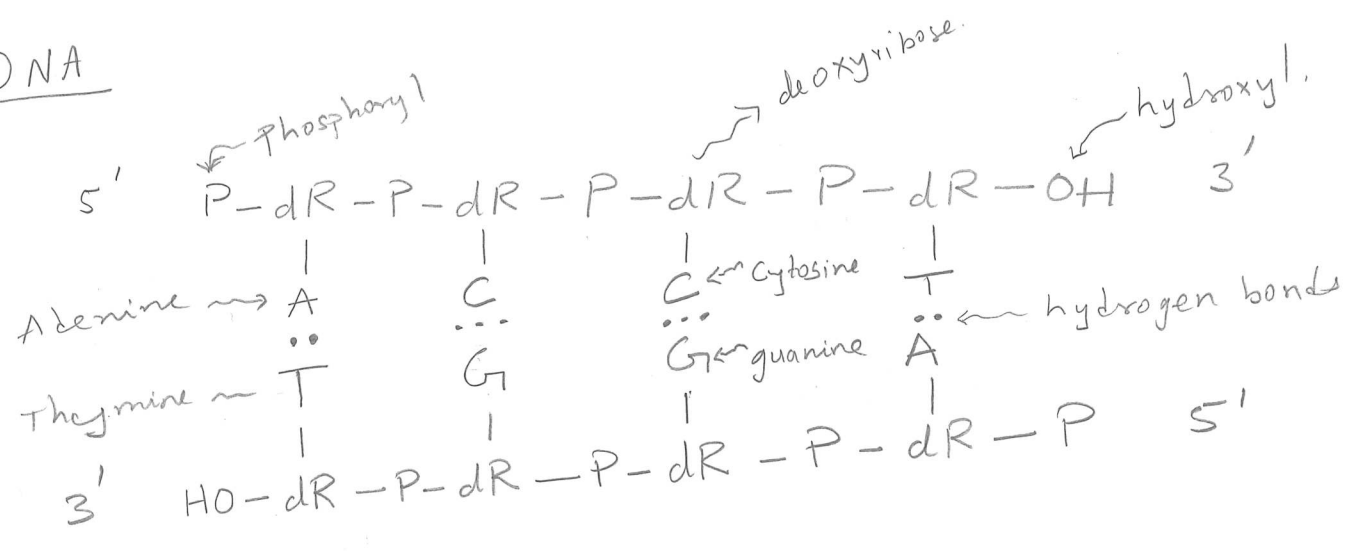


DNA



- Hereditary information of living things is carried by DNA sequences $\in \{A, C, G, T\}^N$ (its genome).
 $N \approx 12M$ For yeast.
 set of nucleotides

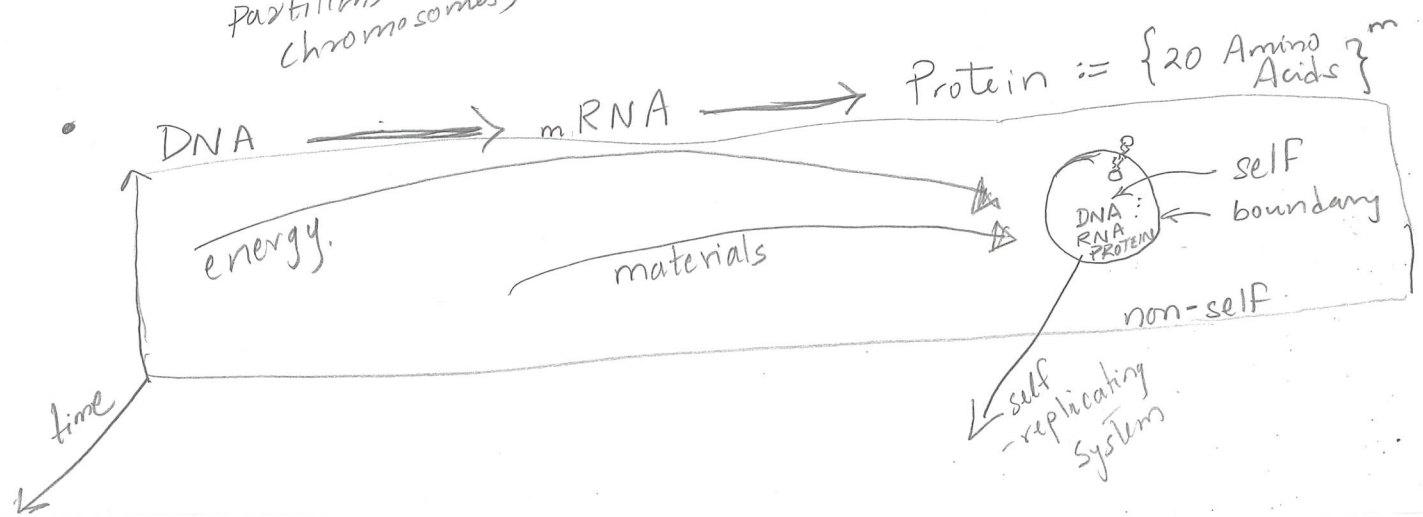
- Double-stranded DNA can replicate! by copying from each strand



- Nucleotide frequencies is non-uniform for Yeast

A	= 0.3090
T	= 0.3078
C	= 0.1917
G	= 0.1913

- organisms have different number of copies of their genome (organized by partitions called chromosomes)
 - 1 (haploid) eg. bacteria
 - 2 (diploid) eg. humans, (most animals)
 - 4 (tetraploid) eg. plants.
 - 6 (hexaploid) eg. wheat
 - > 6 (polyploid) eg. sorghum has 100 chromosomes of 8 types?

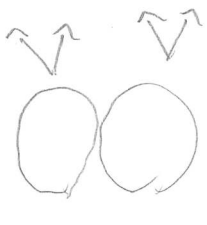
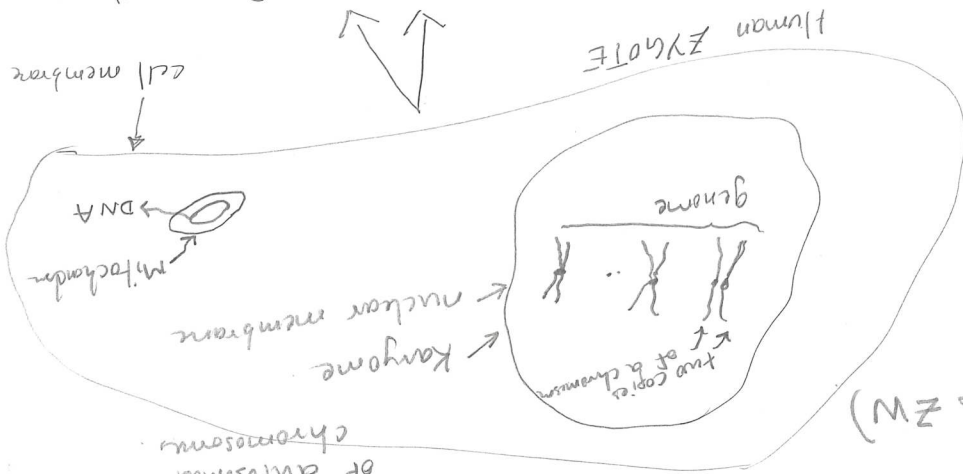


2

Sexual Reproduction in diploids & recombination

humans → 1 sex chromosome, females XX; males XY
 → 22 autosomes, females & males have two copies AA

(Birds: males ZZ; females ZW)



Adult female

Adult male

Recombination between homologous pair of chromosomes



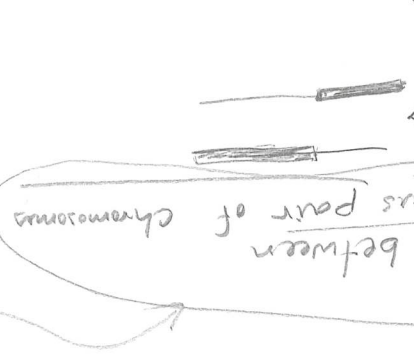
recombination



spERM

diploid zygote

No recombination - Y chromosome & Mitochondrial DNA



egg

mitochondria

Wright-Fisher Model (1931)

(2)

genetic locus := a location in the genome of an organism

Alleles := different versions of the genetic information encoded at a locus.

eg Alleles A and a could represent 'distinct' DNA sequences: $A =$ CTGAAATCGTA A
 $a =$ CTGATATCGTAG G

posn 10193 on chrom 13
X X

or just at a single position on a chromosome = site

$A = G$ site at posn. 15326 on chrom 14
 $a = T$

Diploid organisms have two copies so they will be

$A A$ or $a a$
 $A a$
 (AA) (Aa) (aa)

Three genotypes at a biallelic locus

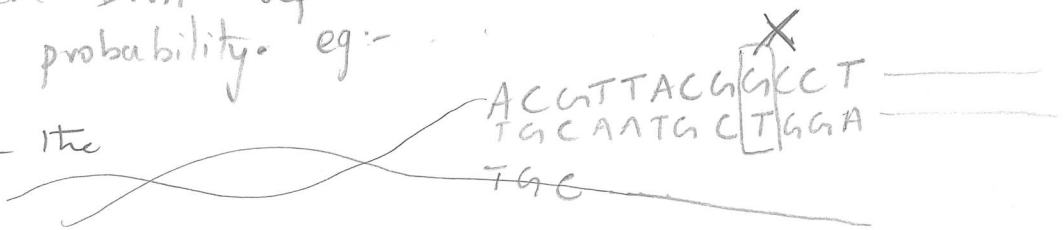
fitness

of an individual is a measure of its ability to survive and to produce offspring.

Mutation

When DNA replicates mistakes ^{or mutations} can be made with small probability. eg:-

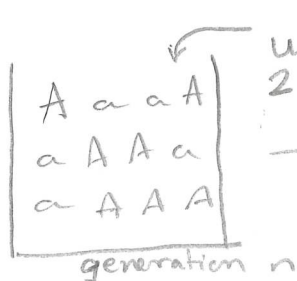
(mutations change the DNA sequence of the original template DNA)



neutral evolution when mutation does not affect fitness.

WF Model

as sampling with replacement from an urn with $2N$ balls.



- nonoverlapping generation (deaths & rebirths every generation)
- random mating.

④ At gen. n i balls (individuals) have allele A and $2N-i$ balls have allele a

$2N$ = number of haploid alleles in a diploid population of size N .

To build up the $(n+1)$ -th gen. choose from the urn at gen n , $2N$ times with replacement

X_n = Number of A 's in gen n

X_n is a Markov chain on $\{0, 1, 2, \dots, 2N\}$

Since $P_r(X_{n+1} = j | X_n = i, X_{n-1} = x_{n-1}, \dots, X_0 = x_0)$

$$= P_r(X_{n+1} = j | X_n = i)$$

$$=: p(i, j)$$

$$= \binom{2N}{j} \left(\frac{i}{2N}\right)^j \left(\frac{2N-i}{2N}\right)^{2N-j}$$

prob. of choosing on A j times

prob. of choosing allele a $2N-j$ times

Binomial $(2N, \frac{i}{2N})$ RV

PMF of

$$\binom{2N}{j} = \frac{(2N)!}{j!(2N-j)!}$$

where,

Binomial coefficient

where $j!$ = $1 \cdot 2 \cdot 3 \cdot \dots \cdot j$ factorial is number of ways of ordering j items in a row.

number of ways of choosing j items out of $2N$ items

Long-Term behaviour of $\{X_n\}_{n=0}^{\infty}$, the W-F Markov chain for the number of A alleles among $2N$ alleles. (5)

Either $X_n \xrightarrow{n \rightarrow \infty} \begin{cases} 0 & \text{with no A alleles} \\ 2N & \text{with no a alleles} \end{cases}$

Because 0 & $2N$ are absorbing states.
(the chain can never leave once it enters either 0 or $2N$)

We say fixation has occurred when X_n enters an absorbing state (whole population is fixated on one allele).

fixation Time τ , i.e. first time the popn is all a's or all A's.

$$\tau := \min \left\{ n : X_n = 0 \text{ or } X_n = 2N \right\}$$

Thm 1 In the W-F model, the prob. of fixation in the all A's state, given you start with i A's is:

$$\Pr(X_\tau = 2N \mid X_0 = i) = \frac{i}{2N}$$

Proof: Since $2N < \infty$, fixation with all A's or a's will eventually occur (i.e. 0 & $2N$ are absorbing states).

Since the expectation of the Binomial(n, p) RV is $n \cdot p$

$$E(X_{n+1} \mid X_n = i) = 2N \left(\frac{i}{2N} \right) = i = X_n$$

Taking expected value on both sides, we get

$$E(X_{n+1}) = E(X_n) \implies \text{so average value of the number of A's remains the same through time.}$$

Now, since $X_n = X_\tau$ when $n > \tau$

$$(\star) \quad i = E(X_n \mid X_0 = i) = E(X_\tau; \tau \leq n \mid X_0 = i) + E(X_\tau; \tau > n \mid X_0 = i)$$

where, $E(X; A) =$ Expected value over the set A. (6)

Now let $n \rightarrow \infty$ with the fact that $|X_n| \leq 2N$

we get that: $E(X_t; \tau \leq n | X_0 = i) \rightarrow E(X_t | X_0 = i)$

and $E(X_t; \tau > n | X_0 = i) \rightarrow 0$

So, from (*)

$$i = E(X_t | X_0 = i) = 2N \cdot P(X_t = 2N | X_0 = i)$$

$$\text{and } \therefore P(X_t = 2N | X_0 = i) = \frac{i}{2N} \quad \square$$

Thus, the prob. of being fixed in all A's state having started with i A alleles is simply $\frac{i}{2N}$, the proportion of A alleles at the start.

Mutation

Now suppose that mutations occur in the population, whereby $A \xrightarrow{\alpha}$ or $a \xrightarrow{\alpha}$ ^{mutates to}

at rate μ . Then from Thm 1 we get Kimura's result:

Thm 2

Under W-F model, the rate of fixation of neutral mutations in a ^{haploid} popⁿ of size $2N$ is the mutation rate μ

Proof: Note that mutations occur to some individual in the population at rate $2N\mu$ and since this is the only individual with this mutation, it goes to fixation (by Thm 1) with prob. $\frac{1}{2N}$.

Heterozygosity is the prob. that two copies of the locus chosen (without replacement) at time n are different:

$$H_n^o := \frac{X_n(2N - X_n)}{\frac{2N(2N-1)}{2}} \leftarrow \begin{matrix} \text{number of choosing} \\ \text{an A and an a allele} \end{matrix}$$

$$= \frac{2X_n(2N - X_n)}{2N(2N-1)} \leftarrow \begin{matrix} \binom{2N}{2} = \frac{2N!}{(2N-2)! \cdot 2!} = \frac{2N(2N-1)}{2} \\ \text{The number of ways of} \\ \text{choosing any two out of } 2N \\ \text{items} \end{matrix}$$

Thm 3

Let $h(n) := E(H_n^o)$, the expected heterozygosity at time n in the W-F model. Then,

$$h(n) = \left(1 - \frac{1}{2N}\right)^n h(0)$$

Proof

Let us label the alleles by $1, 2, \dots, 2N$
 $2N$ copies of the locus. or individuals.

Suppose we pick two individuals at time n labelled by $x_1(0)$ and $x_2(0)$.

Each indiv. $x_i(0)$, $i=1, \dots, 2N$ is a descendant of some individual $x_i(1)$ at time $n-1$, who is a descendant of $x_i(2)$ at time $n-2$, etc.

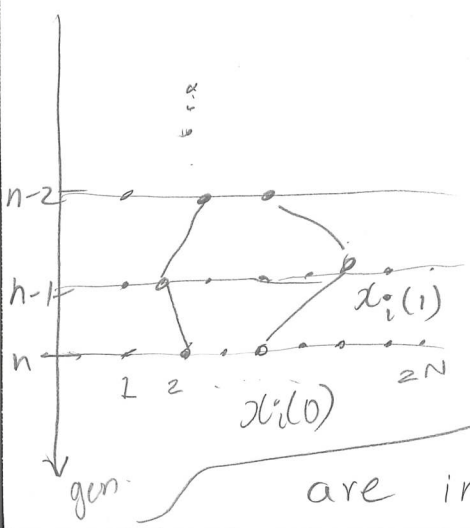
So, $(x_i(m) : 0 \leq m \leq n)$ gives the genealogy or ancestral lineage of $x_i(0)$
(ancestors back in time)

Let $x_1(0)$ and $x_2(0)$ be two randomly chosen indivs.

NOTE If $x_1(m) = x_2(m)$

then $x_1(l) = x_2(l)$ for $m < l \leq n$.

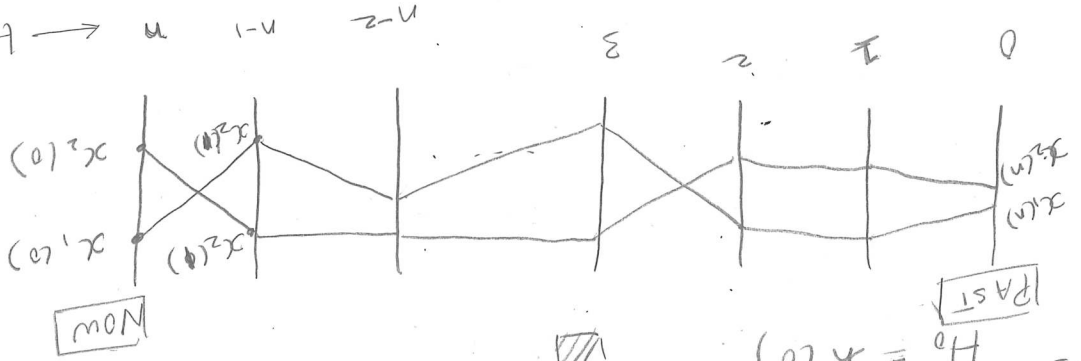
If $x_1(m) \neq x_2(m)$ the parental choices are indep. with $\Pr\{x_1(m+1) \neq x_2(m+1)\} = 1 - \frac{1}{2N}$



For $x_1(n) \neq x_2(n)$, different parents should be chosen for all times $1 \leq m \leq n$,

$$\Pr \{ x_1(m) \neq x_2(m) : 1 \leq m \leq n \} = \left(1 - \frac{1}{2N}\right)^n$$

Finally $x_1(n)$ and $x_2(n)$ are two random indivs from time 0, so the prob they are different is $H_0 = h(0)$



A pair of genealogies remaining distinct for n gens.

The Coalescent (Kingman 1982)

$$1 - x \approx e^{-x}$$

if $2N$ is large

$$h(n) = \left(1 - \frac{1}{2N}\right)^n h(0) \approx e^{-n/2N} h(0)$$

Thm 3 (with $\frac{1}{2N}$ small)

When x is small

Suppose, instead of 2, we sample K random individuals, K will pick the same parent in 2 of the K prev. gen. is

$$\frac{1}{2N} \approx \frac{2}{K(K-1)}$$

Number of ways of picking 2 from K indivs

The prob. that these two indivs will choose the same parent

! it's only approximate because we ignore prob. of $O(\frac{1}{N^2})$ with prob. of three indivs choose same parent

Then the prob. in same parent in 2, we sample K random individuals, K will pick the same parent in 2 of the K prev. gen. is

Thm 4 When measured in units of $2N$ generations, the amount of time during which there are k lineages, t_k , has approximately exponential distribution with rate $k(k-1)/2$. (9)

Proof: $\Pr \{ k \text{ lineages remain distinct for } n \text{ generations} \}$

$$\approx \left(1 - \frac{k(k-1)}{2} \cdot \frac{1}{2N} \right)^n \approx \exp \left(- \frac{k(k-1)}{2} \cdot \frac{n}{2N} \right)$$

Recall $\left[\begin{array}{l} T \sim \text{Exponential}(\lambda) \\ P(T > t) = e^{-\lambda t} \end{array} \right. \begin{array}{l} \uparrow \text{rate parameter} \\ E(T) = \frac{1}{\lambda} \end{array}$

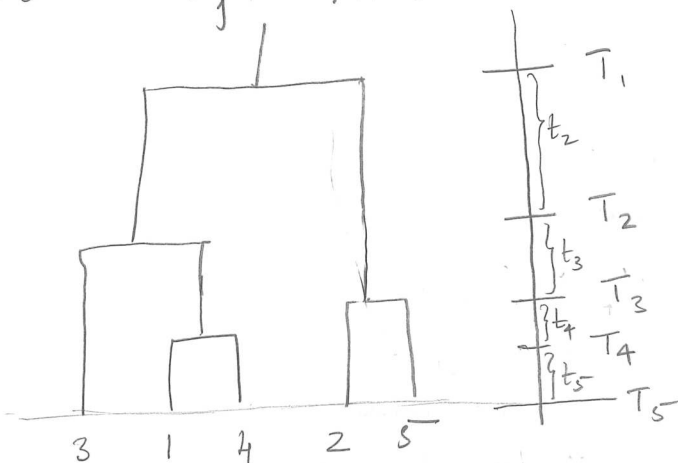
By letting $N \rightarrow \infty$ and expressing time in terms of $2N$ generations, i.e. letting $t = n/2N$, we get

The time to the first coalescence (choosing same ancestor) event is $\text{Exponential} \left(\frac{k(k-1)}{2} \right)$ RV.

Thus k lineages coalesce to $k-1$ lineages at rate $k(k-1)/2$ using continuous time Markov chain (CTMC) terminology. \square

The limit of genealogies in Thm 4 is called the coalescent

Let $T_j =$ first time with j lineages, then we have.



A realisation of the coalescent for a sample of size 5 drawn randomly from a large pop of size $2N$.

Expected Coalescent times

$E(t_2) = 1$, $E(t_3) = \frac{1}{3}$, $E(t_4) = \frac{1}{6}$, $E(t_5) = \frac{1}{10}$

(10) (we use lower case for EVs here for exposition)

$\therefore E(t_k) = \frac{1}{2k(k-1)}$

T_1 = time of the appearance of the Most Recent Common Ancestor (MRCA) of the sample

$= t_n + t_{n-1} + \dots + t_3 + t_2$

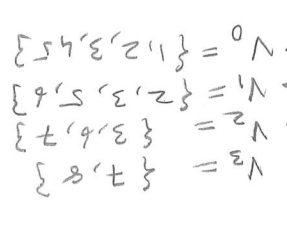
$E(T_1) = E\left(\sum_n^{k=2} t_k\right) = \sum_n \frac{1}{2} k(k-1) = 2 \sum_n \left(\frac{1}{k-1} - \frac{1}{k}\right)$

$= 2 \left(\sum_{k=1}^{n-1} \frac{1}{k} - \sum_{k=2}^n \frac{1}{k} \right) = 2 \left(1 - \frac{1}{n} \right) \rightarrow 2$ for large samples

But $\overline{E(t_2)} = 1 \approx E(T_1)$

Simulating The Coalescent

Let's label internal nodes.



input sample size n

initialize

$V_0 = \{1, 2, 3, \dots, n\}$, $T_n = 0$

for $k = 0, 1, \dots, n-2$ do

pick v_k and j_k from V_k

$V_{k+1} \leftarrow V_k \setminus \{v_k, j_k\} \cup \{n+k+1\}$

In tree connect $v_k \rightarrow n+k+1$ and $j_k \rightarrow n+k+1$

Let $T_{n-k} \sim \text{Exponential}(n-k)$

Let $T_{n-k-1} \leftarrow T_{n-k} + t_{n-k}$

(tree 2)

label branch by smaller number on lower end $1 \leq i \leq 2n-2$

After generating genealogy

⊕ Mutations

μ = mutation rate per gen.

$\theta = 4N\mu$

$X(i) \sim \text{Poi}\left(\frac{\theta}{2}\right)$ (ancestral)

distributed randomly across m sites

ind site $1 \dots m$

(tree 1)

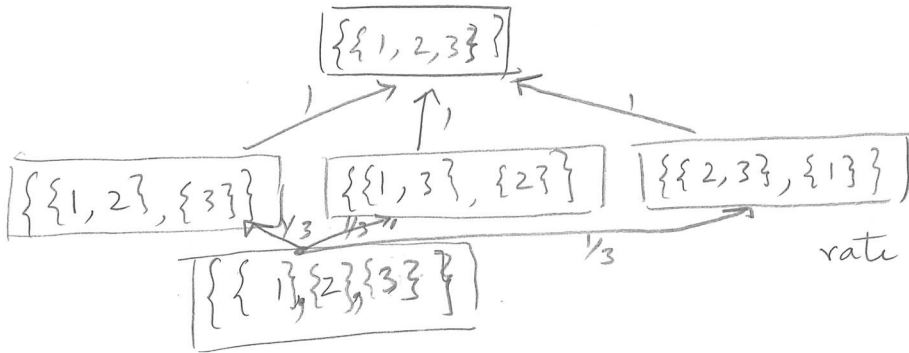
$\text{anc}[i_k] = n+k+1$

$\text{anc}[j_k] = n+k+1$

State space of The coalescent

(11)

ξ_n = Set of all set partitions of $\{1, 2, \dots, n\}$



$$\text{rate} = \binom{3}{2} = \frac{3!}{1!2!} = 3$$

Infinite Sites Model (Kimura, 1969)

Mutations always occur at distinct sites.

(picture DNA sequence as a unit interval in \mathbb{R})

binary incidence matrix

\downarrow [BIM]

0	1	0	...	0
0	0	1	...	1

1 2 3 ... n

T	A	T	...	T
G	C	G	...	C

1 2 3 ... n

ancestral state
↑ individuals
site
chrom. posn.
13 012541
13 012397

derived state
ancestral state

segregating sites

$S_n =$ number of segregating sites = number of site positions where some pair of individual DNA sequences differ.