# Lecture 3

# Models and support

# Bootstrapping

# The bootstrap is non parametric

- Uncertainty in the tree estimate can be inferred indirectly using **bootstrap analysis**

- "Pulling oneself up by one's bootstraps"

- Bootstrap analysis can be performed when using a range of phylogenetic methods:

  - Maximum parsimony

  - Distance-matrix based methods

  - Maximum likelihood
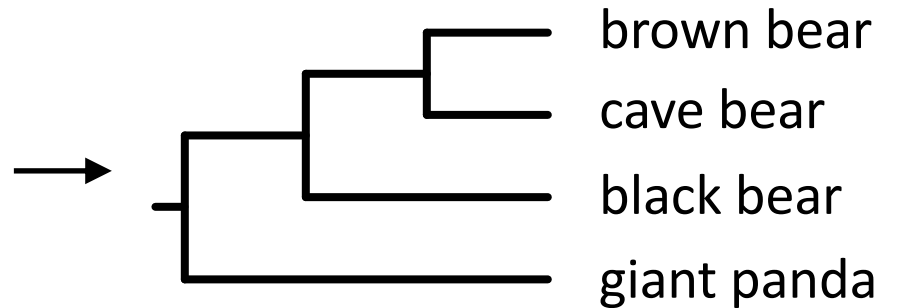
# Bootstrap

brown bear CGTTAGTACACT

cave bear CGATAGTTCACT

black bear CGTTAGTTTACC

giant panda CATTGGTTTACT

Repeat 1000 times

Pseudoreplicate

brown bear ATACTGTCCCT

cave bear ATACTGTCCCA

black bear ACACTGTTCCT

giant panda GTGCTATTCCT



brown bear
cave bear
black bear
giant panda

# Bootstrap



ML tree

0.90
- brown bear
- cave bear
- black bear
- giant panda
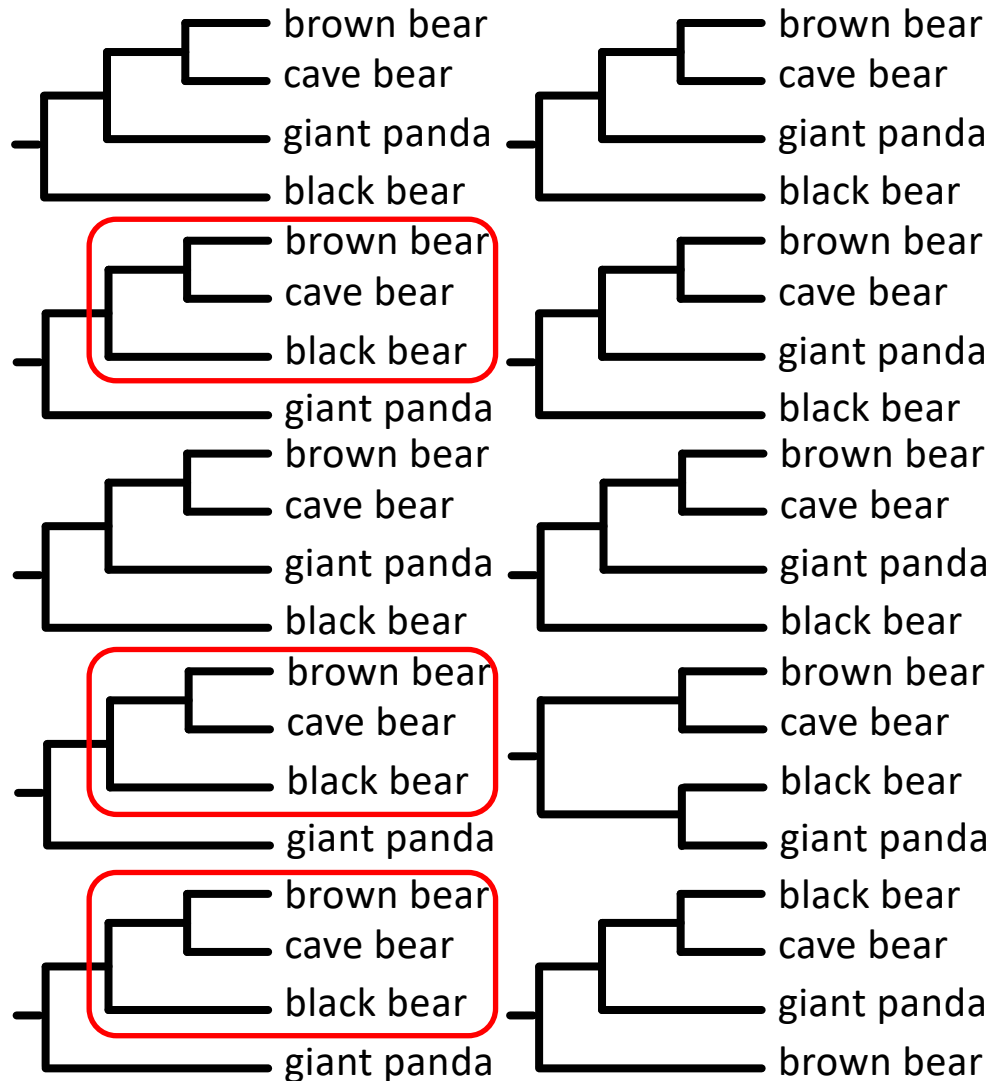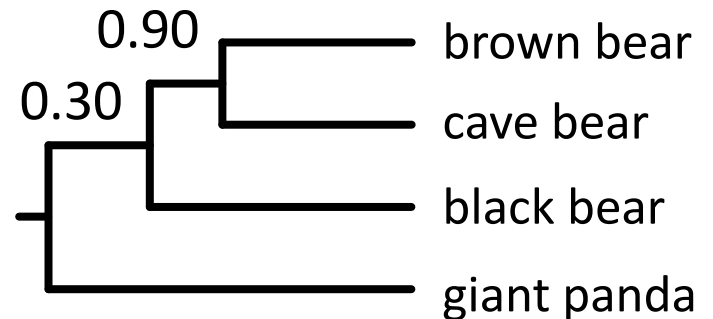
# Bootstrap



ML tree

# Interpreting bootstrap values

- **Felsenstein (1985)**
  The bootstrap gives us a confidence interval that contains *the tree that would be estimated when repeatedly sampling sites from the existing distribution*

- Bootstrap values are **measures of repeatability**

  - High when lots of data are available

  - Has little meaning when genome-scale data are available

Soltis & Soltis (2003) *Stat Sci*

# Popular methods in phylogenetics

1. Maximum parsimony
2. Distance methods
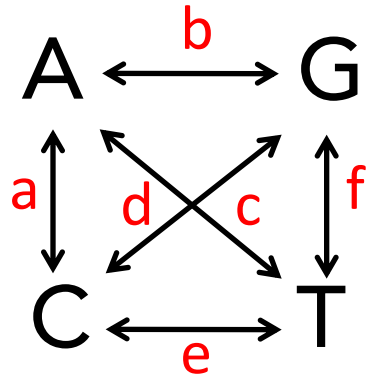3. Maximum likelihood
4. Bayesian inference

Statistical methods

# Substitution models

# DNA substitution models

## Rates matrix  Base frequencies



$$\pi_A + \pi_C + \pi_G + \pi_T = 1$$

### JC
a=b=c=d=e=f

$\pi_A = \pi_C = \pi_G = \pi_T$

### HKY
a=c=d=f, b=e

$\pi_A, \pi_C, \pi_G, \pi_T$
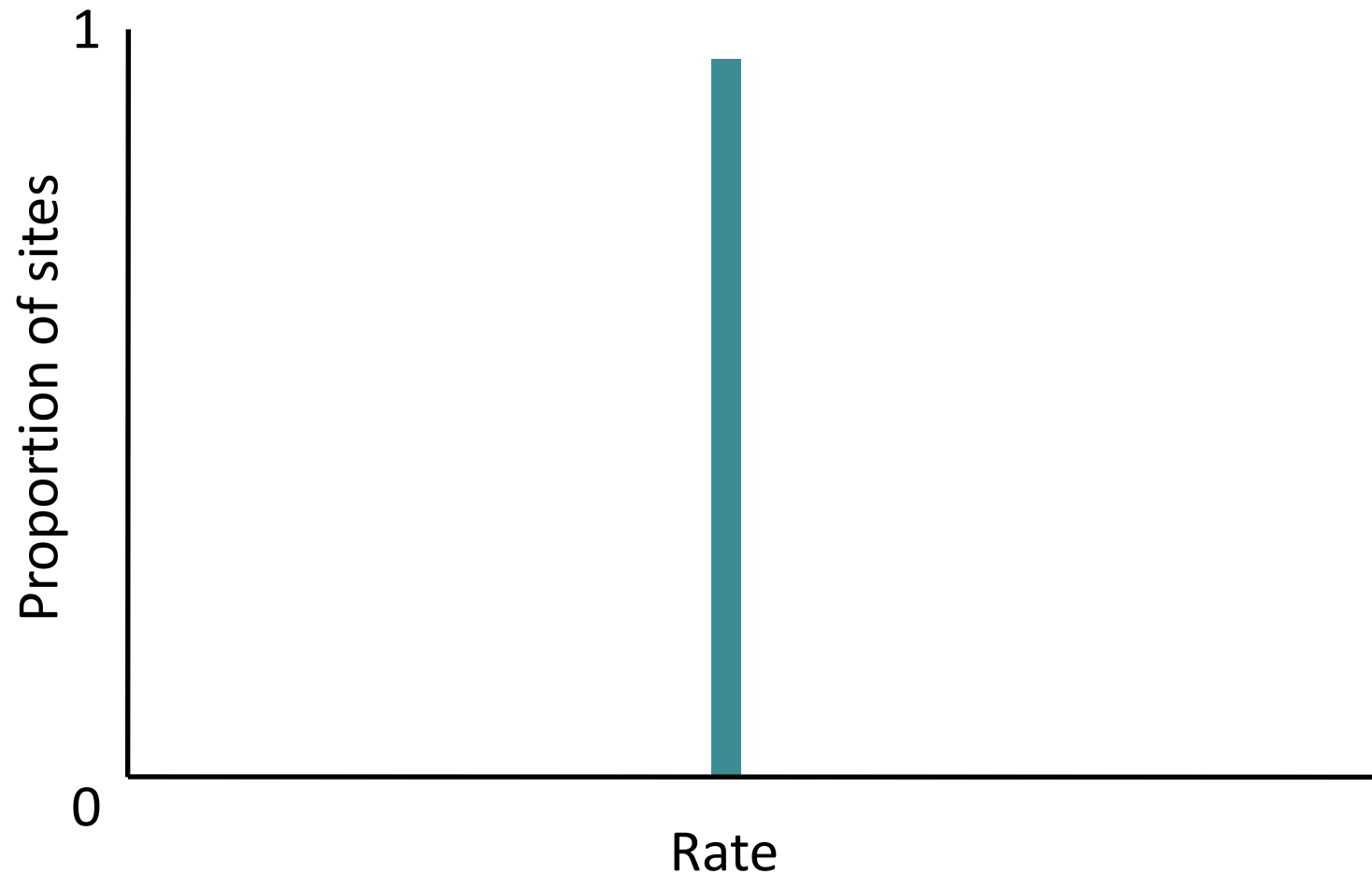
### GTR
a, b, c, d, e, f

$\pi_A, \pi_C, \pi_G, \pi_T$

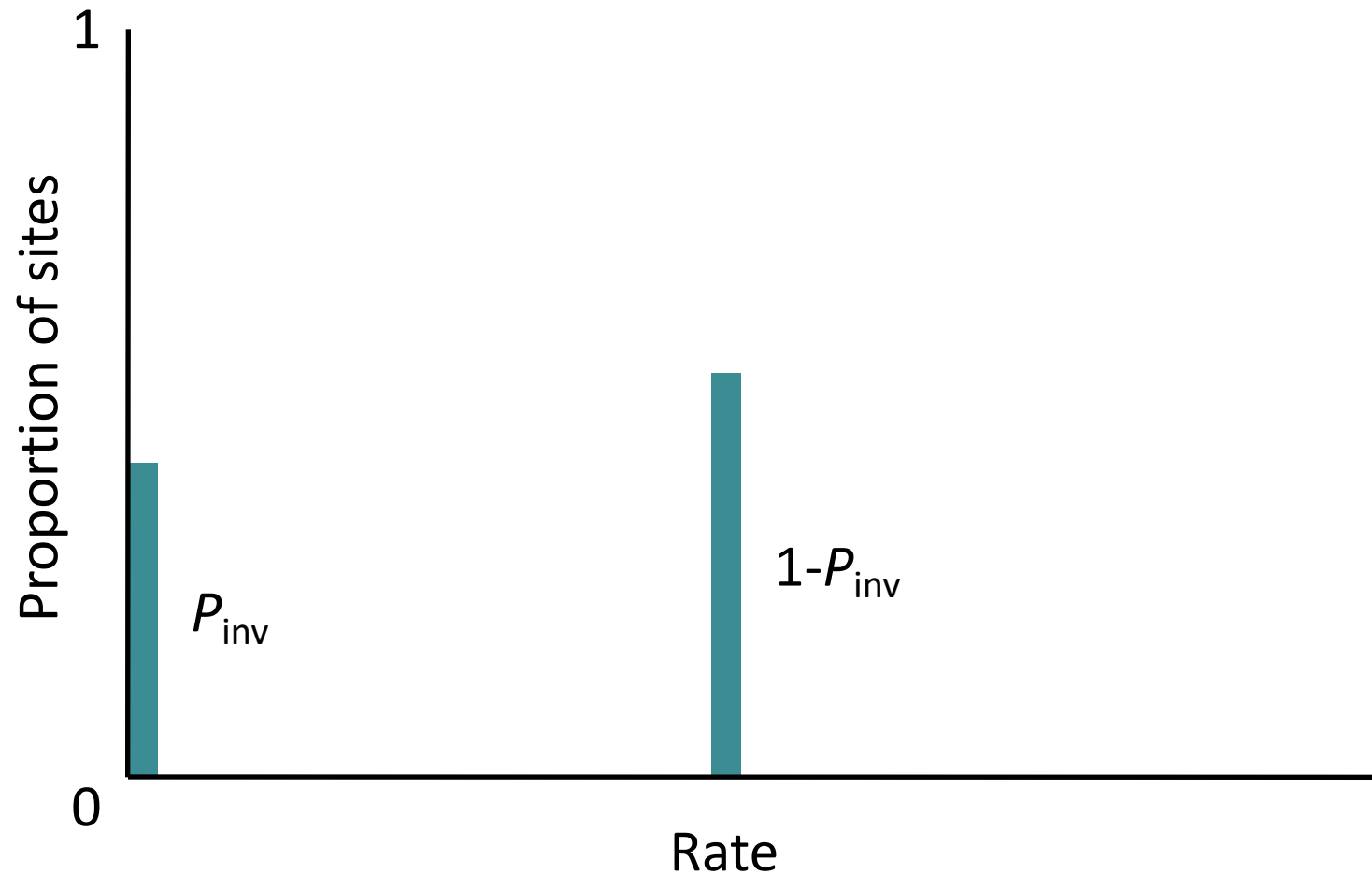# Rate variation across sites



Medium  Slow  Fast

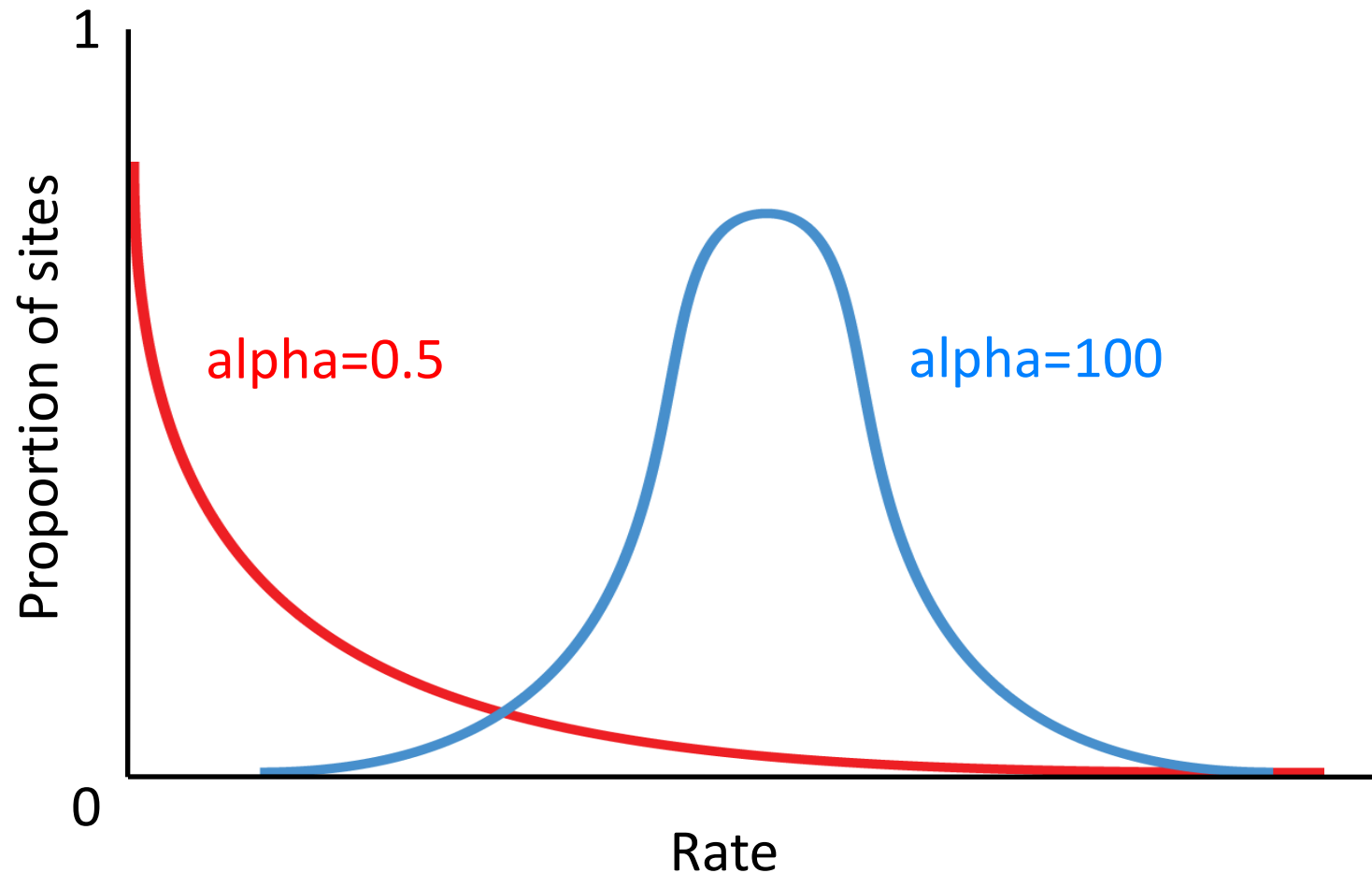# Rate variation across sites

- Identical among all sites

# Rate variation across sites

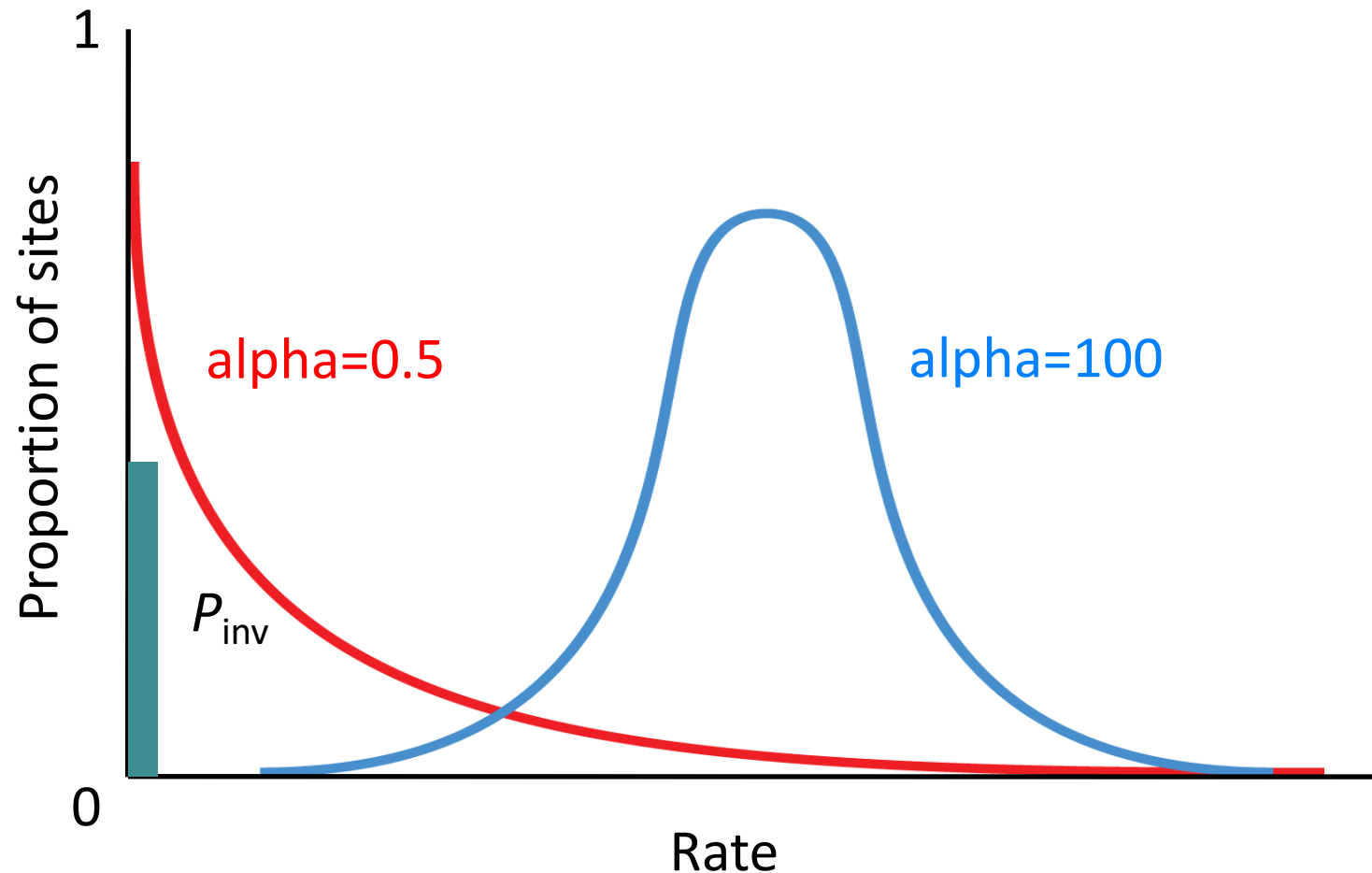- Proportion of invariable sites (**+I** models)

# Rate variation across sites

- Gamma distributed rates across sites (**+G** models)

# Rate variation across sites

- Rates across sites are assumed to follow a gamma distribution and a portion of invariable sites (**+G+I** models)
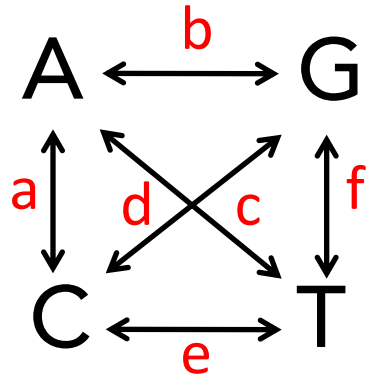


alpha=0.5

alpha=100

$P_{inv}$

Proportion of sites

Rate

# DNA substitution models

## Rates matrix



A $\xleftrightarrow{b}$ G

$a$ $d$ $c$ $f$

C $\xleftrightarrow{e}$ T

## Base frequencies

$\pi_A + \pi_C + \pi_G + \pi_T = 1$

## Site rates

$+ I + G$

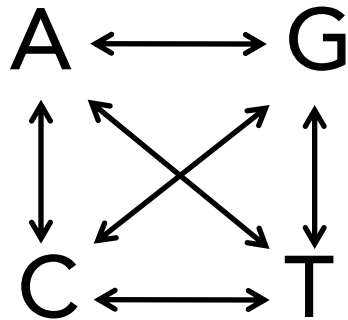| JC | HKY | GTR | GTR+I+G |
|---|---|---|---|
| a=b=c=d=e=f | a=c=d=f, b=e | a, b, c, d, e, f | a, b, c, d, e, f |
| $\pi_A=\pi_C=\pi_G=\pi_T$ | $\pi_A, \pi_C, \pi_G, \pi_T$ | $\pi_A, \pi_C, \pi_G, \pi_T$ | $\pi_A, \pi_C, \pi_G, \pi_T$ I, G |

# DNA substitution models

Rates matrix       Base frequencies       Site rates



$$\pi_A + \pi_C + \pi_G + \pi_T = 1 \qquad + I + G$$

Number of models

$$203 \quad \times \quad 15 \quad \times \quad 4 \quad = 12{,}180$$

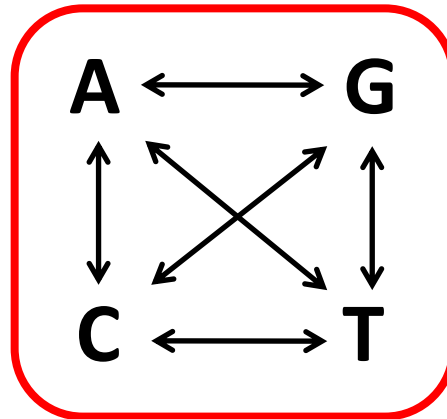In phylogenetics we explore a small portion of these

# Proportion of invariable sites

- Often over-estimated in species-level analyses

- Do not distinguish:

  - Sites that are **invariable** and cannot change

  - Sites that are **constant** and for stochastic reasons do not have any subtitutions

- Little biological meaning

- Site rates can be adequately described using **+G**

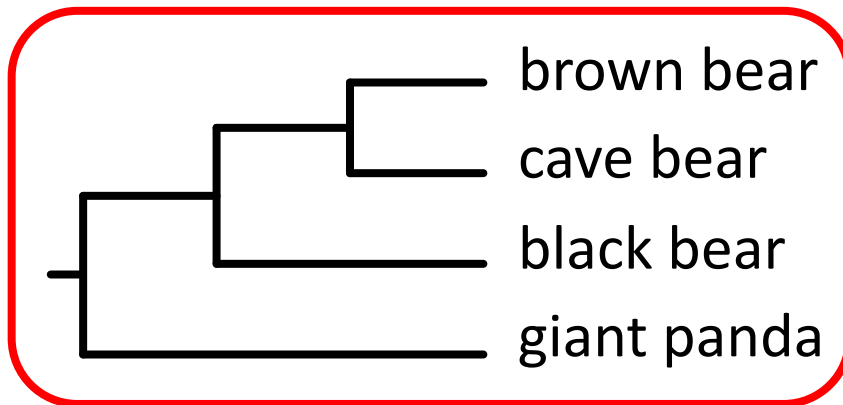We use +G models to account for variable
rates across sites

# Fundamental assumptions

# Amino acid substitution matrices

- Matrix has size 20x20

- Too many parameters to estimate
  - GTR model for DNA: 6 parameters
  - GTR model for proteins: 190 parameters

- Transition probabilities come from vast data sets
  - PAM
  - BLOSUM
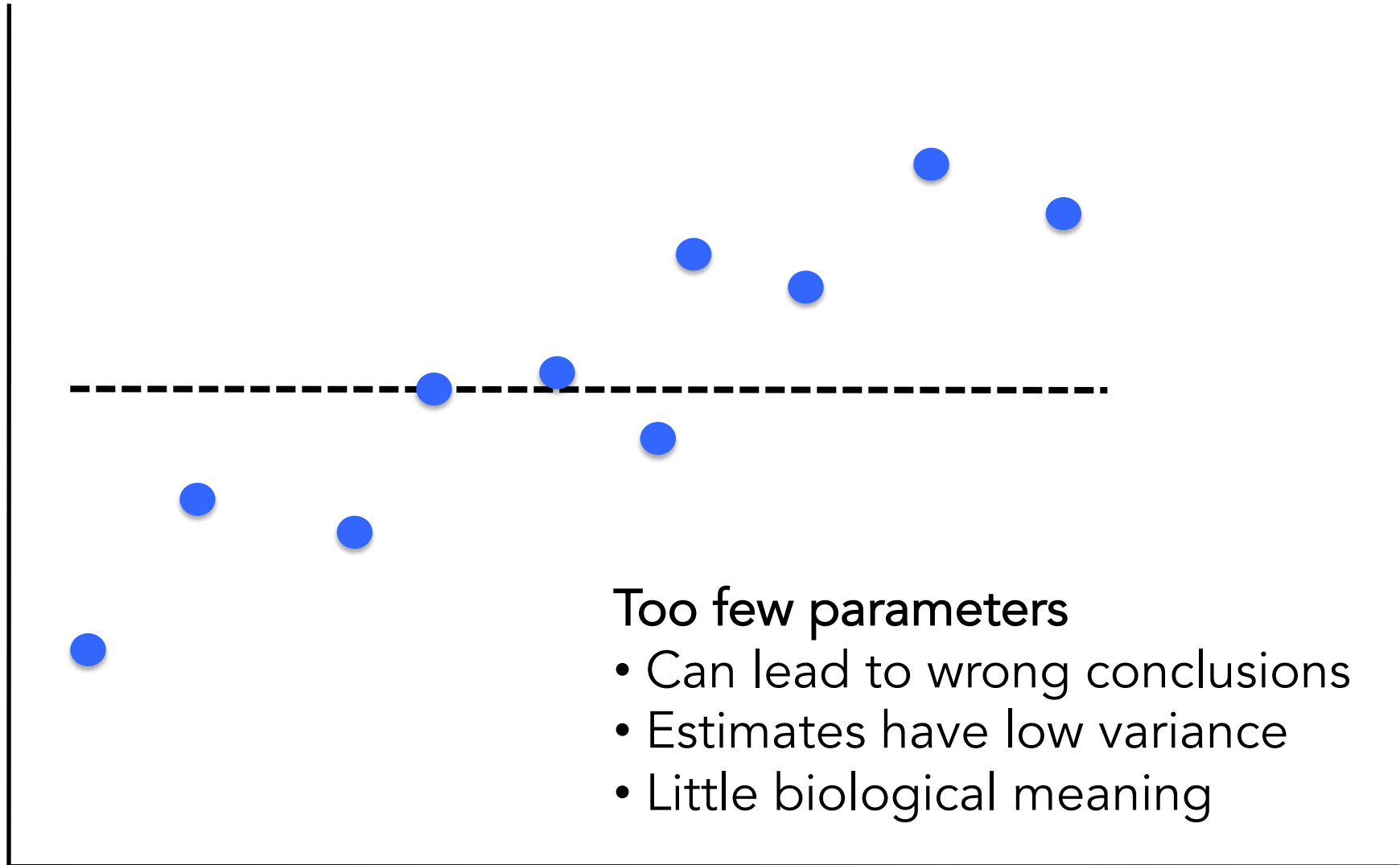  - JTT
  - WAG

# Model selection

# Model selection

1.  **Subjective model selection**

    •   Choosing a model that seems sensible

    •   Balancing the number of parameters against the amount of data available

    •   Biological motivation
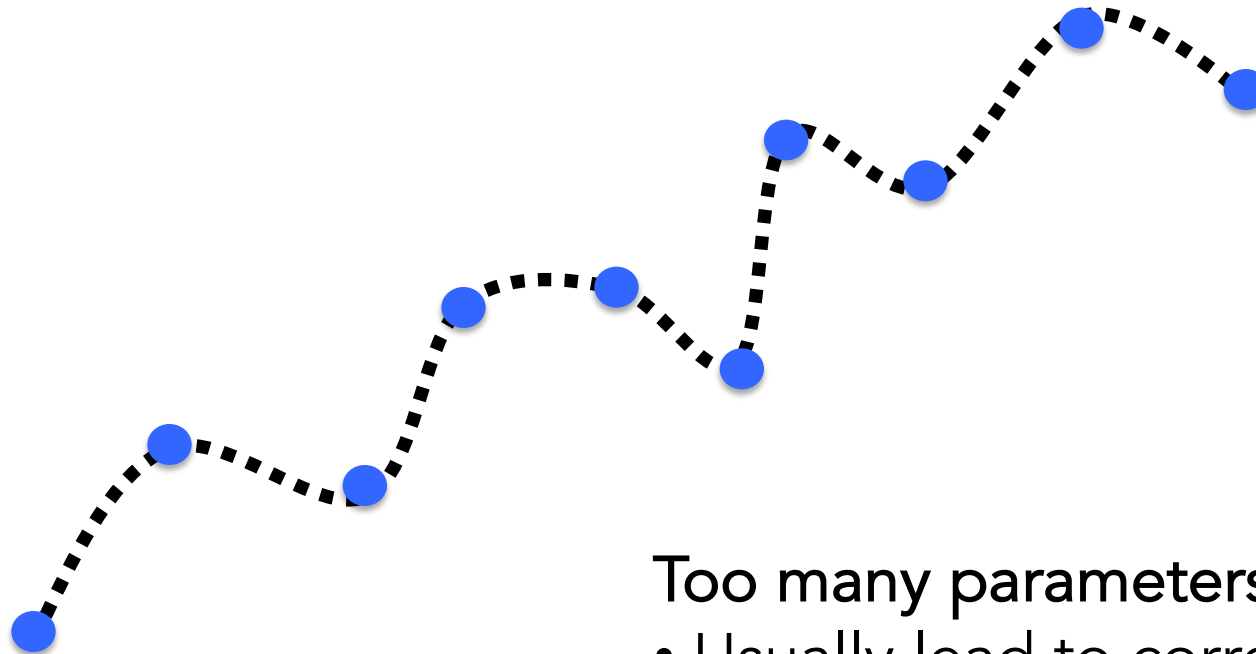

2.  **Objective model selection**

    •   Automated using information theory

    •   Statistical motivation

# Model selection



Too few parameters
- Can lead to wrong conclusions
- Estimates have low variance
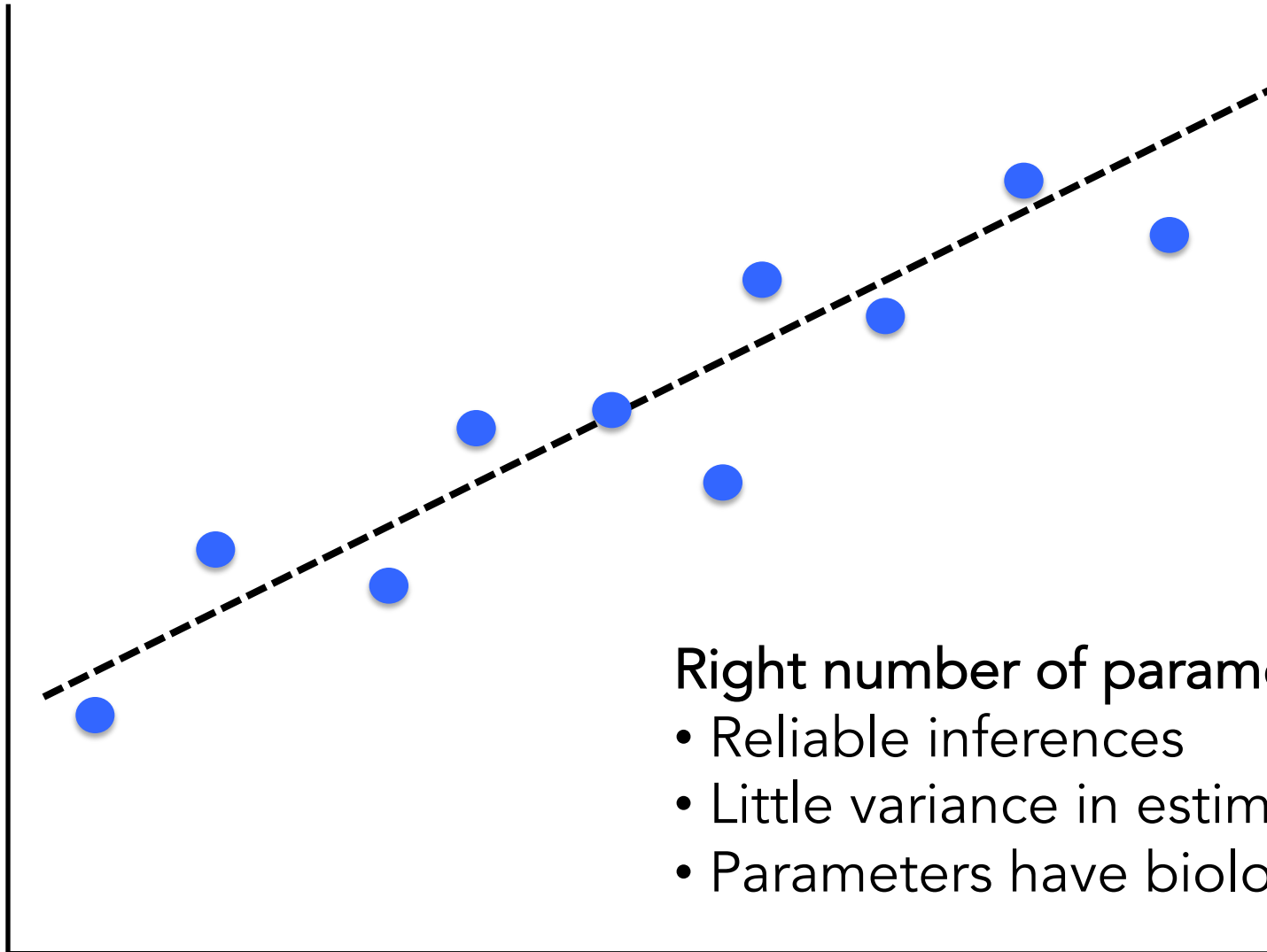- Little biological meaning

# Model selection

**Too many parameters**
- Usually lead to correct conclusions
- High variance in estimates
- Parameters can lack biological meaning
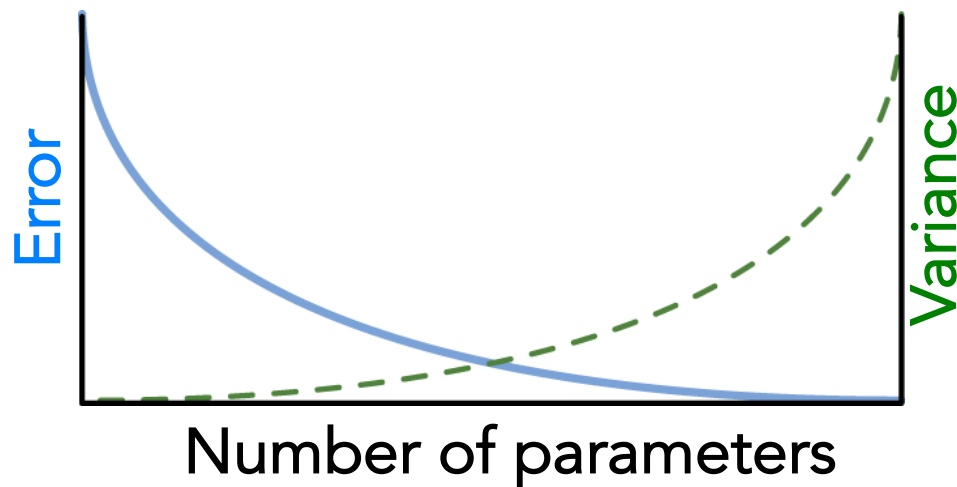
# Model selection



**Right number of parameters**
- Reliable inferences
- Little variance in estimates
- Parameters have biological meaning

# Model selection

- Adding parameters *always* improves model fit

- But adding parameters leads to greater variance in estimates

> Is the cost of additional parameters worthwhile?



Number of parameters

# Model selection

- Likelihood-ratio test (LRT)
  Used for comparing nested models

- Akaike information criterion (AIC)
  $AIC = -2\ln(\text{likelihood}) + 2k$

- Bayesian information criterion (BIC)
  $BIC = -2\ln(\text{likelihood}) + k\ln(n)$

# Substitution models in practice

- The tree topology is highly robust to the model used for inference

- **GTR+G** is acceptable for the majority of data sets

# Useful references

- **Model selection in phylogenetics**
Sullivan & Joyce (2005) *Annual Review of Ecology, Evolution, and Systematics*,
36: 445–466.

- **The effects of partitioning on phylogenetic inference**
Kainer & Lanfear (2015) *Molecular Biology and Evolution*, 32: 1611–1627.