

Brownian motion models, multiple characters, and phylogenies

Bio 550D: Morphometrics in Biology

Joe Felsenstein

13 October 2016

What will approximate change of quantitative characters?

- ... when it occurs by genetic drift of pre-existing alleles?
- ... when it also occurs by mutation to new alleles?
- ... when variable selection affects the alleles at each locus?
- ... when selection is on the fitness based on the whole phenotype?

Approximating genetic drift of two alleles

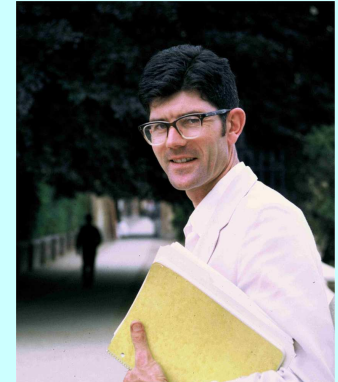
Can we compute transition probabilities for genetic models such as the Wright-Fisher model?

- Can we do this analytically? **No.** (although the right eigenvectors and the eigenvalues are known) the full set of left eigenvalues has never been derived.
- We can take such a model for a given (not-too-big) population size N and compute the transition probability matrix, then either power it up numerically or get its eigenvalues and eigenvectors
- OK, what about the diffusion approximation. Aren't they very close approximations? Yes, they and Kimura (1955a, 1955b) derived transition probabilities for the diffusion process as sums of series in Gegenbauer polynomials. **But** they are difficult to work with.

Edwards and Cavalli-Sforza's approximation



Luca Cavalli-Sforza (and Edwards), 1963



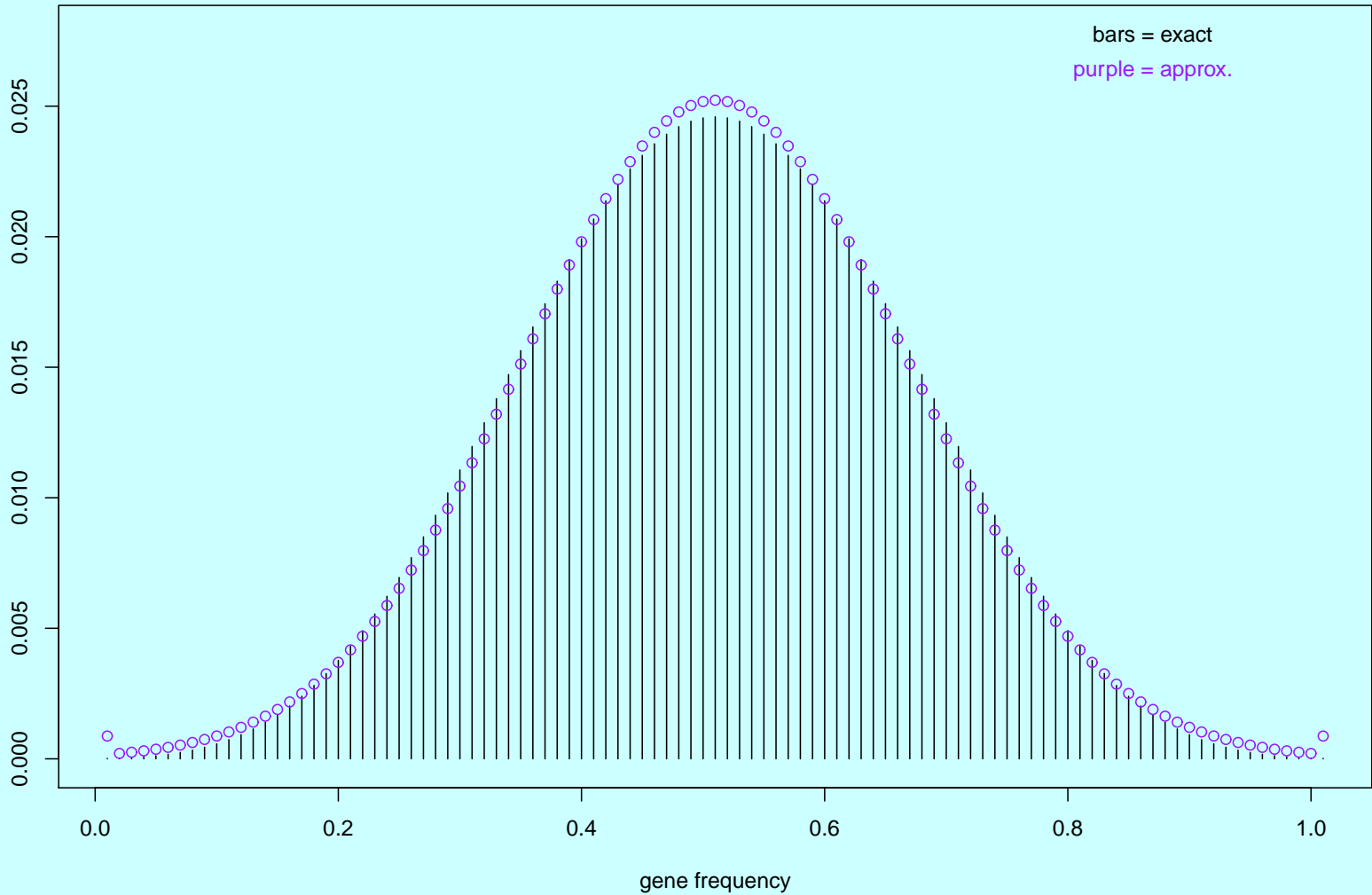
Anthony Edwards, 1970

The expectation of gene frequency change in one generation (under pure genetic drift without mutation) is zero. The variance is the binomial variance

$$E \left[(\Delta p)^2 \right] = \frac{p(1-p)}{2N_e}$$

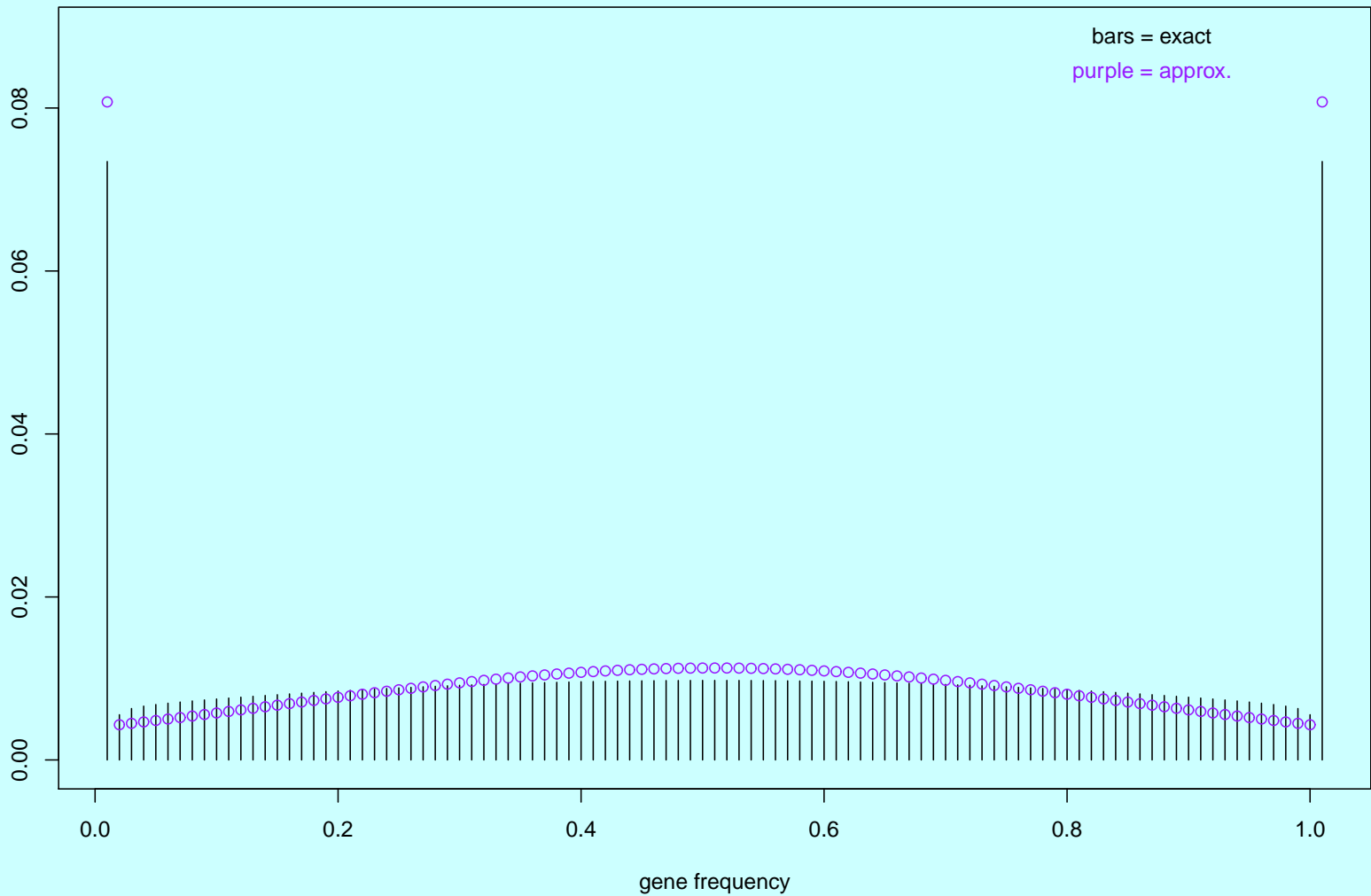
That variance is not constant: it varies with p (in a parabola), but maybe we can roughly approximate it by dealing with the case where all populations have roughly similar gene frequencies, so the variances are nearly the same. Maybe. Roughly.

How good is this?



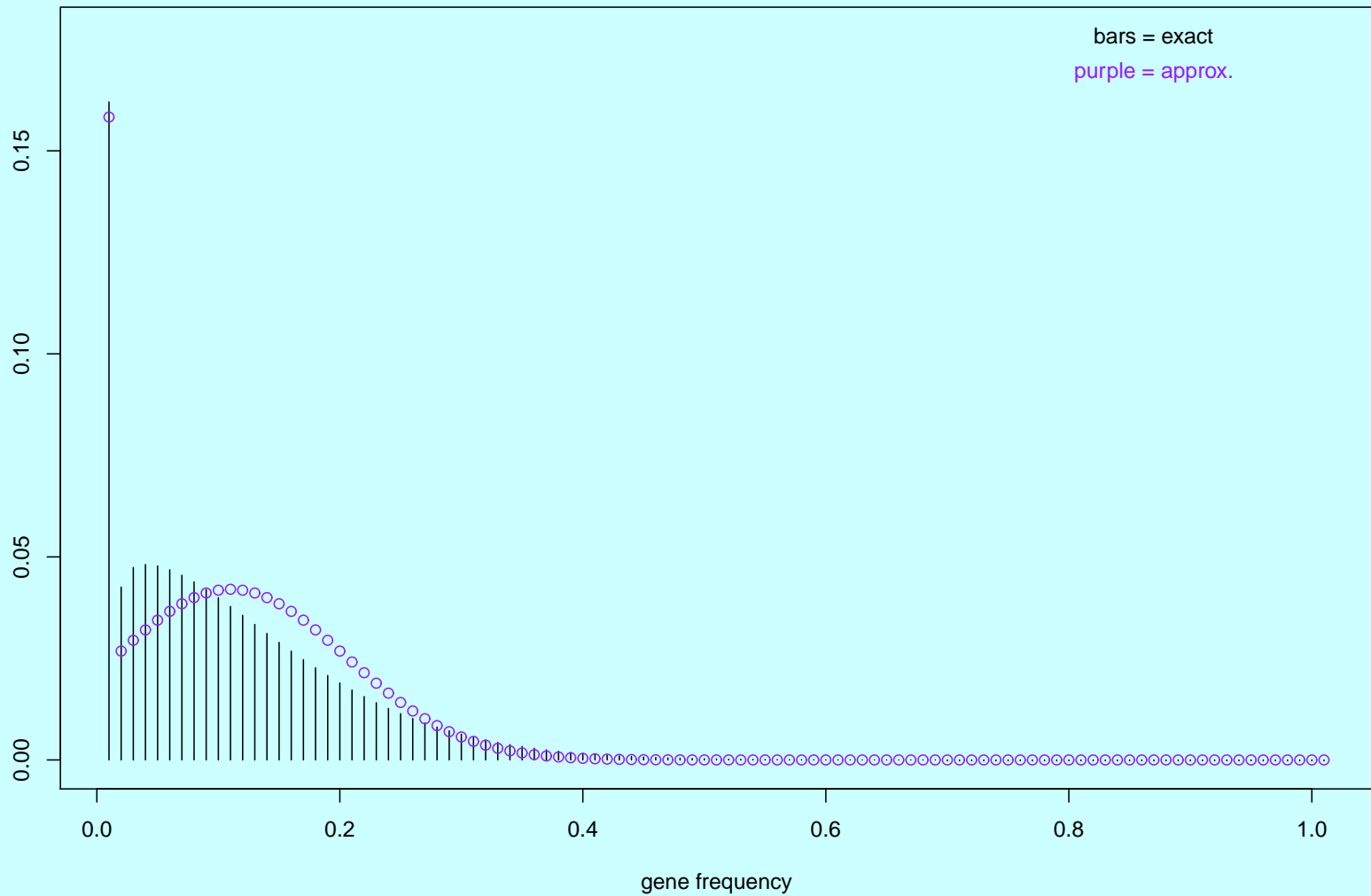
Starting with $p = 0.5$, after 10 generations in a population of size 50.

How good is this?



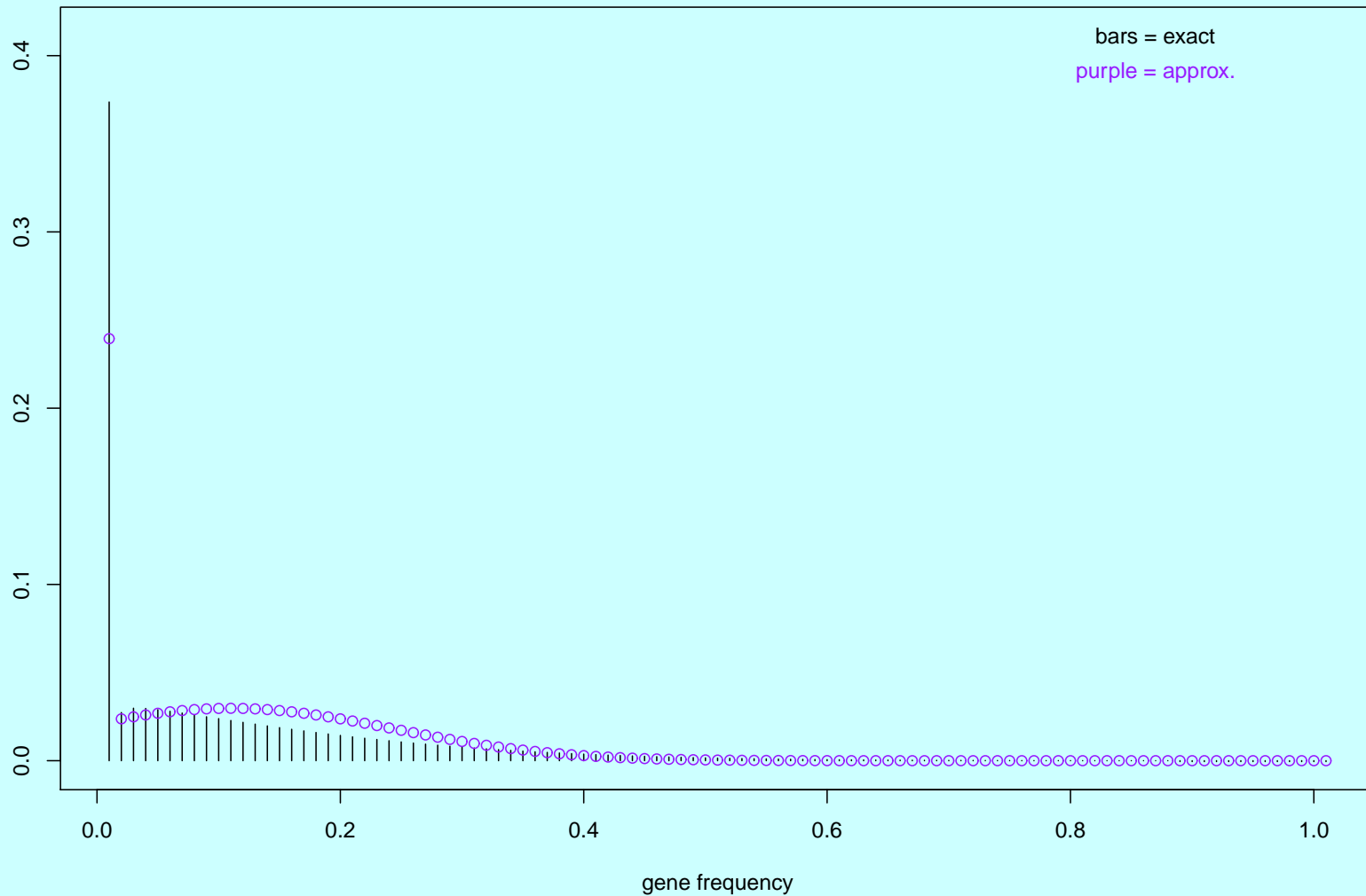
Starting with $p = 0.5$, after 50 generations in a population of size 50.

How good is this?



Starting with $p = 0.1$, after 10 generations in a population of size 50.

How good is this?



Starting with $p = 0.1$, after 20 generations in a population of size 50.

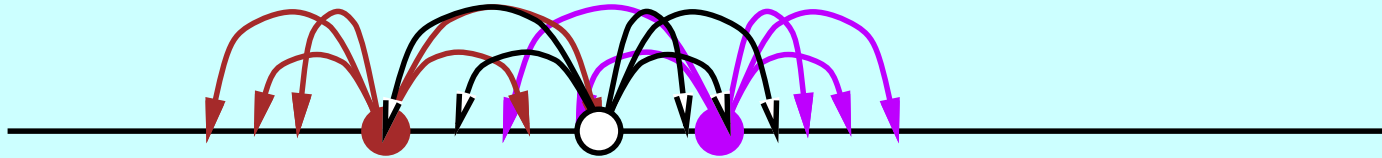
What about a quantitative character?

If a quantitative character is a sum of contributions from a number of loci, then if the individual locus gene frequencies have their change approximated by Brownian Motion, the linear combination will also change by Brownian motion. This works for multiple alleles.

- if there is any dominance, there will be some nonlinearity and the approximation will be less good.
- Epistasis can cause even more trouble.

First discussed by me (Felsenstein, 1973).

But, if there are mutations making incremental changes ..



... as we saw with the discussion of quantitative characters, if a relatively constant genetic variance is maintained, and mutations have additive effects, then genetic drift will cause the mean to change in a random walk close to Brownian Motion.

However, if one approaches some limit where most mutations oppose movement to it, and there are no mutations allowing you to go past that limit, this approximation will be poor.

Brownian motion is mathematically tractable

You can easily compute transition probabilities from one value to another, since the net change after “time” t is normal, with mean zero and variance $\sigma^2 t$, and changes in successive time intervals are independent.

When two lineages share a period of common ancestry, the resulting tip species have phenotypes that covary, the covariance being the variance expected during their shared ancestry.

Covarying character change along a lineage

What is the distribution of changes in multiple characters (say p of them) along a lineage? Simply the appropriate multiple of the infinitesimal rate of change per unit branch length.

If a set of characters \mathbf{x} , changes under covarying Brownian motion, in time t (or a pseudo-time branch-length t) the change will be distributed as

$$\Delta\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{V}t),$$

(where \mathbf{V} is the covariance matrix of the infinitesimal change of the Brownian Motion).

What causes change in quantitative characters?

For neutral mutation and genetic drift, can show that for a quantitative character with additive genetic variance V_A and population size N the genetic (additive) value of the population mean is:

$$\text{Var}(\Delta\bar{g}) = V_A/N$$

If mutation and drift are at equilibrium:

$$E \left[V_A^{(t+1)} \right] = V_A^{(t)} \left(1 - \frac{1}{2N} \right) + V_M$$

In neutral traits additive genetic variance rules

so that

$$E[V_A] = 2NV_M$$

whereby

$$\text{Var}[\Delta\bar{g}] = (2NV_M) / N = 2V_M$$

an analog of Kimura's result for neutral mutation.

Thus to transform characters to independent Brownian motions of equal evolutionary variance, we could use the additive genetic variance V_A .

With selection ... life is harder

There is the “Breeder’s Equation” of Wright and Fisher (1920’s)

$$\Delta z = h^2 S$$

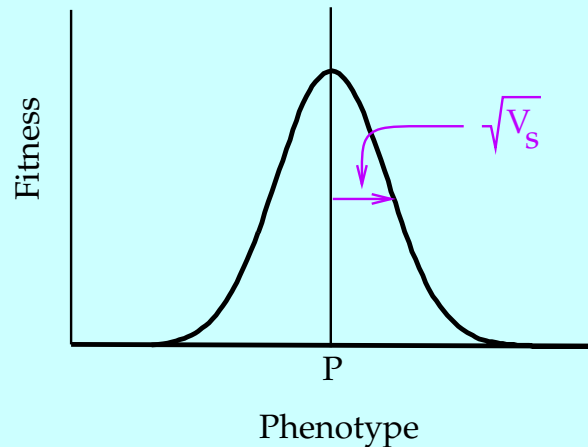


and Russ Lande’s (1976) recasting of that in terms of slopes of mean fitness surfaces:

$$S = V_P \frac{d \log(\bar{w})}{d\bar{x}}$$
$$\Delta z = (V_A/V_P) V_P \frac{d \log(\bar{w})}{d\bar{x}} = V_A \frac{d \log(\bar{w})}{d\bar{x}}$$

Note – it’s heritability times the slope of log of *mean* fitness with respect to *mean* phenotype. There is an exact multivariate analog of this equation.

Selection towards an optimum



If fitness as a function of phenotype is:

$$w(x) = \exp \left[-\frac{(x - p)^2}{2V_s} \right]$$

Then after some completing of squares and integrating, the change of mean phenotype “chases” the optimum:

$$m' - m = \frac{V_A}{V_s + V_P} (p - m)$$

(There is an exact matrix analog of this for multiple characters).

Sources of evolutionary correlation among characters

Variation (and covariation) in change of characters occurs for two reasons:

- **Genetic covariances.** (the same loci affect two or more traits)

Sources of evolutionary correlation among characters

Variation (and covariation) in change of characters occurs for two reasons:

- **Genetic covariances.** (the same loci affect two or more traits)
- **Selective covariances** (Tedin, 1926; Stebbins 1950). The same environmental conditions select changes in two or more traits – even though they may have no genetic covariance.

Sources of evolutionary correlation among characters

Variation (and covariation) in change of characters occurs for two reasons:

- **Genetic covariances.** (the same loci affect two or more traits)
- **Selective covariances** (Tedin, 1926; Stebbins 1950). The same environmental conditions select changes in two or more traits – even though they may have no genetic covariance.

Change of phenotypic means is a result of genetic covariances and selective covariances, where the former affects both response to selection and wandering due to genetic drift.

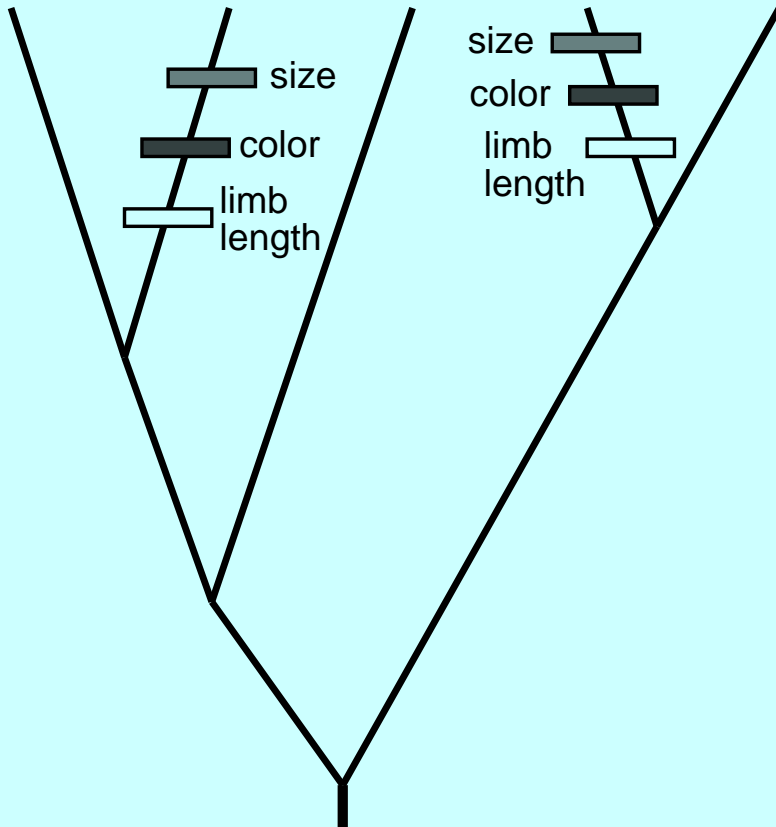
A simple example of selective covariance

covariation due not to genetic correlation
but to covariation of the selection pressure

These are Bergmann's, Allen's and Glogler's Rules
They are presumably not the result of genetic correlations
but result from patterns of selection

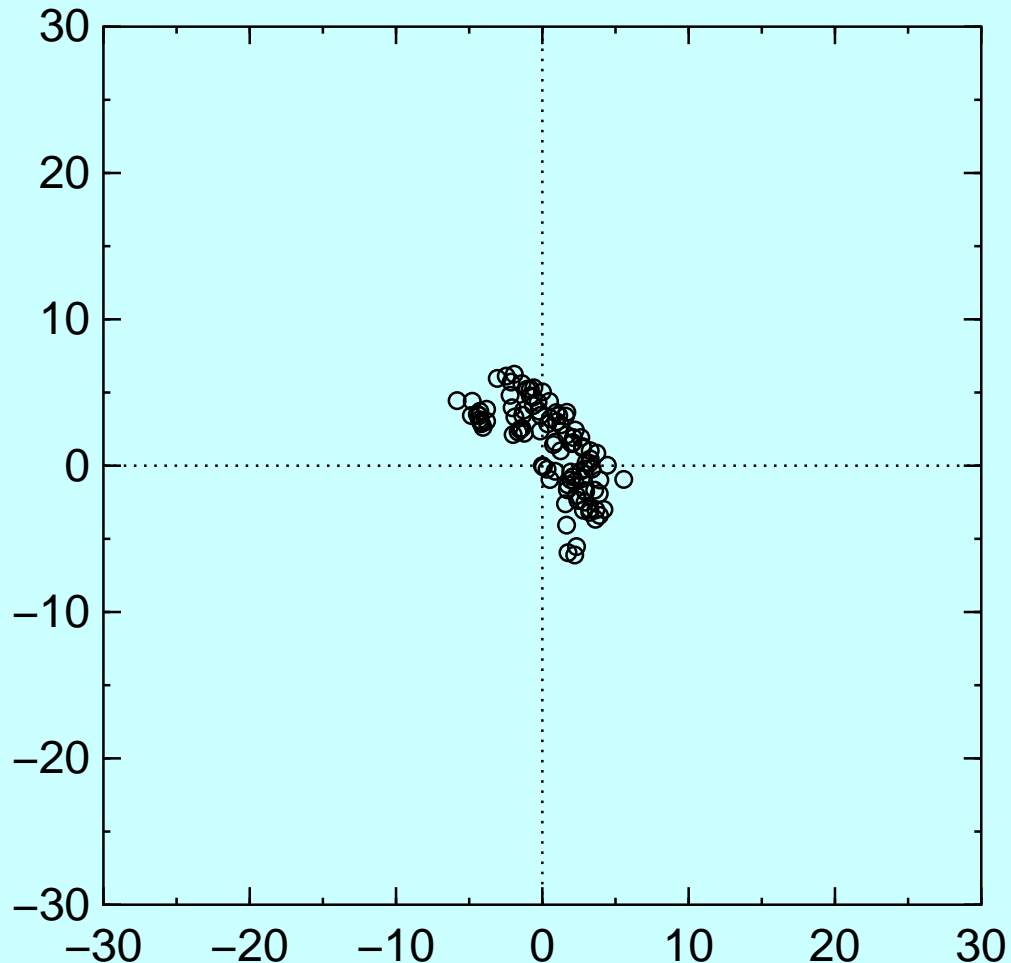
a simple example:

(temperate) (arctic) (temperate) (arctic) (temperate)



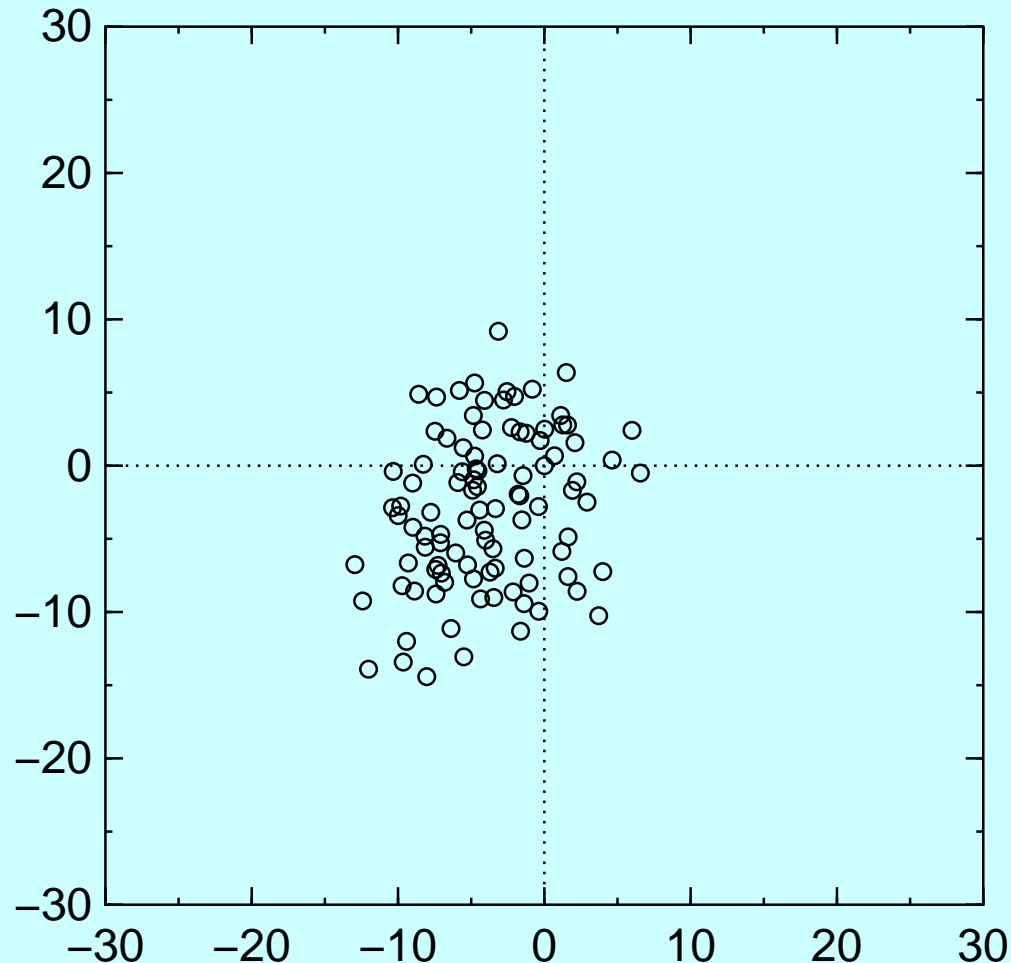
G. L. Stebbins. 1950. *Variation and evolution in plants*. Columbia Univ. Press, New York. page 121

Chasing a peak, simulated with two characters



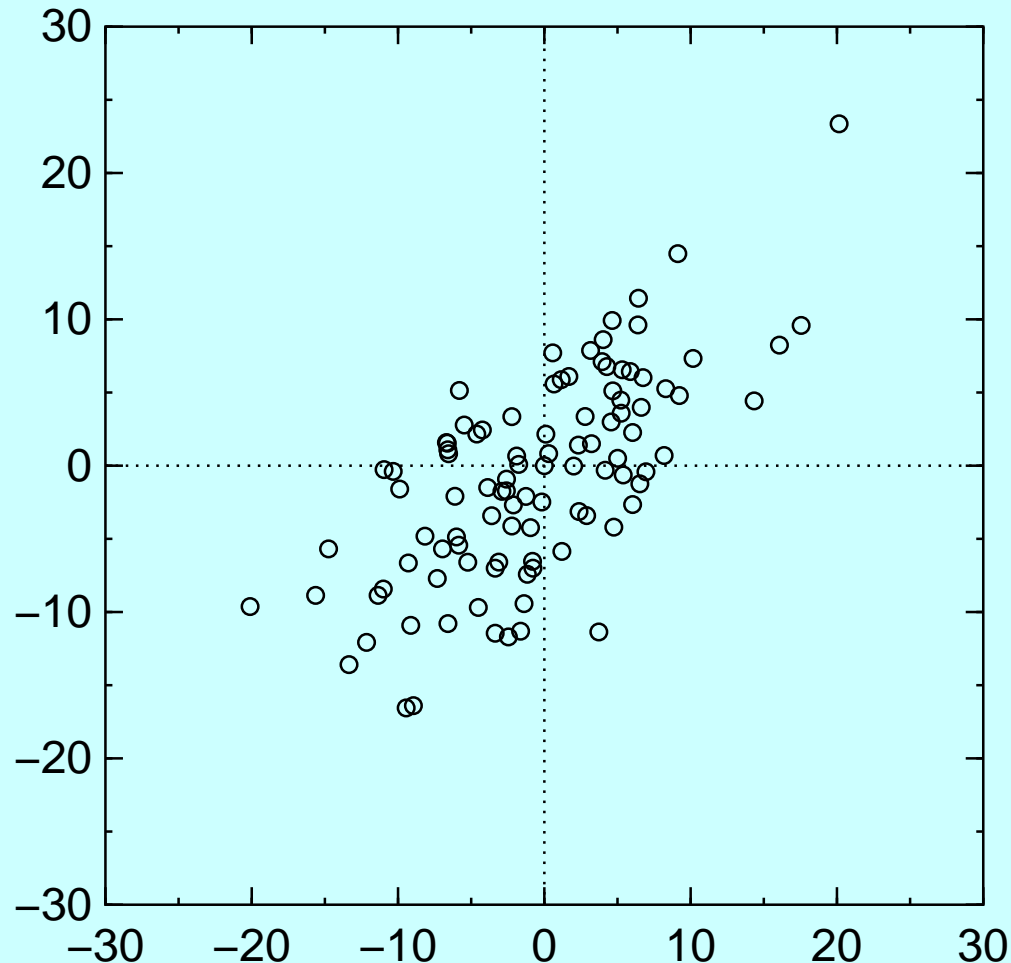
Genetic covariances assumed negative, but the wanderings of the adaptive peaks assumed positively correlated. In the first 100 generations the genetic covariances are most influential.

Chasing a peak, simulated with two characters



Genetic covariances assumed negative, but the wanderings of the adaptive peaks assumed positively correlated. After a while (every 10th generation up to generation 1000), the wanderings of the peaks start to be influential.

Chasing a peak, simulated with two characters



Genetic covariances assumed negative, but the wanderings of the adaptive peaks assumed positively correlated. In the long run (every 100th generation up to generation 10,000) the means go mostly where the peaks go.

A case that has received too little attention

- Suppose characters x and y are genetically correlated.
- and y is under optimum selection, but x is the one we observe.
- What will we see? In effect, the sum (actually, a weighted average) of an Ornstein-Uhlenbeck process and Brownian Motion.
- So Brownian motion restricted in the short run but not in the long run.
- It will look almost like Ornstein-Uhlenbeck Process with an optimum which wanders by Brownian Motion.

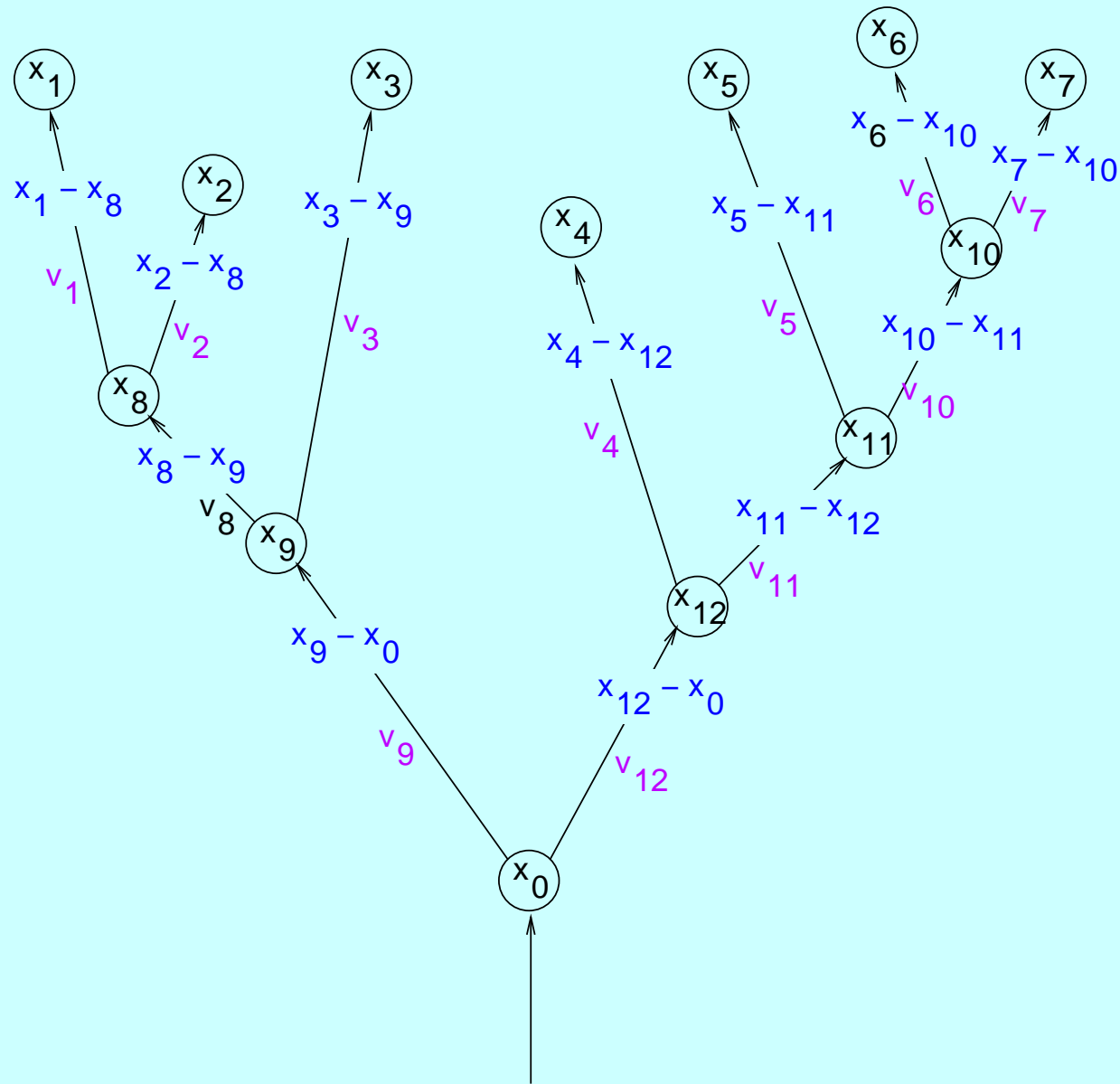
Most models so far do not allow for characters that are observed to covary with those that aren't observed.

A research program?

What we could imagine doing is:

- We might hope to infer additive genetic covariances by doing quantitative genetics breeding experiments to infer them from covariances among relatives, perhaps even in multiple species.
- Infer the covariances of the changes along the phylogeny.
- From them, back-calculate the selective covariances.
- The genetic covariances may also be inferrable from differences between nearby tips on the tree if we do not have breeding experiments.
- There is little or no hope of inferring “selective correlations” more directly without a complete understanding of the functional ecology.

Brownian motion along a tree



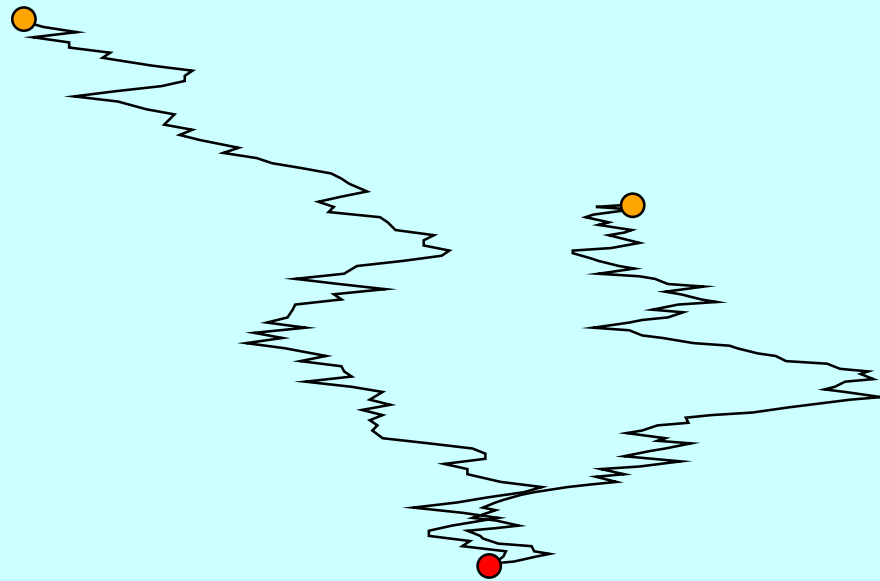
Covariances of species on the tree

$$\begin{bmatrix}
 v_1 + v_8 + v_9 & v_8 + v_9 & v_9 & 0 & 0 & 0 & 0 \\
 v_8 + v_9 & v_2 + v_8 + v_9 & v_9 & 0 & 0 & 0 & 0 \\
 v_9 & v_9 & v_3 + v_9 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & v_4 + v_{12} & v_{12} & v_{12} & v_{12} \\
 0 & 0 & 0 & v_{12} & v_5 + v_{11} + v_{12} & v_{11} + v_{12} & v_{11} + v_{12} \\
 0 & 0 & 0 & v_{12} & v_{11} + v_{12} & v_6 + v_{10} + v_{11} + v_{12} & v_{10} + v_{11} + v_{12} \\
 0 & 0 & 0 & v_{12} & v_{11} + v_{12} & v_{10} + v_{11} + v_{12} & v_7 + v_{10} + v_{11} + v_{12}
 \end{bmatrix}$$

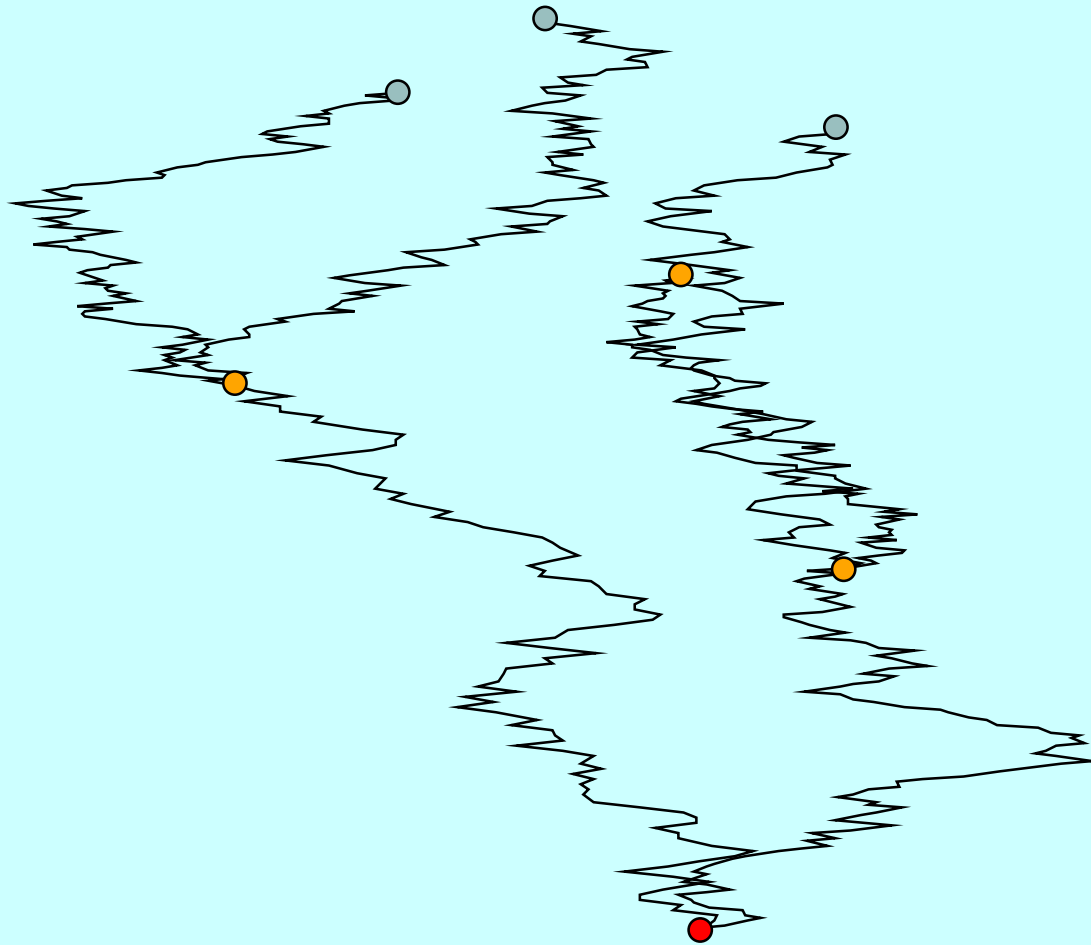
An outcome of Brownian motion on a 5-species tree



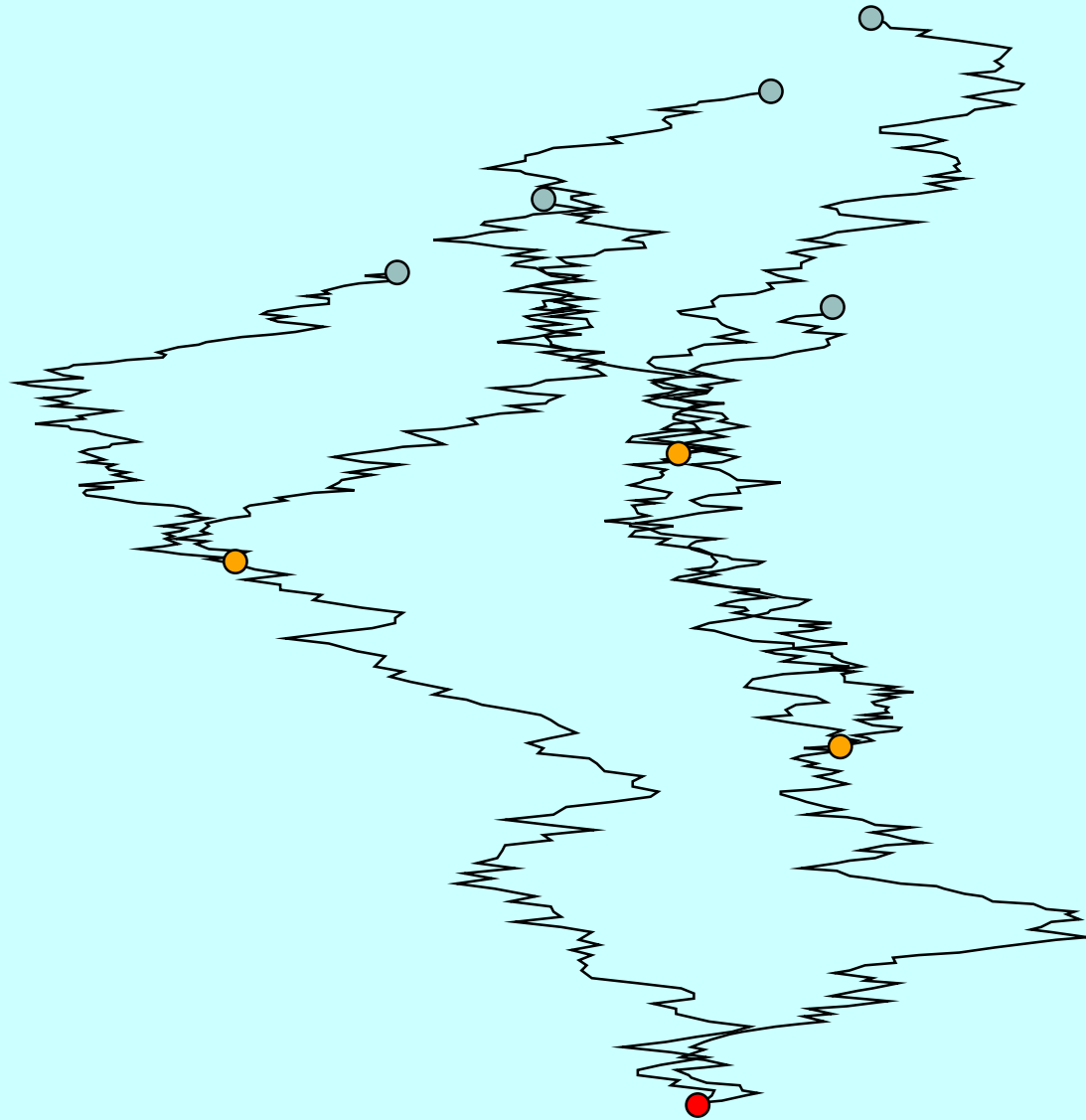
An outcome of Brownian motion on a 5-species tree



An outcome of Brownian motion on a 5-species tree



An outcome of Brownian motion on a 5-species tree



“Pruning” a tree in the Brownian motion case

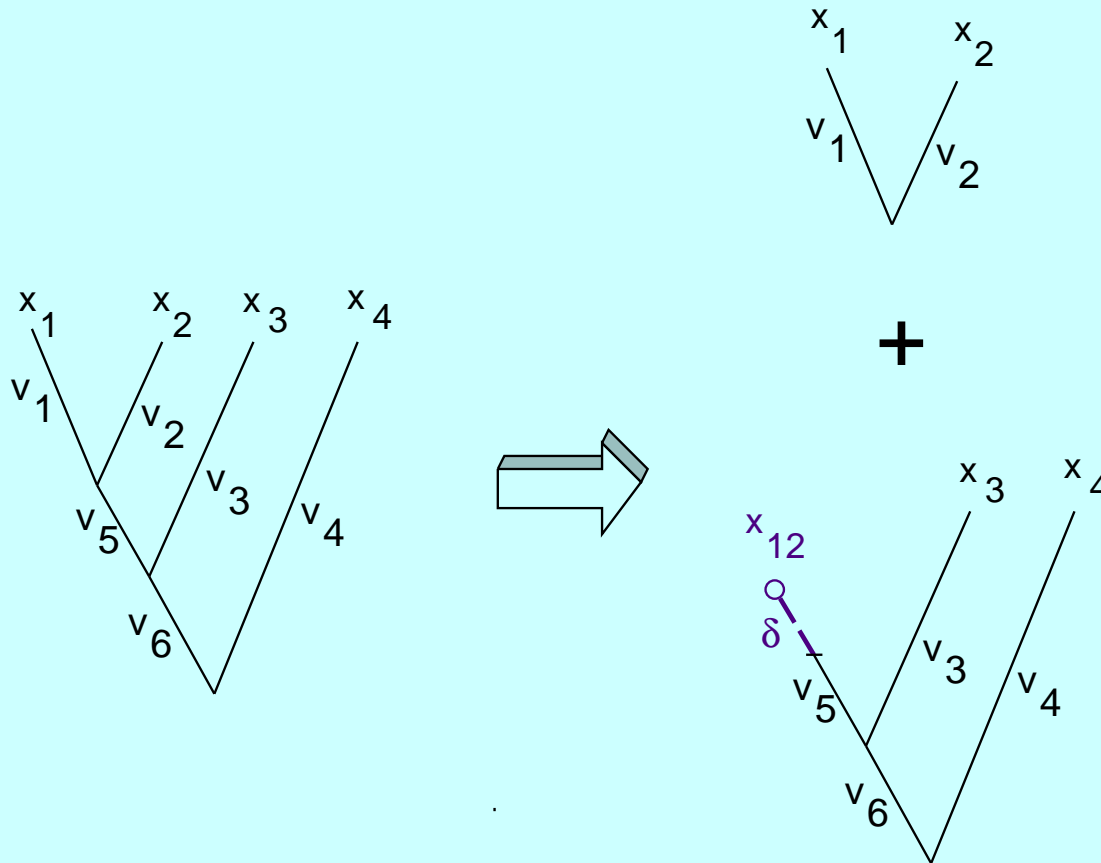
One can take two neighboring tips, and consider their difference $x_1 - x_2$ as well as a weighted average $ax_1 + (1 - a)x_2$. Using weights $a : 1 - a = 1/v_1 : 1/v_2$, the weighted average is independent of the difference, and the difference is also independent of the rest of the tree.

In fact, this weighted average behaves like a tip: Its covariances with the other species are the same as those of x_1 and x_2 . It acts just as if the tree were pruned, cutting off species 1 and 2, leaving a single species whose variance is a bit bigger.

$$\text{Var}[ax_1 + (1 - a)x_2] = v_8 + v_9 + \frac{v_1 v_2}{v_1 + v_2}$$

so in effect, a small extra amount of branch length is added.

“Pruning” a tree in the Brownian motion case

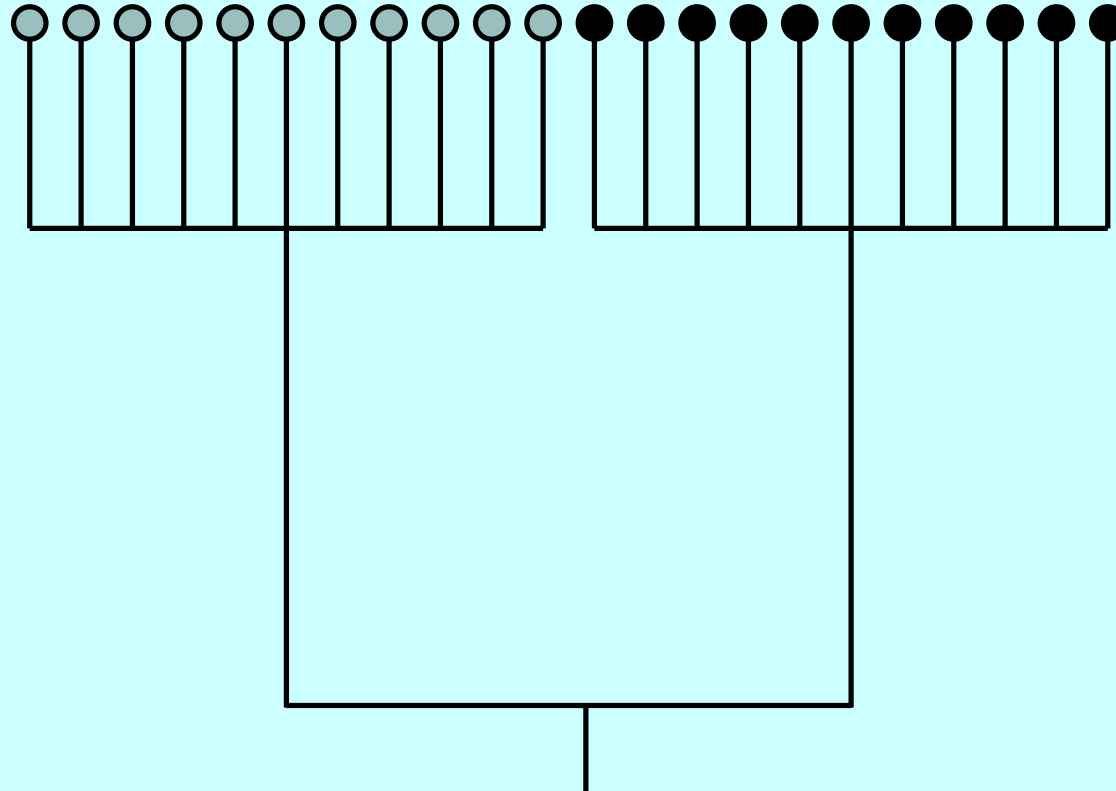


$$\delta = \frac{v_1 v_2}{v_1 + v_2}$$

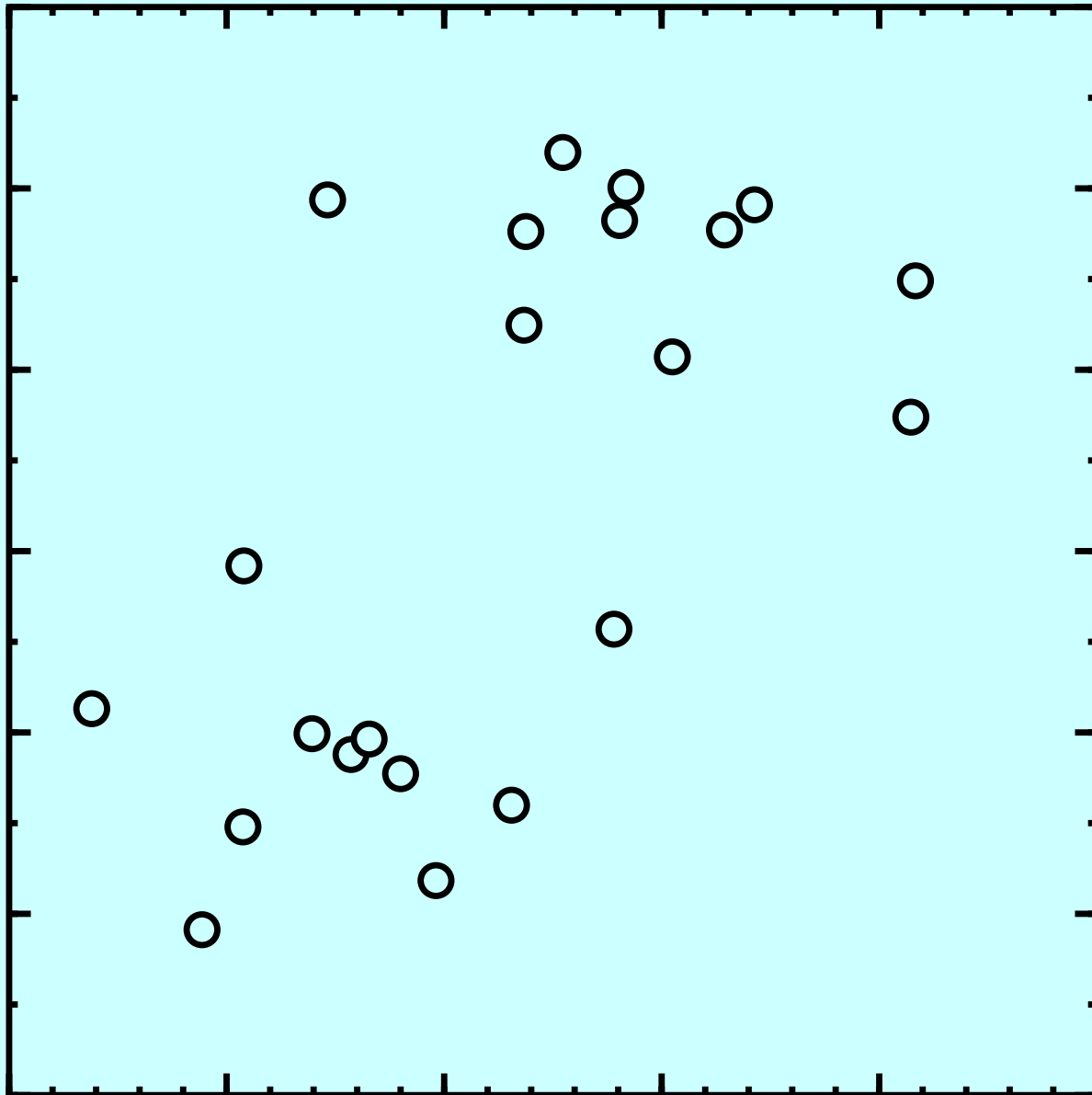
$$x_{12} = \frac{v_2 x_1 + v_1 x_2}{v_1 + v_2}$$

(True in the sense that the log-likelihoods – which are a bit different than the usual likelihoods – add up, since the likelihoods multiply).

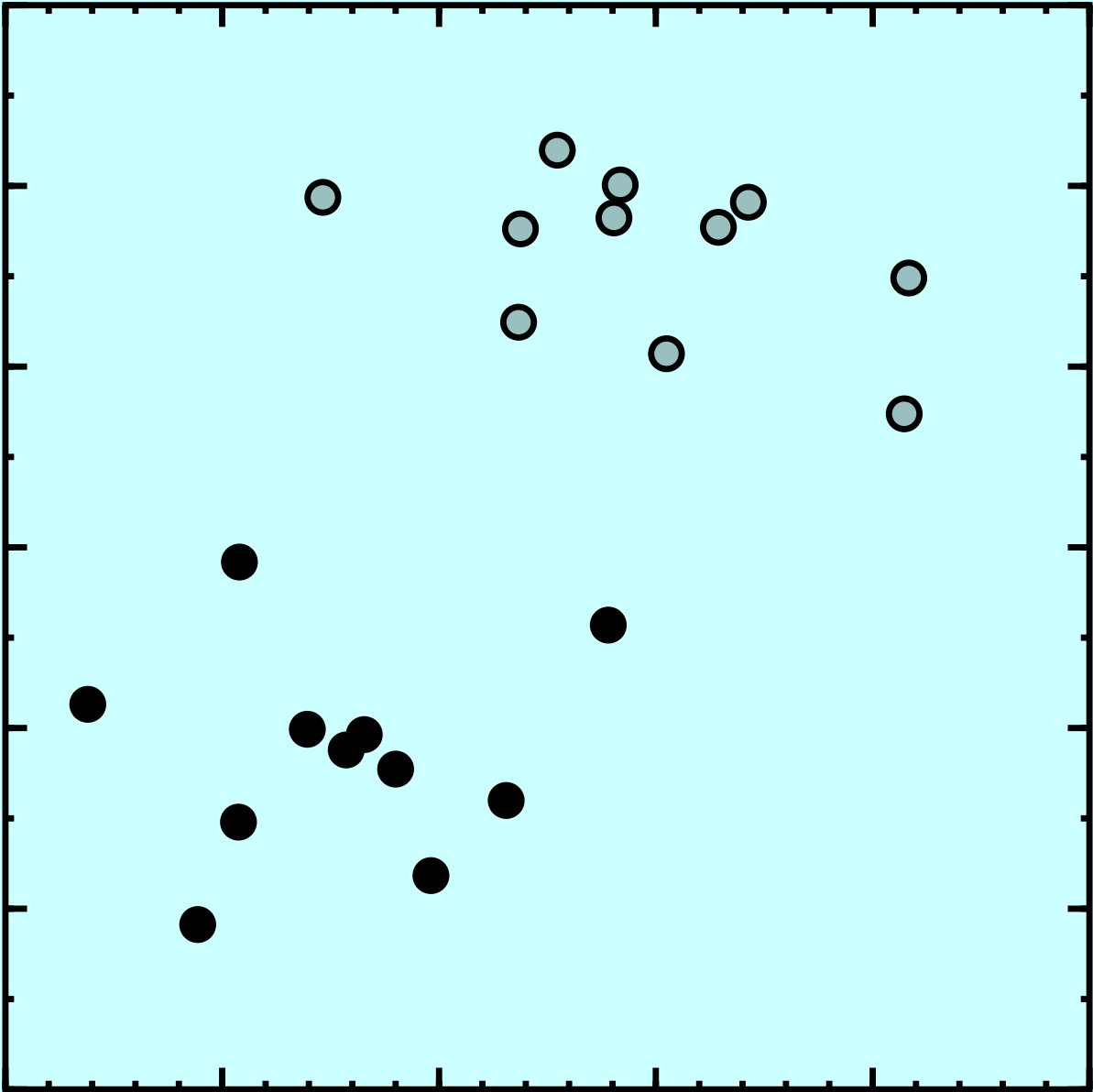
A simple case to show effects of phylogeny



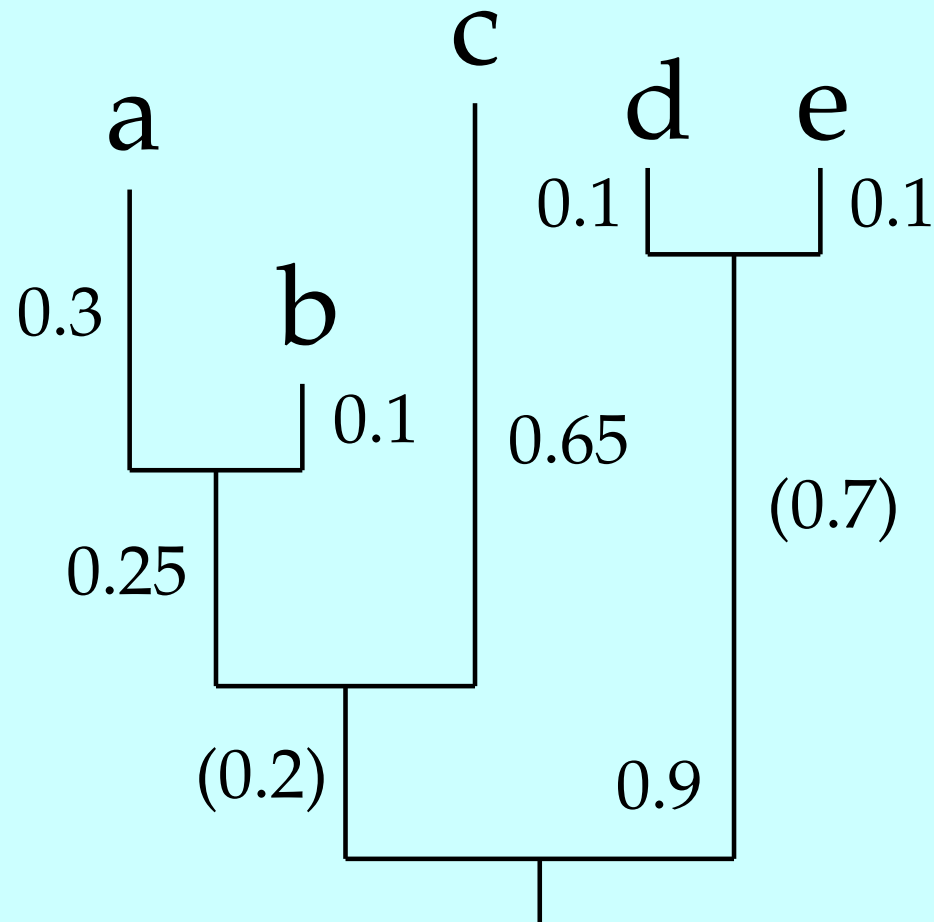
Two uncorrelated characters evolving on that tree



Identifying the two clades



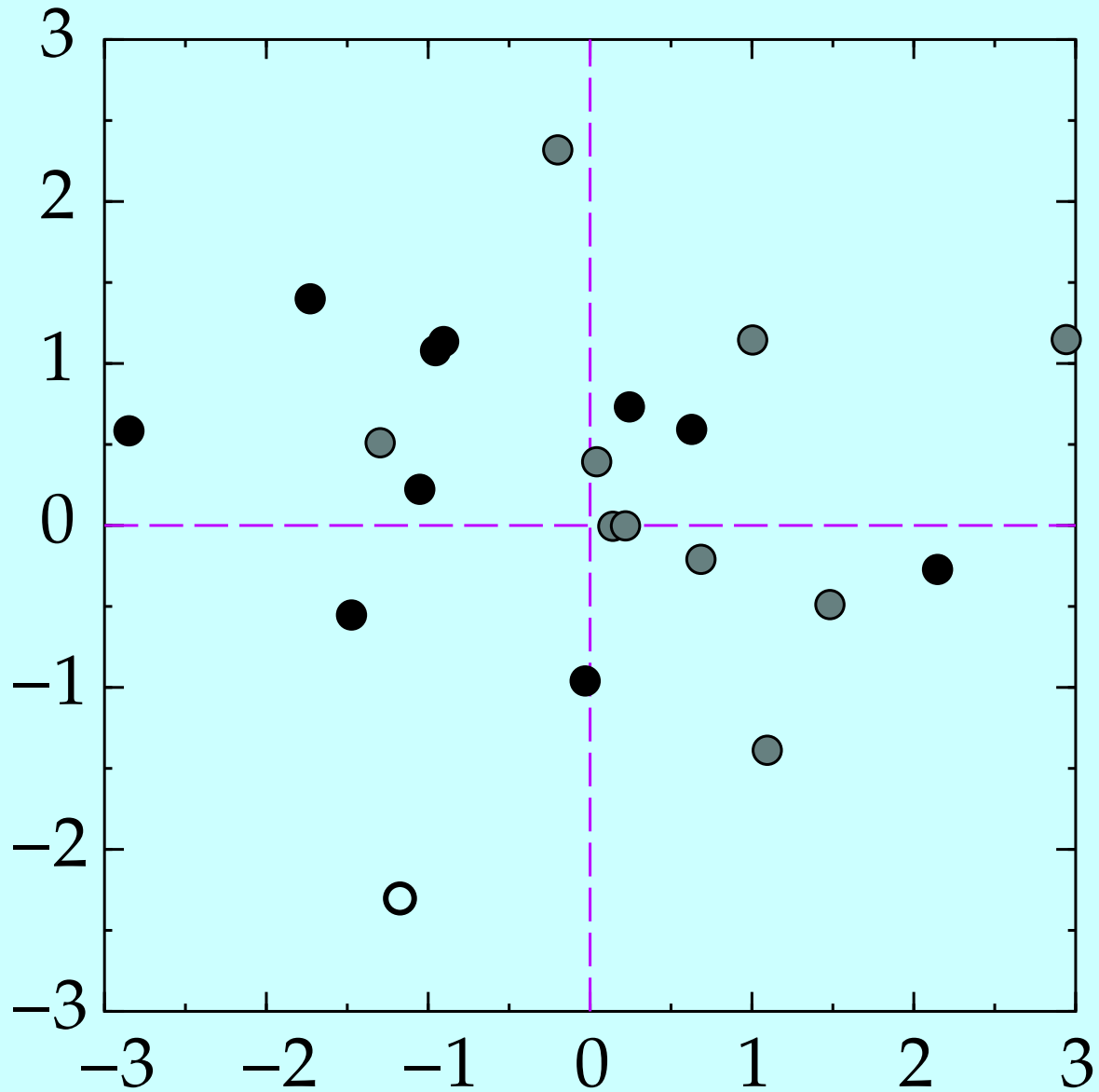
A tree on which we are to observe two characters



Contrasts on that tree

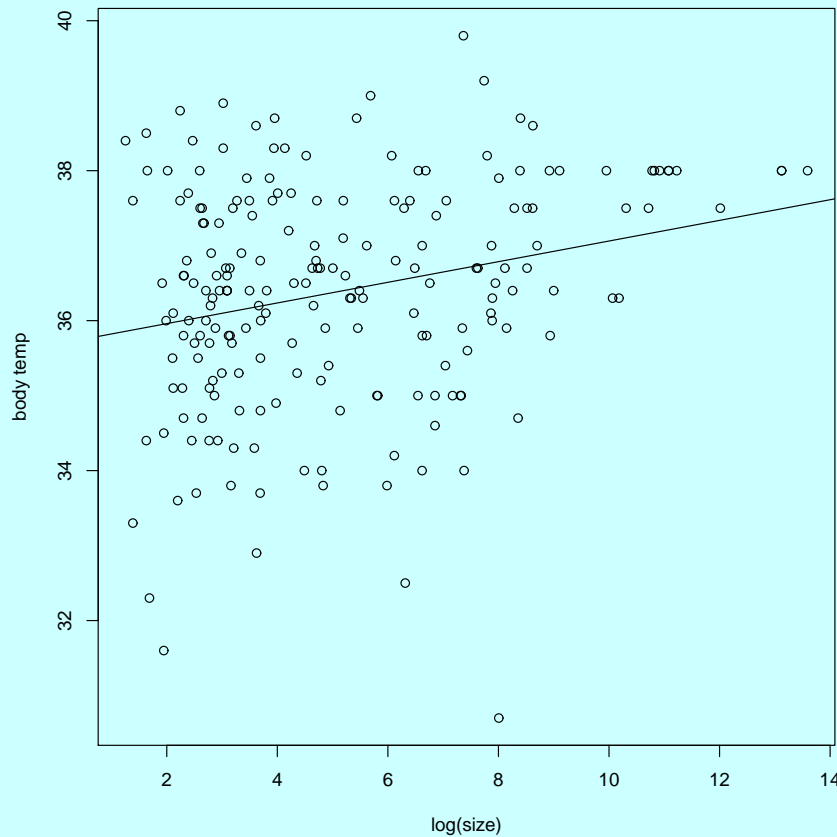
	Contrast	Variance proportional to
y_1	$= x_a - x_b$	0.4
y_2	$= \frac{1}{4} x_a + \frac{3}{4} x_b - x_c$	0.975
y_3	$= x_d - x_e$	0.2
y_4	$= \frac{1}{6} x_a + \frac{1}{2} x_b + \frac{1}{3} x_c - \frac{1}{2} x_d - \frac{1}{2} x_e$	1.11666

Contrasts for the 20-species two-clade example

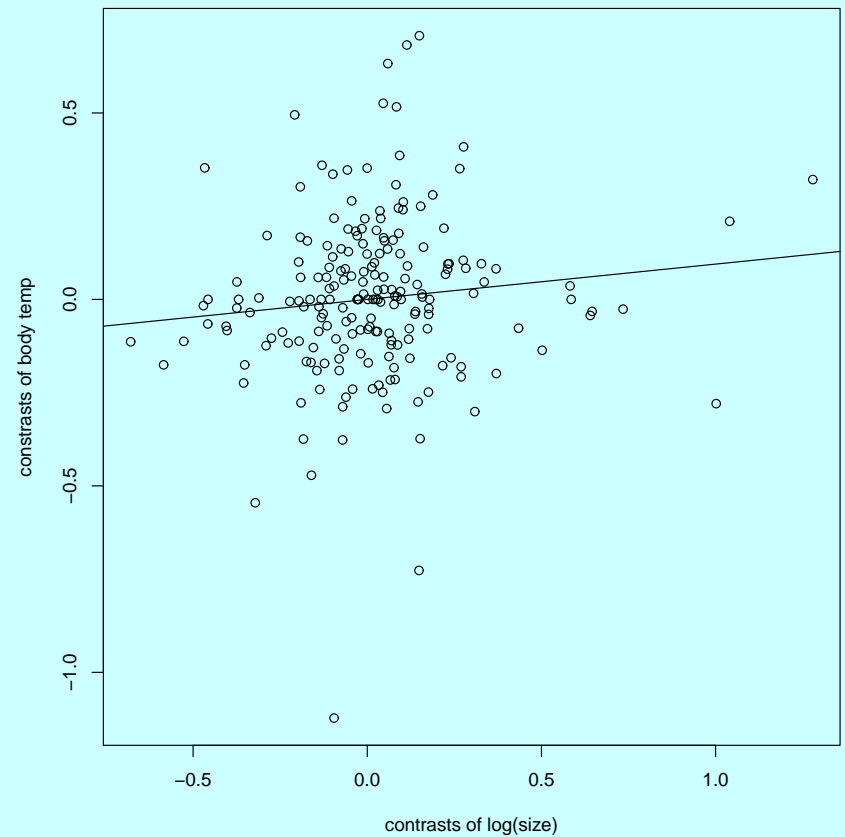


An example: Riek and Geiser, 2013

Alexander Riek and Fritz Geiser. 2013. Allometry of thermal variables in mammals: consequences of body size and phylogeny. *Biological Reviews* 88 (3): 564-572.



body temperature vs. log(body size)
(P for slope $\neq 0$ is 0.000375)



contrasts vs. contrasts
(P for slope $\neq 0$ is 0.116)

When the tree is noisy: Propagating bootstrap sampling

morphological
data

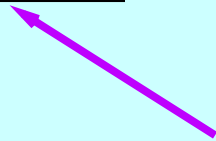
molecular
dataset

Propagating bootstrap sampling

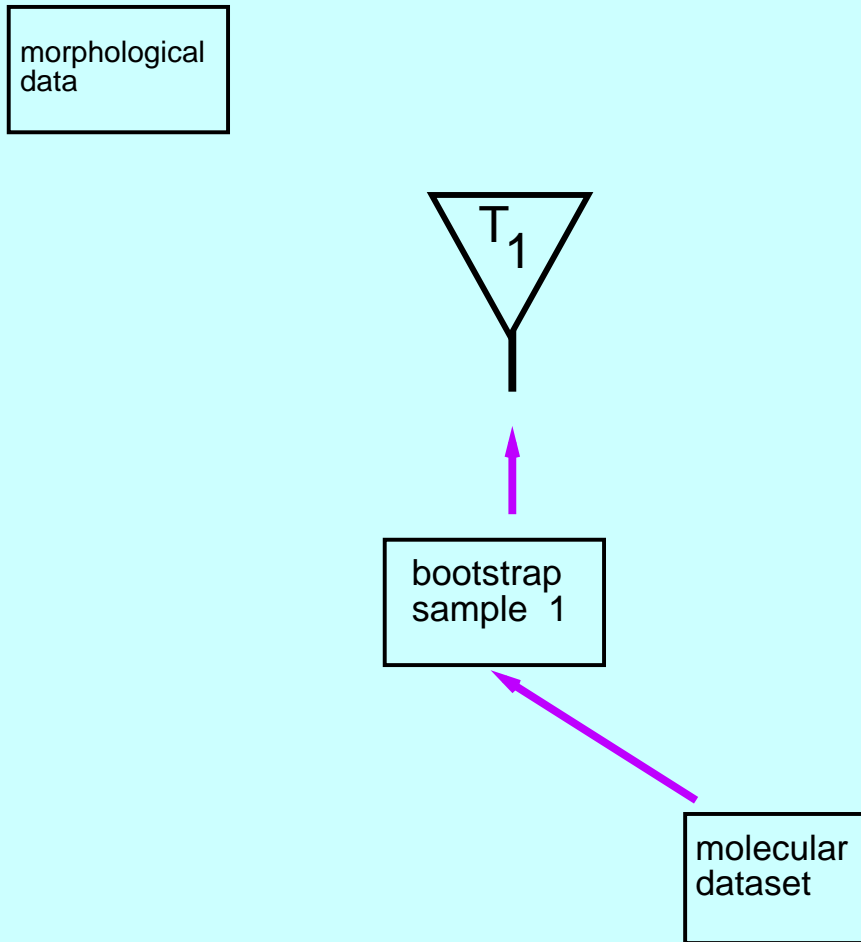
morphological
data

bootstrap
sample 1

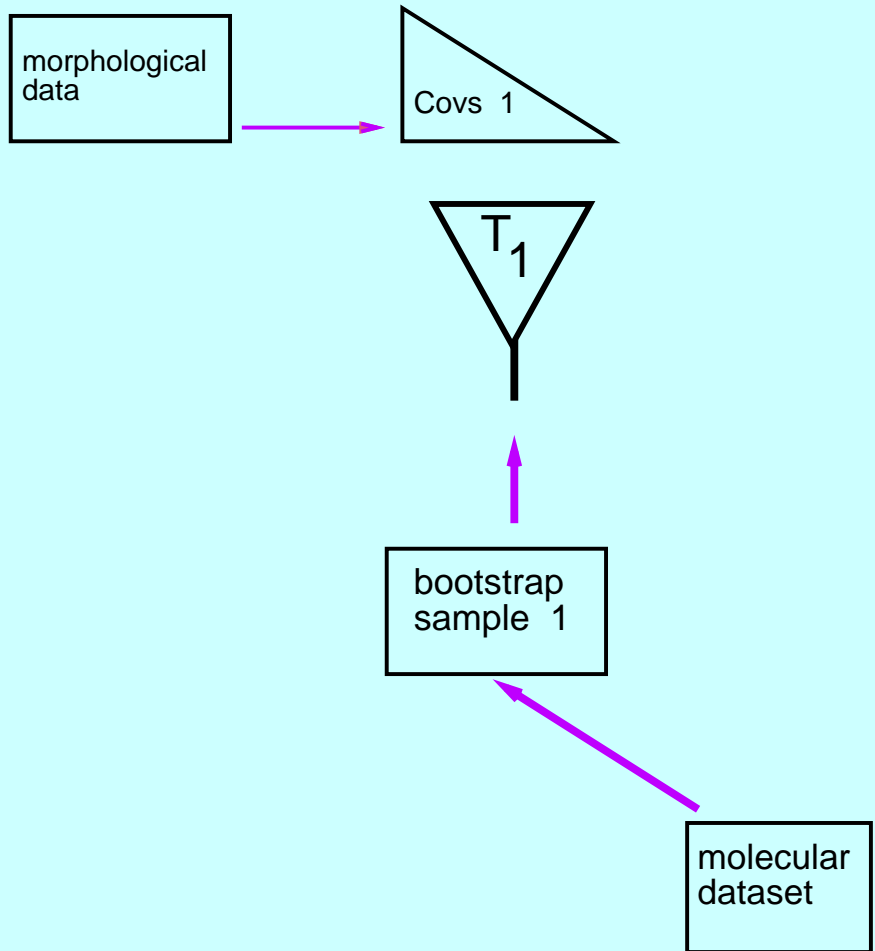
molecular
dataset



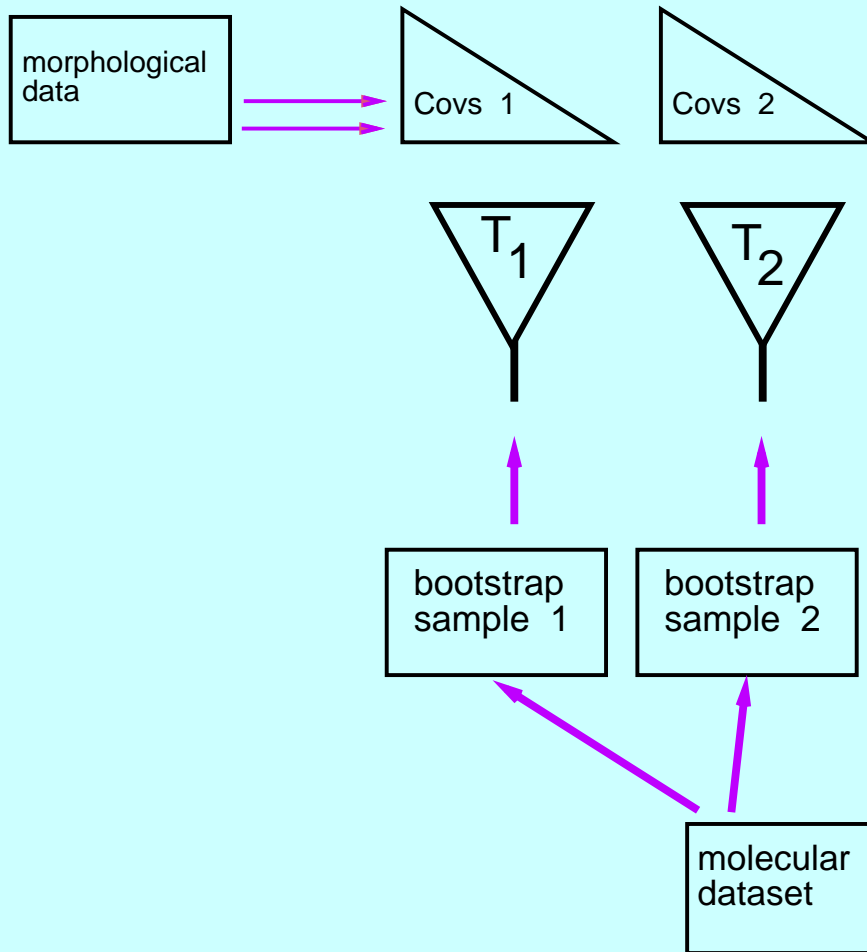
Propagating bootstrap sampling



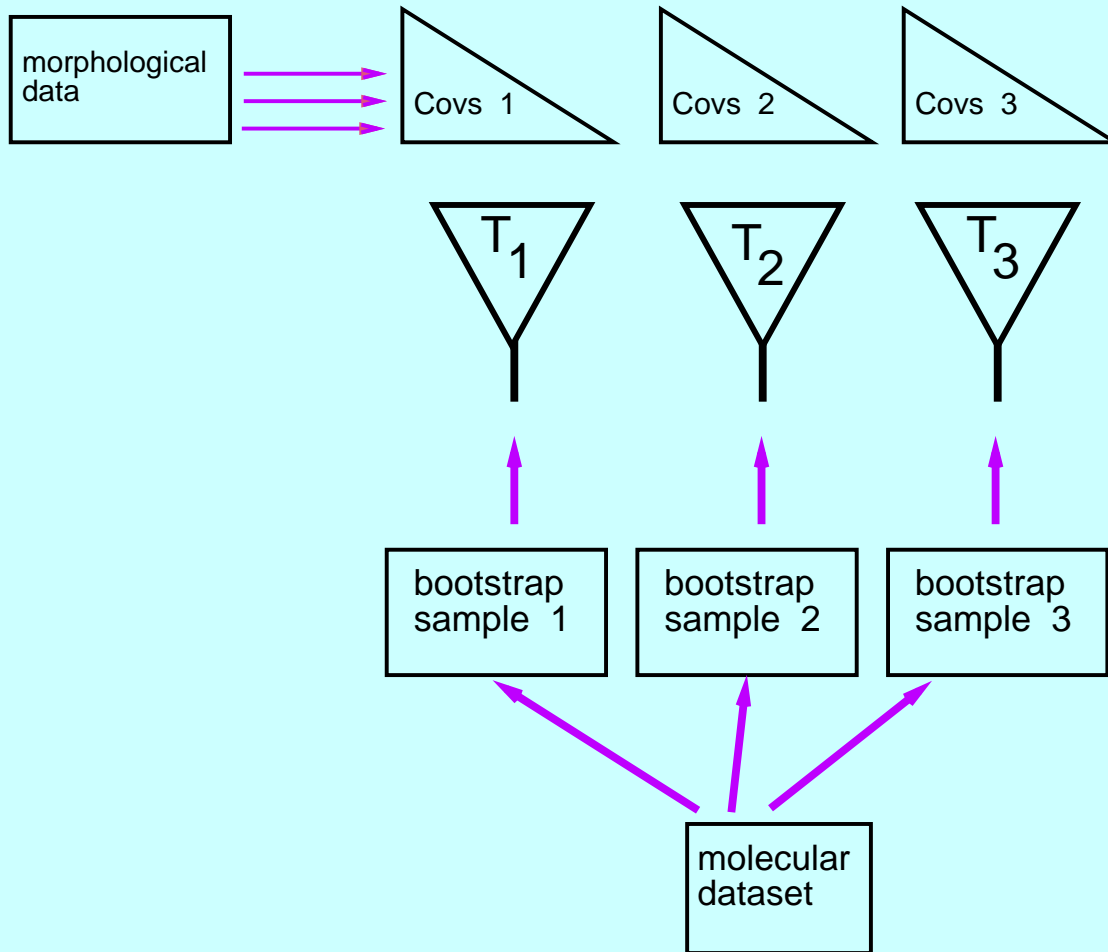
Propagating bootstrap sampling



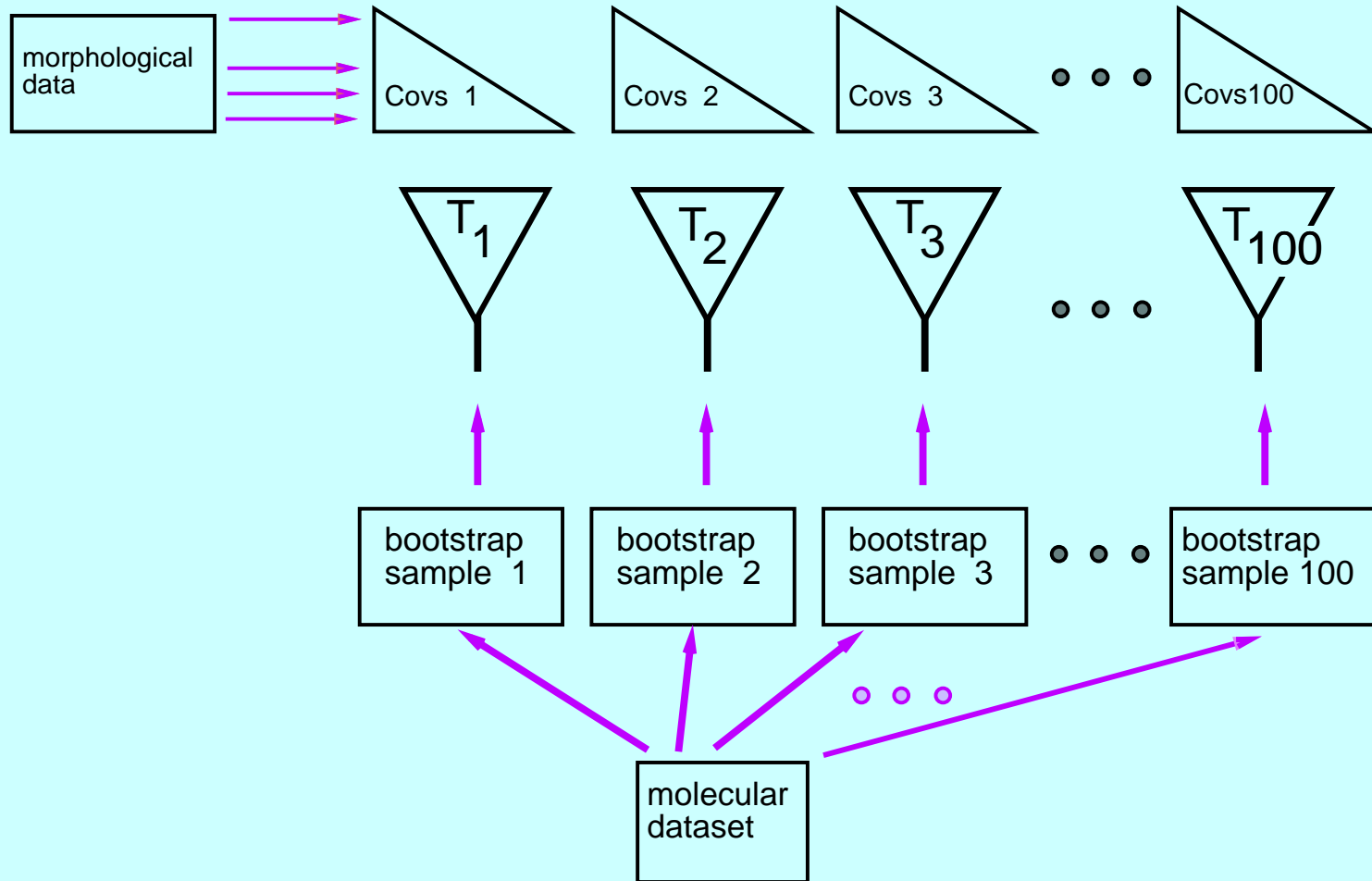
Propagating bootstrap sampling



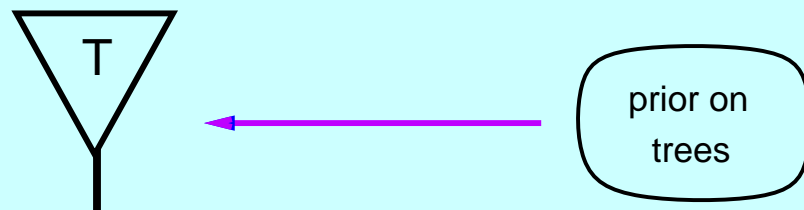
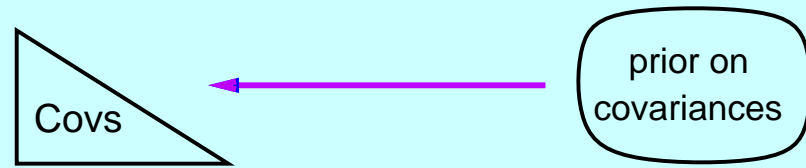
Propagating bootstrap sampling



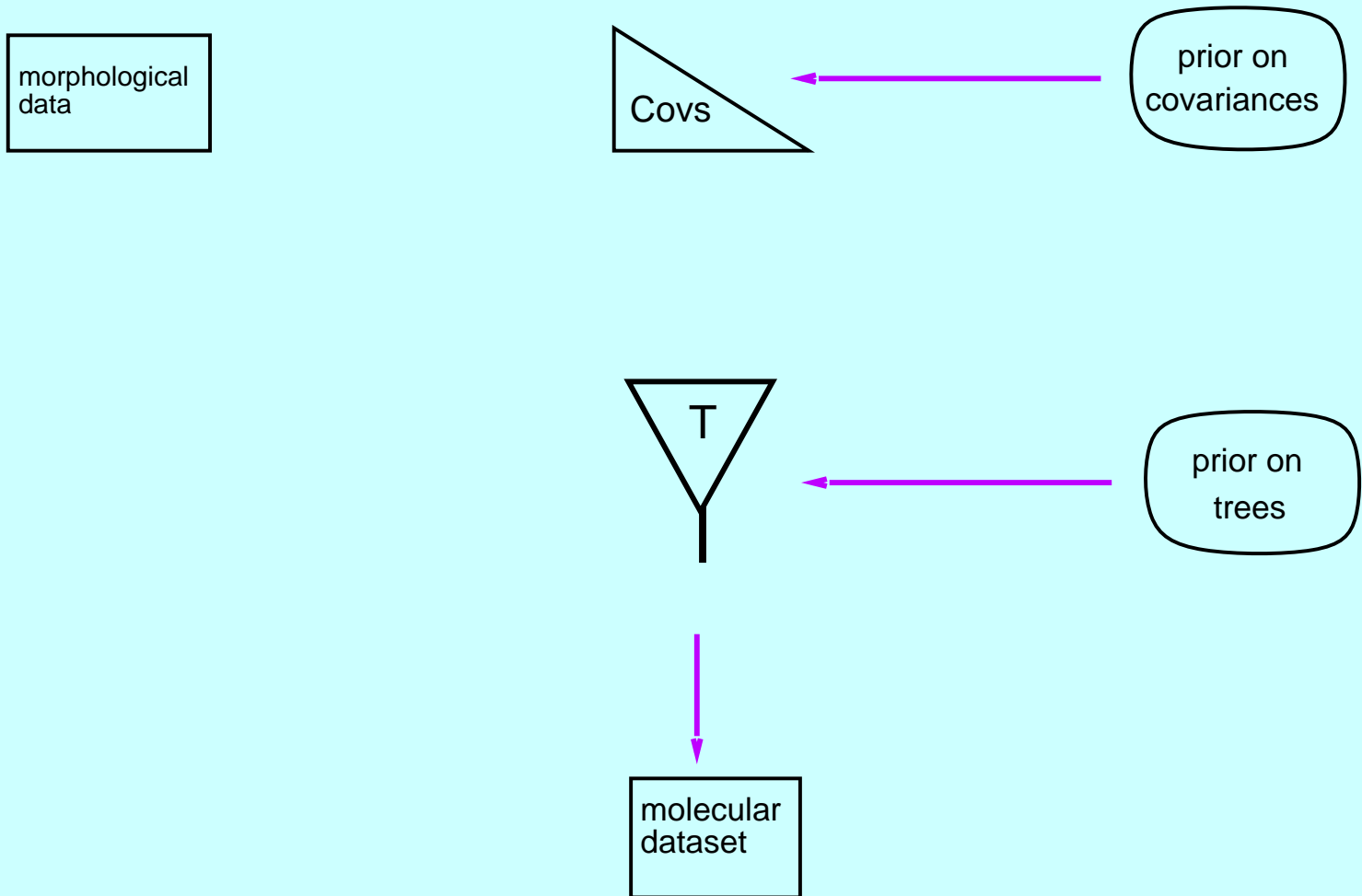
Propagating bootstrap sampling



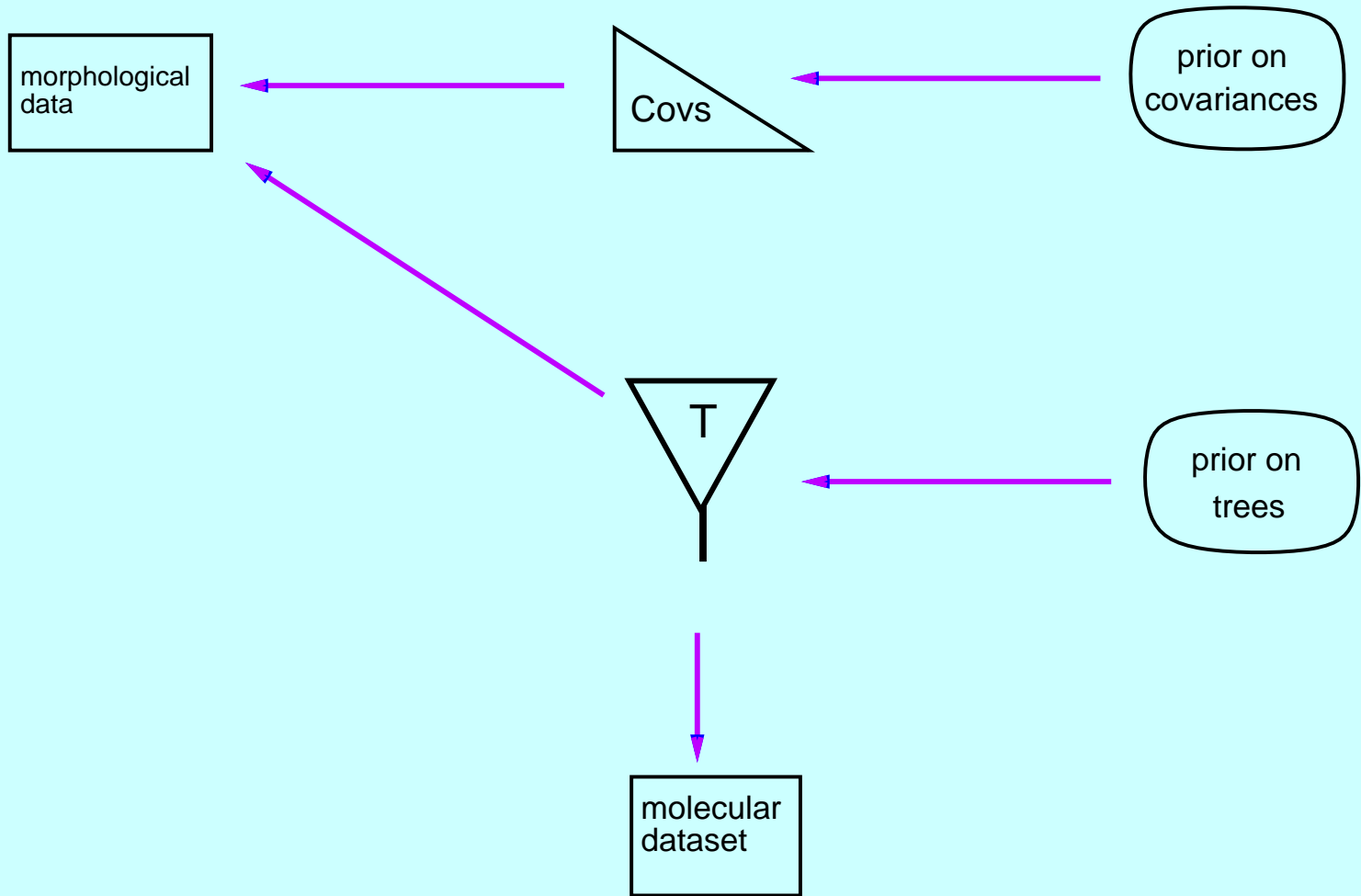
A Bayesian model



A Bayesian model

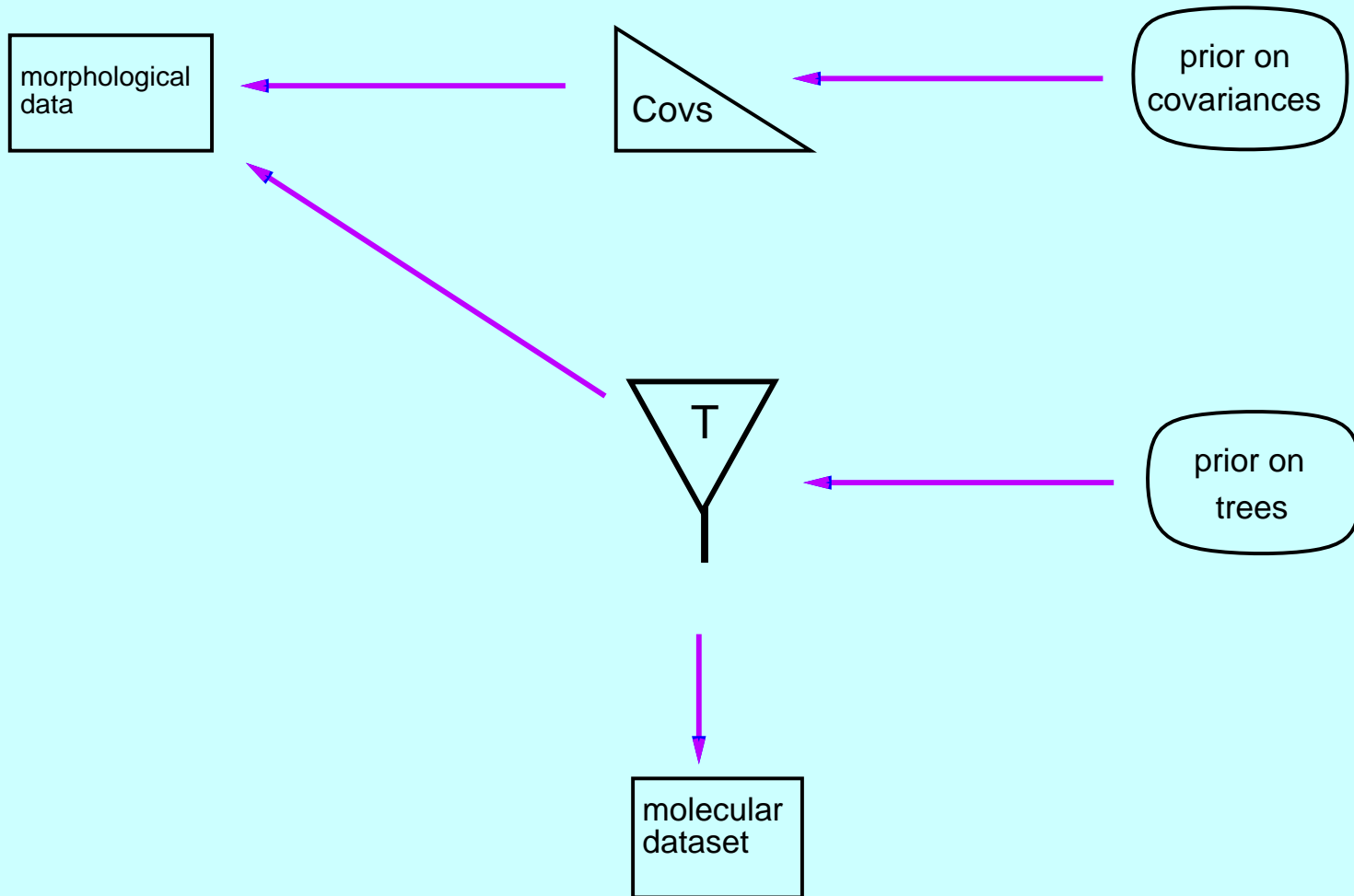


A Bayesian model

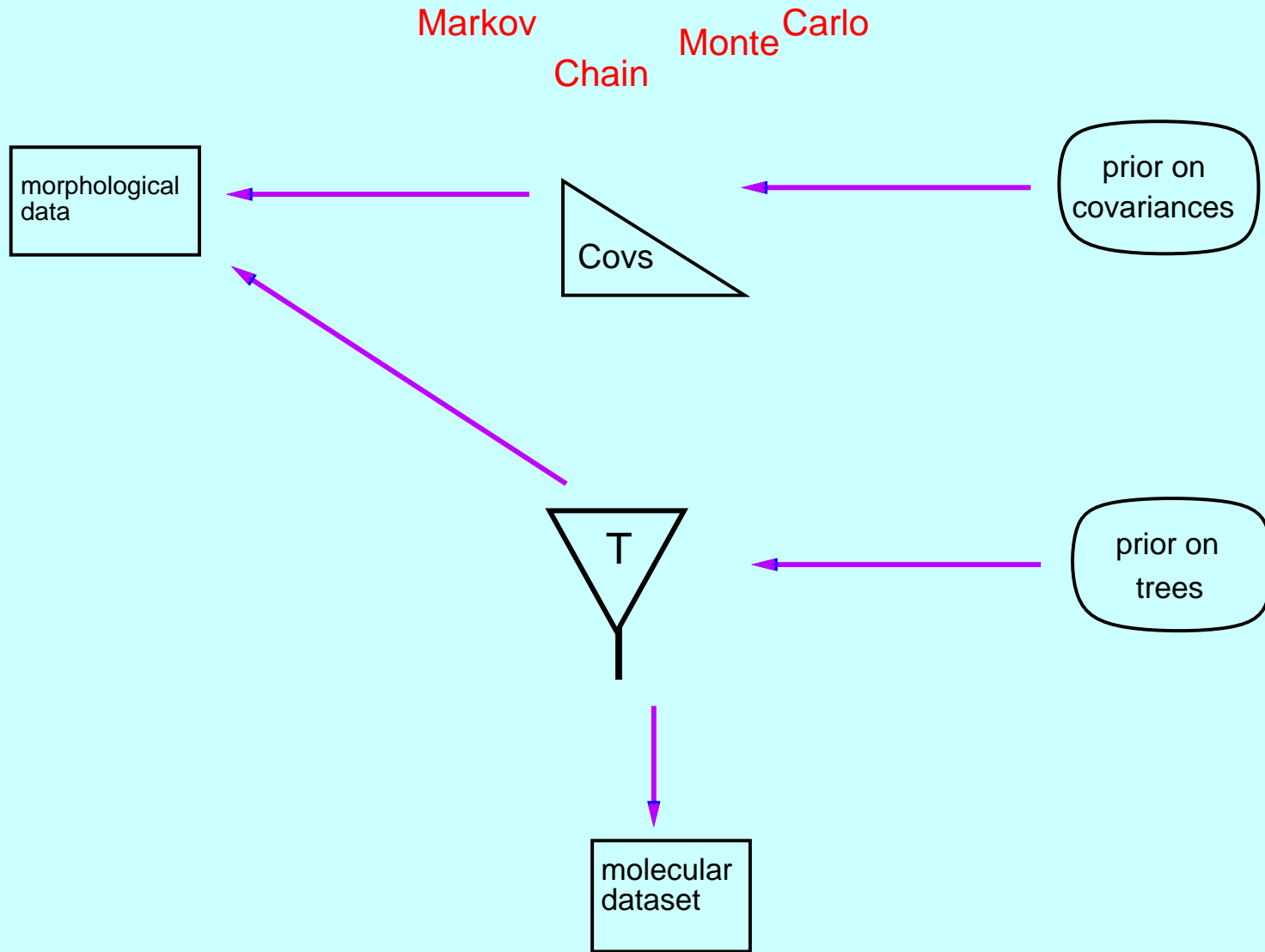


Bayesian MCMC

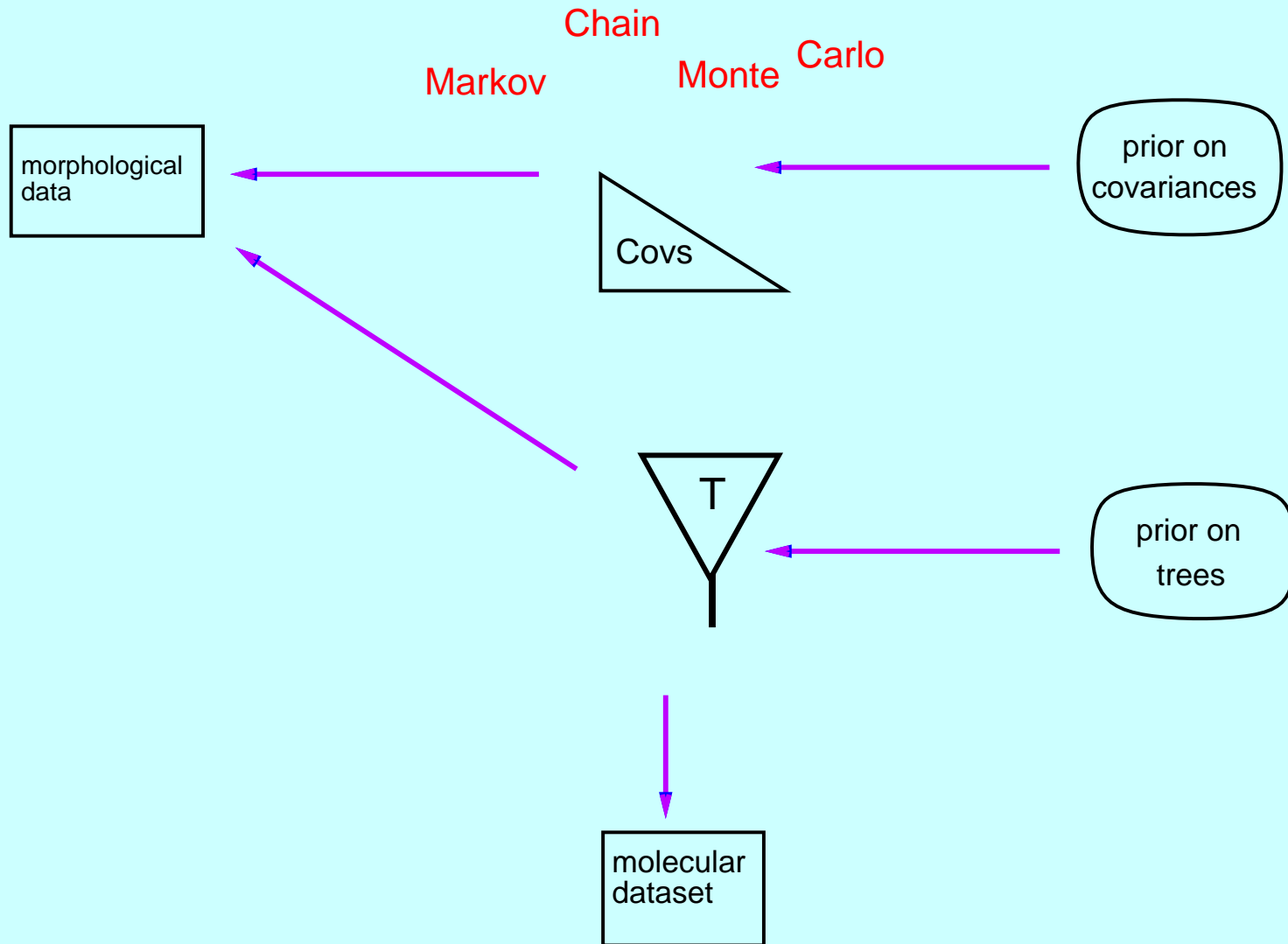
Markov Chain Monte Carlo



Bayesian MCMC

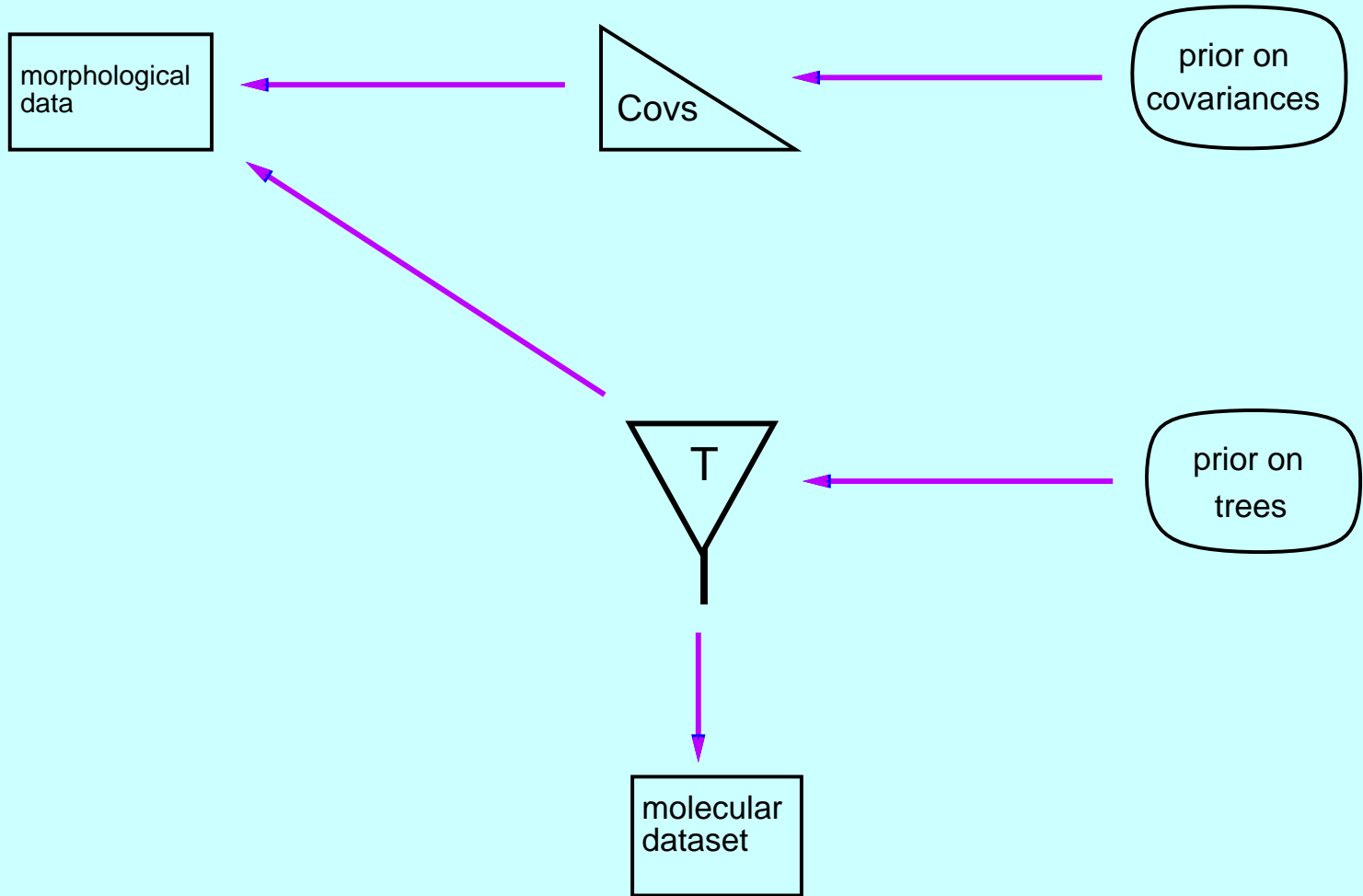


Bayesian MCMC



Bayesian MCMC

Markov Chain Monte Carlo



Some complications

- (As noted above) dealing with uncertainty about the phylogeny

Some complications

- (As noted above) dealing with uncertainty about the phylogeny
- Small sample size from species means their species means are uncertain. Must use a model with another level of variation – within-species phenotypic variation (Ricklefs and Starck, 1996; Ives et al., 2007; Felsenstein, 2008)

Some complications

- (As noted above) dealing with uncertainty about the phylogeny
- Small sample size from species means their species means are uncertain. Must use a model with another level of variation – within-species phenotypic variation (Ricklefs and Starck, 1996; Ives et al., 2007; Felsenstein, 2008)
- Rate of change of morphological characters need not be constant on the molecular tree branch lengths.

Some complications

- (As noted above) dealing with uncertainty about the phylogeny
- Small sample size from species means their species means are uncertain. Must use a model with another level of variation – within-species phenotypic variation (Ricklefs and Starck, 1996; Ives et al., 2007; Felsenstein, 2008)
- Rate of change of morphological characters need not be constant on the molecular tree branch lengths.
- Note – regressions involving contrasts should assume that they all have expectation zero.

Some complications

- (As noted above) dealing with uncertainty about the phylogeny
- Small sample size from species means their species means are uncertain. Must use a model with another level of variation – within-species phenotypic variation (Ricklefs and Starck, 1996; Ives et al., 2007; Felsenstein, 2008)
- Rate of change of morphological characters need not be constant on the molecular tree branch lengths.
- Note – regressions involving contrasts should assume that they all have expectation zero.
- How to infer the effect of an environmental variable when only its present-day values are known but not its values when the past changes were occurring? (note: regressing on the present-day values is generally **wrong**, see paper by Hansen and Bartoszek, *Systematic Biology*, 2012).

Some complications

- (As noted above) dealing with uncertainty about the phylogeny
- Small sample size from species means their species means are uncertain. Must use a model with another level of variation – within-species phenotypic variation (Ricklefs and Starck, 1996; Ives et al., 2007; Felsenstein, 2008)
- Rate of change of morphological characters need not be constant on the molecular tree branch lengths.
- Note – regressions involving contrasts should assume that they all have expectation zero.
- How to infer the effect of an environmental variable when only its present-day values are known but not its values when the past changes were occurring? (note: regressing on the present-day values is generally **wrong**, see paper by Hansen and Bartoszek, *Systematic Biology*, 2012).
- Might be able to assume environment does Brownian motion and infer covariances. Less reason to assume environment does Brownian motion than for characters.

Poor inference of covariation – what to do with that?

- Covariances are hard to infer with only (say) 50 species sampled

Poor inference of covariation – what to do with that?

- Covariances are hard to infer with only (say) 50 species sampled
- ... particularly if they samples are not independent but on a tree

Poor inference of covariation – what to do with that?

- Covariances are hard to infer with only (say) 50 species sampled
- ... particularly if they samples are not independent but on a tree
- ... particularly if the quantitative characters are thresholded

Poor inference of covariation – what to do with that?

- Covariances are hard to infer with only (say) 50 species sampled
- ... particularly if they samples are not independent but on a tree
- ... particularly if the quantitative characters are thresholded
- How do we propagate the resulting uncertainty when biologists want “fly on the wall” certainty?

Poor inference of covariation – what to do with that?

- Covariances are hard to infer with only (say) 50 species sampled
- ... particularly if they samples are not independent but on a tree
- ... particularly if the quantitative characters are thresholded
- How do we propagate the resulting uncertainty when biologists want “fly on the wall” certainty?
- Expanding to more species may put the model at risk

Poor inference of covariation – what to do with that?

- Covariances are hard to infer with only (say) 50 species sampled
- ... particularly if they samples are not independent but on a tree
- ... particularly if the quantitative characters are thresholded
- How do we propagate the resulting uncertainty when biologists want “fly on the wall” certainty?
- Expanding to more species may put the model at risk
- Expanding to more characters just adds new parameters to estimate

References for genetic drift

- Feller, W. 1951. Diffusion processes in genetics. pp. 227-246 in *Proceedings of the 2nd Berkeley Symposium on Mathematical Statistics and Probability*, ed. J. Neyman. University of California Press, Berkeley and Los Angeles.
[Feller's partial solution of the pure drift process for the Wright-Fisher model (and his famous proof that the process converges to the diffusion process)]
- Kimura, M. 1955a. Solution of a process of random genetic drift with a continuous model. *Proceedings of the National Academy of Sciences* **41**: 144-150. [Exact solution in Gegenbauer polynomials for two-allele pure genetic drift in a diffusion process approximation]
- Kimura, M. 1955b. Random drift in a multi-allelic locus. *Evolution* **9**: 419-435.
[The same, for three alleles]

References for the Brownian Motion approximation

- Edwards, A.W. F. and L. L. Cavalli-Sforza. 1964. Reconstruction of evolutionary trees. pp. 67–76 in *Phenetic and Phylogenetic Classification*, ed. V. H. Heywood and J. McNeill. Systematics Association Publ. No. 6, London. **[The first paper on numerical approaches to phylogeny reconstruction; uses parsimony and proposes likelihood for gene frequency trees]**
- Edwards, A.W. F. 1970. Estimation of the branch points of a branching diffusion process. *Journal of the Royal Statistical Society B* **32**: 155–174. **[More detailed consideration of the statistical properties of a maximum likelihood approach to gene frequency phylogenies]**
- Felsenstein, J. 1973. Maximum likelihood estimation of evolutionary trees from continuous characters. *American Journal of Human Genetics* **25**: 471–492. **[REML approach to gene frequency phylogenies, including the contrasts algorithm for rapid computation of likelihood]**
- Nielsen, R., J. L. Mountain, J. P. Huelsenbeck, and M. Slatkin. 1998. Maximum-likelihood estimation of population divergence times and population phylogeny in models without mutation. *Evolution* **52**: 669-677. **[Little-noticed but much more exact method that would require MCMC machinery]**

References on likelihood of Brownian Motion trees

- Thompson, E. A. 1975. *Human Evolutionary Trees*. Cambridge University Press, Cambridge [Thesis monograph on how to infer ML phylogenies from gene frequencies, published because it won a Smith's Prize at Cambridge University]
- Felsenstein, J. 1981. Maximum likelihood estimation of evolutionary trees from continuous characters. *Evolution* 25: 471–492. [Reworks the 1973 paper with more care and some additional algorithmics, including discussion of effect of character covariation]
- Felsenstein, J. 1985. Phylogenies from gene frequencies: A statistical problem. *Systematic Zoology* 34: 300–311. [Shows how gene frequency changes depart from being approximated by Brownian Motion]

References for multivariate Brownian motion

- Felsenstein, J. 1988. Phylogenies and quantitative characters. *Annual Review of Ecology and Systematics* **19**: 445-471. **[Review with mention of usefulness of threshold model]**
- Felsenstein, J. 2002. Quantitative characters, phylogenies, and morphometrics. pp. 27-44 in *Morphology, Shape, and Phylogenetics*, ed. N. MacLeod. Systematics Association Special Volume Series 64. Taylor and Francis, London. **[Review repeating 1988 material and going into some more detail on the question of threshold models.]**
- Felsenstein, J. 2004. *Inferring Phylogenies*. Sinauer Associates, Sunderland, Massachusetts. **[See particularly Chapters 23–25. Mentions issues with multivariate models, and also sample size issues in contrasts method].**
- Lande, R. 1976. Natural selection and random genetic drift in phenotypic evolution. *Evolution* **30**: 314-334. **[Lande's classic paper on drift versus optimum selection]**
- Lande, R. 1979. The quantitative genetic analysis of multivariate evolution, applied to brain-body size allometry. *Evolution* **33**: 402-416. **[Lande on covarying characters]**

References

- Lande, R. 1980. The genetic covariance between characters maintained by pleiotropic mutations. *Genetics* **94**: 203-215. **[Lande model explaining genetic covariance as outcome of mutational covariance and selection]**
- Lynch, M. and W. G. Hill. 1986. Phenotypic evolution by neutral mutation. *Evolution* **40**: 915-935. **[What happens if a quantitative character is selectively neutral?]**
- Stebbins, G. L. 1950. *Variation and Evolution in Plants*. Columbia University Press, New York. **[Describes selective covariance and cites Tedin (1925) for it]**
- Tedin, O. 1925. Vererbung, Variation, und Systematik der Gattung *Camelina*. *Hereditas* **6**: 275-386. **[Original paper noting the existence of selective covariance]**
- Armbruster, W. S. 1996. Causes of covariation of phenotypic traits among populations. *Journal of Evolutionary Biology* **9**: 261-276. **[Good exposition of selective covariance]**

References for phylogenetic comparative methods

- Felsenstein, J. 1985. Phylogenies and the comparative method. *American Naturalist* **125**: 1–5. [Introduces the contrasts method]
- Felsenstein, J. 1988. Phylogenies and quantitative characters. *Annual Review of Ecology and Systematics* [Suggests using bootstrapping to correct comparative methods for uncertainty about the phylogeny] **19**: 445–471.
- Harvey, P. H. and M. D. Pagel. 1991. *The Comparative Method in Evolutionary Biology*. Oxford University Press, Oxford. [The major book introducing statistical phylogenetic comparative methods]
- Grafen, A. 1989. The phylogenetic regression. *Philosophical Transactions of the Royal Society of London, Series B* **326**: 119–157. [Using generalized least squares to evaluate the likelihood for Brownian Motion phylogenies and do comparative methods analysis, without the contrasts methods. In the simplest case, is exactly equivalent to the contrasts method. Discusses ways of coping with unresolved parts of the phylogeny and with varying evolutionary rates.]

References, continued

- Ricklefs, R. E. and J. M. Starck. 1996. Applications of phylogenetically independent contrasts: A mixed progress report. *Oikos* 77: 167–172. **[Pointing put that small sample size within species is a problem for comparative methods]**
- Ives, A. R., P. E. Midford, and T. Garland. 2007. Within-species variation and measurement error in phylogenetic comparative methods. *Systematic Biology* 56: 252-270. **[Taking small sample size into account when we know the within-species phenotypic covariances]**
- Hansen, T. F., and K. Bartoszek. 2012. Interpreting the evolutionary regression: the interplay between observational and biological errors in phylogenetic comparative studies. *Systematic Biology* 61(3): 413 – 425. **[Point out that it is wrong to first remove the effect of one character (or environment) and then assume residuals evolve by Brownian Motion]**
- Felsenstein, J. 2008 Comparative methods with sampling error and within-species variation: contrasts revisited and revised. *American Naturalist* 171: 713–725. **[Inferring both between-species evolutionary covariances and within-species phenotypic variation]**