











































Models of Sequences

- Consists of states (boxes) and transitions (arcs) labelled with probabilities.
- States have probabilities of "emitting" an element of a sequence (or nothing).
- Arcs have transitional probabilities of moving from one state to another.
 - Sum of probabilities of arcs out of a state must be 1
 - Self-loops are allowed.

Markov Chain A sequence is said to be Markovian if the probability of the occurrence of an element in a particular position depends only on the previous elements in the sequence. Order of a Markov chain depends on how many previous elements influence the

- many previous elements influence the probability:
 - **0**th **order**: uniform probability at every position
 - 1st order: probability depends only on the immediately previous position.

Looking for CpG Islands

Example:

- A CpG island in humans refers to the dinucleotide CG and not the basepair CG.
- The C of CpG is generally methylated to inactivate genes hence CpG is found around "start" regions of many genes more often than elsewhere.
- Methylated C is easily mutated into T.





Looking for CpG Islands

- CpG islands are therefore rare in other locations
- CpG islands are generally a few hundred base pairs long
 Questions:
- 1. Given a short DNA fragment, does it come from a **CpG island** or not?
- 2. Given a long unannotated sequence of DNA, how do we find the CpG islands?

Building an HMM for CpG Islands

- A set of human sequences were considered and 48 CpG islands were tabulated.
- Two Markov chain models were built:
 - One for the regions labeled as CpG islands (the '+' model or Model 1)
 - One for the remainder of the sequences (the '-' model or Model 2).

Transition Probabilities The transition probabilities of each model were computed by: $a_{st}^{+} = \frac{c_{st}^{+}}{\sum_{t'} c_{st'}^{+}} \qquad a_{st}^{-} = \frac{c_{st}^{-}}{\sum_{t'} c_{st'}^{-}}$

 C_{st}^+ is the number of times letter t followed letter s in the plus model.

ansition Fre	Marl	kov chain M	Aodel 1	
ansition Fre	guencies wi	1 40 4		
Instition Fre		TININ / N INIT	otivo CoGi	londs in hu
	queneres wi	tmn 48 put	anve CpG i	siands in hu
+	Α	C	G	Т
A	0.180	0.274	4.426	0.120
С	0.171	0.368	0.274	0.188
G	0.161	0.339	4.375	0.125
Т	0.079	0.355	0.384	0.182
	Mar	kov chain encies in N	Model 2 Non CpG is	sland DNA
Trans	ation neque			
Trans		С	G	Т
Trans		C 0.205	G	T 0.210
Trans	A 0.300 0.322	C 0.205 0.298	G - 0:205 - 0.078	T 0.210 0.302
Trans	A 0.300 0.322 0.248	C 0.205 0.298 0.246	G 8:205 0.078 4.298	T 0.210 0.302 0.208



Log Likelihood Ratios $S(x) = \log \frac{P(x \mid the + \text{model})}{P(x \mid the - \text{model})} = \sum_{i=1}^{L} \log \frac{a_{x_i + x_i}^+}{a_{x_i + x_i}^-}$							
A	-0.740	0.417	0.580	-0.803			
С	-0.913	0.302	1.812	-0.685			
G	-0.624	0.461	0.331	-0.730			
	1.170	0.550	0.202	0.670			















Switching between '+' & '-' States

- The maximum scoring path receives a score of 0.0032.
- The most likely state path is found to be C+G+C+G+.
- Given a much longer sequence, the derived optimal path will switch between the CpG and non-CpG states.

HMM for Gene Prediction

- Hidden Markov Models can be applied to predict signals in the gene structure.
- More precisely, HMM can recognize
 - Start codons
 - Stop codons
 - Donor splice sites
 - Acceptor splice sites
 - 5' Poly A site













The Position Weight Matrix

- Due to the high variability of the promoters, exact methods cannot be used for identifying promoter regions.
- Instead we use a pattern search method based on frequencies called the position weight matrix method.
- Use the known promoter regions to construct a table of statistics, where an entry is the frequency of a certain base at a given position.

1110		I DOX	ог <u>н.</u> с	
	A	С	G	Т
1	0.04	0.09	0.07	0.80
2	0.88	0.03	0.01	0.08
3	0.26	0.11	0.12	0.5
4	0.59	0.13	0.16	0.13
5	0.49	0.22	0.12	0.18
6	0.03	0.05	0.02	0.89

















VEIL: Preparing the Training Data They used a set of 570 vertebrate sequences from different species each containing exactly one gene. Quality control step: filter the data by removing pseudogenes, entries with no introns (from cDNA), entries with non-standard splice junctions (introns that do not begin with GT

• The 570 vertebrate sequences contain a total of 2649 exons and 2079 introns.

VEIL: Five-Fold Cross Validation A five-fold cross validation experiment was performed to estimate how well the system would perform when tested on data that was not in the training set. The 570 sequences were randomly partitioned into five sets of 114 sequences each. For each partition, the system is trained on 4 sets and tested on the fifth.

– Combine the results from the five test sets.

and end with AG).













Knowledge Based Gene Prediction

Use characteristics that are known from annotated sequences to predict the genes in unannotated sequences.

- Codon Usage
- Dinucleotide FrequenciesTrinucleotide Frequencies
- Base Composition
- Splice Sites
- CpG Distribution
- Promoters
- Poly A signals
- Transcription Signals
- Translation Signals
- Hexanucleotide
 Frequencies
- Branchpoint Distribution
- Enhancers
- als Transcription Factors
- als Kozak Sequence













- Gene finding is genome-specific : software have to be adapted and trained for each genome.
- The best software for species A (e.g. GenScan for human) is not necessarily the best for species B.

Conclusion (II) The more computers are involved in automating genome annotation, the greater the need for collaboration with biologists. Best methods are based on detailed probabilistic models that include all known effects.

• We are still a long way from having reliable tools for deducing biological functions from sequences.