams and Collyer (2019)

# Significance testing via RRPP (Residual Randomization in Permutation Procedures)

### 1: Fit models

obtain β, and summary stats, SS, MS, $R^2$, F
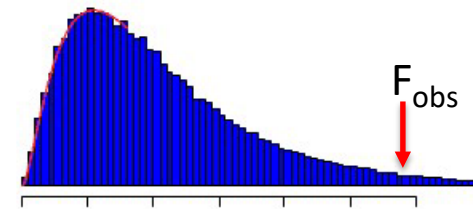
| Full $\mathbf{Y} = \mathbf{X1} + \mathbf{X2} + \mathbf{E}_f$ | Reduced $\mathbf{Y} = \mathbf{X1} + \mathbf{E}_r$ |

### 2: Permute $\mathbf{E_R}$ (residuals of $\mathbf{Y}$)

obtain pseudo-values: $\boldsymbol{y} = \widehat{\mathbf{Y}} + \mathbf{E}_r^*$

### 3: Fit model with $\boldsymbol{y}$, repeat

$F_{obs}$

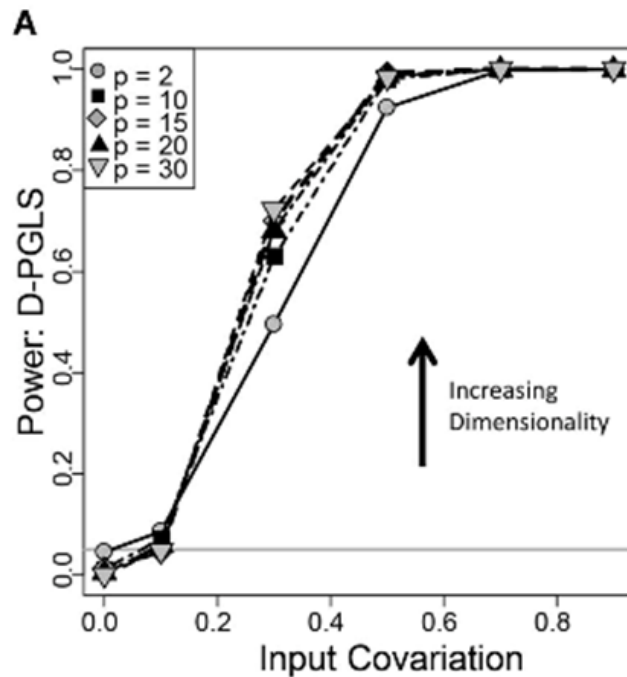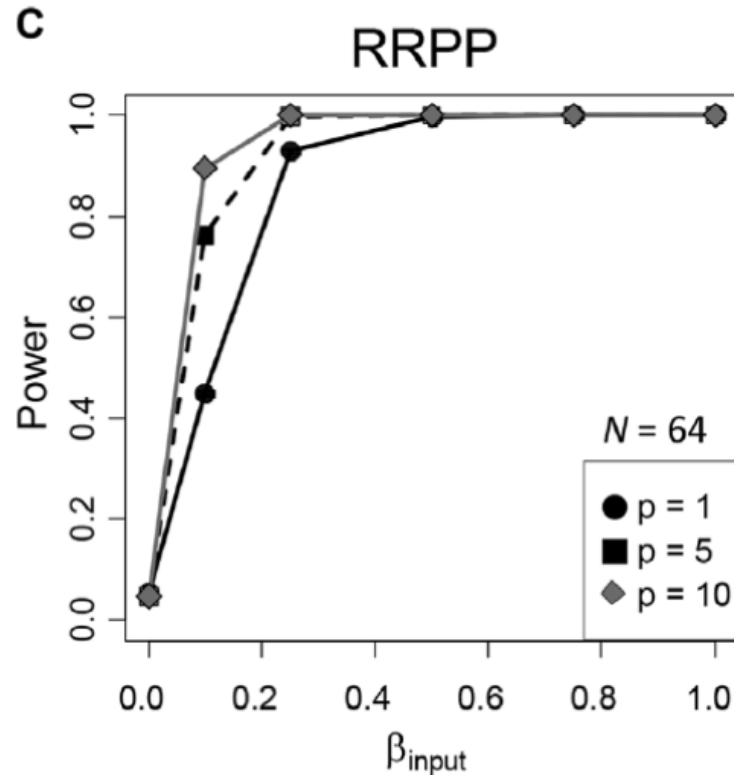### 4: Effect size: $z = \dfrac{(\log(F) - \mu_{\log(Fr)})}{\sigma_{\log(Fr)}}$

*Note: Proper permutation requires identifying correct exchangeable units (Commanges 2003: Adams and Collyer 2018a,b).

Collyer, Sekora, Adams (2015)
Adams & Collyer. (2016)
Adams & Collyer. (2018a); (2018b)

## Breaks Rao's paradox
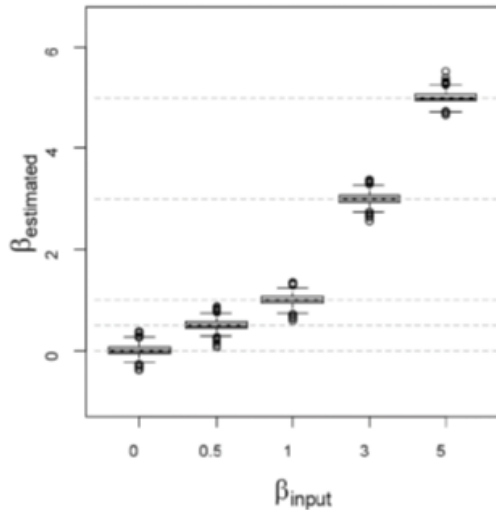


Adams 2014. *Evolution*

Adams and Collyer 2018. *Evolution*

**Displays appropriate type I
error and high power**

## RRPP sampling distribution matches theory (but extends to p>>N)
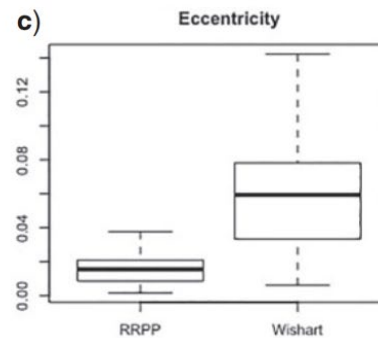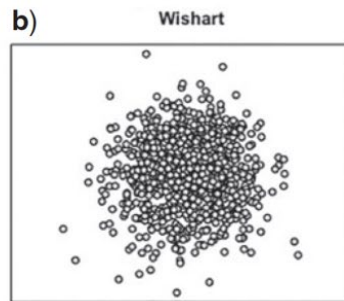


**Correct parameter estimates**

Adams and Collyer (2018a)



**Empirical sampling
distribution matches theory**

Adams and Collyer (2018a)



**Estimated covariance matrices
equivalent to sampling a Wishart
distribution**

Adams and Collyer (2018b)

**Conclusion: RRPP provides analytics for multivariate PCMs (and other applications)**

## Phylo-transform + RRPP enables multivariate PCMs

### 1: PGLS: Phylogenetic ANOVA/Regression



### 2: Phylogenetic PLS
(evolutionary covariation of 2 *SETS* of variables)



### 3: Net Evolutionary Rates



### 4: Phylogenetic Signal



**This facilitates investigations of the macroevolution of shape and other complex phenotypes**

Adams (2014a), 2014b), (2014c)
Adams & Felice (2104)
Adams and Collyer (2015);
Denton and Adams (2015)
Adams and Collyer (2018a), (2018b), (2019)
Collyer, Baken, Adams (2022)

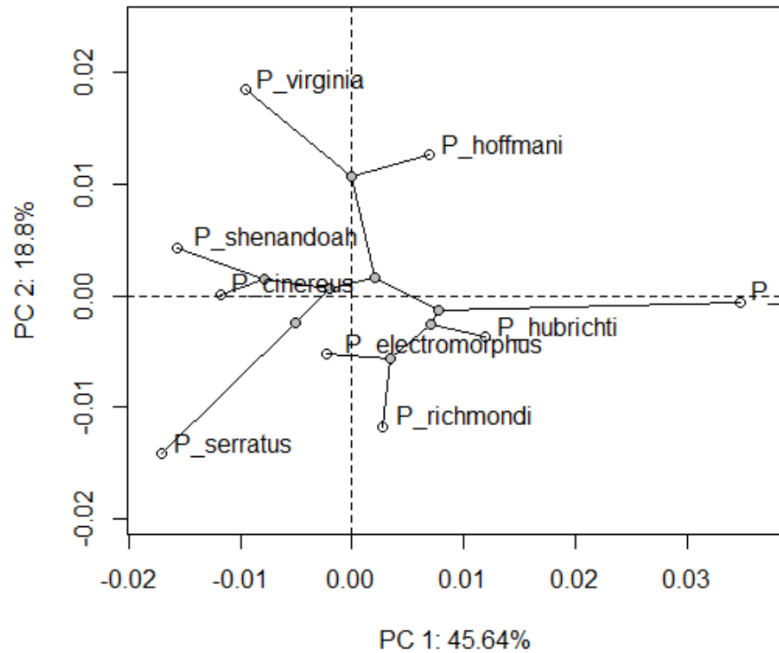# How to Visualize Evolutionary Patterns?

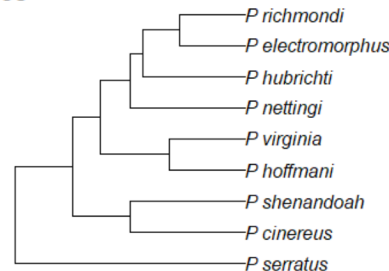## Phylomorphospace (PCA)
Align data to directions of maximal variation

## Phylogenetic PCA (pPCA)
Align data to directions *independent* of phylogenetic signal (1st dimension)



Interpretation can be challenging (e.g., with mixed ecological and phylogenetic signal)

Align data to directions that maximize phylogenetic signal



Data
$\mathbf{P}_{BM} + \mathbf{R}_{noise}$

Data
$\mathbf{P}_{BM} + \mathbf{E}_{ecol} + \mathbf{R}_{noise}$

**PACA reveals phylogenetic trends in data irrespective of other signals!**

Collyer and Adams (2021)

Account for phylogeny during PLS correlation

-PLS of evolutionary covariance (rate) matrix

$$\mathbf{R} = \frac{\left(\mathbf{Y} - E\left(\mathbf{Y}\right)\right)^t \mathbf{C}^{-1} \left(\mathbf{Y} - E\left(\mathbf{Y}\right)\right)}{N - 1}$$

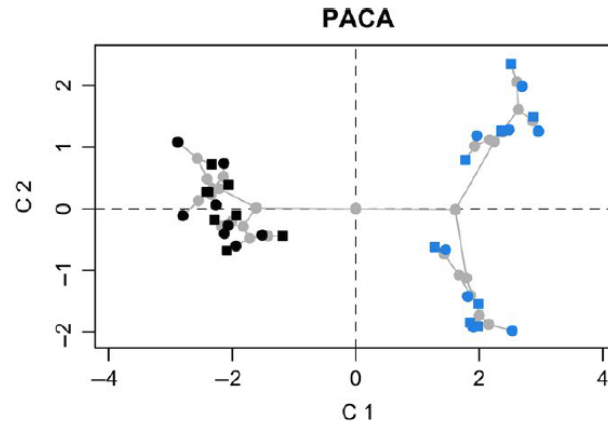$$\mathbf{R}_{12} = \mathbf{U}_{\mathbf{R1}} \mathbf{D} \mathbf{V}^t_{\mathbf{R}2}$$

-Equivalently found from PLS of **PY** (phylo-transformed data)

-Significance found from permutation of phylo-transformed data

Adams and Felice 2014. *PloS One.*
Adams and Collyer. 2018. *Syst. Biol.*

## PLS of cranium vs. mandible in *Plethodon*



a)

b)

c)

cranial shape

cranial shape

mandible shape

mandible shape

PLS Scores (Cranium)

PLS Scores (Mandible)

$r_{PLS}$ = 0.813
$P_{rand}$ = 0.008

Adams and Felice 2014. *PloS One.*

The degree to which phenotypic similarity associates with phylogenetic relatedness

-Blomberg's K: one measure (Adams, 2014 generalized to multivariate)

$$K = \frac{(\mathbf{Y} \text{-} E(\mathbf{Y}))^t (\mathbf{Y} \text{-} E(\mathbf{Y}))}{(\mathbf{Y} \text{-} E(\mathbf{Y}))^t \mathbf{C}^{\text{-}1} (\mathbf{Y} \text{-} E(\mathbf{Y}))} \Bigg/ \frac{tr(\mathbf{C}) - N(\mathbf{1}^t \mathbf{C}^{\text{-}1} \mathbf{1})^{\text{-}1}}{N-1}$$

- Pagel's $\lambda$: a branch-length transformation during logL fitting

$$logL = \log \left[ \frac{\exp\left(-0.5 (\mathbf{Y} - E(\mathbf{Y}))^t \mathbf{V}^{\lambda-1} (\mathbf{Y} - E(\mathbf{Y}))\right)}{\sqrt{(2\pi)^{Np} \times |\mathbf{V}^\lambda|}} \right]$$

Original ($\lambda$=1)    $\lambda$ = 0.5    $\lambda$ = 0 (star)

Both K and $\lambda$ are related to a permutation-based Z-score, which can be used to compare the strength of signal across datasets



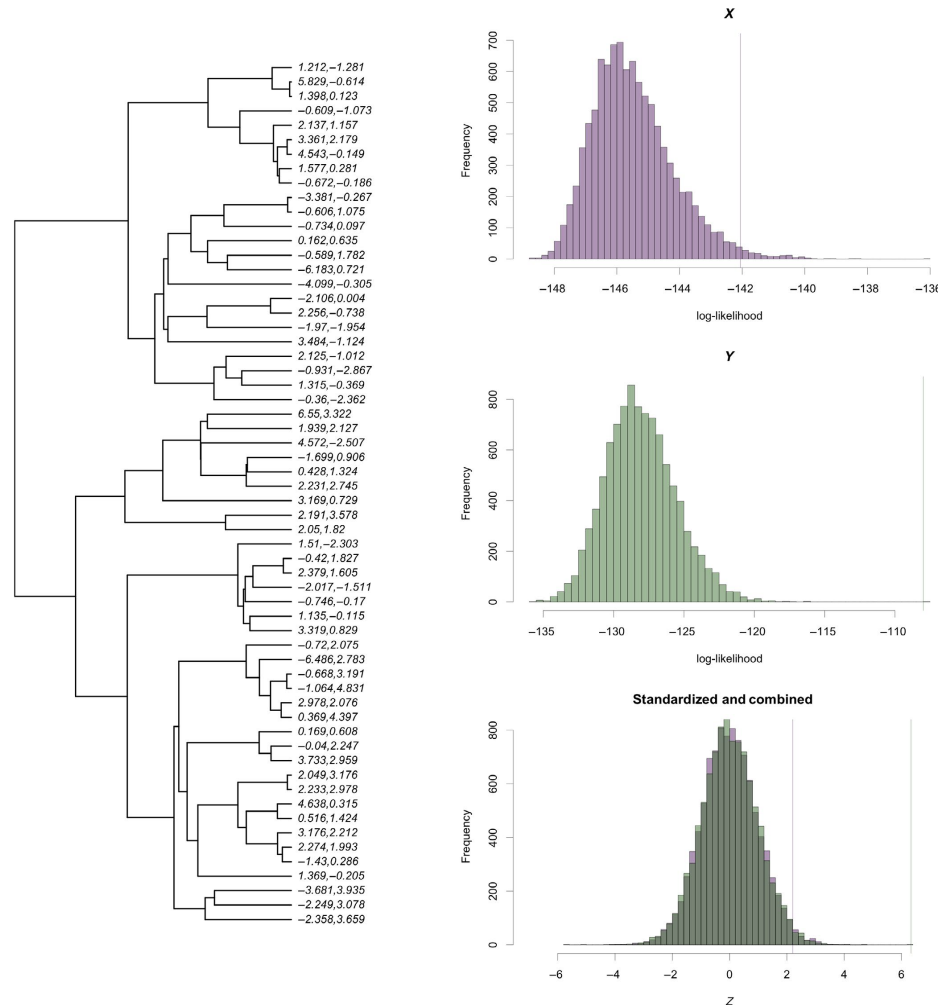FIGURE 1 Plot of phylogenetic tree with $x, y$ values, plus frequency histograms for the RRPP log-likelihood values for two variables, $X$ and $Y$. Vertical lines indicate observed values. In the last panel, histograms have been combined for standardized values

Collyer, Baken, Adams (2022). *MEE.*

What evolutionary model best describes trait variation?

-Fit data to phylogeny under differing evolutionary models



Evolutionary Models

BM          OU          Butler and King 2004

Model comparisons of:
   1: Evolutionary rates (and covariances):  BM1, BMM, etc.
   2: Evolutionary 'modes': BM, OU1, OUM, etc.

Methods for multivariate data:

   1: $logL_{Mult}$ (Revell and Harmon, 2008; Clavel et al. 2015)

   2: $\sigma^2_{mult}$ (Adams, 2014; Denton and Adams, 2015)

   3: $\Sigma logL_{indiv}$ (Ingram & Mahler, 2013; Grundler and Rabosky, 2014; Moen et al. 2016)

   4: PCL (Goolsby, 2016)

## Evolutionary rate for a trait $\sigma^2$: Phylogenetically-standardized variance

-Estimated from data and phylogeny under Brownian motion (see Felsenstein 1973)

$$\sigma^2 = \frac{\left(\mathbf{Y} - E(\mathbf{Y})\right)^t \mathbf{C}^{-1} \left(\mathbf{Y} - E(\mathbf{Y})\right)}{N}$$

Felsenstein 1973. *Am J. Hum. Gen.*

$$\mathbf{R} = \begin{bmatrix} \sigma_1^2 & & \\ \sigma_{21} & \sigma_2^2 & \\ \sigma_{31} & \sigma_{32} & \sigma_3^2 \end{bmatrix}$$

Revell & Harmon 2008. *Ev. Ecol. Res.*

**Is there evidence for multiple evolutionary rates on the phylogeny?**

1: Define 'regimes' for models (BM1, BMM, etc.)



Vs.

2: Estimate $\sigma^2$ (**R** multivariate) and log$L_{mult}$

3: Compare log$L$ (<u>LRT tests, AIC, phylogenetic simulation, etc</u>).

# Type I error of LRT ↑ with p (not useful for high-dimensional data)



Embodiment of 'curse of dimensionality'
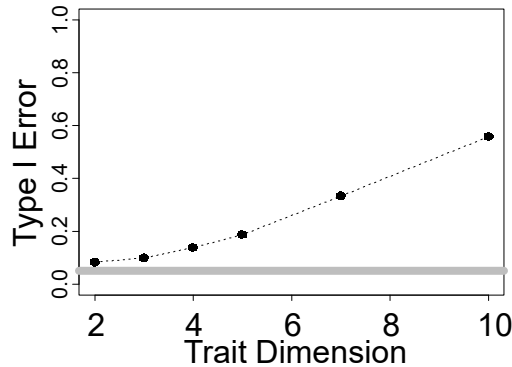
Adams 2014b. *Syst. Biol.*

BM1-simulations
(N=32, p=8)

$BM_1$ vs. $BM_2$
(mvMORPH)

Adams and Collyer 2018. *Syst. Biol.*

# $logL_{Mult}$ cannot be computed when p≥N

For multi-dimensional traits, should have a single rate, not a matrix

# LRT based on the $logL_{mult}$ not a general solution for high-D data

Adams. 2014b. *Syst. Biol.*
Adams and Collyer 2018. *Syst. Biol.*
Adams and Collyer, 2019. *Ann. Rev. Ecol. Evol. Syst.*

Pairwise composite likelihood (PCL) as an alternative to $\log L_{mult}$

1: Define 'regimes' for models (BM1, BMM, etc.)



Vs.

2: Fit $H_0$ and $H_1$ for **PAIRS** of variables; obtain $\log L_{pair}$

3: Sum across $\log L_{pair}$ for overall fit: $\Sigma \log L_{pair}$

4: Simulate data under $H_0$ and compare

## Pairwise composite likelihood to compare BM1 vs. BMM

- Sensitive to *ALL* aspects of multivariate metric spaces



- Arbitrary results
  - Orientation-dependent
  - Cov-Y dependent



Data simulated under BM2
(N=32, p=8) with known
difference in $R_1$ vs $R_2$

- PCL **NOT** useful for comparing evolutionary rates

# Generalize σ² for multidimensional data: **net evolutionary rates**

- •Define 'regimes' for models (BM1, BMM)
- •Phylogenetic transform of data
- •Estimate $\sigma^2_{\text{mult}}$ for BM1, BMM
- •Permute (or simulate), repeat

$$\sigma^2_{mult} = \frac{\mathbf{PD}^t_{\mathbf{U},0}\mathbf{PD}_{\mathbf{U},0}}{N}$$

Adams 2014a. *Syst. Biol.*
Denton and Adams. 2015. *Evol.*

•Method rotation-invariant, and appropriate Type I error/power



• $\sigma^2_{mult}$ **IS** useful for comparing multivariate evolutionary rates!

Adams and Collyer 2018. *Syst. Biol.*

Evolutionary models go beyond Brownian motion

    -BM, OU, EB, ACDC, etc.

    -Fit data to phylogeny under differing evolutionary models

Evolutionary Models



BM      OU      Butler and King 2004

Methods for multivariate data:

    1: $logL_{Mult}$ (Clavel et al. 2015: extending Revell & Harmon, 2008)

    2: $\Sigma logL_{indiv}$ (Ingram & Mahler, 2013; Grundler and Rabosky, 2014; Moen et al. 2016)

    3: PCL (Goolsby, 2016)

1: logL$_{mult}$ (various implementations)

AIC: Model misspecification ↑ with *p* (not useful for high-dimensional data)



AIC from logL$_{Mult}$ not general solution for model comparisons with high-D data

Adams and Collyer 2018. *Syst. Biol.*
Adams and Collyer, 2019. *Ann. Rev. Ecol. Evol. Syst.*

## 2: PCL

- Sensitive to *ALL* aspects of multivariate metric spaces



BM-simulations (p=8) on a 32 species phylogeny
(mean of 100 simulations per scenario)

- High misspecification and arbitrary results



Data simulated under BM1
(N=32, p=8)

- PCL **NOT** useful

Adams and Collyer 2018. *Syst. Biol.*
Adams and Collyer, 2019. *Ann. Rev. Ecol. Evol. Syst.*

3: Evaluate multivariate space dimension by dimension

- Assume trait independence
- Fit evolutionary models separately (on $PPC_1$, $PPC_2$, etc.)
- Obtain $\Sigma logL$ and corresponding AIC to infer best model

(Ingram and Mahler, 2013; Grundler and Rabosky, 2014; Moen et al. 2016)

Mathematical Problems:

- Individual PCs misspecify model
  - EB preferred on lower PCs even for BM data (Uyeda et al. 2015)

- Dimensions not independent evolutionarily (mis-application of Edward's likelihood theorem)
  - Independence when **R** (NOT **S!**) is diagonal
  - ONLY occurs under BM for PPCA
  - For all other models, dimensions evolutionarily correlated
  - Thus, $\Sigma logL \neq LogL_{Mult}$

Adams and Collyer 2018. *Syst. Biol.*

## Consequence: $\Sigma log_{ind}$ greatly supports overly complex models

## Example:  Simulate datasets under BM, infer best model

- 2 or more inferred OU optima = misspecification



Data simulated under BM1
(N=32, p=8)

## Result: > 95% model misspecification!

-NOTE: Comparing observed pattern to set of simulated outcomes post-hoc is not informative,
as one cannot distinguish the 'true' pattern in the observed from the pattern generated by method

## Conclusion: $\Sigma log_{ind}$ methods not reliable

Adams and Collyer 2018. *Syst. Biol.*

Multivariate PCM not trivial

-Algebraic generalizations
  appropriate mathematically

-Useful for hypotheses of:

1: Phylogenetic signal ($K_{mult}$)
2: ANOVA/regression (D-PGLS)
3: Correlation (PPLS)
4: Net evolutionary rates ($\sigma^2_{mult}$)

| Analysis Type | $logL_{Mult}$ | $\Sigma logL$ | PCL | MultG |
|---|---|---|---|---|
| **Phylogenetic Signal** | - | - | - | Yes ($K_{mult}$) |
| **Phylogenetic ANOVA** | - | - | NO | Yes (D-PGLS) |
| **Phylogenetic Regression** | - | - | NO | Yes (D-PGLS) |
| **Phylogenetic Covariation (blocks of variables)** | - | - | NO | Yes (P-PLS) |
| **Comparing Evolutionary Models: BM1 vs BMM** | Limited (when N>>>p) | - | NO | Yes (net rate only) |
| **Comparing Evolutionary Models: BM1 vs BMM vs OU1 vs OUM** | No | No | No | - |

**Current limitation: Brownian motion only**

Adams and Collyer 2018. *Syst. Biol.*
Adams and Collyer, 2019. *Ann. Rev. Ecol. Evol. Syst.*

Multivariate PCM not trivial

-Evolutionary model
 comparisons remain
 a challenge

| Analysis Type | logL$_{Mult}$ | $\Sigma$logL | PCL | MultG |
|---|---|---|---|---|
| **Phylogenetic Signal** | - | - | - | Yes ($K_{mult}$) |
| **Phylogenetic ANOVA** | - | - | NO | Yes (D-PGLS) |
| **Phylogenetic Regression** | - | - | NO | Yes (D-PGLS) |
| **Phylogenetic Covariation (blocks of variables)** | - | - | NO | Yes (P-PLS) |
| **Comparing Evolutionary Models: BM1 vs BMM** | Limited (when N>>>p) | - | NO | Yes (net rate only) |
| **Comparing Evolutionary Models: BM1 vs BMM vs OU1 vs OUM** | No | No | No | - |

**<span style="color:red">Multivariate Ornstein-Uhlenbeck models a particular challenge</span>**

**<span style="color:red">We lack a robust multivariate method for evolutionary model comparisons!</span>**

Adams and Collyer 2018. *Syst. Biol.*
Adams and Collyer, 2019. *Ann. Rev. Ecol. Evol. Syst.*