Generalizing the model

Features of $\begin{cases} GTR \\ GM \end{cases}$ model : independence assumption
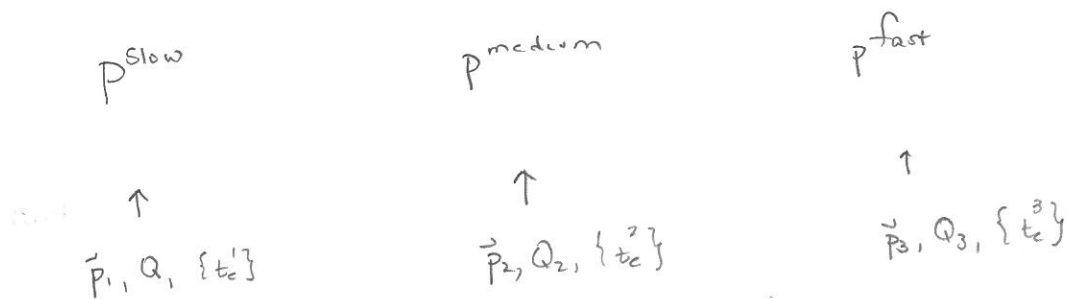
gains you data, but likely unrealistic

identically distributed

If aligned sequences are long, perhaps unrealistic to use <u>same</u> parameters on each site.

A common way to address this is with a MIXTURE MODEL

Simple Example : Fix T.

Suppose you have 3 classes of sites ((slow, medium, fast evolving) and choose GTR parameters for each of them. Compute the expected pattern frequency arrays

$P^{slow}$                    $P^{medium}$                    $P^{fast}$

↑                         ↑                         ↑

$\vec{P}_1, Q_1, \{t_e^1\}$            $\vec{P}_2, Q_2, \{t_e^2\}$            $\vec{P}_3, Q_3, \{t_e^3\}$

Additionally, choose <u>weighting</u> or <u>mixture</u> parameters $\alpha_1, \alpha_2$   (k=2)

then the joint frequency array is

$$P = \alpha_1 P_1 + \alpha_2 P_2 + (1-\alpha_1-\alpha_2) P_3$$

Each $P_i$ is called a MIXTURE COMPONENT and the number of parameters

to infer is        $3(3 + 5 + 2n-3) + 2$                The $\alpha_i$ are called the weights
                                                          or mixing parameters

# component ↑ $\vec{p}$ ↑ Q ↑ b.l.        # of components −1

Note that for a fixed binary tree on n taxa that if a mixture uses K components, then the number of numerical parameters to be inferred increases roughly by a factor of K. I.e. increases a lot.

For both biological and practical reasons, mixture models with fewer parameters are used in practice. Some examples include:

- GTR+I $\equiv$ GTR + Invariable sites model

  Two classes of sites: those that are free to mutate (GTR)

  those that are variable due to perhaps functional constraints (I)

  Parameters: GTR: $\vec{P_a}$, Q          $3+6 = 9$
  
  I: $\vec{P_I}$          $= 3$
  
  $= 1$

  weights: $S$

Excluding branch lengths, this is an 13-parameter model. I.e. 4 additional parameters.

  A variation sets $P_I = P_Q$, i.e. assumes the base distribution is the same over all sites, both variable and invariable.

- GTR + rate variation          (discrete)          Lots of variation

  Assumes K classes of sites but that the mutation rate is scaled depending which class you are in.

Example: GTR + rate variation    $k$ classes    # of classes chosen by user.

Numerical parameters (excluding branch lengths on $T$) are

- GTR parameters              $\vec{p}, Q$                                    used for __all__ sites

- classes weights        $s_1, s_2, \ldots, s_k$    $\sum s_i = 1$    distribution of sites to
                                                              $s_i > 0$          classes

- rates $r_1, r_2, \ldots, r_k$      $r_i \geqslant 0$              scaling rate for
                                                                                    $i$-th class

The pattern frequency array is

$$P = s_1 P_1 + s_2 P_2 + \ldots + s_k P_k$$

Where $P_i$ is the expected pattern freq
array for the $i$-th class

$P_i$ is computed using $\vec{p}, \underbrace{r_i Q}$   and branch lengths $\{t_e\}$

Scaled version of $Q$, with scaling factor the rate $r_i$.
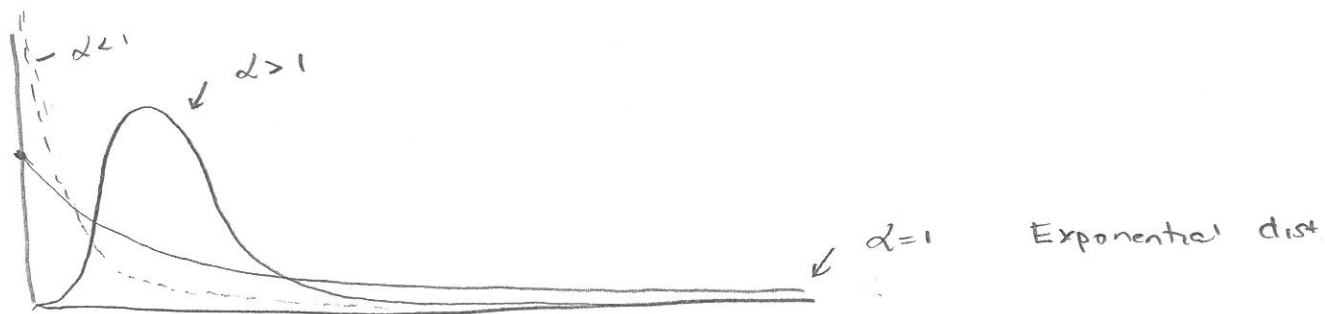
$r_i \geqslant 1$ , $r_i < 1$ , etc.

Visual effect.

Discuss how to use such a model for simulations ...

Variations:   Instead of choosing the rate functions at random,
Choose them from a distribution.   In practice, the $\Gamma$-distribution
is used or in reality the discrete-$\Gamma$.

Gamma distributions in phylogenetics is a 1-parameter family
of distributions with the unknown parameter called $\alpha$ = shape parameter.

The densities for various values of $\alpha$ are shown



$\alpha < 1$
$\alpha > 1$
$\alpha = 1$     Exponential dist

Notice all $r_i > 0$ are possible rates and the shape of the density says

something about the probability of various rates.

　　　See R demo and discuss meaning.

These models are frequently called RATES-ACROSS-SITES models.

In software, a RAS model GTR+$\Gamma$ is implemented as a

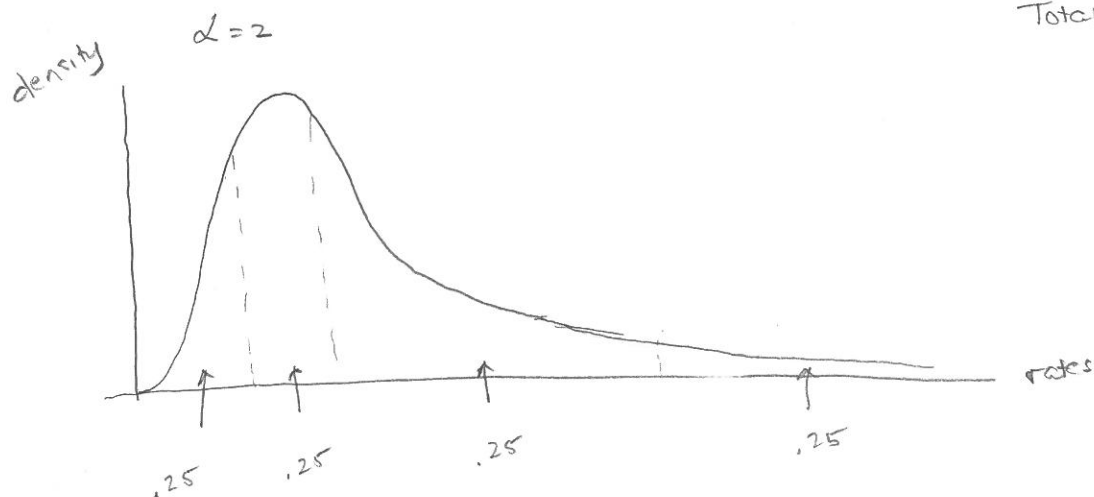GTR+ discrete-$\Gamma$.   i.e.   GTR+$d\Gamma(4)$         Z. Yang   1994   PAML

　　　　　　　　　　　　　　　↑
　　　　　　　　　　　　　# of categories

　　　　　　↑

　　　　discrete.

　　Example: GTR+ $d\Gamma(4)$         Suppose the user choose 4 categories

and GTR parameters $\vec{\pi}$, Q are known as is the $\Gamma$-shape parameter $\alpha$.

density         $\alpha = 2$                              Total area under curve = 1
                                                          Areas $\longleftrightarrow$ probability



　　　　　　　　　　　　　　　　　　　　　　rates

.25     .25        .25              .25

Break the density curve into 4 regions, each of which has area .25

$\Rightarrow$ probability = .25

In each of the four regions, find the mean $r_1, r_2, r_3, r_4$, then

if a site is chosen to be in class $i$, $1 \leq i \leq 4$, the Markov matrix

on a branch of length $t$ is given by

$$M = e^{Q r_i t}$$

ie scale $Q$ by $r_i$ and $t$.

and the expected pattern frequency array is

$$P = s_1 P_1 + s_2 P_2 + s_3 P_3 + (1 - (s_1 + s_2 + s_3)) P_4$$

$$\uparrow r_1 \qquad \uparrow r_2 \qquad \uparrow r_3 \qquad \uparrow r_4$$

$\vec{P}, Q$ for all 4 classes.

Other variants are of course possible.

The most widely used model in practice is

$$GTR + I + \Gamma$$

secretly $GTR + I + d\Gamma(4)$

Parting comments:

In a $GTR + I + \Gamma$ model, a site belongs to 1 category so that rate applies to that site in all parts of the tree. There are no speed ups or slow downs for that site in different parts of the tree.

Looking at the $K = 4$ categories together, the metric trees they infer are all scalings of 1 tree.