

Introduction to Bioinformatics



Biological words

Recap

- p DNA codes information with alphabet of 4 letters: A, C, G, T
- p In proteins, the alphabet size is 20
- p DNA -> RNA -> Protein (genetic code)
 - n Three DNA bases (triplet, codon) code for one amino acid
 - n Redundancy, start and stop codons

Given a DNA sequence, we might ask a number of questions

What sort of statistics should be used to describe the sequence?

```
1 atgagccaag ttccgaacaa ggattcgccgg ggaggataga tcagcgcccc agaggggtga  
61 gtcggtaaag agcattggaa cgtcgagat acaactccca agaaggaaaa aagagaaagc  
21 aagaagcgga tgaatttccc cataaacgcca gtgaaaactct aggaagggga aagagggaaag  
81 ctggaaagaga aggagggggg cgtccccatc cgaggggggcc agggggccaa tggtggggac  
41 actccgggcc gaagggttga gagtaccccca gagggaggaa gccacacggg gtagaaacaga  
01 gaaatcacct ccagaggacc cttcagcga acagagagcg catcgcgaga gggagttagac  
61 catagcgata ggagggatg ctaggagttg ggggagaccg aagcgaggag gaaagcaaag  
21 agagcagcgg ggctagcagg tgggtgttcc gcccccccgag aggggacgag tgaggcttat  
81 cccggggAAC tcgacttatac gtccccacat agcagactcc cggacccccc ttcaaagtga  
41 ccgagggggg tgactttgaa cattggggac cagtggagcc atggatgct cctcccgatt
```

What sort of organism did this sequence come from?

Does the description of this sequence differ from the description of other DNA in the organism?

What sort of sequence is this? What does it do?

Biological words

- p We can try to answer questions like these by considering the *words* in a sequence
- p A *k*-word (or a *k-tuple*) is a string of length *k* drawn from some alphabet
- p A DNA *k*-word is a string of length *k* that consists of letters A, C, G, T
 - n 1-words: individual nucleotides (bases)
 - n 2-words: dinucleotides (AA, AC, AG, AT, CA, ...)
 - n 3-words: codons (AAA, AAC, ...)
 - n 4-words and beyond

1-words: base composition

- Typically DNA exists as *duplex* molecule
(two complementary strands)

5' -GGATCGAAGCTAAGGGCT- 3'

3' -CCTAGCTTCGATTCCCGA- 5'

Top strand: 7 G, 3 C, 5 A, 3 T

Bottom strand: 3 G, 7 C, 3 A, 5 T

Duplex molecule: 10 G, 10 C, 8 A, 8 T

Base frequencies: 10/36 10/36 8/36 8/36

These are something
we can determine
experimentally.

$$\text{fr}(G + C) = 20/36, \text{fr}(A + T) = 1 - \text{fr}(G + C) = 16/36$$

G+C content

- p $\text{fr}(G + C)$, or *G+C content* is a simple statistics for describing genomes
- p Notice that one value is enough characterise $\text{fr}(A)$, $\text{fr}(C)$, $\text{fr}(G)$ and $\text{fr}(T)$ for duplex DNA
- p Is G+C content (= base composition) able to tell the difference between genomes of different organisms?
 - n Simple computational experiment, if we have the genome sequences under study (-> exercises)

G+C content and genome sizes (in megabasepairs, Mb) for various organisms

p	Mycoplasma genitalium	31.6%	0.585
p	Escherichia coli K-12	50.7%	4.693
p	Pseudomonas aeruginosa PAO1	66.4%	6.264
p	Pyrococcus abyssi	44.6%	1.765
p	Thermoplasma volcanium	39.9%	1.585
p	Caenorhabditis elegans	36%	97
p	Arabidopsis thaliana	35%	125
p	Homo sapiens	41%	3080

Base frequencies in duplex molecules

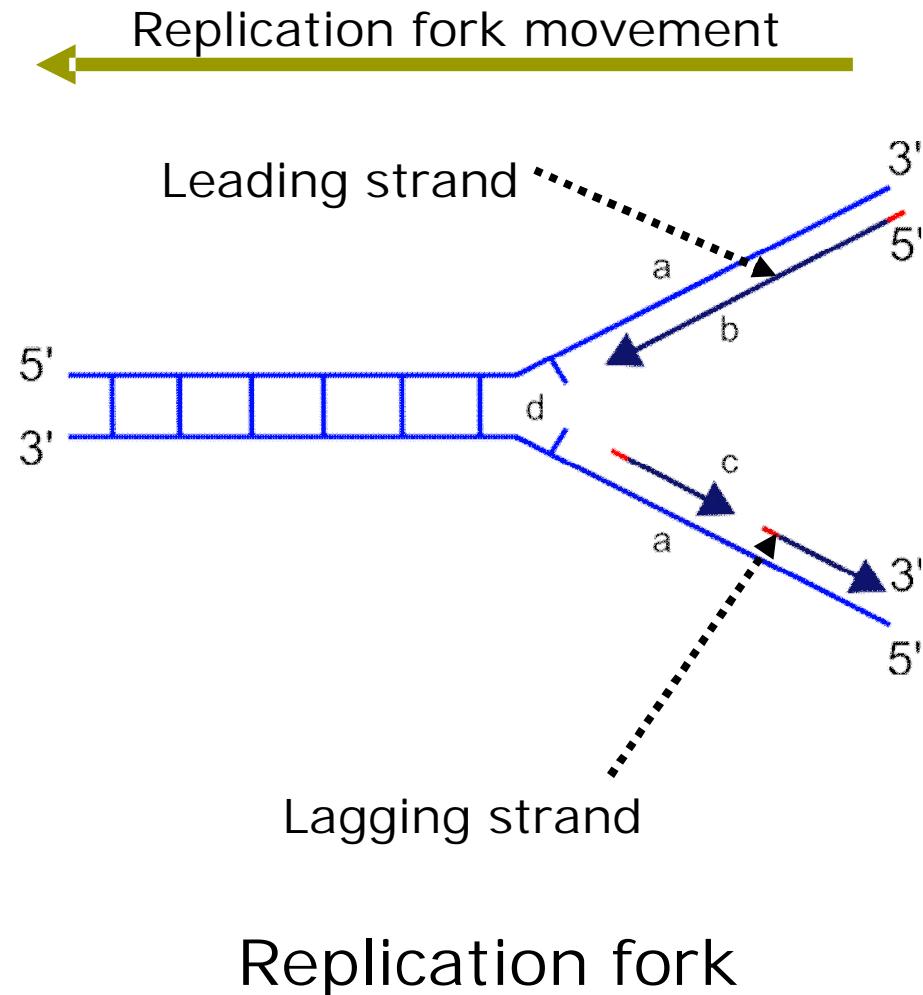
- p Consider a DNA sequence generated randomly, with probability of each letter being independent of position in sequence
- p You could expect to find a uniform distribution of bases in genomes...

5' - . . . GGATCGAAGCTAAGGGCT . . . - 3'
3' - . . . CCTAGCTTCGATTCCCGA . . . - 5'

- p This is not, however, the case in genomes, especially in prokaryotes
 - n This phenomena is called *GC skew*

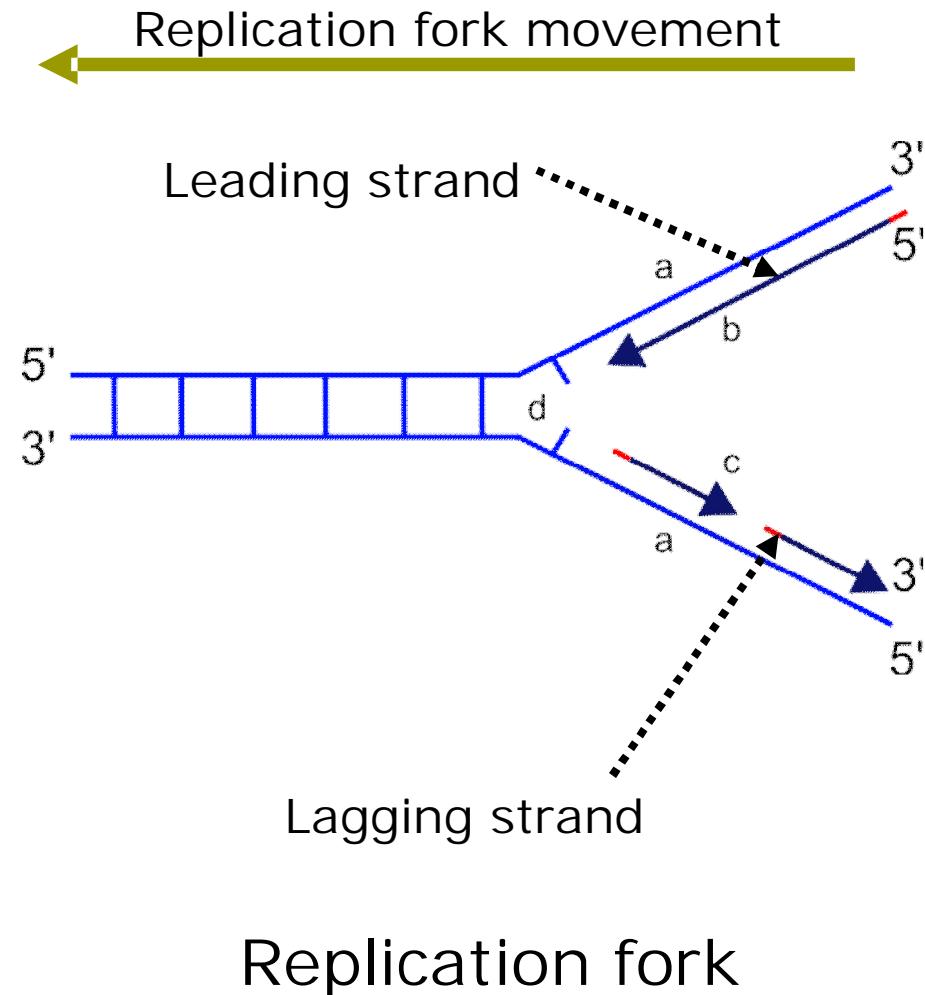
DNA replication fork

- p When DNA is replicated, the molecule takes the *replication fork* form
- p New complementary DNA is synthesised at both strands of the "fork"
- p New strand in 5'-3' direction corresponding to replication fork movement is called *leading strand* and the other *lagging strand*



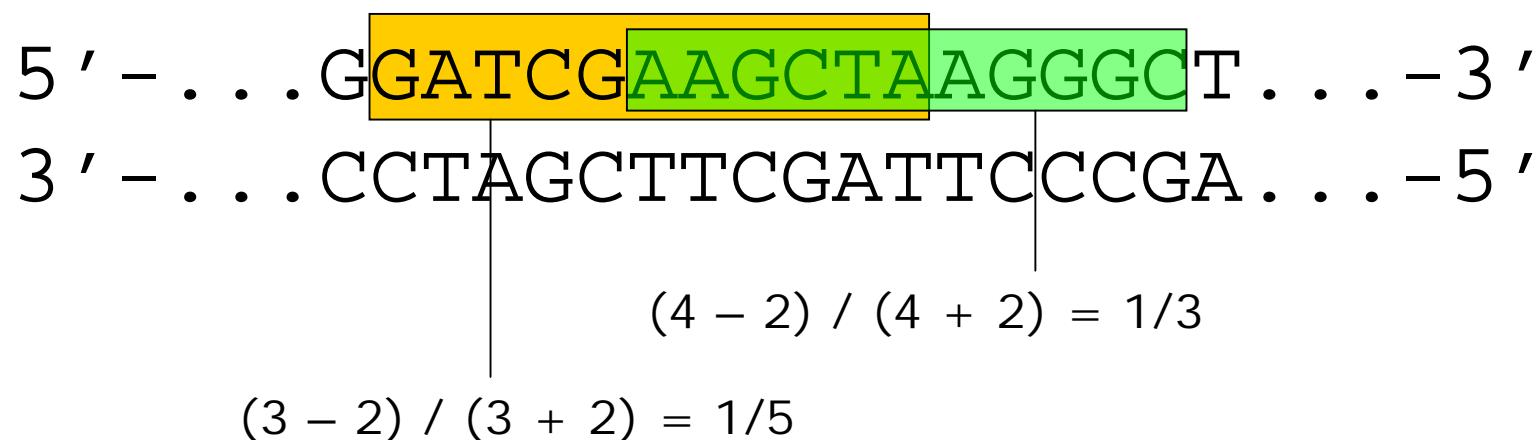
DNA replication fork

- This process has specific starting points in genome (*origins of replication*)
- Observation: Leading strands have an excess of G over C
- This can be described by *GC skew statistics*



GC skew

- GC skew is defined as $(\#G - \#C) / (\#G + \#C)$
- It is calculated at successive positions in intervals (windows) of specific width



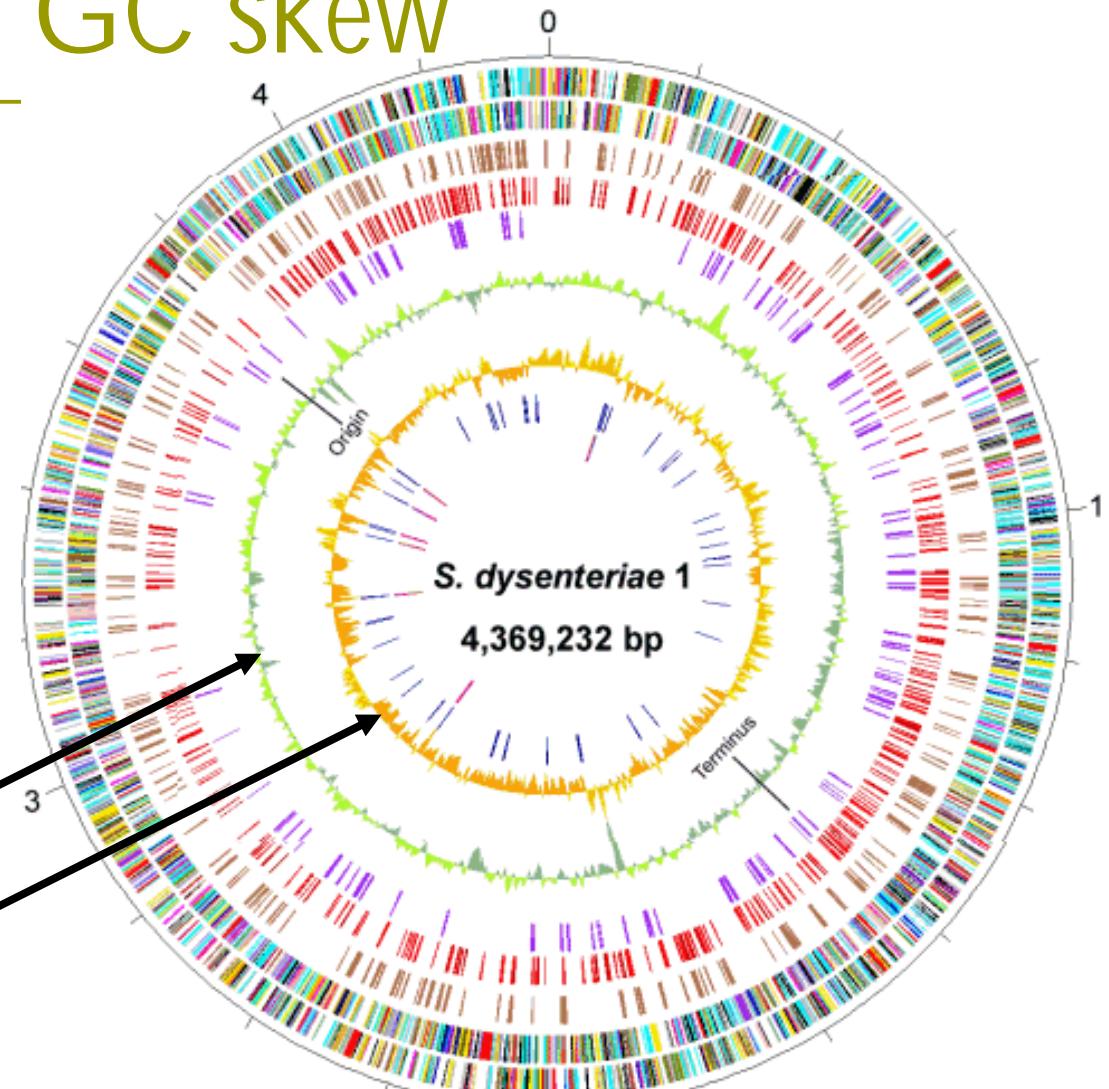
G-C content & GC skew

p G-C content & GC skew statistics can be displayed with a *circular genome map*

(10kb window size)

G+C content

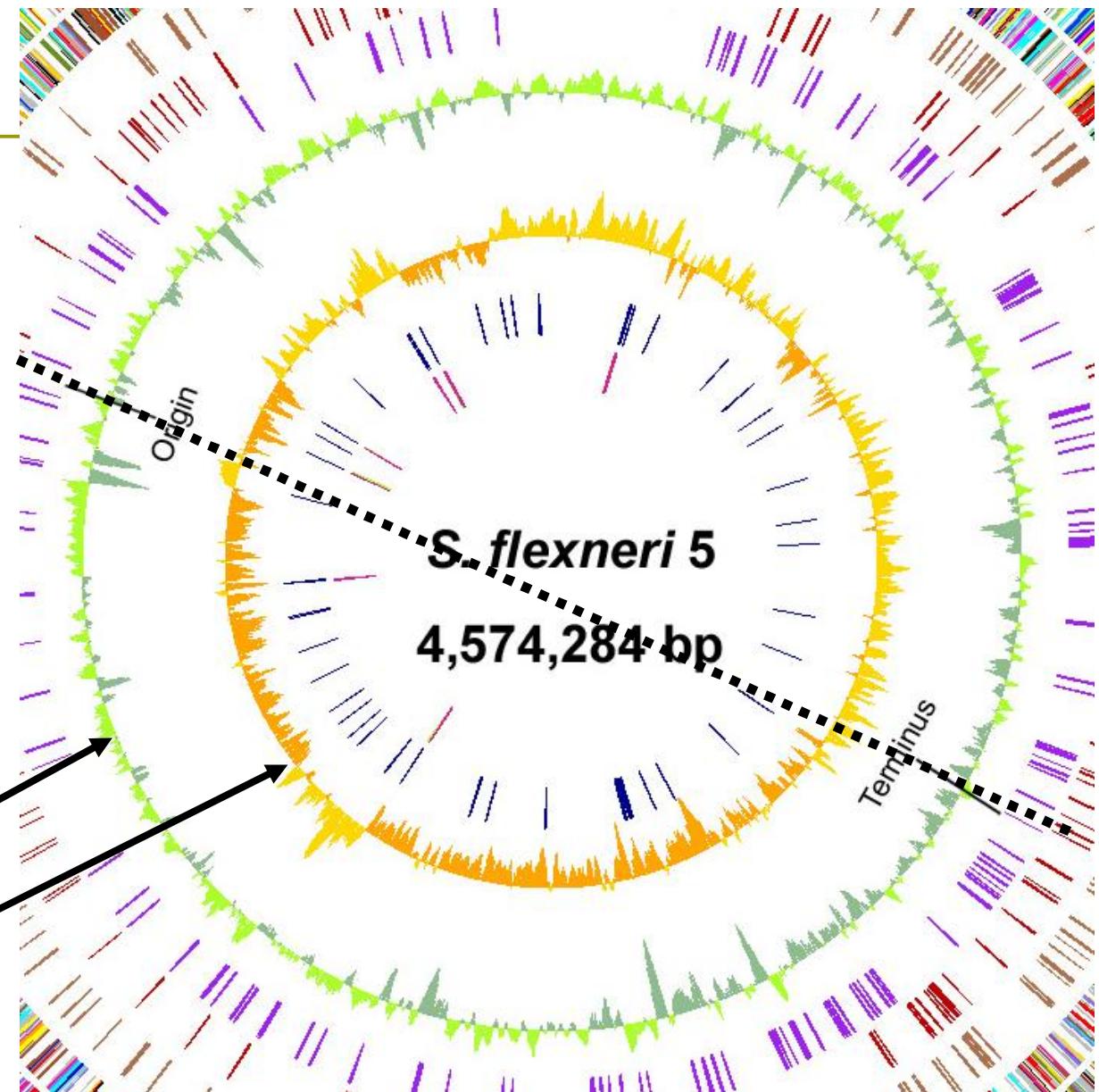
GC skew



Chromosome map of *S. dysenteriae*, the nine rings describe different properties of the genome
http://www.mgc.ac.cn/ShiBASE/circular_Sd197.htm

GC skew

p GC skew often changes sign at origins and termini of replication



Nie et al., BMC Genomics, 2006

2-words: dinucleotides

- Let's consider a sequence L_1, L_2, \dots, L_n where each letter L_i is drawn from the DNA alphabet {A, C, G, T}
- We have 16 possible dinucleotides $L_i L_{i+1}$: AA, AC, AG, ..., TG, TT.

i.i.d. model for nucleotides

- p Assume that bases
 - n occur independently of each other
 - n bases at each position are identically distributed
- p Probability of the base A, C, G, T occurring is p_A , p_C , p_G , p_T , respectively
 - n For example, we could use $p_A=p_C=p_G=p_T=0.25$ or estimate the values from known genome data
- p Probability of $I_i I_{i+1}$ is then $P_{I_i} P_{I_{i+1}}$
 - n For example, $P(TG) = p_T p_G$

2-words: is what we see surprising?

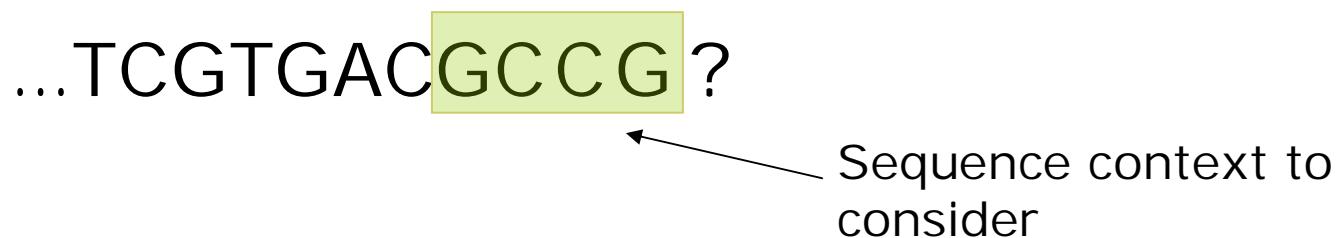
- p We can test whether a sequence is "unexpected", for example, with a χ^2 test
- p Test statistic for a particular dinucleotide r_1r_2 is $\chi^2 = (O - E)^2 / E$ where
 - n O is the observed number of dinucleotide r_1r_2
 - n E is the expected number of dinucleotide r_1r_2
 - n $E = (n - 1)p_{r1}p_{r2}$ under i.i.d. model
- p Basic idea: high values of χ^2 indicate deviation from the model
 - n Actual procedure is more detailed -> basic statistics courses

Refining the i.i.d. model

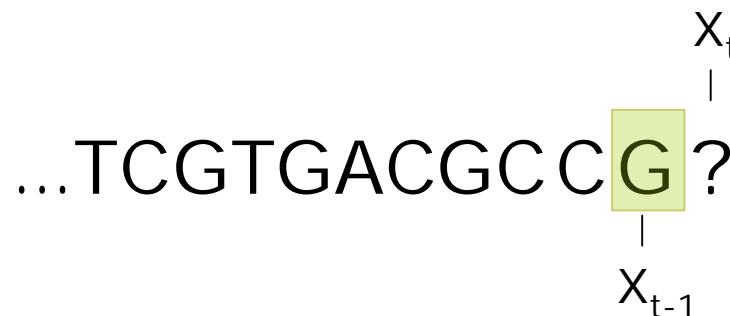
- i.i.d. model describes some organisms well (see Deonier's book) but fails to characterise many others
- We can refine the model by having the DNA letter at some position depend on letters at preceding positions

...TCGTGACGCCG?

Sequence context to consider



First-order Markov chains



- ☐ Lets assume that in sequence X the letter at position t , X_t , depends only on the previous letter X_{t-1} (*first-order markov chain*)
- ☐ Probability of letter j occurring at position t given $X_{t-1} = i$: $p_{ij} = P(X_t = j \mid X_{t-1} = i)$
- ☐ We consider *homogeneous* markov chains: probability p_{ij} is independent of position t

Estimating p_{ij}

- We can estimate probabilities p_{ij} ("the probability that j follows i") from observed dinucleotide frequencies

	A	C	G	T	Frequency of dinucleotide AT in sequence
A	p_{AA}	p_{AC}	p_{AG}	p_{AT}	
C	$p_{CA} + p_{CC} + p_{CG} + p_{CT}$				← Base frequency $fr(C)$
G	p_{GA}	p_{GC}	p_{GG}	p_{GT}	
T	p_{TA}	p_{TC}	p_{TG}	p_{TT}	

...the values $p_{AA}, p_{AC}, \dots, p_{TG}, p_{TT}$ sum to 1

Estimating p_{ij}

$$p_{ij} = P(X_t = j \mid X_{t-1} = i) = \frac{P(X_t = j, X_{t-1} = i)}{P(X_{t-1} = i)}$$

Probability of transition $i \rightarrow j$

$$0.052 / 0.345 \approx 0.151$$

Dinucleotide frequency

$$P(X_{t-1} = i)$$

Base frequency of nucleotide i , $fr(i)$

	A	C	G	T
A	0.146	0.052	0.058	0.089
C	0.063	0.029	0.010	0.056
G	0.050	0.030	0.028	0.051
T	0.086	0.047	0.063	0.140

101

$$P(X_t = j, X_{t-1} = i)$$

	A	C	G	T
A	0.423	0.151	0.168	0.258
C	0.399	0.184	0.063	0.354
G	0.314	0.189	0.176	0.321
T	0.258	0.138	0.187	0.415

$$P(X_t = j \mid X_{t-1} = i)$$

Simulating a DNA sequence

- p From a transition matrix, it is easy to generate a DNA sequence of length n:
 - n First, choose the starting base randomly according to the base frequency distribution
 - n Then, choose next base according to the distribution $P(x_t | x_{t-1})$ until n bases have been chosen

T T C T T C A A

Look for R code in Deonier's book

	A	C	G	T
A	0.423	0.151	0.168	0.258
C	0.399	0.184	0.063	0.354
G	0.314	0.189	0.176	0.321
T	0.258	0.138	0.187	0.415

$$P(X_t = j | X_{t-1} = i)$$

Simulating a DNA sequence

- p Now we can quickly generate sequences of arbitrary length...

Simulating a DNA sequence

	Dinucleotide frequencies	
	Simulated	Observed
aa	0.145	0.146
ac	0.050	0.052
ag	0.055	0.058
at	0.092	0.089
ca	0.065	0.063
cc	0.028	0.029
cg	0.011	0.010
ct	0.058	0.056
ga	0.048	0.050
gc	0.032	0.030
gg	0.029	0.028
gt	0.050	0.051
ta	0.084	0.086
tc	0.052	0.047
tg	0.064	0.063
tt	0.138	0.0140

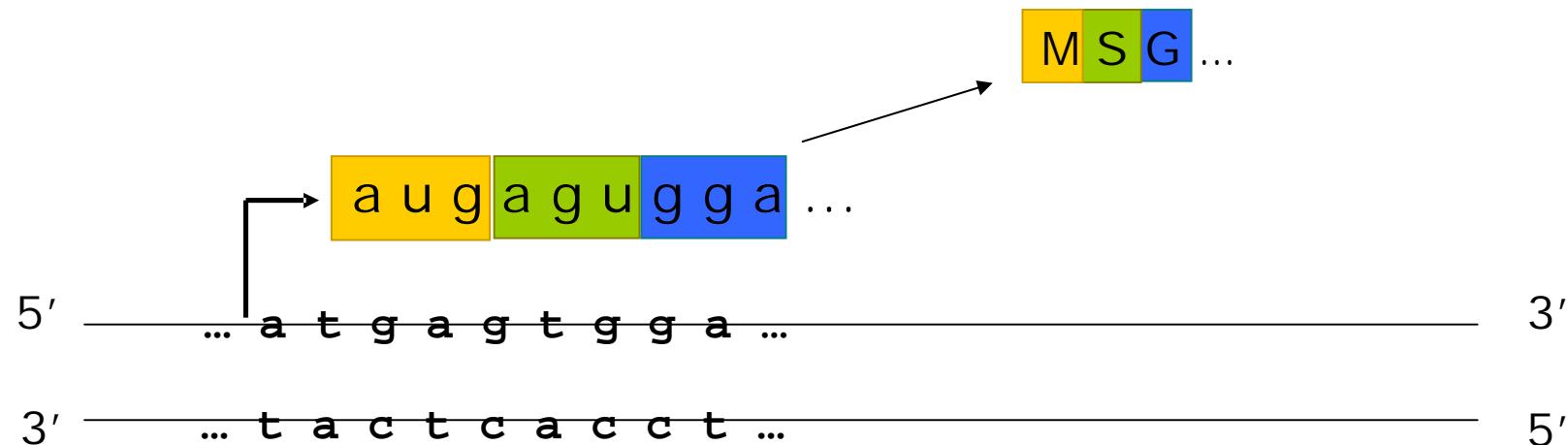
n = 10000

Simulating a DNA sequence

- p The model is able to generate correct proportions of 1- and 2-words in genomes...
 - p ...but fails with $k=3$ and beyond.

3-words: codons

- We can extend the previous method to 3-words
- $k=3$ is an important case in study of DNA sequences because of genetic code



3-word probabilities

- Let's again assume a sequence L of independent bases
- Probability of 3-word $r_1r_2r_3$ at position $i, i+1, i+2$ in sequence L is

$$\begin{aligned} P(L_i = r_1, L_{i+1} = r_2, L_{i+2} = r_3) &= \\ P(L_i = r_1)P(L_{i+1} = r_2)P(L_{i+2} = r_3) \end{aligned}$$

3-words in Escherichia coli genome

Word	Count	Observed	Expected	Word	Count	Observed	Expected
AAA	108924	0.02348	0.01492	CAA	76614	0.01651	0.01541
AAC	82582	0.01780	0.01541	CAC	66751	0.01439	0.01591
AAG	63369	0.01366	0.01537	CAG	104799	0.02259	0.01588
AAT	82995	0.01789	0.01490	CAT	76985	0.01659	0.01539
ACA	58637	0.01264	0.01541	CCA	86436	0.01863	0.01591
ACC	74897	0.01614	0.01591	CCC	47775	0.01030	0.01643
ACG	73263	0.01579	0.01588	CCG	87036	0.01876	0.01640
ACT	49865	0.01075	0.01539	CCT	50426	0.01087	0.01589
AGA	56621	0.01220	0.01537	CGA	70938	0.01529	0.01588
AGC	80860	0.01743	0.01588	CGC	115695	0.02494	0.01640
AGG	50624	0.01091	0.01584	CGG	86877	0.01872	0.01636
AGT	49772	0.01073	0.01536	CGT	73160	0.01577	0.01586
ATA	63697	0.01373	0.01490	CTA	26764	0.00577	0.01539
ATC	86486	0.01864	0.01539	CTC	42733	0.00921	0.01589
ATG	76238	0.01643	0.01536	CTG	102909	0.02218	0.01586
ATT	83398	0.01797	0.01489	CTT	63655	0.01372	0.01537

3-words in Escherichia coli genome

Word	Count	Observed	Expected	Word	Count	Observed	Expected
GAA	83494	0.01800	0.01537	TAA	68838	0.01484	0.01490
GAC	54737	0.01180	0.01588	TAC	52592	0.01134	0.01539
GAG	42465	0.00915	0.01584	TAG	27243	0.00587	0.01536
GAT	86551	0.01865	0.01536	TAT	63288	0.01364	0.01489
GCA	96028	0.02070	0.01588	TCA	84048	0.01812	0.01539
GCC	92973	0.02004	0.01640	TCC	56028	0.01208	0.01589
GCG	114632	0.02471	0.01636	TCG	71739	0.01546	0.01586
GCT	80298	0.01731	0.01586	TCT	55472	0.01196	0.01537
GGA	56197	0.01211	0.01584	TGA	83491	0.01800	0.01536
GGC	92144	0.01986	0.01636	TGC	95232	0.02053	0.01586
GGG	47495	0.01024	0.01632	TGG	85141	0.01835	0.01582
GGT	74301	0.01601	0.01582	TGT	58375	0.01258	0.01534
GTA	52672	0.01135	0.01536	TTA	68828	0.01483	0.01489
GTC	54221	0.01169	0.01586	TTC	83848	0.01807	0.01537
GTG	66117	0.01425	0.01582	TTG	76975	0.01659	0.01534
GTT	82598	0.01780	0.01534	TTT	109831	0.02367	0.01487

2nd order Markov Chains

- p Markov chains readily generalise to higher orders
- p In 2nd order markov chain, position t depends on positions t-1 and t-2
- p Transition matrix:

	A	C	G	T
AA				
AC				
AG				
AT				
CA				
...				

- p Probabilistic models for DNA and amino acid sequences will be discussed in Biological sequence analysis course (II period)

Codon Adaptation Index (CAI)

- p Observation: cells prefer certain codons over others in highly expressed genes
 - n Gene expression: DNA is transcribed into RNA (and possibly translated into protein)

Amino acid	Codon	Predicted	Gene class I	Gene class II	
Phe	TTT	0.493	0.551	0.291	Moderately expressed
	TTC	0.507	0.449	0.709	
Ala	GCT	0.246	0.145	0.275	Highly expressed
	GCC	0.254	0.276	0.164	
Asn	GCA	0.246	0.196	0.240	
	GCG	0.254	0.382	0.323	
	AAT	0.493	0.409	0.172	
	AAC	0.507	0.591	0.828	

Codon frequencies for some genes in E. coli

Codon Adaptation Index (CAI)

CAI is a statistic used to compare the distribution of codons **observed** with the **preferred** codons for highly expressed genes

Codon Adaptation Index (CAI)

- p Consider an amino acid sequence $X = x_1x_2\dots x_n$
- p Let p_k be the probability that codon k is used in highly expressed genes
- p Let q_k be the highest probability that a codon coding for the same amino acid as codon k has
 - n For example, if codon k is "GCC", the corresponding amino acid is Alanine (see genetic code table; also GCT, GCA, GCG code for Alanine)
 - n Assume that $p_{GCC} = 0.164$, $p_{GCT} = 0.275$, $p_{GCA} = 0.240$, $p_{GCG} = 0.323$
 - n Now $q_{GCC} = q_{GCT} = q_{GCA} = q_{GCG} = 0.323$

Codon Adaptation Index (CAI)

p CAI is defined as

$$\text{CAI} = \left(\prod_{k=1}^n p_k / q_k \right)^{1/n}$$

p CAI can be given also in *log-odds* form:

$$\log(\text{CAI}) = (1/n) \sum_{k=1}^n \log(p_k / q_k)$$

CAI: example with an E. coli gene

q_k
 p_k

M	A	L	T	K	A	E	M	S	E	Y	L	...
ATG	GCG	CTT	ACA	AAA	GCT	GAA	ATG	TCA	GAA	TAT	CTG	
1.00	0.47	0.02	0.45	0.80	0.47	0.79	1.00	0.43	0.79	0.19	0.02	
	0.06	0.02	0.47	0.20	0.06	0.21		0.32	0.21	0.81	0.02	
	0.28	0.04	0.04		0.28			0.03			0.04	
	0.20	0.03	0.05		0.20			0.01			0.03	
		0.01						0.04			0.01	
		0.89						0.18			0.89	
ATG	GCT	TTA	ACT	AAA	GCT	GAA	ATG	TCT	GAA	TAT	TTA	
	GCC	TTG	ACC	AAG	GCC	GAG		TCC	GAG	TAC	TTG	
	GCA	CTT	ACA		GCA			TCA			CTT	
	GCG	CTC	ACG		GCG			TCG			CTC	
		CTA						AGT			CTA	
		CTG						AGC			CTG	$\frac{1}{n}$
1.00	0.20	0.04	0.04	0.80	0.47	0.79	1.00	0.03	0.79	0.19	0.89...	
1.00	0.47	0.89	0.47	0.80	0.47	0.79	1.00	0.43	0.79	0.81	0.89	

CAI: properties

- p CAI = 1.0 : each codon was the most frequently used codon in highly expressed genes
- p Log-odds used to avoid numerical problems
 - n What happens if you multiply many values <1.0 together?
- p In a sample of E.coli genes, CAI ranged from 0.2 to 0.85
- p CAI correlates with mRNA levels: can be used to predict high expression levels

Biological words: summary

- p Simple 1-, 2- and 3-word models can describe interesting properties of DNA sequences
 - n GC skew can identify DNA replication origins
 - n It can also reveal *genome rearrangement* events and *lateral transfer* of DNA
 - n GC content can be used to locate genes: human genes are comparably GC-rich
 - n CAI predicts high gene expression levels

Biological words: summary

- p k=3 models can help to identify correct *reading frames*
 - n Reading frame starts from a start codon and stops in a stop codon
 - n Consider what happens to translation when a single extra base is introduced in a reading frame
- p Also word models for $k > 3$ have their uses

Next lecture

- p Genome sequencing & assembly – where do we get sequence data?

Note on programming languages

- p Working with probability distributions is straightforward with R, for example
 - n Deonier's book contains many computational examples
 - n You can use R in CS Linux systems
- p Python works too!

```
#!/usr/bin/env python
```

```
import sys, random
```

```
n = int(sys.argv[1])
```

```
tm = {'a' : {'a' : 0.423, 'c' : 0.151, 'g' : 0.168, 't' : 0.258},  
      'c' : {'a' : 0.399, 'c' : 0.184, 'g' : 0.063, 't' : 0.354},  
      'g' : {'a' : 0.314, 'c' : 0.189, 'g' : 0.176, 't' : 0.321},  
      't' : {'a' : 0.258, 'c' : 0.138, 'g' : 0.187, 't' : 0.415} }
```

```
pi = {'a' : 0.345, 'c' : 0.158, 'g' : 0.159, 't' : 0.337}
```

```
def choose(dist):  
    r = random.random()  
    sum = 0.0  
    keys = dist.keys()  
    for k in keys:  
        sum += dist[k]  
        if sum > r:  
            return k  
    return keys[-1]
```

```
c = choose(pi)  
for i in range(n - 1):  
    sys.stdout.write(c)  
    c = choose(tm[c])  
sys.stdout.write(c)  
sys.stdout.write("\n")
```

Example Python code for generating DNA sequences with first-order Markov chains.

Initialisation: use packages 'sys' and 'random', read sequence length from input.

Transition matrix tm and initial distribution pi.

Function choose(), returns a key (here 'a', 'c', 'g' or 't') of the dictionary 'dist' chosen randomly according to probabilities in dictionary values.

Choose the first letter, then choose next letter according to $P(x_t | x_{t-1})$.