Generalizing the model

Features of $\begin{cases} \text{GTR model} \\ \text{GM} \end{cases}$ : independence assumption

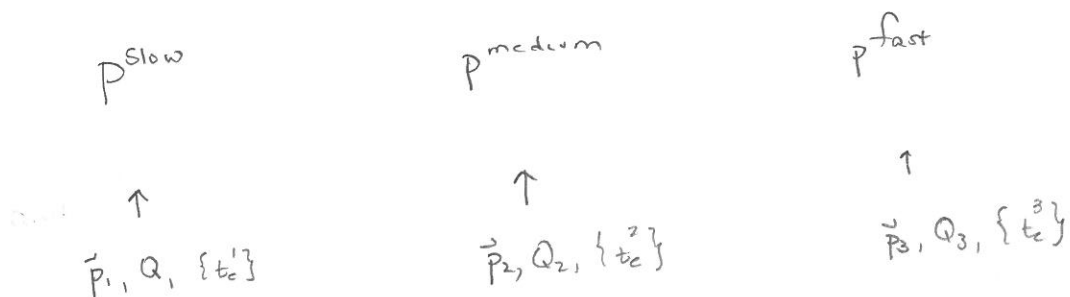gains you data, but likely unrealistic

identically distributed

If aligned sequences are long, perhaps
unrealistic to use _same_ 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

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

$\uparrow$                      $\uparrow$                        $\uparrow$

$\vec{P_1}, Q_1, \{t_e^1\}$          $\vec{P_2}, Q_2, \{t_e^2\}$          $\vec{P_3}, Q_3, \{t_e^3\}$

Additionally, choose _weighting_ or _mixture_ 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       $3(3 + 5 + 2n-3) + 2$

$\underset{\text{# component}}{\uparrow} \underset{\vec{P}}{\uparrow} \underset{Q}{\uparrow} \underset{b.l.}{\uparrow}$     $\underset{\text{# of components} -1}{\uparrow}$

The $\alpha_i$ are called the weights
or mixing parameters

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        ≡    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 : $\vec{P_Q}, Q$                     $3+6 = 9$
                        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.

                                          (discrete)              Lots of variation
- GTR + rate 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

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

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 $r_i$.

$$r_i \geqslant 1, \quad r_i < 1, \quad \text{etc.}$$
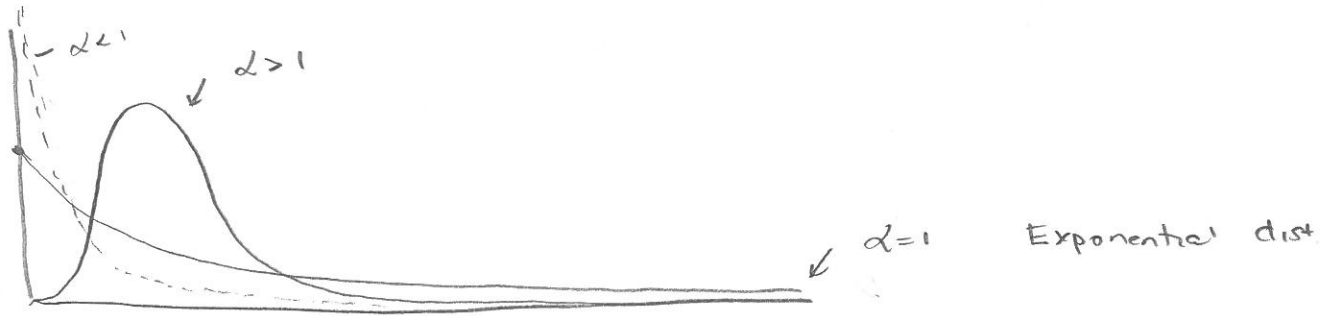
visual effect.

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



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.