

Regulatory sequence analysis

Sequence models

(Bernoulli and Markov models)

Why do we need random models ?

- Any pattern discovery relies on an underlying model to estimate the random expectation.
 - This model can be simple (succession of independent and equiprobable nucleotides) or more elaborate (differences in oligonucleotide composition).
 - The choice of an inappropriate model can lead to false conclusions.
 - In practice, a sequence model can be used to generate random sequences, which will serve to validate some theoretical assumptions.
- Example: comparison of observed and expected occurrences with the binomial distribution, as applied with oligo-analysis :
 - Relies on an assumption that successive oligonucleotides are independent from each other.
 - This is clearly not the case: each k-letter word depends on the k-1 neighbour words on both sides. How far does it affect the conclusions ?
 - We could test it by generating random sequences, counting words, and fitting the distribution of observed occurrences with a binomial distribution.

Probability of a sequence segment

- What is the probability for a given sequence segment (oligonucleotide, “word”) to be found at a given position of a DNA sequence ?
- Different models can be chosen
 - **Bernoulli model**
 - Assumes independence between successive nucleotides.
 - The probability of each residue is fixed a priori (*prior residue probability*)
 - Example: $P(A) = 0.35$; $P(T) = 0.32$; $P(C) = 0.17$; $P(G) = 0.16$
 - Particular case: equiprobable residues
 - $P(A) = P(T) = P(C) = P(G) = 0.25$
 - Simple, but **NOT realistic** !
 - **Markov model**
 - The probability of each residue depends on the ***m*** preceding residues.
 - The parameter ***m*** is called the *order* of the Markov model
 - Remark: a Markov model of order 0 is a Bernoulli model.

- The simplest model : Bernoulli with identically and independently (i.i.d.) distributed nucleotides.

$$p = P(A) = P(C) = P(G) = P(T) = 0.25$$

$$P(S) = p^L$$

- The probability of a sequence
 - Is the product of its residue probabilities (independence)
 - Equiprobability: since all residues have the same probability, it is simply computed as the residue proba (p) to the power of the sequence length (L)
 - S is a sequence segment (e.g. an oligonucleotide)
 - L length of the sequence segment
 - p nucleotide probability
 - $P(S)$ is the probability to observe this sequence segment at given position of a larger sequence
- Example
 - $P(\text{CACGTG}) = 0.25^6 = 2.44 \times 10^{-4}$

Bernoulli model : independently distributed nucleotides

- A more refined model consists in using residue-specific probabilities. The probability of each residue is assumed to be constant on the whole sequence (Bernoulli schema).
- The probability of a sequence is the product of its residue probabilities.
 - $i = 1..k$ is the index of nucleotide positions
 - r_i is the residue found at position i
 - $P(r_i)$ is the probability of this residue
- Example: non-coding sequences in the yeast genome
 - $P(A) = P(T) = 0.325$
 - $P(C) = P(G) = 0.175$
 - $P(CACGTG) = P(C) P(A) P(C) P(G) P(T) P(G)$
 $= 0.325^4 * 0.175^2$
 $= 9.91E^{-5}$

$$P(S) = \prod_{i=1}^L P(r_i)$$

Markov chains and transition matrices

$$P(r_i | S_{i-m,i-1})$$

Transition matrix, order 1

	a	c	g	t
A	P(A A)	P(C A)	P(G A)	P(T A)
C	P(A C)	P(C C)	P(G C)	P(T C)
G	P(A G)	P(C G)	P(G G)	P(T G)
T	P(A T)	P(C T)	P(G T)	P(T T)

Transition matrix, order 2

Pref	A	C	G	T
AA	P(A AA)	P(C AA)	P(G AA)	P(T AA)
AC	P(A AC)	P(C AC)	P(G AC)	P(T AC)
AG	P(A AG)	P(C AG)	P(G AG)	P(T AG)
AT	P(A AT)	P(C AT)	P(G AT)	P(T AT)
CA	P(A CA)	P(C CA)	P(G CA)	P(T CA)
CC	P(A CC)	P(C CC)	P(G CC)	P(T CC)
CG	P(A CG)	P(C CG)	P(G CG)	P(T CG)
CT	P(A CT)	P(C CT)	P(G CT)	P(T CT)
GA	P(A GA)	P(C GA)	P(G GA)	P(T GA)
GC	P(A GC)	P(C GC)	P(G GC)	P(T GC)
GG	P(A GG)	P(C GG)	P(G GG)	P(T GG)
GT	P(A GT)	P(C GT)	P(G GT)	P(T GT)
TA	P(A TA)	P(C TA)	P(G TA)	P(T TA)
TC	P(A TC)	P(C TC)	P(G TC)	P(T TC)
TG	P(A TG)	P(C TG)	P(G TG)	P(T TG)
TT	P(A TT)	P(C TT)	P(G TT)	P(T TT)

- In a Markov model, the probability to find a letter at position i depends on the residues found at the m preceding residues.
- The tables represent the transition matrices for Markov chain models of order $m=1$ (top) and $m=2$ (bottom).
- Each row specifies one **prefix**, each column one **suffix**.
- The values indicate the probability to observe a given residue (suffix r_i) at position (i) of the sequence, as a function of the m preceding residues (the prefix $S_{i-m,i-1}$)
- Particular case
 - A Bernoulli model is a Markov model of order 0.

Markov model estimation (“training”)

- Transition frequencies for a Markov model of order m can be estimated from the frequencies observed for oligomers (k-mers) of length $k=m+1$ in a reference sequence set.
- Example
 - The upper table shows dinucleotide frequencies ($k=2$) computed from the whole set of upstream sequences of the yeast *Saccharomyces cerevisiae*.
 - This table can be used to estimate a Markov model of order $m = k-1 = 1$.

Dinucleotide frequencies		
Sequences	Occurrences	Frequency
S	N(S)	F(S)
AA	526,149	0.112
AC	251,377	0.054
AG	275,056	0.059
AT	414,453	0.088
CA	294,423	0.063
CC	178,324	0.038
CG	146,052	0.031
CT	275,859	0.059
GA	277,343	0.059
GC	184,367	0.039
GG	173,404	0.037
GT	239,569	0.051
TA	369,980	0.079
TC	280,475	0.060
TG	279,932	0.060
TT	521,236	0.111

$$P(r_i | S_{1..m}) = \frac{F_{bg}(r_i | S_{1..m})}{\sum_{j \in A} F_{bg}(r_j | S_{1..m})} = \frac{F_{bg}(S_{1..m} r_i)}{\sum_{j \in A} F_{bg}(S_{1..m} r_j)}$$

$$\begin{aligned}
 P(G|T) &= \frac{F(G|T)}{\sum_{j \in A} F(j|T)} = \frac{F(TG)}{F(T^*)} \\
 &= \frac{0.060}{0.079 + 0.060 + 0.060 + 0.111} \\
 &= \frac{0.060}{0.310} = 0.194
 \end{aligned}$$

Transition matrix, order 1						
Prefix \ Suffix	A	C	G	T	P(Prefix)	N(Suffix)
A	0.359	0.171	0.187	0.283	0.313	1,467,035
C	0.329	0.199	0.163	0.308	0.191	894,658
G	0.317	0.211	0.198	0.274	0.187	874,683
T	0.255	0.193	0.193	0.359	0.310	1,451,623
P(Suffix)	0.313	0.191	0.187	0.310		
N(Suffix)	1,467,895	894,543	874,444	1,451,117		

Examples of transition matrices

- The two tables below show the transition matrices for a Markov model of order 1 (top) and 2 (bottom), respectively.
- Trained with the whole set of **non-coding upstream sequences** of the yeast *Saccharomyces cerevisiae*.
- Notice the high probability of transitions from **AA to A** and **TT to T**.

$$P(r_i | S_{i-m,i-1})$$

Pre/Suffix	A	C	G	T	P(Prefix)
a	0.371	0.165	0.178	0.285	0.321
c	0.327	0.190	0.167	0.316	0.183
g	0.312	0.214	0.189	0.285	0.177
t	0.273	0.179	0.173	0.375	0.320
Sym	1.283	0.748	0.708	1.261	
P(suffix)	0.321	0.183	0.176	0.320	

Prefix/Suffix	A	C	G	T	P(Prefix)
aa	0.416	0.151	0.187	0.246	0.119
ac	0.352	0.181	0.171	0.297	0.053
ag	0.339	0.202	0.193	0.267	0.057
at	0.346	0.166	0.162	0.326	0.092
ca	0.344	0.185	0.180	0.291	0.060
cc	0.305	0.200	0.171	0.324	0.035
cg	0.282	0.232	0.193	0.294	0.031
ct	0.241	0.189	0.184	0.385	0.058
ga	0.411	0.144	0.187	0.257	0.055
gc	0.334	0.192	0.182	0.293	0.038
gg	0.315	0.220	0.194	0.271	0.033
gt	0.307	0.156	0.200	0.338	0.050
ta	0.304	0.184	0.160	0.352	0.087
tc	0.313	0.192	0.152	0.343	0.057
tg	0.300	0.214	0.180	0.307	0.055
tt	0.218	0.194	0.164	0.423	0.120
Sum	5.127	3.000	2.860	5.013	
P(suffix)	0.321	0.183	0.176	0.319	

Pi	a	c	g	t
a	0.371	0.165	0.178	0.285
c	0.327	0.190	0.167	0.316
g	0.312	0.214	0.189	0.285
t	0.273	0.179	0.173	0.375

Pi	a	c	g	t
aa	0.416	0.151	0.187	0.246
ac	0.352	0.181	0.171	0.297
ag	0.339	0.202	0.193	0.267
at	0.346	0.166	0.162	0.326
ca	0.344	0.185	0.180	0.291
cc	0.305	0.200	0.171	0.324
cg	0.282	0.232	0.193	0.294
ct	0.241	0.189	0.184	0.385
ga	0.411	0.144	0.187	0.257
gc	0.334	0.192	0.182	0.293
gg	0.315	0.220	0.194	0.271
gt	0.307	0.156	0.200	0.338
ta	0.304	0.184	0.160	0.352
tc	0.313	0.192	0.152	0.343
tg	0.300	0.214	0.180	0.307
tt	0.218	0.194	0.164	0.423

Markov chains and Bernoulli models

- By extension of the concept of Markov chain, Bernoulli models can be qualified as Markov models of order 0 (the order 0 means that there is no dependency between a residue and the preceding ones).
- The prior probabilities of a Markov model of order $m=0$ can be estimated from the residue of single nucleotides ($k=m+1=1$) in a background sequence set.
- The table below shows the residue frequencies in the genomes of the yeast *Saccharomyces cerevisiae* and the bacteria *Escherichia coli* K12, respectively.
- Notice the strong differences between these genomes.

Markov order 0 = Bernoulli

A	C	G	T	Genome
0.310	0.191	0.191	0.309	<i>Saccharomyces cerevisiae</i>
0.246	0.254	0.254	0.246	<i>Escherichia coli</i> K12

Scoring a sequence segment with a Markov model

- The example below illustrates the computation of the probability of a sequence segment $P(S|B)$ with a background Markov model B of order 2, calibrated from 3nt frequencies on the yeast genome.

CCTACTATATGCCCAGAATT

Background model B

Transition matrix, order 2

Prefix/Suffix	A	C	G	T	P(Prefix)	N(Prefix)
AA	0.388	0.161	0.200	0.251	0.112	525,000
AC	0.339	0.198	0.173	0.290	0.054	251,072
AG	0.345	0.204	0.196	0.255	0.059	274,601
AT	0.311	0.184	0.182	0.323	0.088	413,946
CA	0.347	0.178	0.189	0.286	0.063	293,750
CC	0.341	0.190	0.161	0.309	0.038	178,110
CG	0.293	0.221	0.196	0.290	0.031	145,876
CT	0.229	0.195	0.205	0.371	0.059	275,634
GA	0.394	0.155	0.187	0.264	0.059	277,053
GC	0.330	0.205	0.169	0.297	0.039	184,192
GG	0.318	0.217	0.187	0.277	0.037	173,266
GT	0.285	0.175	0.204	0.336	0.051	239,384
TA	0.300	0.193	0.168	0.339	0.079	369,426
TC	0.313	0.203	0.152	0.332	0.060	280,131
TG	0.302	0.209	0.208	0.282	0.060	279,783
TT	0.210	0.208	0.189	0.392	0.111	520,906
P(Suffix)	0.313	0.191	0.187	0.310		
N(suffix)	1,466,075	893,444	873,260	1,449,351		

Sequence probability given the background model

$$P(S|B) = P(S_{1,m}|B) \prod_{i=m+1}^L P(r_i | S_{i-m,i-1}, B)$$

pos	P(R W)	wR	S	P(S)
1	P(CC)	0.038 cc	CC	3.80E-02
2	P(T CC)	0.309 ccT	CCT	1.17E-02
3	P(A CT)	0.229 ctA	CCTA	2.69E-03
4	P(C TA)	0.193 taC	CCTAC	5.19E-04
5	P(T AC)	0.290 acT	CCTACT	1.50E-04
6	P(A CT)	0.229 ctA	CCTACTA	3.45E-05
7	P(T TA)	0.339 taT	CCTACTAT	1.17E-05
8	P(A AT)	0.311 atA	CCTACTATA	3.63E-06
9	P(T TA)	0.339 taT	CCTACTATAT	1.23E-06
10	P(G AT)	0.182 atG	CCTACTATATG	2.25E-07
11	P(C TG)	0.209 tgC	CCTACTATATGC	4.69E-08
12	P(C GC)	0.205 gcC	CCTACTATATGCC	9.61E-09
13	P(C CC)	0.190 ccC	CCTACTATATGCCC	1.82E-09
14	P(A CC)	0.341 ccA	CCTACTATATGCCCA	6.21E-10
15	P(G CA)	0.189 caG	CCTACTATATGCCCAG	1.17E-10
16	P(A AG)	0.345 agA	CCTACTATATGCCCAGA	4.04E-11
17	P(A GA)	0.394 gaA	CCTACTATATGCCCAGAA	1.59E-11
18	P(T AA)	0.251 aaT	CCTACTATATGCCCAGAAT	4.00E-12
19	P(T AT)	0.323 atT	CCTACTATATGCCCAGAATT	1.29E-12

Background sequences

- The frequencies observed for a k -letter word in a reference sequence set (background sequence) can be used to estimate the expected frequencies of the same k -letter word in the sequences to be analyzed.
- Typical background models:
 - Whole genome
 - But this will bias the estimates towards coding frequencies, especially in microbial organisms, where the majority of the genome is coding.
 - Whole set of intergenic sequences
 - More accurate than whole-genome estimates, but still biased because intergenic sequences include both upstream and downstream sequences
 - Whole set of upstream sequences, same sizes as the sequences to be analyzed
 - Requires a calibration for each sequence size
 - Whole set of upstream sequences, fixed size (default on the web site)
 - Reasonably good estimate for microbes, NOT for higher organisms.

Nucleotide composition of the *Plasmodium* upstream sequences

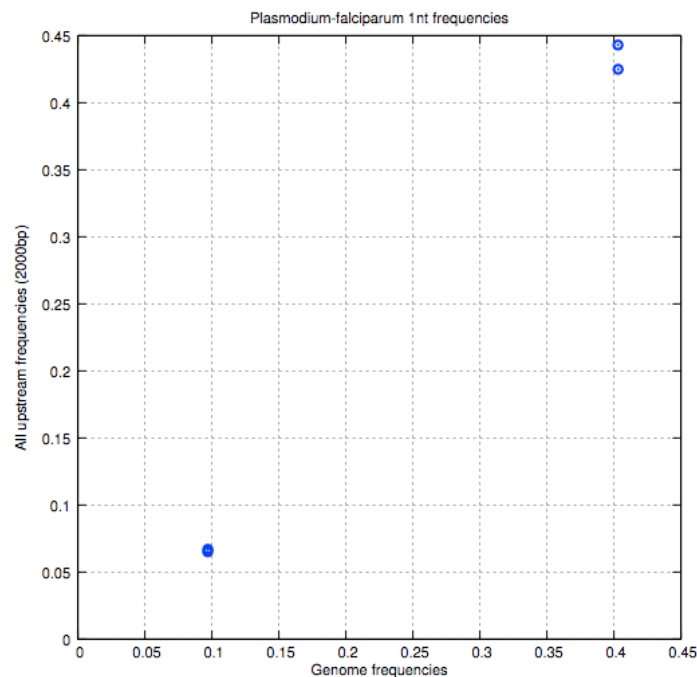
- The genome shows a strong richness in A and T residues (80%AT).
- This enrichment is even stronger in upstream non-coding sequences (86%AT).

Frequencies

Residue	Genome	Upstream (max 2kb)
a	0.40	0.42
c	0.10	0.07
g	0.10	0.07
t	0.40	0.44

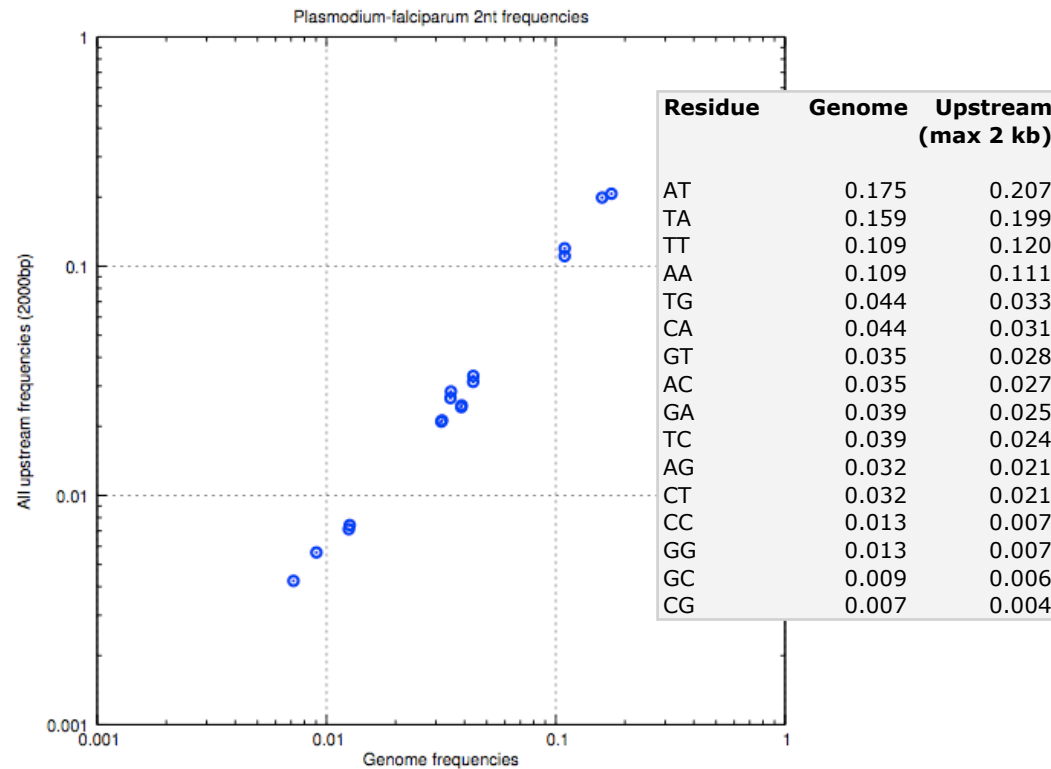
Upstream frequencies Gray scale

a	c	g	t
0.42419	0.06761	0.06951	0.43870



Dinucleotide composition of the *Plasmodium* upstream sequences

- Dinucleotide frequencies reflect the AT-richness.



Markov order m=1
(based of dinucleotides k=2)

pr	a			
	a	c	g	t
a	0.39983	0.06522	0.05243	0.48253
c	0.47917	0.13198	0.06745	0.32140
g	0.36686	0.08719	0.12634	0.41961
t	0.44833	0.05728	0.07758	0.41681

Transition frequencies

Markov order m=3
(based on tetranucleotides k=4)

- On the basis of oligonucleotide frequencies, one can compute Markov models, which indicate the probability to observe a certain residue (suffix) after a certain oligonucleotide (prefix).
- The Markov model can be represented in the form of a transition table.

Markov order m=1
(based on dinucleotides k=2)

Pr	a	c	g	t
a	0.39983	0.06522	0.05243	0.48253
c	0.47917	0.13198	0.06745	0.32140
g	0.36686	0.08719	0.12634	0.41961
t	0.44833	0.05728	0.07758	0.41681

Markov order m=2
(based on trinucleotides k=3)

Pr	a	c	g	t
aa	0.55133	0.05771	0.06855	0.32241
ac	0.59611	0.10571	0.07782	0.22037
ag	0.45875	0.08478	0.16910	0.28737
at	0.60625	0.03631	0.08099	0.27645
ca	0.32382	0.11650	0.04995	0.50973
cc	0.36057	0.16777	0.04837	0.42329
cg	0.32114	0.10703	0.08672	0.48511
ct	0.28437	0.09985	0.07587	0.53991
ga	0.54541	0.06894	0.09101	0.29464
gc	0.49165	0.11509	0.08103	0.31223
gg	0.43450	0.07578	0.16309	0.32663
gt	0.43973	0.06254	0.13930	0.35844
ta	0.26328	0.06301	0.03410	0.63961
tc	0.38894	0.15242	0.05942	0.39921
tg	0.29671	0.08900	0.09461	0.51969
tt	0.29185	0.07502	0.06441	0.56871

Pr	a	c	g	t
aaa	0.60977	0.04557	0.06493	0.27973
aac	0.58765	0.10970	0.06361	0.24944
aag	0.48695	0.07425	0.13023	0.28943
aat	0.58314	0.03951	0.09415	0.28320
aca	0.32624	0.11699	0.04732	0.49945
acc	0.42056	0.13956	0.05166	0.38822
acg	0.31069	0.10224	0.09021	0.47886
act	0.31955	0.08427	0.09022	0.49599
aga	0.59494	0.05904	0.09664	0.24934
agc	0.50259	0.12264	0.07180	0.30297
agg	0.47636	0.07224	0.15773	0.28966
agt	0.45732	0.06399	0.13495	0.34404
ata	0.24700	0.05147	0.02792	0.67361
atc	0.44741	0.14610	0.06464	0.34182
atg	0.31388	0.08453	0.09069	0.51090
att	0.36935	0.07058	0.06213	0.49790
caa	0.50256	0.09844	0.08040	0.31860
cac	0.60914	0.10891	0.07457	0.20738
cag	0.41695	0.10501	0.14426	0.33467
cat	0.52552	0.05702	0.09066	0.32679
cca	0.33910	0.12537	0.05311	0.48242
ccc	0.31962	0.24652	0.03963	0.39424
cgg	0.31798	0.11607	0.07717	0.48888
cct	0.26231	0.10092	0.06513	0.57163
cga	0.49791	0.09494	0.08617	0.32098
cgc	0.48094	0.13045	0.08932	0.29929
cgg	0.39413	0.10204	0.14089	0.36294
cgt	0.40734	0.07709	0.13139	0.38418
cta	0.29871	0.12770	0.05561	0.51798
ctc	0.36301	0.16743	0.05648	0.41309
ctg	0.30288	0.10998	0.09812	0.48912
ctt	0.28491	0.10316	0.07477	0.53716
gaa	0.50960	0.09393	0.09071	0.30076
gac	0.50425	0.12584	0.10360	0.26531
gag	0.48604	0.10232	0.17386	0.23777
gat	0.45644	0.07266	0.12443	0.30647
gca	0.29541	0.15864	0.06259	0.48336
gcc	0.36130	0.16281	0.05062	0.41719
gcg	0.32615	0.14755	0.08999	0.43630
gct	0.30694	0.12070	0.08886	0.48150
gga	0.54253	0.07813	0.09408	0.28520
ggc	0.46171	0.11411	0.08953	0.33466
ggg	0.41850	0.07394	0.23191	0.27560
ggt	0.30290	0.08119	0.14861	0.37731
gta	0.29983	0.11245	0.06728	0.52044
gtc	0.41194	0.17140	0.07744	0.33922
gtg	0.26188	0.12597	0.10396	0.50419
gtt	0.32069	0.12418	0.10655	0.44857
taa	0.46627	0.06216	0.06566	0.40592
tac	0.60619	0.09894	0.07071	0.22426
tag	0.41057	0.09261	0.15341	0.34341
tat	0.63289	0.03019	0.07157	0.28535
tca	0.30623	0.10005	0.04961	0.54411
tcc	0.33077	0.15949	0.04747	0.46227
tcg	0.30935	0.09638	0.06339	0.50889
tct	0.25488	0.10510	0.06881	0.57121
tga	0.50400	0.07174	0.08498	0.33929
tgc	0.49330	0.10821	0.06360	0.31489
tgg	0.39794	0.07761	0.14156	0.38289
tgt	0.44502	0.05733	0.14044	0.35722
tta	0.28876	0.07117	0.03837	0.60169
ttc	0.35959	0.15103	0.05488	0.43451
ttg	0.28338	0.07849	0.09647	0.54166
ttt	0.24799	0.06960	0.06044	0.62198

Markov models show strong variations between organisms

Saccharomyces cerevisiae
(Fungus)

Pr	a	c	g	t
a	0.37000	0.16588	0.17908	0.28504
c	0.32610	0.19058	0.16818	0.31514
g	0.31163	0.21456	0.18957	0.28424
t	0.27256	0.17991	0.17364	0.37389

Escherichia coli K12
(Proteobacteria)

Pr	a	c	g	t
a	0.34491	0.18156	0.17676	0.29677
c	0.30806	0.21557	0.22129	0.25507
g	0.27123	0.25972	0.21545	0.25360
t	0.24080	0.19176	0.21144	0.35599

Mycobacterium leprae
(Actinobacteria)

Pr	a	c	g	t
a	0.23239	0.28694	0.25692	0.22375
c	0.24574	0.24601	0.30574	0.20252
g	0.21748	0.29238	0.25535	0.23479
t	0.18806	0.26081	0.31784	0.23329

Mycoplasma genitalium
(Firmicute, intracellular)

Pr	a	c	g	t
a	0.45565	0.11743	0.13602	0.29091
c	0.39457	0.13008	0.06403	0.41132
g	0.31505	0.18738	0.12047	0.37710
t	0.32450	0.09573	0.11934	0.46044

Bacillus subtilis
(Firmicute, extracellular)

Pr	a	c	g	t
a	0.38159	0.13935	0.18767	0.29139
c	0.33699	0.19499	0.16508	0.30293
g	0.34249	0.18100	0.23541	0.24110
t	0.25122	0.17199	0.19402	0.38278

Plasmodium falciparum
(Apicomplexa, intracellular)

Pr	a	c	g	t
a	0.39821	0.06446	0.05206	0.48527
c	0.47798	0.13336	0.06695	0.32171
g	0.36764	0.08587	0.12431	0.42217
t	0.44739	0.05676	0.07673	0.41912

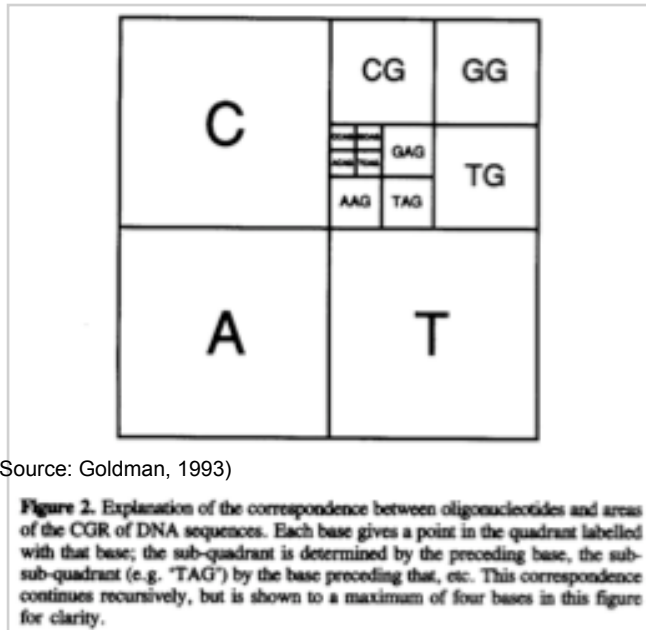
Anopheles gambiae
(Insect)

Pr	a	c	g	t
a	0.34603	0.21388	0.18890	0.25119
c	0.31499	0.21232	0.24159	0.23109
g	0.26036	0.25414	0.20275	0.28275
t	0.20368	0.20710	0.24970	0.33951

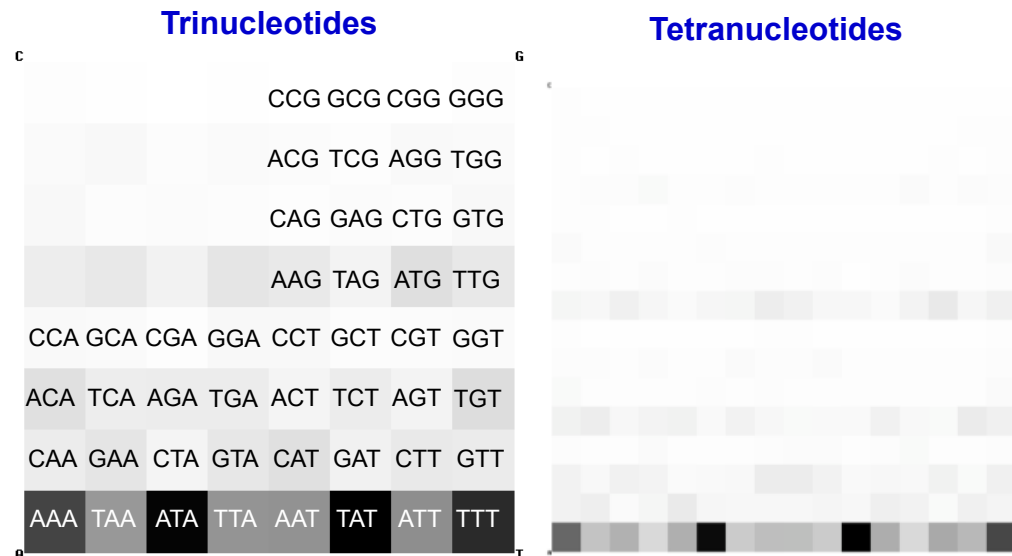
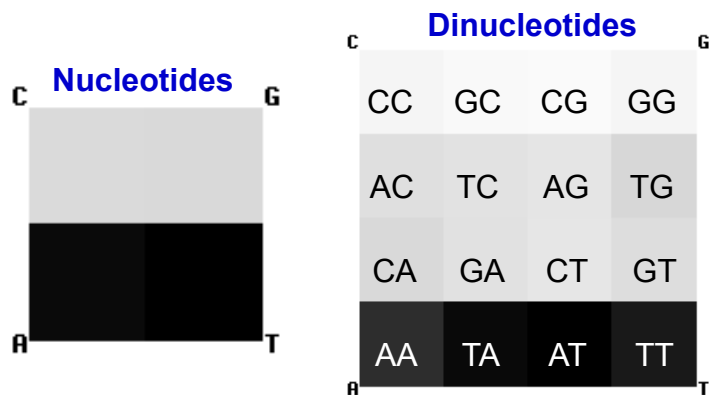
Homo sapiens
(Mammalian)

Pr	a	c	g	t
a	0.29760	0.19031	0.28856	0.22353
c	0.28019	0.30209	0.11692	0.30080
g	0.24408	0.24738	0.30309	0.20545
t	0.18589	0.23061	0.27491	0.30859

Chaos representation - upstream frequencies



- The chaos representation (Jeffrey, 1990) permits to visualize oligonucleotide frequencies and detect enrichment in particular ones.
- Plasmodium upstreall sequences are particularly rich for the following motifs
 - A, T nucleotides
 - Oligonucleotides made of As and Ts only (last row pf each chaos map)
 - Poly-A and poly-T oligos (bottom corners of the maps)
 - (TA){n} motifs (the darkest boxes from dinucleotides to tetranucleotides).



Goldman. Nucleotide, dinucleotide and trinucleotide frequencies explain patterns observed in chaos game representations of DNA sequences. Nucleic Acids Res (1993) vol. 21 (10) pp. 2487-91

Jeffrey. Chaos game representation of gene structure. Nucleic Acids Res (1990) vol. 18 (8) pp. 2163-70

Hexanucleotide frequencies in *Plasmodium* Genome versus upstream (2Kb)

- Hexanucleotides show a very wide range of frequencies in the whole genome (X axis) as well as in the subset of upstream sequences (max 2kb, Y axis).

